

Towards Machine-understanding of Document Images

Thesis submitted in partial fulfillment
of the requirements for the degree of

Doctor of Philosophy
in
Computer Science and Engineering

by

Minesh Mathew

201307668

`minesh.mathew@research.iiit.ac.in`



International Institute of Information Technology

(Deemed to be University)

Hyderabad - 500 032, INDIA

January, 2024

Copyright © Minesh Mathew, 2024

All Rights Reserved

International Institute of Information Technology
Hyderabad, India

CERTIFICATE

It is certified that the work contained in this thesis, titled “**Towards Machine-understanding of Document Images**” by Minesh Mathew, has been carried out under my supervision and is not submitted elsewhere for a degree.

Date

Adviser: Prof. C V Jawahar

To,
My parents, Martil and Mathew

Acknowledgments

I am grateful to my guide, lab mates, and friends for their support throughout my PhD journey.

Abstract

Imparting machines the capability to understand documents like humans do is an AI-complete problem since it involves multiple sub-tasks such as reading unstructured and structured text, understanding graphics and natural images, interpreting visual elements such as tables and plots, and parsing the layout and logical structure of the whole document. Except for a small percentage of documents in structured electronic formats, a majority of the documents used today, such as documents in physical mediums, born-digital documents in image formats, and electronic documents like PDFs, are not readily machine readable. A paper-based document can easily be converted into a bitmap image using a flatbed scanner or a digital camera. Consequently, machine understanding of documents in practice requires algorithms and systems that can process document images—a digital image of a document.

Successful application of deep learning-based methods and use of large-scale datasets significantly improved the performance of various sub-tasks that constitute the larger problem of machine understanding of document images. Deep-learning based techniques have successfully been applied to the detection, and recognition of text and detection and recognition of various document sub-structures such as forms and tables. However, owing to the diversity of documents in terms of language, modality of text present (typewritten, printed, handwritten or born-digital), images and graphics (photographs, computer graphics, tables, visualizations, and pictograms), layout and other visual cues, building generic-solutions to the problem of machine understanding of document images is a challenging task. In this thesis, we address some of the challenges in this space, such as text recognition in low-resource languages, information extraction from historic/handwritten collections and multimodal modeling of complex document images. Additionally, we introduce new tasks that call for a top-down perspective—to understand a document image as a whole, not in parts—of document image understanding, different from the mainstream trend where the focus has been on solving various bottom-up tasks. Most of the existing tasks in Document Image Analysis (DIA) deal with independent bottom-up tasks that aim to get a machine-readable description of certain, pre-defined document elements at various abstractions such as text tokens or tables. This thesis motivates a

purpose-driven DIA wherein a document image is analyzed dynamically, subject to a specific requirement set by a human user or an intelligent agent.

We first consider the problem of making document images printed in low-resource languages machine-readable using an OCR and thereby making these documents AI-ready. To this end, we propose to use an end-to-end neural network model that can directly transcribe a word or line image from a document to corresponding Unicode transcription. We analyze how the proposed setup overcomes many challenges to text recognition of Indic languages. Results of our synthetic to real transfer learning experiments for text recognition demonstrate that models pre-trained on synthetic data and further fine-tuned on a portion of the real data perform as well as models trained purely on real data. For 10+ languages for which there have not been public datasets for printed text recognition, we introduce a new dataset that has more than one million word images in total. We further conduct an empirical study to compare different end-to-end neural network architectures for word and line recognition of printed text.

Another significant contribution of this thesis is the introduction of new tasks that require a holistic understanding of document images. Different from existing tasks in Document Image Analysis (DIA) that attempt to solve independent bottom-up tasks, we motivate a top-down perspective of DIA that requires a holistic understanding of the image and purpose-driven information extraction. To this end, we propose two tasks—DocVQA and InfographicVQA—fashioned along Visual Question Answering (VQA) in computer vision. For DocVQA, we show results using multiple strong baselines that are adapted from existing models for existing VQA and QA problems. For InfographicVQA, we propose a transformer-based, BERT-like model that jointly models multimodal—vision, language, and layout—input. We conduct open challenges for both tasks, attracting hundreds of submissions so far.

Next, we work on the problem of information extraction from a document image collection. Recognizing text from historical and/or handwritten manuscripts is a major challenge to information extraction from such collections. Similar to open-domain QA in NLP, we propose a new task in the context of document images that seek to get answers for natural language questions asked on collections of manuscripts. We propose a two-stage retrieval-based approach for the problem that uses deep features of word images and textual words. Our approach is recognition-free and returns image snippets as answers to the questions. Although our approach is recognition-free and consequently oblivious to the semantics of the text in the documents, it can look for documents or document snippets that are lexically similar to the question. We show that our approach is a reasonable alternative when using text-based QA models is infeasible due to the difficulty in recognizing text in the document images.

Contents

Chapter	Page
1 Introduction	1
1.1 Scope	2
1.2 Applications	4
1.3 Machine understanding of images and text: existing tasks and perspectives	4
1.4 Motivation	8
1.5 Challenges	13
1.6 Problem space	14
1.7 Our major contributions	15
1.8 Re-evaluating thesis locus and contributions at the dawn of an LLM era.	18
1.9 Thesis outline	21
1.10 Publications	22
2 Background	25
2.1 Text recognition in the deep learning era	25
2.1.1 Feature learning using CNNs.	25
2.1.2 Text recognition as a structured prediction using RNNs	27
2.1.2.1 CTC-based text recognition	28
2.1.2.2 Encoder-Decoder Networks for text recognition	29
2.1.2.3 Multidimensional RNN	31
2.2 Question Answering for electronic text	31
2.2.1 Machine Reading Comprehension (MRC)	31
2.2.1.1 BERT for QA/MRC	32
2.2.2 Open domain QA	35
2.3 Layout-aware BERT-like models for DIA tasks	36
2.3.1 LayoutLM	37
2.4 Visual Question Answering (VQA)	38
2.4.1 Grid-based and region-based features for VQA	38
2.4.2 VQA on natural images using VisualBERT	40
2.5 VQA tasks that require reading text on images	41
2.5.1 Scene text VQA	42
2.5.1.1 LoRRA: Look, Read, Reason and Answer	43
2.5.1.2 Multimodal Multi-Copy Mesh (M4C)	43

2.5.2	Specialized VQA datasets where reading text on images is necessary.	44
2.5.2.1	VQA on images of book covers	44
2.5.2.2	VQA on charts	46
2.6	Multimodal QA tasks	46
3	CTC-based Seq2seq OCR for low resource languages	47
3.1	Introduction	47
3.2	Related works	51
3.2.1	First-generation OCRs [1970 - 2000]	51
3.2.2	Second-generation OCRs [2000 - 2012]	52
3.2.3	Third-generation OCRs [2012 - 2022]	53
3.3	Datasets	54
3.3.1	Internal dataset	54
3.3.2	<i>Mozhi</i> dataset	55
3.3.3	Synthetic dataset	57
3.4	Text recognition using CTC transcription	58
3.4.1	Extracting feature sequence	58
3.4.2	Encoder	60
3.4.3	Feature and encoder configurations	60
3.4.4	Decoder	62
3.4.5	Transcription using CTC	63
3.4.6	Training	63
3.4.7	Inference	64
3.5	Experimental setup	64
3.5.1	Implementation details	65
3.5.2	Evaluation	66
3.6	Experiments and results	67
3.6.1	Comparing different feature + encoder configurations	67
3.6.2	Page OCR evaluation on the internal dataset	69
3.6.3	Transfer learning from synthetic to real	71
3.6.4	Evaluating CRNN on <i>Mozhi</i> dataset	73
3.7	Summary	74
4	DocVQA: VQA on business documents	75
4.1	Introduction	75
4.2	DocVQA dataset	77
4.2.1	Data collection	80
4.2.1.1	Document images	80
4.2.1.2	Selection of workers	82
4.2.1.3	Questions and answers	83
4.2.2	Statistics and analysis	86
4.3	Baselines	89
4.3.1	Heuristics and upper bounds	90

4.3.2	VQA models	90
4.3.3	Reading comprehension models	90
4.4	Experiments	91
4.4.1	Evaluation metrics	91
4.4.2	Experimental setup	91
4.4.3	Results	93
4.5	Summary	96
5	InfographicVQA: VQA on infographics	97
5.1	Introduction	97
5.2	Related works	99
5.2.1	Question answering in a multimodal context.	99
5.2.2	Multimodal pretraining for vision and language tasks	100
5.2.3	Infographics understanding.	100
5.3	InfographicVQA	101
5.3.1	Collecting infographics	101
5.3.2	Collecting questions and answers	101
5.3.3	Question-answer types: answer-source, evidence and operation	103
5.3.4	Statistics and analysis of the dataset	104
5.4	VLL-BERT	108
5.4.1	Model architecture	109
5.4.2	Input representation	109
5.4.3	Pretraining	111
5.4.4	Finetuning	112
5.5	Experimental setup	113
5.5.1	Experimental setup for VLL-BERT	113
5.5.1.1	Extracting visual features	113
5.5.2	Experimental setup for heuristics, upper bounds, and human performance.	116
5.5.3	Experimental setup for M4C.	116
5.5.4	Evaluation metrics	116
5.6	Results	116
5.7	Summary	120
6	Recognition-free QA on handwritten document collection	122
6.1	Introduction	123
6.2	Related works	125
6.3	Recognition-free QA	125
6.3.1	Joint embedding of textual words and word images	127
6.3.2	Aggregating word embeddings	127
6.3.2.1	Aggregate by summing (SUM)	127
6.3.2.2	Aggregate using Fisher Vector (FV) framework	128
6.3.3	Document retriever	128

6.3.4	Answer snippet extraction	129
6.4	HW-SQuAD and BenthamQA datasets	129
6.4.1	HW-SQuAD	130
6.4.1.1	Data preprocessing	130
6.4.1.2	Synthetic rendering of SQuAD1.0 passages as handwritten document images	131
6.4.2	BenthamQA	132
6.4.3	Performance evaluation	133
6.5	Experiments and results	135
6.5.1	Implementation details	135
6.5.2	Performance of document retriever	136
6.5.3	Evaluating end-to-end answer snippet extraction	139
6.5.4	Evaluating a recognition-based QA approach on HW-SQuAD and Ben- thamQA	141
6.5.4.1	Transcribing handwritten images using a CRNN OCR	143
6.5.4.2	Evaluating TF-IDF-based document retriever	144
6.5.4.3	Evaluating TF-IDF + BERT full pipeline QA framework	144
6.6	Summary	145
7	Conclusion and future work	147
7.1	Conclusions	147
7.2	Future directions	148
	<i>Appendix A: Appendix to chapter 3</i>	150
A.1	Examples of pages in the internal page OCR dataset	150
A.2	Examples of line images in the Mozhi dataset	150
A.3	Distribution of lengths of words in the Mozhi dataset	150
	<i>Appendix B: Appendix to chapter 4</i>	163
	<i>Appendix C: Appendix to chapter 5</i>	169
	Bibliography	177

List of Figures

Figure	Page
1.1 Different types of document images	3
1.2 Challenges to DIA.	12
1.3 DocVQA vs. existing tasks in DIA.	16
2.1 Character and word classification using CNN.	26
2.2 CTC for Image to text OCR transcription.	29
2.3 Encoder-decoder style models proposed for unconstrained scene text recognition.	30
2.4 Extractive MRC task in SQuAD1.1.	32
2.5 Embeddings used in BERT.	33
2.6 Pretraining and finetuning stages of BERT.	34
2.7 Two-stage pipeline for open domain QA	35
2.8 VQA on natural images.	36
2.9 LayoutLM architecture.	38
2.10 Grid-based features and region-based features for VQA	39
2.11 VisualBERT for Visio-linguistic pretraining.	40
2.12 LoRRA for scene text VQA.	42
2.13 M4C for scene text VQA	44
2.14 VQA on book covers and charts.	45
3.1 Indic languages diversity.	49
3.2 Samples of synthetic word images created for Hindi, Malayalam, and Telugu.	57
3.3 Different CTC-based transcription network configurations.	59
3.4 Architecture of the CNN used in CNN_only and CRNN configurations.	61
3.5 Transfer learning from real to synth by finetuning the synth-only model.	72
4.1 Examples of question-answer pairs in DocVQA.	76
4.2 Distribution of documents by industry, document creation date and document type.	78
4.3 Screenshot of annotation stage 1 - question-answer collection.	80
4.4 Screenshot of annotation stage 2 - data verification.	80
4.5 Screenshot of annotation stage 3 - Reviewing answer mismatch cases.	81
4.6 Top 15 most frequent questions in DocVQA.	82

4.7	Top 15 most frequent answers in DocVQA.	83
4.8	Top 15 non numeric answers in DocVQA.	84
4.9	The 9 question types in DocVQA dataset and share of questions in each type.	84
4.10	Question, answer and OCR tokens statistics in DocVQA	85
4.11	Word clouds of answers and OCR tokens in DocVQA dataset	88
4.12	Distribution of questions by their starting 4-grams.	89
4.13	Qualitative results of baselines on DocVQA.	95
4.14	Performance of baseline models for different question types.	96
5.1	Examples of question-answer pairs in InfographicVQA.	98
5.2	Distribution of questions in validation set of InfographicVQA by QA types.	102
5.3	Word cloud of answers and OCR tokens in InfographicVQA.	105
5.4	Visualization of the top 20 topics in the InfographicVQA dataset.	107
5.5	Starting 4-grams of common questions in InfographicVQA.	108
5.6	Input representation for the VLL-BERT.	109
5.7	Pretraining LayoutLM-based model using MVLM	112
5.8	Finetuning LayoutLM-based model for extractive VQA.	113
5.9	Qualitative results of VQA on InfographicVQA.	120
5.10	Performance of VLL-BERT and other baselines for different question-answer types.	121
6.1	Recognition-free QA on document images collection.	123
6.2	The proposed recognition-free approach to QA on document collection.	126
6.3	Sample documents in the newly introduced HW-SQuAD and BenthamQA datasets.	129
6.4	A question-answer pair and the associated document from the test set of HW-SQuAD dataset.	134
6.5	Number of document proposals(N) vs top-N accuracy	138
6.6	Top-5 accuracy for document retriever for different sizes of FV	139
6.7	Qualitative results of end-to-end answer snippet extraction.	141
6.8	Snippet accuracy vs. number of document proposals.	142
6.9	Snippet accuracy vs. question length.	142
A.1	Assamese, Bangla, Gujarati and Hindi samples from our internal OCR dataset.	151
A.2	Kannada, Malayalam, Manipuri and Marathi samples from our internal OCR dataset.	152
A.3	Odia, Punjabi, Tamil and Telugu samples from our internal OCR dataset.	153
A.4	A sample page image of Urdu from our internal OCR dataset.	154
A.5	Samples of Assamese, Bangla, Gujarati and Hindi text-line images in the newly introduced <i>Mozhi</i> dataset for text recognition.	155
A.6	Samples of Kannada, Malayalam, Manipuri and Marathi text-line images in the newly introduced <i>Mozhi</i> dataset for text recognition.	156
A.7	Samples of Odia, Punjabi, Tamil, Telugu and Urdu text-line images in the newly introduced <i>Mozhi</i> dataset for text recognition.	157

A.8	Distribution of word lengths of Assamese, Bangla and Gujarati samples in the new Mozhi OCR dataset.	158
A.9	Distribution of word lengths of Hindi, Kannada and Malayalam samples in the new Mozhi OCR dataset.	159
A.10	Distribution of word lengths of Manipuri, Marathi and Odia samples in the new Mozhi OCR dataset.	160
A.11	Distribution of word lengths of Punjabi, Tamil and Telugu samples in the new Mozhi OCR dataset.	161
A.12	Distribution of word lengths of Urdu samples in the new Mozhi OCR dataset. .	162
B.1	DocVQA: Examples where BERT QA model answers questions other than ‘running text’ type.	164
B.2	DocVQA: M4C’s performance on questions based on pictures or photographs. .	165
B.3	DocVQA: Contrasting results for similar questions.	166
B.4	DocVQA: Questions based on figures and diagrams.	167
B.5	DocVQA: Impact of OCR errors.	168
C.1	InfographicVQA: A question where color codes and information on a Map are required to arrive at the answer.	170
C.2	InfographicVQA : An example where answer is Question-span	171
C.3	InfographicVQA: An example for multi-Span answer.	172
C.4	InfographicVQA : Example where symbols/markers need to be counted to find an answer.	173
C.5	InfographicVQA : An example where values shown in a bar chart need to be sorted to find the answer.	174
C.6	InfographicVQA: Question requiring arithmetic operation.	175
C.7	InfographicVQA: Performing multiple discrete operations to find the answer. .	176

List of Tables

Table	Page
1.1 Existing tasks that concern with understanding of text, natural images and document images.	6
1.2 Where do our works fit in the existing tasks hierarchy in DIA ?	19
3.1 Statistics of the 1000 pages internal dataset.	55
3.2 Statistics of the new Mozhi dataset	56
3.3 Recognition-only results on the validation splits of the internal dataset.	67
3.4 Recognition-only performance using CRNN on the test splits of the internal dataset.	68
3.5 Performance of our page OCR pipelines compared to other public OCR tools .	70
3.6 Synthetic to real transfer learning results.	71
3.7 CRNN evaluation on Mozhi dataset.	73
4.1 DocVQA annotation instructions.	79
4.2 Results of different upper bounds and heuristic baselines	92
4.3 Results of Scene text VQA models on DocVQA dataset	93
4.4 Performance of BERT QA model on DocVQA.	94
5.1 Summary of VQA and Multimodal QA datasets where text on the images needs to be read to answer questions.	99
5.2 Top 20 topics in InfographicVQA found using LDA.	106
5.3 Statistics of questions, answers, and OCR tokens in InfographicVQA and other similar VQA datasets.	107
5.4 Results of heuristics, upper bounds and human performance evaluation on InfographicVQA.	117
5.5 Performance of different variants of the M4C on InfographicVQA	118
5.6 Performance of VLL-BERT on InfographicVQA.	119
5.7 Statistics of object detections on InfographicVQA.	121
6.1 HW-SQuAD and BenthamQA dataset statistics.	131
6.2 Top-5 accuracy for the document retriever when the aggregation scheme used is SUM.	136

6.3	Top-5 accuracy (%) for the document retriever when FV aggregation is used.	137
6.4	Performance evaluation of end-to-end answer snippet extraction.	140
6.5	Performance of CRNN-based OCR on HW-SQuAD and BenthamQA.	143
6.6	Results of a TF-IDF-based document retriever on transcriptions of the documents in HW-SQuAD and BenthamQA.	144
6.7	Results of text-based full pipeline QA model.	145

Chapter 1

Introduction

What came first, documents or writing? Suzanne Briet, who is called ‘Madame Documentation’ defines a document as ‘*any concrete or symbolic indexical sign[indice], preserved or recorded, towards the ends of representing, of reconstituting, or of proving a physical or intellectual phenomenon.*’ [1]. If we go by Briet’s definition, documents existed long before writing systems were invented. According to recent studies, our species, and possibly Neanderthals, started making non-figurative drawings in caves 64,000 years ago [2]. The next major development in the history of documents is clay tablets engraved with pictograms and ideograms. The Proto-Cunieform system that was used by the Mesopotamian civilization in around 3200 BCE is a good example of this style of proto-writing. The symbols in proto-writing that represented objects and abstract concepts were gradually replaced by symbols that correspond to sounds in spoken language [3]. While the earliest documents were on rocks and clay or wooden tablets, by the beginning of the Common Era(CE), documents in the form of parchments, papyrus, and strolls were common. By 1500 CE, printing presses were a common sight in cities in western Europe. The invention and quick spread of printing technology in Europe led to a printing revolution [4]. This resulted in faster and broader knowledge sharing and played a crucial role in the Renaissance and the scientific revolution.

Half a millennium since the printing revolution, documents continue to act as the primary medium of knowledge dissemination. While we continue to create documents in digital and physical forms, they are primarily made for human consumption. The documents in the physical medium can easily be made digitally accessible by converting them to digital images using a scanner or other imaging devices. But neither these images of physical documents nor the born-digital document images are machine readable—i.e., in a structured format that a computer can readily process. Research on making images of documents machine readable has come a long way, thanks to recent developments in supervised machine learning algorithms. While there has been significant progress in independent information extraction tasks such as

recognizing printed text from document images and extracting regions of interest, research on Artificial Intelligence (AI) that can facilitate a holistic, human-like understanding of document images is in a nascent stage.

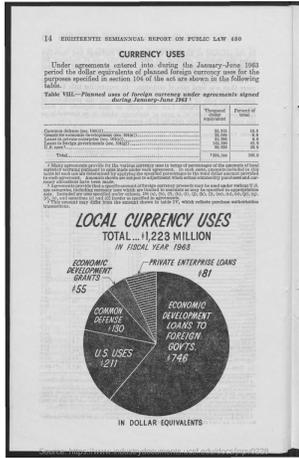
1.1 Scope

A document image can either be a digital image of a physical document or a born-digital document in image format. Different from scene text images that are natural images—or images in the wild—that contain text, a document image predominantly contains a document. The text on these images can either be unstructured text arranged in sentences and paragraphs or structured text that is part of forms, tables, figures, or other visual elements that constitute a document image. For example, in Figure 1.1a, there are i) text that makes up titles, and headers, ii) a table where text is structured in rows and columns, ii) text appearing as sentences or paragraphs, and iv) text used for labelling and annotating a figure.

Most document images created and shared today can be categorized into 4 types. The first type is scanned images of paper documents used in offices. An example of this type is shown in Figure 1.1a. Although the use of paper documents has been declining since the beginning of this century¹², paper documents are still widely used in businesses. In 2020, nearly 100 million metric tons of graphic papers were produced despite the slowdown caused by the COVID-19 pandemic [10]. The second category is document images created as a result of digitization efforts aimed at preserving cultural heritage. Archival and historical documents are being digitized in large amounts around the globe [11, 12, 13, 14, 15, 16]. An image of a manuscript written by Jeremy Bentham that was digitized as part of the Transcribe Bentham [17] project is shown in Figure 1.1d. Digitization options offered by smartphones are immense, and we take photos of paper documents such as handwritten notes and receipts on a daily basis. These images constitute the third category. An example of this type—a photograph of a printed advertisement put up on a wall—is shown in Figure 1.1e. The fourth category is born-digital document images such as visualisations, infographics, digital posters, and screenshots that contain text and/or graphics. One such image is shown in Figure 1.1f. Every time we take a screenshot of a tweet or a web page snippet, we create a document image. Although all social media platforms offer options to share links, we find it often convenient to take screenshots for quick sharing in our chats, social media stories, and statuses.

¹https://www.statista.com/statistics/270317/production-volume-of-paper-by-type/?_ga=2.204976272.1018593204.1622837058-191240632.1618425162

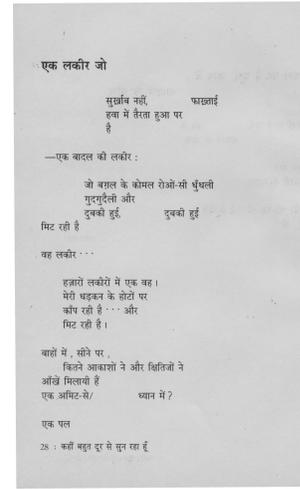
²<https://www.statista.com/statistics/1241313/graphics-paper-production-worldwide/>



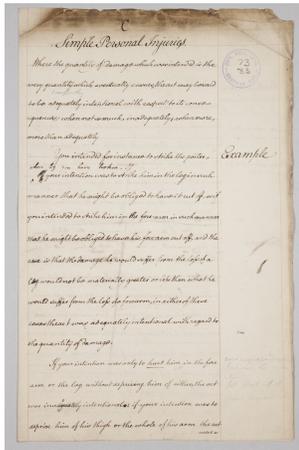
(a)



(b)



(c)



(d)



(e)



(f)

Figure 1.1: **Different types of document images.** Images (a) and (b) are pages of a financial report and a magazine, respectively. Image (c) shows a page from a Hindi book scanned using a flatbed scanner. In (d) is a scanned manuscript from Bentham manuscripts collection. Shown in (e) is a photo of an advertisement in Malayalam captured using a mobile phone camera and (f) shows a born-digital infographic. All these images are taken from datasets created as part of our works. Images (a) and (b) are a part of the DocVQA dataset [5]. Images (c), (d), (e) and (f) are part of IIIT_Hindi_100 [6], BenthamQA [7], IIIT-ILST [8] and InfographicVQA [9], respectively.

On the other hand, electronic documents, unlike the name suggests, are not readily machine-readable. A majority of the electronic documents are in PDF format. Although PDF is the most widely used format to create documents, the format does not preserve the logical structure of the documents. Finding the reading order of a PDF document itself is non-trivial. Con-

sequently, analysis of a born-digital PDF (also called ‘native’ PDF) is nearly as difficult as analysis of scanned PDF [18, 19].

1.2 Applications

With the explosive growth in the creation and dissemination of document images, and non-machine readable electronic documents, the need for algorithms and systems that can read, understand and analyze these documents is ever more imperative. Such solutions have many applications, including some major ones we list below.

1. **Searching and indexing.** Reading—i.e., recognizing text tokens— on document images would allow indexing, searching and retrieving of document images [20, 21]. For example, such an application would allow a user easily find a particular invoice the user is looking for from hundreds of photos of invoices stored in the user’s phone.
2. **Automated data entry.** Transforming data in document images to structured digital formats can facilitate automatic data entry [22, 23, 24, 25]. For example, the manual effort in entering details filled in a paper-based form into a database can be avoided if there are robust text recognition and form understanding solutions.
3. **Assistive systems.** Recognizing text and converting data in document images into structured representations has applications in building assistive systems for visually impaired people. OCR is essential for making non-native/scanned PDFs work with screen readers [26, 27]. Similarly, many books, including textbooks and most of the historical material being digitized are not accessible to the visually impaired since Braille or audiobooks are not available for these books. With the help of reading systems developed for document images and text-to-speech technology, it is possible to generate audiobooks for these books with little or minimum human effort [28, 29, 30].

1.3 Machine understanding of images and text: existing tasks and perspectives

Document Image Analysis (DIA) is a field of study that seeks to build algorithms and systems that can process and analyse document images. At the heart of DIA is developing AI

that can ‘understand’ document images as humans do. Since understanding a document image involves reading and analysing text and interpreting the visual elements, the field of DIA is closely related to Computer Vision (CV) and Natural Language Processing (NLP). In Table 1.1, we list various tasks in NLP, CV, and DIA that are critical to ‘understanding’ of natural language, natural images, and document images, respectively. We arrange the tasks in a 4-level hierarchy reflective of the level of understanding involved in solving a task.

Level 1 tasks, in general, deal with assigning semantic labels or attributes to parts or whole of the given data. In the case of the text (i.e., machine readable natural language text), these tasks deal with assigning semantic labels to words, phrases, or sentences. Tasks such as Named entity recognition [31], Part of speech tagging [32, 33], intent detection [34] and sentiment analysis [35, 36] fall in this category. In the case of natural images and document images, classification and localization tasks fall in this category. For example, image classification [37, 38, 39] and the task of locating and identifying visual objects [40, 41, 42] in natural images are examples of level 1 tasks. Similar tasks in DIA deal with word or line detection [43, 44], text recognition of cropped characters, words or lines [45, 46] and document layout analysis—the problem of localizing and labelling different document sub-structures such as paragraph, title, table and figure [47, 48].

Level 2 tasks model relationships between the atomic entities, and work at the level of groups of these entities. In the case of text, level 2 tasks deal with inferring how two entities in a piece of text are related. Examples include Relation prediction [71] and Coreference resolution [72, 73]. In the case of natural images, visual relationship detection [74] and action recognition in still images [80, 81] are two tasks that involve modelling relationships between objects and other semantic entities in natural images. In the case of document images, tasks of this type include emotion recognition in comics using multimodal context [75], associating related graphics and text elements in charts [77] and associating speech balloons with corresponding speakers in comics [76].

Level 3 tasks across the three fields aim to i) generate a description/summary of the given text/image in natural language or ii) convert the given image or text into a structured machine readable form. In the case of natural language understanding, the most popular task that falls under level 3 is text summarization [61]. Similarly, image captioning [66] and scene graph generation [65] are examples of level 3 tasks for natural images. Tasks in DIA that aim to convert various document substructures to machine readable form, such as Table to HTML [23], Chart to raw table [70], and Form understanding [22, 69], are examples of level 3 tasks for document images. Note that the text summarization task is nothing but asking an AI model to write a summary of one or more documents. This task is inspired by the summary writing

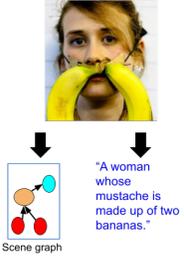
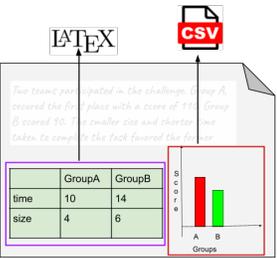
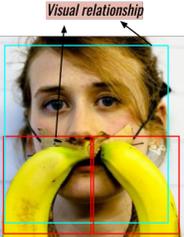
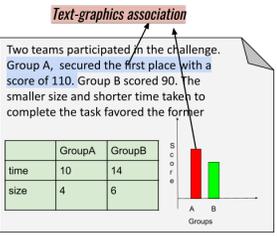
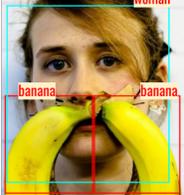
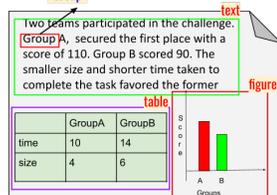
	Text	Natural images	Document images
Level 4	<p>Q: What is the valley called? A: Chetali</p> <p>A girl with silver anklets and eyes prettied with surma came to Chetali's valley to gather flowers. The Chempaka tree stood alone- efflorescent, serene. The flower gatherer reached out and held down a soft twig to pluck the flowers. As the twig broke the Chempaka said, "My little sister you have forgotten me !"</p>	<p>Q: What is the mustache made of? A: bananas</p> 	<p>Q: Which group has highest score? A: Group A</p> 
<ul style="list-style-type: none"> purpose driven interpretation of the data often require reasoning skills and external knowledge 	QA [49, 50, 51] , Dialogue systems [52, 53, 54]	VQA [55, 56], Visual Dialogue [57]	Chart VQA [58, 59], OCR-VQA [60]
Level 3	<p>A girl with silver anklets and eyes prettied with surma came to Chetali's valley to gather flowers. The Chempaka tree stood alone- efflorescent, serene. The flower gatherer reached out and held down a soft twig to pluck the flowers. As the twig broke the Chempaka said, "My little sister you have forgotten me !"</p> <p>"A girl meets her elder sister who now is a tree at the valley of Chetali."</p>	 <p>"A woman whose mustache is made up of two bananas."</p>	
<ul style="list-style-type: none"> model high-level semantics generate a summary or a structured representation for the whole 	Summarization [61, 62, 63]	Scene graph generation [64, 65], Image captioning [66, 67]	Table recognition [68, 23], Form understanding [22, 69] Chart to raw table [70]
Level 2	<p>Coreference</p> <p>A girl with silver anklets and eyes prettied with surma came to Chetali's valley to gather flowers. The Chempaka tree stood alone- efflorescent, serene. The flower gatherer reached out and held down a soft twig to pluck the flowers. As the twig broke the Chempaka said, "My little sister you have forgotten me !"</p>	<p>Visual relationship</p> 	<p>Text-graphics association</p> <p>Two teams participated in the challenge. Group A, secured the first place with a score of 110. Group B scored 90. The smaller size and shorter time taken to complete the task favored the former</p> 
<ul style="list-style-type: none"> model relations between atomic entities identify groups, structures and relations 	Relation prediction [71], Coreference resolution [72, 73]	Visual relationship detection [74]	Emotion recognition in comics [75], Text and graphics association [76, 77]
Level 1	<p>POS = "verb" Named Entity</p> <p>A girl with silver anklets and eyes prettied with surma came to Chetali's valley to gather flowers. The Chempaka tree stood alone- efflorescent, serene. The flower gatherer reached out and held down a soft twig to pluck the flowers. As the twig broke the Chempaka said, "My little sister you have forgotten me !"</p> <p>Sentiment == "Negative"</p>		<p>"Group" text figure</p> <p>Two teams participated in the challenge. Group A, secured the first place with a score of 110. Group B scored 90. The smaller size and shorter time taken to complete the task favored the former</p> <p>table</p> 
<ul style="list-style-type: none"> Assign labels and attributes Segment and recognize 	NER [31], POS tagging [32, 33], Intent detection [34], Sentiment analysis [35, 36]	Object detection & semantic segmentation [40, 42]	Layout analysis [47, 48], Text recognition [78, 46, 79]

Table 1.1: Existing tasks that concern with understanding of text, natural images and document images.

exercises for school children. On the other hand, the image captioning task for natural images is similar to the ‘show and tell’ exercise commonly used in elementary schools.

The most distinguishing characteristic of level 4 tasks is the dynamic interpretation of given data, conditioned on a requirement or a query set by a human user. A natural language question or a dynamic requirement raised during a conversation between a human user and an intelligent system dictates how the image or text needs to be analysed and what response needs to be returned. Question Answering (QA) [49, 50] and Machine Reading Comprehension (MRC) [51] for text and VQA on natural images [55, 56] are examples of level 4 tasks. Chart VQA [58, 59] and VQA on book covers [60] are prior works that deal with VQA on document images. The former is a VQA task for a limited number of standard charts such as bar charts, line plots, and pie charts, and the latter deals exclusively with book covers. From a machine understanding perspective, level 4 tasks that are most challenging are QA/VQA tasks that require reasoning abilities [82, 83].

The advancement made in AI in recent years has primarily been driven by ‘deep learning’ that strives for better performance by employing bigger neural networks and large-scale data. Two fields where AI using deep learning techniques have been utilized effectively are CV and NLP. Following the success of AlexNet [84] for image classification, Convolutional Neural Networks (CNN) was quickly adopted for almost every task that involves images. Similarly, Recurrent Neural Networks (RNN)—particularly the Long Short Term Memory (LSTM) and Gated Recurrent Unit (GRU) variants—emerged as the de facto choice for many sequence learning problems such as speech recognition [85], handwriting recognition [86] and machine translation [87, 88].

Deep learning-based approaches significantly improved performance for almost all tasks for text and natural images that we list under level 1 and level 2 of the tasks hierarchy we discussed above(see section 1.3). These results inspired researchers to attempt newer tasks—that fall in level 3 and level 4 in the hierarchy—that aim for a high-level, human-like understanding of text and images. Such tasks were largely unexplored or lacked large-scale datasets until this period. Popular datasets for tasks such as Question Answering (QA), image captioning and Visual Question Answering (VQA) were introduced around this time. The most widely used dataset for image captioning—MS-COCO captions [89, 90], VQA V1 dataset [55] for VQA, SQuAD [91] for Machine Reading Comprehension (MRC) and WikiQA dataset [92] for open domain QA were all introduced during the 2015-16 period. Robust performance for constituent lower-level tasks is indispensable for most of these higher-level tasks. For example, almost every VQA [93, 94, 56] or captioning model [93, 66] makes use of a CNN pretrained for image classification [95] or object detection model [96] in order to extract features from the images.

Graves et al. [86] were the first to use an end-to-end neural network for unconstrained handwriting recognition. Their approach learns to map directly from handwritten text lines to a sequence of characters without the need for explicit alignment between the characters in the image and characters in the ground truth label sequence. This approach has then been adopted for printed text recognition [97, 98, 99] and scene text recognition [100, 46] in multiple languages. The success of CNN for CV tasks inspired the use of CNN-based methods for many similar tasks in DIA, such as word image representation learning [101, 102], document image classification [103, 104, 105], text detection [106] and layout analysis [47, 48]. Improved performance for text recognition and detection, motivated researchers to study higher-level DIA problems that deal with the machine readable representation of forms, tables and charts [70, 23, 22] (level 3 tasks as per the hierarchy we discuss in section 1.3). These tasks are aimed at extraction of a machine readable representation and/or high-level understanding of only a certain document sub-structure (a table, a form, or a chart) that the task is designed for. Similarly, in the case of level 4 tasks concerning document images, both the tasks expect certain types of document images—charts for chart VQA [58, 59] and book covers for OCR VQA [60]—as input and the datasets contain template questions with limited diversity. In the case of OCR-VQA, the questions are generated using metadata associated with the books, such as book title, author, and book genre. Consequently, the questions do not demand reasoning over the visual, layout, and graphical information on the images.

Document images like the ones shown in Figure 1.1a, Figure 1.1b, and Figure 1.1f are complex multimodal manifestations of data and information. Human-like understanding of these images is an AI-complete problem that goes beyond bottom-up tasks that deal with the detection and recognition of text or visual elements in them. However, we believe that progress made in solving the bottom-up tasks now permits us to explore problems that require a holistic machine understanding of document images. On the other hand, even today, text recognition from documents—the necessary first step to machine understanding—in low-resource languages (e.g., Figure 1.1c and Figure 1.1e) and historical and/or handwritten documents (e.g., Figure 1.1d) remains an open problem.

1.4 Motivation

In this section, the motivation behind the work done as part of this thesis is presented through a conversation between the fictional detective character Sherlock Holmes and his associate Dr. Watson.

221B Baker Street hardly had any visitors since the COVID lockdown. Except for delivery riders of Deliveroo, Mrs. Hudson has not seen anyone going upstairs since the lockdown. Homes found refuge in his Kindle and occasional violin serenade as the regular investigative avenues are temporarily closed.

Holmes has been quite restless since last week. Mrs. Hudson noticed that cigar deliveries are quite frequent. It all began with an email from Holmes's friend, Byomkesh Bakshi. Mr. Bakshi sent him details of a case he is after. The email included a link to a cloud drive full of photographs and documents relevant to the case. A theatre artist and socialite named Maya has gone missing in Kolkata. Mr. Bakshi believes her disappearance has something to do with her grandmother's writings. Her grandmother was a well-known folklore researcher. And Ms. Maya had lately been reading her grandmother's writing and had plans to make some of the folklore stories into plays. Her friends, whom Mr. Bakshi talked to, felt that she was so immersed in those stories that she started hallucinating that some of those folklore characters, particularly "yakshis" and djinns," visited her at night. Holmes instantly developed an interest in Maya's case and started browsing through the documents Mr. Bakshi sent. However, they were not in any particular order. Mr. Bakshi sent him a folder containing hundreds of PDFs, scanned images, and photographs. Holmes suddenly remembered that Dr. Watson has been doing online courses on computing and machine learning. He decided to seek his help and sent the link to Dr. Watson.

It has almost been a week since the details were sent to Watson. Holmes could not make any progress skimming through those documents. Impatient, at 3 in the morning, he lit another cigar and pinged Watson on Discord.

HOLMES: Watson, any updates?

WATSON: I have already sent you a link to a web-based tool where you can search within the documents and even ask questions about the content of the documents.

HOLMES: It is definitely helpful. However, it works only for a small percentage of the documents that Mr. Bakshi sent me.

WATSON: Yes. Those are Ms. Maya's notes she used to make online and some of her blog posts. I have even run an OCR to convert some of the scanned PDFs to electronic text. The results were pretty good. You can query them as well.

HOLMES: Alright! What about her journal? The drive includes photographs of pages from her journal.

WATSON: I am working on those. Those are handwritten. The typical English OCRs won't work well. I am using a handwriting recognition model for it. The results are ok. There are some errors, but you must be able to search for keywords.



HOLMES: If there is an OCR and handwriting recognition system, why are you not OCRing other documents like those old family records and her grandmother's diary?

WATSON: Those family records are at least a hundred years old. They are in English, but the documents are in a pretty bad state. The handwriting recognition systems perform poorly on those.

HOLMES: You've got to do something, Watson. After all, you have been spending so much money on ML courses. I consider you an ML wizard. It will be great if I can query those family records. It is tedious to manually go through thousands of those pages. What about the Question Answering model you were talking about? Forget it! Can you at least help me find a tool that can locate the pages that might be relevant to my queries.?

WATSON: I need to get those documents in an electronic (machine-readable) format using an OCR before you can query them. Let me see if I can build some sort of retrieval system for those old documents.

HOLMES: But how is retrieval different from a QA system?. Both sound the same to me.

WATSON: Retrieval won't get you an exact answer, like in the case of QA. It would instead return a small set of documents that match your query.

HOLMES: Something is better than nothing.

The image of Holmes working on his laptop is created using an AI text-to-image tool hosted by hotpot.ai

HOLMES: The QA model would be of great help in going through those travel bills and receipts. And the medical and lab reports as well. If I am not mistaken, an OCR would yield near perfect results on those receipts and reports. I tested one of those blood reports in an online demo of an OCR, and the result was perfect.

WATSON: I am using the BERT QA model. The reports and receipts have a lot of content in the form of tables, diagrams, and forms. The QA system is made to handle only text in the form of sentences. It cannot handle text in structured form in tables.

HOLMES: Ah, I did not know there was software by the name of a muppet. So this muppet cannot do QA over the receipts and reports? I thought structured content is easier for machines! There is something called VQA that enables you to ask questions about images.

WATSON: Yes. In fact, one of my course assignments was to train a VQA model.

HOLMES: And it never occurred to you to train a VQA model for the reports and receipts? After all, they are images.

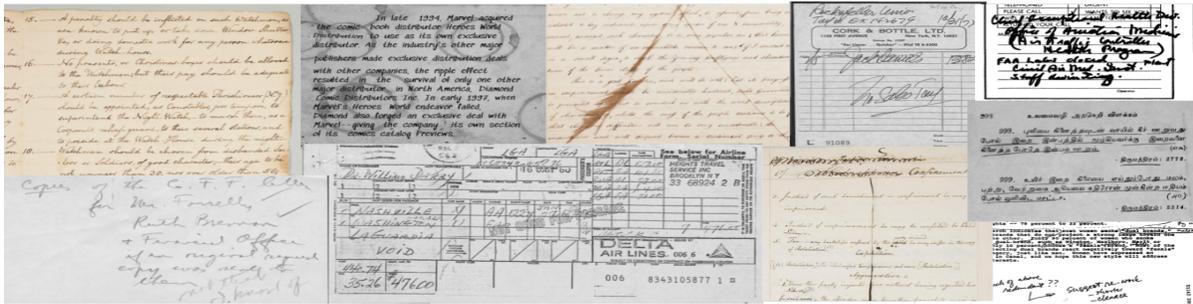
WATSON: There are no datasets that I can use to train a VQA model for images of documents. Additionally, I don't think the VQA models would work out of the box here. These are documents. VQA models typically work by identifying what things and objects they see on the images and how they are related. They are great for QA on photographs of things we see around us. The documents are quite different. The documents carry tables, diagrams, photographs, and what not!

HOLMES: But you will never know. Why don't you try to find a dataset and train on it?

WATSON: Yeah! That is the plan. I wish there was a dataset to train a VQA model on images of documents.

HOLMES: By the way, I need help with those documents in Indian languages. There are many newspapers, magazines, and books that are in Bengali, Urdu, and Malayalam. I was thinking of translating those to English so that I could go through them.

WATSON: But they are scanned images! We need to first run an OCR and then translate. But the free OCRs I tried for Malayalam and Urdu were not good. I thought of training my own OCR models. But there is no data. Unlike English, these scripts are complex. I need to do a bit of research for training OCRs for these languages. Language technologies for non-Latin languages are yet to catch up with what we have for English.



(a)



(b)



(c)

Figure 1.2: **Challenges to DIA.** On one hand, even machine-reading (detection and recognition of text) can be challenging for i) handwritten or historical documents and documents that are degraded as shown in (a), and ii) text in low-resource languages (we show some examples of printed documents in low-resource languages in (b)). On the other hand, multimodal document images like the ones shown in (c) remind us that machine-understanding of document images is a complex cognitive task that goes far beyond machine-reading. All the image snippets shown in (a), (b) and (c) are from multiple datasets that are used in works done as part of this thesis.

1.5 Challenges

In the below, we list challenges to DIA that are particularly critical to the works done as part of this thesis.

1. **Degradation and imaging artifacts.** Document images are either images of physical documents, and born-digital documents. Aging and degradation of physical documents and artefacts induced during imaging pose challenges to DIA [107, 108]. The former is common in historical documents, and the latter is prevalent in the case of images captured using hand-held cameras and smartphones.
2. **Handwriting recognition.** The text on document images can be handwritten, typewritten, printed, or born-digital text. Text recognition of Latin-based languages is largely solved except for handwriting recognition. Owing to large variations in handwriting, offline handwriting recognition systems that are as good as printed OCRs are yet to be developed. We show examples of historical and degraded documents and handwritten text that are difficult to recognize in Figure 1.2a.
3. **Multimodal understanding.** In addition to text, document images contain tables, visualizations, natural images, graphical elements, and other visual cues. Understanding document images is an AI-complete problem that goes beyond extracting machine readable representation of text and visual elements in the documents. In Figure 1.2b, we show examples of complex multimodal document images such as pages of magazines and infographics.
4. **Skewed OCR development** Since most of the documents are rich in text, recognizing the text—extract text in machine readable form—is a critical requirement for making these documents machine readable. Language technologies, including OCR for complex scripts and text in low-resource languages, have traditionally been lagging behind [109]. In Figure 1.2c, we show snippets of documents in multiple Indian languages that fall under more than 2 language families. Scripts of these languages vary widely, and many inherent complexities of the script and language [110] make OCR of these languages difficult compared to English OCR. We believe the skewed development of text recognition models is a significant challenge to universal access to DIA solutions.

1.6 Problem space

This thesis addresses problems in the space of machine understanding of document images. Note that we use the term *machine understanding*, not *machine readable*, as typically used in literature. Gorman and Kasturi, in their 2002 book titled ‘Document Image Analysis a Primer’ state that Document Image Analysis (DIA) refers to ‘*algorithms and techniques that are applied to images of documents to obtain a computer-readable descriptions from pixel data*’ [111]. This observation is in line with the conclusion we draw in section 1.3 that the present-day DIA largely comprises bottom-up tasks designed to extract machine readable description of text and visual elements in document images.

What is the computer-readable or machine readable description for document images like the ones shown in Figure 1.1a, Figure 1.1b or Figure 1.1f? Would recognizing the text in the right reading order and extracting structured raw data presented in tables, tables, and plots yield a computer-readable description of a document image? For example, in the case of Figure 1.1b, we do not know how to convert the photographs and pictures in the image into a computer-readable description. Similarly, we do not know how to preserve the layout of a document image along with data, visual cues, and layout information so that the document in entirety is machine readable. However, the primary objective of making a document image machine readable is to make the machines read, see, and ultimately understand the images like humans do. This is the reason why we consciously use the term ‘machine understanding’ in contrast to ‘machine reading’ as one of the stated objectives of the works reported in this thesis

machine understanding of document images is a vast problem space and our specific areas of interest are listed below.

1. **Making documents in low resource language machine-readable.** One of the problems we are interested in is development of OCRs to make printed documents in low resource languages machine-readable. We want to investigate if data-driven, deep learning techniques can be adopted for OCR of multiple low-resource languages without the need for any script or language-specific modules.
2. **Extracting intended information from document images.** Earlier, we quoted Gorman and Kasturi, who define DIA as a field of study that builds algorithms for obtaining computer-readable description from document images [111]. In the same book, the authors state that ‘*The objective of document image analysis is to recognize the text and graphics components in images, and to extract the intended information as a human would.*’. Are the advancements made in solving various bottom-up DIA tasks equipping us to re-

alise the above-stated objective of DIA? This is one of the research questions we pursue in this thesis.

3. **Multimodal learning for document image understanding.** A complex document image can contain text—both unstructured and structured—, natural images, tables, visualizations, color schemes, and other visual cues organized in a particular layout. Therefore machine understanding of document images is essentially a multimodal problem. In this thesis, we would like to explore multimodal learning involving vision and language for machine understanding of complex document images.
4. **Information extraction from large document collections.** Electronic libraries and digitization projects have been digitizing large volumes of historical documents. Although document retrieval and searching for keywords—both text-based search and keyword spotting—on such collections are well-studied problems, obtaining specific information from a large document collection is mostly unexplored. We study this problem, particularly in the context of document collections where robust text recognition is difficult.

1.7 Our major contributions

Major contributions of this thesis are,

1. **CTC-based transcription OCR for multiple low-resource languages.** Inspired by the success of CTC-based transcription for unconstrained text recognition, we use CTC-based models for OCR of multiple Indian languages. We study the challenges to the recognition of Indic scripts and investigate how segmentation-free recognition using CTC helps to overcome many of these challenges. We empirically compare the performance of four types of CTC-based text recognition models that differ in terms of feature extraction and sequence encoding. The results of this empirical study establish the importance of using learnt features (as opposed to hand-crafted features) and language modelling for text recognition. We also compare word-level and line-level recognition models. We show that our text recognition models perform better than commercial OCR solutions on our internal dataset for 8 out of the 13 languages when used for end-to-end (i.e., document level) text recognition. This study is presented in chapter 3.
2. **New public datasets for printed text recognition in Indian languages** For the problem of printed text recognition of the Indic languages, we introduce a new dataset called

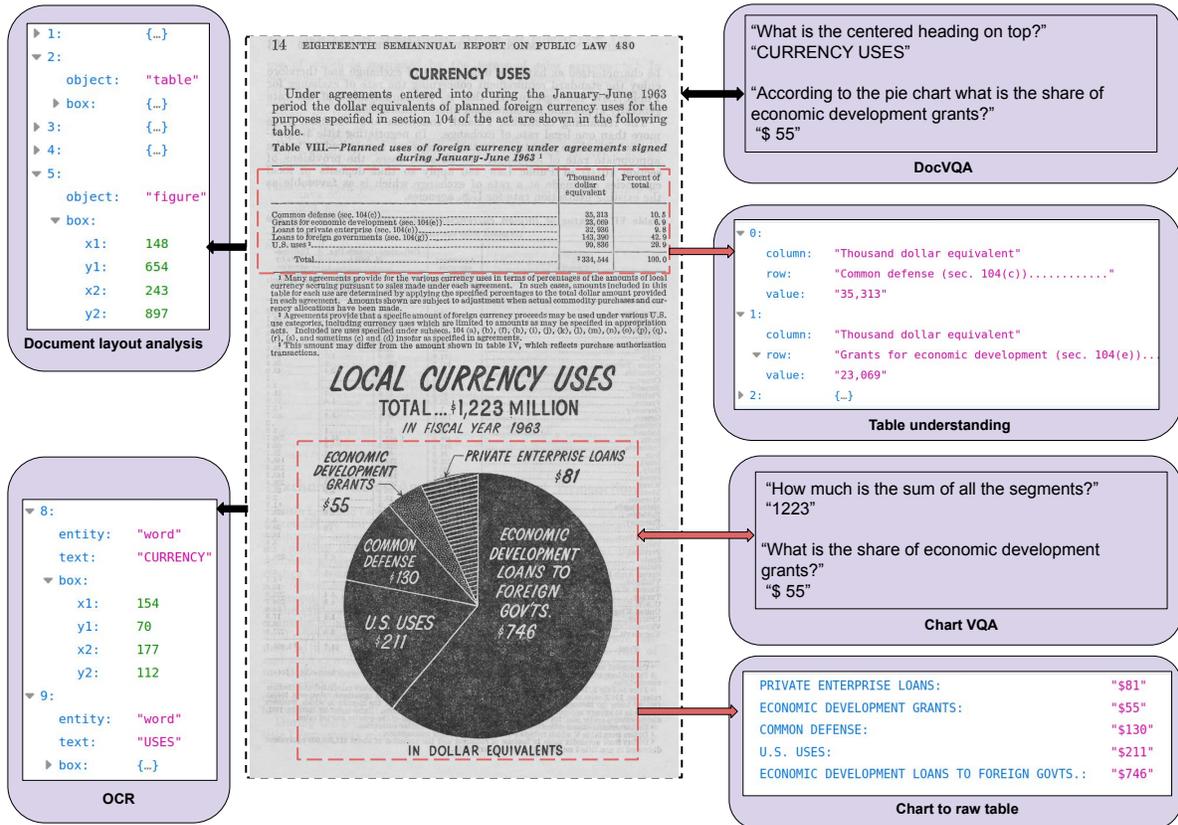


Figure 1.3: **DocVQA vs existing tasks in DIA.** Existing tasks that work at the level of documents detect and label text tokens, semantically meaningful groups of text and graphics, and visual elements like tables and charts. These tasks do not require a holistic understanding of a document image in all its details. Although there are high-level DIA tasks like Chart VQA and Table understanding that deal with high-level human-like understanding of an image of a chart or a table, most document images like the one shown here contain many document objects such as tables, plots and paragraphs within a single image. One of the major contributions of this thesis is our work called DocVQA which introduces a new task and dataset for VQA on generic document images like the one shown here.

‘Mozhi’ that has cropped word and line annotations for 13 Indian languages. The dataset has more than 1.2 million annotated word images in total, making it the largest-ever dataset for printed text recognition in Indian languages. Details of the dataset is presented in chapter 3

3. **Motivating purpose-driven DIA using Document Visual Question Answering.** In section 1.4, we note that existing tasks in DIA largely focus on solving various bottom-up tasks like detection and recognition. Contrary to this mainstream paradigm, we motivate

a top-down approach to DIA by introducing new tasks that use question answering as a proxy task to assess how well machines understand document images. Unlike chart VQA or VQA for book covers (see subsection 2.5.2), the newly introduced tasks deal with QA over generic document images, and the questions require reasoning over text, graphics, layout, and other visual aspects of the document images. Under the umbrella theme of ‘Document Visual Question Answering (DocVQA), we introduce multiple VQA tasks and datasets for document images and document image collections and organize a series of challenges and workshops in top computer vision and document analysis conferences ³.

4. **VQA on industry documents.** We introduce DocVQA, a large scale dataset of 12,767 document images of varied types and content, over which we have defined 50,000 questions and answers. The questions defined are categorised based on their reasoning requirements, allowing us to analyze how DocVQA methods fare for different question types. We evaluate two strong baselines based on state-of-the-art MRC and VQA models on the newly introduced dataset. Our baselines and initial results motivate the simultaneous use of visual and textual cues to answer questions on document images. DocVQA dataset and experimental results are presented in chapter 4.

5. **VQA on complex multimodal images.** In order to motivate the research in machine understanding of complex multimodal images, we propose a VQA task specifically for infographics that are rich in the interplay of textual and visual information. We introduce a new dataset for VQA on infographics, InfographicVQA, comprising 30,035 questions over 5,485 images. Questions in the dataset include questions grounded on tables, figures, and visualizations and questions that require combining multiple cues. Since most infographics contain numerical data, we collect questions that require elementary reasoning skills such as counting, sorting, and arithmetic operations. We believe our dataset is ideal for benchmarking progress of algorithms at the meeting point of vision, language, and document understanding. We propose a multimodal—vision, language, and layout—BERT-like model, called VLL-BERT, for InfographicVQA. Although the proposed approach performs better than strong baselines that were originally developed for MRC or VQA on natural images, there exists a huge gap in performance compared to the human performance on the dataset. InfographicVQA dataset and VLL-BERT are presented in chapter 5.

³<https://www.docvqa.org/>

6. **QA over document images collection.** Similar to open domain QA in NLP/IR, we motivate QA on a collection of document images. We particularly focus our interest on collections of historical documents or manuscripts whose OCR transcription is difficult. Since text-based QA models are not a viable option when OCR transcriptions are absent or noisy, we propose a new QA scheme that returns image snippets as responses to the questions/queries. We introduce new datasets for the task and put forward an evaluation scheme for evaluating answers returned in the form of image snippets. We present a retrieval-based solution for the new task by looking for snippets in the document collection that contain words similar to the words in the question. The new task for QA over document collection, the datasets, and the proposed recognition-free approach are presented in chapter 6.

1.8 Re-evaluating thesis locus and contributions at the dawn of an LLM era.

Recent years have witnessed the application of Transformer-based architectures to a wide variety of tasks in NLP and CV. Large Language Models (LLMs) such as BERT [113] and GPT [114] emerged as foundation models in NLP. Similarly, LayoutLM [115] is a foundation model for document image analysis. LayoutLM extends BERT to document images by incorporating layout information in addition to textual information. LayoutLM is pretrained on millions of document images without any extra annotations but only the OCR transcriptions of the documents. These are called foundation models since such models pretrained in a self supervised manner on a large scale, unlabelled data can be finetuned to perform multiple tasks. For example, a pretrained BERT can be finetuned to perform sentiment analysis and extractive QA. GPT-3 [116] is finetuned for tasks such as translation and abstractive QA. Similarly, the pretrained LayoutLM is finetuned to perform multiple downstream tasks in DIA, such as document image classification and invoice parsing. We used BERT as a text-only baseline for the DocVQA dataset we propose (see subsection 4.3.3). The VLL-BERT model we propose for VQA on infographics (see section 5.4) is inspired by LayoutLM architecture.

Different from foundation models like BERT and LayoutLM which are transformer encoder-only models, models like GPT-3 and PALM [117] are capable of generating text in an auto-regressive manner. What made LLM a buzzword is the recent popularity of their chatbot spin-offs ChatGPT ⁴ and Bard ⁵, respectively. For example, ChatGPT is a GPT-3 that is trained

⁴<https://openai.com/blog/chatgpt>

⁵<https://bard.google.com/>

	Existing DIA tasks	Our works
<p>Level 4</p> <ul style="list-style-type: none"> purpose driven interpretation of the data often require reasoning skills and external knowledge 	<p>Q: Which group has highest score ? A: Group A</p>  <p>Chart VQA [58, 59], OCR-VQA [60]</p>	<p>Q: What date is written on the seal at the top? A: 23 Sep 1970</p>  <p>DocVQA [5]</p> <p>Q: How many companies have more than 10K delivery workers? A: 2</p>  <p>Infographic VQA [9]</p> <p>Q: In which year was John McIntire murdered? A:</p>  <p>Handwritten collection QA [7]</p>
	DocVQA [5], Infographic VQA [9], Handwritten document collection QA [7]	

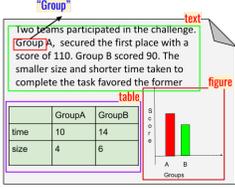
	Existing DIA tasks	Our works									
<p>Level 1</p> <ul style="list-style-type: none"> Assign labels and attributes Segment and recognize 	<p>Group</p> <p>Two teams participated in the challenge. Group A, secured the first place with a score of 110. Group B scored 90. The smaller size and shorter time taken to complete the task favored the former.</p> <table border="1"> <thead> <tr> <th></th> <th>GroupA</th> <th>GroupB</th> </tr> </thead> <tbody> <tr> <td>time</td> <td>10</td> <td>14</td> </tr> <tr> <td>size</td> <td>4</td> <td>6</td> </tr> </tbody> </table>  <p>Layout analysis [47, 48], Text recognition [78, 46, 79]</p>		GroupA	GroupB	time	10	14	size	4	6	<p>लक्का-सुन्नी / 85</p> <p>समेट कर बापस मूँह में डालती है। इसारे पास एक युवा माँ-बाप अपने बेटे के साथ आदसक्रीम खा रहे हैं। बच्चा चकरधिन्नी की तरह ऊपर-नीचे कूद रहा है। ता कुछ पल उसे निमिभेष ताकती है—“जिसके माँ-बाप नहीं होते, उनको बामते हो, क्या कहते है?”</p> <p>“नहीं,” मैं कहता हूँ। मेरी आँखों में प्ल किरकिराती है।</p> <p>“अन्नाथ ! लक्का-सुन्नी भी अन्नाथ हैं !” वह अपनी फीली-फीली नाक में भीतर तक उँगली दे देती है। मैं उसे धीमे से बरजता हूँ। कहता हूँ कि लक्का-सुन्नी भी गन्दी आदत सीख जायेंगे ऐसे।</p> <p>Printed text recognition of Indic languages [112]</p>
	GroupA	GroupB									
time	10	14									
size	4	6									

Table 1.2: Where do our works fit in the existing tasks hierarchy in DIA ? Our works on OCR of Indic languages fall in level 1 of the hierarchy we discuss in Table 1.1. Lack of resources and challenges in the input and labels space make OCR for Indic languages challenging and we address these challenges in our work [6, 112]. Existing QA/VQA tasks (level 4) in DIA focus on specific document entities such as charts [58, 59] or focus on textual content, as in the case of book cover VQA [60]. In contrast, our single document VQA tasks [5, 9] require multimodal reasoning and holistic document understanding. Our work on handwritten document collection QA [7] studies information retrieval from document collections where text recognition and text-based information retrieval are not an option.

further using Reinforcement Learning from Human Feedback (RLHF). Unlike the application specific chatbots that existed before, ChatGPT can be used in a variety of contexts and can be used as a plug-and-play module for a variety of applications. For example chatPDF⁶ is an application built using chatGPT API that allows users to ask questions based on PDFs. Similarly, there are ChatGPT-based tools for creating PPTs⁷ and code review⁸. Although the works done as part of this thesis were carried out before the introduction of these groundbreaking chatbots,

⁶<https://www.chatpdf.com/>
⁷<https://slidesgpt.com/>
⁸<https://github.com/anc95/ChatGPT-CodeReview>

we believe it is important to re-evaluate the positioning and contributions of this thesis in light of these new developments.

In the below, we note down how the thesis contributions align with a new AI paradigm that LLMs open up.

1. **Contributing to data-readiness for LLMs.** ChatGPT like agents allow us to mine information from the electronic documents, ask questions based on the data available and engage with the information in a multitude of ways. However, the engagement is possible only if the data is machine-readable. For example, the ChatPDF application we mentioned above can only converse with native PDFs. In the LLM era, we believe AI-readiness is an important aspect and building OCRs for low resource languages is a critical component of it. Developing OCRs for documents printed in low resource languages make it possible to leverage capabilities of ChatGPT like chatbots for information extraction from scanned document images in these languages.
2. **New benchmarks for the LLM era.** ChatGPT like chatbots are pushing the limits of AI. As we discussed in section 1.3, at the core of this pursuit is to make AI understand text and images as humans do. Benchmarks that help us to evaluate how well the AI systems understand multimodal content is thus necessary. Since document images are inherently multimodal, benchmarks that evaluate machine-understanding of document images serve as an excellent test bed for large multimodal foundation models. We believe the two multimodal document image VQA datasets that we introduce—DocVQA and InfographicVQA—suit the requirement. The two datasets have emerged as popular benchmark datasets for evaluating machine understanding of visually rich documents. For example, Multimodal LLM (MLLM), GPT-4⁹, benchmarks the model’s multimodal capabilities on both DocVQA and InfographicVQA datasets. Similarly, another MLLM called ChatSpot [118] makes use of DocVQA dataset for instruction tuning the model.
3. **Leveraging chain-of-thought prompting for complex multimodal reasoning.** Extracting information from infographics often requires combining multiple modalities and multiple parts of images and performing discrete arithmetic or logical operations on intermediate results. We believe the questions in InfographicVQA are suitable to hone chain of thought prompting [119] for training large multimodal models.

⁹<https://openai.com/research/gpt-4>

1.9 Thesis outline

The thesis is divided into 7 chapters, as given below

- In chapter 1 we discuss the motivation and scope of the thesis and briefly outline our problem space and our areas of interest.
- In chapter 2, we present a summary of the related works from DIA, CV and NLP space.
- In chapter 3, we present our work dealing with printed text recognition of Indic languages. We present an empirical study of various end-to-end deep learning networks for segmentation-free text recognition in 13 Indic languages and introduce new public datasets for the languages that are part of the study.
- In chapter 4, we present our work that proposes a new purpose-driven DIA task called DocVQA that involves VQA on document images. We introduce a new dataset for the task, and report results using multiple strong baselines.
- In chapter 5, we present our work that proposes a VQA task, called InfographicVQA, aimed at machine understanding of infographics that are rich in multimodal content. We propose a transformer-based extractive VQA model that jointly reasons over multimodal input—visual features, text and layout—for the newly introduced task.
- In chapter 6, we present our work on question answering on document collections. We introduce a new QA task that returns image snippets as answers and proposes a retrieval-based, recognition-free approach for the newly introduced task. We demonstrate that the proposed approach can be an effective alternative to recognition-based QA models when robust text recognition from document images is difficult.
- In chapter 7, we present our conclusions and future directions.

1.10 Publications

Works discussed in this thesis have previously been presented as part of the following publications.

Journal publications:

- **Minesh Mathew**, Lluís Gomez, Dimosthenis Karatzas, Ernest Valveny and CV Jawahar - *Asking Questions on Handwritten Document Collections*, International Journal of Document Analysis and Recognition (IJ DAR) , ICDAR-IJ DAR special issue, 2021

Conference publications:

- **Minesh Mathew**, Dimosthenis Karatzas and CV Jawahar - *DOCVQA: A dataset for VQA on document images*, IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) 2021
- **Minesh Mathew**, Viraj Bagal, Ruben Tito, Dimosthenis Karatzas, Ernest Valveny and CV Jawahar - *Infographic VQA*, IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) 2022

arXiv technical reports

- **Minesh Mathew**, and CV Jawahar - *An empirical study of CTC based models for OCR of Indian languages* - arXiv:2205.06740 - 2022

Other publications that are not part of this thesis

Conference publications:

- **Minesh Mathew**, Ajeet Kumar Singh and CV Jawahar - *Multilingual OCR for Indic scripts* , IAPR International Workshop on Document Analysis Systems (DAS) 2016
- Mohit Jain, **Minesh Mathew** and CV Jawahar - *Unconstrained OCR for Urdu using deep CNN-RNN hybrid networks*, Asian Conference on Pattern Recognition (ACPR) 2017
- Kartik Dutta, Praveen Krishnan, **Minesh Mathew** and CV Jawahar - *Towards accurate handwritten word recognition for Hindi and Bangla*, National Conference on Computer Vision, Pattern Recognition, Image Processing, and Graphics (NCVPRIPG) 2017

- Kartik Dutta, Praveen Krishnan, **Minesh Mathew** and CV Jawahar - *Offline handwriting recognition on Devanagari using a new benchmark dataset*, IAPR International Workshop on Document Analysis Systems (DAS) 2018
- Kartik Dutta, Praveen Krishnan, **Minesh Mathew** and CV Jawahar - *Improving CNN-RNN hybrid networks for handwriting recognition*, International Conference on Frontiers in Handwriting Recognition (ICFHR) 2018
- Kartik Dutta, Praveen Krishnan, **Minesh Mathew** and CV Jawahar - *Towards spotting and recognition of handwritten words in indic scripts*, International Conference on Frontiers in Handwriting Recognition (ICFHR) 2018
- Kartik Dutta, **Minesh Mathew**, Praveen Krishnan and CV Jawahar - *Localizing and recognizing text in lecture videos*, International Conference on Frontiers in Handwriting Recognition (ICFHR) 2018
- Deepayan Das, Jerin Philip, **Minesh Mathew** and CV Jawahar - *A cost efficient approach to correct OCR errors in large document collections*, International Conference on Document Analysis and Recognition (ICDAR) 2019
- Sangeeth Reddy, **Minesh Mathew**, Lluís Gomez, Marçal Rusiñol, Dimosthenis Karatzas and CV Jawahar - *Roadtext-1k: Text detection & recognition dataset for driving videos*, International Conference on Robotics and Automation (ICRA) 2020
- Yash Khare , Viraj Bagal, **Minesh Mathew**, Adithi Devi, U Deva Priyakumar and CV Jawahar - *MMBERT: multimodal BERT pretraining for improved medical VQA*, International Symposium Biomedical Imaging (ISBI) 2021
- Sergi Garcia Bordils, George Tom, Sangeeth Reddy, **Minesh Mathew**, Marçal Rusiñol, CV Jawahar and Dimosthenis Karatzas - *Read While You Drive-Multilingual Text Tracking on the Road*, IAPR International Workshop on Document Analysis Systems (DAS) 2022
- Soumya Jahagirdar, **Minesh Mathew**, Dimosthenis Karatzas and CV Jawahar - *Watching the News: Towards VideoQA Models that can Read*, IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) 2023
- George Tom, **Minesh Mathew**, Sergi Garcia Bordils, Dimosthenis Karatzas and CV Jawahar - *Reading Between the Lanes: Text VideoQA on the Road*, International Conference on Document Analysis and Recognition (ICDAR) 2023

Workshop publications

- **Minesh Mathew**, Mohit Jain and CV Jawahar - *Benchmarking scene text recognition in Devanagari, Telugu and Malayalam*, Multilingual OCR (MOCR) Workshop, International Conference on Document Analysis and Recognition (ICDAR) 2017
- Vinitha VS, **Minesh Mathew** and CV Jawahar - *An Empirical Study of Effectiveness of Post-processing in Indic Scripts*, Multilingual OCR (MOCR) Workshop, International Conference on Document Analysis and Recognition (ICDAR) 2017
- Mohit Jain, **Minesh Mathew** and CV Jawahar - *Unconstrained scene text and video text recognition for arabic script*, International Workshop on Arabic Script Analysis and Recognition (ASAR) 2017

Challenge reports:

- Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluís Gomez, Marçal Rusinol, **Minesh Mathew**, CV Jawahar, Ernest Valveny, Dimosthenis Karatzas - *ICDAR 2019 Competition on Scene Text Visual Question Answering*, International Conference on Document Analysis and Recognition (ICDAR) 2019
- **Minesh Mathew**, Ruben Tito, Dimosthenis Karatzas, R Manmatha and CV Jawahar - *Document visual question answering challenge 2020*, IAPR International Workshop on Document Analysis Systems (DAS) *short paper* 2020
- Ruben Tito*, **Minesh Mathew***, Dimosthenis Karatzas, R Manmatha and CV Jawahar - *ICDAR 2021 competition on Document Visual Question Answering*, International Conference on Document Analysis and Recognition (ICDAR) 2021; [* - equal contribution]
- George Tom, **Minesh Mathew**, Dimosthenis Karatzas and CV Jawahar - *ICDAR 2023 Competition on RoadText Video Text Detection, Tracking and Recognition*, International Conference on Document Analysis and Recognition (ICDAR) 2023

Chapter 2

Background

In this section, we briefly summarize literature from DIA, CV and NLP that are related to the works done as part of the thesis.

2.1 Text recognition in the deep learning era

Text recognition is the problem of recognizing text—extracting machine-readable text—from a word or line image. Traditional text recognition approaches typically involved segmentation of word or line images into patches of symbols, extracting hand-crafted features from the patches and symbol-level classifiers. Feature engineering for such models usually depends on the script in which the text is written and text-modality—printed, handwritten, or scene text. Such modular approaches resulted in the compounding of errors. Additionally, the use of language or text-modality-specific feature engineering prevented the development of robust models for different types of text modalities and languages.

The key characteristics of text recognition methods developed during the deep learning era are i) the use of CNNs for feature learning and ii) the use of RNNs for modelling intra-word or intra-line relationships. The former avoided the need for feature engineering, and the latter enabled us to do unconstrained text recognition—i.e, the set of possible words that need to be recognized are not limited by a fixed vocabulary—without the need for sub-word or sub-line segmentation. Below, we discuss notable developments during this period.

2.1.1 Feature learning using CNNs.

Feature learning using CNNs is part of almost all of the text recognition approaches proposed since the onset of deep learning. Early text recognition works using CNN used it to

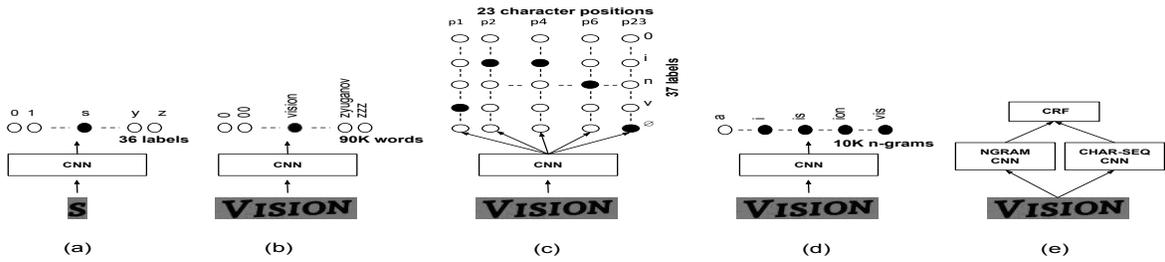


Figure 2.1: **Character and word classification using CNN.** A CNN-based character or symbol classifier, as used in many of the early text recognition works that use a CNN [121], is shown in (a). Others are different models proposed by Jaderberg et al. [124, 125] for whole word recognition. (b) is a multi-class classification model for words. (c) shows a network that has different classification heads corresponding to different character positions within the image. (d) shows a CNN that is used in a multi-label classification setting. It is trained to detect whether the text in the given word image has any of the predefined n-grams. Once the presence of certain n-grams is detected using the model, the word with the nearest—in terms of Euclidean distance—n-gram encoding is picked from a lexicon. In (e), models shown in (c) and (d) are combined using a CRF.

learn features of individual characters or digits. One of the earliest CNN that was proposed in the literature, called LeNet-5 [120], is a network that was made for recognizing digits. By 2010, many works were introduced that use a CNN for the classification of isolated characters or digits [121, 122, 123]. A schematic of CNN-based classifier used for character image classification is shown in Figure 2.1(a).

In contrast to the early works that use a bottom-up approach to word recognition (i.e., predictions are made at the character level, and these predictions are further used to derive a word level transcription), later works follow a top-down approach. In the top-down approach, character patches are not separated out to recognize text in a given word image. Instead, input to the CNN-based classifier is an image of a word. Jaderberg et al. [124] proposed three different encoding schemes to model the output space when the input to the CNN is a word image. In the first encoding scheme, text recognition is modelled as a multi-class classification problem (see Figure 2.1(b)). The classes are 90K words drawn from a fixed vocabulary. In the second approach, the last fully-connected layer of a CNN is connected to 23 separate classification heads, each having 36 output units corresponding to 26 English characters, 10 digits, and a ‘null’ (\emptyset) class. The null class corresponds to the case where no character needs to be emitted. The 23 classification heads independently predict the probability of a label at a particular position in the word. Authors call this approach ‘CHAR CNN’. In the third scheme called, ‘NGRAM CNN’, units in the final classification layer correspond to a fixed set of n-grams.

Activations at the classification layer are thresholded to binary activations indicating the presence or absence of any of the n-grams in the word image. Two types of decoding shall be used to arrive at the target word from the final binary activation vector. In the naive decoding, the word in fixed size lexicon with the nearest n-gram encoding is found. An alternative, that yielded comparable results to the naive decoding, was a linear Support Vector Machine (SVM) that was trained to map from the n-gram encoding to the lexicon. In a later work, authors proposed a new model that combines a CHAR CNN and an NGRAM CNN using a Conditional Random Field (CRF) [125]. The unary potentials of the CRF are provided by the CHAR-CNN, and the NGRAM CNN provides the higher order terms. The character predictions are position-dependent, while the n-gram predictions are not. The objective here is to maximize the CRF score, preserving the consistency of the positional order of the character predictions. Network parameters of both the CHAR CNN and the NGRAM CNN are jointly optimized by minimizing a structured output loss modelled by the CRF.

The above approaches that use a CNN for whole word recognition report results on standard scene text benchmarks without using any dataset-specific lexicon. This marked a shift from most of the earlier approaches that relied on dataset-specific lexicons to make the final prediction. The multi-class classification model that uses a large (90K words) lexicon yields the best results among the 4 different schemes. However, this method cannot recognize a word that is not in the 90K lexicon. The same is the case with the NGRAM CNN since it needs a lexicon (the same 90K lexicon is used) to derive a target word from the n-grams that the network predicted as being present in the given word image. The CHAR CNN and the joint model are capable of unconstrained recognition, since they are not bound by a lexicon. These models can recognize words of length up to a max length of 23 characters.

2.1.2 Text recognition as a structured prediction using RNNs

Structured output problems are a class of problems where the output is not a class label or a number but an object, like a sequence, that possesses an inherent structure [126]. In most cases, both the input and output possess a structure, and these structures are somehow related. Typically, the objective in such problems is to model the structures in both the input and output domain and to find a mapping from the input to the output. For example, many popular problems in NLP, such as Machine Translation (MT), Automatic Speech Recognition (ASR) and Text to Speech (TTS) are structured prediction problems. The input and output in all these problems have a sequential structure.

Recognizing text on a word or a line image is a sequence-to-sequence (seq2seq) structured prediction problem. A sequence of image features from the word or line image needs to be

mapped to a sequence of character labels or Unicode points. The sequential nature of the text recognition motivated researchers to make use of RNNs for modelling the features extracted from the word or line image. RNN-based text recognition models in the literature generally fall into two categories: models that use Connectionist Temporal Classification (CTC) and those that follow an encoder-decoder architecture.

2.1.2.1 CTC-based text recognition

Standard loss functions used with neural networks require a target label for each input instance to compute the loss. Such loss functions are not suitable for seq2seq problems where input and output sequences are not aligned. Graves et al. proposed the CTC algorithm [85] that allows learning a mapping from an input sequence of features to an output sequence of class labels without explicit alignment between the two sequences. Seq2seq models using CTC typically comprise three modules: i) a feature extraction block, ii) a neural network—most often an RNN is used—that outputs a probability distribution over output classes for each input instance—and iii) CTC transcription.

Graves et al. [86] were the first to use CTC for the problem of text recognition. They showed that a bi-directional LSTM (BLSTM) network could be trained end-to-end to map from a sequence of features to a sequence of English characters using CTC. A schematic of this approach is shown in Figure 2.2a. This approach improves word accuracy for offline handwriting recognition on the IAM dataset from 64.5% to 74.1%, over a Hidden Markov Model(HMM)-based model.

The CTC-based transcription approach was then adapted for printed text recognition in multiple languages, such as Indic languages [97, 99, 110, 127] and Fraktur [98]. The fact that the CTC-based approach does not require sub-line segmentation has been hugely beneficial for scripts like Indic scripts and Arabic, where segmenting word or line images to isolated symbols is a challenging task [110]. In [100], authors employed a similar approach for scene text recognition, using gradient-based features extracted from scene text word images in a sliding window manner.

A novel approach combining a CNN for feature learning and an RNN for sequence modelling was introduced in [46]. This hybrid network called CRNN, is end-to-end trainable using CTC. This approach is shown in Figure 2.2b. Since the RNN expects a sequence of fixed-size vectors, the 3D feature maps from the CNN need to be mapped to a sequence of features. In CRNN, this is achieved by adjusting pooling window sizes so that the height of the feature map eventually becomes 1.

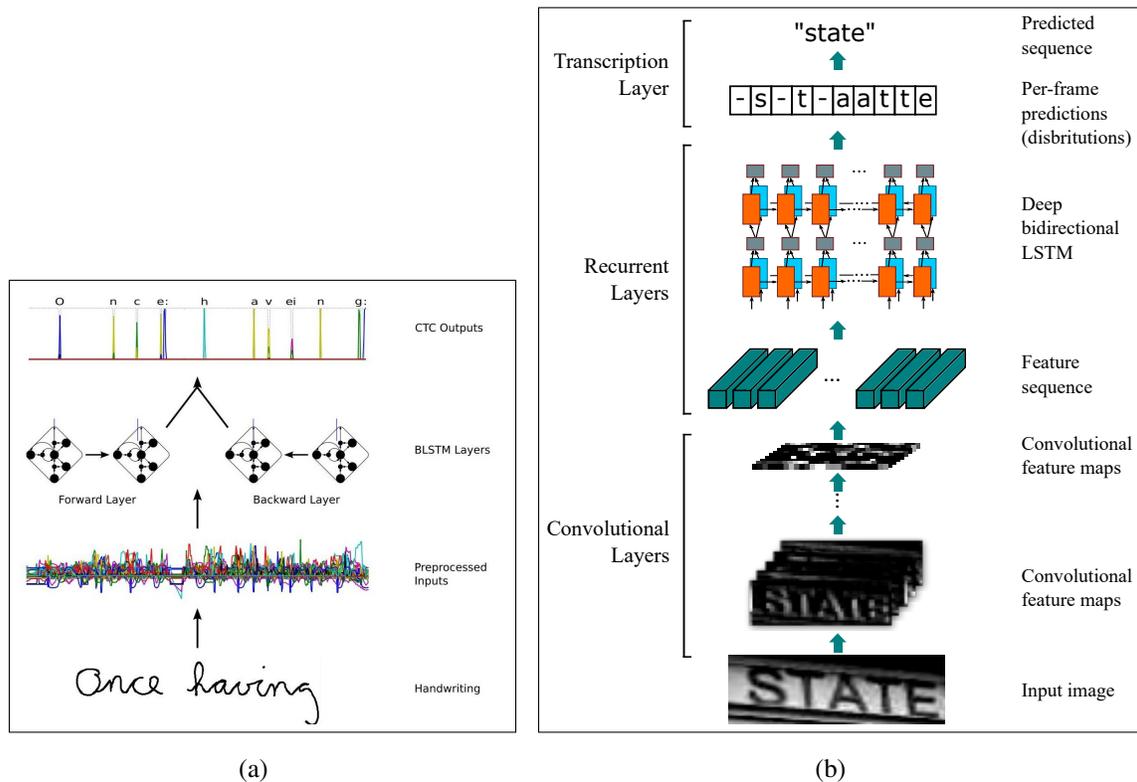


Figure 2.2: **CTC for Image to text OCR transcription.** On the left is A BLSTM + CTC network that Graves et al. [86] proposed for offline handwriting recognition. This approach uses handcrafted features extracted from the handwritten line image in a sliding window manner. Shown on the right is a hybrid CNN-RNN network proposed by Shi et al. [46] which is trained end-to-end using CTC to generate text transcription directly from the given input image. The figure on the left is from [86] and the figure on the right is from [46]

2.1.2.2 Encoder-Decoder Networks for text recognition

Encoder-decoder neural networks became popular with the success of using them for machine translation [87, 88]. In the case of a naive encoder-decoder network that uses an RNN at both the encoder and decoder, the encoder compresses the entire input sequence into a fixed dimensional context vector. When used for machine translation, it compresses the source sentence into a fixed-size vector, irrespective of the length of the sentence or the amount of information present in it. In addition to this, in the case of structured input, the structure is lost once the input is encoded. These two drawbacks of naive encoder-decoder networks are addressed by introducing an attention mechanism between the encoder and decoder. Here a single context vector is replaced by a set of context vectors. Each vector corresponds to a particular spatial, temporal, or spatio-temporal location of the input. During training, the attention model learns

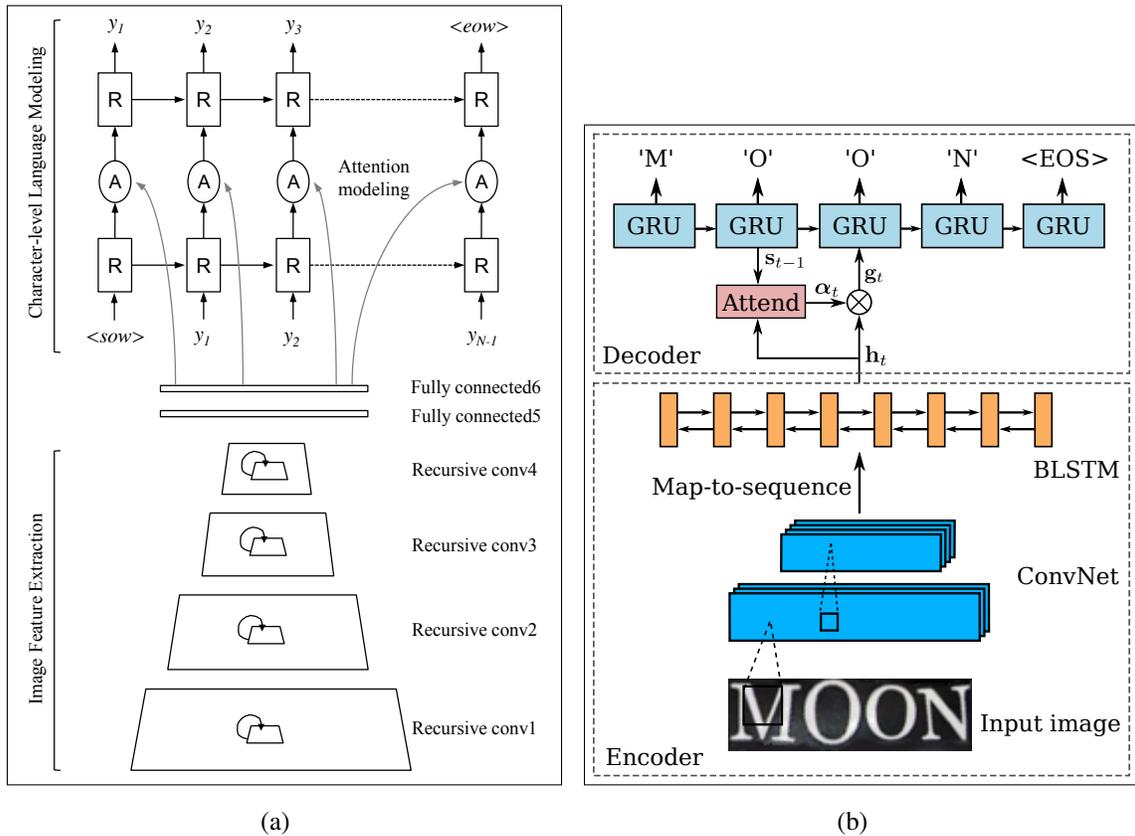


Figure 2.3: **Encoder-decoder style models proposed for unconstrained scene text recognition.** (a) is an image captioning style network named R^2AM proposed in [129] where the encoder is a recursive CNN and the decoder is a bi-directional LSTM . Another model called RARE proposed in [130] is shown in (b). It uses a CNN and bidirectional-LSTM as the encoder and a GRU at the decoder. The figures (a) and (b) are from [129] and [130], respectively.

how relevant each ‘part’ of the input is for each output. This is achieved by learning a set of attention weights [128].

For scene text recognition, Lee et al. [129] proposed an encoder-decoder network that uses a CNN encoder and an RNN decoder. This approach can be interpreted as a conditional RNN language model (RNN-LM). Character level RNN-LM models the conditional probability of the next character, provided the previous characters. Such a model shall also be used to generate text. Lee et al. use an RNN-LM where the text generation is conditioned on image features extracted from a scene text word image. They use a CNN to learn the image features. The authors demonstrate multiple ways in which the image features shall be fed to the RNN. The best-performing model, called R^2AM , feeds the image features through an attention mechanism to the second layer of a two-layer RNN-LM at all time-steps. We show this architecture

in Figure 2.3a. In RARE [130], the encoder-decoder framework for scene text recognition uses a CNN and a bi-directional RNN at the encoder side. The decoder is a single layer RNN-LM that generates a sequence of characters conditioned on the output from the encoder, aided by an attention mechanism. A schematic of this model is shown figure Figure 2.3b.

2.1.2.3 Multidimensional RNN

Graves et al. [131] proposed the multi-dimensional RNN (MDRNN) [131] that can model dependencies in more than one spatio-temporal dimension. Although MDRNN has been proven successful for text recognition tasks [132], a later study by Puigcerver argued that computationally expensive MDRNNs are not essential for obtaining similar performances [133]. He showed that the features extracted by MDRNN layers are visually similar to those extracted by a CNN. He further argues that two-dimensional long-term dependencies that are modelled by MDRNN layers may not be essential for the problem of text recognition.

2.2 Question Answering for electronic text

Our works that deal with question answering on document images are related to tasks that deal with question answering for unstructured electronic text. In this section, we discuss major tasks, datasets, and methods used for QA for electronic text. Depending on whether the questions are asked given a specific context of one or a few passages or not, there are two types of QA tasks: Machine Reading Comprehension and Open-domain QA.

2.2.1 Machine Reading Comprehension (MRC)

In MRC, the task is to answer a natural language question given a passage as the context. Although early works in MRC go back to the late 1970s, MRC tasks gained momentum in the latter half of the last decade. The success of deep learning for a wide variety of NLP tasks and the introduction of multiple large-scale datasets such as SQuAD 1.1 [91] and MS-MARCO [134] gave a fillip to research in MRC. Liu et al. [51], who review the early developments in the MRC space, observe that there is a 10-fold increase in the number of publications related to MRC in top conferences and journals during the 2015-18 period. MRC in SQuAD is formulated as an extractive task—i.e., the answer is extracted as a span from the given passage. In Figure 2.4, we show example QA pairs in SQuAD 1.1 dataset. On the other hand, in datasets

like MS-MARCO [134], MRC is an abstractive task where the answer needs to be generated, not extracted.

Transformer [135]-based architectures like BERT [113] and XLNet [136] that capitalize on large-scale pretraining report performance surpassing human performance on extractive MRC benchmarks like SQuAD 1.1 [91]. Below, we discuss the BERT architecture and how it is pretrained and finetuned for MRC.

2.2.1.1 BERT for QA/MRC

Bidirectional Encoder Representations from Transformer (BERT) is a method of pretraining language representations from unlabelled text using transformers. These pretrained models can then be used for downstream tasks by using an output head suitable for the task at hand. In

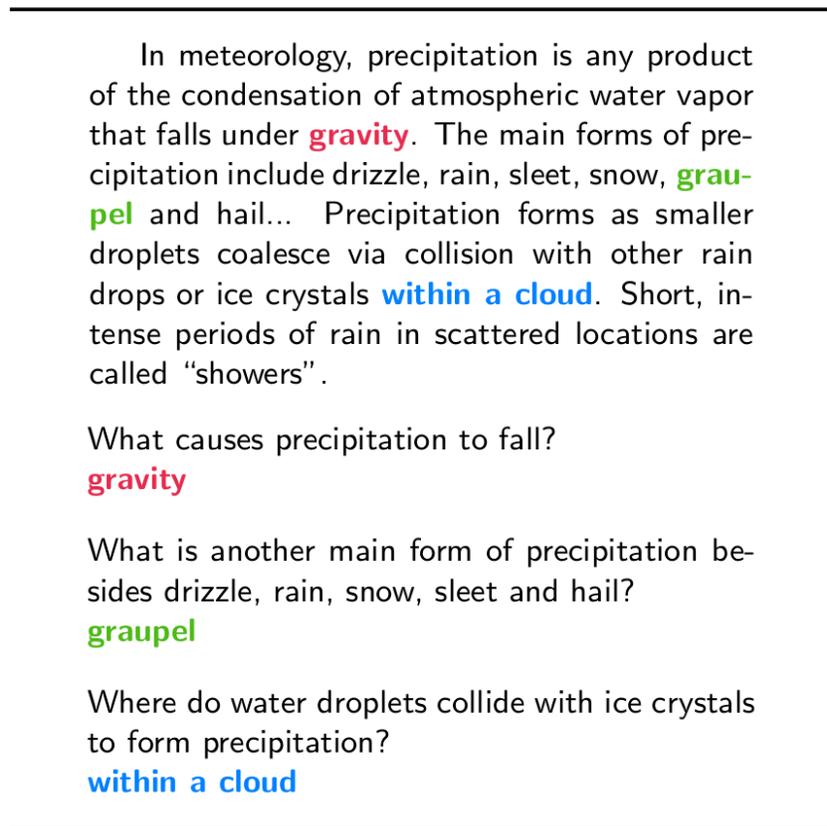


Figure 2.4: **Extractive MRC task in SQuAD1.1.** An example passage in the SQuAD 1.1 dataset and QA pairs annotated on the passage. The task is an extractive QA task since the answer to any question in the dataset is a span—a sequence of continuous text tokens—from the respective passage. The figure is from [91]

the case of extractive QA, this is an output layer that predicts the start and end indices of the answer.

BERT’s model architecture is essentially a bidirectional transformer encoder. In the original work, they work with two configurations of the architecture: BERT_{BASE} and BERT_{LARGE}. BERT_{BASE} has 12 transformer blocks in the encoder and the hidden size is 768. The total number of parameters in it is 110M. BERT_{LARGE} has 24 Transformer blocks and a hidden size of 1024. This variant has 340M parameters.

Input to BERT is a sequence of text tokens. A tokenizer based on the word piece model [137] is used. Every input sequence starts with a special [CLS] token. In the case of tasks where there is more than one type of sequence, a special [SEP] token is used to separate the two sequences. For example, in the case of MRC, the sequence of tokens from the question is separated from the sequence of tokens from the passage using the [SEP] token. The two sequences are also differentiated by adding a learnt segment embedding. For example, in the case of MRC, all the question tokens will have a segment embedding different from the segment embedding for tokens in the passage. The information about the position of a token in the sequence is fed to the model using a learnt position embedding. Thus, each token in the input sequence is the sum of the token embedding, segment embedding, and position embedding.

BERT has two training stages. In the pretraining stage, the model is trained on large amounts of language data that do not have any extra annotations. Two pretraining tasks are used in BERT. The first pretraining task is called ‘Masked LM’. In this task, some of the tokens in the input sequence are masked and the model is made to predict the masked tokens. In the second pretraining task, input has a pair of sequences, each separated by a [SEP] token. The task is to predict if the second sequence is a continuation of the first in the pretraining corpus. For

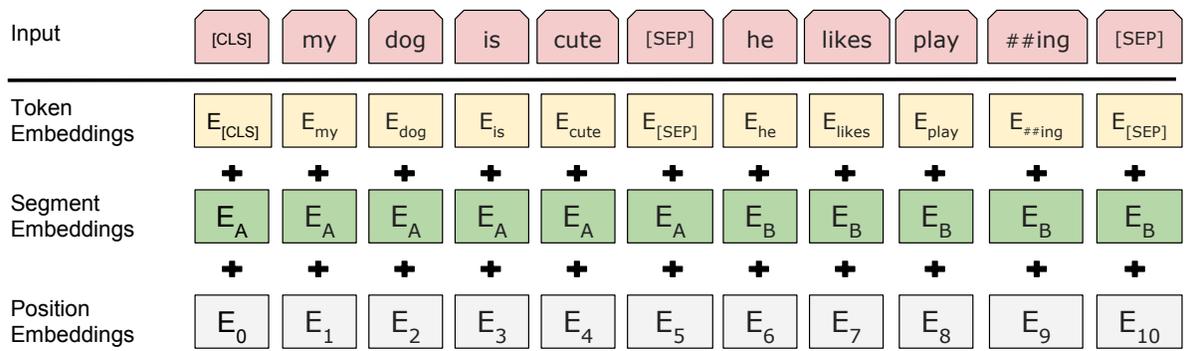


Figure 2.5: **Embeddings used in BERT.** The input sequence is tokenized using Word piece [137] tokenizer. Each input token is represented by the sum of its token embedding, segment embedding, and position embedding. The figure is from [113]

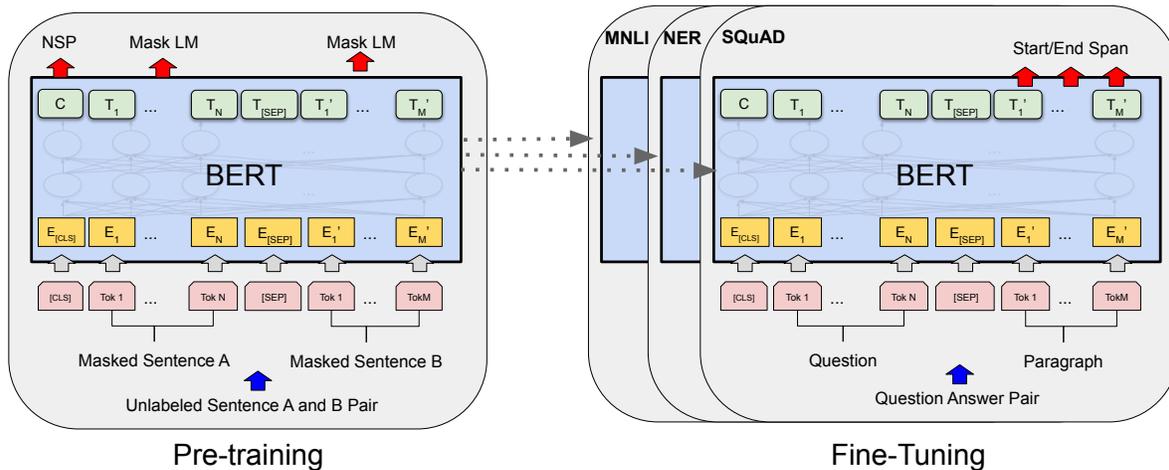


Figure 2.6: **Pretraining and finetuning stages of BERT.** During pretraining, the Transformer encoder block is trained on tasks such as Masked Language Modelling (MLM) and Next Sentence Prediction (NSP). The pretrained model is then finetuned separately for each of the downstream tasks. Shown here is the finetuning for the MRC on SQuAD1.1, where the output layer is modified to predict the start and end tokens of the answer span. The figure is from [113]

example, if you think in terms of linguistic sentences, this task makes the model predict if the second sentence follows the first in the training corpus or not. Therefore this task is called ‘Next Sentence Prediction’ (NSP) in the original work.

The second stage of the BERT training is the finetuning stage. In this stage, the pretrained BERT is trained in a supervised manner for the downstream task at hand. For MRC, input comprises two sequences—question tokens and tokens from the passage—separated by a [SEP] token. The output head for MRC acts as a span classifier. This head has two learnable weight vectors, S and E , for start and end tokens of the span, respectively. Both vectors have the same size as the hidden size of the model.

At the time of finetuning for the MRC, dot product $T_i \cdot S$ is computed for each of the tokens, where T_i is the final hidden vector for the token at position i . Softmax normalization is applied to the dot product values for all the T_i s yielding a probability vector where the value at position i in the vector is the probability of token at position i being the start token of the answer span. Using the ground truth information, cross entropy between the probability vector and a one-hot vector—indicating which token is the actual start token—is calculated. Cross-entropy for the end token is also computed similarly. The training objective is to minimize the sum of the two cross-entropies.

At the time of inference, we pick tokens at positions i and j that maximize $T_i \cdot S + T_j \cdot E$ for $j > i$, as the start and end tokens of the answer span, respectively.

Open-domain QA

SQuAD, TREC, WebQuestions, WikiMovies

Q: How many of Warsaw's inhabitants spoke Polish in 1933?

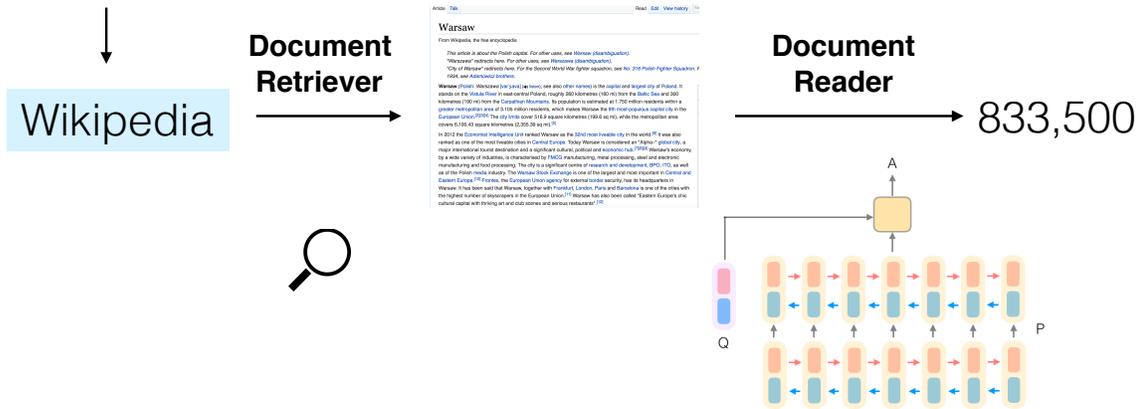


Figure 2.7: **Two-stage pipeline for open domain QA** Chen et al. [138] proposed a two-stage pipeline for open-domain QA. In the first stage, a TF-IDF-based ranked retrieval is used to retrieve a few documents (usually 5 or 10) that are most relevant to the question being asked. In the second stage, authors employ an extractive MRC model that uses an LSTM-based network as proposed in [139]. The figure is from [138]

2.2.2 Open domain QA

We have seen that MRC is a QA task where the questions are based on a specific context, usually one or a few passages. On the other hand, open domain QA, as the name suggests, involves QA where no specific context is given. In most cases, open domain QA is modelled as a task wherein the question’s answer needs to be found from a large electronic document collection, say the entire Wikipedia.

Existing solutions to open domain QA typically follows a two-stage solution. In the first stage, documents that are relevant to the question are retrieved from the document collection. This stage is usually called ‘Document Retriever’. In the second stage, a ‘Document Reader’ returns a textual answer using the retrieved documents as context. IR-based document retrieval techniques that do not involve any machine learning are typically used in the retriever stage. For QA tasks where answers can be extracted, deep learning based MRC models are employed for the second stage. In Figure 2.7, we show an example of the two-stage solution named DrQA, proposed by Chen et al. [138].



Q: What is the mustache made of?

A: [bananas](#)



Q: Where will the driver go if turning right? a) Onto 24 3/4 Rd. b) Onto 25 3/4 Rd. c) Onto 23 3/4 Rd. d) Onto Main Street.

A: [a\) Onto 24 3/4 Rd.](#)

Figure 2.8: **VQA on natural images.** On the left is an example from the VQA v1 dataset [55]. Answers in the dataset are free-form natural answers. More than 97% of the answers in the dataset have less than three words in it, and the share of unique answers is quite less. For this reason, VQA on the VQA v1 dataset is often modelled as a classification problem. On the right is an example from the Visual7W dataset [140] that features multiple-choice questions. In addition to the questions where answers are textual, the dataset also includes what the authors call ‘pointing questions’, where answers are a cropped region from the given image.

2.3 Layout-aware BERT-like models for DIA tasks

The success of transformer-based models for text understanding inspired the use of similar models for DIA. LayoutLM [115] and LAMBERT [141] incorporate layout information into the BERT architecture by using embeddings of the 2D positions of the text tokens in the image.

Additionally, there have been many newer or concurrent works that deal with a joint understanding of the text, image, and layout in document images [142, 143, 144, 145]. These models build on Transformer encoder layers and leverage large-scale pretraining on unlabelled data.

The LayoutLM serves as the basis for the VLL-BERT we propose for extractive VQA on our InfographicVQA dataset presented in chapter 5. Below, we discuss the details of the LayoutLM architecture and how it is trained.

2.3.1 LayoutLM

A schematic of the LayoutLM is shown in Figure 2.9. LayoutLM is a pretrained transformer-based model for document image understanding. It is based on the BERT architecture and is designed to leverage both textual and layout information for various DIA tasks.

As discussed in subsection 2.2.1.1, input to BERT is the sum of text, position, and segment embeddings for each token. Since the LayoutLM is trained on text tokens spotted on document images, in addition to the aforementioned embeddings, a 2D position embedding is also used. Similar to how position embedding captures the order of a token in a sequence, 2D position embedding captures the spatial location of the token on the document image plane. Corresponding to the 4 coordinates— x_1, y_1, x_2, y_2 —of the bounding box of a token, 4 different embeddings are learnt but only two embedding tables are used one each for x and y dimensions.

In addition to the 2D position embeddings, LayoutLM incorporates visual information at the time of finetuning the model for various tasks. LayoutLM uses a late fusion of visual features with other embeddings. That is, visual features (aka image embeddings) are not a part of the input representation, but are added to the attention-enriched representations outputted by the transformer stack. Image embeddings are not part of the self attention. Image embeddings used by LayoutLM are region-based features (subsection 2.4.1) corresponding to each OCR token. For the [CLS] token, region-based feature corresponding to the entire document image is used. A Faster-RCNN [96] trained on Visual Genome dataset [146] for object detection is used for extracting the region-based features.

The model is pretrained on IIT-CDIP Test Collection [147]. While pretraining, similar to BERT, the model is trained in a self-supervised manner for Masked Visual Language Modelling (MVLM). MVLM is similar to MLM in BERT. The only difference is that 2D position embeddings are also a part of the input representation. The model is trained to predict the masked text token, given its 2D and 1D position and the surrounding context. In addition to the self-supervised MVLM task, LayoutLM is pretrained to perform Multi-label Document Classification (MDC) task. Document tags in IIT-CDIP Test Collection are used for training the model to perform MDC. The pretrained model is finetuned to perform multiple DIA tasks such as document image classification, form understanding, and receipt understanding.

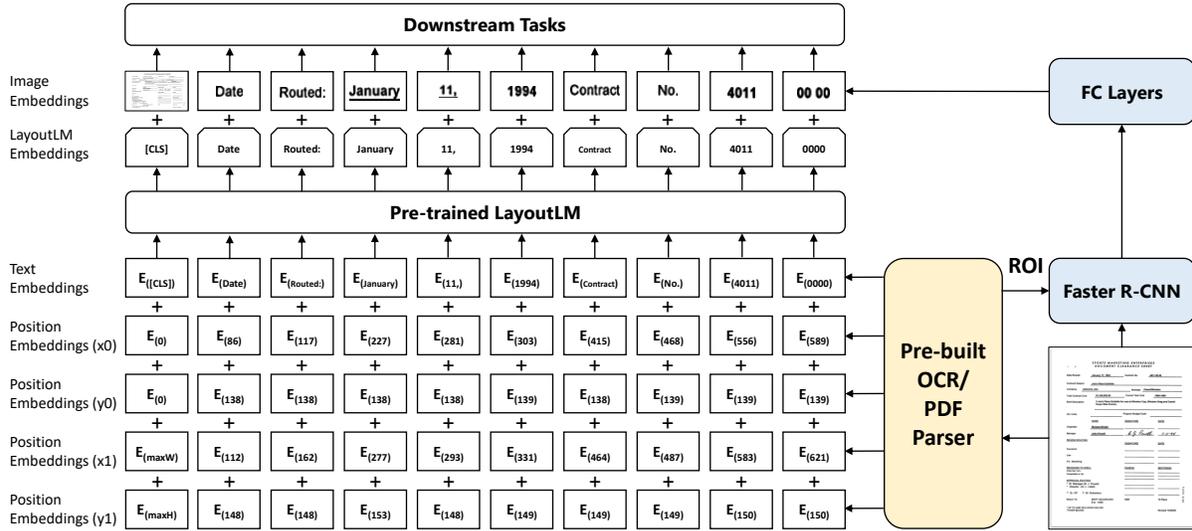


Figure 2.9: **LayoutLM architecture.** LayoutLM is essentially a BERT extended to incorporate layout information. A 2D position embedding that represents the location of a text token in the image plane is added in addition to the embeddings that are already included in the input representation of BERT. While finetuning, image embeddings—region-based features for the corresponding text token—are added to the embeddings obtained from the LayoutLM model. The figure is from [115]

2.4 Visual Question Answering (VQA)

Visual Question Answering (VQA) aims to answer a question asked on a given image. Similar to QA in electronic text, research in VQA gained momentum in the second half of the last decade owing to i) the success of deep learning techniques for various tasks in CV and NLP and ii) the introduction of large-scale datasets for the problem such as VQA V1 [55], Visual7W [140] and Visual Genome VQA [146]. Most of these early datasets contain natural images. In Figure 2.8, we show an example each from VQA V1 and Visual7W datasets.

2.4.1 Grid-based and region-based features for VQA

The early methods for VQA typically employ a CNN to encode the image, an RNN to encode the question, and eventually learn a joint embedding of the two. Many of these models employ question-guided attention over the visual features from an image encoder [56]. The image encoder used is a CNN that is pretrained to perform image classification. For the given image, a feature map—usually the last feature map—from the CNN is used for extracting visual features. If the feature map is of size $w \times h \times c$, each location in the spatial extent of the feature corresponds to a patch of the input image. The patch that corresponds to a particular

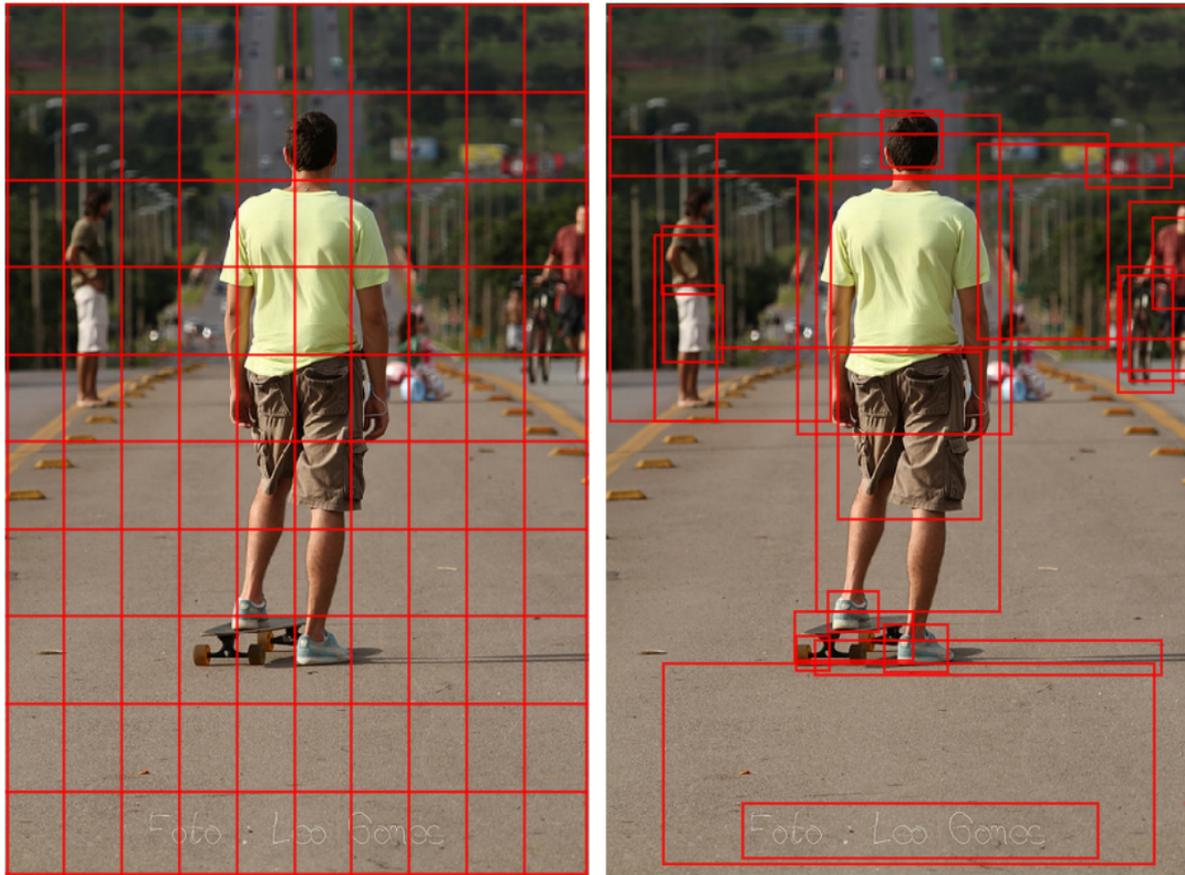


Figure 2.10: **Grid-based features and region-based features for VQA** In the case of grid-based features (left), the VQA models make use of features from a CNN feature map. Locations on the spatial extent of this feature map correspond to a uniform grid on the input image. In the case of region-based features (right) introduced by Anderson et al. [93], each visual feature extracted from the image corresponds to the bounding box of an object detected on the image. Figure is from [93]

location in the feature map is the receptive field of the convolutional filters that computed that features at that location. Therefore the feature map corresponds to a uniform grid of image regions that are of same size. Consequently, these features are referred as grid-based features in the literature. In other words, a feature map of size $w \times h \times c$ can be thought of as $w \cdot h$ feature vectors each of size c and each representing a patch of the input image.

Anderson et al. [93] observed the VQA models that attend over the grid-based features ignore the semantics of the image. They proposed to use an object detection model to detect objects in the given image and extract features corresponding to the detected objects. Typically, ROI pooled features corresponding to the detected object instances are extracted from the last or penultimate fully connected layer of the box head. These features are called region-based

features. Different from the top-down attention over the grid-based features, attention over the region-based features is a bottom-up attention since the locations to attend to are semantically meaningful regions within the image. The bottom-up attention has then been used for almost all of the VQA and image captioning works that followed, such as Pythia [148], LoRRA [149], M4C [94] and VisualBERT [150].

2.4.2 VQA on natural images using VisualBERT

Following the success of BERT [113]-like models for NLP tasks, multiple works were proposed extending BERT to the Vision + Language space. Models like VL-BERT [151], VisualBERT [150], and UNITER [152] show that pretraining of a BERT-like architecture jointly on vision and language inputs achieve SoTA performances on various downstream tasks, including VQA on natural images. In the following, we briefly discuss one of such architectures called Visual BERT. We pick this model for discussion since the model design follows naturally from the BERT.

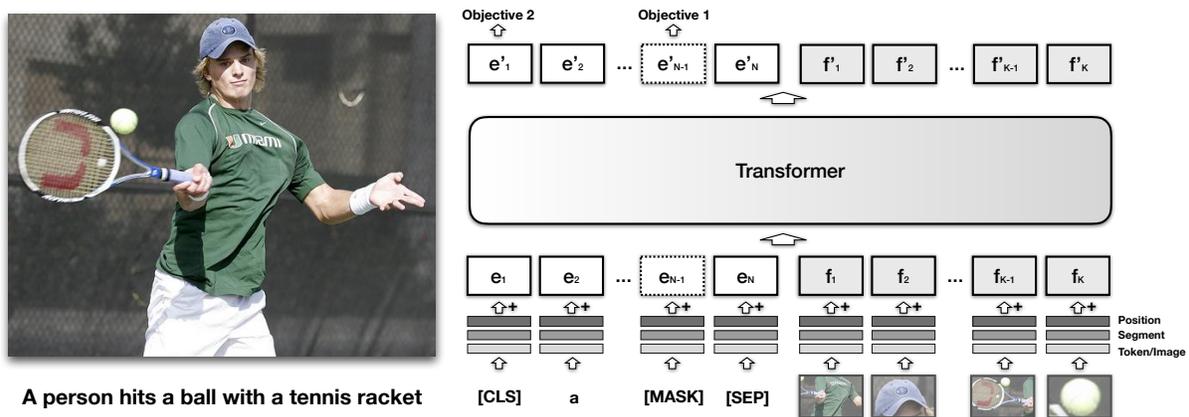


Figure 2.11: **VisualBERT for Visio-linguistic pretraining.** Similar to BERT, VisualBERT [150] builds an encoder architecture that jointly models vision and language using self attention mechanism in Transformers. As shown in the figure, at the time of pretraining, the input to the model comprises embeddings of the text tokens in the caption and visual embeddings (region-based features) of the visual tokens—objects detected in the image using an object detection model like Faster-RCNN. During finetuning for VQA, similar to BERT QA model, the tokens from the question are used in the input along with the visual tokens. The output head used during finetuning for VQA is a classification layer that picks an answer from a fixed vocabulary of answers. The figure is from [150]

VisualBERT, as the name suggests, is directly inspired by BERT. It is essentially a BERT that incorporates changes to accommodate visual inputs. Similar to BERT, input to VisualBERT is a sequence of tokens. But the tokens include both textual and visual tokens. The visual tokens are the visual objects detected on an image. Similar to bottom-up attention we discuss in subsection 2.4.1, region-based features corresponding to objects detected on the image are used as embeddings for the visual tokens. Each token is represented by the sum of three embeddings: i) text token embedding for text tokens and visual feature embedding for the visual tokens, ii) segment embedding and iii) position embedding as used in BERT.

The Visual BERT is pretrained on image + caption pairs from MS-COCO [89] using two pretraining tasks: MLM and sentence-image prediction. Similar to BERT, for MLM, some text tokens are masked, but no visual tokens are masked. For the sentence-image prediction task, an image is coupled with two captions. Half of the time, both captions are for the same image, and for the rest of the cases, one of the captions is not related to the image. The task is to distinguish the two cases. A schematic of the model, showing the pretraining setup, is shown in Figure 2.11

While finetuning for VQA on natural images, authors train the model on the VQA V2 dataset [153]. An answer classification head is used to predict the correct answer from 3000+ answers. Following the common practice, the most frequent answers in the train split are used as the answer vocabulary.

Grid-based features from a Faster-RCNN trained on Visual Genome for object detection are used for visual features of the visual tokens.

2.5 VQA tasks that require reading text on images

In the previous section, we talked about VQA in general and the early datasets that were introduced for the problem. In this section, we discuss VQA tasks and datasets where answering the questions requires reading and comprehending textual content in the images. Gurari et al. [154] show that in a goal-oriented VQA setting where visually impaired individuals ask questions on images they take, answering a good number of questions require the ability to read and interpret text on the images. In fact, the early VQA datasets such as, VQA v1 and Visual7W, included questions that require recognizing and reading text on the images. For example, in the case of the example from Visual7W shown in Figure 2.8, the answer is the name of a road displayed on a traffic information board. However, a number of text-based questions in these datasets is not many, and the models developed for VQA on these datasets often ignored the textual cues from the image.

Realising the importance of scene text in understanding natural images, the computer vision community introduced two new datasets—ST-VQA [155] and TextVQA [149]—that introduce a new VQA task that requires reading text in natural images to answer questions. We shall refer to this type of VQA as ‘Scene text VQA’ in the rest of the discussion. Around the same time, there have been works on VQA on special categories of images such as charts and book covers where reading the text on the images is necessary to answer the questions. A brief summary of these datasets and some of the popular methods for such VQA tasks are discussed below.

2.5.1 Scene text VQA

VQA datasets like VQA V1 or Visual7W focused on asking questions concerning visual aspects such as people, objects, and relationships between various semantic entities seen in the image. Some of the questions in these datasets required reading text on the images to answer the questions, but there were not many. Around 2019, ST-VQA and Text VQA datasets were introduced, where reading text on the images is critical to answering almost all questions in the datasets.

While TextVQA sourced images from OpenImages [156], images in Scene text VQA are from various scene text detection and recognition datasets. In ST-VQA, all questions can be answered purely using text present in the images. On the other hand, in TextVQA, while almost every question requires reading text, there is no guarantee that answers are composed purely of text tokens present in the image.

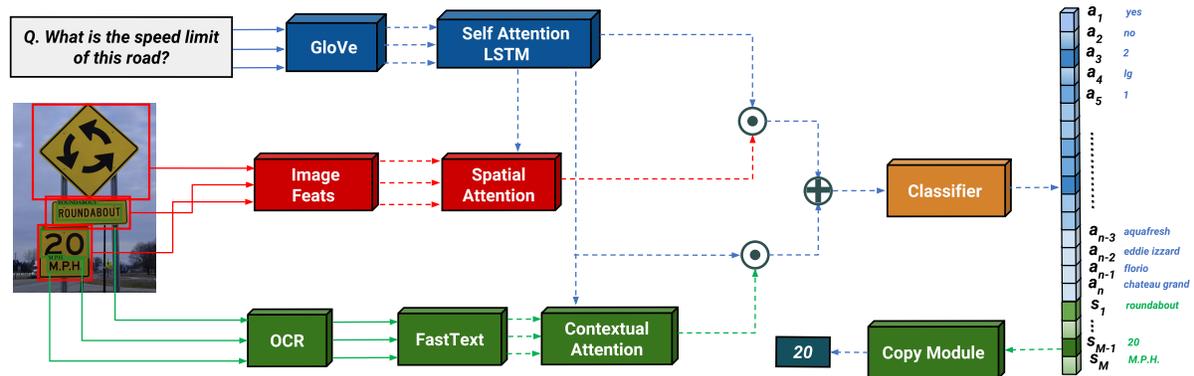


Figure 2.12: **LoRRA for scene text VQA.** In LoRRA [149] a question guided attention is used to compute attention weights for visual features and embeddings of OCR tokens. The end-to-end trainable network uses a classification layer that picks an answer either from a fixed vocabulary or from the OCR tokens spotted on the image. Figure is from [149]

Below, we discuss two important multimodal architectures that were proposed for TextVQA and Scene text VQA tasks.

2.5.1.1 LoRRA: Look, Read, Reason and Answer

LoRRA extends the original bottom-up and top-down attention [93] scheme with additional bottom-up attention over OCR tokens from the images. In LoRRA, tokens in a question are first embedded using a pretrained word embedding model (GloVe [157]), and then these tokens are encoded using an LSTM [158] encoder. The model uses both grid-based features and region-based features (see subsection 2.4.1). For grid features, a ResNet-152 [95], which is pretrained on ImageNet [38] for image classification is used. For region-based features, a Faster R-CNN trained on Visual Genome dataset [146] for object detection is used. OCR tokens from the image are embedded using a pretrained word embedding (fastText [159]). An attention mechanism is used to compute the attention weighed average of the image features as well the OCR tokens' embeddings. These averaged features are combined and fed into an output module. The final classification layer of the model predicts an answer either from a fixed vocabulary (made from answers in the train set) or copy an answer from a dynamic vocabulary which essentially is the list of OCR tokens in an image. Here the copy mechanism can copy only one of the OCR tokens from the image. Consequently, it cannot output an answer that comprises two or more OCR tokens.

2.5.1.2 Multimodal Multi-Copy Mesh (M4C)

M4C uses a multimodal transformer and an iterative answer generation decoder to yield state-of-the-art results on Text VQA, ST-VQA, and OCR-VQA [60] datasets. OCR-VQA is a specialized VQA dataset for images of book covers. We discuss this dataset in subsection 2.5.2.1. In M4C, tokens in the question are embedded using a pretrained BERT. Images are represented using (i) region-based features of the objects detected using an object detection model—Faster-RCNN pretrained on Visual Genome and (ii) location information - bounding box coordinates of the detected objects. Each OCR token recognized from the image is represented using (i) a pretrained word embedding (fastText [159]), (ii) visual feature for the OCR token, (iii) PHOC [160] representation of the token and (iv) bounding box coordinates of the OCR token. Visual features for OCR tokens are extracted using the same faster-RCNN that is used to detect objects on the images. For the bounding box corresponding to each OCR token, a visual feature is extracted using ROI pooling. Then these feature representations of the three entities (question tokens, objects, and OCR tokens) are projected to a common, learned embedding space. Then a stack of Transformer [135] layers is applied over these features in the

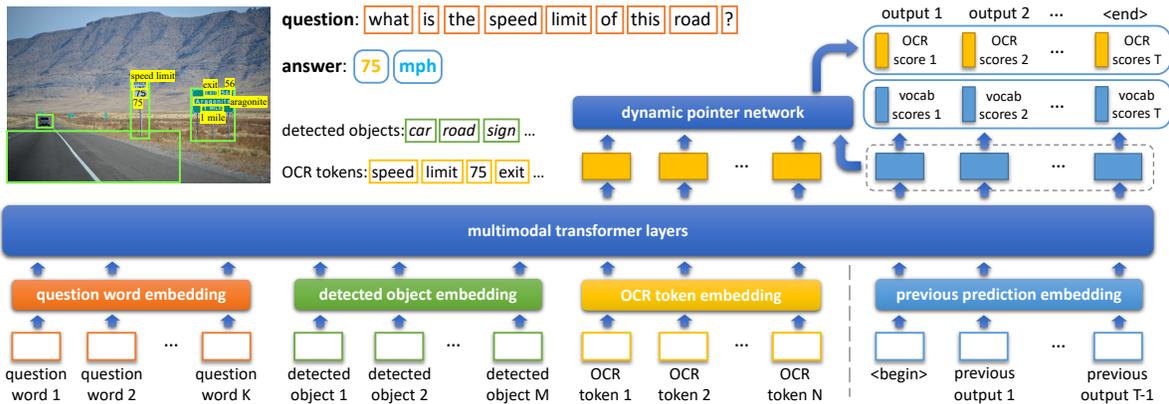


Figure 2.13: **M4C for scene text VQA.** M4C [94] uses a transformer block that jointly models relationships between the question tokens, detected objects, and OCR tokens. The model outputs answers in an auto-regressive fashion using both a fixed vocabulary of the most common answers and a dynamic vocabulary comprising OCR tokens spotted on the image. Figure is from [94]

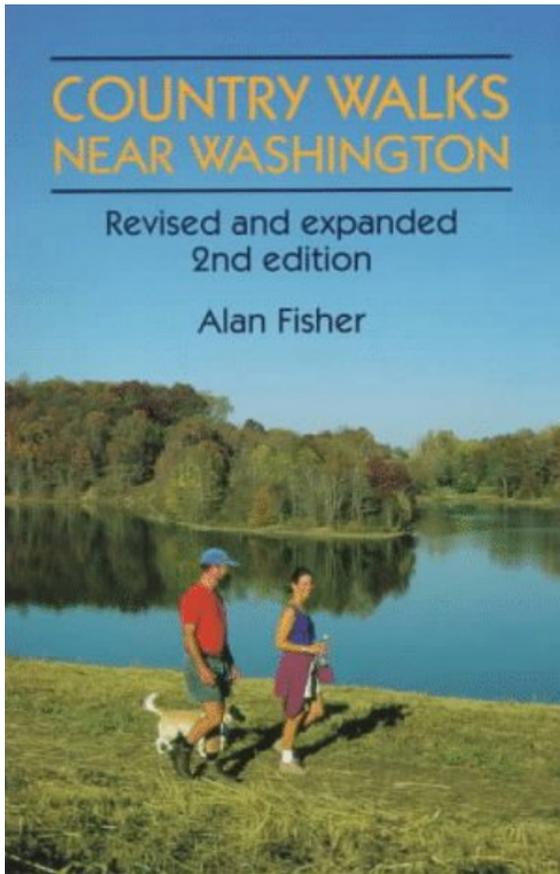
common embedding space. The multi-head self attention in transformer layers enables both inter-entity and intra-entity attention. Finally, answers are predicted through iterative decoding in an auto-regressive manner. The fixed vocabulary used here is made up of the most common answer words in the train split. Note that in this case, the fixed vocabulary comprises answer words, not answers, as in the case of LoRRA. At each step in the decoding, the decoded word is either an OCR token from the image or a word from the fixed vocabulary of common answer words.

2.5.2 Specialized VQA datasets where reading text on images is necessary.

In addition to VQA on natural images, there are a few specialized VQA tasks that require reading text on the images in order to answer the questions. We briefly discuss these datasets below.

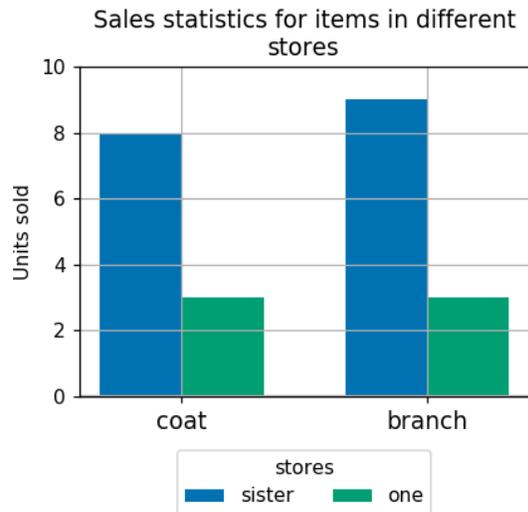
2.5.2.1 VQA on images of book covers

OCR-VQA [60] comprises more than 1 million question-answer pairs over 207K+ images of book covers. The questions in this dataset are domain-specific and generated based on template questions and answers extracted from available metadata. More than 50% of the questions have



Q: Who wrote this book?

A: Alan Fisher



Q: Which item sold the most number of units summed across all the stores?

A: a) branch

Figure 2.14: **VQA on book covers and charts.** Shown on the left is an example from the OCR-VQA dataset [60] that comprises questions and answers asked on images of book covers. The example on the right is from a VQA dataset for bar charts called DVQA [58]. The book covers in OCR-VQA are sourced from the internet and the chart images in DVQA are synthetically generated using a chart plotting library. Questions in both datasets are synthetically generated using a few question templates. Owing to the nature of the images in both datasets, reading text is critical to answering most of the questions.

answers that are not scene text instances, including 40% binary (yes/no) questions and 10% questions about book genres, for example.

2.5.2.2 VQA on charts

The DVQA dataset presented by Kafle et al. [58] comprises synthetically generated images of bar charts and template questions defined automatically based on the bar chart metadata. The dataset contains more than three million question-answer pairs over 300,000 images. They proposed a method called SANDY for DVQA. It is a modified version of the stacked attention network (SAN) [161] that uses off-the-shelf OCR and dynamic encoding to represent out-of-vocabulary words. They observe that contrary to VQA on natural images, small, localised, changes in the chart image might have a big influence on the VQA performance. This stands true with document images in general, as changing a single element (word, check box, etc.) in a document image might completely alter the meaning of the data presented in the document.

Another dataset for VQA on charts, FigureQA [59], comprises over one million yes/no questions, defined on over 100,000 images. The dataset comprises bar charts, line plots and pie charts. Similar to DVQA, images are generated using a chart plotting library, and question-answer pairs are created using a finite number of question templates. The authors propose a relation network [162]-based model as a strong baseline for the task.

2.6 Multimodal QA tasks

Textbook Question Answering (TQA) [163] and RecipeQA [164] deal with QA in a multimodal context. For TQA, contexts are textbook lessons, and for RecipeQA, contexts are recipes containing text and images. TQA is probably the first work that deals with QA in a multimodal context. Multiple choice questions are asked on lessons from science textbooks, and each question is provided with passages from the lesson and associated diagrams. Similarly, in RecipeQA, multiple-choice questions are asked about recipes. Each recipe includes instructions (machine readable text) and images depicting these instructions. Contrary to VQA datasets that involve reading text, text in multimodal QA datasets is not embedded in the images but provided in machine readable form as a separate input.

Chapter 3

CTC-based Seq2seq OCR for low resource languages

In section 1.5 we note that research and development of DIA for Indic languages has been lagging behind Latin-based languages. Deep learning based techniques that avoid the need for script-specific handcrafted features and segmentation-free transcription have become a boon for Indic languages OCR. Approaching text recognition as a structured prediction using CTC is the most commonly used approach for segmentation-free OCR.

This chapter presents a comprehensive empirical study of various neural network models that use CTC for transcribing step-wise predictions in the neural network output to a Unicode sequence. The study is conducted for 13 Indic languages, using an internal dataset that has around 1000 pages per language. We study the choice of line vs word as the recognition unit, and use of synthetic data for training. We compare our models with publicly available OCR tools for end-to-end document image recognition. By tying our recognition model with an existing text detection model, we evaluate end-to-end text recognition performance. We demonstrate that our end-to-end pipeline outperforms the public OCR engines for 8 out of the 13 languages. We also introduce a new public dataset called *Mozhi* for word and line recognition in Indic languages. The dataset contains more than 1.2 million annotated word images (around 120 thousand text lines) across the 13 languages. This work has been previously presented as part of a technical report [112].

3.1 Introduction

Optical Character Recognition (OCR) is generally used as an umbrella term for the process and technology involved in converting text present in a document image to machine readable text. End-to-end OCR generally involves two steps: i) text detection: locating the regions where text tokens are present in an image, and ii) text recognition: transcribing text present in

a line or word region identified in the detection step to machine readable sequence of characters or Unicode points.

Commercial OCR engines for Latin-based languages began to appear during the mid-1960s. These early OCR engines predominantly used template-matching techniques for recognizing isolated characters [165]. This was followed by machines that could recognize machine-printed and handwritten numerals. Simultaneously, Toshiba launched the first automatic letter-sorting machine for postal code numbers. During the 1970-80 period, several OCR engines were developed that could recognize printed and handwritten English characters [165, 166]. Presently there are several publicly available OCR systems for Latin scripts such as Tesseract¹, ocropus², ABYY FineReader³ and Google cloud vision OCR⁴ that can perform OCR even on documents with complex layouts with sub 1% character error rates.

The challenges in text recognition vary based on the language/script the text is written in, how the text is rendered—handwritten, printed, or typewritten—and how the document is imaged— scanned, captured using a mobile device, or born-digital. This work deals with the recognition of text printed in Indic languages. Our focus is on text recognition alone. That is, we assume that cropped word or line images are provided.

This chapter presents our efforts to build OCR for multiple Indic languages. By Indic languages, we refer to languages that are spoken in the Indian subcontinent. Our work deals with text recognition in 13 Indic languages. All 13 languages are granted official language status in India and have more than a million speakers in India [167]. The languages are Assamese, Bangla, Gujarati, Hindi, Kannada, Malayalam, Manipuri, Marathi, Odia, Punjabi, Tamil, Telugu and Urdu. Some of these languages are used widely in other countries in the Indian subcontinent. 99% of the people (around 164 million people) in Bangladesh speak Bangla⁵. According to the 2017 census of Pakistan⁶, more than 38% of the population (around 80 million) speak Punjabi and nearly 7% speak Urdu (around 15 million). Tamil, a Dravidian language that is the primary language in the Indian state of Tamil Nadu, is spoken by more than 15% of the population in Sri Lanka⁷. Many of these languages share common linguistic and grammatical structures. But the script remains very different except for a few languages. Among the 13 languages, Hindi and Marathi use Devanagari script, and Bangla, Assamese and Manipuri use Bangla script. Others have their own unique scripts. Thus our study deals with printed text

¹<https://opensource.google.com/projects/tesseract>

²<https://github.com/ocropus/ocropus>

³<https://pdf.abbyy.com>

⁴<https://cloud.google.com/vision/docs/ocr>

⁵<http://www.bbs.gov.bd/site/page/47856ad0-7e1c-4aab-bd78-892733bc06eb/Population-and-Housing-Census>

⁶<https://www.pbs.gov.pk/content/final-results-census-2017>

⁷<http://www.statistics.gov.lk/Population/StatisticalInformation/CPH2011>

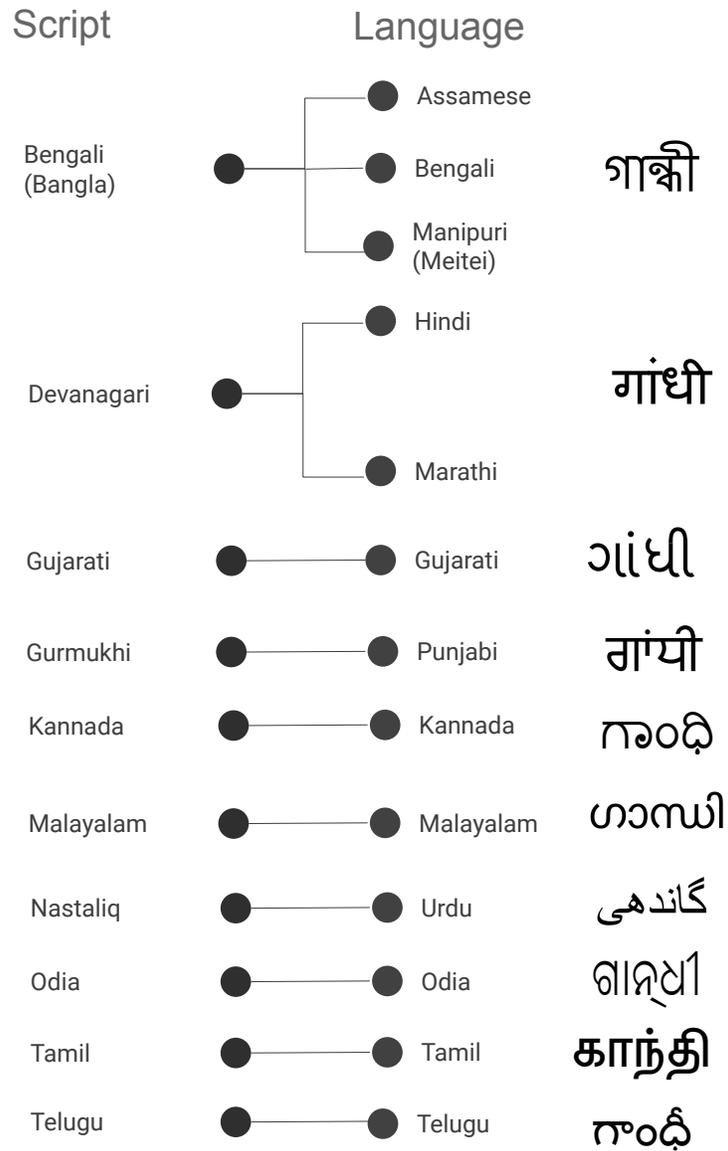


Figure 3.1: **Indic languages diversity.** We study printed text recognition of 13 languages (10 different scripts) that are used in the Indian subcontinent. Although these languages essentially use a common alphabet, the script used to write these languages are different except for a few languages that share a common script, such as Hindi and Marathi. In the last column, we show how the name “Gandhi” is written in all 10 scripts.

recognition of 13 Indic languages that use 10 different scripts. Figure 3.1 shows how the name *Gandhi* is written in the 10 scripts.

Though there had been many attempts to develop OCRs for Indic scripts from the 1970s to the beginning of the last decade [168, 169, 170], methods that can scale across languages and yield reasonable results over a wide variety of documents were not devised. Inherent challenges with the scripts and languages and lack of large-scale annotated data hindered the development of Indic languages OCR for decades. Following the success of Connectionist Temporal Classification (CTC) for speech transcription, CTC has been adapted for recognition of handwritten text [86], printed text [98, 97] and scene text [100, 46]. Most popular open source OCR tools such as Tesseract, EasyOCR ⁸ and ocrpy use a CTC-based model for text recognition. The primary reason for the popularity of this approach is that text on a word or line images can be recognized without the need for sub-word segmentation.

Segmenting a word into sub-word units is much more difficult for Indic languages than for English [171]. Another challenge in developing recognizers for Indic languages is the complex relationships between atomic units of the script (visual), language (text), and machine representation. In a script, the atomic unit is an isolated symbol (a glyph), and from the language perspective, the atomic unit is an *akshara* or an orthographic syllable. And for the machine representation of text, the atomic unit is a Unicode point. An *akshara* can be a combination of multiple glyphs in the script. Similarly, an *akshara* is often represented by a sequence of Unicode. Splitting text at *aksharas* and mapping from *aksharas* to the corresponding Unicode sequence requires knowledge of the language and script [171, 6]. For these reasons, approaching text recognition as a sequence modeling problem using CTC has become the de-facto choice for OCR of Indic languages [171, 99, 127]. This approach can directly map a sequence of features from the word or line image to a target Unicode sequence and an explicit alignment between the feature sequence and the output Unicode sequence is not required during training.

In this work, we present a comprehensive empirical study of the various design choices involved in building a CTC-based printed text recognition model for Indic languages. Major contributions of this chapter are i) we study the challenges to text recognition of Indic scripts and investigate how the seq2seq transcription using CTC helps to overcome some of these challenges, ii) for 13 Indic languages, we empirically compare the performance of four types of CTC-based text recognition that differs in terms of feature extraction and sequence encoding, iii) investigate the effectiveness of synthetic data as an alternative to real training data, iv) compare performance of our recognition model against publicly available OCR tools and show that our model yields better or comparable performance for 8 of the 13 languages, and v)

⁸<https://github.com/JaidedAI/EasyOCR>

introduce the largest ever dataset for printed text recognition of Indic languages that has more than 1.2 million word images (nearly 100 thousand word images per language).

3.2 Related works

This section summarizes previous works that deal with printed text recognition of Indic scripts. We organize the section into three subsections that reflect the evolutionary progress made in this space over the years.

3.2.1 First-generation OCRs [1970 - 2000]

First-generation OCRs [168, 169] for Indic languages typically follow a template-matching style approach for character matching and use intuitive features such as shape and water reservoir. Pal and Chaudhuri [172] provide an excellent summary of methods developed in this period. Sinha and Mahabala [168] use a syntactic pattern analysis system with an embedded picture language to recognize the Devanagari script. Structural representations for each symbol of the Devanagari script are stored beforehand in terms of primitives and their relationships. The input word is first digitized, cleaned, thinned, and segmented (segmented to isolated symbols/glyphs) and labeled (local feature extraction process). Recognition involves searching primitives on the labeled pattern based on the stored description. Contextual constraints are also used to arrive at the correct recognition. Chaudhuri and Pal [169] propose an approach that uses stroke-based features and a tree-based classifier for the recognition of printed Bengali script. The classifier is followed by a template-matching approach to improve accuracy. The character unigram statistics are used to make the tree classifier efficient. Several heuristics are also used to speed up the template-matching approach.

Antani and Agnihotri [173] created a dataset for Gujarati from scanned images and various sources on the internet. They use invariant moments and raw (regular) moments as features. They use K-Nearest Neighbour (KNN) and minimum hamming distance classifiers. Negi et al. [174] use connected component analysis to extract isolated symbols from Telugu words. The segmented symbols are then matched against a stored template bank using fringe distance as the distance measure to perform the classification. Pal and Sarkar [175] propose a system to recognize Urdu script using a combination of topological, contour, and water reservoir-based features and a tree-based classifier.

Lehal and Singh [176] developed a Gurumukhi OCR during this period. Their work uses connected component analysis to extract sub-word components from word images. They use

two types of features called primary and secondary features. Primary features include features like the number of junctions and loops and their positions. Secondary features include the number of end-points and locations, and the nature of profiles of different directions. For classification, a multistage classification scheme is used by combining a binary tree classifier and a nearest neighbour classifier. They improve the performance using a post-processor that makes use of statistical information of Punjabi syllable combinations and certain heuristics.

3.2.2 Second-generation OCRs [2000 - 2012]

The second-generation of OCRs started using statistical and data-driven features such as Discrete Cosine Transform (DCT) and Principal Component Analysis (PCA). This period is also characterized by the use of Machine learning models like Support Vector Machines (SVM) and Artificial Neural Networks (ANN). A comparison and evaluation of state-of-the-art OCR systems developed during this period are presented in [170].

A Tamil and English bilingual text recognition system introduced by Aparna et al. [177] uses geometric moments and DCT coefficients to classify sub-word symbols. For the classification of the symbols, a nearest neighbour classifier is used. Ashwin and Sastry [178] proposed to use SVM for the recognition of Kannada characters. To capture the shapes of the Kannada characters, they extract structural features that characterize the distribution of foreground pixels in the radial and angular directions. Kumar and Ramakrishnan [179] use coefficients of the DCT, Discrete Wavelet Transform (DWT) and Karhunen-Louve Transform (KLT) as features for Kannada OCR. They use both K-NN and Radial Basis Function(RBF) neural network for classification. Kunte and Samuel [180] develop a Kannada OCR that uses Hu's invariant moments and Zernike moments as features. They use an RBF neural network for classification.

Jawahar et al. [181] proposed a bilingual OCR for Hindi and Telugu. Their model uses PCA for feature extraction and SVM for the classification of isolated characters. They evaluate the performance of this system on nearly 200K character images of Hindi and Telugu and report a classification accuracy of 96.7%. Ghosh et al. [182] modified an existing Bangla OCR model to recognize Assamese characters. Their model uses an SVM classifier followed by a spell checker. Rasagna et al. [183] developed a multi-font Telugu OCR using the Histogram of Oriented Gradients (HOG) features and an SVM classifier. They report more than 96% character accuracy on a dataset that has more than 145,000 Telugu character samples in 359 classes and 15 fonts. Neeba and Jawahar [184] empirically evaluated different types of features and character classification schemes for the problem of character classification. Classifiers studied in this work are K-NN, decision trees, Multilayer Perceptron (MLP), Convolutional Neural Network (CNN) and SVM.

Motivated by the success of Hidden Markov Models (HMM) for continuous speech recognition, BBN BYBLOS OCR systems were introduced in the late 1990s [185, 186]. The HMM-based approach followed in THE BBN BYBLOS system does not require word or character level segmentation, and training is language independent. Natarajan et al. [187] extended this approach for Devanagari text recognition. This is one of the first segmentation-free recognizers for an Indic language.

3.2.3 Third-generation OCRs [2012 - 2022]

Third-generation OCRs for Indic scripts are primarily driven by segmentation-free approaches that directly generate a sequence of labels given a word or line image. Sankaran et al. [97] were the first to adopt CTC-based sequence modelling for the problem of printed text recognition of an Indic language. They use an RNN encoder and CTC transcription to map from a sequence of features extracted from a Devanagari word image (i.e., recognition unit is a word) to a sequence of class labels. Profile-based features [188] computed from a 25 x 1 sliding window is used. The model uses *aksharas* as the output classes and employs a rule-based *akshara* to Unicode mapping. They extend this approach in [171] wherein feature sequence from the word image is directly mapped to the Unicode sequence and thereby avoiding the need for a rule based mapping from *aksharas* to Unicode. The CTC-based transcription approach came as a boon for Indic scripts since sub-word segmentation has always been a challenge for most of the Indic scripts. In addition to it, transcribing the word image directly to machine readable form (Unicode sequence) avoided the need to write language-specific rules to map from latent output classes (for example, classes corresponding to the possible set of *aksharas* used in [97]) to a valid Unicode sequence. Krishnan et al. [127] use profile-based features and CTC-based model similar to the one in [171] for the recognition of 7 Indian languages. They evaluate their approach on a test set comprising thousands of document images per language. Their results demonstrate that a unified framework that uses CTC-based transcription works well for the recognition of multiple Indian languages without requiring language/script-specific modules.

Similar to the works discussed above, Hasan et al. [99] proposed using an RNN+CTC model to recognize printed text in Urdu. This work directly outputs a Unicode sequence given an image of a text line (i.e., the recognition unit is a line). Raw pixels extracted from the line image using a 30 x 1 sliding window form the input feature sequence. Our previous work [6] in this space deals with multilingual OCR for 12 Indian languages and English. Our work reports the result of a two-stage system—a script identification module and a recognition module—on thousands of pages per language. Similar to Hasan et al. [99], we used raw pixel values as features and an RNN+CTC classifier. We demonstrated that the RNN+CTC model outper-

forms multiple free and commercial OCR engines for Devanagari text recognition on a newly introduced public dataset.

Chavan et al. [189] conducts a comparative study by evaluating the performance of RNN and multidimensional RNN (MDRNN) [131] encoders when used with the CTC transcription. They use HOG (Histogram of Gradients) features with the RNN encoder and raw pixels for the MDRNN. This study concludes that the MDRNN encoder performs better compared to the RNN encoder. An RNN+CTC transcription model is proposed for the recognition of Bengali script in [190]. This work reports 99%+ character/symbol accuracy for a test set that has word images rendered using more than 20+ fonts. Kundaikar and Pawar [191] study the robustness of CTC-based Devanagari OCR to font and font size variations. Dwivedi et al. [192] proposed to use an encoder-decoder model for the recognition of Sanskrit. The proposed solution achieves under 3% character/symbol error rate on a new private dataset of Sanskrit line images.

Most of the word-based and line-based recognition models for Indic languages that were developed recently rely on CTC-based transcription. In this chapter, we conduct a comprehensive empirical study of CTC-based transcription for both line and word recognition by comparing different types of encoders and features.

3.3 Datasets

This section presents details of the datasets we use in this study. We use three datasets: i) an internal dataset that has nearly 1000 pages per language, ii) a new public dataset of cropped words and lines and the corresponding ground truth transcriptions, and iii) a dataset of synthetic word images for the experiments involving synthetic data.

3.3.1 Internal dataset

Our internal dataset is a subset of a multilingual dataset created for Indic languages by a consortium of universities in India [193]. We shall refer to the original dataset as ‘Indic consortium dataset’. The Indic consortium dataset is annotated with bounding boxes and Unicode transcriptions of all words and lines in the page [194, 195]. For Urdu, owing to the difficulty in annotating words, only line-level annotations were added. The pages are taken from multiple books and scanned using a flatbed scanner. The pages are scanned in 300 PPI for Assamese, Manipuri, and Urdu and in 600 PPI for the rest. The pages mostly contain a single column of text arranged in paragraphs in simple layouts. There are a few pages that contain multi-column text. An insignificant number of pages contain graphics or tables.

Language	Train				Test			
	Books	Pages	Lines	Words	Books	Pages	Lines	Words
Assamese	14	901	23,811	186,471	5	98	3,038	33,636
Bangla	11	900	24,490	199,890	2	100	3,120	26,307
Gujarati	17	899	23591	186555	8	99	3141	32626
Hindi	27	899	23752	195111	6	101	2416	26131
Kannada	22	899	24270	184323	5	101	3448	17269
Malayalam	24	900	24383	183462	7	100	3054	14796
Manipuri	22	897	23466	183304	3	101	2890	25834
Marathi	17	898	23982	191496	3	102	3140	26709
Odia	14	898	24054	192494	3	102	2959	30582
Punjabi	24	899	23725	194900	8	101	2683	32158
Tamil	18	900	24129	181238	5	100	2869	14638
Telugu	24	899	23596	181083	4	101	2791	16748
Urdu	8	866	23250	-	1	93	1829	-

Table 3.1: **Statistics of the internal dataset used in our study.** Each page is annotated with word and line bounding boxes and the corresponding ground truth text transcriptions. For Urdu, word-level annotations are not available.

From the Indic consortium dataset, for 13 languages (the ones shown in Figure 3.1), we sample 1000 pages per language to form our internal dataset. For each language, the pages are split approximately in an 80:10:10 ratio into train, validation, and test splits. The splits are done in such a way that no two splits have pages from the same book. Statistics of the dataset are presented in Table 3.1. Examples of page images in the internal dataset are shown in section A.1 of Appendix A.

3.3.2 *Mozhi* dataset

To the best of our knowledge, there are no large-scale public datasets for the problem of printed text recognition of Indic languages. Most of the early works in this space use a dataset of cropped characters or isolated symbols since these works deal with the classification of dis-

Language	Train		Validation		Test	
	Lines	Words	Lines	Words	Lines	Words
Assamese	9566	79959	1196	9945	1196	10146
Bangla	7579	80113	948	9787	947	10113
Gujarati	8632	79910	1080	10016	1079	10090
Hindi	6525	79762	816	10114	816	10173
Kannada	13462	80085	1683	10088	1683	9838
Malayalam	15112	80146	1889	9893	1889	9980
Manipuri	9765	79691	1221	10254	1221	10061
Marathi	8380	80151	1048	10005	1048	9855
Odia	8260	79945	1033	10089	1033	9994
Punjabi	6726	79931	841	10036	841	10038
Tamil	16074	80022	2010	10021	2009	9974
Telugu	12722	80337	1591	9811	1590	9876
Urdu	9100	-	1138	-	1137	-

Table 3.2: **Statistics of the new *Mozhi* dataset.** It is a public dataset for printed text recognition of cropped words and lines. It has more than 1.2 million words annotated in total. For Urdu, only cropped lines are annotated.

joint characters [180, 181]. Later works that make use of line or word-level annotated data use either internal collections [97, 127, 6, 189, 196] or large-scale synthetically generated samples [192, 191, 99] to train their models. Some of the recent works, including ours, have introduced public datasets containing real samples for a few Indic languages like Hindi and Urdu. They contain a limited number of samples meant only for the evaluation of the models [6, 196]. To the best of our knowledge, except for these two datasets, there are no other public benchmarks of real samples for word, line or page-level printed text recognition of Indic languages. In an attempt to address the scarcity of annotated data and to standardize the data used to train and evaluate printed text recognition models for Indic languages, we introduce a new public dataset named *Mozhi* (meaning “language” or “word” in Tamil and Malayalam) for all 13 languages we study in this work. The *Mozhi* dataset is created by sampling random line-level samples from the Indic consortium dataset. We sampled the line segments randomly



Figure 3.2: Samples of synthetic word images created for Hindi, Malayalam, and Telugu.

so that we can release the dataset without copyright infringement. The dataset contains both line and word-level annotations. For all the 13 languages, cropped line images and their corresponding ground truth text annotations are provided. Word images cropped from these text lines and the corresponding word level ground truths are also included for all languages except Urdu. The dataset has 1.2 million word annotations in total (nearly 100,000 per language), making it the largest ever public dataset of real word images for text recognition in Indic languages. For each language, the line-level data is split randomly in an 80:10:10 ratio to train, validation, and test splits, respectively. Words that are cropped from line images in train split of lines form the train split for word-level recognition. Similarly, for validation and test splits.

Examples of line images in Mozhi dataset and distribution of lengths of words in the dataset are given in section A.2 and section A.3 of Appendix A.

3.3.3 Synthetic dataset

With the advent of deep-learning-based data-driven methods, reliance on large-scale data to build machine learning models has only increased [38, 197, 84]. Deeper networks have more learnable parameters and require larger amounts of data to generalize well. Manually annotating data to train these models is often a tedious and expensive task. An alternative to using real data is to generate synthetic samples. This approach has successfully been used for many Computer Vision problems [124, 198, 199]. The successful adaption of modern machine learning models for the problem of text recognition—whether it is printed, handwritten, or scene text—would not have been possible without large-scale synthetic datasets [200, 201, 124, 202, 8, 203]. Synthetic samples generated by font rendering have been used in many works, including ours, that deal with OCR of Indic languages. [202, 99, 191, 192, 8]

In this work, we investigate the effectiveness of synthetic data as an alternative to real data for training a text recognition model. To this end, we render synthetic word images for the same

set of words that make up the train split of the internal dataset for the specific language. We compare the performances when the recognition model is trained on real and synthetic datasets. We conduct this study for three languages— Hindi, Telugu, and Malayalam. Freely available Unicode fonts are used to render the synthetic word images. The number of unique fonts used for Hindi, Malayalam, and Telugu are 97, 19 and 62, respectively. For rendering text onto images, we use the `convert` tool of ImageMagick ⁹, which makes use of Pango ¹⁰ text layout engine and Cairo ¹¹ graphics library. In order to mimic word images from pages of books that our real data is composed of, we generate images whose background is always lighter (higher intensity) compared to the foreground. Each word is rendered as an image using a random font. Font size, font stylings such as bold and italic, foreground and background intensities, kerning, and skew are varied for each image to generate a diverse set of samples. A random one-fourth of the images are smoothed using a Gaussian filter with a standard deviation(σ) of 0.5. Finally, all the images are resized to a height of 32 while keeping the original aspect ratio. Samples from our synthetic dataset are shown in Figure 3.2

3.4 Text recognition using CTC transcription

Given an input image I of a word or a line, the task of text recognition aims to output the text present on the image in machine readable form. We model the task as a sequence modelling problem using CTC. Input is a sequence of features $\mathbf{x} = x_1, x_2, \dots, x_T$ where $x_t \in \mathbb{R}^D$ are extracted from image I . Output is a sequence of class labels $\mathbf{l} = l_1, l_2, \dots, l_N$, where $l_n \in L$, where L is the output alphabet, i.e., the set of unique class labels. In our case, L is the set of all Unicode code points we are interested in recognizing.

For the below discussion, we use an encoder-decoder style interpretation of the CTC as given in [204].

3.4.1 Extracting feature sequence

Graves et al. [85] first used CTC to transcribe speech to text. In their work, features are extracted along the time axis of the speech signal in a sliding window manner. They use a window of size 10 milliseconds (ms) and a step size of 5 ms. A fixed-size feature vector is extracted at each instance of the sliding window. Each individual unit of the input sequence is

⁹<https://imagemagick.org>

¹⁰<https://pango.gnome.org/>

¹¹<https://www.cairographics.org/>

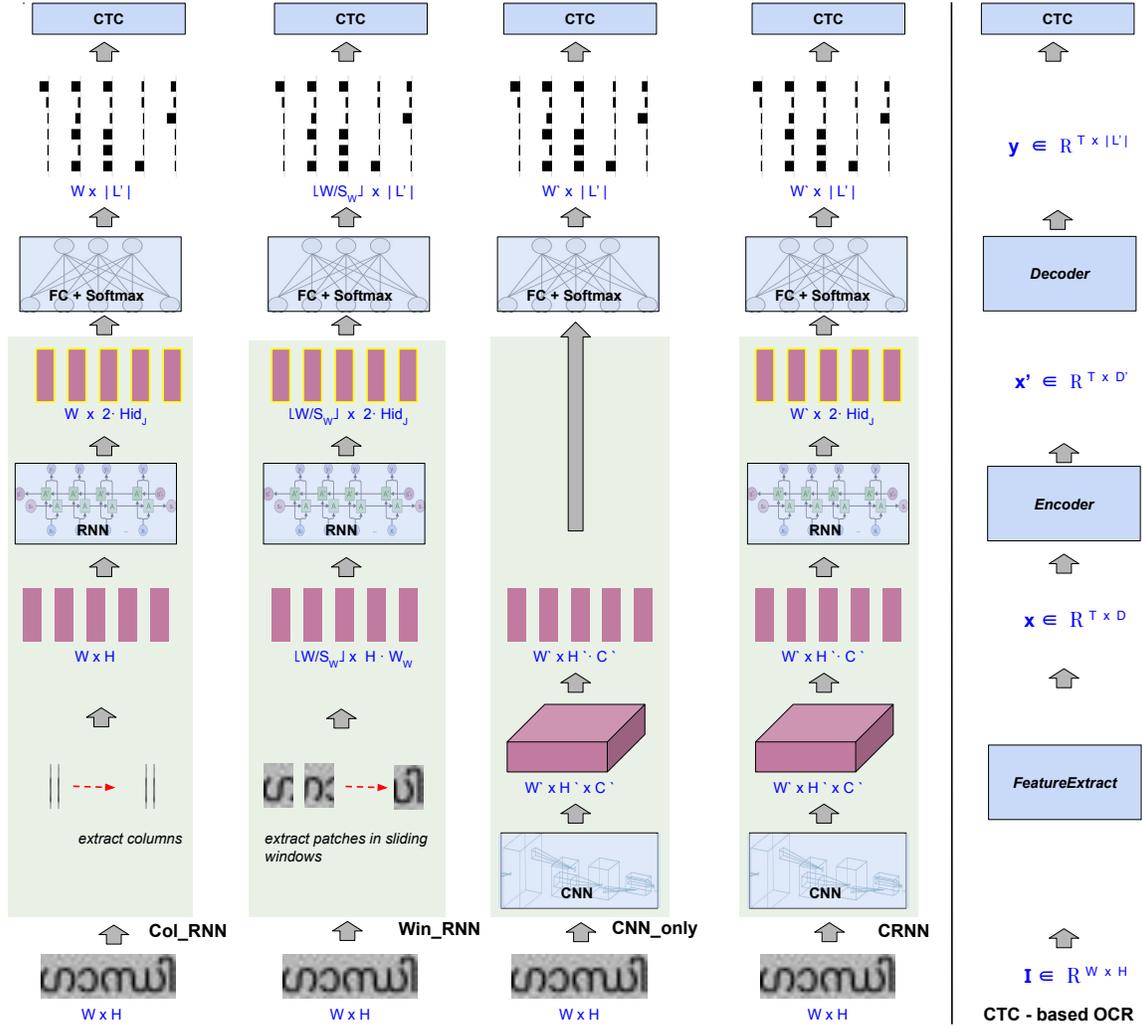


Figure 3.3: **Different CTC-based transcription network configurations.** We study four CTC based text recognition approaches — Col_RNN, Win_RNN, CNN_only and CRNN, that differ in terms of the features extracted and the kind of sequence encoding used. W and H are width and height of the input image I respectively. $|L'|$ is the number of class labels including the *blank* label. Hid_j denotes the number of hidden units in the last layer of the RNN. For the Win_RNN, W_W and S_W are width and step size of the sliding window respectively.

referred to as a ‘time-step’ or a ‘frame’. Unlike a 1D time-varying signal like speech, a grey-scale image is a 2D scalar-valued spatial signal. Therefore, in order to form a 1D sequence of features, methods that use CTC to transcribe text from images conventionally extract features along the horizontal axis of the image[171, 99, 46]. We follow the same approach. That is, feature vectors in the input sequence x correspond to a sequence of horizontal segments of the

given image. Similar to the original work [85], an instance of the input sequence is referred to as a ‘time-step’ or a ‘frame’. The horizontal extent of a frame varies based on how the features are extracted. The feature sequence, \mathbf{x} , is extracted in the same direction as the script is written. That is, for all the languages except Urdu, features are extracted from left to right, and for Urdu features are extracted in the opposite direction.

In summary, given an image $I \in \mathbb{R}^{W \times H}$ (we work with greyscale images), a feature sequence is extracted as follows:

$$\mathbf{x} \in \mathbb{R}^{T \times D} = \text{FeatureExtract}(I) \quad (3.1)$$

3.4.2 Encoder

The job of the sequence encoder is to take the input sequence \mathbf{x} and map it to the encoded representation $\mathbf{x}' \in \mathbb{R}^{T \times D'}$ where D' is the encoding size; i.e., the fixed size to which each feature vector is encoded into.

i.e.,

$$\mathbf{x}' \in \mathbb{R}^{T \times D'} = \text{Encoder}(\mathbf{x}) \quad (3.2)$$

3.4.3 Feature and encoder configurations

Below, we discuss 4 different types of feature and encoder combinations we study in this work.

Col_RNN: In this case, a frame or time step of the input corresponds to a single column of the image. Features are nothing but normalized pixel values (normalized in 0–1 range) of each column. This approach has been used in previous works, including ours, that deal with printed text recognition of Indic languages [6, 127]. From the given image $I \in \mathbb{R}^{W \times H}$, feature sequence $\mathbf{x} \in \mathbb{R}^{W \times H}$ is extracted, and an RNN is used to encode \mathbf{x} to $\mathbf{x}' \in \mathbb{R}^{W \times D'}$.

Win_RNN: Win_RNN uses an approach similar to the original work [85] to extract features in a sliding window manner. A sliding window of size $W_w \times H$ is moved across the image and at each step t , pixel values of the columns in the window are stacked to form the feature vector x_t . If the sliding window is moved with a step size S_w , then \mathbf{x} will be of shape $\lfloor W/S_w \rfloor \times (H \cdot W_w)$. Further, \mathbf{x} is encoded using an RNN to $\mathbf{x}' \in \mathbb{R}^{\lfloor W/S_w \rfloor \times D'}$. Col_RNN is in fact, a Win_RNN with $W_w = S_w = 1$.

CNN_only: Unlike the above two approaches that use pixel intensities as features, a CNN is used to extract features. Given an image I , CNN outputs a feature map $F \in \mathbb{R}^{W' \times H' \times C'}$. F

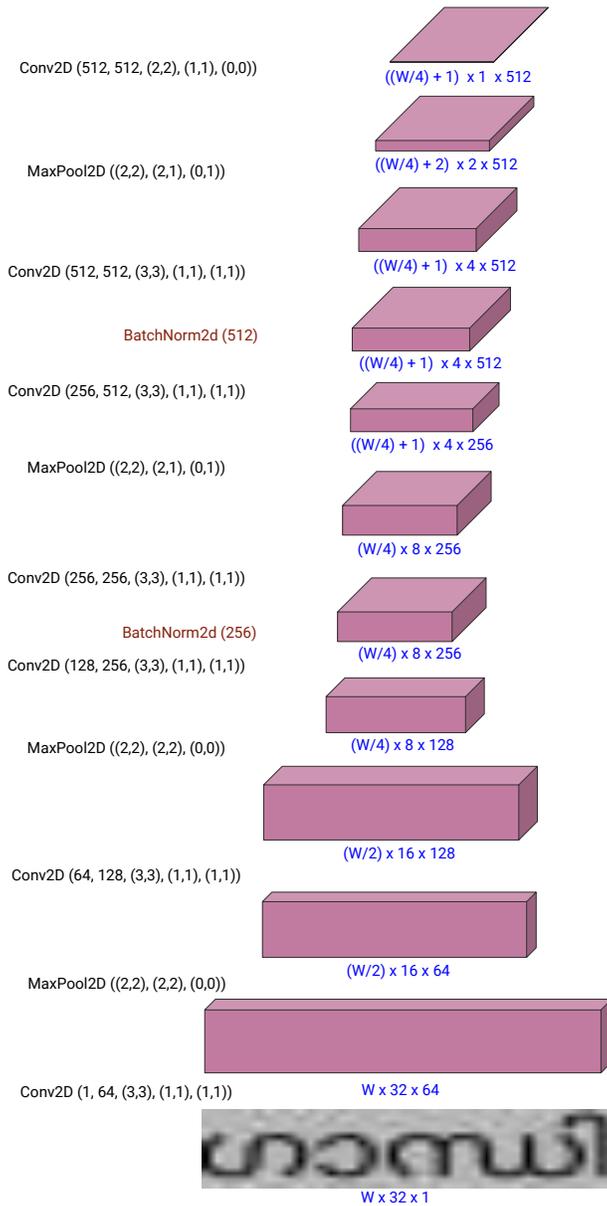


Figure 3.4: **Architecture of the CNN used in CNN_only and CRNN configurations.** This architecture is the same as the one used in the original CRNN paper [46]. The size in *width* \times *height* \times *num.ofchannels* format is written below the image and feature maps. Conv2D, MaxPool2D, and BatchNorm2D denote 2D convolution, 2D max pooling, and 2d Batch Normalization operations respectively. The parameters for these operations are shown in the same format as the one used in PyTorch.

is then reshaped to $W' \times (H' \cdot C')$ to form the sequence of features \mathbf{x} . In other words, from the feature map F , a sequence of W' feature vectors, each of size $H' \cdot C'$ are formed. In this

configuration, there is no dedicated sequence encoder. Or the encoder is an identity mapping. That is, in this configuration, $\mathbf{x}' = \mathbf{x}$.

CRNN: In this configuration, features are extracted using a CNN, and an RNN is used to encode the sequence of features. This approach called, "CRNN" was proposed by Shi et al. [46] for English scene text recognition. Similar to the CNN_only configuration, a CNN is used to obtain a feature sequence $\mathbf{x} \in W' \times (H' \cdot C')$. \mathbf{x} is then encoded using an RNN to form the final encoded sequence $\mathbf{x}' \in \mathbb{R}^{W' \times D'}$. All configurations we discuss above, but CNN_only, use an RNN at the encoder, and the size of the encoding D' depends on the number of units in the last recurrent layer. If a J layer bi-directional RNN is used, then $D' = 2 \cdot Hid_J$, where Hid_J is the size of the hidden state of the last layer of the RNN.

The RNN encoder used in the above configurations is essentially encoding a 1D sequence of features. Or in other words, the RNN captures long-term dependencies along the horizontal axis only. An obvious choice to capture dependencies along multiple axes is MDRNN. However, as discussed in subsection 2.1.2.3, the MDRNN encoder is empirically found to be an overkill for the text recognition problem. For this reason, we do not experiment with MDRNNs in this work. Instead, our CNN_only and CRNN configurations discussed above use convolutional layers to model dependencies along both axes of an image.

3.4.4 Decoder

The encoded features \mathbf{x}' are projected to the size equivalent to the number of the output classes using a linear projection layer followed by Softmax normalization. This step can be viewed as the decoding part of the CTC as interpreted in [204]. The original output alphabet L is augmented by adding an extra label \sim for blank. That is, $L' = L \cup \sim$. A blank label corresponds to the case when we want to assign no label for an input. The result of Softmax normalization at a time-step can be interpreted as class conditional probabilities at the time-step. Or in other words, Softmax yields the posterior distribution over the classes.

In summary, given the sequence of encoded features \mathbf{x}' ,

$$\mathbf{y} \in \mathbb{R}^{T \times L'} = Decoder(\mathbf{x}') \tag{3.3}$$

where each $y_t \in R^{L'}$ represents activations at time step t . Thus y_t^k is a score indicating the probability of k^{th} label at time step t .

3.4.5 Transcription using CTC

The objective of the CTC transcription is to find the most probable sequence of class labels given \mathbf{y} . Let L'^T be the set of sequences of length T over the alphabet L' . An element of L'^T is called a *path* and denoted by π . CTC assumes that the target label sequence's length is always shorter or equal to the length of the input sequence (T). Therefore, we consider all length T sequences over the alphabet L' .

If we assume that network prediction at a particular time step is independent of the predictions at other timesteps, the probability of a path is the product of probabilities of individual labels in the path, at respective time steps. That is, probability of a path π , given input sequence \mathbf{x} is:

$$p(\pi|\mathbf{x}) = \prod_{t=1}^T y_t^{\pi^t} \quad (3.4)$$

where $y_t^{\pi^t}$ is the probability of t^{th} label of the path π .

A many-to-one sequence to sequence mapping \mathcal{B} is defined from the set of all possible paths to the set of all possible labellings whose length is at most T . i.e., $\mathcal{B} : L'^T \mapsto L^{\leq T}$. Note that the paths are defined over the augmented alphabet L' , and the labellings are defined over the original alphabet L . \mathcal{B} maps a path π to a labelling \mathbf{l} by removing the blank labels and the repeated labels. For example, the path “g~~aa~nd~hh~ii” is mapped to the labelling “gandhi”.

Given an input feature sequence \mathbf{x} , the conditional probability of a labelling \mathbf{l} is the sum of probabilities of all paths in L'^T that maps to \mathbf{l} . That is,

$$p(\mathbf{l}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{l})} p(\pi|\mathbf{x}) \quad (3.5)$$

Explicitly computing the summation in Equation 3.5 is difficult since there are a large number of paths that map to a given labelling. Inspired by the forward-backward algorithm for HMMs [205], Graves et al. [85] proposed a dynamic programming algorithm for the efficient computation of Equation 3.5.

3.4.6 Training

Let the training dataset be $S = \{I_i, \mathbf{l}_i\}$ where I_i is a word or line image and \mathbf{l}_i is the corresponding ground truth labelling. The objective function for training the encoder-decoder neural network for CTC transcription is based on the principle of Maximum Likelihood. Minimizing

the objective function maximizes the log likelihoods of the ground truth labelling. Therefore the objective function used is,

$$\mathbb{O} = - \sum_{I_i, \mathbf{l}_i \in S} \log p(\mathbf{l}_i | I_i) \quad (3.6)$$

The above objective function can be optimized using gradient descent and backpropagation.

3.4.7 Inference

At the time of inference, given an input sequence \mathbf{x} , our CTC-based classifier needs to output the labelling \mathbf{l}^* that has the highest probability as defined in Equation 3.5. Thus the CTC-based classifier can be expressed as a function $h(\mathbf{x})$, where

$$h(\mathbf{x}) = \arg \max_{\mathbf{l} \in L^{\leq T}} p(\mathbf{l} | \mathbf{x}) \quad (3.7)$$

Similar to the HMMs, this step, where the most probable labelling is found, is called *decoding*. But there is no tractable algorithm for finding the labelling \mathbf{l} that maximizes $p(\mathbf{l} | \mathbf{x})$. Graves et al. [85] proposed two approximate methods instead — best path decoding and prefix search decoding. In this work, we use the former. ‘Best path decoding’, as the name suggests, output the labelling corresponding to the most probable path as the most probable sequence. i.e.,

$$h(\mathbf{x}) \approx \mathcal{B}(\pi^*) \quad (3.8)$$

where π^* is a path formed by concatenating the most probable labels in each time step. Note that best path decoding is an approximation, and there is no guarantee that it will always yield the most probable labelling.

3.5 Experimental setup

Details of the steps taken to preprocess the data, hyperparameters of the encoder and decoder, and specifics of the network training are presented in this section. We also summarize the evaluation metrics we use to evaluate the text recognition performance in both recognition-only and end-to-end settings.

3.5.1 Implementation details

In all experiments, input images of cropped words or lines are converted to greyscale and resized to a height of 32 pixels while keeping the original aspect ratio. There is no separate validation split for the internal dataset. Therefore, for all languages, we form a validation split by taking random 5% pages from each of the books in the train split. Thus validation split of the internal dataset has pages similar to the pages in the train split, while the test split has pages from a different set of books.

While training our models on word or line samples from the internal dataset, the output alphabet L for a language is the set of unique Unicode points found in the respective train split for that language. Similarly, while working with the Mozhi dataset, the output alphabet for a language is the set of Unicode points in the train split for that language in the Mozhi dataset. For any language, the alphabet used for the line recognition model will have only one extra label — the label corresponding to white space — compared to the alphabet used for the word recognition model of the same language. Word images in the synthetic dataset for a language (see subsection 3.3.3) are rendered using the same set of words in the train split for that language in the internal dataset. Therefore, the output alphabet used while training on synthetic data for a language is the same as the alphabet used for the language while training on the internal dataset.

For Win_RNN, the sliding window width W_W is 20, and the step size W_S is 5. The RNN we use in Col_RNN, Win_RNN and CRNN has the number of layers $J = 2$. We use a bi-directional LSTM with 256 hidden units per direction in each layer. Therefore the size of the output of the RNN at each time step is $2 \times 256 = 512$. The CNN block in CNN_only and CRNN models have the exact same architecture (shown in Figure 3.4) as in the original CRNN paper [46]. Our models are implemented in PyTorch [206]. We build on an existing third-party implementation of CRNN architecture¹². All our models are trained on a single Nvidia GeForce 1080 Ti GPU. While training on the internal dataset, we train all the models for 15 epochs, and the CRNN models trained on the Mozhi dataset are trained for 30 epochs. A batch size of 64 and 16 is used for word and line recognition models, respectively. We use RMSProp [207] as the optimizer. For Col_RNN and Win_RNN, a learning rate of $10e - 03$ is used. For CNN_only and CRNN variants, a slower learning rate of $10e - 04$ was found to be better for faster convergence. The checkpoint that yields the highest Character Accuracy (refer to subsection 3.5.2) on the respective validation data is saved for evaluation on the test split.

¹²<https://github.com/Holmeyoung/crnn-pytorch>

3.5.2 Evaluation

We require to evaluate text recognition in three scenarios, i) word OCR: recognition of a cropped word image, ii) line OCR: recognition of a cropped line image, and iii) page OCR: end-to-end text recognition where input is a document image. In all three cases, our primary evaluation metric is Character Accuracy (CA), which is based on the Levenshtein distance between predicted and ground truth strings.

For a formal definition of CA, let us denote the predicted text for a word/line/page as l_i and the corresponding ground truth as g_i . If there are N such samples, CA is defined as

$$CA = \frac{\sum_i \text{len}(g_i) - \sum_i LD(l_i, g_i)}{\sum_i \text{len}(g_i)} \times 100 \quad (3.9)$$

where len is a function that returns the length of the given string, and LD is a function that computes the Levenshtein distance between the given pair of strings. Note that Character Error Rate (CER), which is another commonly used metric for OCR evaluation, is essentially $100 - CA$.

For word OCR and line OCR, in addition to CA, we report Sequence Accuracy (SA). It is the percentage of samples for which the prediction is fully correct (i.e., $LD(l_i, g_i) = 0$). In the case of a word recognition model, SA is the same as ‘word accuracy’ or ‘accuracy’ as used in scene text recognition literature. For a line recognition model, SA is nothing but the percentage of text lines that are transcribed correctly.

For evaluating page-level OCR where the input is a document image, we use a standard OCR evaluation toolkit. A modern port¹³ of the original ISRI Analytic Tools for OCR Evaluation [208] is used. Using the ISRI toolkit, we report Character Accuracy (CA) and Word Accuracy (WA). In ISRI toolkit, CA is computed in the same manner as given in Equation 3.9. Word accuracy is computed by aligning the sequence of words in the prediction l_i and the sequence of words in ground truth g_i , and finding the Longest Common Sub-sequence (LCS) of the two. For a set of pages,

$$WA = \frac{\sum_i \text{len}(LCS(l_i, g_i))}{\sum_i \text{len}(g_i)} \times 100 \quad (3.10)$$

where len returns the number of words in a given sequence of words.

¹³<https://github.com/eddieantonio/ocreval>

Language	Word								Line							
	Col_RNN		Win_RNN		CNN_only		CRNN		Col_RNN		Win_RNN		CNN_only		CRNN	
	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA
Assamese	98.6	95.4	97.6	92.9	98.3	96.0	99.0	96.5	99.1	78.8	98.1	65.5	99.0	73.8	99.2	80.8
Bangla	99.1	97.0	98.3	94.5	99.2	97.3	99.4	97.9	99.1	73.7	98.4	59.6	99.2	73.9	99.4	79.7
Guajrati	96.2	92.4	95.1	89.5	96.2	90.9	96.5	93.9	96.5	66.3	94.6	49.9	96.0	62.2	96.9	67.4
Hindi	97.6	95.1	96.3	92.3	97.4	94.2	98.2	96.3	98.8	64.9	97.8	48.8	98.8	64.3	98.9	66.9
Kannada	97.4	88.9	96.4	84.7	96.7	85.8	97.7	90.7	97.4	49.2	96.4	38.2	97.0	42.9	97.5	49.4
Malayalam	99.5	96.6	99.3	95.6	98.0	83.7	99.7	97.7	99.5	84.8	99.3	80.5	98.5	49.7	99.7	86.9
Manipuri	98.6	95.4	97.8	92.8	98.2	93.1	99.0	96.9	99.4	79.9	98.7	67.6	99.4	79.2	99.5	80.9
Marathi	99.0	96.2	98.5	94.2	98.9	95.0	99.2	96.9	99.1	71.2	98.4	55.0	99.0	67.1	99.1	71.7
Odia	96.8	93.5	95.7	90.8	96.9	93.7	97.2	94.8	97.9	73.9	96.8	60.2	97.9	70.5	98.1	74.2
Punjabi	99.1	97.7	98.4	96.4	99.2	97.8	99.5	98.7	99.1	76.6	98.3	62.5	99.2	78.7	99.3	79.9
Tamil	97.9	91.0	97.4	88.4	97.3	87.2	98.0	91.8	96.3	43.8	95.9	40.7	96.2	41.2	96.5	45.0
Telugu	96.3	91.4	95.3	86.8	96.4	92.0	96.8	93.6	96.5	68.9	95.2	50.5	96.7	68.4	97.0	75.0
Urdu	-	-	-	-	-	-	-	-	93.9	23.2	75.8	4.2	91.9	17.0	93.5	24.1

Table 3.3: **Recognition-only results on the validation splits of an internal dataset.** The left half shows results when the recognition unit is a word (word OCR) and the other half shows results of the line OCR setting. Note that we train a separate model for each language in each setting. In both settings, we report Character Accuracy (CA) and Sequence Accuracy (SA) for all four model configurations. The numbers shown in bold are the best CA and SA among the four feature + encoder configurations.

3.6 Experiments and results

In this section, we present the details of the experiments we conduct and report and discuss the results.

3.6.1 Comparing different feature + encoder configurations

We evaluate the performance of the 4 different feature + encoder configurations (see subsection 3.4.3) on the internal dataset for word and line recognition performance. The results of this experiment are shown in Table 3.3. Models are trained on the train split and evaluated on the validation split of the respective language.

Language	Word		Line	
	CA	SA	CA	SA
Assamese	99.1	96.9	99.4	73.3
Bangla	98.9	96.7	98.7	75.4
Guajrati	97.2	93.0	97.5	53.1
Hindi	97.4	94.0	98.0	51.2
Kannada	94.2	86.5	93.7	49.8
Malayalam	99.3	94.8	99.3	77.7
Manipuri	98.1	93.7	98.7	63.2
Marathi	99.6	97.9	99.5	81.7
Odia	97.8	94.8	98.0	61.4
Punjabi	99.0	97.8	99.2	78.8
Tamil	95.4	84.5	95.9	41.2
Telugu	99.0	94.7	98.9	69.4
Urdu	-	-	93.5	7.5

Table 3.4: **Recognition-only performance using CRNN on the test splits of the internal dataset.** Here we use the CRNN checkpoint that yielded the highest CA on validation split.

Note that each CA and SA pair in Table 3.3 correspond to a CTC-based network that was trained separately for a certain combination of language, recognition unit (line or word) and feature + encoder configuration (Col_RNN, Win_RNN, CNN_only or CRNN). In all cases except for Urdu line recognition, CRNN performs the best among the 4 configurations. CRNN performing better than Col_RNN and Win_RNN substantiate the superiority of features learnt using a CNN compared to handcrafted features like normalized pixel values.

Similarly, improved performance for CRNN compared to CNN_only configuration validates the need for modelling long-term dependencies in the word or line images. Unlike fully connected networks, neurons in successive layers in a CNN ‘see’ only a local region of the input feature map. In order to build a CNN where a neuron in the last layer has a receptive field covering the entire input, we need to stack a large number of convolutional layers. The 7-layer CNN we use is not deep enough to model long term dependencies along the horizontal axis. This is compensated by using a sequence encoder (a bi-directional LSTM) that efficiently mod-

els long-term dependencies in both forward and backward directions along the horizontal axis of the input.

Since CRNN works the best except for Urdu line-level recognition, we only evaluate the CRNN configuration on the test set. These results are reported in Table 3.4. The Col_RNN configuration that performs the best for Urdu line recognition on the validation split yields a CA of 92.0 and SA of 3.6 on the Urdu test split.

3.6.2 Page OCR evaluation on the internal dataset

In the page OCR setting, input to the OCR is a document image and the expected output is the transcription of the textual content in the image. For a page OCR, a typical approach is to use a page segmentation step that detects lines or words followed by a word or line-level text recognition model. Finally, text transcriptions for individual lines or words are combined to form a page-level transcription. In a typical scenario, end-to-end OCR involves layout analysis to identify different document layout objects and techniques to identify the reading order. This work focuses on text recognition (i.e., recognizing text present in a given word or line image). Detecting words or lines on document images is beyond the scope of this work. Therefore we build an end-to-end page OCR pipeline where text detection is done using existing methods or tools, and recognition of words or lines is done using our CRNN models. Once we have transcriptions for individual words or lines, we concatenate them in the same reading order as found by the page segmentation tool. In order to assess the upper bound on the end-to-end text recognition performance of our CRNN model, we evaluate an end-to-end pipeline where the gold standard line or word detections are used. We further compare the results of our end-to-end OCR with two publicly available OCR tools.

We use two public end-to-end OCR tools. Tesseract and Google cloud vision OCR’s DOCUMENT_TEXT_DETECTION API. The former is an open-source OCR while the latter is a commercial, cloud-based solution. Tesseract version 5.1.4 is used. We used page segmentation mode (`psm`) 3 of Tesseract, which does automatic page segmentation but without automatic script detection. We used model checkpoints provided in the official Tesseract repository ¹⁴. For Manipuri, since there is no trained model available, we used the model trained for Bengali since both languages use the same script. At the time we used Google cloud vision OCR, Manipuri was not supported, and Urdu and Odia OCRs were in the experimental stage.

While building end-to-end pipelines that use our recognition model with automatic text detection, we try out text detections from the following: i) line and word detections from Tesseract, ii) line and word detections from Google cloud vision OCR, and iii) line detections

¹⁴<https://github.com/tesseract-ocr/tessdata>

Language	GT detection + our CRNN				Automatic detection + our CRNN										End-to-end OCR tools			
	GT word		GT line		Tesseract word		Tesseract line		Google word		Google line		Scale space line		Tesseract		Google	
	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA	CA	SA
Assamese	99.3	97.0	99.4	97.2	94.4	94.5	96.8	94.5	94.6	92.0	98.7	95.7	98.6	96.5	92.7	91.2	90.0	86.0
Bangla	99.1	97.3	99.0	96.8	97.7	96.2	98.6	96.3	91.8	92.0	96.3	92.5	96.4	94.9	93.5	96.2	84.0	91.3
Gujarati	98.0	93.7	97.7	91.9	91.7	81.9	93.6	88.4	88.6	74.3	93.1	79.6	75.2	67.7	96.9	92.4	93.0	95.2
Hindi	98.1	96.0	98.0	95.6	94.6	91.5	95.5	93.2	92.4	87.6	95.7	92.3	96.1	93.7	95.0	93.3	95.2	97.3
Kannada	95.6	89.2	95.9	86.4	70.0	61.1	70.6	62.1	72.5	63.9	72.2	64.2	67.2	64.6	94.9	85.1	85.7	84.6
Malayalam	99.4	98.0	99.3	97.9	98.1	91.6	98.8	96.7	96.3	83.2	97.7	91.9	98.2	96.3	96.2	78.7	88.0	74.8
Manipuri	98.4	94.7	98.7	94.9	95.9	89.2	97.4	93.7	86.9	64.0	96.4	87.5	98.0	93.8	90.9	80.6	85.7	77.4
Marathi	99.6	98.2	99.5	98.0	96.3	96.1	97.2	96.2	97.4	93.2	98.1	95.6	86.5	82.9	97.9	97.4	98.3	98.4
Odia	98.6	95.4	98.0	94.5	95.4	89.2	96.9	93.3	86.6	67.1	96.0	89.5	98.3	94.8	94.0	83.6	92.6	90.0
Punajbi	99.2	98.3	99.3	97.9	94.7	91.6	96.0	95.1	91.5	85.0	97.6	95.3	96.5	95.7	93.2	89.8	92.7	96.7
Tamil	96.1	85.6	96.5	85.4	92.4	80.0	93.6	83.4	88.0	60.6	93.7	79.1	93.3	82.0	79.3	42.4	92.5	93.1
Telugu	99.1	95.1	98.9	94.0	89.3	83.2	89.2	83.5	91.4	71.9	96.3	86.0	83.8	79.9	93.7	79.3	94.2	89.2
Urdu	-	-	94.7	81.5	-	-	88.9	74.4	-	-	90.0	68.8	56.4	45.5	68.3	26.2	92.7	85.7

Table 3.5: **Performance of our page OCR pipelines compared to other public OCR tools.** In this setting, we evaluate text recognition in an end-to-end manner on the test split of the internal dataset. For text detection, either gold standard word/line bounding boxes or automatic text detection tools are used. Under ‘End-to-end OCR tools ’we show results of Tesseract and Google Cloud Vision OCR. Given a document image, these tools output a transcription of the page along with the bounding boxes of the lines and words detected. Under ‘GT detection + our CRNN’, we show results of end-to-end pipeline where gold standard word and line detections are used. For instance, ‘GT Word’ means we used ground truth (GT) word bounding boxes and the CRNN model trained for recognizing words, for that particular language. Under ‘Automatic detection + our CRNN’ we show results of the following end-to-end pipelines: i) Tesseract word detections + our CRNN word model, ii) Tesseract line detections + our CRNN line model, iii) Google OCR word detections + our CRNN word model, iv) Google OCR line detections + our CRNN line model and v) Scale space [209] line detections + our CRNN line model.

using a third party implementation [209] of the scale space method proposed by Manmatha et al. [210]. The Tesseract and Google OCR that we use for text detection is the same as the ones we use in the end-to-end setting. Along with the text transcriptions, these tools provide the bounding boxes of detected lines and words. These detections are used with our recognition model to build an end-to-end OCR. The parameters for scale space method are set as instructed in [209]. We use different parameter values for different languages. Values of the parameters are determined based on the average height and aspect ratio of the text lines in the train split.

The lines detected using the scale space technique are ordered in the default reading order. The default order may not be correct in the case of complex layouts and multi-column text. Since our internal dataset comprises mostly document images with single-column text, using default reading order to order line detections works for most of the images. We found the word detections using scale space method highly noisy with many under segmentation cases. For this reason, we do not try out an end-to-end pipeline that uses word detections from the scale space method.

Results of all the end-to-end evaluations are reported in Table 3.5.

Model	Trained on	finetuned on	Hindi		Malayalam		Telugu	
			CA	SA	CA	SA	CA	SA
real-only	real	NA	97.4	94.0	99.3	94.8	99.0	94.7
synth-only	synthetic	NA	93.3(↓ 4.2%)	84.9(↓ 9.7%)	98.2(↓ 1.1%)	88.4(↓ 6.7%)	95.5(↓ 3.5%)	80.0(↓ 15.5%)
synth + 0.1 real	synthetic	10% real	96.9(↓ 0.5%)	92.8(↓ 1.2%)	99.2(↓ 0.1%)	94.2(↓ 0.6%)	98.7(↓ 0.3%)	93.1(↓ 1.6%)
synth + 0.5 real	synthetic	50% real	97.4(↑ 0.0%)	94.0(↑ 0.0%)	99.3(↑ 0.0%)	94.7(↓ 0.1%)	99.0(↑ 0.0%)	94.7(↑ 0.0%)

Table 3.6: **Synthetic to real transfer learning results.** The numbers reported are evaluation results on the word-level test split of the internal dataset for the respective language. ‘real-only’ and ‘synth-only’ are models trained exclusively on real and synthetic samples, respectively. We finetune the synth-only models using varying amounts of real data. Synth + 0.1 real is the case when the synth-only model is finetuned with a random 10% of the real samples. Similarly, synth + 0.5 real is a model where the synth-only model is finetuned with half of the real training data. The values shown in brackets next to accuracy values are relative gain or drop in accuracy compared to the corresponding accuracy for real-only model. ↓ indicates drop in accuracy and ↑ is used when there is no change in accuracy. For all languages, the synth-only model can be finetuned for performance on par with the real-only model using only half of the real training data.

3.6.3 Transfer learning from synthetic to real

We investigate the effectiveness of synthetic data for training word recognition models. The results of these experiments are reported in Table 3.6. For three languages—Hindi, Malayalam and Telugu—we train word-level CRNN models on the respective synthetic datasets (see subsection 3.3.3). These models are called ‘synth-only’ since they are trained purely on synthetic data. The validation data used for these experiments contains not synthetic but real samples. The validation data we use is the same as the validation data used for the respective languages for training word-level models on the internal dataset. The model checkpoint that yields the

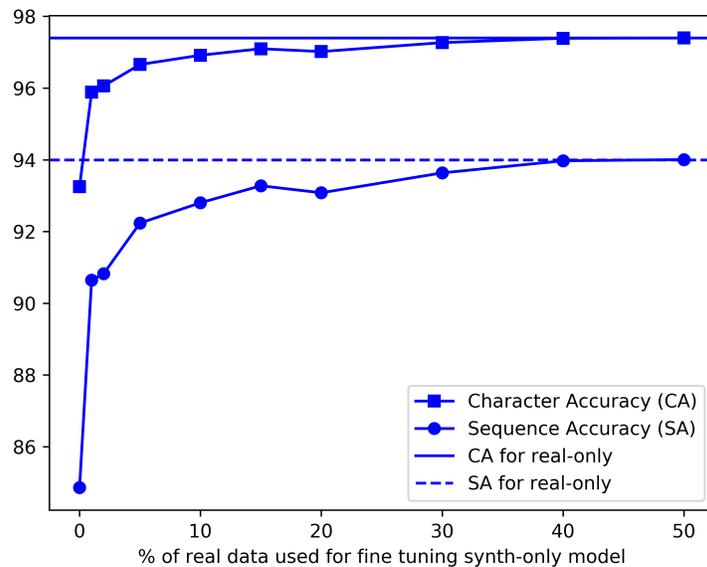


Figure 3.5: **Transfer learning from real to synth by finetuning the synth-only model.** The plot shown above is for Hindi. The model that uses no real data for finetuning is called ‘synth-only’. That is, it is trained purely on synthetic data. This model is then finetuned with real data that is randomly sampled from the train split of the internal dataset. With finetuning on just 1% of the real data, the CA jumps from 93.3 to 95.9, and the SA jumps from 84.9 to 90.6. With 40% or more real data, the performance gain from the finetuning starts saturating and matches the performance of the real-only model.

best CA on the validation data is saved for evaluation on the test split. The saved checkpoint is then evaluated on the respective word-level test split of the internal dataset. In Table 3.6, a ‘real-only’ model for a language is a CRNN model that is trained on the word samples in the train split of the internal dataset for the language. Thus the results of real-only models we report in the table are the same as the numbers we report for the three languages for word-level CRNN on validation and test splits in Table 3.3 and Table 3.4, respectively. As expected, for all languages, the synth-only model has a lower performance than the real-only model. But the drop in performance is not large. The relative drop in CA, compared to the corresponding real-only model, is not more than 5% for any language.

We conduct transfer learning experiments where we finetune the synth-only models using real data. The results of these experiments are shown in the last two rows of Table 3.6. ‘synth + 0.1 real’ is a model where we finetune the synth-only model using a random 10% of the real samples from the train split of the internal dataset. Similarly, ‘synth + 0.5 real’ is finetuned using a random 50% of the internal dataset’s train split. With finetuning on 50% of the real

data, the CA is the same as that of the real-only models for all the languages. In Figure 3.1, we show how increasing the amount of real data used for finetuning improves the performance of the synth-only model for Hindi. Both CA and SA saturate and matches the performance of the real-only model with 40% of the real data.

3.6.4 Evaluating CRNN on Mozhi dataset

Results of word and line recognition on the new Mozhi benchmark dataset using CRNN are reported in Table 3.7. For a language, for both word and line, we train a CRNN from scratch on the respective train split of the Mozhi dataset.

Language	Validation				Test			
	Word		Line		Word		Line	
	CA	SA	CA	SA	CA	SA	CA	SA
Assamese	98.0	95.1	98.7	76.9	98.9	96.2	99.2	76.8
Bangla	99.2	97.0	98.4	69.5	99.0	96.9	98.1	68.4
Guajrati	98.3	95.3	97.8	61.4	98.0	94.9	97.4	63.1
Hindi	98.2	95.9	98.8	61.8	98.1	95.5	98.8	63.5
Kannada	96.6	88.4	97.1	53.2	97.1	88.7	97.5	53.9
Malayalam	99.6	97.2	99.5	86.0	99.5	97.3	99.5	87.3
Manipuri	98.4	95.8	99.1	80.3	98.4	95.9	99.2	79.4
Marathi	99.2	96.7	99.2	73.5	99.3	97.0	99.3	73.8
Odia	98.1	85.2	98.9	74.4	97.5	94.3	98.8	73.1
Punjabi	99.4	98.6	99.4	81.9	99.2	98.2	99.3	79.7
Tamil	98.2	92.0	98.5	70.1	98.0	91.6	98.3	68.1
Telugu	99.2	96.1	98.9	74.2	99.1	95.4	98.9	71.7
Urdu	-	-	93.6	26.5	-	-	93.8	24.2

Table 3.7: **CRNN evaluation on Mozhi dataset.** For each language, we train both word and line level CRNN models on the respective train split of the Mozhi dataset. The models are trained from scratch.

3.7 Summary

In this chapter, we conduct an empirical study of different CTC-based models for word and line recognition for 13 Indic languages. Our study concludes that CRNN, which uses a CNN for feature representation and a dedicated RNN-based sequential encoder, works the best. Using existing text detection tools and our recognition models, we build page OCR pipelines and show that our approach works better than two popular OCR tools for most languages. We create font-rendered synthetic word image samples and train word recognition models for 3 languages. We conduct a transfer learning experiment where we analyze how the performance of models trained purely on synthetic data improves when finetuned on real data. Our results show that models pretrained on synthetic data and then finetuned with half of the real training data perform equally well as the models trained on the whole of the real data. The results suggests that font-rendered synthetic samples are a good alternative to real data to train text recognition models for low-resource Indic languages. We also introduce a new public dataset for cropped word/line recognition in 13 Indic languages, which has more than 1.2 million annotated words in total.

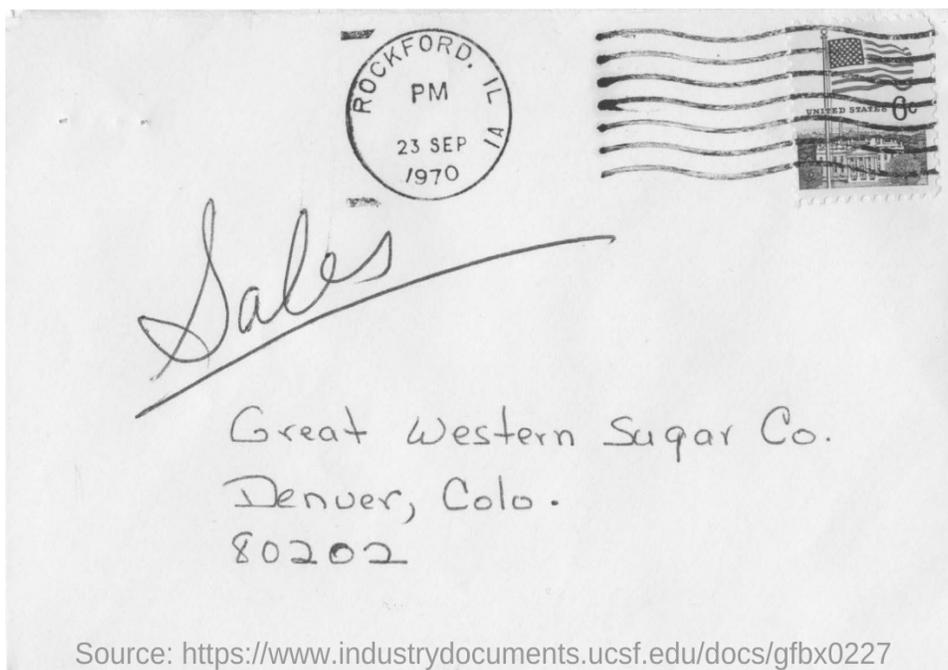
Chapter 4

DocVQA: VQA on business documents

In the previous chapter we presented an empirical study comparing different feature and encoder configurations for CTC-based printed text recognition of Indic languages. We believe our study and the new benchmark dataset would broaden the language coverage of DIA and make it more universal. In this chapter and the subsequent chapters, we present our works that contribute to deepening the machine understanding capabilities of DIA. As discussed in section 1.3, question answering tasks and similar tasks that require human-like reasoning skills gained a lot of traction since the second half of last decade. However, QA/VQA is not actively studied in the context of document images. Previous datasets that deal with VQA on document images have template-based questions, and the images are limited to certain types of document images, such as book covers, and charts subsection 2.5.2. In this chapter, we present a new task and dataset for VQA on document images from a wide variety of business documents. The dataset consists of 50,000 questions defined on 12,000+ document images. Detailed analysis of the dataset in comparison with similar datasets for VQA and MRC is presented. We report several baseline results by adopting existing VQA and MRC models for the DocVQA. Although the existing models perform reasonably well on certain types of questions, there is a significant performance gap compared to human performance (94.36% accuracy). The models need to improve specifically on questions where understanding the structure of the document is crucial. This work has previously been presented in our publication [5].

4.1 Introduction

As discussed in section 1.4, current research in DIA is generally focused on information extraction tasks that aim to convert the information in document images into machine readable form, such as character recognition [211], table extraction [212] or key-value pair ex-



Q: Mention the ZIP code written?

A: 80202

Q: What date is seen on the seal at the top of the letter?

A: 23 sep 1970

Q: Which company address is mentioned on the letter?

A: Great western sugar Co.

Figure 4.1: **Examples of question-answer pairs in DocVQA.** Answering questions in the new dataset require models not just to read text but interpret it within the layout/structure of the document.

traction [213]. Such algorithms tend to be designed as task-specific blocks, blind to the end-purpose the extracted information will be used for. Progressing independently in such information extraction processes has been quite successful, although it is not necessarily true that holistic document image understanding can be achieved through a simple constructionist approach, building upon such modules. The scale and complexity of the task introduce difficulties that require a different point of view.

In this chapter, we introduce Document Visual Question Answering (DocVQA) as a high-level task dynamically driving DAR algorithms to interpret document images conditionally. By doing so, we seek to inspire a “purpose-driven” point of view in DIA research. In the case of Document VQA, as illustrated in Figure 4.1, an intelligent reading system is expected to

respond to ad-hoc requests for information, expressed in natural language questions by human users. To do so, reading systems should not only extract and interpret the textual (handwritten, typewritten or printed) content of the document images but exploit numerous other visual cues, including layout (page structure, forms, tables), non-textual elements (marks, tick boxes, separators, diagrams) and style (font, colours, highlighting), to mention just a few.

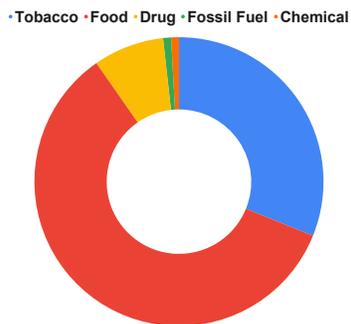
Departing from generic VQA (see section 2.4) and scene text VQA (see subsection 2.5.1) approaches, the document images warrant a different approach to exploit all the above visual cues, making use of prior knowledge of the implicit written communication conventions used, and dealing with the high-density semantic information conveyed in such images. Answers in the case of document VQA cannot be sourced from a closed dictionary, but they are inherently open-ended. Previous approaches to bringing VQA to the domain of document images have either focused on specific document entities such as charts or on specific collections such as book covers (subsection 2.5.2). In contrast to such approaches, we recast the problem to its generic form and put forward a large scale, varied collection of real documents.

The main contributions of this work are the following:

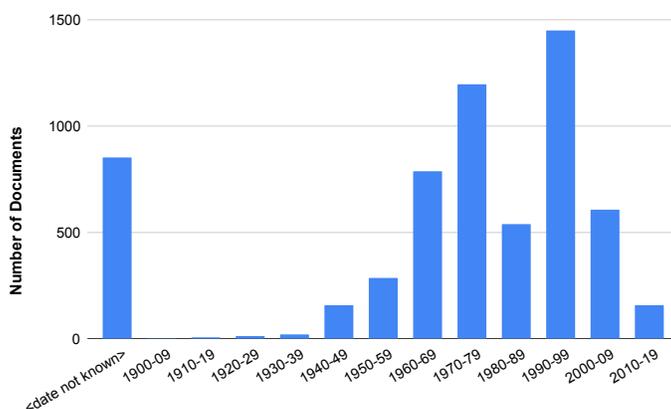
- We introduce DocVQA, a large-scale dataset of 12,767 document images of varied types and content, over which we have defined 50,000 questions and answers. The questions defined are categorised based on their reasoning requirements, allowing us to analyze how DocVQA methods fare for different question types.
- We define and evaluate various baseline methods over the DocVQA dataset, ranging from simple heuristic methods and human performance analysis that allow us to define upper-performance bounds given different assumptions to state-of-the-art scene text VQA models and MRC models.

4.2 DocVQA dataset

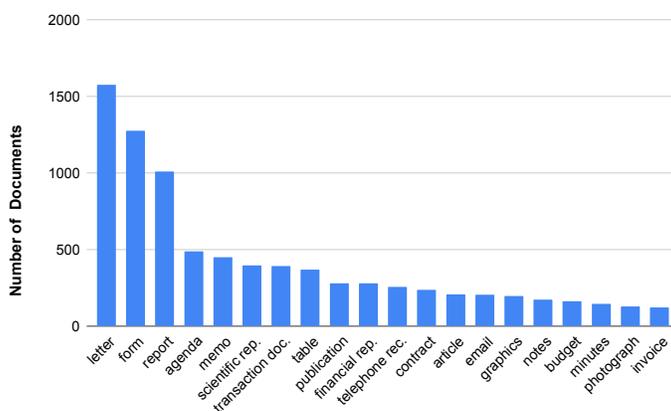
In this section, we explain the data collection and annotation process and present statistics and analysis of DocVQA.



(a) Industry-wise distribution of the documents.



(b) Year wise distribution of the documents.



(c) Various types of documents used.

Figure 4.2: **Distribution of documents by industry, document creation date and document type.** Document images we use in the dataset come from 6071 documents from the UCSF industry documents library, spanning many decades, of a variety of types and originating from 5 different industries.

Stage 1 instructions	Stage 2 instructions
<ol style="list-style-type: none"> 1. You need to add questions and corresponding answers based on the given image. 2. Make sure the questions you ask can be answered purely based on the content in the image 3. Try to frame more questions based on tables, figures and forms. 4. You are allowed to reject the image if, <ol style="list-style-type: none"> (a) the image is of bad resolution or mostly contains illegible text. (b) the content is almost entirely in a non-English language. (c) you are unable to frame any question. 5. The text which forms your answer must be <ol style="list-style-type: none"> (a) found verbatim in the image as a contiguous sequence of tokens in the reading order 	<ol style="list-style-type: none"> 1. You need to enter answers for the questions shown based on the given image 2. if you cannot find the answer to the question based on the image, flag the question as "can't answer" 3. For each question, add appropriate question type(s) from the below list. It is possible that a question can have more than one question type associated with it. <ol style="list-style-type: none"> (a) Handwritten: The question is based on a content that is handwritten (b) Form: The question is based on content that is presented like a form. (c) Layout: If the question is about headings, page numbers or anything concerning the location or layout of text or any other document entity (d) Table/list: The question is grounded on data presented in the form of a table or a list (e) Running text: The question is grounded on the information presented in the form of sentences or paragraphs. (f) Photograph: The question is grounded on a picture of a person or a thing. (g) Figure: The question is grounded on a diagram, visualization, figure or schematic. (h) Yes/No: The answer to the question is yes/no. (i) Other: If none of the above types is valid for the question, assign "Other"

Table 4.1: DocVQA annotation instructions.

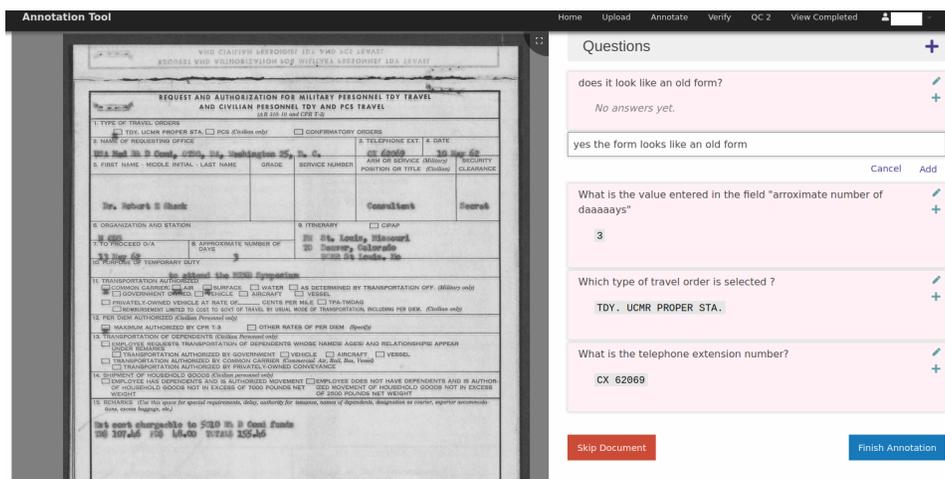


Figure 4.3: Screenshot of annotation stage 1 - question-answer collection. Questions and answers are collected for a given document image. An annotator can add up to 10 questions for an image. The image can be skipped if it is not possible to frame questions on it.

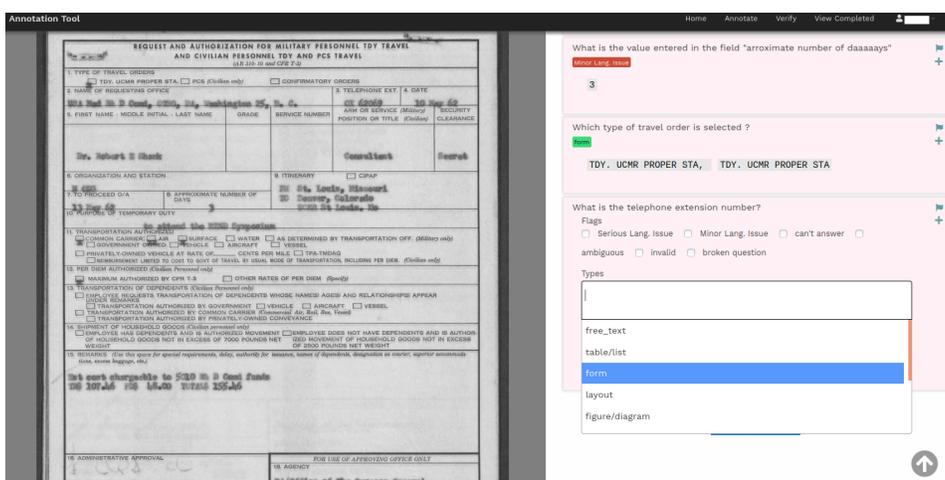


Figure 4.4: Screenshot of annotation stage 2 - data verification. For each question shown annotators have to (i) enter answer(s) (answer(s) from the first stage are not shown) and (ii) Tag the question with one or more question types from the 9 question types shown in a drop-down (question types assigned to a question are shown in green highlight color.) or (iii) flag/ignore the question by selecting the checkbox corresponding to one of the reasons such as “invalid question”, “Serious lang. issue” etc. (the reasons chosen for flagging a question are shown in red highlight color)

4.2.1 Data collection

4.2.1.1 Document images

Images in the dataset are sourced from documents in UCSF Industry Documents Library¹. The library is a digital archive of documents from various industries concerning public health.

¹<https://www.industrydocuments.ucsf.edu/>

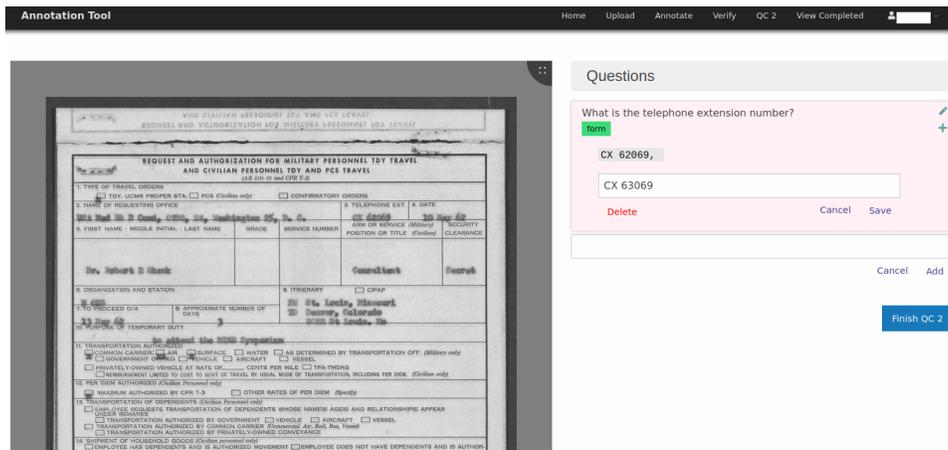


Figure 4.5: Screenshot of annotation stage 3 - reviewing answer mismatch cases. If none of the answers entered in the first stage for a question match with any of the answers entered in the second stage, the question is sent for review in the third stage. In this stage, the worker is allowed to edit question as well as answers or add new answers before accepting the question.

As of June 2020, the library hosts over 15 million documents. The documents are organized under different industries and further under different collections. We downloaded documents from different collections, and hand-picked pages from these documents for use in the dataset. The majority of documents in the library are binarized, which has taken on a toll on the image quality. We tried to minimize binarized images in DocVQA since we did not want poor image quality to be a bottleneck for VQA. We also prioritized pages with tables, forms, lists and figures over pages that only have running text.

The final set of images in the dataset is drawn from pages of 6,071 industry documents. We made use of documents from as early as 1900 to as recent as 2018. (Figure 4.2b). Most of the documents are from the 1960-2000 period, and they include typewritten, printed, handwritten and born-digital text. There are documents from all 5 major industries for which the library hosts documents — tobacco, food, drug, fossil fuel and chemical. We use many documents from food and nutrition-related collections, as they have a good number of non-binarized images. The majority of the documents we use are from food-related collections since a good share of them are not binarized. See Figure 4.2a for industry-wise distribution of the 6071 documents used. The dataset comprises a wide variety of document types, as shown in Figure 4.2c.

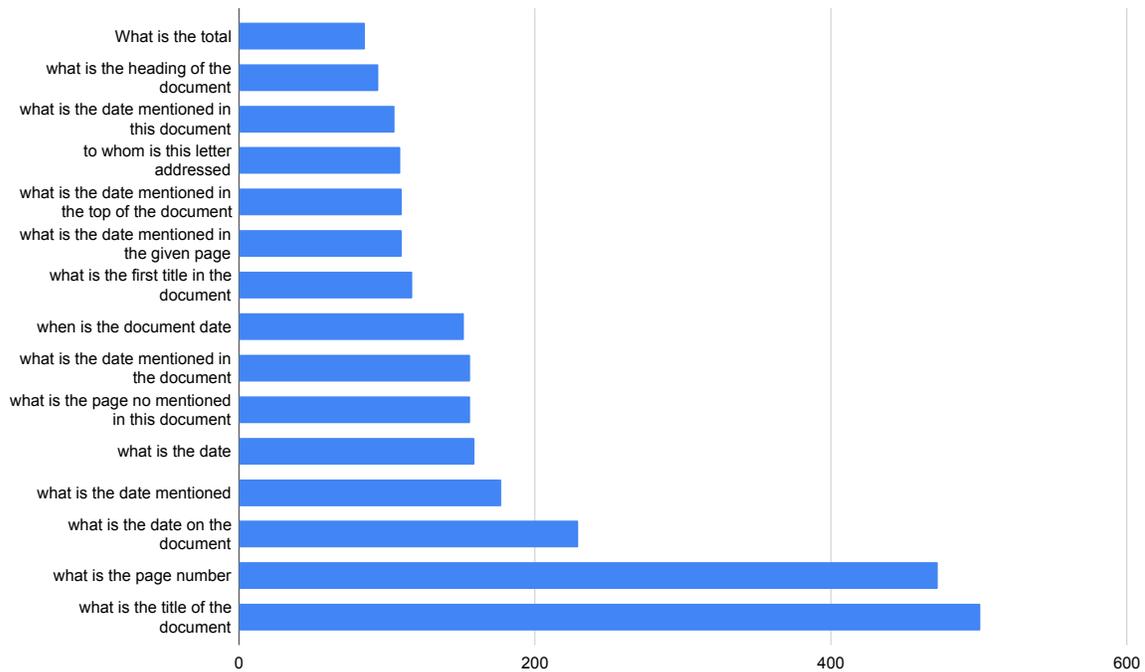


Figure 4.6: Top 15 most frequent questions in DocVQA and their frequencies.

4.2.1.2 Selection of workers

Initially, we hosted a pilot annotation on a crowd-sourcing platform for collecting question-answer pairs. But more than 40% of the question-answer pairs added during the pilot annotation were noisy. We realized that some of the requirements were not easy to understand from written instructions. For example, the kind of question-answers that are allowed—only extractive questions—and assigning question types are better understood when explained using examples. Consequently, we decided to use an internal web-based tool for the annotation and hired workers with whom we could interact closely.

To select the workers, we reached out to individuals looking for annotation-type jobs through mailing lists and other online groups. Interested applicants were invited to join a 90-minute webinar explaining the process and all the requirements. Table 4.1 lists the written annotation instructions that were handed out to the applicants. During the webinar, we explained each of the instructions with many examples of the accepted types of questions. Following the webinar, the applicants were asked to take an online quiz to assess how well they understood the process and the policies. Based on the quiz scores we selected 45 individuals for the annotation. The selected workers were called for another round of webinar, where we discussed the answer key for the quiz and clarified their doubts. The workers were added to an online forum so that

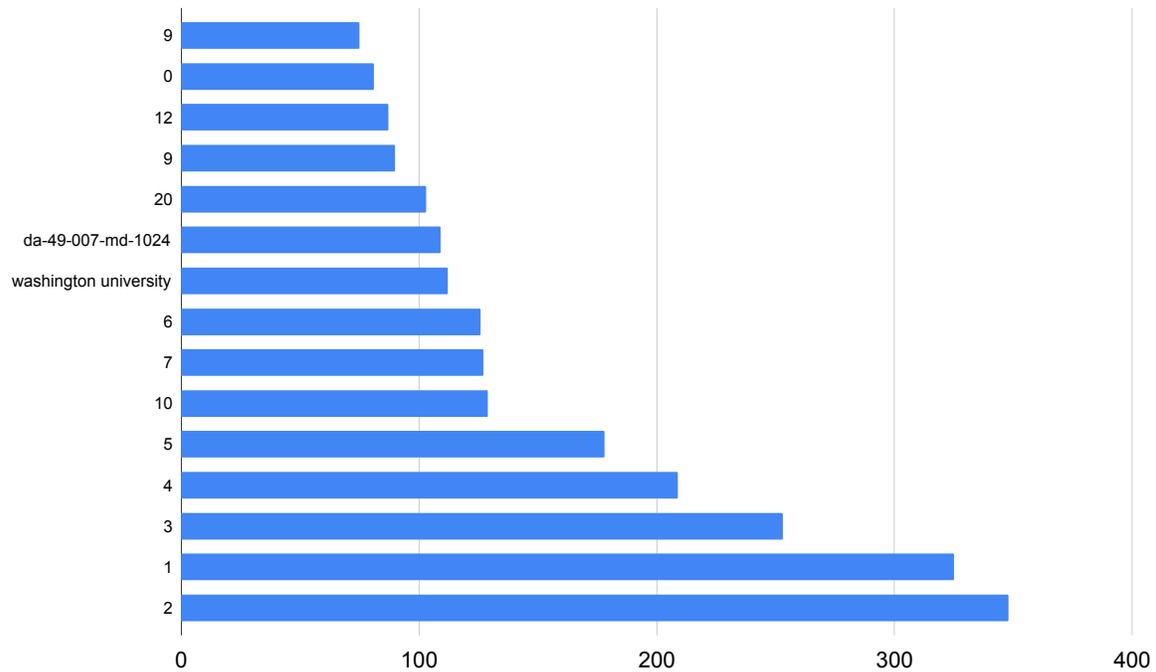


Figure 4.7: Top 15 most frequent answers in DocVQA and their frequencies.

they could post their queries related to the annotation in the forum. They were encouraged to post questions with screenshots whenever in doubt. They would keep the particular image in pending and move on to other images in the queue until one of the authors give a reply to the question they raised. This way, we could reduce annotation errors drastically. Every time a worker is in doubt, they were advised to keep the document in pending, post their question in the forum and move to another image. The image kept in pending will be cleared later when one of the authors replies to the question.

4.2.1.3 Questions and answers

The annotation process was organized into three stages. In stage 1, workers were shown a document image and asked to define at most 10 question-answer pairs on it. They were given the option to skip an image if it is of bad quality or if the worker could not frame any question. We encouraged the workers to add more than one ground truth answer per question when warranted. Figure 4.3 shows a screenshot from stage 2 of the annotation.

The second annotation stage aims to verify the data collected in the first stage. We made sure that the second stage was done by a worker different from the one who collected questions and

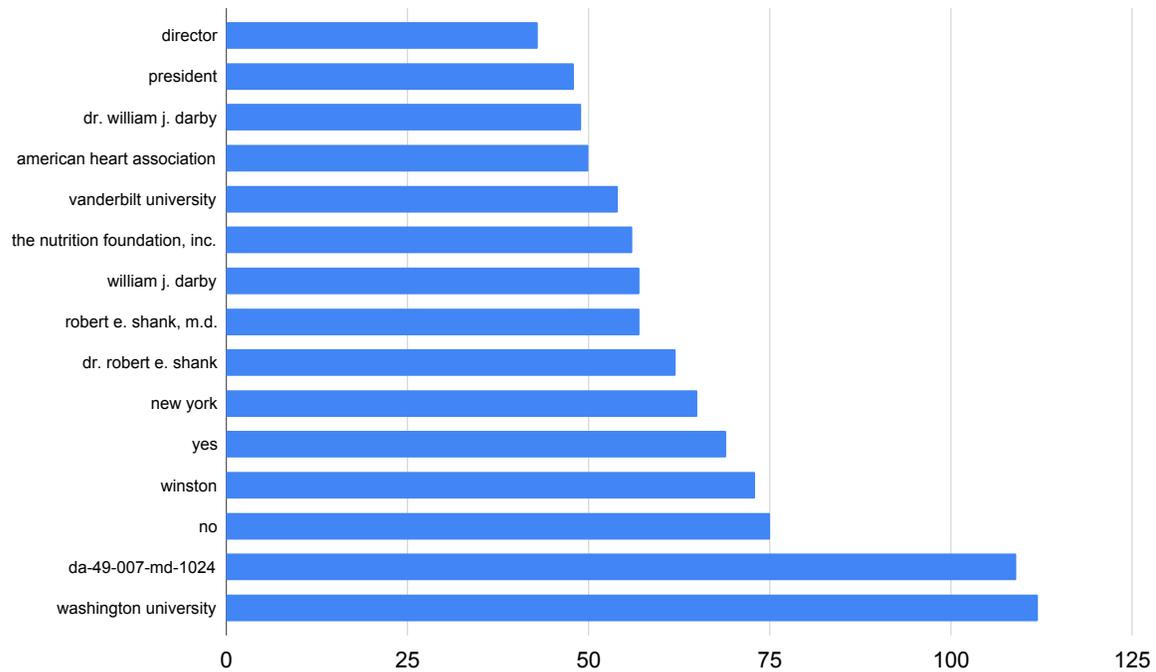


Figure 4.8: Top 15 non numeric answers in DocVQA and their frequencies.

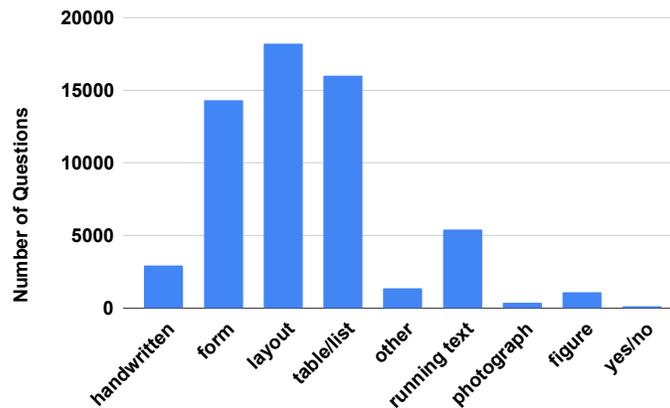
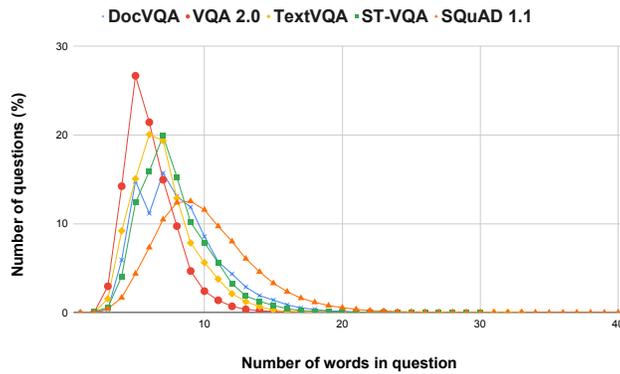
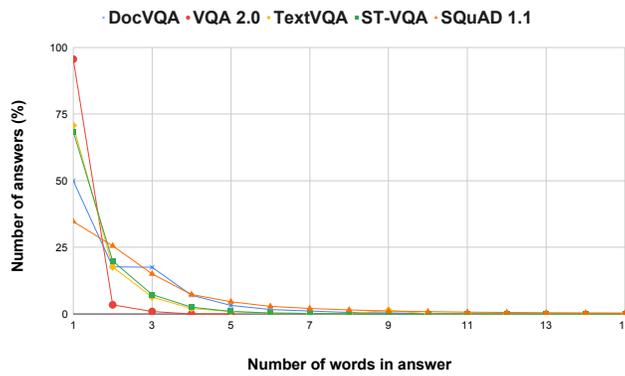


Figure 4.9: **The 9 question types in DocVQA dataset and share of questions in each type.** For each question in the dataset we collect additional meta data called question types that indicate the nature of the question and the kind of information the question is grounded on.

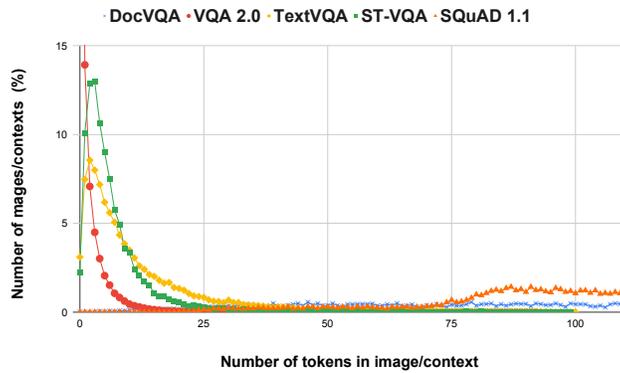
answers on the same image in the first stage. Here a worker was shown an image and questions defined on it in the first stage (but not the answers from the first stage), and was required to enter answers for the questions. In this stage, workers were also required to assign one or



(a) Questions with a particular length.



(b) Answers with a particular length.



(c) Images/contexts with a particular length of text tokens

Figure 4.10: Question, answer and OCR tokens’ statistics compared to similar datasets from VQA—VQA 2.0 [153], ST-VQA [155], TextVQA [149]—and SQuAD 1.1 [91] reading comprehension dataset.

more question types to each question. The different question types in DocVQA are discussed in subsection 4.2.2. During the second stage, if the worker finds a question inapt owing to

language issues or ambiguity, an option to flag the question was provided. Such questions are not included in the dataset. Figure 4.4 shows a screenshot from stage 2 of the annotation.

If none of the answers entered in the first stage match exactly with any of the answers from the second stage, the particular question is sent for review in the third stage. In the third stage, the question and all the answers that were entered during the first and second annotation stages are shown. Questions and answers are editable, and the reviewer either accepts the question-answer (after editing if necessary) or ignores it. The third stage review is done by the authors themselves. Figure 4.4 shows a screenshot of this stage. For a question, we retain all the unique answers (i.e, unique strings after converting all answers to lower case) . Hence a question can have more than 1 valid answer.

4.2.2 Statistics and analysis

The DocVQA comprises 50,000 questions framed on 12,767 images. The data is split randomly in an 80-10-10 ratio to train, validation and test splits. The train split has 39,463 questions and 10,194 images, the validation split has 5,349 questions and 1,286 images, and the test split has 5,188 questions and 1,287 images.

As mentioned before, questions are tagged with question type(s) during the second stage of the annotation process. Figure 4.9 shows the 9 question types and the number of questions in each type. A question type signifies the type of data where the question is grounded. For example, ‘table/list’ is assigned if answering the question requires an understanding of a table or a list. If the information is presented like a form—i.e., in key:value format—the ‘form’ type is assigned. ‘Layout’ is assigned for questions that require spatial/layout information to find the answer. For example, questions asking for a title or heading require one to understand the structure of the document. If a question is based on a picture of a person or a thing, “photograph” type is assigned. Similarly, a question that is based on a figure/plot/visualization is tagged as of type “figure”. Any question that has a yes or no answer is assigned “yes/no” type. If the answer for a question is based on information in the form of sentences/paragraphs, the type assigned is ‘running text’. For all questions where the answer is based on handwritten text, the ‘handwritten’ type is assigned. If none of the above-mentioned question types is valid for a question, the workers were asked to assign the type for the question as ”other”. It must be noted that a question can have more than one question type assigned to it. For example, if a question is based on a handwritten form, both “form” and “handwritten” types are assigned. In some cases, multiple question types are assigned when it is difficult to decide what document entity the question is based on. For example, there are many instances where a form looks very

similar to tabular data. In such cases, we advised the workers to assign both “table/list” and “form” types.

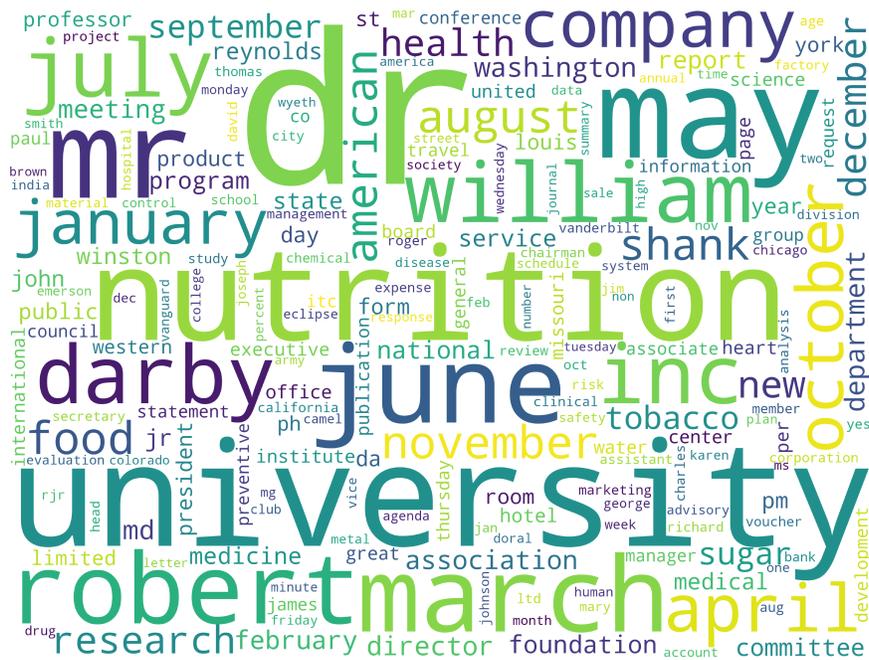
In the following analysis, we compare statistics of questions, answers, and OCR tokens with other similar datasets for VQA —VQA 2.0 [153], TextVQA [149] and ST-VQA [155] and SQuAD 1.1 [91] dataset for reading comprehension. Statistics for other datasets are computed based on their publicly available data splits. For statistics on OCR tokens, for DocVQA, we use OCR tokens generated using Microsoft’s Read OCR². For VQA 2.0, TextVQA and ST-VQA, we use OCR tokens made available by authors of LoRRA [149] and M4C [94] as part of the MMF framework [214]. For statistics involving question or answer tokens, tokens are generated by splitting question or answer at white space.

Figure 4.10a shows the distribution of question lengths for questions in DocVQA compared to other similar datasets. The average question length is 8.12, which is the second highest among the compared datasets. The minimum question length is 1, and the longest question has 29 words. In DocVQA, 35,362 (70.72%) questions are unique, which is lower compared to other datasets except VQA 2.0. Figure 4.6 shows the top 15 most frequent questions and their frequencies. There are questions repeatedly being asked about dates, titles, and page numbers. For example “What is the title of the document” is asked at least 500 times. A sunburst of the first 4 words of the questions is shown in Figure 4.12. Here, for better visualization, we only show words that occur at least 80 times. It can be seen that a large majority of questions start with “what is the”, asking for the date, title, total, amount, or name.

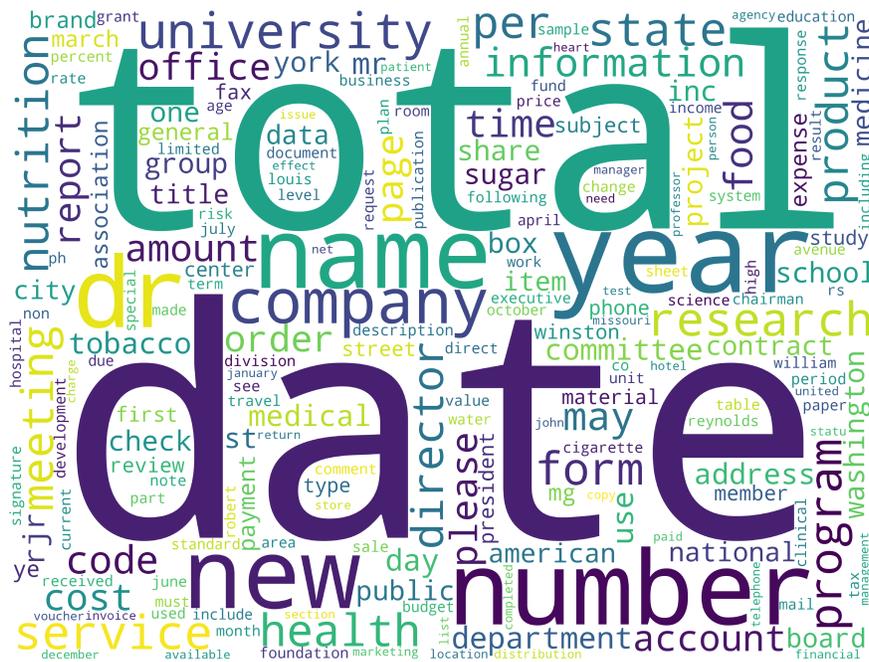
The distribution of answer lengths is shown in Figure 4.10b. We observe in the figure that both DocVQA and SQuAD 1.1 have a higher number of longer answers compared to the VQA datasets. The average answer length is 2.17. The longest answer in DocVQA has 30 words in it. 63.2% of the answers are unique, which is second only to SQuAD 1.1 (72.5%). The top 15 answers in the dataset are shown in Figure 4.7. The most common answer is an answer for only 0.7% of questions in the dataset. We observe that almost all of the top answers are numeric values, which is expected since there are a good number of document images of reports and invoices. In Figure 4.8, we show the top 15 non-numeric answers. These include named entities such as names of people, institutions, and places. The word cloud in Figure 4.11a shows frequent words in answers. The most common words are the names of people and the names of calendar months.

In Figure 4.10c, we show the number of images (or ‘context’s as called in the case of SQuAD 1.1) containing a particular number of text tokens. The average number of text tokens in an image or context is the highest in the case of DocVQA (182.75). It is considerably higher

²<https://learn.microsoft.com/en-us/azure/cognitive-services/computer-vision/overview-ocr>



(a) Word cloud of tokens in answers in DocVQA



(b) Word cloud of OCR tokens in DocVQA

Figure 4.11: Word clouds of answers and OCR tokens in DocVQA dataset

compared to SQuAD 1.1, where contexts are usually small paragraphs whose average length is

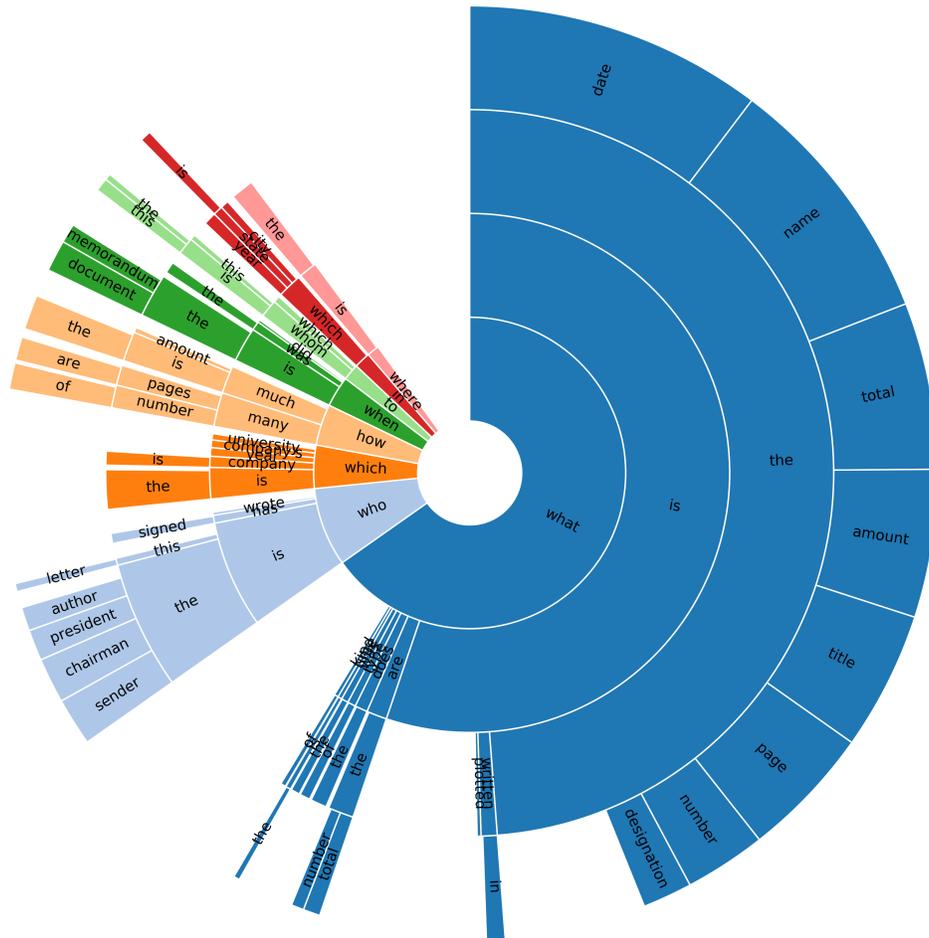


Figure 4.12: **Distribution of questions by their starting 4-grams.** Most questions aim to retrieve common data points in documents such as date, title, total mount, and page number.

117.23. In the case of VQA datasets that comprise real-world images average number of OCR tokens is not more than 13. The word cloud in Figure 4.11b shows the most common words spotted by the OCR on the images in DocVQA. There is a high overlap between common OCR tokens and words in answers.

4.3 Baselines

In this section, we explain the baselines we use, including heuristics and trained models.

4.3.1 Heuristics and upper bounds

The heuristics we evaluate are: (i) **Random answer:** measures performance when we pick a random answer from the answers in the train split. (ii) **Random OCR token:** performance when a random OCR token from the given document image is picked as the answer. (iii) **Longest OCR token** is the case when the longest OCR token in the given document is selected as the answer. (iv) **Majority answer** measures the performance when the most frequent answer in the train split is considered as the answer.

We also compute the following upper bounds: (i) **Vocab UB:** This upper bound measures the performance upper bound one can get by predicting correct answers for the questions, provided the correct answer is present in a vocabulary of answers, comprising all answers which occur more than once in the train split. (ii) **OCR substring UB:** is the upper bound on predicting the correct answer provided the answer can be found as a substring in the sequence of OCR tokens. The sequence is made by serializing the OCR tokens recognized in the documents as a sequence separated by space, in top-left to bottom-right order. (iii) **OCR subsequence UB:** the upper bound on predicting the correct answer provided the answer is a subsequence of the OCR tokens' sequence.

4.3.2 VQA models

In subsection 2.5.1. we discussed two scene text VQA models—LoRRA and M4C—that have the capability to read and reason over text spotted on the images. We evaluate the models on DocVQA. In our experiments, we use original LoRRA and M4C models and a few variants of these models. Document images in DocVQA usually contain a higher number of text tokens compared to images in scene text VQA datasets. Hence we try out larger dynamic vocabularies (i.e., more OCR tokens are considered from the images) for both LoRRA and M4C. For both models, we also evaluate performance when no fixed vocabulary is used. Since the notion of visual objects in real world images is not directly applicable in the case of document images, we also try out variants of LoRRA and M4C where region features are omitted.

4.3.3 Reading comprehension models

In addition to the VQA models, which can read text, we try out extractive QA/MRC models from NLP. In particular, we use the BERT [113] QA model(subsection 2.2.1.1).

4.4 Experiments

In this section, we explain evaluation metrics and our experimental settings and report the results of experiments.

4.4.1 Evaluation metrics

Two evaluation metrics we use are Average Normalized Levenshtein Similarity (ANLS) and Accuracy (Acc.). ANLS was originally proposed for the evaluation of VQA on ST-VQA [215]. The Accuracy metric measures the percentage of questions for which the predicted answer matches exactly with any of the target answers. The accuracy metric awards a zero score even when the prediction differs slightly from the target answer. Since no OCR is perfect, we propose to use ANLS as our primary evaluation metric so that minor answer mismatches stemming from OCR errors are not severely penalized.

The ANLS score is defined as:

$$\text{ANLS} = \frac{1}{N} \sum_{i=0}^N \left(\max_j s(a_{ij}, o_{q_i}) \right) \quad (4.1)$$

$$s(a_{ij}, o_{q_i}) = \begin{cases} (1 - NL(a_{ij}, o_{q_i})) & \text{if } NL(a_{ij}, o_{q_i}) < \tau \\ 0 & \text{if } NL(a_{ij}, o_{q_i}) \geq \tau \end{cases}$$

In Equation 4.1 N is the number of questions and M is the number of ground truth answers per question. The ground truth answers are denoted as a_{ij} . Index i in the range $\{0, \dots, N\}$ corresponds to the question, and index j indicates the index of the particular ground answer for a question. o_{q_i} is the predicted answer for the i^{th} question q_i . $NL(a_{ij}, o_{q_i})$ is the Normalized Levenshtein distance between the strings a_{ij} and o_{q_i} after converting both the strings to lower case. The distance is a value between 0 and 1. The similarity score $s(a_{ij}, o_{q_i})$ for a_{ij} and o_{q_i} is set to $1 - NL(a_{ij}, o_{q_i})$ if the normalized distance is less than a threshold $\tau = 0.5$. Otherwise, the similarity score is set to zero. The rationale for the threshold is that it is unlikely that the mismatch is caused by OCR error if a predicted answer has a normalised edit distance of more than 0.5 to a ground truth answer.

4.4.2 Experimental setup

For measuring human performance, we collect answers for all questions in the test split, with the help of a few volunteers from our institution. In all our experiments, including heuristics

Baseline	val		test	
	ANLS	Acc.	ANLS	Acc.
Human	-	-	0.981	94.36
Random answer	0.003	0.00	0.003	0.00
Random OCR token	0.013	0.52	0.014	0.58
Longest OCR token	0.002	0.05	0.003	0.07
Majority answer	0.017	0.90	0.017	0.89
Vocab UB	-	31.31	-	33.78
OCR substring UB	-	85.64	-	87.00
OCR subsequence UB	-	76.37	-	77.00

Table 4.2: **Results of different heuristic baselines and upper bounds.** Predicting random answers or majority answer do not even yield 1% accuracy. Answers are a substring of the serialized OCR output in more than 85% of the cases.

and trained baselines, the OCR tokens we use are extracted using a commercial OCR application. For the heuristics and upper bounds, we use a vocabulary of 4,341 answers which occur more than once in the train split.

For LoRRA and M4C, we use official implementations available as part of the MMF framework. The training settings and hyperparameters are the same as the ones reported in the original works. Grid features and region features for LoRRA and region features for M4C are extracted using the models originally used in the respective works. The fixed vocabulary we use for LoRRA is the same as the vocabulary we use for computing vocabulary-based heuristics and upper bounds. For M4C, the fixed vocabulary we use is a vocabulary of the 5,000 most frequent words from the answers in the train split of DocVQA.

For QA using BERT, three pretrained BERT models³ from the Hugging face library [216] are used. The models we use are bert-base-uncased, bert-large-uncased-whole-word-masking, and bert-large-uncased-whole-word-masking-finetuned-squad. We abbreviate the model names as bert-base, bert-large, and bert-large-squad, respectively. Among these, bert-large-squad is a pretrained model, which is also finetuned on SQuAD 1.1 for question answering. In the case of extractive QA or MRC datasets, ‘contexts’ on which questions are asked are passages of electronic text. But in DocVQA, ‘contexts’ are document images. Hence to finetune the BERT

³https://huggingface.co/transformers/pretrained_models.html

Method	Region features	Fixed vocab.	Dynamic vocab. size	val		test	
				ANLS	Acc.	ANLS	Acc.
LoRRA [149]	✓	✓	50	0.110	7.22	0.112	7.63
	✓	✗	50	0.041	2.64	0.037	2.58
	✗	✓	50	0.102	6.73	0.100	6.43
	✓	✓	150	0.101	7.09	0.102	7.22
	✓	✓	500	0.094	6.41	0.095	6.31
M4C [94]	✓	✓	50	0.292	18.34	0.306	18.75
	✓	✗	50	0.216	12.44	0.219	12.15
	✗	✓	50	0.294	18.75	0.310	18.92
	✗	✓	150	0.352	22.66	0.360	22.35
	✗	✓	300	0.367	23.99	0.375	23.90
	✗	✓	500	0.385	24.73	0.391	24.81

Table 4.3: **Performance of the VQA models which are capable of reading text— LoRRA [149] and M4C[94] on DocVQA.** Detection of visual objects and their features (bottom-up attention), which is a common practice in VQA is ineffective in the case of DocVQA.

QA models on DocVQA, we need to prepare the data in SQuAD style format where the answer to a question is a ‘span’ of the context, defined by start and end indices of the answer. To this end, we first serialize the OCR tokens recognized on the document images to a single string, separated by space, in top-left to bottom-right order. To approximate the answer spans, we follow an approach proposed in TriviaQA [217], which is to find the first match of the answer string in the serialized OCR string.

The bert-base model is finetuned on DocVQA on 2 Nvidia GeForce 1080 Ti GPUs, for 2 epochs, with a batch size of 32. We use Adam optimizer [218] with a learning rate of $5e - 05$. The bert-large and bert-large-squad models are finetuned on 4 GPUs for 6 epochs with a batch size of 8, and a learning rate of $2e - 05$.

4.4.3 Results

The results of the heuristic baselines and upper bounds are reported in Table 4.2. We can see that none of the heuristics get even a 1% accuracy on the validation or test splits. *OCR*

substring UB yields more than 85% accuracy on both validation and test splits. However, the substring match in all cases need not be an actual answer match. For example, if the answer is “2” which is the most common answer in the dataset, it will match with a “2” in “2020” or a “2” in “2pac”. This is the reason why we evaluate the *OCR subsequence UB*. An answer is a sub sequence of the serialized OCR output for around 76% of the questions in both validation and test splits.

The results of our trained VQA baselines are shown in Table 4.3. The first rows for both methods report the results of the original model proposed in those works. In the case of LoRRA, the original setting proposed by the authors yields the best results compared to the variants we try out. With no fixed vocabulary, the performance of the model drops sharply, suggesting that the model primarily relies on the fixed vocabulary to output answers. Larger dynamic vocabulary results in a slight performance drop, suggesting that incorporating more OCR tokens from the document images does little help. Unlike LoRRA, M4C benefits from a larger dynamic vocabulary. Increasing the size of the dynamic vocabulary from 50 to 500 improves the ANLS by around 50%. In the case of LoRRA, although region features help, the improvement is insignificant. For M4C, the setting where region features are omitted performs slightly better compared to the original setting. These results can be attributed to the fact that the region features used by LoRRA and M4C come from an object detection model trained on natural images. The document images in DocVQA are significantly different from natural images, and the definition of objects in natural images makes little sense in the context of document images.

The results of different variants of the BERT QA model are reported in Table 4.4. We observe that all BERT models perform better than the best VQA baseline using M4C (last row in 4.3). The best-performing model out of all the baselines analysed is the bert-large-

Pretrained model	Finetuned on DocVQA	val		test	
		ANLS	Acc.	ANLS	Acc.
bert-base	✓	0.556	45.6	0.574	47.6
bert-large	✓	0.594	49.28	0.610	51.08
bert-large-squad	✗	0.462	36.72	0.475	38.26
bert-large-squad	✓	0.655	54.48	0.665	55.77

Table 4.4: **Performance of BERT QA model on DocVQA.** A BERT_{LARGE} model which is finetuned on both SQuAD 1.1 [91] and DocVQA performs the best.

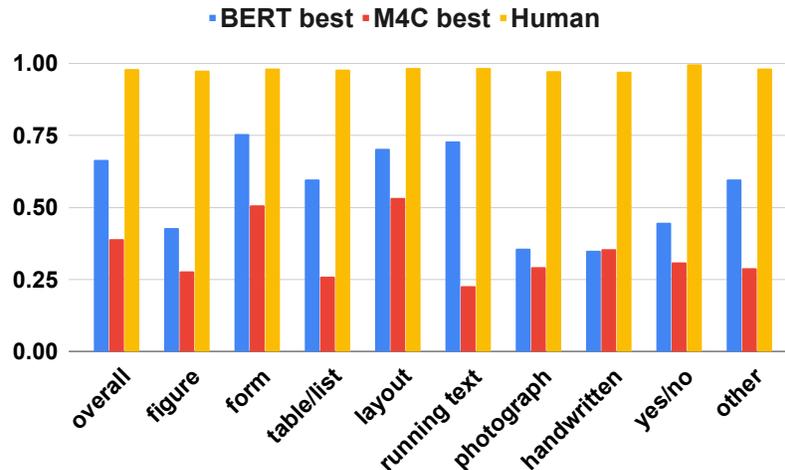


Figure 4.14: **Performance of baseline models for different question types.** Best baselines from VQA space and reading comprehension space pitted against the human performance for different question types. We need models which can understand figures and text on photographs better. We need better handwriting recognizers too!

4.5 Summary

We introduce a new data set and an associated VQA task with the aim to inspire a “purpose-driven” approach in document image analysis and recognition research. Our dataset comprises a variety of business documents that feature structured and unstructured text, tables, figures and photographs. However existing text-based VQA models perform poorly on DocVQA as they are developed exclusively for natural images. We believe DocVQA could drive methods that use low-level cues (text, layout, arrangements) and high-level goals (purpose, relationship, domain knowledge) in solving problems of practical importance.

Chapter 5

InfographicVQA: VQA on infographics

In the previous chapter, we presented a novel VQA task called DocVQA, where images comprise pages from business documents. Our experiments on DocVQA show that a text-only baseline using BERT correctly answers nearly 55% of the questions. Additionally, results of the DocVQA challenge [219] suggest that a good share of the questions in the DocVQA can be answered without using any visual information. The findings drove us to create a DocVQA++, where visual cues are indispensable for answering the questions. To this end, we present InfographicVQA, a new dataset that comprises a diverse collection of infographics along with question-answer annotations. Infographics communicate information using a combination of textual, graphical and visual elements. We believe infographics that feature complex multimodal interactions and rich in information, are an excellent benchmark for multimodal reasoning. We curate the dataset with an emphasis on questions that require elementary reasoning and arithmetic skills. We propose a multimodal BERT-like architecture for VQA on InfographicVQA. Although our approach performs better than the baselines, there is a significant performance gap compared to near-perfect human performance on the dataset. This work on VQA on infographics has previously been presented in our publication [9].

5.1 Introduction

Infographics are documents created to convey information in a compact manner using a combination of textual and visual cues. The presence of the text, numbers, and symbols, along with the semantics that arises from their relative placements, make infographics understanding a challenging problem. Machine understanding of infographics requires methods to jointly reason over the document layout, text, graphical elements, data visualisations, color schemes and visual art, among others. Motivated by the multimodal nature of infographics, and the human

centered design, we propose a Visual Question Answering (VQA) approach to infographics understanding.



How many companies have more than 10K delivery workers?

Answer: 2

Evidence: Figure

Answer-source: Non-extractive Operation: Counting Sorting

Who has better coverage in Toronto - Canada post or Amazon?

Answer: canada post

Evidence: Text

Answer-source: Question-span Image-span Operation: none

In which cities did Canada Post get maximum media coverage?

Answer: vancouver, montreal

Evidence: Text Map

Answer-source: Multi-span Operation: none

Figure 5.1: **Example image from InfographicVQA along with questions and answers.** For each question, source of the answer, type of evidence the answer is grounded on, and the discrete operations required to find the answer are shown.

In this work, we introduce a dataset for VQA on infographics named InfographicVQA. The dataset comprises 30,035 questions over 5,485 images. An example from our dataset is shown in Figure 5.1. Questions in the dataset include questions grounded on tables, figures and visualizations and questions that require combining multiple cues. We believe our dataset is ideal for benchmarking the progress of algorithms at the meeting point of vision, language and document understanding.

We propose a multimodal BERT-like model called VLL-BERT (vision, language and layout BERT) for VQA on the new dataset. LayoutLM (subsection 2.3.1) and VisualBERT(subsection 2.4.2 serve as inspiration for our model. Despite the fact that our method outperforms a

strong baseline like M4C, there remains a sizable performance difference when compared to near-perfect human performance on the dataset.

5.2 Related works

Dataset	Images	Synthetic Images	Template questions	Text type	# Images	# Questions	Answer type
TQA [163]	Science diagrams	✗	✗	MR	1K	26K	MCQ
RecipeQA [164]	Culinary pictures	✗	✓	MR	251K	37K	MCQ
ST-VQA [155]	Natural images	✗	✗	ST	23K	31K	Ex
TextVQA [149]	Natural images	✗	✗	ST	28K	45K	Ex, SAb
OCR-VQA [60]	Book covers	✗	✓	BD	207K	1M	Ex, Y/N
DVQA [58]	Bar charts	✓	✓	BD	300K	3.4M	Ex, Nm, Y/N
FigureQA [59]	Charts - 5 types	✓	✓	BD	120K	1.5M	Y/N
LEAF-QA [220]	Charts - 4 types	✓	✓	BD	250K	2M	Ex, Nm, Y/N
VisualMRC [221]	Webpage screenshots	✗	✗	BD	10K	30K	Ab
DocVQA [5]	Industry documents	✗	✗	Pr, Tw, Hw, BD	12K	50K	Ex
InfographicVQA	Infographics	✗	✗	BD	5.4K	30K	Ex, Nm

Table 5.1: **Summary of VQA and Multimodal QA datasets where text on the images needs to be read to answer questions.** Text types are: Machine Readable (MR), Scene Text (ST), Born Digital (BD), Printed (Pr), Handwritten (Hw), and Typewritten (Tw). Answer types are: Multiple Choice (MCQ), Extractive (Ex), Short abstractive (SAb), Abstractive (Ab), Yes/No (Y/N), and Numerical (Nm)—when answer is numerical and not extracted from image or question.

5.2.1 Question answering in a multimodal context.

Textbook Question Answering (TQA) and RecipeQA deal with QA in a multimodal context (see section 2.6). Contrary to InfographicVQA and other datasets mentioned below, text in these two datasets is not embedded on the images but provided separately in machine readable form. ST-VQA and TextVQA (see subsection 2.5.1) comprise images captured in the wild with sparse text content. InfographicVQA has born-digital images with an order of magnitude more text tokens per image, richer in layout and in the interplay between textual and visual elements. Compared to the OCR-VQA dataset (subsection 2.5.2.1 which focuses exclusively on book covers and contains text-centric template questions, questions in InfographicVQA require multimodal reasoning and the dataset contains visually rich images. Charts in chart VQA

tasks (subsection 2.5.2.2) are a type of infographics. But these datasets contain a few types of standard charts and contain only template questions. Compared to these, InfographicVQA contains complex infographics that go beyond the definition of a "chart" and questions are manually annotated.

VisualMRC [221] is a document VQA dataset that was introduced concurrently to InfographicVQA. VisualMRC is an abstractive VQA (answers cannot be directly extracted from text in the images or questions) benchmark where images are screenshots of web pages. Compared to VisualMRC, InfographicVQA is an extractive VQA (answers are extracted as 'span'(s) of the question or text present in the given image), except for questions that require certain discrete operations resulting in numerical non-extractive answers. InfographicVQA can be seen as a natural extension of DocVQA to a new category of images that are information rich and more challenging in terms of layout complexity and vision-language interplay. Table 5.1 presents a high-level summary of the QA/VQA datasets related to InfographicVQA.

5.2.2 Multimodal pretraining for vision and language tasks

Following the success of BERT [113], there have been multiple works extending it to the Vision-Language space. Models like VisualBERT(subsection 2.4.2), VL-BERT [151], ViLBERT [222], LXMERT [223] and UNITER [152] show that pretraining of BERT-like architectures on vision and language inputs achieve SoTA performances on downstream tasks such as VQA on natural images. Similarly, BERT inspired similar models for DIA tasks where 2D position embeddings of OCR tokens are used in addition to the embeddings of the text tokens. LAMBERT [141] and LayoutLM [115] are two prior works that fall into this category. LayoutLM is a basis for the model we propose for InfographicVQA. We discuss details of LayoutLM architecture in subsection 2.3.1.

5.2.3 Infographics understanding.

Bylinskii et al. [224] and Madan et al. [225] looked at generating textual and visual tags from infographics. Landman uses an existing text summarization model to generate captions for infographics [226]. But the model uses only text recognized from infographics to generate the captions, and layout/visual information is not considered. These three works use Visually29K dataset that comprises images from a single infographics website. MASSVIS [227] is a collection of infographics created to study infographics from a cognitive perspective. As observed by Lu et al. [228], it is a specialized collection focusing on illustrations of scientific procedures and statistical charts, therefore not representative of general infographics. To

summarize, existing datasets containing infographics are either specialized collections or infographics collected from a single source. In contrast, the InfographicVQA dataset comprises infographics drawn from hundreds of different sources, with diverse layouts and designs, and without any topic specialization.

5.3 InfographicVQA

Details of the data collection and analysis of the images, questions, and answers are presented here.

5.3.1 Collecting infographics

Images in the dataset are sourced from the internet. Initially, we downloaded more than 10K images for the search query “infographics” using Google and Bing image search engines. The downloaded images were first de-duped using a Perceptual Hashing approach implemented in *imagededup* library ¹. This step helped us to remove nearly 2000 duplicates. In the second round of de-duplication, we compared the images using the Jaccard similarity of the text tokens spotted on the images. For recognizing text tokens on the images, we used Amazon Textract OCR ². After the two rounds of de-duplication, around 7K images were left. These were added to the annotation system for question-answer annotation.

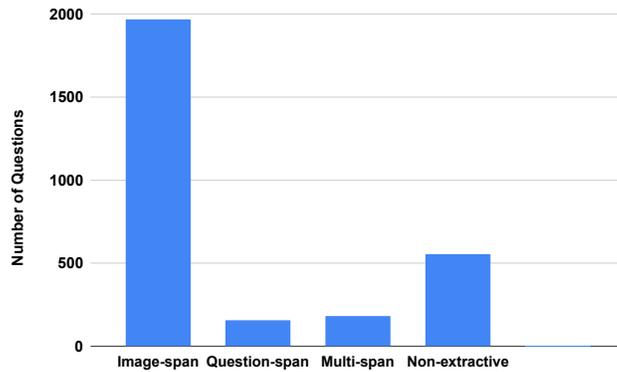
5.3.2 Collecting questions and answers

The annotation tool used for InfographicVQA is the same as the one used for DocVQA annotation (see subsection 4.2.1). Minor changes were incorporated into the tool to support the collection of new question-answer types that we discuss in subsection 5.3.3. Workers who had been part of the DocVQA annotation were invited for InfographicVQA annotation. We conducted an online webinar to explain the process to the prospective workers using examples from pilot annotation that the authors themselves did. Following the webinar, a quiz was conducted to ensure that the workers understood the annotation guidelines correctly. Similar to DocVQA annotation (subsection 4.2.1), the annotation of InfographicVQA involved three stages. The process followed in the three stages is same as how we annotated the DocVQA

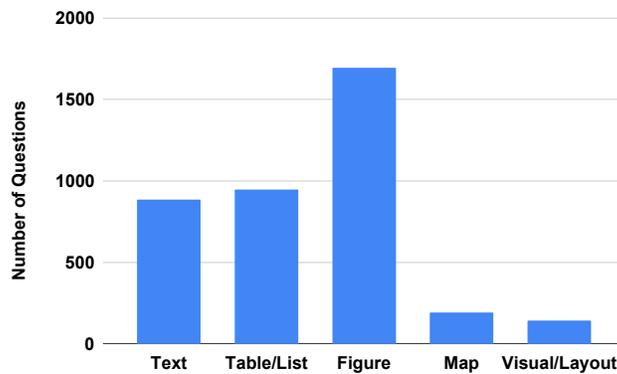
¹<https://github.com/idealo/imagededup>

²<https://aws.amazon.com/textract>

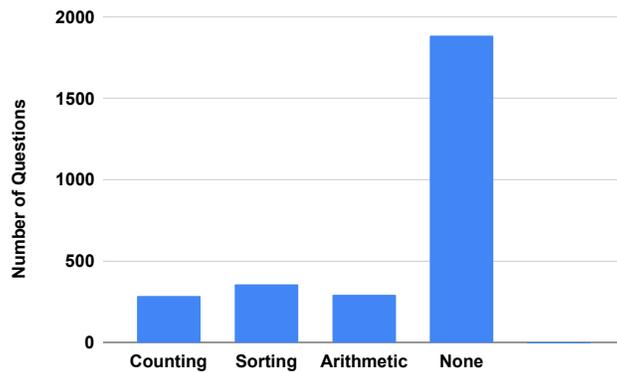
dataset (subsection 4.2.1.3). The only major difference compared to DocVQA annotation is the question-answer types—similar to question types in DocVQA—that were collected.



(a) Answer-sources and their counts



(b) Evidence types and their counts



(c) Operation types and their counts

Figure 5.2: Count of questions in validation set by their Answer-source, (Figure 5.2a), Evidence required to answer (Figure 5.2b) and the discrete Operation performed to find the answer (Figure 5.2c).

5.3.3 Question-answer types: answer-source, evidence and operation

In the second stage, in addition to answering questions collected in the first stage, we instructed the workers to add question-answer types (QA types). QA types are a set of category labels assigned to each question-answer pair. DocVQA and VisualMRC have QA types that indicate the kind of document entity a question is based on. In InfographicVQA, we collect QA types under three categories: Answer-source, Evidence, and Operation.

There are four types of Answer-source—Image-span, Question-span, Multi-span, and Non-extractive. Akin to the definition of ‘span’ in SQuAD and DocVQA, an answer is considered Image-span if it is a single span (a contiguous sequence) of text extracted verbatim in the reading order from text present in the image. Similarly, when the answer is a span from the question, it is labelled as Question-span. In Figure 5.1, the answer to the second question is a span of text found both in the image and the question. Consequently, both Image-span and Question-span are valid answer sources for the question. A Multi-span answer is composed of multiple spans of text from the image. We instructed our workers to enter Multi-span answers by separating each individual span by a comma and a white space. For example, in Figure 5.1, for the last question, the answer is the names of two cities that do not appear in a contiguous sequence of text. Hence it is a Multi-span answer. For a Multi-span answer, any order of the individual spans is a valid answer. In the above example, both “Vancouver, Montreal” and “Montreal, Vancouver” are valid answers. Since such answers are unordered lists, we consider all permutations of the list as valid answers for the question at evaluation time. The ‘Non-extractive’ type is assigned when the answer is not an extracted one. We instructed the workers not to add Non-extractive questions unless the answer is a numerical value.

The inclusion of Question-span, Multi-span, and numerical Non-extractive answers in InfographicVQA is inspired by QAs collected in the DROP dataset [82]. We see this as a natural next step in VQA involving text, different from the purely extractive QA setting in datasets like DocVQA and ST-VQA and abstractive question answering in VisualMRC where automated evaluation is difficult. Allowing non-extractive answers only when the answer is a numerical value makes sure that such answers are short and unique, giving no room for variability. This ensures that automatic evaluation metrics can be used effectively even though answers are generated, not extracted. Near perfect human performance while using automatic evaluation metrics (Table 5.4) validates that answers in InfographicVQA are unique with minimal variability when answered by different individuals.

The Evidence type indicates the kind of evidence behind the answer. Types of evidence are Text, Figure, Table/List, Map, and Visual/Layout. For example, Map is used if the question is based on data shown on a geographical map. Visual/Layout type is added when evidence

is based on the visual or layout aspect of the image. For example, questions such as “What is the color of the hat – brown or black?” or “What is written at the top left corner” fall in this category. Sometimes it is difficult to discern evidence for a question-answer pair. For example, for the first question in Figure 5.1, although the evidence type assigned by the worker is ‘Figure’, it could even be ‘Table/List’ since the visualization looks like a table.

The operation type captures the kind of discrete operation(s) required to arrive at an answer — Counting, Arithmetic or Sorting. Figure 5.2 shows the distribution of questions in the validation split based on Answer-source, Evidence, and Operation. As evident from Figure 5.1, a question can have multiple types of answer sources, evidence, or operation and many questions do not require any of the specified discrete operations to find the answer. For these reasons, counts in plots shown in Figure 5.2 do not add up to 100%.

5.3.4 Statistics and analysis of the dataset

The InfographicVQA dataset has 30,035 questions and 5,485 images in total. These images are from 2,594 distinct websites. The data is split randomly into 23,946 questions and 4,406 images in train, 2,801 questions and 500 images in validation, and 3,288 questions and 579 images in test splits.

To analyze the topics covered by images in InfographicVQA, we used the Latent Dirichlet Allocation (LDA) [230] model. LDA implementation in Gensim library [231] is used. Since our dataset comprises images, the text recognized from the images using the Textract OCR is used for the topic modelling. Table 5.2 shows that images in InfographicVQA cover a wide range of topics such as energy, war, health, and social media. In Figure 5.4, we show a visualization of the top 20 topics in the dataset, visualized using the pyLDAvis tool ³.

We report basic statistics of questions, answers, and OCR tokens in InfographicVQA and similar datasets in Table 5.3. Table 5.3 shows that InfographicVQA has the highest percentage of unique questions and the highest average question length compared to similar datasets. Figure 5.5 shows a sunburst of the common questions in the dataset. There are a good number of questions asking for “How many...” or percentages. This is expected since infographics carry a lot of numerical data.

Table 5.3 shows that the percentage of questions having unique answers in InfographicVQA is the lowest (48.84%) among the compared datasets. InfographicVQA has the shortest average answer length (1.60) compared to other document VQA datasets—VisualMRC (9.55) and DocVQA (2.49)—and slightly longer than scene text VQA datasets such as ST-VQA (1.56) and TextVQA (1.51). It can be seen in Figure 5.3a that the most common answers are num-

³<https://github.com/bmabey/pyLDAvis>

No.	Topic
1	cost lead increase system non risk energy reduce cause clean
2	war violence symptom domestic potential die injury mil acquire birth
3	health person white black police department doctor respiratory smith officer
4	child food water parent potential eat drink essential green sugar already
5	death woman age man old adult love likely statistic rate
6	country high account say month report change global survey event
7	social medium job value program find direct authority salary candidate
8	first purchase call sport still house kid name bring early
9	case university point physical idea language mass brain thought presentation
10	fire act min sunday encounter concentration daily active th monthly
11	paper common check photo add type virus print christmas present
12	game mobile internet app olympic london medal online device mm_mm
13	public right patient human goal influence earth plant face individual
14	help free american likely provide need support contact tip hand
15	company school design content employee college technology create offer audience
16	new state top city rank york art west east california
17	business customer service population sale product small software increase investment
18	force industry car line waste register decrease driver victim throw
19	year world people day make time com average number source
20	user use facebook share site video post google search worldwide

Table 5.2: **Top 20 topics in InfographicVQA found using LDA.** We used text tokens spotted on the images for topic modelling.

The number of average text tokens in InfographicVQA images is 217.89. Compared to ST-VQA and TextVQA, this value is much higher and it is even higher than the average number of OCR tokens in document images in DocVQA. This is expected since infographics in the dataset are often originally made to be used as posters or leaflets and generally contain more text tokens compared to a page from a business document. The word cloud of OCR tokens spotted in the images is shown in Figure 5.3b. Since infographics are rich in numerical data,

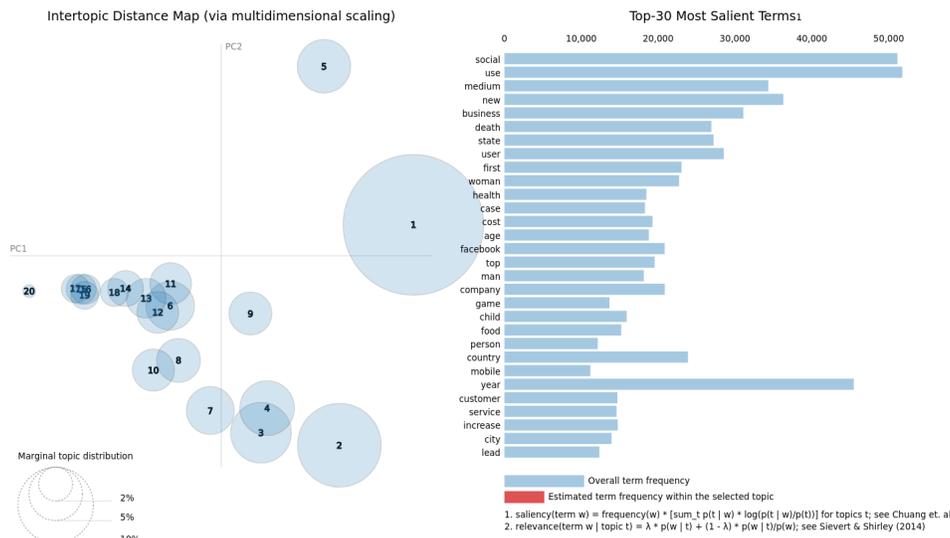


Figure 5.4: **Visualization of the top 20 topics in InfographicVQA dataset.** We used LDA to find the topics of infographics in the dataset. On the left is an inter-topic distance map where each circle represents a topic. The area of a circle is proportional to the overall relevance of the topic. Distance between two topics is computed using Jensen–Shannon divergence. The distances are then mapped to a two dimensional space using Multidimensional Scaling [229]. On the right, we show the top 30 most salient terms(most prevalent terms in the entire corpus) among the text present in the images. This diagram was created using the pyLDAvis tool.

Dataset	Questions		Answers		OCR tokens
	%Unique	Avg. len	%Unique	Avg. len	Avg. len
ST-VQA	84.84	8.80	65.63	1.56	7.52
TextVQA	80.36	8.12	51.74	1.51	12.17
VisualMRC	96.26	10.55	91.82	9.55	151.46
DocVQA	72.34	9.49	64.29	2.43	182.75
InfographicVQA	99.11	11.54	48.84	1.60	217.89

Table 5.3: Statistics of questions, answers, and OCR tokens in InfographicVQA and other similar VQA datasets.

numbers dominate the word cloud. It is evident from the two word clouds that the common answers and common text tokens in the infographics are similar.



Figure 5.6: **Input representation for the VLL-BERT.** Input to the model comprises text tokens from the question, OCR tokens, and visual tokens. $E()$ denotes a learnt embedding function. In the figure, we show a case where there are K question tokens, N OCR tokens, and M visual tokens. x_1 , y_1 , x_2 , and y_2 are the 4 coordinates of the bounding box of an OCR token or a visual token. Visual tokens correspond to local image regions.

5.4.1 Model architecture

Our model backbone is a stack of transformer layers like in BERT. Specifically, we use the $BERT_{BASE}$ architecture.

5.4.2 Input representation

Input to the model consists of three types of tokens: i) tokens from the question, ii) tokens from the text spotted on the image using an OCR, and iii) visual tokens. The visual tokens we use are either grid-based or region-based. The two styles of visual tokens are inspired by the use of grid-based and region-based features that are popular in the VQA literature(subsection 2.4.1). In the case of grid-based visual tokens, the tokens correspond to different spatial locations on a convolutional feature grid which in turn correspond to a uniform grid of image

patches. In the case of region-based visual tokens, each token corresponds to the spatial extent of an object detected on the image using an object detector.

The input stream starts with a special [CLS] token, followed by the question tokens, OCR tokens, and visual tokens. Question and OCR tokens are separated by a special [SEP] token. Similarly, OCR tokens and visual tokens are also separated by a [SEP] token. A token in the input stream is represented by a representation that is the sum of the following embeddings.

Text embedding. Following the original setting in BERT, for text embedding, we use the WordPiece embedding [137] that uses a vocabulary comprising 30,000 tokens. These embeddings are of size 768. For the question and the OCR tokens, WordPiece embeddings of the respective text tokens are used. For all the visual tokens, we use a special text token: [VisTok]. Consequently, all the visual tokens get the same text embedding; a WordPiece embedding for the token [VisTok].

Segment embedding differentiates different segments in the input. Question tokens, OCR tokens, and visual tokens are given segment ids of 0, 1, and 2, respectively.

1D position embedding is used to indicate the order of a token within the input stream. For the question tokens, the sequential order of the tokens is used. For the OCR tokens, we use the default reading order by serializing the tokens in top-left to bottom-right order. For the visual tokens, as there is no 1D order for the visual tokens, all of them are given the same 1D position and consequently have the same 1D position embedding.

2D position embedding is used to capture the 2D position of the OCR tokens and visual tokens in the image plane. For this, we follow the same process as LayoutLM. Given the bounding box coordinates (x_1, y_1, x_2, y_2) of an OCR token or a visual token, we embed all 4 coordinate values using 4 separate embedding layers. Bounding box coordinates for the grid-based visual tokens are found by estimating the spatial extent of the image patch that each spatial location on the feature map is based on. If the image is of size $W \times H$ and the feature map is of size $w \times h \times c$, then $w \cdot h$ feature vectors each of size c correspond to a uniformly spaced grid over the input image. For example, bounding box coordinates for the first visual token—the one corresponding to the top left location in the grid—will be $(0, 0)$ and $(W/w, H/h)$. Similarly, coordinates for the last visual token—the one corresponding to the bottom right location in the grid—will be $((w - 1)W/w, (h - 1)H/h)$ and (W, H) . x_1 and x_2 share the same embedding table and y_1 and y_2 share another common embedding table. Similar to LayoutLM, the coordinate values are normalized to lie in the range 0–1000 before embedding. For the [CLS] token, we use the 2D embedding corresponding to the coordinates $(0, 0, 1000, 1000)$. For the question tokens, all four 2D position embeddings are set to 0s.

Visual embeddings / Visual features. For a token, visual embedding is a visual feature (VF) corresponding to the token. Unlike the other embeddings that are learnt using an embedding

table, the visual embedding is extracted using a separate network. The visual features we use as visual embedding for the tokens are of two styles—grid-based and region-based.

In either case, there are no visual embeddings for the question tokens and [SEP] tokens. They are set to zeros. When the grid-based visual tokens are used, visual features for the visual tokens are the grid-based features corresponding to each token. Assume the feature map from the CNN used for the grid-based features is of the shape $w \times h \times c$. Then each location on the spatial extent of the feature map is a visual token and the c dimensional feature vector at the location is the visual embedding/visual feature for the token. The visual feature for the [CLS] token is taken as the average of visual embeddings for all the visual tokens. We refer to it as ‘Grid avg.’ in further discussion. Visual embeddings for OCR tokens are set to zero when the grid-based visual tokens are used.

When region-based visual tokens are used, visual embedding for the visual tokens are ROI pooled feature corresponding to each visual token. Remember in this case a visual token corresponds to an object detected on the image using an object detection model. The ROI pooled feature is extracted from the box head of an object detection model like Faster-RCNN. Similarly, for the OCR tokens, ROI pooled features corresponding to the bounding box of the OCR token are extracted. For the [CLS] token, the ROI pooled feature corresponding to a bounding box spanning the entire image is used. In the following discussion, we will refer to it as ‘Full Img.’

In Figure 5.6, we show how an input stream comprising question tokens, OCR tokens, and visual tokens are represented as a sum of text, layout—1D and 2D position—, visual and segment embeddings.

5.4.3 Pretraining

At the time of pretraining, similar to how Large Language Models(LLM) like BERT are trained, we train the model in a self-supervised manner. Following the LayoutLM, we use Masked Visual-Language Model (MVLM) task for pretraining with a masking probability of 0.15. The training objective is to predict the masked token. Whenever masking, we replace an OCR token with the [MASK] token 80% of the time, with a random token 10% of the time, and keep it unchanged 10% of the time. Note that we do not mask or zero out the positional information or the visual information for the token that is being masked. A schematic of the pretraining setup is shown in Figure 5.7.

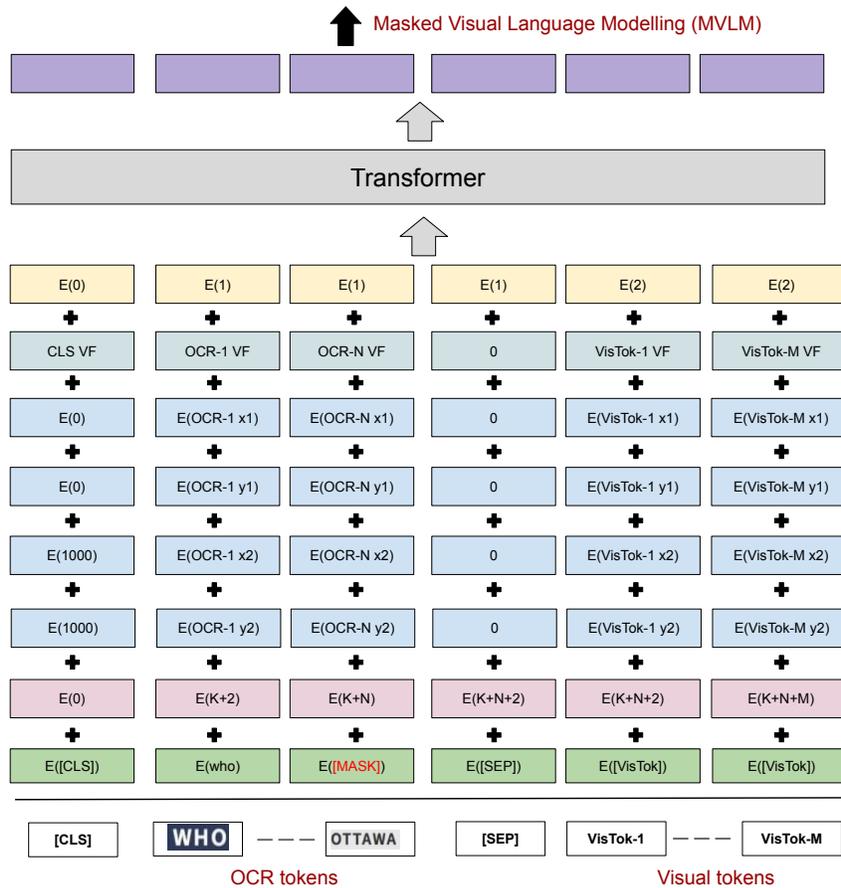


Figure 5.7: **Pretraining LayoutLM-based model using MVLM** We pretrain the model using Masked Visual Language Modelling (MVLM). MVLM aims to predict a masked OCR token given other OCR tokens and visual tokens. In this case, the OCR token "ottawa" is picked for masking. For this token, text embedding of a special [MASK] token is used. Other embeddings remain the same.

5.4.4 Finetuning

Similar to the output head of the BERT QA model (subsection 2.2.1.1, we use an output head that predicts the answer as a single continuous span of the OCR tokens spotted on the image. Therefore, our model can only handle questions whose Answer-source is Image-span. Nearly 70% of questions in validation and test split are of this type. Extending this model to include the questions with other Answer-sources—Question-Span, Multi-Span, and Non-Span—is an interesting direction for future work. Figure 5.8 depicts the finetuning setup.

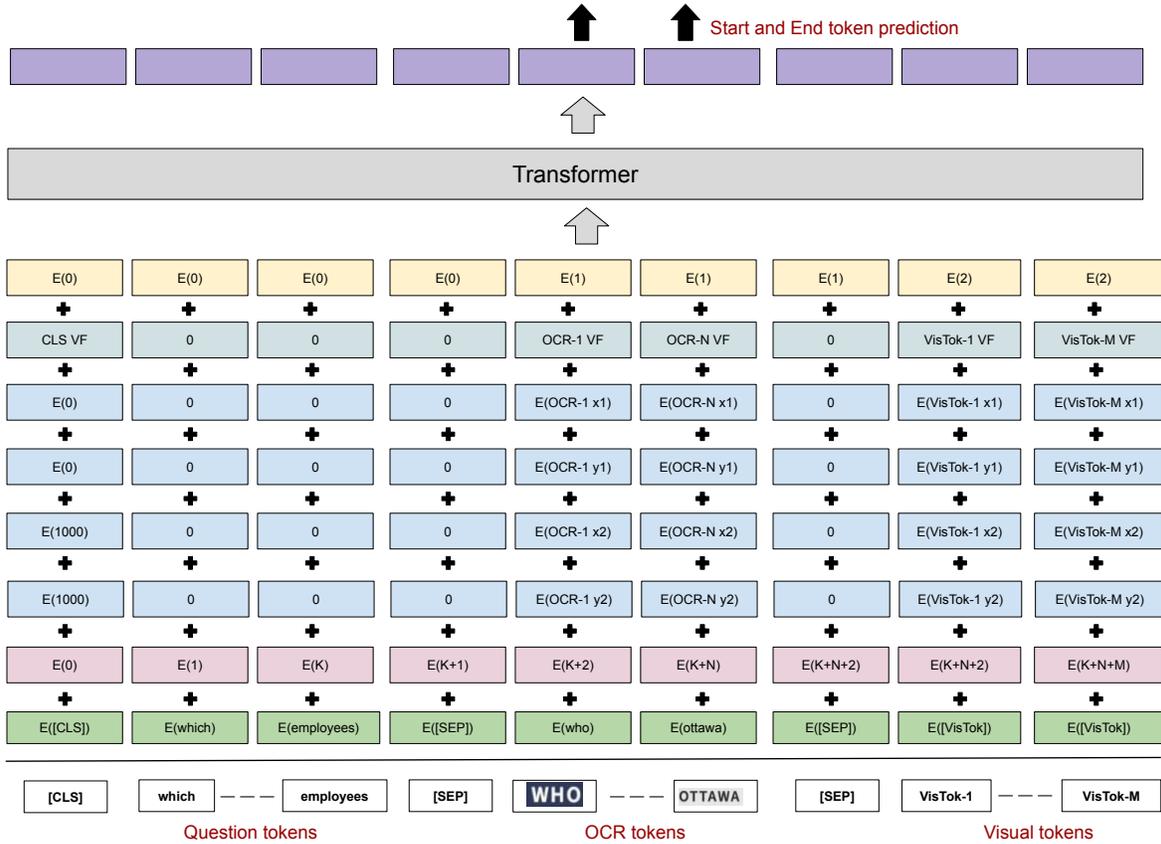


Figure 5.8: **Finetuning LayoutLM-based model for extractive VQA.** At the time of finetuning, tokens from both the question and the image are used. The training objective is to predict the Start and End tokens of the answer among the OCR tokens. This finetuning setup is similar to finetuning BERT for MRC(subsection 2.2.1.1).

5.5 Experimental setup

In this section, we present details of the experiments we conduct, implementation details of models we use, and the evaluation metrics.

5.5.1 Experimental setup for VLL-BERT

5.5.1.1 Extracting visual features

Following Anderson et al.’s work [93] that proposed to use attention over both the region-based and grid-based features, models like LoRRA (subsection 2.5.1.1) followed suit to utilize both the grid-based and the region-based features. However, the recent works on VQA,

such as M4C and multimodal pretrained models like VisualBERT (for a list of models similar to VisualBERT, see subsection 5.2.2), use only grid-based features. All these models make use of the ROI pooled features from a Faster-RCNN [96] trained for object detection on Visual Genome [146]. LayoutLM, which is a pretrained model for DIA tasks, uses region-based features as ‘image embedding’ of the OCR tokens. Although document images are quite different from natural images, authors of LayoutLM use a Faster-RCNN trained on Visual Genome dataset, for region-based features.

Inspired by the use of the region-based features in the aforementioned works, we first decided to follow the same approach in VLL-BERT. In addition to the region-based features from an object detection model for natural images, we also try out region-based features using a model that is trained for Document Layout Analysis (DLA). In DLA, the objective is to detect layout objects of interest, such as titles, text blocks, and tables. Our initial experiments showed that VLL-BERT which uses region-based features performed poorer compared to the VLL-BERT variant which doesn’t use any visual features. Consequently, we investigated other means to incorporate visual information into the VLL-BERT model. Ultimately as presented in subsection 5.4.2, in addition to region-based features, we also experimented with grid-based features.

For grid-based features and region-based features following are the 5 different visual feature (VF) sources that we experiment with.

1. **VG-regions.** Similar to M4C, VisualBERT and other recent works, we use a Faster-RCNN trained on Visual Genome (VG). We used the scripts and pretrained Faster-RCNN checkpoint used in the MMF framework [214]. As done in M4C, for VLL-BERT, we extract ROI pooled features of size 2048 from the penultimate fully-connected layer of the box head of this network. But unlike in M4C, we do not finetune the last fully-connected layer of the box head while pretraining or finetuning the VLL-BERT.
2. **DLA-regions.** It is a Mask-RCNN trained on document images in PubLayNet [48] for DLA. PubLayNet dataset has annotations for 5 document layout object classes—text, title, list, table, and figure. We used a publicly available Detectron2 ⁴-based implementation ⁵ of Masked-RCNN trained on PubLayNet. We use the pretrained checkpoint provided in the repository. ROI pooled features of size 1024 from the last fully-connected layer of the box head of the model are used as the visual features.

⁴<https://github.com/facebookresearch/detectron2>

⁵<https://github.com/hpanwar08/detectron2>

3. **ImageNet-grid.** Grid-based features from a ResNet-152 trained for image classification on the ImageNet dataset are used. We used the pretrained checkpoint available in PyTorch. The feature map from which we extract the grid features is of size $7 \times 7 \times 2048$
4. **RVL-CDIP-grid.** We finetuned the ‘Imagenet-grid’ model on RVL-CDIP dataset [232] for document image classification. The RVL-CDIP is a document image classification dataset containing 400,000 images in 16 document classes. The images in RVL-CDIP are a subset of the IIT-CDIP [147] collection.
5. **Visually29K-grid.** We finetuned a pretrained ResNet-152—the one that is used for ImageNet-grid—for multi-label image classification on Visually29K [224]. In Visually29K, each infographic is annotated with a category label and one or more tags. There are a total of 26 different categories and 391 different tags. While finetuning, we replace the last fully-connected layer with a new fully-connected layer that has 391 units and uses Sigmoid activation (since the model is being finetuned for multi-label classification) for the last layer. The model is finetuned using binary cross-entropy loss.

VLL-BERT is implemented in PyTorch. In all our experiments, we start from a pretrained checkpoint of LayoutLM. Specifically, we use the ‘layoutlm-base-uncased’ checkpoint⁶. Text, segment, and 1D and 2D position embeddings are of size 768 (see subsection 5.4.2). Since the visual features we use is of a different size—2048 in the case of VG-regions, 2048 for the grid-based features, and 1024 in the case of DLA-regions—than the common embedding size of 768, we use a linear layer to project the visual features to this common embedding size. This layer is initialized from scratch. For pretraining on InfographicVQA using MVLm, we used Adam optimizer with a learning rate of $2e - 5$. We used a batch size of 4. For finetuning, we used a batch size of 8 and a learning rate of $1e - 5$. For pretraining—i.e., in-domain pretraining on InfographicVQA—and finetuning, no additional data other than train split of InfographicVQA is used. Since we train VLL-BERT for extractive VQA, answers need to be mapped as spans of OCR tokens. To map answers in the InfographicVQA train split to SQUAD-style spans, we follow the same approach we used for DocVQA(see subsection 4.4.2). We take the first subsequence match of an answer in the serialized transcription as the corresponding answer span. This way, we find approximate spans for 52% of questions in the train split. The rest of the questions are not used for finetuning the model.

⁶<https://huggingface.co/microsoft/layoutlm-base-uncased>

5.5.2 Experimental setup for heuristics, upper bounds, and human performance.

On DocVQA, we evaluated various heuristic baselines and upper bounds as presented in subsection 4.3.1. We evaluate them on InfographicVQA as well. For evaluating human performance, all questions in the test split of the dataset are answered by a volunteer. For Vocab UB and heuristics involving a vocabulary, we use a vocabulary of 5,000 most common answers in the train split.

5.5.3 Experimental setup for M4C.

In addition to VLL-BERT, we evaluate M4C on InfographicVQA. We use the official implementation of the model [214]. The training parameters and other implementation details are the same as the ones used in the original paper. As done in the original M4C, the fixed vocabulary used with the model is created from the 5,000 most common words among words from answers in the train split. The original M4C uses region-based features from a Faster-RCNN trained on Visual Genome. We call this feature extractor VG-regions in subsection 5.5.1.1. In addition to it, we experiment with region-based features from DLA-regions as well. We also evaluate an M4C variant without visual features and one with features drawn from only one region that covers the entire image extent.

5.5.4 Evaluation metrics

The evaluation setup is same as the evaluation for DocVQA (see subsection 4.4.1). We use two metrics: ANLS and accuracy.

5.6 Results

The results of heuristic baselines, upper bounds, and human performance are shown in Table 5.4. Human performance is comparable to human performance on DocVQA. This suggests that humans can answer questions on infographics with the same level of easiness as answering questions on business documents such as letters and reports. Vocab + OCR UB suggests that more than three-quarters of questions have their answers present as a span of the OCR tokens or in a vocabulary of the most common answers in the train split.

Baseline	ANLS		Accuracy(%)	
	val	test.	val	test
Human performance	-	0.980	-	95.70
Random answer	0.006	0.005	0.00	0.00
Random OCR token	0.011	0.014	0.29	0.49
Majority answer	0.041	0.035	2.21	1.73
Vocab UB	-	-	53.16	51.34
OCR subsequence UB	-	-	53.95	56.96
Vocab + OCR subsequence UB	-	-	76.71	77.4

Table 5.4: **Results of heuristics and upper bounds and human performance evaluation.** Heuristics yield near-zero results. From the numbers for different upper bounds, it is evident that more than 75% of the answers are either present in the present vocabulary or found as a sub sequence of the OCR tokens spotted on the respective image.

We show results using M4C in Table 5.5. The first row in the table corresponds to the settings used in the original M4C. The other rows show the results of different ablations and variants of the original setting. It can be seen that the region-based features from VG-regions and DLA-regions do little help for VQA on InfographicVQA. At the same time, using region-based feature corresponding to the entire image (i.e., the ‘Full Img.’) works better than the variant that does not use any visual features. With more OCR tokens taken as part of the input to the M4C, the performance is steadily increasing. In the original M4C, the maximum number of OCR tokens that are fed to the model is 50. However, on average, InfographicVQA images have a lot more OCR tokens. Therefore, including more OCR tokens helps with the performance.

The results of extractive VQA using VLL-BERT are given in Table 5.6. The first row in the table is essentially a LayoutLM finetuned on the InfographicVQA. The second row shows the case when the LayoutLM is pretrained and finetuned on InfographicVQA. It can be seen that the in-domain pretraining significantly improves the performance. The rest of the rows show results of VLL-BERT that use visual features, either region-based or grid-based. Similar to M4C, VLL-BERT variants that use region-based features using VG-regions or DLA-regions perform poorer compared to the VLL-BERT which does not use any visual features. Although the grid-based features from ImageNet-grid and RVL-CDIP-grid do little to help in improving

VF source	Finetune VF source	Region type	# regions	# OCR tokens	val		test	
					ANLS	Acc.	ANLS	Acc.
VG-regions	✓	objects	100	50	0.11	4.81	0.12	4.87
VG-regions	✓	objects	20	50	0.11	4.82	0.12	4.87
VG-regions	✗	objects	20	50	0.13	4.89	0.13	4.89
VG regions	✗	objects	20	300	0.13	4.90	0.13	5.08
VG-regions	✗	None	0	300	0.14	5.86	0.14	6.58
VG regions	✗	Full Img.	1	300	0.14	5.93	0.15	6.64
DLA-regions	✗	objects	20	50	0.11	4.86	0.13	5.02
DLA-regions	✗	objects	20	300	0.13	5.95	0.14	6.50
DLA-regions	✗	None	0	300	0.14	5.90	0.14	6.39
DLA-regions	✗	Full Img.	1	300	0.14	5.97	0.14	6.42

Table 5.5: **Performance of different variants of the M4C.** The setting used in the original M4C is the one shown in the first row. The column ‘VF source’ indicates the source of the visual features (VF) used. ‘Finetune detector’ denotes the case when features from the penultimate fully-connected layer of the box head of the detector are used for VQA while the last fully-connected layer is finetuned along with the M4C model. This is the default setting in M4C. In our experiments, we get better results without finetuning the detector. In the column ‘Region type’, ‘objects’ means grid-based features corresponding to the objects detected on the images are used. ‘Full Img.’ stands for the case where only one region-based feature corresponding to the bounding box covering the entire infographic is considered. Using no visual features and only the Full Img’s feature gives better results than using the region-based features corresponding to the objects detected on the image. The original M4C considers only the first 50 OCR tokens. We experiment with up to 300 OCR tokens since the number of OCR tokens per image is much higher for InfographicVQA compared to the scene text VQA datasets, as shown in Table 5.3.

the performance, the grid-based features from Visually29K-grid yield the best among all the variants we try out.

We believe the region-based visual features we use impart little information since the object detectors we use—a detector trained for detecting objects on natural scene images and another for document layout analysis—are not suitable for infographics. This is evident from Table 5.7. Both detectors detect only a few instances of objects on infographics.

In-domain pretraining	VF source	Visual tokens	OCR VF	[CLS] VF	val		test	
					ANLS	Acc.	ANLS	Acc.
✗	✗	✗	✗	✗	0.21	13.40	0.23	15.30
✓	✗	✗	✗	✗	0.25	18.14	0.27	19.74
✓	DLA-regions	✗	✗	Full Img.	0.25	18.56	0.26	19.16
✓	VG-regions	✗	✗	Full Img.	0.22	16.47	0.24	16.51
✓	VG-regions	✓	✗	Full Img.	0.19	10.20	0.19	10.60
✓	DLA-regions	✓	✗	Full Img.	0.23	16.14	0.25	17.58
✓	VG-regions	✓	✓	Full Img.	0.18	9.90	0.17	10.30
✓	DLA-regions	✓	✓	Full Img.	0.24	16.14	0.25	17.61
✓	ImageNet-grid	✓	✗	Grid avg.	0.21	16.22	0.24	16.21
✓	RVL-CDIP-grid	✓	✗	Grid avg.	0.20	15.40	0.20	16.18
✓	Visually29K-grid	✓	✗	Grid avg.	0.28	19.91	0.29	21.10

Table 5.6: **Performance of VLL-BERT on InfographicVQA.** The first column indicates if the VLL-BERT is pretrained using MVLM on InfographicVQA. For different visual feature sources that we show in the second column, refer to subsection 5.5.1.1. The ‘OCR VF’ column indicates if visual embeddings are used for the OCR tokens. When visual features are not used, they are set to zeros. Similarly, the ‘CLS VF’ column shows what visual feature is used for the [CLS] token. When region-based visual features are used, the visual feature for the [CLS] token is the ROI pooled feature corresponding to the full image (abbreviated as ‘Full Img.’). In the case of grid-based features, the visual feature for [CLS] is computed as the average of the visual features for all the visual tokens. We refer to it as ‘Grid avg.’.

In Figure 5.10, the performance of our trained baselines on the test split is compared against the upper bounds and human performance. The M4C and VLL-BERT models used for this comparison are the variants that yield the best ANLS on the test data. A few qualitative results of our experiments are shown in Figure 5.9. More qualitative examples are shown in Appendix C.

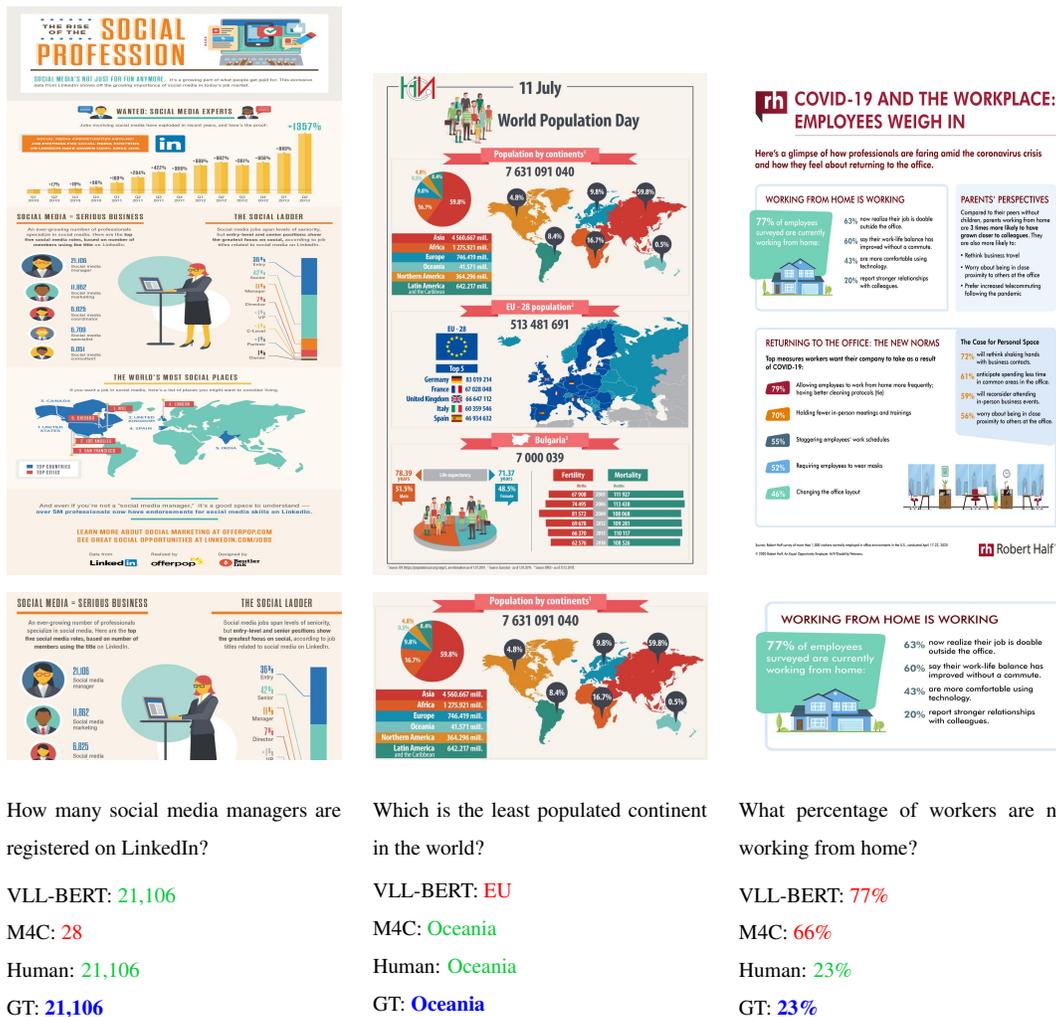


Figure 5.9: **Qualitative results.** Shown in the first row are the original infographics on which the questions are asked. The second row shows the crops of the respective original images that are relevant to the question. For the leftmost question, VLL-BERT gets the answer correct, while M4C fails. In the case of the second question, where evidence is a Table/List and Sorting is involved, M4C answers the question correctly. In the case of the last question that requires the subtraction of a number from another, neither M4C, nor VLL-BERT gets the answer correct.

5.7 Summary

We introduce a new dataset for VQA on infographics with the goal of developing a VQA benchmark that necessitates non-trivial multimodal reasoning in order to respond to the questions. We present a comprehensive comparison of the new dataset with other similar datasets.

Detector	TextVQA		DocVQA		InfographicVQA	
	Avg.	<2 det.(%)	Avg.	<2 det.(%)	Avg.	<2 det.(%)
VG-regions	28.8	0.0	4.1	43.9	7.4	23.9
DLA-regions	1.0	97.9	4.7	0.0	2.9	43.4

Table 5.7: **Statistics of object detections using two detectors – VG-regions and DLA-regions.** DLA-regions is trained for detecting document layout objects in PubLayNet, and VG-regions is an object detection model trained on natural images in Visual Genome. Avg. shows the average number of detections per image. ‘<2 det.(%)’ is the percentage of images on which the number of detected objects is less than 2.

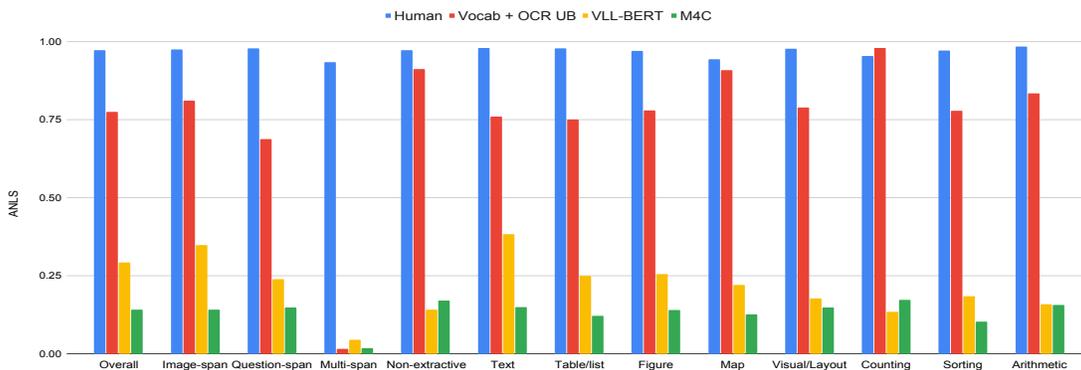


Figure 5.10: **Performance of VLL-BERT and other baselines for different question-answer types.**

Inspired by multimodal transformer-based models, we propose a BERT-like model for extractive VQA on InfographicVQA. The proposed model, called VLL-BERT, takes in textual, layout, and visual features as its input. Additionally, we evaluate different variants of M4C and many heuristic baselines and upper bounds on the dataset. Unlike images in DocVQA, infographics in InfographicVQA are significantly different from generic document images like the ones in popular DIA datasets such as IIT-CDIP, RVL-CDIP and PubLayNet. Infographics are not like natural images either. Due to this distribution shift, visual feature extractors(both grid-based and region-based) trained on natural images or generic document images were proven to be ineffectual in extracting visual features from Infographics.

Chapter 6

Recognition-free QA on handwritten document collection

DocVQA and InfographicVQA adhere to the mainstream VQA setting where a question is asked on a single image and a textual answer is returned. Recognizing text on the given image and/or understanding its semantics is essential for building solutions for the two tasks. Consequently, noisy OCR can act as a bottle neck for recognition-based solutions to VQA on document images. In this chapter, we investigate the possibility of extracting information from document images in a recognition-free manner. To this end, we propose a new recognition-free QA task. We prefer to call this a QA task not a VQA, as the handwritten document images we work with are text-only pages with little visual content. Additionally all the images contain text in the form of paragraphs in a single column layout, making the layout analysis trivial. Unlike typical QA and VQA formulations where the answer is a short text, we aim to locate a document snippet where the answer lies. The proposed approach works without recognizing the text in the documents. We argue that the recognition-free approach is suitable for handwritten documents and historical collections where robust text recognition is often difficult. At the same time, for human users, document image snippets containing answers act as a valid alternative to textual answers. The proposed approach uses an off-the-shelf deep embedding network that can project both textual words and word images into a common subspace. This embedding bridges the textual and visual domains and helps us retrieve document snippets that potentially answer a question. We evaluate results of the proposed approach on two new datasets: (i) HW-SQuAD: a synthetic, handwritten document image counterpart of SQuAD1.0 dataset, and (ii) BenthamQA: a smaller dataset containing QA pairs defined on documents from the Bentham manuscripts collection. We also present a thorough analysis of the proposed recognition-free approach compared to a SoTA recognition-based approach. This work has previously been presented in our publication [7].

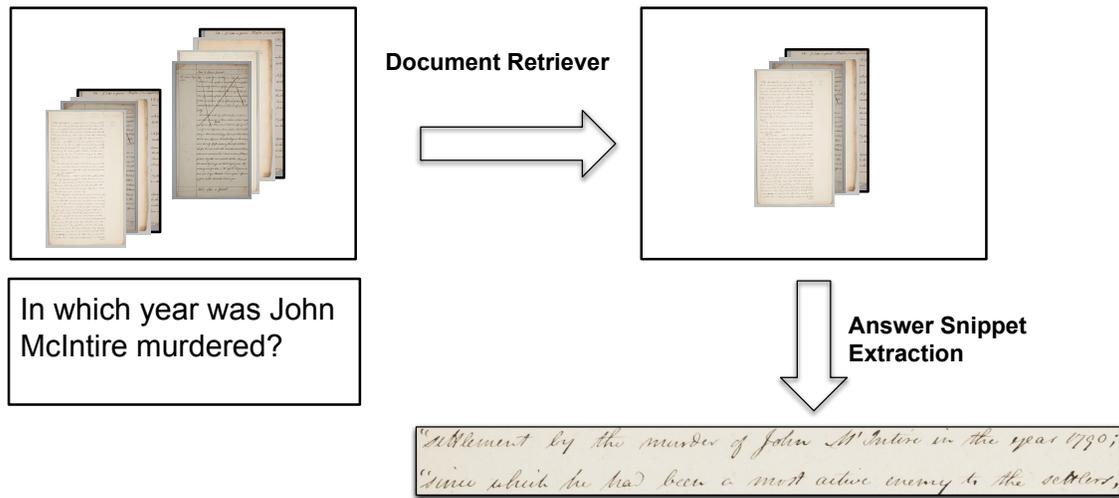


Figure 6.1: This work addresses the problem of Question Answering on handwritten document collection. Given a natural language question, a document snippet that answers the question is returned.

6.1 Introduction

Although deep learning-based techniques have led to considerable performance improvement in handwriting recognition [101, 233, 234], it is still far from OCR performance on printed documents. Owing to variability in handwriting and prevalence of degradation, text recognition from manuscripts—particularly historical manuscripts—is a challenging problem. Recognizing text on a new manuscript collection using a generic OCR that is not finetuned or retrained on the target collection, often result in noisy OCR transcriptions. A recent study by Bazzo et al. [235] observes that even a 5% word error significantly impacts information retrieval from document images that are automatically transcribed using an OCR. In another study based on data in a large digital library, Chiron et al. [236] observe that a significant number of user queries are affected by OCR errors. They also observe that at least 15% of the OCR errors are for named entities that dictionary-based error correction approaches cannot correct. This motivates us to devise a retrieval-based solution for QA that does not require text recognition and extracts a relevant region from an image in the collection that answers the question.

We propose to return answers to the questions as image/document snippets. Figure 6.1 depicts an example case where the answer to a natural language question on a document collection is an image snippet/crop from one of the documents in the collection. Although this approach does not yield a textual answer to the question, it would still meet the requirement of human users looking for information from document collections. This is in line with works in VQA,

which argue that visual evidence is as important as coming up with a textual answer [237, 238]. For example, in the STE VQA dataset [238], in addition to the textual answer for a question, a visual answer is also provided. The visual answer is a bounding box representing the image region on which the question is grounded. Moreover, returning image snippets as answers is similar to extractive QA tasks in NLP/IR, where an answer is not generated but returned as a span (a sequence of contiguous text tokens) of the given context [91, 239].

We propose a method that performs QA in the image space without any explicit text recognition from the documents in the collection. We use an end-to-end deep embedding network [101] to extract feature embeddings of both textual words in questions and word images in documents. These representations are aggregated for each question, document, or snippet to form question vectors, document vectors, and snippet vectors, respectively. This way, documents best matching a question or snippets best matching a question can be retrieved using the nearest neighbour search.

Major contributions of this chapter are:

- Formulate the problem of information extraction from a document image collection as a question answering task. This is motivated by VQA in Computer Vision and QA in NLP. The question is defined as a natural language query resulting in an answer in the form of image snippets.
- Introduce two new datasets: HW-SQuAD and BenthamQA, for QA on handwritten document collection. HW-SQuAD is curated using an existing dataset for electronic text — SQuAD1.0 [91]. We render passages in SQuAD1.0 as images using handwritten fonts. For BenthamQA, we annotate questions and answers on handwritten manuscripts from Transcribe Bentham Project [17].
- Propose an evaluation protocol to evaluate QA, where the answer is a document image snippet.
- A two-stage solution for QA on document collection and evaluate it on the newly introduced datasets. Our method is segmentation-based (requires word and line bounding boxes), recognition-free (does not require to recognize text in the documents), and lexicon-free (does not depend on a fixed vocabulary). We compare our approach against a recognition-based approach and analyze how the two approaches fare when text recognition is noisy.

6.2 Related works

Work presented in this chapter is related to QA and VQA works in NLP and Computer Vision. We discuss important tasks, datasets and popular methods for QA on electronic text in section 2.2. The main difference between QA on electronic text and the problem we address in this chapter is that, text, in our case is not in machine readable form. We address the problem of answering questions asked on a set of document images without any given transcription.

As noted in section 2.2, most of the early VQA datasets and methods disregard text present in the images, and the problem is often modelled as a multi-class classification. Questions in these datasets typically focus on visual aspects, such as objects, attributes, and relationships. On the other hand, scene text VQA, requires reading text on natural images(see subsection 2.5.1). Our work is different from the scene text VQA on two accounts: (i) these datasets contain images “in the wild” which are drawn from popular scene text datasets or datasets like OpenImages [156] that predominantly have scattered text tokens compared to the handwritten document images we consider, and (ii) almost all VQA problems, including the ones involving text on the images, are formulated as QA on a single image, while the proposed QA task is for a collection of document images.

Our works on single document VQA—DocVQA and InfographicVQA—that we present in the previous two chapters are closely related to the work presented in this chapter. While DocVQA and InfographicVQA stick to the standard VQA setting where the answer is textual, work done in this chapter propose to respond to natural language queries by providing “visual answers” in the form of image snippets. HW-SQuAD and BenthamQA datasets contain handwritten and/or historical documents whose OCR transcription is generally noisy. Compared to multimodal document images in DocVQA and InfographicVQA that feature tables, text, figures and photographs, the document images in datasets for the proposed recognition-free QA task contain only handwritten or historical manuscripts that predominantly contain text in the form of sentences and paragraphs.

6.3 Recognition-free QA

In the following discussion, we use the term ‘document’ to refer to a handwritten document image. ‘Word’ refers to a textual word in questions and a word image in documents. We define a document snippet or a ‘snippet’ as a horizontal slice of one or more contiguous text lines from a document.

The proposed QA method returns the answer to a given question as a snippet extracted from one of the queried documents. Given a question like *In which year was John McIntire*

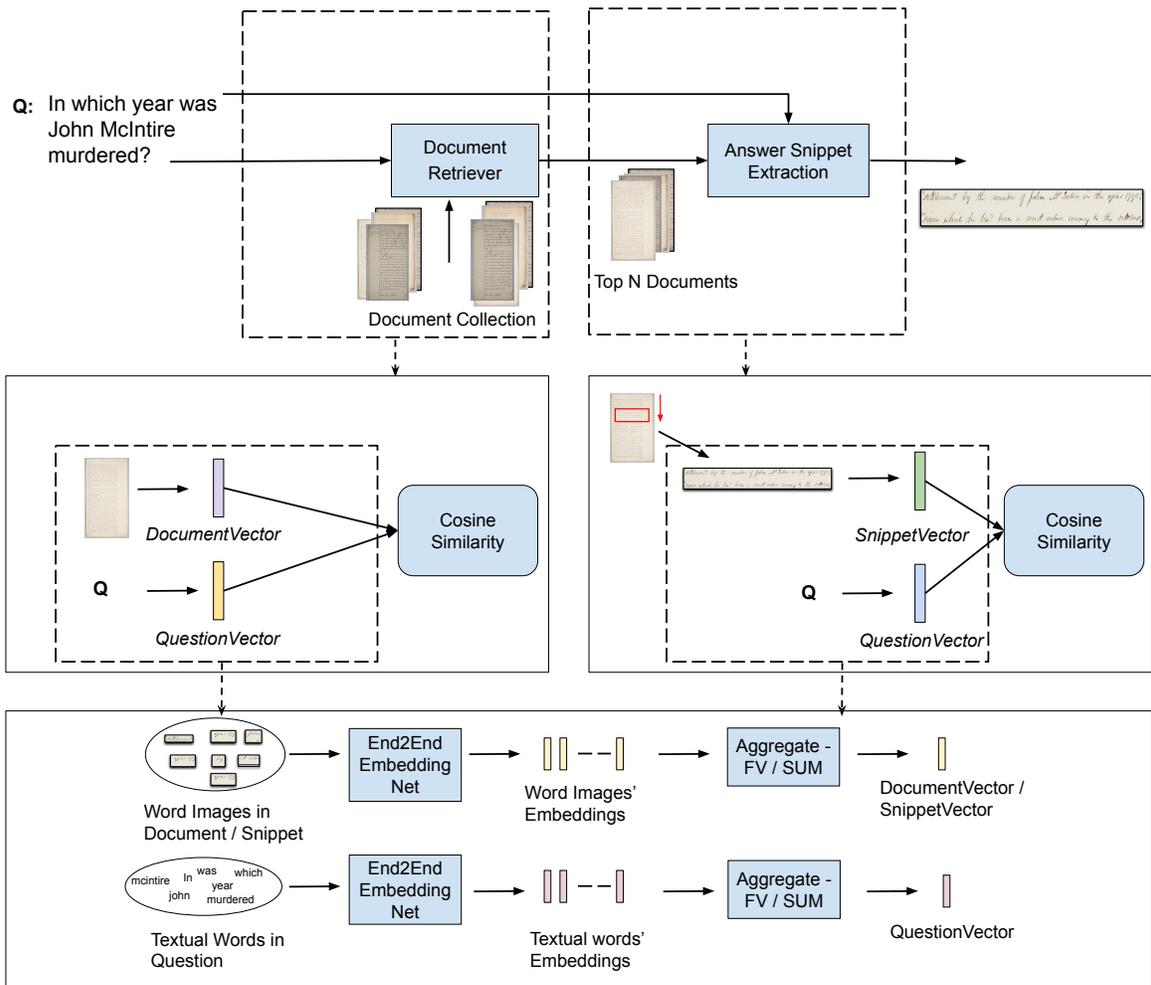


Figure 6.2: **Proposed recognition-free approach to QA on document collection.** Our approach works by aggregating embeddings of words in a question, a handwritten document or a document snippet to a *QuestionVector*, *DocumentVector* or *SnippetVector* respectively. We use an end-to-end embedding network [101] for the embeddings. These vector representations help us to retrieve documents similar to a question or snippets similar to a question in a Vector Space model. The snippets are extracted in a two-step process. In the first step, a small number of relevant documents are retrieved. In the second step, a document snippet answering the question is extracted from one of the previously retrieved documents.

murdered? (see Figure 6.1) we want to find an embedding/representation for the question. We aim to find a snippet in the new embedding space that can be the most probable answer using a simple nearest neighbour search. This approach leads to the following technical challenges:

- We need an embedding where both text (questions) and image (document/snippet) modalities can be embedded.
- The embedding of a question must be independent of its length and its paraphrasing.
- In the joint embedding space, a question’s matching snippet must have its embedding similar to the question.

We describe a solution that meets the above requirements below. Figure 6.2 shows a schematic representation of the proposed solution.

6.3.1 Joint embedding of textual words and word images

In order to come up with a joint space where questions and snippets can match, we first embed the constituent words of a question or snippet into a joint space. Later these individual embeddings are aggregated to form a single global embedding for the question or the snippet. For the joint embedding of words, we use an end-to-end embedding network introduced in [101], which simultaneously learns both text and image embeddings using a multi-task loss. It is an end-to-end trainable network which does not require to embed image and text in separate steps and later combine them. The network learns a feature space where words that are lexically similar lie closer to each other. This facilitates word spotting in this joint space in both query-by-string and query-by-example settings.

6.3.2 Aggregating word embeddings

The answer extraction process we discuss below requires us to represent questions, documents, and snippets as vectors of a fixed dimension. This is achieved by aggregating the embeddings of the constituent words of a question/document/snippet into a global representation whose size is independent of the number of words being aggregated. Stop words in questions and documents/snippets are removed prior to the aggregation of word embeddings (see auto6.5.1 for details). Two different approaches are tried out for aggregating the word embeddings and they are discussed below.

6.3.2.1 Aggregate by summing (SUM)

A set of M word embeddings $\{x_1, x_2, \dots, x_M\}$ of size D_w are summed together to form an aggregate vector of the same size.

6.3.2.2 Aggregate using Fisher Vector (FV) framework

The Fisher Vector (FV) framework introduced in [240] has been widely used to aggregate local image descriptors to a global descriptor in image retrieval and classification problems [241, 242, 243]. Given a set of word embeddings $X = \{x_1, x_2, \dots, x_M\}$ we assume that these continuous word embeddings have been generated by a Gaussian Mixture Model (GMM). We denote the parameters of the GMM, $\lambda = \{w_i, \mu_i, \Sigma_i, i = 1 \dots K\}$, where K is the number of Gaussians. w_i , μ_i and Σ_i are, respectively, the mixture weight, mean vector and covariance matrix of the Gaussian i . We assume that the covariance matrices are diagonal and hence it is denoted by the variance vector, σ_i^2 . The GMM, u_λ is estimated/trained offline from a representative set of samples using Maximum Likelihood (ML) estimation.

For the i^{th} Gaussian in the GMM, let $\mathcal{G}_{\mu,i}^\lambda$ be the gradient with respect to the mean (μ_i) and $\mathcal{G}_{\sigma,i}^\lambda$ be the gradient with respect to the standard deviation (σ_i). FV of the given set of embeddings X will be the gradient vector \mathcal{G}_X^λ which is obtained by concatenating $\mathcal{G}_{\mu,i}^\lambda$ and $\mathcal{G}_{\sigma,i}^\lambda$ for $i = 1 \dots K$. Following the standard practice, we do not consider gradients with respect to the mixture weights, w_i s, since it adds little discriminative information [244, 243]. Finally, we apply Power normalization and L2 normalization to the FV as suggested in [242]. Power normalization proposed by Perronnin et al. [243] helps to “unsparisfy” the FV as it becomes sparser with more Gaussians. To each dimension of the FV, power normalization applies the following function $f(z) = sign(z)|z|^\alpha$. $0 \leq \alpha \leq 1$ is the power normalization parameter and usually set to $1/2$.

6.3.3 Document retriever

Following the standard practice in QA systems for electronic text [138], we first select a small set of documents that are likely to contain the answer. We call this step “Document Retriever”. Given a question and a document collection of M documents, the Document Retriever generates N proposals ($N \ll M$) for the ‘target document’, i.e., the document where the potential answer to a given question lies. For this step, the question and the documents are summarized into a *QuestionVector* and *DocumentVectors* of fixed size using one of the aggregation methods mentioned above. Subsequently, target document proposals are generated by retrieving N documents best matching the question by computing cosine similarity between the *QuestionVector* and the *DocumentVectors*.

6.3.4 Answer snippet extraction

Once a set of document images are retrieved—as proposals for the target document—by the Document Retriever, all possible snippets from these proposals are considered as ‘candidate snippets’. Since our approach assumes that word and line segmentations are known, the candidates are extracted in a sliding window manner, scanning text lines in a document from top to bottom. For instance, if the size of answer snippets is set as 2 text lines, we group each pair of consecutive text lines as a snippet. Then, candidate snippets are converted into *SnippetVectors* using one of the aggregation methods proposed in subsection 6.3.2. The aggregation method used for this step could be different from the one used with the document retriever. The best aggregation scheme for each step shall be found empirically by evaluating different options on a validation set. Finally, similar to Document Retriever, cosine similarity between the *QuestionVector* and *SnippetVectors* are computed and the top matching snippet is returned as the answer to the given question.

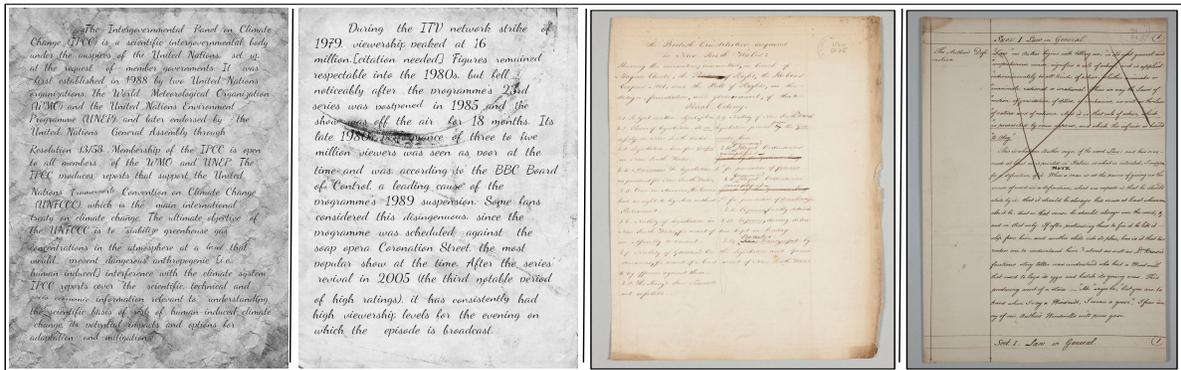


Figure 6.3: Sample documents in the newly introduced HW-SQuAD and BenthamQA datasets. The first two are synthetically rendered images in HW-SQuAD. The synthetic samples have diverse fonts, backgrounds, textures, image artifacts, and inter-word and inter-line spacing. The last two are real handwritten manuscripts that are part of the BenthamQA dataset.

6.4 HW-SQuAD and BenthamQA datasets

The annotation of a QA dataset on a handwritten document collection requires considerable human effort. Moreover, modern deep learning-based supervised models benefit from large amounts of training data. As an alternative to costly human annotation, we build a dataset using an existing QA dataset for electronic text by reusing the questions and answers. Since this

dataset is a derivative of the existing SQuAD1.0 dataset, we term it HW-SQuAD. We use HW-SQuAD for experimenting with different ablations of the proposed approach. For benchmarking the approach on a real collection of handwritten manuscripts, we create another dataset called BenthamQA. Images in BenthamQA are scanned images of handwritten manuscripts from the 18th and 19th centuries.

6.4.1 HW-SQuAD

SQuAD1.0 [91] is a popular benchmark dataset used by NLP and IR communities for QA and MRC on electronic text. It contains 100,000+ question–answer pairs on 23,000+ passages drawn from over 500 Wikipedia articles. Every question is associated with a passage, and the answer to the question is marked as a ‘span’ of contiguous words in the passage. The train set of the dataset has only one answer mapped per question. But in the development (public test) set, each question is provided with one or more target (ground truth) answers. The analysis presented in [91] shows the dataset has diverse questions in terms of i) syntactic and lexical variations and ii) the extent of reasoning required to answer the questions.

The passages need to be rendered as document images to reuse SQuAD1.0 for our purpose. Motivated by the success of synthetic datasets for OCR problems [124, 101], we synthetically render the passages in SQuAD1.0 as handwritten document images. In the following, we describe how we use SQuAD1.0 data to create the new HW-SQuAD dataset.

6.4.1.1 Data preprocessing

Samples in SQuAD1.0 are first subjected to standard pre-processing steps¹ to tokenize questions and passages, and to generate question–passage correspondences in terms of token offsets. Non-ASCII tokens in the questions and the passages are then converted to the nearest ASCII transliterations using the `Unidecode` library². This was required since most of the handwritten fonts we use do not support characters outside the ASCII character set. Following the pre-processing step, a few questions—corresponding to the answers for which mapping the character offsets to the token offsets fail—are excluded. The original train set of SQuAD1.0 is split randomly into two subsets in HW-SQuAD — 95% in train and 5% for validation (val). We use the original development set (public test set) of SQuAD1.0 to create the test set in HW-SQuAD.

¹<https://github.com/facebookresearch/DrQA>

²<https://pypi.org/project/Unidecode/>

	HW-SQuAD			BenthamQA
	train	val	test	
Number of questions	67887	7578	9477	200
Number of documents	17007	1889	2067	338
Average number of words per document	116.6	117.2	122.8	202.0
Average number of words per question	10.1	10.1	10.2	18.0
Average number of words per answer	2.4	2.4	2.7	3.9

Table 6.1: HW-SQuAD and BenthamQA dataset statistics.

6.4.1.2 Synthetic rendering of SQuAD1.0 passages as handwritten document images

We render a passage as a single-column document image using a handwriting-style font. The font used is randomly sampled from over 100 handwriting-style fonts downloaded from Google Fonts³.

The first step is to render each word in a passage in SQuAD1.0 onto a transparent background. The font size is set randomly from a range of 28–52 pts. The intensity of the foreground (i.e., the pixels making up the text strokes) is varied randomly in the 0–50 range (0 being the darkest and 255 being the brightest). This means all our documents have darker text on lighter backgrounds, which is the case for document images in general. For a random 15% of the word images, we induced degradation by eroding the images using the `erode()` function available in OpenCV [245].

In the next step, all the word images are pasted onto a background image in the same order as in the original passage using alpha blending. The background image is randomly sampled from a set of 20+ manuscripts like textures downloaded from the internet. While pasting the words onto the background image, we break the lines after every few words. The number of words in a line is varied randomly, with a minimum of 5 and a maximum of 7. For a document, the fixed inter-word spacing is the average per-character width (width of a word image divided by the number of characters in it) of the document. Similarly, the fixed inter-line spacing used for a document is the average height of word images. For a random 15%, we used a different spacing than the fixed value. In such cases, we multiply the fixed value by a random multiplier in the range of 0.9–2.5 for inter-word spacing and a random multiplier in the range of 0.9–1.3 for inter-line spacing. The borders on all 4 sides of the document are also set as a function of

³<https://fonts.google.com/>

the inter-word spacing used in the document. The border on each side is the fixed inter-word spacing for the document multiplied by a random value in the range 1.5–5.

Once all word images were pasted onto the background image, we subjected the image to a random amount of skew. The skew angle we used is a random value in the range of -5–5. Then a random 15% of the images are re-sampled and brought back to the original dimensions to induce artifacts resulting from the re-sampling of images. For re-sampling, we used a random scaling factor in the 0.6–1.4 range.

6.4.2 BenthamQA

Compared to HW-SQuAD which contains synthetically rendered images, BenthamQA contains real historical manuscripts. We decided to use a manuscript collection with word and line bounding box annotation already available so that QA methods that require gold-standard annotations for words or lines can use them. We used manuscripts from the ImageCLEF 2016 Bentham Handwritten Retrieval dataset [246]’s dev split. The manuscripts are at least 200 years old. At first, duplicate and near duplicate images were removed by comparing the images using Jaccard Similarity of their gold-standard transcriptions. 338 images were left after the de-duplicaton.

To annotate QAs on the manuscripts, we used a hosted version of the Haystack annotation tool⁴. The gold-standard transcriptions of the manuscripts were shown to the annotators to add questions and answers. Similar to the typical extractive QA annotation, the answers were marked as spans of the transcripts. Once the textual spans are annotated, the equivalent visual region on the document image that encompass the textual answer are marked. Since we know which words constitute the textual answer and the bounding boxes of these words are known (since we have ground truth for locations and transcriptions of the words in the manuscripts), marking the answers’ visual regions can be done without any manual effort. With the help of two volunteers, 200 question–answer pairs based on 94 different document images were collected. There were 224 images that were not annotated with any question. We keep these images also in the dataset as distractors for the QA task. In other words, we have 200 questions on 338 images and 94 images have at least one question based on them.

Basic statistics of questions, answers and documents in both HW-SQuAD and BenthamQA are shown in Table 6.1. Questions in BenthamQA are almost double the length of questions in HW-SQuAD. This is expected since the SQuAD1.0 dataset was annotated for MRC, where each question is asked in the context of a single passage. A detailed discussion on this aspect is given in section 6.3, when we analyze the performance based on the length of the questions.

⁴<https://annotate.deepset.ai>

While annotating BenthamQA, annotators were instructed not to ask questions specific to a particular document, such as *What is the title of the document?* or *What is the name of the city mentioned in the first paragraph?* This made the annotators ask longer questions. For example, take this question; *In one of the notes where the author talks about “Genus” and “Species” in the context of law, he makes a comparison with a flower. Which flower is it?’. The question gives context that helps find the target manuscript from the collection and then talks about the context within the manuscript. On average, document images in BenthamQA have more words in them since these documents are manuscript pages compared to passages rendered as document images in HW-SQuAD. We show two documents each from HW-SQuAD and BenthamQA in Figure 6.3.*

6.4.3 Performance evaluation

QA and MRC benchmarks use evaluation metrics like F1 score and Exact Match [91, 138]. On the other hand, standard VQA benchmarks, particularly those that do not involve text, use a metric called VQA Accuracy. It treats a predicted answer as correct if at least 3 of the 10 annotators entered the exact same answer [55]. ST-VQA and DocVQA use ANLS subsection 4.4.1. Since all these metrics are conceptualized for textual answers, they cannot be used for a QA task where answers are document snippets. Hence a new evaluation scheme is devised where the predicted answer snippet is evaluated against a ground truth defined in the image space.

Wang et al. [238] propose a VQA task where models need to predict a textual answer and a bounding box as ‘evidence’. The ground truth ‘evidence’ is used to assess the predicted ‘evidence’. Inspired by evaluation in object detection benchmarks, the authors propose to use Intersection Over Union (IOU) as a measure to evaluate the predicted box. We considered a similar evaluation scheme by treating the minimum rectangle bounding the answer words as the ground truth (GT) box. Methods that work without recognizing the text in the documents, such as the method we propose in this work, do not understand the semantics of the text and must be relying on the similarity between words in question and the context where the answer lies. It is unlikely that such approaches can narrow down to a tight box around the answer words. This would result in low IOU scores when pitted against a GT box which tightly bounds the answer. An alternative is to check if the GT box is inside the predicted snippet. However, a method can always output bigger snippets that contain the answer words and satisfy the condition. This motivated us to develop an evaluation scheme that accommodates methods that return answers as short snippets containing the exact answer and some context. At the same time, we need to make sure that larger snippets with too much context are discouraged.

Q: Telnet used what interface technology?

A1: host interface to X.25 and terminal interface to X.29

A1: X.25

A1: ARPANET

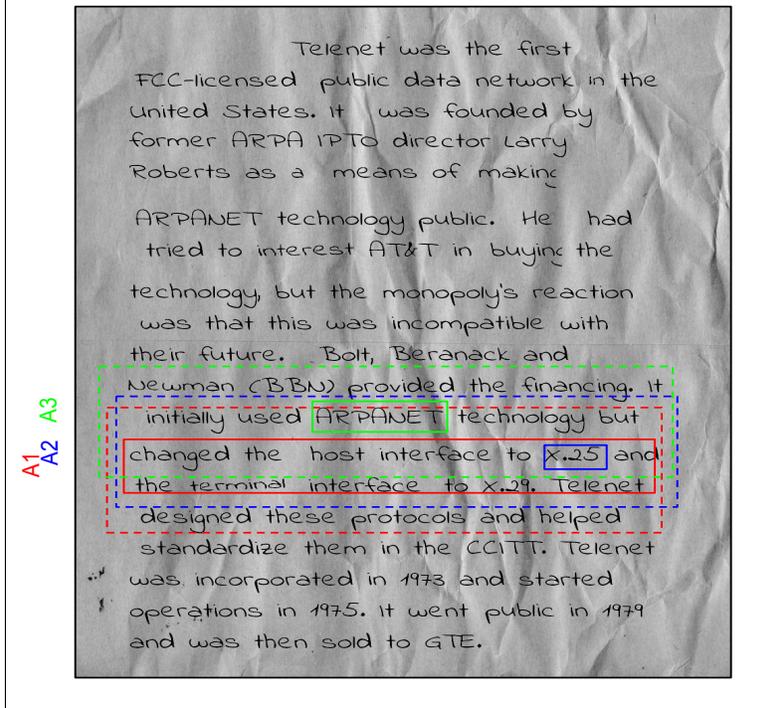


Figure 6.4: A question–answer pair and the associated document from the test set of HW-SQuAD dataset. To evaluate answer snippet extraction, we define ground truth in the image space in terms of a Large Box (LB), which includes the lines where the answer is, plus a line above and below it and a Small Box (SB), which tightly encloses the words constituting the answer. In this case, there are three ground truth answers. For visualization, they are denoted as A1, A2, and A3 and drawn using red, blue and green colours, respectively. Note: the margins of the boxes shown in the figure are changed a bit to show the overlapping boxes without clutter.

We treat the predicted answer snippet (or Answer Box; AB for short) as a correct prediction if: (i) AB is inside a Large Box (LB), which is the region in the image enclosing the text lines where the words making up the ground truth answer lie and one line above and one line below it and (ii) the smallest box containing the words constituting the textual answer (i.e., a Small Box or SB) is inside the AB . In short, the condition for an answer snippet to be a correct prediction is a double inclusion criterion: $AB \subset LB \wedge SB \subset AB$. To understand this better,

we show LB and SB for an example case in Figure 6.4. To incorporate the double inclusion criterion into a score, we define Double Inclusion Score (DIS) as:

$$DIS = \frac{AB \cap SB}{|SB|} \times \frac{AB \cap LB}{|AB|}$$

For a question, if the predicted answer snippet has a $DIS > 0.8$ with any of the ground truth answers, we consider it a correct prediction. Accuracy is then calculated as the percentage of questions for which we have a correct prediction. Note that the HW-SQuAD and BenthamQA datasets can also be used for recognition-based QA where textual answers are expected. In such a setting, it is ideal to use the standard metrics for text-based QA, such as Exact Match (EM) score and F1 score [91]. In subsection 6.5.4, we report results on both the datasets in the recognition-based setting where textual answers are expected.

6.5 Experiments and results

In this section, we present experimental settings, the results of our experiments, and a discussion of the results.

6.5.1 Implementation details

Questions are split into constituent words, and stop words are removed using NLTK [247]. Words from the documents are extracted using gold-standard word bounding boxes available with the datasets. In practice, word bounding boxes are obtained using a document segmentation algorithm, and the segmentations need not be perfect. Hence results we report below do not factor in possible segmentation errors in such scenarios. We train a Convolutional Neural Network (CNN)-based binary classifier, to filter out word images of stop words. The feature extraction block of this network follows the same architecture as the one in CRNN [46] network. Feature maps after the last convolutional layer are mapped to a binary classification layer. The classifier is trained using synthetic, handwritten word images in the HW-SYNTH dataset [203] and word images from the IAM dataset [248].

After stop words removal, words from both the questions and documents are fed to the end-to-end embedding network to represent them as word embeddings. We use a pretrained model⁵ for this network, made publicly available by the authors. The pretrained model is trained on HW-SYNTH and IAM datasets. The network output is an L2 normalized embedding vector $x \in \mathbb{R}^{2048}$. We use Principal Component Analysis (PCA) [249] to map x to a $x' \in \mathbb{R}^{Dw}$ in

⁵<https://github.com/kris314/e2eEmbed>

Embedding size (D_w)	HW-SQuAD		BenthamQA
	val	test	
2048	26.0	22.1	23.5
200	21.8	-	-
128	20.5	-	-
64	16.3	-	-

Table 6.2: Top-5 accuracy for the document retriever when the individual embeddings are added together to obtain document and question representations. That is the aggregation scheme used is SUM.

a lower dimensional space. PCA rotation matrices are learnt on the training set of the HW-SQuAD. In our experiments, particularly while using the FV aggregation scheme, we work with these dimensionality reduced vectors.

Parameters of the GMMs are estimated offline using embeddings of around 1.2 million words (after filtering stop words) drawn from documents in the train set of HW-SQuAD. Since words in questions are almost a subset of the words in the documents, and word images add more variability than textual words, we use only words from documents to fit the GMMs. To get candidate snippets from the document proposals generated by the document retriever (see subsection 6.3.4), we use sliding windows of size 2 and a step size of 1 over the text lines, scanning from top to bottom. For example, from a document proposal having l text lines, we have $l - 1$ candidate snippets, each one enclosing two contiguous text lines.

6.5.2 Performance of document retriever

Since the performance of the answer snippet extraction step (see subsection 6.3.4) is dependent on the quality of the document proposals generated by the document retriever (see subsection 6.3.3), we first evaluate the document retriever. Top- N accuracy—the percentage of questions for which the target document is in top N proposals—is used as the evaluation metric. Table 6.2 reports top-5 accuracy for the proposals retrieval on both the datasets for different word embedding sizes (D_w) when the aggregation scheme used is SUM. We appreciate that the original embedding vector (without any dimensionality reduction) of size 2048 yields the best results in this case.

		HW-SQuAD						BenthamQA
		val			test			
D_w	K	32	64	128	256	512	512	512
	<i>Gradients w.r.t to both mean and SD are considered i.e. $FV = \mathcal{G}_{\mu,i}^\lambda s$ and $\mathcal{G}_{\sigma,i}^\lambda s$; $FV = 2KD_w$</i>							
	64	14.6	15.6	16.2	-	-	-	-
	128	10.8	14.4	17.5	-	-	-	-
	200	6.8	12.9	15.4	-	-	-	-
<i>Only gradients w.r.t to mean are considered i.e - $FV = \mathcal{G}_{\mu,i}^\lambda s$; $FV = KD_w$</i>								
	64	31.9	39.5	40.9	44.0	48.2	-	-
	128	36.6	43.2	45.7	48.3	50.4	46.5	55.5
	200	30.1	40.9	42.0	44.7	46.5	-	-

Table 6.3: **Top-5 accuracy (%) for the document retriever when FV aggregation is used.** The top half shows the results when FV is a concatenation of gradients with respect to means ($\mathcal{G}_{\mu,i}^\lambda$) and standard deviations of the Gaussians. ($\mathcal{G}_{\sigma,i}^\lambda$). The bottom half shows results when FV formed is a concatenation of gradients with respect to means alone. In all the cases the resulting FVs are normalized using power normalization and L2 normalization.

The performance of the document retriever using FV aggregation for different word embedding sizes and the number of Gaussians in GMM (K) is shown in Table 6.3. The top half of the table shows results when the FV is a concatenation of gradients with respect to both the mean ($\mathcal{G}_{\mu,i}^\lambda$) and standard deviation ($\mathcal{G}_{\sigma,i}^\lambda$) of the Gaussians. Thus the size of the resulting FV ($|FV|$) is $2KD_w$, since each gradient vector is of size D_w (see subsection 6.3.2.2). We repeat the same set of experiments using FVs where only the gradients with respect to the means of the Gaussians are concatenated to form the FV. These results are shown in the bottom half of the table. This setting yields much better performance, consistent across different values of K and D_w . It also helps in faster computations since $|FV|$ becomes half in this case. For this reason, experiments considering gradients with respect to both the mean and standard deviation of Gaussians are not conducted for higher values of K . On the test set of the HW-SQuAD and on BenthamQA, we report results only for the setting which yields the best performance on the validation (val) set of HW-SQuAD.

The best performance on the validation set is for $K = 512$ and dimensionality-reduced word embeddings of size $D_w = 128$. Embeddings of size 128 give better results compared to

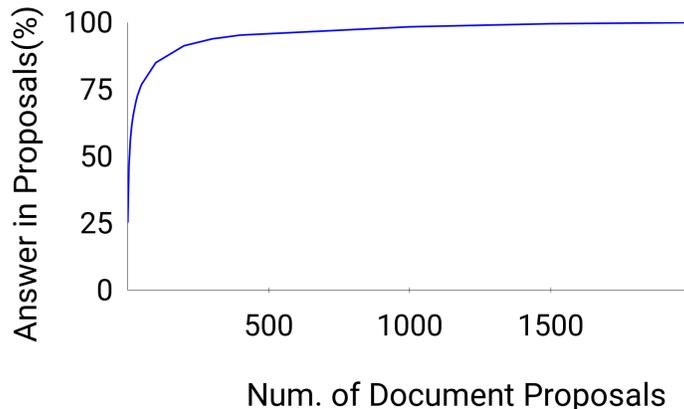


Figure 6.5: **Number of document proposals (N) vs top-N accuracy.** The higher the number of document proposals for the target document, the better the chance of finding the answer to a question in the proposals. When the number of proposals is more than 2000, the answer is found within the proposals for almost all the questions. The given plot is for the validation set of HW-SQuAD.

a larger embedding size of 200 for all values of K . This is in line with the finding in [244] that dimensionality reduction of the local features using PCA improves the overall performance of classification and retrieval. In [244] the authors cite two possible reasons for this: (i) decorrelated data can be fitted more accurately when the covariance matrices of the GMM are diagonal, and (ii) GMM estimation is less noisy if only the ‘stronger’ components are considered. In Table 6.3, it can be seen that the larger the number of Gaussians K , the better the retrieval performance. However, for higher values of K dimensionality of the FV becomes very large and it is computationally expensive to work with such large vectors in practical scenarios.

Although we apply L2 normalization and power normalization to the fisher vectors as a default choice, we study how power normalization helps in improving the results. Figure 6.6 depicts the effect of power normalization of the FVs for the document retriever step. The chart plots top-5 accuracy for the document retriever on the HW-SQuAD validation set as a function of $|FV|$. We use dimensionality-reduced embeddings of size $D_w = 128$ for different values of K . It can be seen that power normalization of the FVs consistently improves performance. The impact is more prominent as $|FV|$ increases. This trend is in line with the observation in previous studies [244, 242] that power normalization helps in improving the quality of fisher vectors for classification/retrieval problems.

Figure 6.5 shows how the number of questions for which the target document is found within the proposals grows as a function of the number of proposals. With the number of proposals > 2000 , every question’s target document is found among the proposals.

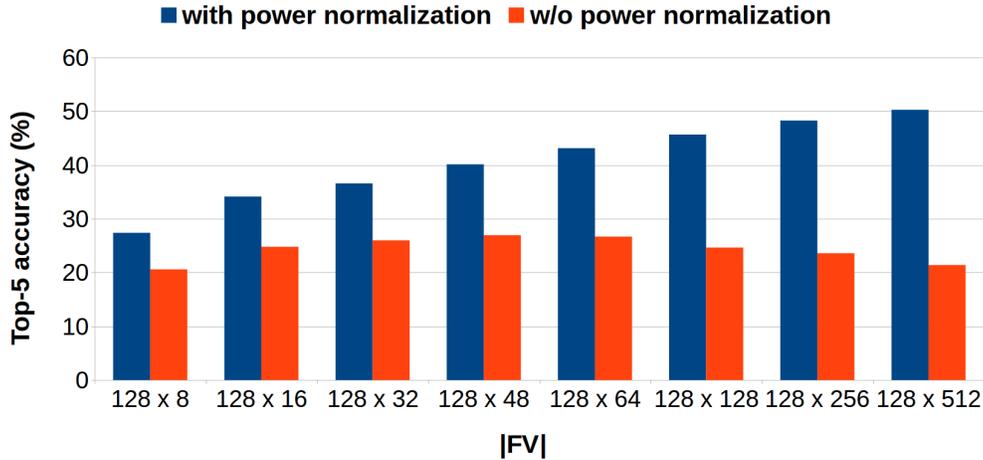


Figure 6.6: **Top-5 accuracy for document retriever for different sizes of FV , with and without power normalization.** Embedding size used in this experiment, $D_w = 128$. Power normalization yields improved performance for FV of all sizes, but the effect is prominent for larger FVs.

6.5.3 Evaluating end-to-end answer snippet extraction

Next, we evaluate our full pipeline, which retrieves relevant documents first and then extracts the answer snippet. Qualitative results of snippet extraction using the full pipeline are shown in Figure 6.7.

Table 6.4 shows results for different aggregation schemes used for the snippet extraction step with our best setting for document retriever. In addition to the accuracy measure using DIS, we use a metric that measures the F1 score at the level of text lines.

Since our approach generates candidate snippets as a set of two consecutive text lines, we believe it is meaningful to define a metric at the level of text lines. We define Precision (P) as the ratio of the number of text lines common to the answer snippet and the target answer to the number of text lines in the answer snippet. Recall (R) is the ratio of the number of lines common to the answer snippet and target answer to the number of text lines in the target answer. F1 is computed as the Harmonic mean of P and R. This metric makes sense only when the answer snippets are defined in terms of text lines. This is the reason why we do not define this metric in subsection 6.4.3, where we define a generic evaluation scheme for the proposed QA task where answers are document snippets.

We evaluate performance on the test set of HW-SQuAD when the target document for each question is made available to the answer extraction step directly, bypassing the proposal generation step. This is equivalent to an MRC task since the document containing the answer

Aggregation scheme for the answer extraction step	HW-SQuAD				BenthamQA	
	val		test		Acc.	F1(lines)
	Acc.	F1(lines)	Acc.	F1(lines)		
FV ($D_w = 128, K = 64$)	15.1	12.5	14.44	11.7	11.0	10.4
FV ($D_w = 128, K = 512$)	10.6	9.9	-	-	-	-
SUM ($D_w = D_o = 2048$)	15.9	12.7	15.9	12.7	17.5	15.1

Table 6.4: **Performance evaluation of end-to-end answer snippet extraction.** In all cases, the document retriever generates 5 document proposals. The aggregation scheme used for document retriever is FV ($D_w = 128, K = 512$). Acc. stands for the snippet extraction accuracy in percentage based on the DIS score proposed in subsection 6.4.3 for a DIS threshold of 0.8. F1 (lines) is the F1 score in percentage for the text lines.

is given, and the only job left is to extract the answer snippet. In this case, using the SUM aggregation, answers are found with an accuracy of 43%.

Figure 6.8 shows how the performance varies when the number of proposals generated varies. When only one document proposal is generated, snippet extraction accuracy is 7.6%, and it goes up to 16.72% when the number of proposals is 25 and then starts dropping. When all the documents in the test set are considered as proposals for every question in the set, the accuracy is 12.7%.

Chen et al. [138] note that open-domain QA performance on the SQuAD1.0 dataset is significantly affected by the nature of questions. The dataset was originally created for MRC and the questions are asked in the context of a short paragraph. Information sought in the question is ambiguous in many cases if the context is not given. For example, questions like *What day was the game played on?* or *Why did he walk?* which are part of the development set of SQuAD1.0 cannot be answered unless the associated passage or document is given. In order to study the impact of such questions on QA on HW-SQuAD, we analyze in Figure 6.9 how the performance varies when the question length varies (question length is measured in terms of the number of non-stop words). The trend seen in the plot is that performance improves with more words in the question. This could mean that, with more words in the question, questions become less ambiguous and hence easier to locate the right answer snippet.

In our experiments, two aspects make the experimental setting (see subsection 6.5.1) advantageous for HW-SQuAD, compared to the BenthamQA dataset. Firstly, the off-the-shelf, end-to-end embedding model, which we use for word embeddings, is trained on the HW-SYNTH dataset, whose synthetic word images are similar to the word images in the HW-

When is the oldest recorded incident of civil disobedience?

civil disobedience took place during the Roman Empire[citation needed]. Unarmed Jews gathered in the

Why was the Merit network formed in Michigan?

Michigan's public universities, was formed in 1966 as the Michigan Educational Research

In every quasi jury, how many classes of quasi jurors are there?

Art. 1. In every Quasi Jury are two classes of Quasi Jurors the ordinary, and the select.

What caused scarcity of wheat in public stores towards the end of the year 1799?

"The scarcity which they themselves had occasioned, had raised the price of wheat to £1.10 per bushel; a shameful extortion!"

Figure 6.7: **Qualitative Results of end-to-end answer snippet extraction.** The first row shows a success and failure case from the test set of HW-SQuAD. In the case of the second question (the one on the right) although the extracted snippet talks about the formation of the Merit network, it does not answer the question. The second row shows examples from the BenthamQA dataset. The left one is a success. In the case of the right one, our model returns a snippet talking about wheat scarcity, but it does not say anything about the reason for the scarcity.

SQuAD dataset. Both HW-SYNTH and HW-SQUAD use handwritten fonts to render text and there could be many common fonts between the two datasets. Secondly, PCA rotation matrices and GMMs used for the FV aggregation scheme are learnt on the train split of the HW-SQuAD dataset alone. Despite the disadvantages, the document retriever and the answer extraction steps perform better on BenthamQA. We attribute this to i) the smaller size of the dataset, which makes the proposal generation easier with a lesser number of distractors, and ii) longer, unambiguous questions in BenthamQA, as discussed in subsection 6.4.2.

6.5.4 Evaluating a recognition-based QA approach on HW-SQuAD and BenthamQA

To study how recognition-based QA models work on HW-SQuAD and BenthamQA, we evaluate a full pipeline QA framework on the text transcriptions of the datasets. In NLP/IR domain, the standard approach to QA over a document collection is to use a document retriever first and then use a document reader that extracts final answers. The document retriever is

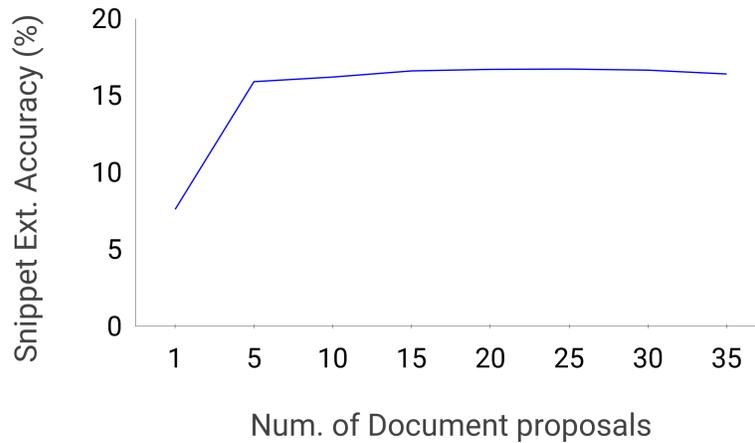


Figure 6.8: **Snippet accuracy vs number of document proposals.** The plot shows the performance of answer snippet extraction step on the test set of HW-SQuAD as a function of the number of document proposals generated. The performance flattens when the number of proposals > 20 , suggesting that with more proposals, answer extraction becomes difficult.

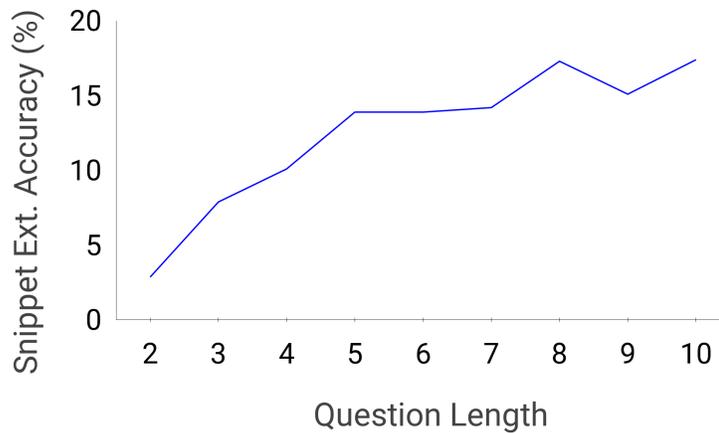


Figure 6.9: **Snippet accuracy vs question length.** Longer questions are likely to carry more information helpful in finding the right answer snippet when there are multiple candidates with similar content. Better performance for longer questions substantiate this.

similar in function to the one we use, i.e., to generate document proposals for the document which is likely to contain the answer. The document reader— typically an extractive QA model—extracts an answer using the document proposals as its context. In our experiments,

we use a TF-IDF-based document retriever, the same as the one used by Chen et al. [138] in their work on recognition-based full pipeline QA. They observe that a simple TF-IDF-based retriever performs better than other similar models such as, Elasticsearch [250]. For the document reader part, we use a BERT [113] QA model. The specific model we use is a BERT_{LARGE} model finetuned for QA on the SQuAD1.0 dataset. The exact model name in Transformers model zoo⁶ is ‘bert-large-uncased-whole-word-masking-finetuned-squad’. Note that HW-SQuAD, which is a derivative of SQuAD1.0, has little domain gap with HW-SQuAD data, provided the text transcriptions are good.

OCR	Training data	Word accuracy(%)	
		HW-SQuAD test	BenthamQA
SynthIam	HW-SYNTH [203] + IAM [248]	44.87	13.72
SynthIamHwSq	HW-SYNTH + IAM + HW-SQuAD	97.85	23.17

Table 6.5: Performance of a CRNN [46] OCR model on the documents in the newly introduced datasets. Without in-domain training data, the OCR transcriptions are quite noisy. Compared to HW-SQuAD, performance on BenthamQA is much worse.

6.5.4.1 Transcribing handwritten images using a CRNN OCR

To recognize text in the handwritten documents in HW-SQuAD and BenthamQA, we used a CRNN [46] OCR model. By the OCR model, we refer to a word recognition model, and it is assumed that segmented word images are available. Since the gold-standard word bounding boxes are available for both the datasets, we use them directly to crop the word images. We trained two OCRs: i) SynthIam: trained on 9 million synthetic handwritten word images in HW-SYNTH [203] and the train split of IAM dataset with real handwritten images, and ii) SynthIamHwSq: trained on HW-SYNTH, IAM train split and train split of HW-SQuAD. Except for different training data, both the OCRs are the same in terms of network architecture and training setting.

In Table 6.5, we show the performance of both the OCRs on the test set of HW-SQuAD and BenthamQA. We report word accuracy, which is the percentage of words for which the recognized text matches exactly with the ground truth. The 9 million word images in HW-SYNTH are synthetically generated using handwriting fonts, it is likely that a lot of these fonts are used to render document images in HW-SQuAD. This must be the reason why SynthIam performs

⁶https://huggingface.co/transformers/pretrained_models.html

significantly better on HW-SQuAD compared to its performance on BenthamQA. On BenthamQA, although both OCRs yield low word accuracy, SynthIamHwSq is significantly better than SynthIam. We believe better performance with SynthIamHwSq is primarily due to the paper/manuscript style backgrounds used while rendering document images in HW-SQuAD. This is in contrast to IAM and HW-SYNTH, where backgrounds are solid colors. Some of the background images and textures used for HW-SQuAD resemble the background of historical manuscripts. The diversity in backgrounds and textures for pages in HW-SQuAD must have helped the SynthIamHwSq model to better recognize the word images in BenthamQA.

6.5.4.2 Evaluating TF-IDF-based document retriever

Table 6.6 shows the performance of the TF-IDF-based document retriever for different types of transcriptions. The results are directly comparable with our recognition-free document retriever (subsection 6.3.3) since the task and evaluation metrics are the same. As expected, with higher OCR accuracy, the recognition-based approach is better at retrieving documents. However, when OCR transcriptions are noisy, the proposed recognition-free approach is better at retrieving documents. In the case of BenthamQA, TF-IDF-based retrieval on transcriptions from SynthIamHwSq results in a top-5 accuracy of 32%. At the same time, the recognition-free approach using FV yields a top-5 accuracy of 55.5%, as reported in Table 6.3.

6.5.4.3 Evaluating TF-IDF + BERT full pipeline QA framework

We build a text-based QA pipeline, which uses the TF-IDF-based document retriever to generate document proposals and a BERT_{LARGE}-based extractive QA model. We discuss BERT QA model in subsubsection 2.2.1.1. We use the BERT_{LARGE} architecture. The results are shown in Table 6.7. In the table, “F1” stands for the F1 score which is the most commonly used evaluation metric for QA and MRC benchmarks like SQuAD1.0. It measures the average

Transcriptions	top-5 Accuracy(%)	
	HW-SQuAD test	BenthamQA
SynthIam	47.93	15.76
SynthIamHwSq	86.10	32.00
Gold-standard	90.2	98.5

Table 6.6: Results of a TF-IDF-based document retriever on transcriptions of the documents in HW-SQuAD and BenthamQA.

Transcription	HW-SQuAD test			BenthamQA		
	F1	EM	Snippet Acc.	F1	EM	Snippet Acc.
SynthIam	23.05	13.54	19.13	1.05	0.34	0.82
SynthIamHwSq	65.24	55.31	59.26	3.00	1.5	2.46
Gold-standard	76.82	70.73	74.76	78.41	66.00	72.85

Table 6.7: **Results of text-based full pipeline QA model.** F1 score, Exact Match (EM) percentage and Snippet Accuracy (Snippet Acc.) for full pipeline QA using the transcribed text from the documents in the datasets. F1 is the standard evaluation metric used for evaluating extractive QA. For calculating Snippet Acc., we map the textual answers to document snippets and evaluate them using the DIS score proposed in subsection 6.4.3.

overlap between the predicted answer and the ground truth answer [91]. Do not confuse this F1 score, with the F1 score for text lines we compute to evaluate snippet extraction performance in Table 6.4. Exact Match (EM) is the percentage of questions for which an answer matches exactly with at least one of the ground truth answers. The predicted answers are mapped to a corresponding document snippet and evaluated for snippet extraction performance as well. Since the predicted answers are extracted as spans from the given documents, we take the text lines where the span lies as the answer snippet and evaluate snippet accuracy by calculating the DIS score as given in subsection 6.4.3. These numbers are shown under “Snippet Acc.” in the table. Comparing the snippet extraction accuracy in Table 6.4 with the Snippet Acc. in Table 6.7, it is evident that the proposed recognition-free approach works better when robust OCR transcriptions are not available.

Results in Table 6.7 suggest that highly noisy OCR transcriptions lead to sub-standard QA performance even when state-of-the-art NLP/IR QA models are used. Although some works address recognition-based document retrieval in the presence of noisy OCR outputs [251, 252], we do not explore this direction since our approach focuses on QA without explicit OCR.

6.6 Summary

In this chapter, we have introduced the problem of QA on handwritten document collections and presented two new datasets — HW-SQuAD and BenthamQA. The problem has proven to be challenging, and we expect it to attract the interest of the research community. We also have presented a new method for QA that is segmentation-based, recognition-free, and lexicon free.

We have presented extensive experiments with different aggregation schemes that allow us to find answers as nearest neighbours of a given question representation in the space of possible answer snippets. The results demonstrate that the proposed solution would perform better than recognition-based document retrieval and QA models that use noisy OCR transcriptions.

Chapter 7

Conclusion and future work

In this thesis, we presented our works dealing with machine-understanding of document images. We believe our works expand the applicability of DIA by expanding the language coverage and launching new lines of inquiry that call for a holistic, human-like understanding of document images. In this chapter, we present a summary of the thesis reiterating our major contributions and list potential directions for future research.

7.1 Conclusions

Works presented in this thesis contribute to DIA on two fronts. It makes DIA more universal by developing OCRs for low-resource Indic languages and releasing the largest ever public benchmarks for these languages. Secondly, different from the bottom-up approaches in DIA that are focused on machine reading, we motivate a purpose-driven DIA with machine-understanding as the core objective.

In section 1.5, we listed the slower development of DIA technologies for non-Latin-based languages as one of the hurdles to the development of universal DIA solutions. We believe our work on printed text recognition of 13 Indic languages (chapter 3) contributes to the development of robust DIA solutions in Indic languages. The combined number of speakers of Indic languages is more than 1.5 billion. Developing robust OCRs and other DIA solutions for these languages can contribute significantly to streamlining office work in governments and businesses, preserving the the cultural heritage of the region, and most importantly improving the AI-readiness of data in these languages.

In chapter 4 we introduced the novel problem of VQA on business documents called DocVQA. We showed that neither models for VQA on natural images nor MRC models work well for VQA on document images. DocVQA has emerged as one of the most popular tasks in

the DIA space, attracting thousands of downloads for the dataset and models developed for it ¹. Followed by DocVQA, we introduced a more challenging document images VQA task on infographics called InfographicVQA (chapter 5). Infographics are rich in the interplay of visual, textual and layout elements, and machine understanding of infographics requires strong multi-modal reasoning capabilities. This is validated by the poor performance of existing vision and language models on the InfographicVQA dataset. The results also highlight the need for better visual encoders for images like infographics that are quite different from natural images and generic document images. We conducted open challenges on both DocVQA and InfographicVQA [219, 253]. The challenges remain open for submissions and the two datasets have evolved as an important benchmark not only in DIA but also in Vision and language research space.

In section 1.5, while discussing challenges to DIA in general, we note that text recognition from handwritten and/or historical documents remains an open problem. Often, these document images are part of a collection that has historical or cultural significance. To be able to query these collections and extract information is an important requirement. However, poor text recognition performance remains a bottleneck. To motivate study on information extraction from such collections, we propose a new QA task that returns answers as image snippets. We presented this work in chapter 6. Since recognition-based approaches are ineffective when OCR transcriptions are noisy, we propose a recognition-free approach that relies on cross-model retrieval to retrieve a document snippet that best matches the question. Although the retrieval is based on the lexical similarity between the question and the snippet—and does not take into account the semantics—, the proposed approach works better than SoTA recognition-based models when the OCR transcription is noisy.

7.2 Future directions

The following are some of the prospective directions that the thesis opens up:

- **AI-readiness of low-resource languages** OCR technology for non-Latin scripts have been gaining momentum since the last decade thanks to deep learning-based text recognition models that are largely language-agnostic. Our work presented in chapter 3 validates the robustness of CTC-based transcription for segmentation-free OCR of 12 Indic languages. However, there are many more Indic languages that lack functional OCRs and text recognition benchmark datasets. Although the script is different for most of the

¹https://huggingface.co/models?pipeline_tag=document-question-answering

Indic languages, the languages share many linguistic properties. We believe the common linguistic heritage can be exploited to build OCRs for these languages. Rather than building large-scale annotated datasets for each of these languages, cross-lingual domain adaptation and use of synthetic datasets need to be explored.

- **Multi-page DocVQA.** In both DocVQA and InfographicVQA, an answer for a question needs to be found from a single page or image. A natural extension of the single-page VQA is the multi-page VQA. This is similar to open-domain QA [50] in NLP. Multi-page VQA would not be as trivial as recognizing text on these images and then employing an existing open domain QA model. We are talking about a collection of multimodal document images, i.e., document images that feature not just text but tables, visualizations, figures and forms. One of the challenges in creating such a dataset is the difficulty in annotating question answers on document collections. One possibility is to take existing questions in DocVQA and pose them over the whole set of document images in the dataset. However, many of the questions in DocVQA are ambiguous in the context of the whole images in the dataset. For example, questions in DocVQA such as "What is the title of the document?" or "Who sent this letter?" make little sense when asked on a document collection that contains many pages with titles and many letters. A possible way to annotate a multi-page VQA dataset is to follow the way MS-MARCO [134] is annotated. That is, to collect a set of questions that users ask over a document collection and then employ workers to find and mark answers for each question in one of the pages (or images) in the collection or flag the question as unanswerable.
- **Generative document VQA.** Another important direction to pursue in both single document and multi-document VQA is generative/abstractive VQA—a VQA task where answers are not extracted—from text present in the images—but generated. In InfographicVQA, there are some questions that require generated answers, such as questions that require counting and arithmetic operations. Abstractive answers in InfographicVQA are always short, numerical answers. Recent works such as VisualMRC [221] deal with abstractive VQA on single document images. However, document images in VisualMRC are born-digital, Wikipedia screenshots. They are significantly different from the majority of the documents that are created and disseminated in real life, such as letters, reports, newspapers and periodicals.

Appendix A

Appendix to chapter 3

A.1 Examples of pages in the internal page OCR dataset

In chapter 3, in order to conduct page-level OCR experiments an internal dataset that comprises 1000 pages per language is used (see subsection 3.3.1). This dataset is a subset of the Indic consortium dataset. In figures Figure A.1– Figure A.4, we show examples of pages in the internal dataset.

A.2 Examples of line images in the Mozhi dataset

In chapter 3, we introduced a new public dataset called *Mozhi* (see subsection 3.3.2) for line-level and word-level text recognition in 13 Indian languages. In figures Figure A.5– Figure A.7, we show examples of line images in the dataset.

A.3 Distribution of lengths of words in the Mozhi dataset

In figures Figure A.8– Figure A.12, we show the distribution of word lengths for all languages in the *Mozhi* dataset.

৪৬ জনা সৃষ্টি অনা পুৰস

কিঃ মংগলবোৰে মংগলবোৰে বেখোন অফলা ভাণিহে? গানে শ্ৰোতা মুখখন মেৰি হাঁহি এটি উলিয়াই বোলে, 'কথা শুনাও। আমি ইয়াতে এখন মংগলবৰীয়া হাট পাহিৰে মুছিয়ে।' গাঁৱৰ বাইৰখন সুবিধা হয়— সপ্তাহে দুবাৰ আশুতো, কচুতো বেচি কেলকণ চেনেকৰ কিনিব পাৰে, যোৰহাটলৈ বাটকুৰি বাই নগালেও হয়।

এইবোৰ কেইবা সপ্তাহো গ'ল। মংগলবোৰে মংগলবোৰে মনুৰে সমাগম বাঢ়িহে আহিছে— হাটলৈ হ'লে নাম-গোছেই নাই। বাবুতো সোদায় মংগলবোৰ আহে, সোদায়োম কথা-নকতা হয়। হেঁচাপো, হাঁহিউকনি, থকাথুকাও হয়। দুৰ-বিদূৰৰ গাঁৱৰ মনুৰে বাস্তৱ কোমল চাউল ওৰৰ টোপোনা লৈ বাটকুৰি বাই আহে। মোৰ সন্দেহ আকৌ যুগুৰনি লাগিল— বোলাৰ কথাটো কিঃ

বৰশীউৰি দুটা কোমল মাটিত টানকৈ গোল মাৰি থৈ এদিন মনৰ উত্থুৰুনিৰে থাকিব নোৱাৰি গ'লো আশুৰি মনুৰে অহালৈকে উৰুই গছৰ ওপটো লাগিল। গৌৰত বাবুজনে হাতমেনাৰ পৰা উলিয়াই সোণাবীন তৰুই ভূমি কিবা নিজে আক মনুৰোৰে ভূমি বাকি মংগলবৰীয়া চলচলীয়া ৰূপৰ টকা দি বেৰাহোৰে কিনিছে। দুবুৰে বুজিলো আশীম অৰ্থাৎ কনিৰ বানসায়। গামক ঘূৰি আহোঁতে লগ পাই আকৌ মুছিলো, বোলা 'মংগলবোৰে মংগলবোৰে অহা এই বাবুজন আশীম কোচা চকৰবীৰি বিয়া মেৰি?' গানে বোলে, 'হেঁতৰে মনুৰোৰে বৰ কথা সুনি থাকুৰক। তই বৰশী বোবা মনুৰে, বৰশী বাই থাক। এইবোৰ কথা কেতেই? মই ইতিহাসখন চুপে-খুপাতে তামলে অকল উলিয়াই দিলত পিছে গামক মুখৰ হাঁহি বিৰিঙন। মুখৰ আঁধৰুটো দি কথা মুছিলে ধৰিলে। ক'লে বোলে, 'তই বেখোন জানিবিহে।' কৈয়ে পোহৰ দে।

'মংগলবৰীয়া আশীম বোবা মংগলবোৰক আঁজিকালি আশীম, কিনিব নোৱাৰা হ'ল— জাহে জাহে ভলশিমাৰে নিজেই কনি থাকে, হাতে-ভৰিৰে ধৰি বাইজক আশীম নাথালো কাহুতি কৰে, সেইবোৰো কামসৰ দিদি বুজবাজ নলে আশুৰি গ'লে মানপিতা কৰে। পুণ্ডিত ভলশিমাৰেবক হাজোতৰ ভাষাটো— পিছে ইহঁত উপক্ৰমৰ গুৰি নিটো নটো এজাক ওলায়। গাৰ্হাৰজন শিখ বাটগী পিন্ধা এই ভলশিমাৰেবোৰে আমাৰ বাইজপলক হ'লে বৰ বিপদত পেলাইছে সেই।

আজি এৰাহমান আপুৰে কুপ্পনীয়ে কনিৰ মহলাবোৰৰ দৰ্ভাত তলা লগালে— বোলে কনি খোৰাটো সোৱৰ কথা। ইফালে আমাৰ কাঠিচোকো বাইজক মৰৰ মিলিল। মই বোলা কুপ্পনীয়ে বাইজপলক এটাকৈ গছত তুলি ওৰি কটো কিয়? কিবা বোলানে তোমালোকেই বহুতোমা? আমি উলিয়ালাই— বাইজক মাজত চলিল। নহ'লে আমাৰ বাইজে এইমোৰা বহুৰ সোৱাল পালেহেঁতেন কথা পৰা? পিছে সেইয়া বাইজক লতা লাগিল, চিকিৰা এটা মনুৰিলে গাটো বাইজাই আক মনুৰো কিনিবিন কৰা হ'ল, সেয়াটো কাৰোবাৰ কাৰোবাৰ তেনে হ'ল— এতিয়া তোমালোকেই ফেছাই বোলে

(a) Assamese

আনমনা হয়ে যাচ্ছিলে। ঠিক বুঝতে পারছিলাম হোমার বই পাঠিয়ে গেছে। তখনই আমার মনে হল, হোমার জন্য কিছু করতে হবে। আমরা একসঙ্গে মাতাল হবো, নাচবো, ফুর্তি করব। এর মধ্যে বিশপনে কী আছে? টান্নি ঘীর হয়ে এল। তারপর একটা বা চকচকে বার কাম হোমারের সামনে থাকল। বনি রায় দেখলেন, এ হচ্ছে একেবারে নইলই হেঁচুটুটুদের বেলেম্পানো করার হাত জায়গা। একটা নির্জন পরিষ্কার মধ্যে টান্নি ভাড়া দিচ্ছে বনি মেয়েটির হাতে হাত ধরে ঢুক সেলেন ভিতরে। নীল ফিফটিটা নিশ্চিত খোমছে কাঠেপাঠে। বনি রায় বাড় খোলালেন না। আলো এত সুন্দর যে, বাইরে থেকে ভিতরে ঢুকলে যুঁশুটি অন্ধকার বলে মনে হয়। মেয়েটি হাত ধরে টানে নিয়ে না গলে বরকে কিছুক্ষণ হাতছাড়াতে হত চারপাশে।

কোলে একটা ফাঁকা একটেরে কিউবিকল। মেয়েটি খুব কাছ থেকে শরীরে শরীর লাগিয়ে বলল, হোমারীয় শব্দে নেই বললেই হয়। কী যাবে? পুরসেবা তো হুইকিই যায়। আমি বাবো তোদাকা। বনি শুধু বললেন, এনিথিং ইউ সে। কাচের দরজাটা ঠেলে দুজন লোক ঘরে ঢুকল। চারদিকে তাকাল। তারপর আঁধা অন্ধকারে কোথাও বলে গেল। মেয়েটি মূদু করে বলল, হুইস এ মাড জয়েক। বাতেরে দিকে এ জায়গাটা একেবারে ফ্রোই হয়ে যায়। হ্যাঁ, এ হচ্ছে বৌবনের জায়গা আমার মতো বুড়াদের নয়। মেয়েটি বনির পাশে একটা ঠোনা দিয়ে বলল, খুশি মোটেই বুড়া নও। বোলা, আমি টয়লেট থেকে আসছি। চিকা উঠে যেতেই বনি রায় দুটো জিনিস লক্ষ করলেন। চিকা তার হ্যাণ্ডব্যাগ সঙ্গে নিয়ে গেছে। সেটা স্বাভাবিকও হতে পারে। মেয়েদের অনেক সাঙ্কশোজের জিনিসও হ্যাণ্ডব্যাগে থাকে। বিতীয়ত চিকা ফ্লিকসেনে অভরি দিয়ে যায়নি। টয়লেট কোন দিকে তা বনি রায় জানেন না, চিকা গেল ভদনদিকে। বহিবে বোলাদের দরজাও ওই দিকেই। বনি রায় পরিষ্কৃতি বুঝে নিলেন চোখের পলকে। ঠিক কম্পিউটারের মতোই। নির্জন এই রোস্তারী তীর করখানা হয়ে উঠতে পারে, যদি সত্যক না হন তিনি। বনি হঠাৎ উঠে নড়ালেন। আর তৎক্ষণাৎ বুঝতে পারলেন, পিছনের ৫০

(b) Bangla

૫૪ કિલો : વાસ્તી □ ૧૬૦

સૌથીઊંચા	૧૯૪૮	૧૯૫૨	
	૩૨	૨૯	૫૫ સુધી
			૧૬૪ % ઊંચે,
ઊંચા વાસ્તી	૬૮	૬૩	૩૦ % ૫૫ સુધી
			૩૧ %
ઘેરેજી વીજળીની	૧૯૪૮	૧૫	૫૨
	૨૬.૩૫	૧૨.૬૭	૧૦.૩૦
			૮.૮૭

જાંબાં વીજી કોન્યાના વનર વાપરતા. ઘરે નહીં મીંડે નથી? છે. મુક્તિ પહોંચ ન શકું.

Rate per unit સૌથીઊંચા ૯૦૦ Yen
ઘેરના ૩૩૦૦ Yen

૧૦૦૦ કિલોવૉટ ૧ કલાક ૫૫ રૂ ૧ મુક્તિ મીંડે શાંકાંકાં પાનેજી મુક્તિ પહોંચ મીંડેક વુર માટે ઊંચીકિં ઘેરેથી માટે. આંવેરેકાં માટે. જંતુ મારવા માટે.

વેશ્યાંકાં પ્રાન્યાંકાં વાસ્તીથી વીજળી. જ્યાંનો નમુનો. ઠેરથી યાવ છે? નથી જાણી.

૨. કાપવાથી વન

ઘાથયાળ	૮ કલાકમાં	૩૦ મીંડે
ધંર	૮ કલાકમાં	૪૦ yards.

૩. રેશમ

Mutual Aid Teams		૧
ઘેરકાં	૪૮૫ ૧૦૦	૫૨૫ ૧૨૬૦
મીંડેકાં	૧૦૦	૨૧૩ %

પાને ને વાણુ કસાં. છેડે કાંરેવાળાં ઘાં. મલવની વુશ. જ્યાં માંડ જ્યાં. ઘરે ઘરેકાં ઘા આપે છે. ઘેરોમાં કાંકાં પાંડાં ૫૨ જ્યાંની મિલકત. ઘાંકાં જ્યાંની ઘાં. પાંડાં ૫૨. કાંકાં ઘેરોમાં વાંકાં જ્યાં. ઘાં. ૨ ઘાંકાં ઘેરોમાં સળીઓમાં. પાંડાં. સળીમાં મુકવી માટે. રેશમ મુકવા. ઘેરોમાં સાર કરે છે. પહોંચાં ઘાથથી. ઘરે મર્યાનાં કાપકાં જુદા પાંડાં એ રીતના ઘાથનથી. કેરોમાં ઘે રીત

પહોં માટે. તરુનો છેડે મીં. ૮ લેવા કરી ઘેરે કાંકે ને રીલ બનાવે. કેરોમાં પાંચ. ઉત્કાંકાં ૪૮૫ ૫૨ ૪ જુદું વધુ. પાંડાં નહીં તે માટે પાંડાંથી કાંકાં છે. ડાંકાંની મર્યાં આંકાં. વીજળીથી ૪૮. ૫૦. ૫૧. ૫૨. ૫૩. વાસ્તી. પહોંચાં ૩ મુક્તિ ૨ ઉપર. ઘરે ૪. ૧૨ ઉપર. ઘરેવાળું કાપકાં.

(c) Gujarati

ई बेटे का नाम रखा गुरविवन्दर इसका पश्चात छोटी मनप्रीत ने एक के बाद एक चार बेटियों को जन्म दिया। गुरप्रीत को बहुत चाह थी, एक और बेटे की। सो चार बेटियों के बाद एक और बेटा जना मनप्रीत ने तो परिवार मानो सम्पूर्ण हो गया।

किलकारियों से गुजरी घर की दीवारें पहले से अधिक मजबूत होने लगी। किसी को झुक दाखिल कराना है या कोई बीमार है, डाक्टर को दिखाना है सब गुरप्रीत का जिम्मा था। दो मासों की शीतल, सौहार्दपूर्ण छत्रछाया में बच्चे पल कर बढ़े हो गए, पता ही न चला।

हरजोत जवान हुआ तो उसके दिल में भी विवाह की लालसा तिर उठाने लगी। परन्तु एक अपाठिज को कौन अपनी बेटा देता? पेशा ही तो सब कुछ नहीं।

बैसे हरजोत प्रतिदिन दुकान पर जाता। सरदार या गुरविवन्दर के साथ पीछे स्कूटर पर बैठता तो गुरप्रीत सहारा देते हुए अपने बेटे को निरीह नजरों से टटोलीत। उसकी आंखों में तैरते कोमल सपनों को समझने की कोशिश करती। साय दुकान से वापस आ अपने साथ के लड़कों की बीवियों और बच्चों का जिक्र उड़ता तो एक हीनता का भाव उसके चेहरे पर उमर आता।

मा ने ताड़ लिया। बस तभी से जुट गई बहू की ललाशा में आखिर गुरप्रीत ही तो थी जिसकी कोशिशों से आज घर-आंगन महक रहा था। बेटे का घर बसाना कौन मुश्किल काम था। सुन्दर, पढ़ी-लिखी एवं सुरील लड़की खोज लार्ई हरजोत के वास्ते। गरीब घर की थी तो क्या हुआ। अपने घर में किसी चीज की कमी थोड़े ही थी। इन्दरजीत के मां-बाप दूर शर्त पर राजी हो गए थे कि विवाह का समस्त प्रबन्ध हरजोत के घरवाले ही करेंगे। गहने और कपड़े-लत्ते भी गुरप्रीत ने ही बनाने की हामी मरी थी।

सिक्कड़ते वाक्ये : 17

(d) Hindi

Figure A.1: Assamese, Bangla, Gujarati and Hindi samples from our internal OCR dataset.

ನಕ್ಕು ನಗಿಸು

ಒಮ್ಮೆ ನಕ್ಕು ನನ್ನ ನಗಿಸು ;
ಬಿಮ್ಮೆನುಕಿಹುದೆನ್ನ ಮನಸು,
ಮುನ್ನಿಲ್ಲದ ಚಿನ್ನ ಗನಸು
ನಿನ್ನ ಹಾನ ಕೈಗೂಡಿತು.

ತಣ್ಣನೆ ಕಡಲಂತೆ ನಗು,
ಬಣ್ಣದ ಮುಗುಳಂತೆ ನಗು,
ಹುಣ್ಣೆಮೆಯುರುಳಂತೆ ನಗು ,
ನನ್ನೆಂಬುಮೆಯ ಹೆಣ್ಣೆ ನಗು .

ಅತ್ತು ಸತ್ತುಹೋಕ ಬದುಕ -
ನೆತ್ತಿ ತುಂಬಿ ಜೀವಕಳೆಯ
ಮುತ್ತಿನಿಟ್ಟು ರೆಪೆಯತ್ತಿ
ಕಣ್ಣನೊಳಗೆ ದೀಪಹಚ್ಚಿ .
ಎದೆಯ ಹೊರೆಯಿಸಿಲಿ, ಮೇಲೆ
ಬದುಕ ನಗುವ ಬದುಕ ಬೆಲೆ ,
ನಕ್ಕಲೆ ನಗಬಲ್ಲುದಿಳೆ.

ಕೆನ್ನೆ ಕೆನ್ನೆಯೊಂದು ಮಾಡಿ
ನಗುವ ನನ್ನ ನಿನ್ನ ನೋಡಿ
ಕನ್ನಡಿಯಲ್ಲಿ ಒಂದು ಜೋಡಿ
ನಕ್ಕುಬಿಡಲಿ, ಒಮ್ಮೆ ನಗು.

ನಕ್ಕುಬಿಡು, ಚಿಕ್ಕಮುತ್ತು
ಸುರಿದುಹೋಗಲಿ ;
ಪಕ್ಕನೊಮ್ಮೆ ನಕ್ಕುಬಿಡು.
ಚಿಂತೆ ತೊಲಗಲಿ.

೨೪

(a) Kannada

೧೫೫

ಸ್ವಾಹಿರಿಯಾಚಾರ್ಯರ ಜೀವಿತ, ಉಗ್ರಹೃದಯವಿಭೂಷಿತ ಅನುಭವಗಳಿಂದ ಸಾಕಾರವಾಗಿ ಬರೆಯಲ್ಪಟ್ಟಿರುವ ಈ ಕೃತಿ, ಉತ್ತಮವಾಗಿ ಬರೆಯಲ್ಪಟ್ಟಿದೆ. ಅದರಲ್ಲಿ ಅನೇಕ ಅಂಶಗಳನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಅದರಲ್ಲಿ ಅನೇಕ ಅಂಶಗಳನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಅದರಲ್ಲಿ ಅನೇಕ ಅಂಶಗಳನ್ನು ಸೂಚಿಸಲಾಗಿದೆ.

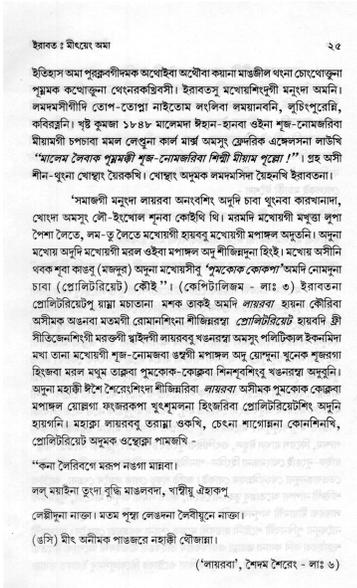
ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ.

ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ.

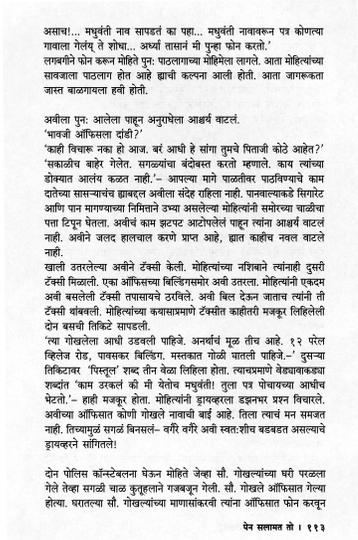
ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ.

ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ. ಇಲ್ಲಿಯೂ ವ್ಯುತ್ಪತ್ತಿಯನ್ನು ಸೂಚಿಸಲಾಗಿದೆ.

(b) Malayalam



(c) Manipuri



(d) Marathi

Figure A.2: Kannada, Malayalam, Manipuri and Marathi samples from our internal OCR dataset.



Figure A.4: A sample page image of Urdu from our internal OCR dataset.

“গেজ কাপোৰেৰে দুটা ফুক হৈ কাপোৰ থাকিবহে।” মাষ্টৰে ক’লে, “এটা বাৰু পিচত
প্ৰতিবিন্দুটোৰ দৰেই মাথোঁ এটি শব্দ- এছাটি বতাহে ভাঙি চূৰ্ণ কৰি পেলোৱাৰ
পাইছোনে? অথচ মোৰ নাটকত অভিনয়ৰ প্ৰথম সুযোগ পোৱা ভাস্কৰলৈ
ভকতি মিনতি তুতি নজানো পামৰ-মতি।। ৯১৪

(a) Assamese

কাণ্ডজ্ঞান ততটা নাই, তাহাতে তাহাৰ বৰ্তমান মুগ্ধ অবস্থায় তাহাৰ সাংসাৰিক
যাওয়া দড়িৰ খাটে শুয়ে পড়ল ফেৰ। হিংসাগুলি নাচতে থাকল পাশে। ক্ৰোধগুলি
তাকে সামান্য নাড়িয়ে বললেন—“মা, কোনও কাৰণে অভিমান কৰেছ ? ঘৰে চলো। মায়ের জন্যে
আমার মা কোনোদিন ইন্ধুলে পড়েননি, নিজের চেষ্টায় বাংলাটা আয়ত্ত

(b) Bangla

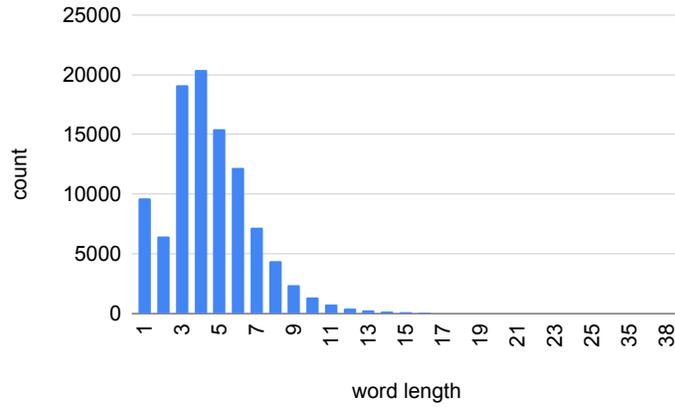
নে पिता आव्या त्यां अत्रो ज्युं के दीकरीनुं आजनुं रुदन कोर्ष अजल
कृष्ण कथामृत” सौ जिज्ञासु साधके जइर वांयवुं. रामकृष्ण परमहंस साक्षात्कारी संत उता.
अने महादेवनाथ जेवा समर्थ यरित्रकार होय तो केवुं श्रवंत श्रवनयरित्र
पाटली वय्ये आंधी इमालनी गणाओणीमां अधर रामेलो उतो. स्टीमरना

(c) Gujarati

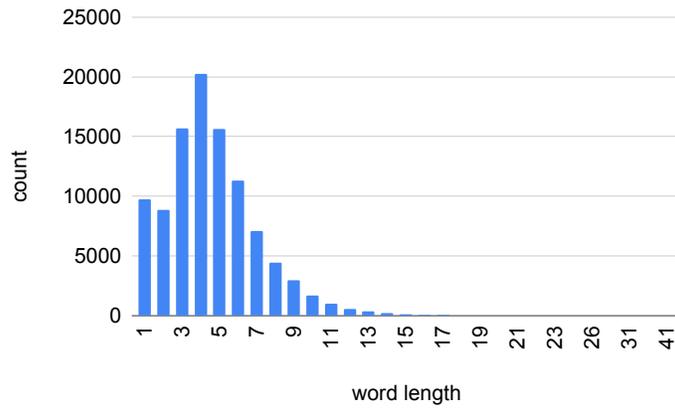
अन्यथा वह मुझे अवश्य सूचित करता कि महिला बिना किसी
“हुजूर.... हुजूर किसी कोने में छिप गया होगा सांप। भीतर तो जरूर है।
मगर यह अस्वस्थ शरीर कब तक साथ देता ! अधिक परि-
हाँ, वह एक नर्स ही थी। उसी तरह के कपड़े, सफेद जूता, मोजा, सफेद

(d) Hindi

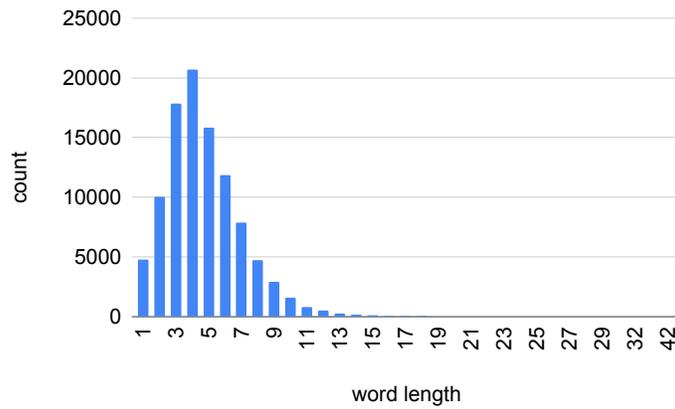
Figure A.5: Samples of Assamese, Bangla, Gujarati and Hindi text-line images in the newly introduced *Mozhi* dataset for text recognition.



(a) Assamese

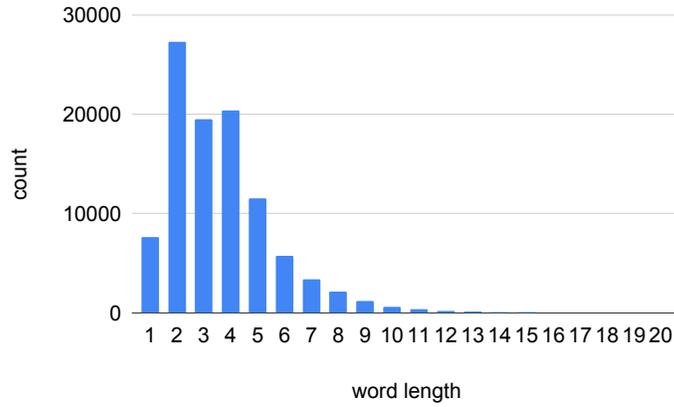


(b) Bangla

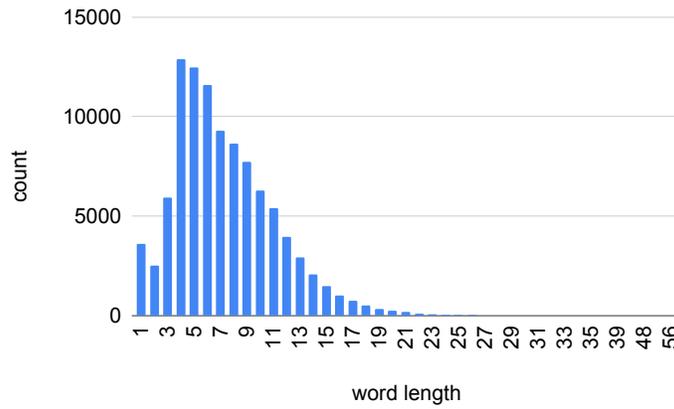


(c) Gujarati

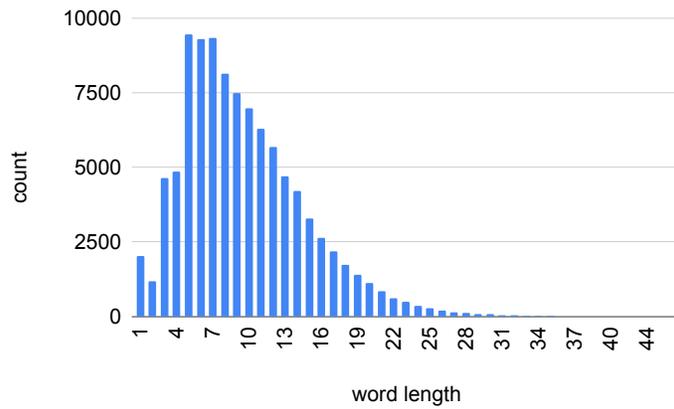
Figure A.8: Distribution of word lengths of Assamese, Bangla and Gujarati samples in the new Mozhi OCR dataset.



(a) Hindi

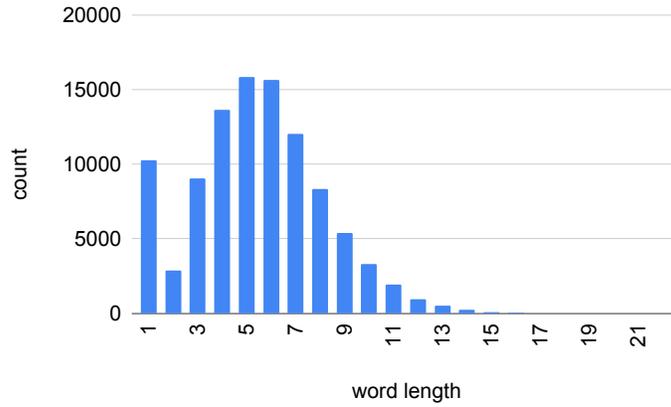


(b) Kannada

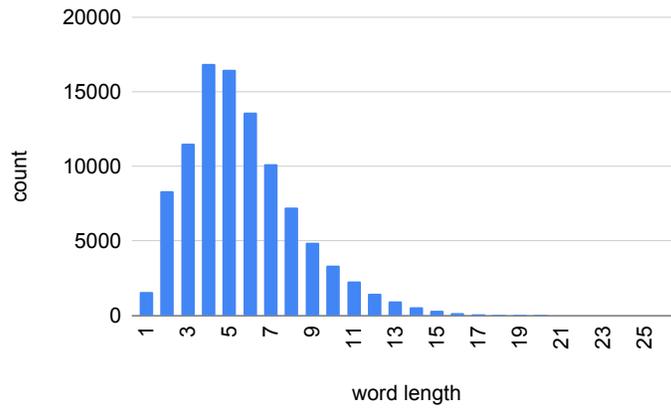


(c) Malayalam

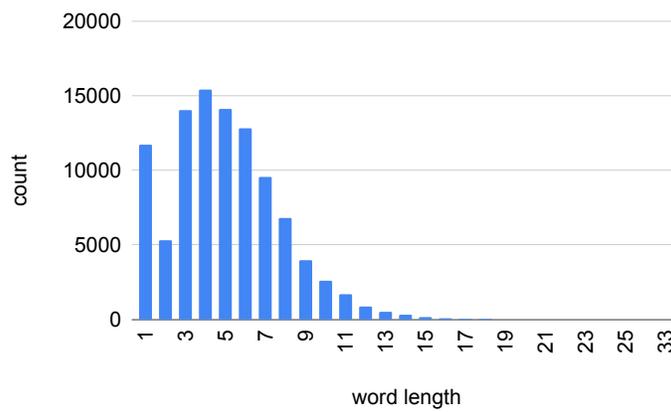
Figure A.9: Distribution of word lengths of Hindi, Kannada and Malayalam samples in the new Mozhi OCR dataset.



(a) Manipuri

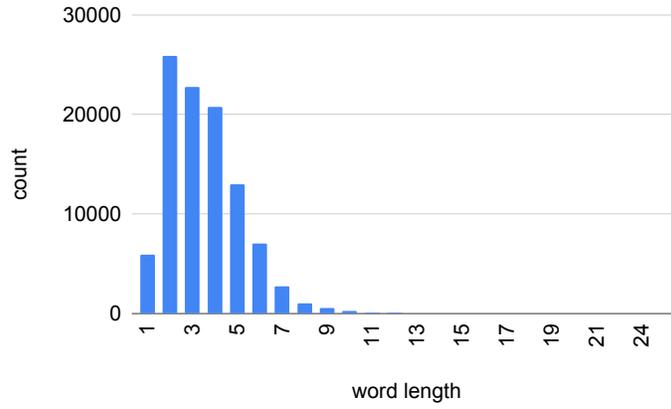


(b) Marathi

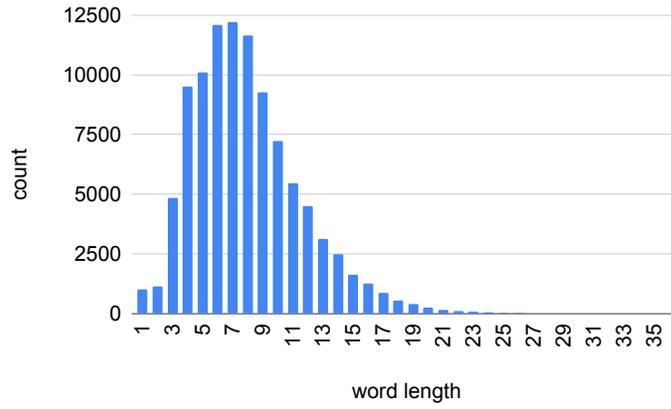


(c) Odia

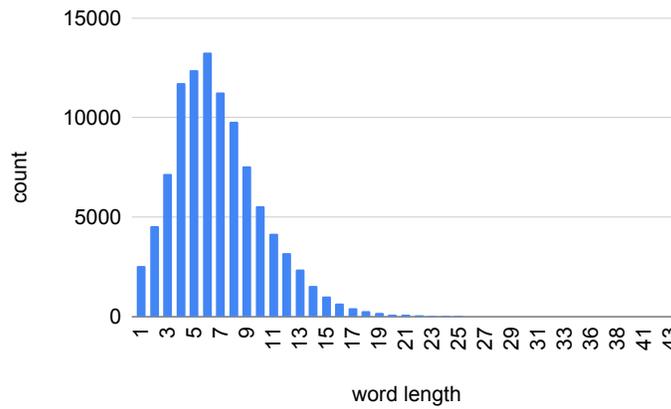
Figure A.10: Distribution of word lengths of Manipuri, Marathi and Odia samples in the new Mozhi OCR dataset.



(a) Punjabi



(b) Tamil



(c) Telugu

Figure A.11: Distribution of word lengths of Punjabi, Tamil and Telugu samples in the new Mozhi OCR dataset.

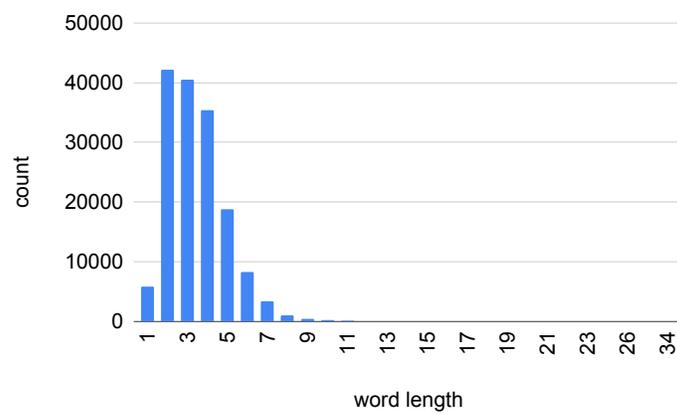


Figure A.12: Distribution of word lengths of Urdu samples in the new Mozhi OCR dataset.

Appendix B

Appendix to chapter 4

In Figure B.1– Figure B.5, we show qualitative examples of performance of baseline models–BERT and M4C, on the DocVQA dataset introduced in chapter 4.

ROCHE, A. P.
370 56 0985

Privileged Communication

4) D₀. This will be done at the Brehm Laboratory, Wright State University, under the direction of Dr. Thomas O. Tierman. Unit charges per test (3 samples per test at \$35 per sample) are \$105 for the -04 year, increasing by 10% per year. There is a unit cost at Webb Associates of \$20 per person (test) for provision of receptacles, handling and transporting the specimens. These unit costs are based on information in letters from Dr. Tierman (pp. 25-26) and Dr. Webb (p. 23).

	Unit Cost	-04		-05		-06	
		N	Total	N	Total	N	Total
D ₀ (Brehm Lab.)	\$105*	220	\$23,100	156	\$16,018	94	\$11,943
Handling & transport (Webb Assoc)	20	220	4,400	156	3,120	94	1,880
Compensation to participants	5	220	1,100	156	780	94	470
TOTALS			\$28,600		\$21,918		\$14,293

TOTALS FOR ALL YEARS \$64,811.
*increased 10% annually.

5) Fat cell size and number. The unit cost is \$209 with a 10% cost adjustment for each additional year. This work will be done at Mt. Sinai School of Medicine under the direction of Dr. Jerome Kistritz. The cost includes provision of personnel, use of equipment, supplies, preparation of reports, relevant consultation and interpretation. This unit cost is described in a letter on p. 27. The total costs are shown below. There is a unit cost at Fels of \$3.00 to cover the cost of disposable syringes (one for local anesthetic and one to obtain fat at each examination), local anesthetic, band-aids and shipping.

	Unit Cost	-04		-05		-06	
		N	Total	N	Total	N	Total
Fat cell size (Mt. Sinai)	\$209*	220	\$45,980	156	\$35,864	94	\$23,772
Biopsy and handling (Fels)	3.00	220	660	156	468	94	282
Compensation to participants	3.00	220	660	156	468	94	282
TOTALS			\$47,300		\$36,800		\$24,336

TOTAL FOR ALL YEARS \$108,436
*The unit cost has been increased 10% annually in accordance with union contracts.

Page 17 January, 1980

Source: <https://www.industrydocuments.ucsf.edu/docs/hm2227>

island desserts (continued from page 7)

HAWAIIAN FRUIT CAKE

1 cup seedless raisins
½ cup seeded raisins
2/3 cup diced citron
¼ cup diced candied orange peel
1¼ cup diced pineapple
1/8 cup chopped dates
½ cup candied cherries
½ cup dried candied lemon peel
1½ cup chopped macadamia nuts
1½ cups shredded coconut
1 tablespoon brandy
1 tablespoon sherry
½ teaspoon ginger juice
1/3 cup flour
1 cup shortening
1 cup brown sugar
4 eggs
½ teaspoon cinnamon
½ teaspoon nutmeg
½ teaspoon soda
¼ cup guava jelly

Preheat oven to 375 degrees. Grease two 8½ x 2½ inch loaf pans. Line with foil or brown paper. Combine fruits, nuts and coconut. Sprinkle with brandy, sherry and ginger juice. Stir in ½ cup of the flour, cream shortening and sugar. Add eggs and beat well. Stir remaining flour with cinnamon, nutmeg, cloves and soda; stir into the creamed mixture. Add jelly and mix well. Stir in fruit mixture. Pour into prepared pans and bake for 3 hours. Make 2 - 2 lb. cakes.

MALASADAS

1 package yeast
½ cup warm water
1 tablespoon sugar
6 cups flour
½ cup sugar
2 cups warm milk
1/8 lb. melted butter
8 eggs slightly beaten

Dissolve yeast in water and 1 tablespoon sugar. Measure dry ingredients in a large bowl, add melted butter, beaten eggs, dissolved yeast and warm milk, (added slowly). Mix well to form soft dough. Cover and place in warm area. Let stand until it rises to double in bulk. Form into small balls and drop into hot oil and cook until brown. Roll in granulated sugar. Serve warm. Yield: Approx. 2½ dozens.

BAUBIA

2 cups coconut milk (frozen)
¼ cup water
4-6 tablespoons sugar
6 tablespoons cornstarch

Melt coconut milk in a double boiler. Combine and stir the above ingredients until smooth. Cook and stir over a low heat until it has thickened completely. Increase the heat slightly and stir the pudding vigorously to prevent it from burning. Remove pudding and pour into a 1½ inch deep cake pan. Let it cool till set, then cut into 2-inch squares and serve.

Board a sailing ship out of the past.
Explore the endless night skies...
Witness authentic chants and dances...
And discover the life of ancient Hawaii.

Bishop Museum's
PASSPORT TO POLYNESIA

SEE YOUR HOTEL TRAVEL DESK OR HERITAGE THEATRE, KING'S ALLEY 922-3368

FOR FREE CATALOG SEND STAMPED SELF-ADDRESSED ENVELOPE TO THE POLY SHOP, P.O. Box 2401, Honolulu, Hawaii 96825

EXOTIC

THE DIAMOND PALACE

NEW LOCATION! World Square, Kalia-King
Located at the International Market Place • Waikiki Beach • Honolulu, Hawaii

46 Latitude 20/November 1978

are shown below. There is a unit cost at Fels of \$3.00 to cover the cost of disposable syringes (one for local anesthetic and one to obtain fat at each examination), local anesthetic, band-aids and shipping.

	Unit Cost	-04		-05		-06	
		N	Total	N	Total	N	Total
Fat cell size (Mt. Sinai)	\$209*	220	\$45,980	156	\$35,864	94	\$23,772
Biopsy and handling (Fels)	3.00	220	660	156	468	94	282

island desserts (continued from page 7)

HAWAIIAN FRUIT CAKE

1 cup seedless raisins
½ cup seeded raisins
2/3 cup diced citron
¼ cup diced candied orange peel
1¼ cup diced pineapple
1/8 cup chopped dates

MALASADAS

1 package yeast
½ cup warm water
1 tablespoon sugar
6 cups flour
½ cup sugar
2 cups warm milk

Q: What is the total cost for Fat cell size (Mt. Slnai) in the -05 year ?

Question type: table/list

GT: \$35,864

M4C best: 4400

BERT best: \$35 , 864

Human: \$35,864

Q: What is the first recipe on the page?

Question type: table/list

GT: hawaiian fruit cake

M4C best: island desserts (continued from cake

BERT best: hawaiian fruit cake

Human: hawaiian fruit cake

Figure B.1: DocVQA: Examples where BERT QA model [113] answers questions other than 'running text' type. On the left is a question based on a table and for the other question one needs to know the 'first recipe' out of the two recipes shown. For the first question the model gets the answer correct except for an extra space, and in case of the second one the predicted answer matches exactly with the ground truth answer.

16. Courses in which you wish to enroll: (Please read instructions carefully and check the appropriate boxes).

A.M. Schedule
 Fundamentals of Biostatistics M-T-W-Th-F.
 If you have selected this course, do not select the following morning courses:

Alternate
 Principles and Methods of Epidemiologic Research T-Th-S
 or
 Epidemiology of Occupational Hazards T-Th-S
 and
 Epidemiology of Cancer M-W-F
 or
 Epidemiology of Cardiovascular Diseases M-W-F

P.M. Schedule
 Fundamentals of Epidemiology M-T-Th-F
 If you have selected this course, do not select the following afternoon courses:

Alternate
 Health Services Planning and Evaluation T-Th
 or
 Epidemiology of Injuries T-Th
 and
 Infectious Disease Epidemiology M-W-F
 or
 Epidemiology of Nutritional Diseases and Abnormalities M-W-F

17. Please make the following room reservation for me:
 Middlebrook Hall (Two in a room, twin beds) Arrive _____ Depart _____

18. Signature of Applicant *Grande H. Hwang* Date *4/9/76*

19. I approve of this application _____
 Department, Chairman or Advisor

Send this form with check for \$25.00, made payable to the University of Minnesota, to: Dr. Leonard M. Schuman, Program Director, Epidemiology Summer Session, Division of Epidemiology, Room 1-117 Unit A, Health Sciences Building, University of Minnesota, Minneapolis, Minnesota 55455.

APPLICATIONS MUST BE RECEIVED BY MAY 17.

 **Yourself Grow!**

The University of Minnesota adheres to the principle that all persons shall have equal opportunity and access to facilities in any phase of University activity without regard to race, creed, color, sex, age or national origin. Under this principle, educational, cultural, social, housing, extra curricular and employment opportunities are available to all on an equal basis.

Fulfilling the potential to win

Inspired individuals share their story—what drove them to succeed against the odds and how they realized their goals.

"The launch of three new brands has been the joint effort of everyone, right across the company. From buying the lead, to the blend, the R&D, the processing, the packaging and finally the marketing, sales and distribution all had to work in tandem, to achieve one goal. It's really true, success is about teamwork, it is about five fingers coming together to form a fist."

-GPI grew by 3 times the industry rate of growth in 2002-03

-Successful launch of Tipper and Prince



"I am really proud to be part of the team that worked on re-launching Prince in the market. We knew that we had a winning brand on our hands. We were so clear about the aim, we knew that we had to find a way to achieve it. That's why I believe that when you have the passion to succeed, nothing can stand in your way. It's a human condition, using the forces you need."

The passion to succeed - it's a value that runs right across the organization. Every individual is determined to realize his or her potential to be a winner.

Source: <https://www.industrydocuments.ucsf.edu/docs/znbx0223>



Q: What is written inside logo in the bottom of the document?

Question type: layout

GT: let yourself grow!

M4C best: yourself grow!

BERT best: < no prediction >

Human: let yourself grow!



Q: What Tobacco brand of GPI is shown in the picture?

Question type: photograph

GT: Prince

M4C best: prince

BERT best: < no prediction >

Human: prince

Figure B.2: **DocVQA: M4C [94]’s performance on questions based on pictures or photographs.** Here we show two examples where the best variant of the M4C model outperform the BERT best model in answering ‘layout’ type questions seeking to read what is written in a logo/pack. The BERT model doesn't make any predictions for the questions.



Report on Corporate Governance

- Major accounting entries based on exercise of judgement by management
 - Significant adjustments, if any, arising out of audit
 - Compliance with Accounting Standards
 - Compliance with Stock Exchange and legal requirements concerning financial statements
 - Related party transactions
 - Qualifications, if any, in draft audit report
 - Report of the Directors & Management Discussion and Analysis;
- (d) Reviewing with the management, external and internal auditors, the adequacy of internal control systems and the Company's statement on the same prior to endorsement by the Board;
- (e) Reviewing the adequacy of the internal audit function, including the structure of the internal audit department, staffing and seniority of the official heading the department, reporting structure, coverage and frequency of internal audit;
- (f) Reviewing reports of internal audit, including that of wholly owned subsidiaries, and discussion with internal auditors on any significant findings and follow-up thereon;
- (g) Reviewing the findings of any internal investigations by the internal auditors and the executive management's response on matters where there is suspected fraud or irregularity or failure of internal control systems of a material nature and reporting the matter to the Board;
- (h) Discussion with the external auditors, before the audit commences, on nature and scope of audit, as well as after conclusion of the audit, to ascertain any areas of concern and review the comments contained in their management letter;
- (i) Reviewing the Company's financial and risk management policies;
- (j) Looking into the reasons for substantial defaults, if any, in payment to shareholders (in case of non-payment of declared dividends) and creditors;
- (k) Considering such other matters as may be required by the Board;
- (l) Reviewing any other areas which may be specified as role of the Audit Committee under the Listing

Agreement, Companies Act and other statutes, as amended from time to time.

Composition

The Audit Committee presently comprises four Non-Executive Directors, three of whom are Independent Directors. The Chairman of the Committee is an Independent Director. The Executive Director representing the Finance function, the Chief Financial Officer, the Head of Internal Audit and the representative of the Statutory Auditors are Invitees to the Audit Committee. The Head of Internal Audit is the Co-ordinator and the Company Secretary is the Secretary to the Committee. The representatives of the Cost Auditors are invited to meetings of the Audit Committee whenever matters relating to cost audit are considered. All members of the Committee are financially literate; three members, including the Chairman of the Committee, have accounting and financial management expertise.

The names of the members of the Audit Committee, including its Chairman, are provided under the section 'Board of Directors and Committees' in the Report and Accounts.

Meetings and Attendance

Details of Audit Committee Meetings during the financial year

During the financial year ended 31st March, 2014, eight meetings of the Audit Committee were held, as follows:

Sl. No.	Date	Committee Strength	No. of Members present
1	6th May, 2013	6	4
2	17th May, 2013	6	5
3	25th July, 2013	6	6
4	28th August, 2013	5	5
5	23rd September, 2013	5	5
6	25th October, 2013	5	5
7	17th January, 2014	5	5
8	31st March, 2014	5	3

Source: <https://www.industrydocuments.ucsf.edu/docs/tmbx0223> ITC Limited REPORT AND ACCOUNTS 2014 15

Meetings of the Audit Committee were held, as follows:

Sl. No.	Date	Committee Strength	No. of Members present
1	6th May, 2013	6	4
2	17th May, 2013	6	5
3	25th July, 2013	6	6
4	28th August, 2013	5	5
5	23rd September, 2013	5	5
6	25th October, 2013	5	5
7	17th January, 2014	5	5
8	31st March, 2014	5	3

Q: What was the committee strength for the first meeting?

Question type: table/list

GT: 6

M4C best: 6

BERT best: 6

Human: 6

Q: What was the committee strength for the last meeting?

Question type: table/list

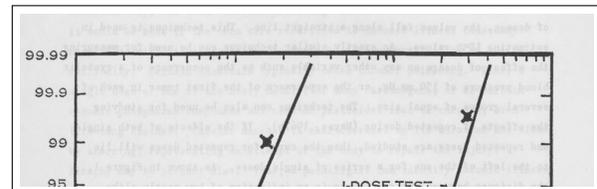
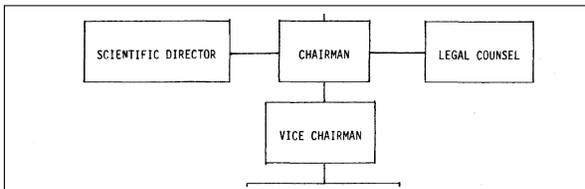
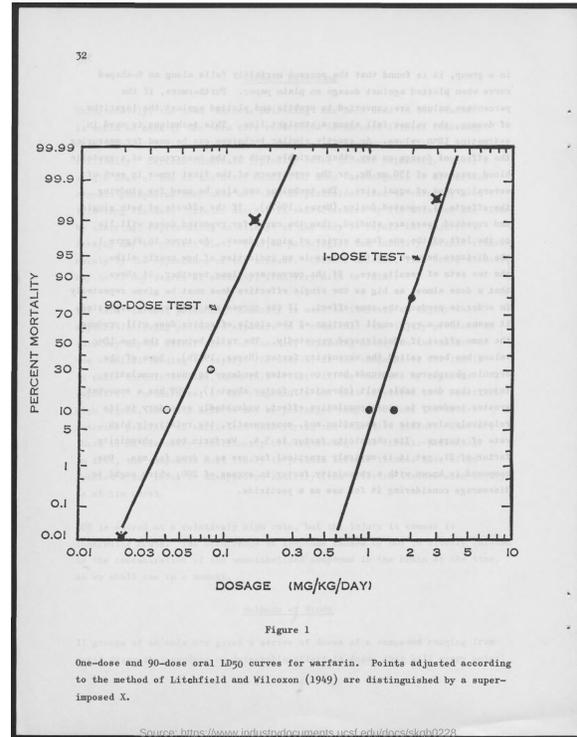
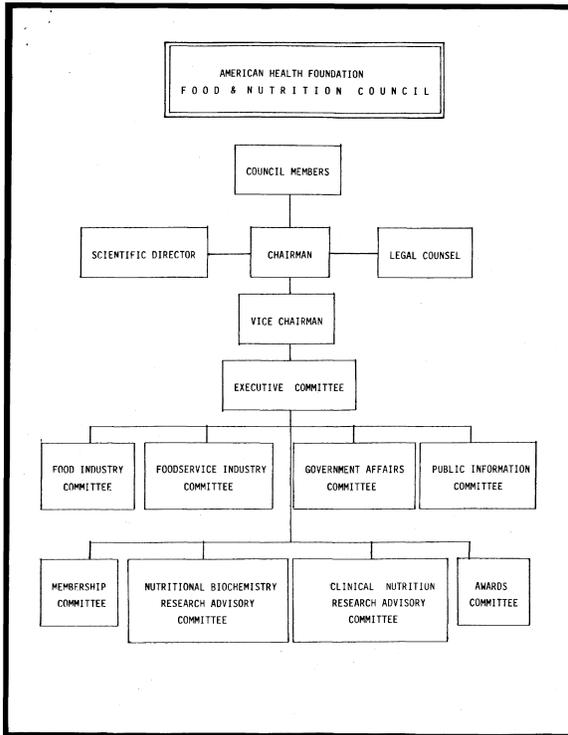
GT: 5

M4C best: 6

BERT best: 6

Human: 5

Figure B.3: DocVQA: Contrasting results for similar questions. Both questions ask for 'committee strength' for a particular meeting (first or last). Both models get the answer right for the first one only. This suggests that the models' predictions are not backed by a proper reasoning/grounding in all cases.



Q: What is the position above "vice chairman" ?

Question type: figure

GT: chairman

M4C best: legal counsel

BERT best: legal counsel

Human: chairman

Q: What is the highest value shown on the vertical axis?

Question type: figure

GT: 99.99

M4C best: 50

BERT best: 32

Human: 99.99

Figure B.4: **DocVQA: Questions based on figures and diagrams.** In case of the question on the left, one needs to understand an organizational hierarchy diagram. For the second question, one needs to know what a 'vertical axis' is, and then find the largest value. Both the models fail to answer the questions.

Q: What is the name of the passenger?

Question type: form, handwritten

GT: dr. william j. darby

M4C best: larry

BERT best: larry

Human: dr. william j. darby

Q: What is the date present in the memo ?

Question type: form, handwritten

GT: 1/7/77

M4C best: 1 7 7 7

BERT best: 1 / 7

Human: 1/7/77

Figure B.5: **DocVQA: Impact of OCR errors.** Here the models are able to ground the questions correctly on the relevant information in the image, but failed to get the answers correct owing to the OCR errors. In case of the question on the left, even the answer entered by the human volunteer is not exactly matching with the ground truth. In case of the second question, OCR has split the date into multiple tokens due to over segmentation, resulting in incorrect answers by both the models.

Appendix C

Appendix to chapter 5

In Figure C.1– Figure C.7, we show qualitative examples of performance of VLL-BERT model on InfographicVQA (see chapter 5), covering different QA types in the dataset.

Gender in the Global Research Landscape

Elsevier's comprehensive report on research performance through a gender lens, *Gender in the Global Research Landscape*, spans 20 years, 12 geographies, and 27 disciplines. This global study draws upon data and analytics, a unique gender disambiguation methodology, and involvement of global experts.

FOCUS ON COMPUTER SCIENCES

Women Men

Researchers

Proportion of women and men among researchers in Computer Sciences 2011-15



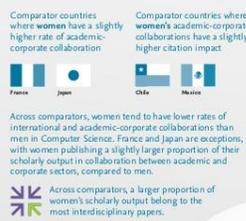
Leadership

Appearance of authors in first or corresponding position for papers in Computer Science 2011-15



Research Collaboration

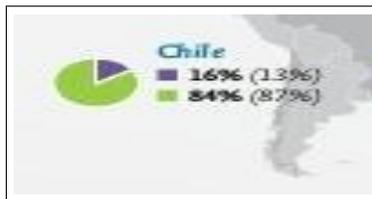
Collaboration and Field-Weighted Citation Impact in Computer Science 2011-15



*Field-Weighted Citation Impact normalizes the data to a world average of 1 in order to account for different citation rates and practices across various fields, topics, and ages. Access the full report at: www.elsevier.com/research-intelligence/essence/industry/industry-report

Elsevier is a registered trademark of Elsevier B.V. | RCUK Online and the RCUK symbol are trademarks of RCUK Intellectual Property UK, used under license. © 2014 Elsevier B.V.

Powered by Scopus



Q: what percent of researchers in Chile were men in the duration of 2011-15?

GT: [84%, 84]

VLL-BERT: 23%

M4C: 23%

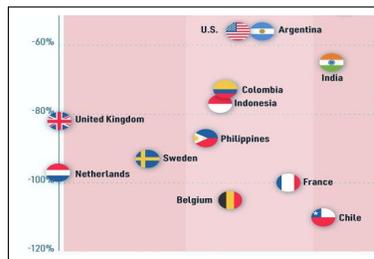
Human: 84%

Answer-source: Image-span

Evidence: Map

Operation: none

Figure C.1: **InfographicVQA: A question where color codes and information on a Map are required to arrive at the answer.** To answer this question, models require to understand that the blue color correspond to women and then pick the number corresponding to the blue color from the data given for Chile. Both the baseline models get this question wrong. Note that here there are two valid answers (GT), one added during the first stage of annotation and the other during the second stage.



Q: Which quadrant does the country India fall into, blue, pink, or gray?

GT: pink

VLL-BERT: country

M4C: pink

Human: pink

Answer-source: Question-span

Evidence: Figure Visual/Layout

Operation: none

Figure C.2: **InfographicVQA : An example for Question-span** To answer this question, a model needs to first locate "India" on the image and then identify the background color there. M4C gets this question correct.



Q: How many championships has Kobe Bryant won?

GT: 5

VLL-BERT: 5

M4C: 5

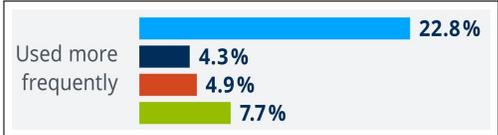
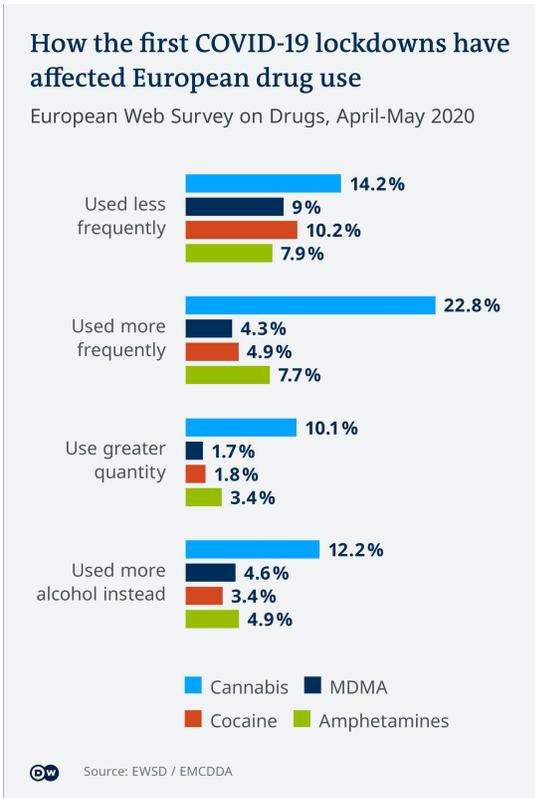
Human: 5

Answer-source: Non-extractive

Evidence: Figure

Operation: Counting

Figure C.4: **InfographicVQA** : Example where symbols/markers need to be counted to find an answer. Both the models get the answer correct for this question that require one to count the yellow squares next to "CHAMPIONSHIPS".



Q: Which drug was used more frequently during lockdown, MDMA, Cocaine, Cannabis, or Amphetamines?

GT: **cannabis**

Answer-source: **Question-span Image-span**

VLL-BERT: **cannabis**

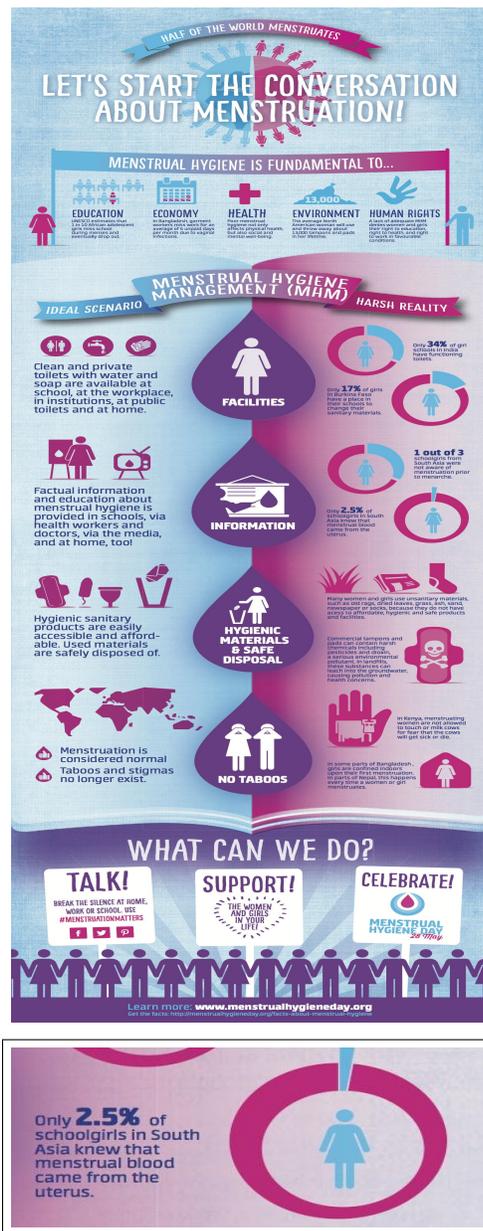
Evidence: **Figure**

M4C: **cocaine**

Operation: **Sorting**

Human: **cannabis**

Figure C.5: **InfographicVQA : An example where values shown in a bar chart need to be sorted to find the answer.** In this question, answer is a span of question (Question-span) and a span of the text on the image (Image-span) as well. The largest among the given items is explicit in the bar chart representation. Alternatively the same can be found by finding the largest by comparing the numbers. Hence 'Sorting' is added as the Operation.



Q: What % of schoolgirls in South Asia do not know that menstrual blood comes from the uterus?

GT: [97.5, 97.5%]

Answer-source: Non-extractive

VLL-BERT: 2.5%

Evidence: Text

M4C: 25

Operation: Arithmetic

Human: 97.5%

Figure C.6: **InfographicVQA: Question requiring arithmetic operation.** To answer this question, the given percentage value needs to be subtracted from 100. Both the models fail to get the answer correct.



Q: Playing against which country did he reach the most number of his milestone runs?

GT: **sri lanka**

VLL-BERT: **bangladesh**

M4C: **pakistan**

Human: **sri lanka**

Answer-source: **Image-span**

Evidence: **Text Figure**

Operation: **Counting Sorting**

Figure C.7: **InfographicVQA: Performing multiple discrete operations to find the answer.** Here the context required to find the answer spans the entire image. Hence we do not show a crop of the image in the inset. This question requires a model to do Counting — count number of milestone runs scored against each country and then perform Sorting — find the country against which the player scored most milestone runs.

Bibliography

- [1] S. Briet and L. Martinet, *What is documentation?: English translation of the classic French text*. Scarecrow Press, 2006.
- [2] D. L. Hoffmann, C. D. Standish, M. García-Diez, P. B. Pettitt, J. A. Milton, J. Zilhão, J. J. Alcolea-González, P. Cantalejo-Duarte, H. Collado, R. De Balbín, *et al.*, “U-th dating of carbonate crusts reveals neandertal origin of iberian cave art,” *Science*, vol. 359, 2018.
- [3] A. Robinson, *Writing and script: a very short introduction*. Oxford University Press, 2009, vol. 208.
- [4] E. L. Eisenstein, E. Elizabeth Lewisohn, *et al.*, *The printing revolution in early modern Europe*. Cambridge University Press, 2005.
- [5] M. Mathew, D. Karatzas, and C. Jawahar, “Docvqa: A dataset for vqa on document images,” in *WACV*, 2021.
- [6] M. Mathew, A. K. Singh, and C. V. Jawahar, “Multilingual OCR for indic scripts,” in *DAS*, 2016.
- [7] M. Mathew, L. Gomez, D. Karatzas, and C. Jawahar, “Asking questions on handwritten document collections,” *IJDAR*, vol. 24, 2021.
- [8] M. Mathew, M. Jain, and C. V. Jawahar, “Benchmarking scene text recognition in devanagari, telugu and malayalam,” in *MOCR Workshop, ICDAR*, 2017.
- [9] M. Mathew, V. Bagal, R. Tito, D. Karatzas, E. Valveny, and C. Jawahar, “Infographicvqa,” in *WACV*, 2022.
- [10] “Global paper and paperboard production volume from 2010 to 2020, by type,” https://www.statista.com/statistics/270317/production-volume-of-paper-by-type/?_ga=2.204976272.1018593204.1622837058-191240632.1618425162, accessed on 04 June 2022.
- [11] V. Ambati, N. Balakrishnan, R. Reddy, L. Pratha, and C. Jawahar, “The digital library of india project: Process, policies and architecture,” in *ICDL*, 2006.
- [12] “Project gutenber,” <https://gutenberg.org>, accessed on 04 June 2022.
- [13] L. Vincent, “Google Book Search: Document Understanding on a Massive Scale,” in *ICDAR*, 2007.
- [14] J. Purday, “Europeana: Digital access to europe’s cultural heritage,” *Alexandria*, vol. 23, 2012.
- [15] M. Moyle, J. Tonra, and V. Wallace, “Manuscript transcription by crowdsourcing: Transcribe bentham,” *Liber Quarterly*, vol. 20, 2010.

- [16] G. Michalek, “The Heinz Electronic Library Interactive On-Line System (HELIOS): Building A Digital Archive Using Imaging, OCR, and Natural Language Processing Technologies,” 1995.
- [17] T. Causer and V. Wallace, “Building A Volunteer Community: Results and Findings from Transcribe Bentham,” *Digital Humanities Quarterly*, vol. 6, 01 2012.
- [18] D. Tkaczyk, P. Szostek, M. Fedoryszak, P. J. Dendek, and L. Bolikowski, “Cermine: Automatic extraction of structured metadata from scientific literature,” *IJDAR*, vol. 18, 2015.
- [19] C. Clark and S. K. Divvala, “Looking beyond text: Extracting figures, tables and captions from computer science papers,” in *AAAI Workshop: Scholarly Big Data*, 2015.
- [20] D. Doermann, “The indexing and retrieval of document images: A survey,” *CVIU*, vol. 70, 1998.
- [21] L. Vincent, “Google book search: Document understanding on a massive scale,” in *ICDAR*, 2007.
- [22] G. Jaume, H. K. Ekenel, and J.-P. Thiran, “Funsd: A dataset for form understanding in noisy scanned documents,” in *ICDAR Workshops*, 2019.
- [23] S. S. Paliwal, V. D. R. Rahul, M. Sharma, and L. Vig, “TableNet: Deep learning model for end-to-end table detection and tabular data extraction from scanned document images,” in *ICDAR*, 2019.
- [24] “Amazon textract form data,” <https://docs.aws.amazon.com/textract/latest/dg/how-it-works-kvp.html>, accessed on 04 June 2022.
- [25] “Amazon textract tables,” <https://docs.aws.amazon.com/textract/latest/dg/how-it-works-tables.html>, accessed on 04 June 2022.
- [26] T. Kanahori and M. Suzuki, “Scientific pdf document reader with simple interface for visually impaired people,” in *ICCHP*, 2006.
- [27] S. K. Kane, B. Frey, and J. O. Wobbrock, “Access lens: a gesture-based screen reader for real-world documents,” in *SIGCHI Conference on Human Factors in Computing Systems*, 2013.
- [28] C. Ponsard, R. Ramdoyal, and D. Dziamski, “An ocr-enabled digital comic books viewer,” in *ICCHP*, 2012.
- [29] A. Domale, B. Padalkar, R. Parekh, and M. Joshi, “Printed book to audio book converter for visually impaired,” in *Texas Instruments India Educators’ Conference*, 2013.
- [30] M. Suzuki, T. Kanahori, N. Ohtake, and K. Yamaguchi, “An integrated ocr software for mathematical documents and its output with accessibility,” in *ICCHP*, 2004.
- [31] S. Ruder, “NLP Progress: Named entity recognition,” http://nlpprogress.com/english/named_entity_recognition.html, accessed on 07 June 2022.
- [32] —, “NLP Progress: Part-of-speech tagging,” http://nlpprogress.com/english/part-of-speech_tagging.html, accessed on 07 June 2022.
- [33] M. Banko and R. C. Moore, “Part-of-speech tagging in context,” in *COLING*, 2004.
- [34] S. Ruder, “NLP Progress: Intent detection and slot filling,” http://nlpprogress.com/english/intent_detection_slot_filling.html, accessed on 07 June 2022.

- [35] —, “NLP Progress: Sentiment analysis,” http://nlpprogress.com/english/sentiment_analysis.html, accessed on 07 June 2022.
- [36] A. Yadav and D. K. Vishwakarma, “Sentiment analysis using deep learning architectures: A review,” *Artif. Intell. Rev.*, vol. 53, 2020.
- [37] A. Krizhevsky, “Learning multiple layers of features from tiny images,” 2009.
- [38] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [39] W. Rawat and Z. Wang, “Deep convolutional neural networks for image classification: A comprehensive review,” *Neural computation*, vol. 29, 2017.
- [40] Z. Zou, Z. Shi, Y. Guo, and J. Ye, “Object detection in 20 years: A survey,” *arXiv preprint arXiv:1905.05055*, 2019.
- [41] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, “Object detection with deep learning: A review,” *TNNLS*, vol. 30, 2019.
- [42] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez, “A review on deep learning techniques applied to semantic segmentation,” *arXiv preprint arXiv:1704.06857*, 2017.
- [43] T. Grüning, G. Leifert, T. Strauß, J. Michael, and R. Labahn, “A two-stage method for text line detection in historical documents,” *IJDAR*, vol. 22, 2019.
- [44] N. Stamatopoulos, B. Gatos, G. Louloudis, U. Pal, and A. Alaei, “ICDAR 2013 handwriting segmentation contest,” in *ICDAR*, 2013.
- [45] G. Cohen, S. Afshar, J. Tapson, and A. Van Schaik, “Emnist: Extending mnist to handwritten letters,” in *IJCNN*, 2017.
- [46] B. Shi, X. Bai, and C. Yao, “An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition,” *TPAMI*, vol. 39, 2016.
- [47] X. Yang, E. Yumer, P. Asente, M. Kraley, D. Kifer, and C. Lee Giles, “Learning to extract semantic structure from documents using multimodal fully convolutional neural networks,” in *CVPR*, 2017.
- [48] X. Zhong, J. Tang, and A. J. Yepes, “Publaynet: largest dataset ever for document layout analysis,” in *ICDAR*, 2019.
- [49] S. Ruder, “NLP Progress: Question answering,” http://nlpprogress.com/english/question_answering.html, accessed on 07 June 2022.
- [50] F. Zhu, W. Lei, C. Wang, J. Zheng, S. Poria, and T.-S. Chua, “Retrieving and reading: A comprehensive survey on open-domain question answering,” *arXiv preprint arXiv:2101.00774*, 2021.
- [51] S. Liu, X. Zhang, S. Zhang, H. Wang, and W. Zhang, “Neural machine reading comprehension: Methods and trends,” *Applied Sciences*, vol. 9, 2019.
- [52] S. Ruder, “NLP Progress: Dialogue,” <http://nlpprogress.com/english/dialogue.html>, accessed on 07 June 2022.

- [53] J. Ni, T. Young, V. Pandelea, F. Xue, and E. Cambria, “Recent advances in deep learning based dialogue systems: A systematic survey,” 2021.
- [54] H. Chen, X. Liu, D. Yin, and J. Tang, “A survey on dialogue systems: Recent advances and new frontiers,” *SIGKDD Explor. Newsl.*, vol. 19, 2017.
- [55] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, “Vqa: Visual question answering,” in *ICCV*, 2015, pp. 2425–2433.
- [56] Q. Wu, D. Teney, P. Wang, C. Shen, A. Dick, and A. van den Hengel, “Visual question answering: A survey of methods and datasets,” *CVIU*, vol. 163, 2017.
- [57] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra, “Visual Dialog,” in *CVPR*, 2017.
- [58] K. Kafle, B. Price, S. Cohen, and C. Kanan, “DVQA: Understanding data visualizations via question answering,” in *CVPR*, 2018.
- [59] S. E. Kahou, V. Michalski, A. Atkinson, Á. Kádár, A. Trischler, and Y. Bengio, “FigureQA: An annotated figure dataset for visual reasoning,” *arXiv preprint arXiv:1710.07300*, 2017.
- [60] A. Mishra, S. Shekhar, A. K. Singh, and A. Chakraborty, “OCR-VQA: Visual question answering by reading text in images,” in *ICDAR*, 2019.
- [61] S. Ruder, “NLP Progress: Summarization,” <http://nlpprogress.com/english/summarization.html>, accessed on 07 June 2022.
- [62] R. Nallapati, B. Zhou, C. Gulcehre, B. Xiang, *et al.*, “Abstractive text summarization using sequence-to-sequence rnns and beyond,” *arXiv preprint arXiv:1602.06023*, 2016.
- [63] W. S. El-Kassas, C. R. Salama, A. A. Rafea, and H. K. Mohamed, “Automatic text summarization: A comprehensive survey,” *Expert Systems with Applications*, vol. 165, 2021.
- [64] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei, “Scene graph generation by iterative message passing,” in *CVPR*, 2017.
- [65] X. Chang, P. Ren, P. Xu, Z. Li, X. Chen, and A. G. Hauptmann, “A comprehensive survey of scene graphs: Generation and application,” *TPAMI*, 2021.
- [66] M. Z. Hossain, F. Sohel, M. F. Shiratuddin, and H. Laga, “A comprehensive survey of deep learning for image captioning,” *ACM Computing Surveys (CSUR)*, vol. 51, 2019.
- [67] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: Lessons learned from the 2015 MSCOCO image captioning challenge,” *TPAMI*, vol. 39, 2016.
- [68] X. Zhong, E. ShafieiBavani, and A. J. Yepes, “Image-based table recognition: data, model, and evaluation,” *arXiv preprint arXiv:1911.10683*, 2019.
- [69] Z. Huang, K. Chen, J. He, X. Bai, D. Karatzas, S. Lu, and C. V. Jawahar, “ICDAR 2019 competition on scanned receipt ocr and information extraction,” in *ICDAR*, 2019.

- [70] K. Davila, B. U. Kota, S. Setlur, V. Govindaraju, C. Tensmeyer, S. Shekhar, and R. Chaudhry, “Icdar 2019 competition on harvesting raw tables from infographics (chart-infographics),” in *ICDAR*, 2019.
- [71] S. Ruder, “NLP Progress: Relation prediction,” http://nlpprogress.com/english/relation_prediction.html, accessed on 07 June 2022.
- [72] V. Ng and C. Cardie, “Improving machine learning approaches to coreference resolution,” in *ACL*, 2002.
- [73] S. Ruder, “NLP Progress: Coreference resolution,” http://nlpprogress.com/english/coreference_resolution.html, accessed on 07 June 2022.
- [74] R. Yu, A. Li, V. I. Morariu, and L. S. Davis, “Visual relationship detection with internal and external linguistic knowledge distillation,” in *CVPR*, 2017.
- [75] N.-V. Nguyen, X.-S. Vu, C. Rigaud, L. Jiang, and J.-C. Burie, “Icdar 2021 competition on multimodal emotion recognition on comics scenes,” in *ICDAR*, 2021.
- [76] C. Rigaud, N. Le Thanh, J.-C. Burie, J.-M. Ogier, M. Iwata, E. Imazu, and K. Kise, “Speech balloon and speaker association for comics and manga understanding,” in *ICDAR*, 2015.
- [77] W. Huang, C. L. Tan, and W. K. Leow, “Associating text and graphics for scientific chart understanding,” in *ICDAR*, 2005.
- [78] R. Smith, “An overview of the Tesseract OCR engine,” in *ICDAR*, 2007.
- [79] N. Nayef, Y. Patel, M. Busta, P. N. Chowdhury, D. Karatzas, W. Khelif, J. Matas, U. Pal, J.-C. Burie, C.-I. Liu, and J.-M. Ogier, “ICDAR 2019 robust reading challenge on multi-lingual scene text detection and recognition — rrc-mlt-2019,” in *ICDAR*, 2019.
- [80] D. Girish, V. Singh, and A. Ralescu, “Understanding action recognition in still images,” in *CVPR Workshops*, 2020.
- [81] B. Yao, X. Jiang, A. Khosla, A. L. Lin, L. Guibas, and L. Fei-Fei, “Human action recognition by learning bases of action attributes and parts,” in *ICCV*, 2011.
- [82] D. Dua, Y. Wang, P. Dasigi, G. Stanovsky, S. Singh, and M. Gardner, “DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs,” in *NAACL-HLT*, 2019.
- [83] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, *et al.*, “Training verifiers to solve math word problems,” *arXiv preprint arXiv:2110.14168*, 2021.
- [84] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NeurIPS*, 2012.
- [85] A. Graves, S. Fernández, F. J. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *ICML*, 2006.
- [86] A. Graves, M. Liwicki, S. Fernandez, R. Bertolami, H. Bunke, and J. Schmidhuber, “A Novel Connectionist System for Unconstrained Handwriting Recognition,” *TPAMI*, 2009.

- [87] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [88] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” *NeurIPS*, 2014.
- [89] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *ECCV*, 2014.
- [90] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick, “Microsoft coco captions: Data collection and evaluation server,” *arXiv preprint arXiv:1504.00325*, 2015.
- [91] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, “Squad: 100,000+ questions for machine comprehension of text,” *arXiv preprint arXiv:1606.05250*, 2016.
- [92] Y. Yang, W.-t. Yih, and C. Meek, “Wikiqa: A challenge dataset for open-domain question answering,” in *EMNLP*, 2015.
- [93] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, “Bottom-up and top-down attention for image captioning and visual question answering,” in *CVPR*, 2018.
- [94] R. Hu, A. Singh, T. Darrell, and M. Rohrbach, “Iterative answer prediction with pointer-augmented multi-modal transformers for textvqa,” in *CVPR*, 2020.
- [95] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” 2015.
- [96] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *NeurIPS*, 2015.
- [97] N. Sankaran and C. Jawahar, “Recognition of Printed Devanagari text using BLSTM neural network,” in *ICPR*, 2012.
- [98] Thomas M. Breuel and Adnan Ul-Hasan and Mayce Ibrahim Ali Al Azawi and Faisal Shafait, “High-Performance OCR for Printed English and Fraktur Using LSTM Networks,” in *ICDAR*, 2013.
- [99] Adnan Ul-Hasan and Saad Bin Ahmed and Sheikh Faisal Rashid and Faisal Shafait and Thomas M. Breuel, “Offline Printed Urdu Nastaleeq Script Recognition with Bidirectional LSTM Networks,” in *ICDAR*, 2013.
- [100] B. Su and S. Lu, “Accurate scene text recognition based on recurrent neural network.” in *ACCV*, 2014.
- [101] P. Krishnan, K. Dutta, and C. Jawahar, “Word spotting and recognition using deep embedding,” in *DAS*, 2018.
- [102] S. Sudholt and G. A. Fink, “PHOCNet: A Deep Convolutional Neural Network for Word Spotting in Handwritten Documents,” in *ICFHR*, 2016.
- [103] L. Kang, J. Kumar, P. Ye, Y. Li, and D. Doermann, “Convolutional neural networks for document image classification,” in *ICPR*, 2014.
- [104] G. Csurka, D. Larlus, A. Gordo, and J. Almazan, “What is the right way to represent document images?” *arXiv preprint arXiv:1603.01076*, 2016.

- [105] A. Das, S. Roy, U. Bhattacharya, and S. K. Parui, "Document image classification with intra-domain transfer learning and stacked generalization of deep convolutional neural networks," in *ICPR*, 2018.
- [106] S. Long, X. He, and C. Yao, "Scene text detection and recognition: The deep learning era," *IJCV*, vol. 129, 2021.
- [107] T. Kanungo, R. M. Haralick, and I. Phillips, "Global and local document degradation models," in *ICDAR*, 1993.
- [108] A. Sulaiman, K. Omar, and M. F. Nasrudin, "Degraded historical document binarization: A review on issues, challenges, techniques, and future directions," *Journal of Imaging*, vol. 5, p. 48, 2019.
- [109] A. Magueresse, V. Carles, and E. Heetderks, "Low-resource languages: A review of past work and future challenges," *arXiv preprint arXiv:2006.07264*, 2020.
- [110] N. TS, "Word recognition in indic scripts," Ph.D. dissertation, IIIT Hyderabad, 2014. [Online]. Available: <https://cvit.iiit.ac.in/research/thesis/thesis-students/word-recognition-of-indic-scripts>
- [111] R. Kasturi, L. O'gorman, and V. Govindaraju, "Document image analysis: A primer," *Sadhana*, vol. 27, 2002.
- [112] M. Mathew and C. Jawahar, "An empirical study of ctc based models for ocr of indian languages," *arXiv preprint arXiv:2205.06740*, 2022.
- [113] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *ACL*, 2019.
- [114] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, *et al.*, "Improving language understanding by generative pre-training," 2018.
- [115] Y. Xu, M. Li, L. Cui, S. Huang, F. Wei, and M. Zhou, "Layoutlm: Pre-training of text and layout for document image understanding," in *ACM SIGKDD*, 2020.
- [116] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, "Language models are few-shot learners," *NeurIPS*, 2020.
- [117] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, *et al.*, "Palm: Scaling language modeling with pathways," *arXiv preprint arXiv:2204.02311*, 2022.
- [118] L. Zhao, E. Yu, Z. Ge, J. Yang, H. Wei, H. Zhou, J. Sun, Y. Peng, R. Dong, C. Han, and X. Zhang, "Chatspot: Bootstrapping multimodal llms via precise referring instruction tuning," 2023.
- [119] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *NeurIPS*, 2022.
- [120] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, 1989.
- [121] Y. LeCun, K. Kavukcuoglu, and C. Farnet, "Convolutional networks and applications in vision," in *ISCAS*, 2010.

- [122] D. C. Ciresan, U. Meier, L. M. Gambardella, and J. Schmidhuber, “Convolutional neural network committees for handwritten character classification,” in *ICDAR*, 2011.
- [123] T. Wang, D. J. Wu, A. Coates, and A. Y. Ng, “End-to-end text recognition with convolutional neural networks,” in *ICPR*, 2012.
- [124] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman, “Synthetic data and artificial neural networks for natural scene text recognition,” in *NeurIPS Deep Learning Workshop*, 2014.
- [125] —, “Deep structured output learning for unconstrained text recognition,” in *ICLR*, 2015.
- [126] G. H. Bakir, T. Hofmann, B. Schölkopf, A. J. Smola, B. Taskar, and S. V. N. Vishwanathan, Eds., *Predicting Structured Data*, ser. Neural Information Processing. MIT Press, 2007.
- [127] P. Krishnan, N. Sankaran, A. K. Singh, and C. V. Jawahar, “Towards a robust OCR system for indic scripts,” in *DAS*, 2014.
- [128] Z. Niu, G. Zhong, and H. Yu, “A review on the attention mechanism of deep learning,” *Neurocomputing*, vol. 452, 2021.
- [129] C. Lee and S. Osindero, “Recursive recurrent nets with attention modeling for OCR in the wild,” in *CVPR*, 2016.
- [130] B. Shi, X. Wang, P. Lyu, C. Yao, and X. Bai, “Robust scene text recognition with automatic rectification,” in *CVPR*, 2016.
- [131] A. Graves, S. Fernández, and J. Schmidhuber, “Multi-dimensional recurrent neural networks,” in *ICANN*, 2007.
- [132] Alex Graves and Jürgen Schmidhuber, “Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks,” in *NIPS*, 2008.
- [133] J. Puigcerver, “Are Multidimensional Recurrent Layers really necessary for handwritten text recognition?” in *ICDAR*, 2017.
- [134] T. Nguyen, M. Rosenberg, X. Song, J. Gao, S. Tiwary, R. Majumder, and L. Deng, “MS MARCO: A Human Generated MACHine Reading COMprehension Dataset.” *CoRR*, vol. abs/1611.09268, 2016.
- [135] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017.
- [136] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, “Xlnet: Generalized autoregressive pretraining for language understanding,” in *NeurIPS*, 2019.
- [137] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, *et al.*, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv preprint arXiv:1609.08144*, 2016.
- [138] D. Chen, A. Fisch, J. Weston, and A. Bordes, “Reading Wikipedia to Answer Open-Domain Questions,” in *ACL*, 2017.

- [139] K. M. Hermann, T. Kocisky, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, and P. Blunsom, “Teaching machines to read and comprehend,” *NeurIPS*, vol. 28, 2015.
- [140] Y. Zhu, O. Groth, M. Bernstein, and L. Fei-Fei, “Visual7w: Grounded question answering in images,” in *CVPR*, 2016.
- [141] Ł. Garncarek, R. Powalski, T. Stanisławek, B. Topolski, P. Halama, M. Turski, and F. Graliński, “LAMBERT: Layout-aware language modeling for information extraction,” in *ICDAR*, 2021.
- [142] Y. Xu, Y. Xu, T. Lv, L. Cui, F. Wei, G. Wang, Y. Lu, D. Florencio, C. Zhang, W. Che, M. Zhang, and L. Zhou, “LayoutLMv2: Multi-modal pre-training for visually-rich document understanding,” in *ACL — IJCNLP*, 2021.
- [143] S. Appalaraju, B. Jasani, B. U. Kota, Y. Xie, and R. Manmatha, “Docformer: End-to-end transformer for document understanding,” in *CVPR*, 2021.
- [144] R. Powalski, Ł. Borchmann, D. Jurkiewicz, T. Dwojak, M. Pietruszka, and G. Pałka, “Going full-tilt boogie on document understanding with text-image-layout transformer,” in *ICDAR*, 2021.
- [145] Y. Li, Y. Qian, Y. Yu, X. Qin, C. Zhang, Y. Liu, K. Yao, J. Han, J. Liu, and E. Ding, “Structext: Structured text understanding with multi-modal transformers,” in *ACM Multimedia*, 2021.
- [146] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei, “Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations,” *IJCV*, 2017.
- [147] D. Lewis, G. Agam, S. Argamon, O. Frieder, D. Grossman, and J. Heard, “Building a test collection for complex document information processing,” in *ACM SIGIR*, 2006.
- [148] Y. Jiang, V. Natarajan, X. Chen, M. Rohrbach, D. Batra, and D. Parikh, “Pythia v0. 1: the winning entry to the vqa challenge 2018,” *arXiv preprint arXiv:1807.09956*, 2018.
- [149] A. Singh, V. Natarjan, M. Shah, Y. Jiang, X. Chen, D. Parikh, and M. Rohrbach, “Towards vqa models that can read,” in *CVPR*, 2019.
- [150] L. H. Li, M. Yatskar, D. Yin, C. Hsieh, and K.-W. Chang, “VisualBERT: A Simple and Performant Baseline for Vision and Language,” *ArXiv*, vol. abs/1908.03557, 2019.
- [151] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, “VI-bert: Pre-training of generic visual-linguistic representations,” *ArXiv*, vol. abs/1908.08530, 2020.
- [152] Y. Chen, L. Li, L. Yu, A. E. Kholly, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, “UNITER: learning universal image-text representations,” *CoRR*, vol. abs/1909.11740, 2019.
- [153] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, “Making the v in vqa matter: Elevating the role of image understanding in visual question answering,” in *CVPR*, 2017.
- [154] D. Gurari, Q. Li, A. J. Stangl, A. Guo, C. Lin, K. Grauman, J. Luo, and J. P. Bigham, “Vizwiz grand challenge: Answering visual questions from blind people,” in *CVPR*, 2018.

- [155] A. F. Biten, R. Tito, A. Mafla, L. Gomez, M. Rusinol, E. Valveny, C. Jawahar, and D. Karatzas, “Scene text visual question answering,” in *ICCV*, 2019.
- [156] I. Krasin, T. Duerig, N. Alldrin, V. Ferrari, S. Abu-El-Haija, A. Kuznetsova, H. Rom, J. Uijlings, S. Popov, S. Kamali, M. Mallocci, J. Pont-Tuset, A. Veit, S. Belongie, V. Gomes, A. Gupta, C. Sun, G. Chechik, D. Cai, Z. Feng, D. Narayanan, and K. Murphy, “Openimages: A public dataset for large-scale multi-label and multi-class image classification.” *Dataset available from <https://storage.googleapis.com/openimages/web/index.html>*, 2017.
- [157] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *EMNLP*, 2014.
- [158] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comput.*, 1997.
- [159] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching word vectors with subword information,” *ACL*, vol. 5, 2017.
- [160] J. Almazán, A. Gordo, A. Fornés, and E. Valveny, “Word spotting and recognition with embedded attributes,” *TPAMI*, 2014.
- [161] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, “Stacked attention networks for image question answering,” in *CVPR*, 2016.
- [162] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap, “A simple neural network module for relational reasoning,” in *NeurIPS*, 2017.
- [163] A. Kembhavi, M. Seo, D. Schwenk, J. Choi, A. Farhadi, and H. Hajishirzi, “Are You Smarter Than A Sixth Grader? Textbook Question Answering for Multimodal Machine Comprehension,” in *CVPR*, 2017.
- [164] S. Yagcioglu, A. Erdem, E. Erdem, and N. Ikingler-Cinbis, “RecipeQA: A Challenge Dataset for Multimodal Comprehension of Cooking Recipes,” in *EMNLP*, 2018.
- [165] S Mori and C.Y Suen and K Yamamoto, “Historical review of OCR research and development,” in *Proceedings of IEEE*, 1992.
- [166] G Nagy, “At the frontiers of OCR,” in *Proceedings of IEEE*, 1992.
- [167] “Census 2011,” <https://censusindia.gov.in/2011-Common/CensusData2011.Html>, accessed on 1 November 2021.
- [168] R. M. K. Sinha and H. Mahabala, “Machine recognition of devnagari script,” in *IEEE Trans. on Systems, Man and Cybernetics*, 1979.
- [169] B. B. Chaudhuri and U. Pal, “A complete printed bangla ocr system,” *Pattern Recognition*, vol. 31, pp. 531–549, 1998.
- [170] D. Arya, T. Patnaik, S. Chaudhury, C. V. Jawahar, B.B.Chaudhuri, A.G.Ramakrishna, C. Bhagvati, and G. S. Lehal, “Experiences of Integration and Performance Testing of Multilingual OCR for Printed Indian Scripts,” in *MOCR Workshop, ICDAR*, 2011.

- [171] N. Sankaran and C. V. Jawahar, "Devanagari Text Recognition:A Transcription Based Formulation," in *ICDAR*, 2013.
- [172] U. Pal and B. B. Chaudhuri, "Indian Script Character Recognition: A Survey ," in *Pattern Recognition*, 2004.
- [173] S. Antanani and L. Agnihotri, "Gujarati Character Recognition," in *ICDAR*, 1999.
- [174] A. Negi, C. Bhagavati, and B. Krishna, "An OCR system for Telugu," in *ICDAR*, 2001.
- [175] U. Pal and A. Sarkar, "Recognition of printed Urdu script," in *ICDAR*, 2003.
- [176] G S Lehal and C Singh, "A Gurumukhi Script Recognition System," in *ICPR*, 2000.
- [177] K. Aparna and A. Ramakrishnan, "A complete Tamil Optical Character Recognition system," in *Document Analysis System*, 2002.
- [178] T. Ashwin and P. S. Sastry, "A font and size independent OCR system for printed kannada documents using support machines," *Sadhana*, 2002.
- [179] B. Vijay Kumar and A. G. Ramakrishnan, "Machine Recognition of Printed Kannada Text," in *Document Analysis Systems V*, 2002.
- [180] R. Sanjeev Kunte and R. Sudhaker Samuel, "A simple and efficient optical character recognition system for basic symbols in printed kannada text," *Sadhana*, vol. 32, no. 5, 2007.
- [181] C. V. Jawahar, MNSSK Pavan Kumar and S. S. Ravikiran, "A Bilingual OCR system for Hindi-Telugu Documents and its Applications," in *ICDAR*, 2003.
- [182] Ghosh, Subhankar and Bora, P. K. and Das, Sanjib and Chaudhuri, B. B., "Development of an Assamese OCR Using Bangla OCR," in *Proceeding of the Workshop on Document Analysis and Recognition*, 2012.
- [183] Venkat Rasagna, Jinesh K.J. and C.V. Jawahar, "On Multifont Character Classification in Telugu," in *International Conference on Information Systems for Indian Languages (ICISIL)*, 2011.
- [184] N. N.V. and C. Jawahar, "Emperical evaluation of Character Classification schemes," in *International Conference on Advances in Pattern Recognition*, 2009.
- [185] J. Makhoul, R. Schwartz, C. R. Christopher Lapre, and I. Bazzi, "Language-independent and segmentation-free techniques for optical character recognition," in *Document Analysis Systems II*, J. J. Hull and S. L. Taylor, Eds., 1996, pp. 67–82.
- [186] Z. Lu, R. M. Schwartz, P. Natarajan, I. Bazzi, and J. Makhoul, "Advances in the BBN BYBLOS OCR System," in *ICDAR*, 1999.
- [187] P. Natarajan, E. MacRostie, and M. Decerbo, "The BBN Byblos Hindi OCR system," in *DRR*, 2005.
- [188] Toni M. Rath and R. Manmatha, "Features for Word Spotting in Historical Manuscripts," in *ICDAR*, 2003.
- [189] V. Chavan, A. Malage, K. Mehrotra, and M. K. Gupta, "Printed text recognition using blstm and mdlstm for indian languages," in *ICIP*, 2017.
- [190] D. Paul and B. B. Chaudhuri, "A blstm network for printed bengali ocr system with high accuracy," *ArXiv*, vol. abs/1908.08674, 2019.

- [191] T. Kundaikar and J. D. Pawar, "Multi-font devanagari text recognition using lstm neural networks," in *International Conference on Sustainable Technologies for Computational Intelligence*, 2020.
- [192] A. Dwivedi, R. Saluja, and R. K. Sarvadevabhatla, "An ocr for classical indic documents containing arbitrarily long words," in *CVPR Workshops*, 2020.
- [193] C. Jawahar, A. Kumar, A. Phaneendra, and K. Jinesh, "Building data sets for indian language ocr research," *Guide to OCR for Indic Scripts: Document Recognition and Retrieval*, pp. 3–25, 2010.
- [194] A. Kumar and C. Jawahar, "Content-level annotation of large collection of printed document images," in *ICDAR*, 2007.
- [195] G. Harit, K. Jinesh, R. Garg, C. Jawahar, and S. Chaudhury, "Managing multilingual ocr project using xml," in *International Workshop on Multilingual OCR*, 2009, pp. 1–10.
- [196] M. Jain, M. Mathew, and C. V. Jawahar, "Unconstrained ocr for urdu using deep cnn-rnn hybrid networks," in *ACPR*, 2017, p. 6.
- [197] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International Journal of Computer Vision*, 2010.
- [198] A. Rozantsev, V. Lepetit, and P. Fua, "On rendering synthetic images for training an object detector," *CVIU*, 2015.
- [199] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez, "The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes," in *CVPR*, 2016.
- [200] K. P. Sankar, C. V. Jawahar, and R. Manmatha, "Nearest Neighbor based Collection OCR," in *Document Analysis Systems*, 2010.
- [201] J. A. Rodríguez-Serrano, F. Perronnin, J. Lladós, and G. Sánchez, "A similarity measure between vector sequences with application to handwritten word image retrieval," in *CVPR*, 2009.
- [202] N. Sabbour and F. Shafait, "A segmentation free approach to arabic and urdu ocr," *Proceedings of SPIE - The International Society for Optical Engineering*, vol. 8658, 2013.
- [203] P. Krishnan and C. V. Jawahar, "Generating synthetic data for text recognition," 2016.
- [204] A. Hannun, "Sequence modeling with ctc," *Distill*, 2017, <https://distill.pub/2017/ctc>.
- [205] L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, 1989.
- [206] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *NeurIPS*, 2019.
- [207] Hinton, "Neural Networks for Machine Learning, Lecture 6," http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf, 2012, accessed on 10 November 2021.

- [208] S. V. Rice and T. A. Nartker, “The isri analytic tools for ocr evaluation version 5.1,” <https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.216.9427>, 1996, accessed on 3 December 2021.
- [209] H. Scheidl, “WordDetector,” <https://github.com/githubharald/WordDetector>, 2021, accessed on 10 November 2021.
- [210] R. Manmatha and N. Srimal, “Scale space technique for word segmentation in handwritten documents,” in *Scale-Space Theories in Computer Vision*, 1999.
- [211] D. Doermann, K. Tombre, *et al.*, *Handbook of document image processing and recognition*. Springer, 2014.
- [212] I. Kavasidis, C. Pino, S. Palazzo, F. Rundo, D. Giordano, P. Messina, and C. Spampinato, “A saliency-based convolutional neural network for table and chart detection in digitized documents,” in *ICIAP*, 2019.
- [213] R. B. Palm, O. Winther, and F. Laws, “Cloudscan—a configuration-free invoice analysis system using recurrent neural networks,” in *ICDAR*, 2017.
- [214] A. Singh, V. Goswami, V. Natarajan, Y. Jiang, X. Chen, M. Shah, M. Rohrbach, D. Batra, and D. Parikh, “MMF: A multimodal framework for vision and language research,” <https://github.com/facebookresearch/mmf>, 2020.
- [215] A. F. Biten, R. Tito, A. Mafla, L. Gómez, M. Rusiñol, M. Mathew, C. V. Jawahar, E. Valveny, and D. Karatzas, “ICDAR 2019 competition on scene text visual question answering,” *CoRR*, vol. abs/1907.00490, 2019.
- [216] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, and J. Brew, “Huggingface’s transformers: State-of-the-art natural language processing,” *ArXiv*, vol. abs/1910.03771, 2019.
- [217] M. Joshi, E. Choi, D. Weld, and L. Zettlemoyer, “TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension,” in *ACL*, 2017.
- [218] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *ICLR*, 2015.
- [219] M. Mathew, R. Tito, D. Karatzas, R. Manmatha, and C. V. Jawahar, “Document visual question answering challenge 2020,” in *DAS short papers*, 2020.
- [220] R. Chaudhry, S. Shekhar, U. Gupta, P. Maneriker, P. Bansal, and A. Joshi, “Leaf-qa: Locate, encode attend for figure question answering,” in *WACV*, 2020.
- [221] R. Tanaka, K. Nishida, and S. Yoshida, “VisualMRC: Machine reading comprehension on document images,” in *AAAI*, 2021.
- [222] J. Lu, D. Batra, D. Parikh, and S. Lee, “ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks,” 2019.
- [223] H. Tan and M. Bansal, “Lxmert: Learning cross-modality encoder representations from transformers,” *arXiv preprint arXiv:1908.07490*, 2019.

- [224] Z. Bylinskii, S. Alsheikh, S. Madan, A. Recasens, K. Zhong, H. Pfister, F. Durand, and A. Oliva, “Understanding infographics through textual and visual tag prediction,” *ArXiv*, vol. abs/1709.09215, 2017.
- [225] S. Madan, Z. Bylinskii, M. Tancik, A. Recasens, K. Zhong, S. Alsheikh, H. Pfister, A. Oliva, and F. Durand, “Synthetically trained icon proposals for parsing and summarizing infographics,” *arXiv preprint arXiv:1807.10441*, 2018.
- [226] N. Landman, “Towards abstractive captioning of infographics,” Master’s thesis, Massachusetts Institute of Technology, Massachusetts Institute of Technology, 2018.
- [227] M. A. Borkin, A. A. Vo, Z. Bylinskii, P. Isola, S. Sunkavalli, A. Oliva, and H. Pfister, “What Makes a Visualization Memorable?” *IEEE Transactions on Visualization and Computer Graphics*, vol. 19, 2013.
- [228] M. Lu, C. Wang, J. Lanir, N. Zhao, H. Pfister, D. Cohen-Or, and H. Huang, “Exploring Visual Information Flows in Infographics,” in *ACM CHI*, 2020.
- [229] M. A. A. Cox and T. F. Cox, *Multidimensional Scaling*. Springer Berlin Heidelberg, 2008, pp. 315–347.
- [230] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent dirichlet allocation,” *J. Mach. Learn. Res.*, vol. 3, p. 993–1022, 2003.
- [231] R. Řehůřek and P. Sojka, “Software Framework for Topic Modelling with Large Corpora,” in *LREC 2010 Workshop on New Challenges for NLP Frameworks*, 2010.
- [232] A. W. Harley, A. Ufkes, and K. G. Derpanis, “Evaluation of deep convolutional nets for document image classification and retrieval,” in *ICDAR*, 2015.
- [233] C. Wigginton, C. Tensmeyer, B. L. Davis, W. A. Barrett, B. L. Price, and S. Cohen, “Start, Follow, Read: End-to-End Full-Page Handwriting Recognition,” in *ECCV*, 2018.
- [234] L. Kang, P. Riba, M. Rusiñol, A. Fornés, and M. Villegas, “Pay Attention to What You Read: Non-recurrent Handwritten Text-Line Recognition,” 2020.
- [235] G. T. Bazzo, G. A. Lorentz, D. Suarez Vargas, and V. P. Moreira, “Assessing the Impact of OCR Errors in Information Retrieval,” in *Advances in Information Retrieval*, J. M. Jose, E. Yilmaz, J. Magalhães, P. Castells, N. Ferro, M. J. Silva, and F. Martins, Eds., 2020.
- [236] G. Chiron, A. Doucet, M. Coustaty, M. Visani, and J.-P. Moreux, “Impact of OCR Errors on the Use of Digital Libraries: Towards a Better Access to Information,” in *ACM/IEEE JCDL*, 2017.
- [237] K. J. Shih, S. Singh, and D. Hoiem, “Where To Look: Focus Regions for Visual Question Answering,” 2016.
- [238] X. Wang, Y. Liu, C. Shen, C. C. Ng, C. Luo, L. Jin, C. S. Chan, A. van den Hengel, and L. Wang, “On the General Value of Evidence, and Bilingual Scene-Text Visual Question Answering,” 2020.
- [239] A. Trischler, T. Wang, X. Yuan, J. Harris, A. Sordoni, P. Bachman, and K. Suleman, “Newsqa: A machine comprehension dataset,” *CoRR*, vol. abs/1611.09830, 2016.
- [240] T. S. Jaakkola and D. Haussler, “Exploiting Generative Models in Discriminative Classifiers,” in *NeurIPS*, 1999.

- [241] F. Perronnin and C. R. Dance, “Fisher Kernels on Visual Vocabularies for Image Categorization.” in *CVPR*, 2007.
- [242] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek, “Image Classification with the Fisher Vector: Theory and Practice,” *Int. J. Comput. Vision*, 2013.
- [243] F. Perronnin, J. Sánchez, and T. Mensink, “Improving the Fisher Kernel for Large-scale Image Classification,” in *ECCV*, 2010.
- [244] H. Jegou, F. Perronnin, M. Douze, J. Sánchez, P. Perez, and C. Schmid, “Aggregating Local Image Descriptors into Compact Codes,” *TPAMI*, 2012.
- [245] G. Bradski, “The OpenCV Library,” *Dr. Dobbs’s Journal of Software Tools*, 2000.
- [246] M. Villegas, J. Puigcerver, A. Toselli, J.-A. Sánchez, and E. Vidal, “Overview of the ImageCLEF 2016 Handwritten Scanned Document Retrieval Task,” in *CLEF*, 2016.
- [247] E. Loper and S. Bird, “NLTK: The Natural Language Toolkit,” in *ACL ETMTNLP*, 2002.
- [248] U.-V. Marti and H. Bunke, “The IAM-database: an English sentence database for offline handwriting recognition,” *IJDAR*, 2002.
- [249] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2007.
- [250] C. Gormley and Z. Tong, *Elasticsearch: The Definitive Guide*. O’Reilly Media, Inc., 2015.
- [251] H. Cao, V. Govindaraju, and A. Bhardwaj, “Unconstrained handwritten document retrieval,” *IJDAR*, 2011.
- [252] Y. Fataïcha, M. Cheriet, J. Y. Nie, and C. Y. Suen, “Retrieving Poorly Degraded OCR Documents,” *IJDAR*, vol. 8, 2006.
- [253] R. Tito, M. Mathew, C. Jawahar, E. Valveny, and D. Karatzas, “Icdar 2021 competition on document visual question answering,” in *ICDAR*, 2021.