

Matching Handwritten Document Images

Praveen Krishnan and C.V Jawahar

CVIT, IIIT Hyderabad, India

praveen.krishnan@research.iiit.ac.in, jawahar@iiit.ac.in

Abstract. We address the problem of predicting similarity between a pair of handwritten document images written by potentially different individuals. This has applications related to matching and mining in image collections containing handwritten content. A similarity score is computed by detecting patterns of text re-usages between document images irrespective of the minor variations in word morphology, word ordering, layout and paraphrasing of the content. Our method does not depend on an accurate segmentation of words and lines. We formulate the document matching problem as a structured comparison of the word distributions across two document images. To match two word images, we propose a convolutional neural network (CNN) based feature descriptor. Performance of this representation surpasses the state-of-the-art on handwritten word spotting. Finally, we demonstrate the applicability of our method on a practical problem of matching handwritten assignments.

Keywords: Handwritten word spotting, CNN features, plagiarism detection

1 Introduction

Matching two document images has several applications related to information retrieval like spotting keywords in historical documents [8], accessing personal notes [22], camera based interface for querying [45], retrieving from video databases [27], automatic scoring of answer sheets [40], and mining and recommending in health care documents [25]. Since OCRs do not reliably work for all types of documents, one resorts to image based methods for comparing textual content. This problem is even more complex when considering unconstrained handwritten documents due to the high variations across the writers. Moreover, variable placement of the words across the documents makes a rigid geometric matching ineffective. In this work, we design a scheme for matching two handwritten document images. The problem is illustrated in Fig. 1(a). We validate the effectiveness of our method on an application, named as measure of document similarity (MODS).¹ MODS compares two handwritten document images and provides a normalized score as a measure of similarity between two images.

¹ In parallel to measure of software similarity (MOSS) [36], which has emerged as the de facto standard across the universities to compare software solutions from students.

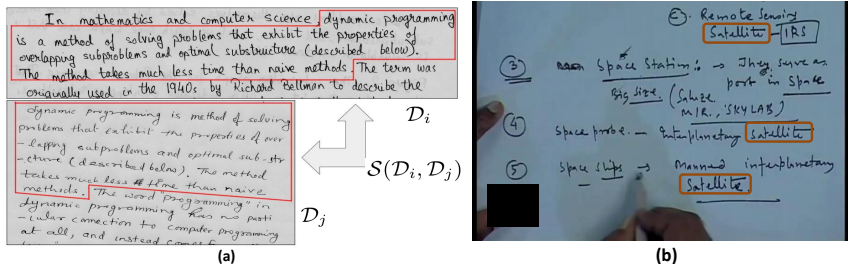


Fig. 1. (a) Given two document images D_i and D_j , we are interested in computing a similarity score $S(D_i, D_j)$ which is invariant to (i) writers, (ii) word flow across lines, (iii) spatial shifts, and (iv) paraphrasing. In this example, the highlighted lines from D_i and D_j have almost the same content but they widely differ in terms of spatial arrangement of words. (b) Query-by-text results on searching with “satellite” on an instructional video. The spotted results are highlighted in the frame.

Text is now appreciated as a critical information in understanding natural images [14, 26, 48]. Attempts for wordspotting in natural images [48] have now matured to end-to-end frameworks for recognition and retrieval [14, 16, 47]. Natural scene text is often seen as an isolated character image sequence in arbitrary view points or font styles. Recognition in this space is now becoming reliable, especially with the recent attempts that use CNNs and RNNs [14, 42]. However, handwritten text understanding is still lacking in many aspects. For example, the best performance on the word spotting (or retrieval) on the popular IAM data set [24] is an mAP of 0.55 [3]. In this work, we improve this to 0.80. We achieve this with the help of synthetic handwritten data that now enables the exploitation of deep learnt representations for handwritten data.

Word Spotting. Initial attempts for matching handwritten words were based on DTW [32] and HMM [9, 34] over variable length feature representations. Although these models were flexible, they were not really scalable. Many approaches such as [2, 29, 35] demonstrated word spotting using fixed length representation based on local features such as SIFT and HOG in a bag of words (BOW) framework. Most of these works employed better feature representations such as Fisher vectors [2, 29], latent semantic indexing [35] and techniques such as query expansion and re-ranking for enhancing the performance. However, the applicability of these methods are still limited for multi-writer scenarios. Recently, Almazán *et al.* [3] proposed a label embedding and attributes learning framework where both word images and text strings are embedded into a common subspace with an associated metric to compare both modalities.

Matching documents. Matching textual documents is a well studied problem in text processing [23] with applications in plagiarism detection in electronic documents [31]. For softwares, MOSS [36] provides a solution to compare two programs and is robust against a set of alterations e.g., formatting and changes in variable names. However, when the documents are scanned images, these methods can not be directly applied. There have been some attempts [4, 17] to find

Query Top ranked retrieval							Scope
							(a) [2,14,17,19] & Sec. 3
							(b) This work

Fig. 2. Word spotting vs. normalized word spotting. (a) shows the conventional word spotting task while (b) extends the task to retrieve semantically similar words using a normalized representation. Here we deal with popular inflectional ending present due to agglutinative property of a language.

duplicate and near duplicates in multimedia databases. However, they are not directly applicable to documents where the objective is to compare images based on the textual content. For printed documents, matching based on geometry or organization of a set of keypoints has been successful [10, 44, 46]. This works well for duplicate as well as cut-and-paste detection in printed documents. However, due to unique set of challenges in handwritten documents such as wide variation of word styles, the extraction of reliable keypoints with geometric matching is not very successful. Other major challenges include paraphrasing of the textual content, non-rigidity of word ordering which leads to word overflows across lines. In our proposed method, we use locality constraints to achieve invariance to such variations. We also extend the word spotting to take care of the popular word morphological variations in the image space as shown in Fig. 2(b). The proposed features can associate similarity between word images irrespective of word morphological variations due to changes in tense and voice of the sentence construction. In the context of retrieval systems it improves overall recall and helps in matching documents in a semantic space.

Contributions. In this work, we compute a similarity score by detecting patterns of text re-usages across documents written by different individuals irrespective of the minor variations in word forms, word ordering, layout or paraphrasing of the content. In the process of comparing two document images, we design a module that compares two handwritten words using CNN features and report a 56% error reduction in word spotting task on the challenging dataset of IAM [24] and pages from George Washington (GW) collection [9]. We also propose a normalized feature representation for word images which is invariant to different inflectional endings or suffixes present in words. The advantage of our matching scheme is that it does not require an accurate segmentation of the documents. To calibrate the similarity score with that of human perception, we conduct a human experiment where a set of individuals are advised to create similar documents with natural variations. Our solution reports a score that match the human evaluation with a mean normalized discounted cumulative gain ($nDCG$) of 0.89. Finally, we demonstrate two immediate applications (i) searching handwritten text from instructional videos, and (ii) comparing handwritten assignments. Fig. 1(a,b) shows a sample result from these applications.

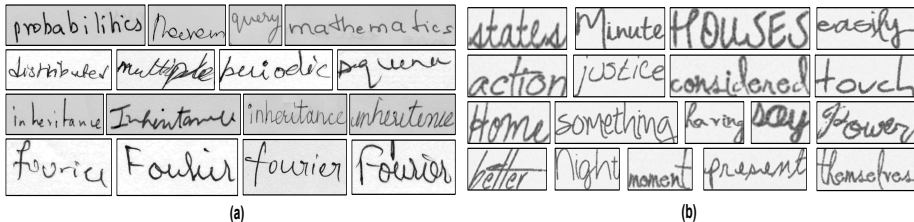


Fig. 3. (a) The top two rows show the variations in handwritten images, the bottom two rows demonstrate the challenges of intra class variability in images across writers. (b) Sample images from the IIIT-HWS dataset created as part of this work to address the lack of training data for for learning complex CNN networks.

2 CNN features for handwritten word images

The proposed document image matching scheme employs a discriminative representation for comparing two word images. Such a representation needs to be invariant to (i) both inter and intra class variability across writers, (ii) presence of skew, (iii) quality of ink, and (iv) quality and resolution of the scanned image. Fig. 3(a) demonstrates the challenges in matching across writers and documents. The top two rows show the variations across images in which some are even hard for humans to read without enough context of nearby words. The bottom two rows show different instances of same word written by different writers, e.g., “inheritance” and “Fourier” where one can clearly notice the variability in shape for each character in the word image. In this work we use convolutional neural networks (CNN) motivated by the recent success of deep neural networks [6, 15, 18, 39, 43] and the availability of better learning schemes [12, 13]. Even though CNN architectures such as [19, 38] were among the first to show high performing classifier for MNIST handwritten digits, application of such ideas for unconstrained continuous handwritten words or documents has not been demonstrated possibly due to the lack of data, and also the lack of appropriate training schemes.

2.1 IIIT-HWS dataset

To address the lack of data for training handwritten word images, we build a synthetic handwritten dataset of 1 million word images. We call this dataset as IIIT-HWS. Some of the sample images from this dataset are shown in Fig. 3(b). Note that these images are very similar to natural handwriting. IIIT-HWS dataset is formed out of 750 publicly available handwritten fonts. We use a subset of Hunspell dictionary and pick a unique set of 10K words for this purpose. For each word, we randomly sample 100 fonts and render its corresponding image. During this process, we vary the following parameters: (i) kerning level (inter character space), (ii) stroke width, and (iii) mean foreground and background pixel distributions. We also perform Gaussian filtering to smooth the final rendered image.

Moreover, we prefer to learn a case insensitive model for each word category, hence we perform three types of rendering, namely, all letters capitalized, all letters lower and only the first letter in caps.²

2.2 HWNet architecture and transfer learning

The underlying architecture of our CNN model (HWNet) is inspired from [15]. We use a CNN with five convolutional layers with 64, 128, 256, 512 and 512 square filters with dimensions: 5, 5, 3, 3 and 3 respectively. The next two layers are fully connected ones with 2048 neurons each. The last layer uses a fully connected (FC) layer with dimension equal to number of classes, 10K in our case, and is further connected to the softmax layer to compute the class specific probabilities. Rectified linear units are used as the non-linear activation units after each weight layer except the last one, and 2×2 max pooling is applied after first, second, and fourth convolutional layers. We use a stride of one and padding is done to preserve the spatial dimensionality. We empirically observed that the recent approach using batch normalization [13] for reducing the generalization error, performed better as compared to dropouts. The weights are initialized randomly from normal distribution, and during training the learning rate is reduced on a log space starting from 0.1. The input to the network is a gray scale word image of fixed size 48×128 . HWNet is trained on the IIIT-HWS dataset with 75-15-10% train-validation-test split using a multinomial logistic regression loss function to predict the class labels, and the weights are updated using mini batch gradient descent algorithm with momentum.

Transfer learning. It is well-known that off-the-shelf CNNs [7, 33] trained for a related task could be adapted or fine-tuned to obtain reasonable and even state-of-the-art performance for new tasks. In our case we prefer to perform a transfer learning from synthetic domain (IIIT-HWS) to real world setting. Here we use popular handwritten labeled corpora such as IAM and GW to perform the transfer learning. It is important to keep the learning rates low in such setting, else the network quickly unlearns the generic weights learned in the initial layers. In this work, we extract the features computed from the last FC layer to represent each handwritten word image.

3 Normalized word spotting

Word spotting [3, 22] has emerged as a popular framework for search and retrieval of text in images with applications in retrieving text from historical manuscripts, handwritten documents where the performance of optical character recognition (OCR) is still limited. It is typically formulated as a retrieval problem where the query is an exemplar image (query-by-example) and the task is to retrieve all word images with similar content. It uses holistic word image representation [2,

² More details on dataset, codes and trained CNN models are available at: <http://cvit.iiit.ac.in/research/projects/cvit-projects/matchdocimgs>

22] which does not demand character level segmentation and the retrieval is performed using nearest neighbor search. Fig. 2(a) shows a word spotting result which retrieves similar word images for the query “looked”. In this work, our interest lies in finding the document similarity between a pair of handwritten documents written by different writers in an unconstrained setting. We observe that such a problem can be addressed in a word spotting framework where the task would be to match similar words between a pair of documents using the proposed CNN features for handwritten word images.

HWNet provides a generic representation for word spotting by retrieving word images with the exact content written. While addressing the larger problem of document retrieval, on similar lines of a text based information retrieval pipeline, we relax this constraint and prefer to retrieve not just similar or exact words but also their common variations. These variations are observed in languages due to morphology. In English, we observe such variations in the form of inflectional endings (suffixes) such as “-s (plural), -ed (past tense), -ess (adjective), -ing (continuous form)” etc. These suffixes are added to the root word, and thereby resulting in a semantically related word. A stemmer, such as the Porter stemmer [30] can strip out common suffixes which generates a normalized representation of words with common roots. We imitate the process of stemming in the visual domain by labeling the training data in terms of root words given by the Porter stemmer, and use the HWNet architecture to learn a normalized representation which captures the visual representation of word images along with the invariance to its inflectional endings. We observe that such a network learns to give less weights to popular word suffixes and gives a normalized representation which is better suited for document image retrieval tasks. Fig. 2(b) shows the normalized word spotting results obtained using the proposed features that includes both “similar” and “semantically-similar” results, e.g., “look”, “looks”, “looking” and “looked”.

4 Measure of document similarity (MODS)

Matching printed documents for retrieving the original documents and detecting cut-and-paste for finding plagiarism were attempted in the past by computing interest points in word images and their corresponding matches [10, 44]. However, handwritten documents have large intra class variability to reliably detect interest points. In addition, the problem of word-overflow in which words from the right end of the document overflow and appear on the left end of the next line make the matching based on rigid geometry infeasible. We state our problem as follows: *given a pair of document images, compute a similarity score by detecting patterns of text re-usages between documents irrespective of the minor variations in word morphology, word ordering, layout and paraphrasing of the content.* Our matching scheme is broadly split into two stages. The first stage involves segmentation of document into multiple possible word bounding boxes while the later stage computes a structured document similarity score which obeys loose word ordering and its content.

4.1 Document segmentation

A document image contains structured objects. The objects here are the words and structure is the order in which words are presented. Segmentation of a handwritten document image into constituent words is a challenging task because of the unconstrained nature of documents such as variable placements of page elements (figures, graphs and equations), presence of skewed lines, and irregular kerning. Most of the methods such as [11, 41] are bottom-up approaches with tunable parameters to arrive at a unique segmentation of words and lines. Considering the complexity of handwritten documents, we argue that a reasonably practical system, should work with multiple possible lines and word segmentation proposals with a high recall. We use a simple multi-stage bottom-up approach similar to [20] by forming three sets of connected components (CCs) on the binarized image based on its sizes. CCs in s_1 set contains area less than 0.1μ , s_3 set contains CCs having area large than $\mu + 2\sigma$ while remaining CCs are categorized as s_2 . Here μ, σ is mean and standard deviation respectively. The small (s_1), medium (s_2) and large (s_3) CC sets are assumed to be punctuation, actual characters and high probable line merge respectively. We associate each component in s_2 with its adjacent component if the cost given by: $Cost(i, j) = OL(i, j) + D(i, j) + \theta(i, j)$, is above a certain threshold. Here i, j are two components, OL is the amount of overlap in y -axis which is given by intersection over union, D is the normalized distance between the centroids of the i^{th} and j^{th} component, and $\theta(i, j)$ gives the angle between the centroids of the components. After the initial assignment, we now associate the s_3 components by checking whether these components intersects in the path of detected lines. In such a case, we slice the component horizontally and join it to the top and the bottom line respectively. Finally the components present in s_1 are associated with nearest detected lines. Given the bounding boxes of a set of CCs and its line associations, we analyse the inter CC spacing and derive multiple thresholds to group it into words. This results in multiple word bounding box hypotheses with a high recall. Minor reduction in the precision at this stage is taken care by our matching scheme.

4.2 SWM matching

We first define a similarity score between a pair of documents as the sum of word matches (SWM). We use l_2 normalized CNN descriptors of the corresponding words images w_k and w_l and compute the l_2 distance d_{kl} . We define the document similarity as the symmetric distance between the best word matches across the documents as follows:-

$$\mathcal{S}_N(\mathcal{D}_i, \mathcal{D}_j) = \frac{1}{|\mathcal{D}_i| + |\mathcal{D}_j|} \left(\sum_{w_k \in \mathcal{D}_i} \min_{w_l \in \mathcal{D}_j} d_{kl} + \sum_{w_l \in \mathcal{D}_j} \min_{w_k \in \mathcal{D}_i} d_{lk} \right). \quad (1)$$

This is a normalized symmetric distance where $|\mathcal{D}_i|$ is the number of words in the document \mathcal{D}_i . In order to reduce the exhaustive matches, we use an approximate nearest neighbor search using KD trees.

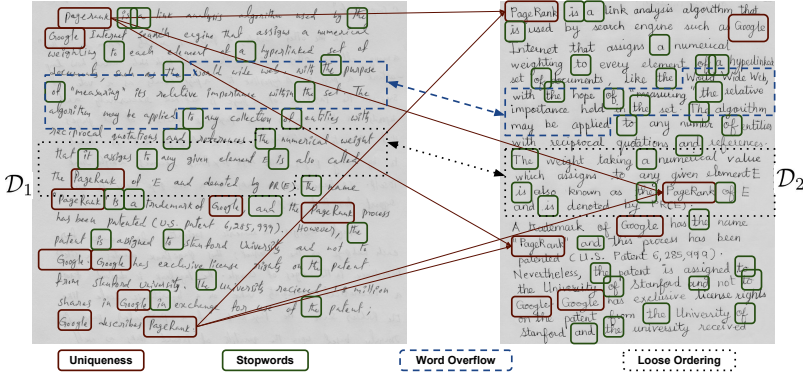


Fig. 4. A few major challenges of the matching process between a pair of documents \mathcal{D}_1 and \mathcal{D}_2 . (i) Finding a unique match of each potential word, (ii) removal of stopwords, (iii) invariance to word overflow problems, and (iv) exploiting the loose ordering of words in matching.

4.3 MODS matching

The problem of document matching and devising a scheme to compute similarity score is a challenging task. This problem along with the challenges is illustrated in Fig. 4. We address these problems along with their solution at two levels: (i) individual word matches, and (ii) bringing locality constraints.

Word matches. (i) Alternations: In general, the pair of documents of interest need not have the same content and hence, not all words need to have a correspondence in the second image. We enforce this with a simple threshold γ on the distance used for matching. (ii) Stopwords: The presence of stopwords in documents acts as a noise which corrupts the matching process for any IR system due their high frequency. In Fig. 4 we show some of these words in dark green boxes. We observed that the trained HWNet is reasonably robust in classifying stopwords due to their limited number and increased presence in training data. Therefore, we could take the softmax scores (probabilities) from last layer of HWNet and classify a word image as a stopword if the scores of one of stopword classes is above a certain threshold.

Locality constraints. The following three major challenges are addressed using locality constraints in the matching process. We first list the challenges and later propose the solution given by MODS. (i) Uniqueness: Though a word in the first image can match with multiple images in the second image, we are interested in a unique match. In Fig. 4 the highlighted words in dark red such as “Google” and “PageRank” occur at multiple places in both documents but the valid matches needs to be unique that obeys the given locality. (ii) Word overflow: As we deal with documents of unconstrained nature, similar sentences across different documents can span variable number of lines, a property of an individual writing style. In terms of geometry of position of words this results in a major shift of words (from right extreme to the left extreme). One such

pair of occurrence is shown in Fig. 4 as blue colored dashed region. We refer to this problem as word overflow. (iii) Loose ordering: Paraphrasing of the words as shown in the Fig. 4 as black dashed rectangle, is a common technique to conceal the act of copying where one changes the order of the words keeping the semantics intact.

We observe that the most informative matching words are the ones which preserve the consistency within a locality. We enforce locality constraints by splitting the document into multiple overlapping rectangular regions. The idea is to find out the best matching pairs of regions within two documents and associate them with individual word matches. For finding the cost of associating two rectangular regions, we formulate the problem as a weighted bipartite graph matching where the weights are the cosine distances of word images in feature space. We use the popular Hungarian algorithm to compute the cost of word assignments, which leads to a one to one mapping of word images between a pair of regions. The score computed between a pair of rectangular regions denoted as p and q from documents \mathcal{D}_i and \mathcal{D}_j respectively is given by:

$$Score(p) = \max_{q \in R(\mathcal{D}_j)} \left(\frac{\sum_{(k,l) \in Matches(p,q)} (1 - d_{kl})}{\max(|p|, |q|)} \right), \forall p \in R(\mathcal{D}_i), \quad (2)$$

where, $R(\mathcal{D}_j)$ denotes the set of all rectangular regions in a document image and $|\cdot|$ denotes number of words in the region. The function $Matches(p, q)$ returns the assignments given by the Hungarian algorithm. Finally, the normalized MODS score for a pair of documents is defined as follows:

$$\mathcal{S}_M(\mathcal{D}_i, \mathcal{D}_j) = \frac{\sum_{p \in R(\mathcal{D}_i)} Score(p)}{\max(|\mathcal{D}_i|, |\mathcal{D}_j|)}. \quad (3)$$

5 Experiments

In this section, we empirically evaluate the proposed CNN representation for the task of word spotting on standard datasets. We validate the effectiveness of these features on newer tasks such as retrieving semantically similar words, searching keywords from instructional videos and finally demonstrate the performance of the MODS algorithm for finding similarity between documents on annotated datasets created for this purpose.

5.1 Word-spotting

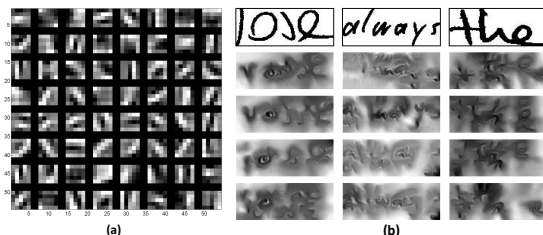
We perform word spotting in a *query-by-example* setting. We use IAM [24] and George Washington [9] (GW) dataset, popularly used in handwritten word spotting and recognition tasks. In case of the IAM dataset, we use the standard partition for training, testing, and validation provided along with the corpus. For GW dataset, we use a random set of 75% for training and validation, and the remaining 25% for testing. Each word image except the stop words in the test

Table 1. Quantitative results on word spotting using the proposed CNN features along with comparisons with various existing hand designed features on IAM and GW dataset.

Dataset	DTW[3]	SC-HMM[34]	FV[3]	EX-SVM[2]	KCCA[1]	KCSR[3]	Ours
GW	0.6063	0.5300	0.6272	0.5913	0.8563	0.9290	0.9484
IAM	0.1230	-	0.1566	-	0.5478	0.5573	0.8061

corpus is taken as the query to be ranked across all other images from the test corpus including stop-words acting as distractions. The performance is measured using the standard evaluation measure namely, mean Average Precision (mAP). HWNet architecture is fine-tuned using the respective standard training set for each test scenario. Table 1 compares the proposed features from state-of-the-art methods on these datasets. The results are evaluated in a case-insensitive manner as used in previous works [2, 3]. The proposed CNN features clearly surpasses the current state-of-the-art method [3] on IAM and GW, reducing the error rates by $\sim 56\%$ and $\sim 27\%$ respectively. This demonstrates the invariance of features for both multi-writer scenario (IAM) and historical documents (GW). Some of the qualitative results are shown in the top three rows of Fig. 6(a). One can observe the variability of each retrieved result which demonstrates the robustness the proposed features.

Visualizations. Fig. 5 shows the visualization of the trained HWNet architecture using popular schemes demonstrated in [18, 21]. Fig. 5(a) visualizes the weights of the first layer which bears a resemblance to Gabor filters and detects edges in different orientations. Fig. 5(b) demonstrates the visualization from a recent method [21] which inverts the CNN encoding back to image space and arrives at possible images which have high degree of probability for that encoding. This gives a better intuition of the learned layers and helps in understanding the invariances of the network. Here, we show the query images on the first row and its reconstruction in the following rows. One can observe that in almost all reconstructions there are multiple translated copies of the characters present in the word image along with some degree of orientations. Similarly, we can see the network is invariant to the first letter being in capital case (see Label: “the” at Col:3, Row:4) which was part of the training process. The reconstruction of the first image (see Label: “rose” at Col:1, Row:1) shows that possible recon-

**Fig. 5.** Visualization: (a) The weights of first layer of HWNet. (b) Four possible reconstructions [21] of three sample word images shown in columns. These are re-constructed from the representation of final layer of HWNet.

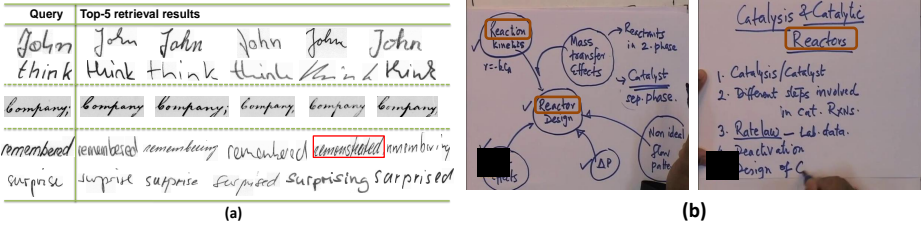


Fig. 6. Qualitative results: (a) Query-by-example results for the task of word spotting results on IAM and GW dataset. The bottom two rows shows results from normalized feature representation where one can observe we are also able to retrieve words with related meanings. (b) Query-by-text results on searching with “reactor” on an instructional video. The top two results are shown along with the spotted words which are highlighted in the frame.

struction images includes Label: “rose” (Col:1, Row:2) and “jose” (Col:1, Row:3) since there is an ambiguity in the query image.

5.2 Enhancements and applications

We now analyse the performance of the normalized features for retrieving semantically similar words which has not been yet attempted in handwritten domain and plays an important role in matching similar documents. We also demonstrate an application of MODS framework in a collection of instructional videos by retrieving relevant frames corresponding to user queries.

Normalized word spotting. Table 2 shows the quantitative results of the normalized (CNN_{Norm}) features which are invariant to common word inflectional endings and thereby learn features for stem or the root part of the word image. For this experiment, we update the evaluation scheme (ref. as inexact) to include not only similar word images but also the word images having common stem. We use Porter stemmer [30] for calculating the stem of a word. Table 2 also compares the performance of CNN features used in Sec. 5.1 and validate it over inexact evaluation. Here we obtain a reduced mAP of 0.7170 whereas using the normalized features, we improve the mAP to 0.7443. We also observe that using normalized features for exact evaluation results in a comparable performance (0.7955) which motivates us to use them in document similarity problems. In Fig. 6(a), the bottom two rows shows qualitative results using these normalized features. The retrieval results for query “surprise” contains the word “surprised”, “surprising” along with the keyword “surprise”.

Evaluation	CNN	CNN_{Norm}
Exact	0.8061	0.7955
Inexact	0.7170	0.7443

Table 2. Word spotting results using normalized features and its comparisons with exact features.

Table 3. Quantitative evaluation of various matching schemes on HW-DocSim dataset. We compare the performance of proposed MODS framework using CNN features over baseline methods such as NN, BOW, and embedded attributes proposed in [3].

Method	NN	BOW	SWM	MODS	SWM	MODS
Feature	Profile	SIFT	KCSR [3]		CNN	
$nDCG@99$	0.5856	0.6128	0.7968	0.8444	0.8569	0.8993
AUC	0.5377	0.4516	0.8231	0.8302	0.9465	0.9720

Searching in instructional videos. To demonstrate the effectiveness and generalization ability of the proposed CNN features we performed an interesting task of searching inside instructional videos where the tutor write handwritten text to aid students in the class. We conducted the experiment in a query-by-text scenario where the query text is synthesized into a word image using one of the fonts used in the IIT-HWS dataset. We took five popular online course videos from NPTEL [28] on different topics from YouTube and manually extracted frames containing textual regions. For each frame, we obtained multiple segmentation output from the proposed segmentation method. For evaluation, we handpicked 20 important queries and labeled the frames containing them. We obtained a frame level mAP of 0.9369 on this task. Fig. 6(b) shows the top-2 matching frames for the query “reactor” along with the spotted words. One can observe that along with retrieving exact matches, we also retrieve similar keywords such as “Reactors”, and “Reaction”.

5.3 HW-DocSim dataset and evaluations

We start with the textual corpus presented in [5] for plagiarism detection. The corpus contains plagiarized short answers to five unique questions given to 19 participants. Hence the corpus contains around 100 documents of which 95 were created in a controlled setting while five were the original answers (source document) which were given to participants to refer to and copy. There are four types or degree of plagiarism introduced in this collection: (i) *near copy*, where the content is an exact copy from different parts from the source; (ii) *light revision*, where the content is taken from source but with slight revisions such as replacing words with synonyms, (iii) *heavy revision*, which includes heavy modification such as paraphrasing, combining or splitting sentences and changing the order; and (iv) *non-plagiarized*, where the content is prepared independently on the same topic. For the task of generating handwritten document images, we included a total of 24 students and asked them to write on plain white sheets of paper. For each document we use a separate student to avoid any biases in writing styles. To keep the content close to its natural form, we did not mention any requirements on spacing between words, and lines, and did not put any constraints on the formatting of text in the form of line breaks and paragraphs. In case of mistakes, the written word was striked out and writing was continued.

Evaluation methodology. To evaluate the performance, we took all source-candidate document pairs and computed their similarity scores. Here we only verify whether the document is similar (plagiarized) or not while discarding the amount of plagiarism. The performance is measured using area under the ROC curve (AUC) by sorting the scores of all pairs. In another experiment, we compute graded similarity measure in accordance to each source document posed as a query which expects the ranking according to the degree of copying. Here we use normalized discounted cumulative gain ($nDCG$), a measure used frequently in information retrieval when grading is needed. Here the query is presented as the *source* document and the target documents are all documents present in the corpus. The discounted cumulative gain (DCG) at position p is given as $DCG_p = \sum_{i=1}^p (2^{rel_i} - 1) / (\log_2(i + 1))$ where rel_i is the ground truth relevance for the document at rank i . In our case, the relevance measures are represented as: 3 - *near copy*, 2 - *light revision*, 1 - *heavy revision*, and 0 - *not copied*. The normalized measure $nDCG$ is defined as $DCG_p / IDCG_p$, where $IDCG$ is the DCG measure for ideal ranking. $nDCG$ values scale between 0.0 – 1.0 with 1.0 for ideal ranking.

Results. We now establish two baselines for comparison. Our first approach uses a classical visual bag of words (BOW) approach computed at the interest points. The BOW representation has been successfully used in many image retrieval tasks including the document images [37, 49]. We use SIFT descriptors, quantized using LLC and represented using a spatial pyramid of size 1×3 . Our second baseline (NN) uses the classical word spotting scheme based on profile features similar to [32]. While the first one is scalable for large datasets, the second one is not really appropriate due to the time complexity of classical DTW. In both these methods, the best match is identified as the document which has most number of word/patch matches. Table 3 reports the quantitative evaluation for various matching schemes along with the baselines. The proposed MODS framework along with CNN features performs better in both evaluation measures consistently. Using SWM word matching scheme over the proposed CNN features, we achieve an $nDCG$ score of 0.8569 and AUC of 0.9465. This is further improved in the MODS, which incorporates loose ordering and is invariant to word overflow problems. Note that in both cases (SWM and MODS), the stopwords are removed as preprocessing. We also evaluate our framework with the state-of-the-art features proposed in [3] and observe a similar trend which validates the effectiveness of MODS. Fig. 7 shows some qualitative results of matching pairs from HW-DocSim dataset.

5.4 Human evaluations

To validate the performance of the system on an unrestricted collection, we introduce HW-1K dataset which is collected from the real assignments of a class as part of an active course. The dataset contains nearly 1K handwritten pages from more than 100 students. The content in these documents varied from text,

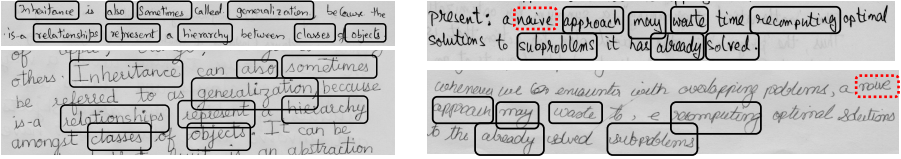


Fig. 7. Qualitative results of the MODS matching algorithm from HW-DocSim dataset. Here we show two sample matching pairs in two columns. The top region is taken from source and bottom one is plagiarized. The highlighted words in rectangle have been correctly matched along with few words which remain undetected.

figures, plots and mathematical symbols. Most of the documents follow a complex layout with misalignment in paragraphs, huge variations in line and word spacing and a high degree of skewness over the content.

We perform a human evaluation where we picked a set of 50 assignment images written by different students, and gathered the top-1 similar document image present in the corpus using MODS. We then ask five humans evaluators to give a score to the top-1 retrieval on a likert scale of 0 – 3 where 0 is “*very dissimilar*”, 1 is “*similar only for few word matches*”, 2 is “*partially similar*” and 3 is “*totally similar*”. Here, the scale-1 refers to the case where the document pair refers to the same topic. Thus there could be individual word matches but the text is not plagiarized. The average agreement to the human judgments as evaluated for the top-1 similar document is reported at 2.356 with 3 as the best score.

6 Discussions

We propose a method which estimates a measure of similarity for two handwritten documents. Given a set of digitized handwritten documents, we estimate a ranked list of similar pairs that can be used for manual validation, as in the case of MOSS and deciding the amount of plagiarism. Our document similarity score is computed using a CNN feature descriptor at the word level which surpasses the state-of-the-art results for the task of word spotting in multi-writer scenarios. We believe that with an annotated, larger set of natural handwritten word images, the performance can be further improved. We plan to use weakly supervised learning techniques for this purpose in the future.

Throughout this work, we characterize the document images with textual content alone. Many of the document images also have graphics. Our method fails to compare them reliably. On a qualitative analyses of the failures, we also find that the performance of matching mathematical expressions e.g., equations and symbols is inferior to the textual content. We believe identifying regions with graphics and applying separate scheme for matching such regions can further enhance the performance of our system.

Acknowledgment. Praveen Krishnan is supported by TCS Research Scholar Fellowship.

References

1. Almazán, J., Gordo, A., Fornés, A., Valveny, E.: Handwritten word spotting with corrected attributes. In: ICCV (2013)
2. Almazán, J., Gordo, A., Fornés, A., Valveny, E.: Segmentation-free word spotting with exemplar SVMs. PR (2014)
3. Almazán, J., Gordo, A., Fornés, A., Valveny, E.: Word spotting and recognition with embedded attributes. PAMI (2014)
4. Chum, O., Philbin, J., Zisserman, A.: Near duplicate image detection: min-hash and tf-idf weighting. In: BMVC (2008)
5. Clough, P.D., Stevenson, M.: Developing a corpus of plagiarised short answers. LREC (2011)
6. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR (2009)
7. Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.: Decaf: A deep convolutional activation feature for generic visual recognition. In: ICML (2014)
8. Fernández-Mota, D., Manmatha, R., Fornés, A., Lladós, J.: Sequential word spotting in historical handwritten documents. In: DAS (2014)
9. Fischer, A., Keller, A., Frinken, V., Bunke, H.: Lexicon-free handwritten word spotting using character HMMs. PRL (2012)
10. Gandhi, A., Jawahar, C.V.: Detection of cut-and-paste in document images. In: ICDAR (2013)
11. Gatos, B., Stamatoopoulos, N., Louloudis, G.: ICDAR2009 handwriting segmentation contest. IJDAR (2011)
12. Hinton, G.E., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.R.: Improving neural networks by preventing co-adaptation of feature detectors. CoRR (2012)
13. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. CoRR (2015)
14. Jaderberg, M., Simonyan, K., Vedaldi, A., Zisserman, A.: Reading text in the wild with convolutional neural networks. IJCV (2014)
15. Jaderberg, M., Simonyan, K., Vedaldi, A., Zisserman, A.: Synthetic data and artificial neural networks for natural scene text recognition. CoRR (2014)
16. Jaderberg, M., Vedaldi, A., Zisserman, A.: Deep features for text spotting. In: ECCV (2014)
17. Ke, Y., Sukthankar, R., Huston, L.: An efficient parts-based near-duplicate and sub-image retrieval system. In: ACM Multimedia (2004)
18. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NIPS (2012)
19. Lecun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. Proceedings of the IEEE (1998)
20. Louloudis, G., Gatos, B., Pratikakis, I., Halatsis, C.: Text line and word segmentation of handwritten documents. PR (2009)
21. Mahendran, A., Vedaldi, A.: Understanding deep image representations by inverting them. In: CVPR (2015)
22. Manmatha, R., Han, C., Riseman, E.M.: Word spotting: A new approach to indexing handwriting. In: CVPR (1996)
23. Manning, C.D., Raghavan, P., Schütze, H.: Introduction to Information Retrieval. Cambridge University Press (2008)

24. Marti, U., Bunke, H.: The IAM-database: an English sentence database for offline handwriting recognition. *IJDAR* (2002)
25. Milewski, R., Govindaraju, V.: Handwriting analysis of pre-hospital care reports. In: *CBMS* (2004)
26. Mishra, A., Alahari, K., Jawahar, C.V.: Top-down and bottom-up cues for scene text recognition. In: *CVPR* (2012)
27. Mishra, A., Alahari, K., Jawahar, C.V.: Image retrieval using textual cues. In: *ICCV* (2013)
28. NPTEL: <http://nptel.ac.in/> (2016), [Online; accessed 10-March-2016]
29. Perronnin, F., Rodríguez-Serrano, J.A.: Fisher kernels for handwritten word-spotting. In: *ICDAR* (2009)
30. Porter, M.F.: An algorithm for suffix stripping. *Program* (1980)
31. Potthast, M., Hagen, M., Beyer, A., Busse, M., Tippmann, M., Rosso, P., Stein, B.: Overview of the 6th international competition on plagiarism detection. In: *CLEF* (2014)
32. Rath, T.M., Manmatha, R.: Word spotting for historical documents. *IJDAR* (2007)
33. Razavian, A.S., Azizpour, H., Sullivan, J., Carlsson, S.: CNN features off-the-shelf: An astounding baseline for recognition. In: *CVPR* (2014)
34. Rodríguez-Serrano, J.A., Perronnin, F.: A model-based sequence similarity with application to handwritten word spotting. *PAMI* (2012)
35. Rusiñol, M., Aldavert, D., Toledo, R., Lladós, J.: Efficient segmentation-free keyword spotting in historical document collections. *PR* (2015)
36. Schleimer, S., Wilkerson, D.S., Aiken, A.: Winnowing: Local algorithms for document fingerprinting. In: *SIGMOD* (2003)
37. Shekhar, R., Jawahar, C.V.: Document specific sparse coding for word retrieval. In: *ICDAR* (2013)
38. Simard, P.Y., Steinkraus, D., Platt, J.C.: Best practices for convolutional neural networks applied to visual document analysis. In: *ICDAR* (2003)
39. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: *ICLR* (2015)
40. Srihari, S.N., Collins, J., Srihari, R.K., Srinivasan, H., Shetty, S., Brutt-Griffler, J.: Automatic scoring of short handwritten essays in reading comprehension tests. *Artif. Intell.* (2008)
41. Stamatopoulos, N., Gatos, B., Louloudis, G., Pal, U., Alaei, A.: ICDAR 2013 handwriting segmentation contest. In: *ICDAR* (2013)
42. Su, B., Lu, S.: Accurate scene text recognition based on recurrent neural network. In: *ACCV* (2014)
43. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: *CVPR* (2015)
44. Takeda, K., Kise, K., Iwamura, M.: Real-time document image retrieval for a 10 million pages database with a memory efficient and stability improved LLAH. In: *ICDAR* (2011)
45. Takeda, K., Kise, K., Iwamura, M.: Real-time document image retrieval on a smart-phone. In: *DAS* (2012)
46. Vitaladevuni, S.N.P., Choi, F., Prasad, R., Natarajan, P.: Detecting near-duplicate document images using interest point matching. In: *ICPR* (2012)
47. Wang, K., Babenko, B., Belongie, S.J.: End-to-end scene text recognition. In: *ICCV* (2011)
48. Wang, K., Belongie, S.J.: Word spotting in the wild. In: *ECCV* (2010)
49. Yalniz, I.Z., Manmatha, R.: An efficient framework for searching text in noisy document images. In: *DAS* (2012)