

# Word Spotting in Silent Lip Videos

Abhishek Jha<sup>1</sup>

Vinay P. Namboodiri<sup>2</sup>

C. V. Jawahar<sup>1</sup>

<sup>1</sup> Center for Visual Information Technology, KCIS, IIT Hyderabad, India

<sup>2</sup> Department of Computer Science and Engineering, IIT Kanpur, India

{abhishek.jha@research, jawahar@}.iiit.ac.in, vinaypn@iitk.ac.in

## Abstract

Our goal is to spot words in silent speech videos without explicitly recognizing the spoken words, where the lip motion of the speaker is clearly visible and audio is absent. Existing work in this domain has mainly focused on recognizing a fixed set of words in word-segmented lip videos, which limits the applicability of the learned model due to limited vocabulary and high dependency on the model’s recognition performance.

Our contribution is two-fold: 1) we develop a pipeline for recognition-free retrieval, and show its performance against recognition-based retrieval on a large-scale dataset and another set of out-of-vocabulary words. 2) We introduce a query expansion technique using pseudo-relevant feedback and propose a novel re-ranking method based on maximizing the correlation between spatio-temporal landmarks of the query and the top retrieval candidates. Our word spotting method achieves 35% higher mean average precision over recognition-based method on large-scale LRW dataset. Finally, we demonstrate the application of the method by word spotting in a popular speech video (“The great dictator” by Charlie Chaplin) where we show that the word retrieval can be used to understand what was spoken perhaps in the silent movies.

## 1. Introduction

Parsing information from videos has been explored in various ways in computer vision. Recent advances in deep learning have facilitated many such tasks. One such parsing requirement is of reading lips from videos. This has applications in surveillance or aiding improvements in speech recognition in noisy outdoor settings. Solving this problem has been attempted using methods based on recurrent neural networks (RNN) [28] and spatio-temporal deep convolutional networks [29]. However, for practical applications, recognizing lip motion into words is still in its nascent stages, with state of the art models [36] being limited to a constrained vocabulary. In this paper, we adopt a

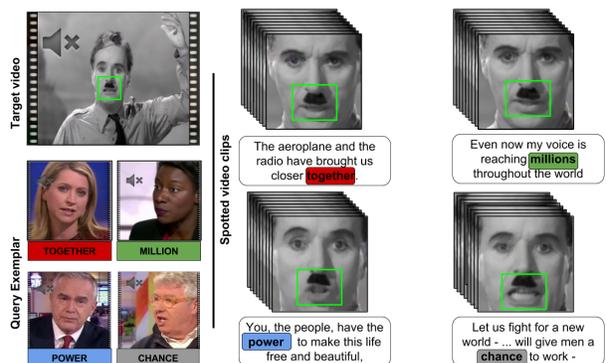


Figure 1. Example of word spotting in black and white Charlie Chaplin silent video: (left) target is the silent video and queries are the exemplars spoken by different people;(right) retrieved video clip segments where the words ‘together’, ‘million’, ‘power’ and ‘chance’ are present.

recognition-free ‘word-spotting’ approach that does not suffer from the vocabulary limitations. Unlike text documents, where the performance in character recognition [43], word recognition [14] and spotting research [37] has seen a great boost in the post deep learning era, this approach has been rarely pursued for lipreading task.

Training a lipreader requires careful word level annotation, which is expensive even for a small vocabulary set. Although progress in speech recognition [42] has resulted in better audio-to-text prediction and can be used for annotation, such methods are often prone to changes in accent and presence of noise in the audio channel. Lipreader’s performance is also susceptible to similar sounding words [36]. In recognition-based retrieval, we use a lipreader to predict the word spoken in a video clip. Evidently, if the word is wrongly predicted due to variations in visual appearance, it would never appear in the top results. In contrast, for recognition-free retrieval, the ‘word spotting’ i.e. matching of words is based on the feature representation of the target word without explicitly predicting the word itself. It intrinsically compares the features of the target word with the query word. Hence, even if the target word is misclassified it appears in the top results.

We are motivated by the fact that for handwritten documents word spotting has shown better performance for retrieving target words in different handwriting styles than word recognition [34]. Likewise, we show that recognition-free retrieval can also be useful for spotting words when target words come from a different source than the data used for training a lipreader, like archaic black and white documents films in Figure 1. We further investigate the applicability of recognition-free pipeline for out-of-vocabulary word spotting, for a different domain of data with respect to what has been used for training the lipreading model. Figure 1 shows few sample results of our pipeline for spotting different query words in the black and white video clip in four spoken sentences.

We further show that the word spotting performance can be improved by using a novel re-ranking method for top-k retrieval candidates. We also adapt the standard pseudo-relevance feedback query expansion method for lipreading task. Our pipeline takes silent speech videos as input and retrieves a queried word that is provided again as a video clip from the target input dataset. The target video is first densely segmented into ‘word proposal clips’, where these clips may or may not contain any word. Any ‘word proposal clip’ is considered a spotted word if the similarity measure between the query and the target ‘word proposal clip’ is greater than a particular threshold.

We show improvement in word spotting on a standard large scale lip video dataset Lipreading in the wild (LRW) [11], and another standard dataset GRID corpus [13] for showing domain invariance. We also assess our pipeline’s performance in a popular speech video by Charlie Chaplin: “*The great dictator*”.

## 2. Related Work

Research in visual speech recognition has been pursued for at least two decades [4, 7, 26] with earlier approaches focusing mainly on handcrafted features and HMM-based speech recognizers [5, 27, 33]. Some of these approaches have been thoroughly reviewed in [24, 44]. Wand *et al.* [39] showed word level lipreading using an LSTM [28] stacked over two-layered neural network on GRID corpus dataset [13]. Recently, Chung and Zisserman [11] have used multiple lipreading models that fuses the temporal sequence at different layers of underlying VGG-M model [8] to classify the input video clip into 500 words. Assael *et al.* [2] uses a Connectionist Temporal Classification (CTC) [22] to show one of the best results on GRID corpus [13].

Lipreading involves modeling temporal sequences of lip video clips into phonemes or characters, hence better sequence learning models using deep networks proved to be pivotal in lipreading research. Chung *et al.* [10] have proposed Watch Listen Attend and Spell (WLAS) architecture that leverages attention model [3] for doing character level

prediction of input lip videos. They provide the best results on Lipreading in the Wild (LRW) dataset and GRID corpus [13]. They however use a much larger Lipreading Sentences (LRS) dataset that is not widely available [10] for pretraining, hence making it a data intensive model that is not accessible. In a recent work, Stafylakis and Tzimiropoulos [36] trained a model entirely on LRW dataset to give state-of-the-art result for word level prediction. This model consisted of three parts: a spatio-temporal convolutional front-end, followed by a Resnet-34 [25], and a bidirectional LSTM [23] at the end. Since this model has been trained to classify lip videos into one of 500 word classes, it does not address out-of-vocabulary words. Our pipeline employs recognition architectures based on [11] and [36] as feature extractors to show how recognition-free leverages these features spaces for improved retrieval performance.

Initial work in word spotting appeared in speech recognition community, majority relying on HMMs [21, 35]. Kernel machines and large margin classifiers introduced by Keshet *et al.* [30] in discriminative supervised setting resulted in an improvement over the previous methods. Post deep learning, RNNs with CTC objective functions gave a major improvement over the HMMs [16] for modeling temporal audio speech signals. Unlike audio speech, visual speech is spatio-temporal signal. Hence, our choice of feature extractors contain VGG-M [8] and Resnet-34 [25] modules for modeling facial features, and uses LSTM and temporal convolution for modeling temporal information.

Word spotting is a well defined problem in document analysis and retrieval [20]: hand writing recognition [17, 19, 34, 37], word image retrieval [32], scene-text [40] etc. Although a large corpus of work exists for word spotting for documents, images and audio speech, the visual speech domain has been largely ignored. The work that is closest to our approach is by Wu *et al.* [41]. In their approach, the authors use geometric and appearance based features to build their word spotting pipeline and they rely on the knowledge of optimal handcrafted feature. In our work, though we also adopt a recognition-free retrieval approach, we do so using recognition-based features and show that the recognition-free approach improves on the recognition-based approach. We further also improve the base recognition-free pipeline by using query expansion and re-ranking extensions. We benchmark our work on standard datasets.

## 3. Proposed Method

In this section, we will discuss the individual components of our proposed word spotting pipeline and move along to develop a holistic overview of the method.

### 3.1. Recognition-free Retrieval

Recognition-based retrieval relies on recognizing words in lip videos by completely depending on the lipreading

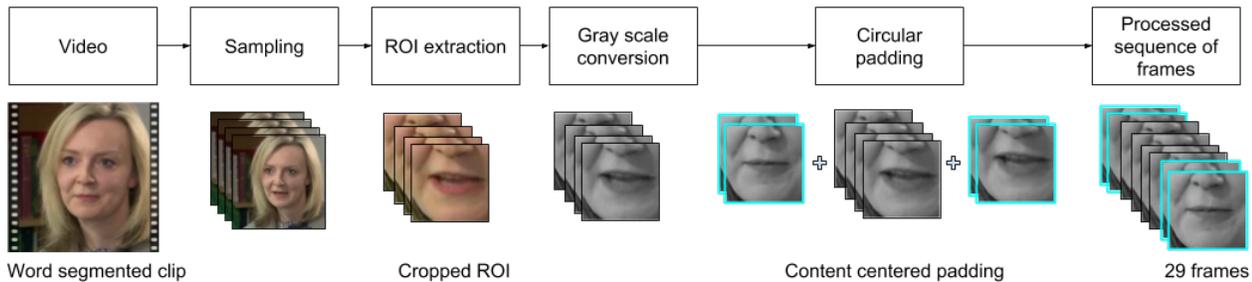


Figure 2. Preprocessing: The pipeline which takes a variable length word clip and converts it into a fixed length sequence of frames.

model. During testing a video clip containing a word is classified as one of the word in the vocabulary it is trained on. Moreover, modeling a lipreader with open vocabulary is an active area of research.

Retrieving a word from a set of candidate silent videos without directly recognizing each candidate words being spoken is recognition-free retrieval or word spotting. This opens up an opportunity to use a sub-performing lipreader with incorrect word recognition. In a recognition-free setup, the user formulates a query and a rank list is computed based on its distance from all the clips in the target corpus (retrieval set), such that most similar candidate is given the highest rank. Since word spotting systems rely heavily on computing similarity, the quality of the feature representation is more important than the classification of input clips.

Word spotting based on the modality of query are of two types: query by string (QbS) where the input query is string and the retrieval is video, and query by exemplar (QbE), where query is video and retrieval is also video. In this work, our query will be through exemplar.

### 3.2. Preprocessing

We use the recognition models as described in [10] and [36] as feature extractors. These models takes inputs as a fixed length input of spatial dimension  $225 \times 225$  and  $112 \times 112$  respectively with a sequence length of 29 frames. The feature extractors are trained on LRW [11] dataset which consists of fixed length video clips of 29 frames and 1.16 sec duration, with actual word at the center. Hence it is required to preprocess the input videos (other than that of LRW) before feeding them to the feature extractors. As shown in Figure 2, the preprocessing step proceeds by just sampling the input video at 25 frames per second, then converting the sampled frames to grayscale. Since words can be of different length we circular pad grayscale sequence of frames on both the side such that the actual content is at the center of the sequence. Circular padding of length 2 for a sequence:  $\{1, 2, 3, 4, 5\}$  on both sides gives  $\{4, 5, 1, 2, 3, 4, 5, 1, 2\}$ .

### 3.3. Video Features

Our first feature extractor only uses the visual stream of the WLAS architecture and hence called Watch, Attend and Spell (WAS) model [10]. Chung *et al.* [10] train WLAS model on LRS dataset [10] and fine tune it on LRW dataset [11]. As LRS dataset [10] is not yet publicly available, we trained our WAS model entirely on LRW dataset. WAS contains two modules: a VGG-M convolution module and an attention-based sequence to sequence LSTM module, followed by 28 neurons with softmax non-linearity. Our output sequence for a lip video clip is maximum 20 character long, 28 dimensional(D) (A to Z, *eos*, *padding*) ground truth (GT) word label. Using early stopping we achieve a word accuracy of 53%.

We also employ another network ‘N3’ as described by Stafylakis and Tzimiropoulos [36] for feature extraction. This network is composed of three modules: A layer of 3D convolutions followed by three dense layers (fully connected layers), and finally a temporal convolution layer. The final layer has 500 neurons with softmax non-linearity. The classification accuracy of this model is 69.7%. We will address this model as CMT in this paper.

In both the feature extractors, the choice of features are the softmax scores or the probabilities of a lip videos belonging to different words in the vocabulary, instead of sparsely belonging to only one word. We also experimented with the output of the last dense layer as feature representation for the input video, and found softmax scores to be empirically better.

### 3.4. Overall Pipeline

In this section, we propose a pipeline for spotting words in silent lip videos. In order to demonstrate generic nature of our pipeline, we first train our two different feature extractors on LRW dataset. We project the query set, consisting of preprocessed annotated video clips, and retrieval set video clips which do not have any labels into the feature space. The label of the query is assigned to a particular candidate clip in the retrieval set, only if the mean similarity score of that candidate with all the same label queries

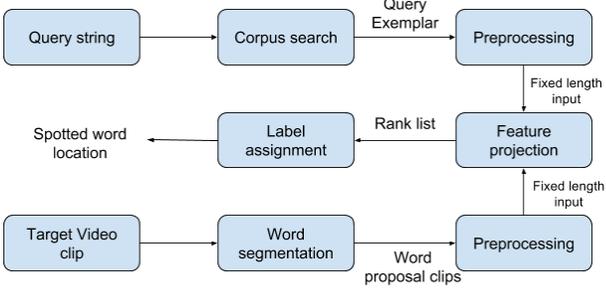


Figure 3. Overall pipeline: First a string is searched in an annotated corpus to formulate an exemplar which is then preprocessed, and projected into feature space. Target video is then segmented into word clips, either using given time-stamp or dense segmentation, preprocessed and projected in the same feature space. A ranking is computed based on the cosine similarity between query exemplar and the word proposal clips. Label is transferred based on majority voting, as discussed later in Subsection 3.4

is greater than a threshold, otherwise it is assumed the candidate word proposal clip does not contain a full word. In Figure 3 we show our overall pipeline.

More precisely, if  $q_i^c$  is the feature representation of  $i^{th}$  query belonging to label  $c$  and  $r_j$  is feature representation of the  $j^{th}$  word proposal clip, the similarity score between the two is given by  $n_{ij}^c$  in Equation 1.

$$n_{ij}^c = \frac{(q_i^c)^T \cdot r_j}{\|q_i^c\| \cdot \|r_j\|} \quad (1)$$

The average similarity between all the queries  $q^c$  belonging to label  $c$  and the candidate  $r_j$  is given by  $s_j^c$  in the below Equation 2.

$$s_j^c = \frac{\sum |q^c| n_{ij}^c}{|q^c|} \quad (2)$$

Finally, the label assignment for the candidate  $r_j$  is  $c$  if the mean similarity score between all the queries belonging to label  $c$ , i.e.  $s_j^c$ , is greater than  $\rho$ . Otherwise, we consider the word proposal clip is either noise or does not contain the whole word, as represented by  $\phi$ .

$$label_{r_j} = \begin{cases} c & \text{if } s_j^c > \rho \\ \phi & \text{otherwise} \end{cases} \quad (3)$$

Hence, these word proposal clips are spotted as word  $c$  using the queries  $q_i^c$  in the target video. We can further use enhancements over this pipeline to improve the retrieval performance, which we will discuss in the next section.

## 4. Enhancements

In this section, we discuss a query expansion technique to search videos with semantic relevance to the given query, followed by re-ranking method to improve ordering of top-k results.

### 4.1. Query Expansion and Re-ranking

Query expansion, in image retrieval [1], has been widely used to improve retrieval performance by increasing the recall and obtain additional documents which might get missed with the original query. Similar to documents, we first feed a *seed* query to our retrieval system which gives us a ranked list of all the candidates from the retrieval set. From this set, top-k candidates are selected to construct a new query based on the weighted sum of the query and top-k candidates feature vectors as the pseudo-relevance feedback to improve the retrieval results.

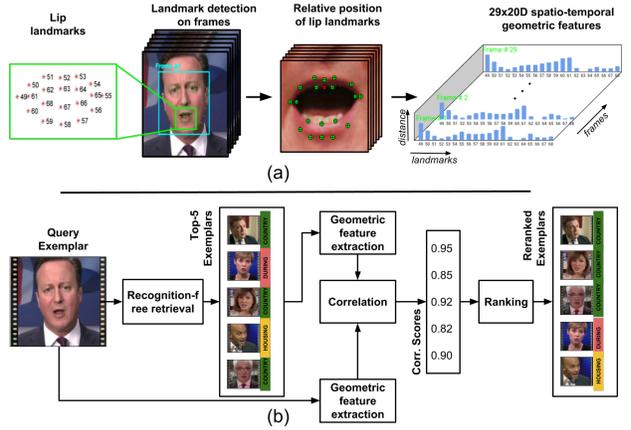


Figure 4. Re-ranking using geometric cues of lip video: (a) shows method of extracting spatio-temporal feature using lip landmarks of each frame of the video clip; (b) shows re-ranking of top-5 retrieved candidates based on the correlation between spatio-temporal features of top-5 candidates and that of the query.

Re-ranking is used to improve the ranking of top retrieval results for a given query. Some of the prominent re-ranking method [15, 38] relies on geometrical consistency between query and its top retrieval candidates. Fergus *et al.* [15] uses RANSAC [18] to re-rank top results from Google Image search engine. Unlike images, lip videos are temporal in nature with each word consisting of a specific set of phonemes. To adapt such a method for lip videos, we extract spatio-temporal features. Out of total 68 facial landmarks [31], we first compute the distance between all the 20 landmarks associated with lip and the lip-central landmark (landmark no. 63), as shown by ‘red’ color landmark in Figure 4(a). Both landmark no. 63 and 67, being in the center, are clearly visible for different head poses and hence can be chosen for computing distances. However, on an average, the motion of the upper lip is lesser than the lower lip for most of the word utterances, makes landmark 63 more stable and a better choice.

This geometric feature extraction results in a 20D spatial feature for each frame, or  $20 \times 29D$  spatio-temporal feature for the video clip. We then re-rank our candidate using their temporal lip landmark correlation with the query lip video,

as shown in Figure 4(b). Using recognition-free retrieval top-k candidates are selected for a given query. Spatio-temporal features for both top-k candidates and query are extracted. The correlation of landmark of the lip region of these top-k candidates with the query is computed, the re-ranking is done in the order of decreasing correlation.

## 5. Experiments

### 5.1. Datasets

**Lipreading in Wild (LRW)** [11] has 500 words classes with 1000 clips for training, 50 for testing, and 50 for validation for each of the words, which has been curated from BBC news videos. Each word clip is of length 1.16 second duration containing 29 frames. We use the LRW to train both feature extractors. The proposed retrieval pipeline only uses the test set for querying and validation set for retrieval, since training set has been used to train feature extractors.



Figure 5. Random frames from LRW dataset (top row), GRID corpus (middle row) and Charlie Chaplin “The great dictator” speech video (bottom row).

**GRID corpus** [13] contains 1000 phrases, spoken by each of 33 speakers. Each phrase has a fixed syntax containing 6 words: *command*(4) + *color*(4) + *preposition*(4) + *letter*(25) + *digit*(10) + *adverb*(4); an example of which is ‘put red at G 9 now’. We use speakers 10-19, similar to [39], in our experiment. For showing domain invariance, we randomly sample 1000 phrases from these speakers to create our query set. Similarly, we sample another 1000 phrases from the same speakers to create our retrieval set. All the speech videos are word segmented and preprocessed before feeding to feature extractors.

For qualitative results we show lipreading on **Charlie Chaplin’s** famous “The great dictator” speech video. We only use the video, without audio cues for our experiment. The video is segmented into sentence level video clips using the timestamps provided by Youtube subtitles, which also gives the ground truth annotations. The retrieval corpus is made by densely segmenting these sentence videos into word proposal clips. Randomly selected frames from these three datasets are shown in Figure 5.

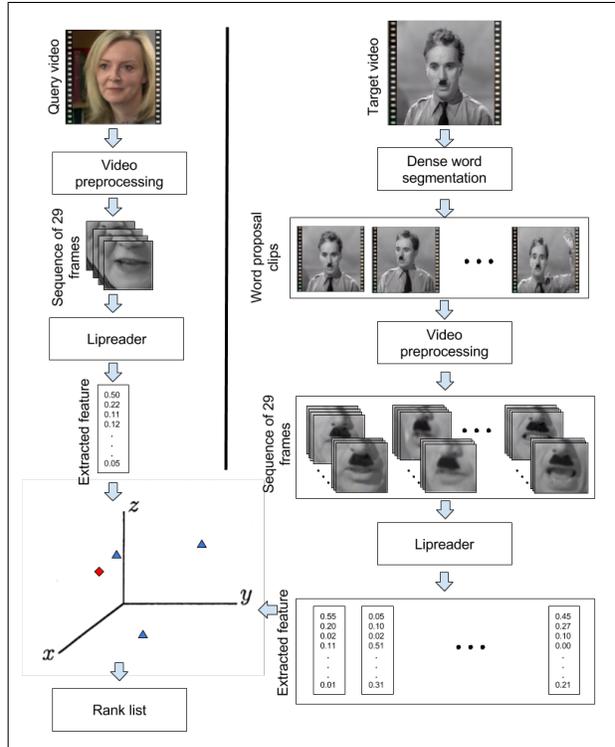


Figure 6. Word spotting in Charlie Chaplin video: (left) a query exemplar with known annotation is preprocessed into fixed length input and fed to the feature extractor. (right) the Charlie Chaplin video is first densely segmented into word proposal clips and fed to the feature extractor. All the word proposal clips and query exemplar is projected into feature space and ranking is computed based on cosine similarity.

### 5.2. Implementation

For WAS, we use the pretrained VGG-M model from Chung and Zisserman [12], and only train attention sequence-to-sequence LSTM module, while freezing the weights of VGG-M module. We use the LRW training set for training our model, with validation set used for parameter tuning. The network has been trained with batch size 64, cross-entropy loss and SGD optimizer. Initial learning rate was set to 0.1 with a decay of 0.01% every two iterations. No data augmentation was used.

For training CMT, we follow the similar procedure as mentioned in Stafylakis and Tzimiropoulos [36] to train our model end-to-end. Again the batch size of 64 was taken with cross-entropy loss and SGD optimizer was used. Initial learning rate was set to  $3e^{-3}$  with exponential decay in learning rate when the validation loss does not decrease for 2 epochs. We also perform data augmentation with random cropping of 4 pixels around the lip region of interest (ROI), and horizontally flipping all frames of randomly chosen input clips. For both the networks, WAS and CMT, early stopping was employed if validation accuracy failed to improve over 3 consecutive epochs. We implement both the

networks in Keras deep learning library [9].

Word spotting on LRW dataset has been shown considering LRW test set as query set and LRW validation set as retrieval set. Here, we want to assign label to the query video clips, considering we know the GT label for retrieval set. Both the query and retrieval set are first preprocessed, as discussed in Section 3.2. Since all the video clips are 29 frames long, circular padding is not required during preprocessing. After feature extraction, the query is searched in the retrieval set, the candidate with highest cosine similarity is ranked highest. To transfer word label from retrieval set the query, we take the majority vote of top-5 candidates in the retrieval set.

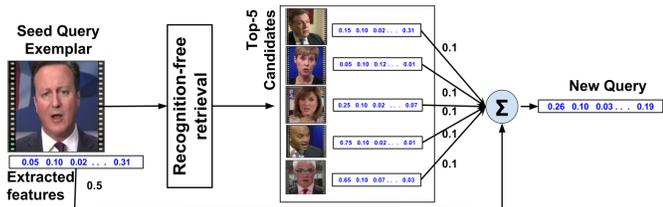


Figure 7. Formulation of new query: The weighted sum of the feature representation of *seed* query and its top-5 retrieved candidates becomes the new query.

During query expansion, we first search a *seed* query in the retrieval set to get top-5 candidates. The ‘New query’ is the weighted sum of the top-5 candidates with weights 0.1 each and *seed* query with weight 0.5, as shown in Figure 7. This query is then used to retrieve a new set of candidates which becomes our final retrieval for the *seed* query.

For each query video coming from LRW test set, we retrieve top-10 candidates from LRW validation set using recognition-free retrieval. For Re-ranking, we then extract spatio-temporal feature for both query video and its top-10 retrieval candidates using DLib [31] and OpenCV [6] libraries. Correlation between spatio-temporal features of query and candidates were computed and were used to re-rank the top-10 candidates. This method proves to be effective in refining the search results for our retrieval pipeline.

For showing word spotting in Charlie Chaplin video, as shown in Figure 6, the sentence videos are densely segmented into fixed length (29 frames) word proposal clips by taking stride of 3 frames. We spot the words in retrieval corpus consisting of these clips. Since the segmentation is dense there will be very few word proposal clips which will entirely cover actual words spoken in the video. As discussed in Section 3.4, we calculate the average similarity score between all the query exemplars coming from LRW validation set belonging to a particular word label and a word proposal clip from Charlie Chaplin video. If the average similarity is more than a threshold ( $\rho$ ), we assign the word label to the word proposal clip. We empirically selected the value of  $\rho = 0.3$  for this experiment.

### 5.3. Baselines

We compare our pipeline with recognition-based retrieval. WAS [10], in the original paper, was first pretrained on LRS dataset, and later fine-tuned on LRW dataset, gives a word accuracy of 76.2%. Our WAS model trained solely on LRW dataset gives the word accuracy of 53%. The recognition-based baseline of our WAS is given in Table 1, column 1. Another lipreader CMT, gives the word accuracy of 69.7%. The recognition-based baseline is given in Table 1 column 3.

For GRID corpus we do not fine-tune our LRW trained base feature extractors on GRID corpus. The recognition-based baseline for the domain-invariance out-of-vocabulary retrieval is shown in Table 3, column 1 and 3.

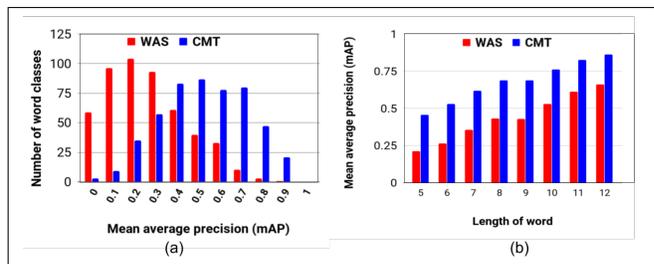


Figure 8. (a) number of words below a certain mAP for WAS and CMT based pipeline: y-axis is the number of words, and x-axis is the mAP; (b) variation of mean average precision (mAP) with the length of the word for CMT and WAS based pipeline: y-axis is average mAP and x-axis is word length in LRW vocabulary

### 5.4. Evaluation metric

For search based applications, the most important performance factor is: how many good results are in the top search results. Hence, Precision@K, which measures the precision at fixed lower levels of retrieval results, makes sense as an important performance metric. It considers the number of desirable results out of the top-k retrieval results without taking into account the overall rank ordering of the search results.

Recall@K is another important evaluation metric that we show, which is the number of desired results retrieved among top-k search results, with respect to the total number of available positive results.

While Precision@K and Recall@K give specific insights into the performance of the retrieval system, both measure performance for a fixed number of retrievals (K) and are insensitive to the overall rank ordering of the search results. We therefore also report the Mean Average Precision (mAP) for our retrieval system. mAP provides a measure of the quality of retrieval across different recall levels. mAP has been shown to have especially good discrimination and stability, and is one of the most standard evaluation measures for word spotting.

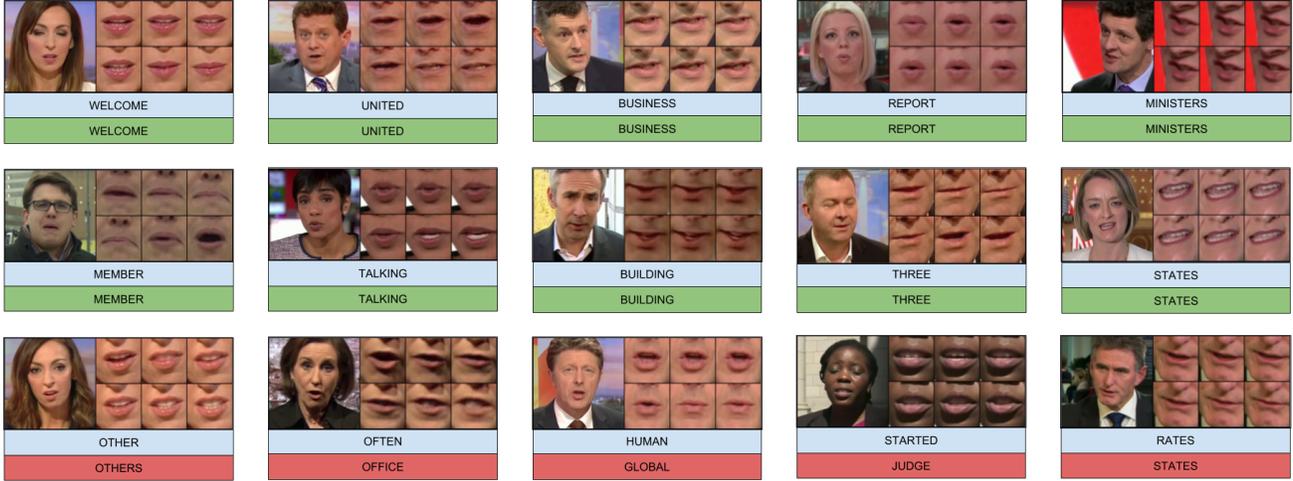


Figure 9. Qualitative results on LRW dataset: Each image depicts the central frame of the query video clip (left) and a sequence of lip ROIs of 6 consecutive frames around central frame, shown in raster order (right); (middle) blue boxes are the ground truths; (bottom) green boxes are correct predictions while the red ones are incorrect predictions. Label is propagated to a query based on the majority label present in the top-5 retrieval candidates.

## 6. Results

### 6.1. Comparison with Baseline Methods

Recognition-free retrieval or word spotting on LRW dataset when the base lipreader is WAS gives an absolute improvement of 35.9% over the recognition-based baseline of mAP 0.23; Table 1, column 2. Similarly, for recognition-free retrieval using CMT lipreader there is an improvement of 49.67% over the recognition-based baseline of mAP 0.38; Table 1, column 4. For recognition-free retrieval using WAS (in red) and CMT (in blue) feature extractor, Figure 8(a) shows the number of words below a certain mAP value. The variation of average mAP with the length of the words in the LRW vocabulary is shown in Figure 8(b). It can be seen that the average mAP value increases with the increase in word length. The qualitative results for word-spotting on LRW dataset using CMT features can be seen in Figure 9.

	WAS		CMT	
	RB (BL)	RF (ours)	RB (BL)	RF (ours)
mAP	0.2317	<b>0.3149</b>	0.3807	<b>0.5698</b>
P@10	0.2928	0.4566	0.3253	0.6519
R@10	0.0586	0.0913	0.0651	0.1304
% imp.in mAP	–	<b>35.90</b>	–	<b>49.67</b>

Table 1. Retrieval performance for LRW dataset: Left two columns show recognition-based (RB) baseline and recognition-free (RF) performances for WAS features; right two columns show the similar results for CMT features. Across columns (first row) mAP is mean average precision, (second row) P@10 is precision at 10, (third row) R@10 is recall at 10, and (last row) % imp.in mAP is percentage mAP improvement of recognition-free retrieval over baseline.

Query expansion on LRW dataset using two lipreaders: WAS and CMT give a mAP of 0.3146 and 0.5722 respectively; Table 2, column 2 and 5. Although the mAP results are comparative to the recognition-free method, we see an overall increase in recall@10. Also, re-ranking using spatio-temporal cues improves the retrieval performance for WAS and CMT, giving a mAP of 0.3179 and 0.5709 respectively; Table 2, column 3 and 6.

	WAS			CMT		
	RF	QExp	ReR	RF	QExp	ReR
mAP	0.3149	0.3146	<b>0.3179</b>	0.5698	0.5722	<b>0.5709</b>
P@10	0.4566	0.4591	0.4566	0.6519	0.6572	0.6519
R@10	0.0913	<b>0.0918</b>	0.0913	0.1304	<b>0.1314</b>	0.1304

Table 2. Different recognition-free performance for LRW dataset: Left three columns are recognition-free (RF), query expansion (QExp) and re-ranking (ReR) performances for WAS features; right three columns show similar results for CMT features. Across columns (first row) mAP is mean average precision, (second row) P@10 is precision at 10, and (last row) R@10 is recall at 10.

Charlie Chaplin “*The great dictator*” speech video, contains 39 words from LRW vocabulary. It has a total of 54 spoken sentences, out of which 33 sentences actually contains LRW vocabulary words. Hence, the query set contains 50 exemplars, from LRW validation set, belonging to each of these 39 common vocabulary words. Using our CMT based recognition-free pipeline we were able to correctly spot instances of 13 instances of the common vocabulary words in 11 sentences. Whereas on using recognition-based pipeline, only 6 instances of common vocabulary words in 6 sentences are correctly predicted. The qualitative results can be seen in Figure 10, where we spot the sentences which



Figure 10. Qualitative results on Charlie Chaplin “*The great dictator video*”: Each image is one of the frames in the sentence clips extracted from the speech video. The top text box in blue color contains the subtitles with **bold** text showing the common LRW vocabulary word present in the subtitle. The bottom text box shows the correctly spotted word.

contain the query words.

## 6.2. Domain Invariance

Domain Invariance provides us the robustness of the pipeline for target data distribution different from the one it is trained on. GRID corpus contains 51 words with only 1 common word available in LRW dataset vocabulary. Hence this experiment also shows out-of-vocabulary retrieval performance of the proposed pipeline.

On GRID corpus, the recognition-based baseline is 0.033 (mAP) for WAS features and 0.06 (mAP) for CMT features, while the recognition-free performance is 0.068 (mAP) for WAS and 0.177 (mAP) for CMT; Table 3, column 2. This signifies the utility of recognition-free retrieval for out-of-vocabulary words when the underlying lipreader is constrained by vocabulary size.

	WAS		CMT	
	RB (BL)	RF(ours)	RB (BL)	RF(ours)
mAP	0.033	<b>0.068</b>	0.060	<b>0.177</b>
P@10	0.034	0.219	0.224	0.322
R@10	0.002	0.016	0.019	0.020
% imp.in mAP	–	<b>106</b>	–	<b>195</b>

Table 3. Domain invariance results on Grid corpus dataset (for both WAS and CMT): Left column has recognition-based (RB) baseline performance and right has our recognition-free (RF) performance where (first row) mAP is mean average precision, (second row) P@10 is precision at 10, (third row) R@10 is recall at 10, and (last row) % imp.in mAP is the percentage mAP improvement of our proposed method over baseline.

## 6.3. Discussions

Many conclusions can be drawn from the result presented in Subsection 6.1. Recognition-free retrieval per-

formed better than recognition-based counterpart for spotting words in LRW dataset. From Figure 8(b), we see that quality of retrieval improves when the length of word increases, as longer the word is more the number of phonemes it contains, and less is the chance of it being similar to other words. Errors in similar sounding words are more likely, as can also be seen in Figure 9.

Performance of recognition-based retrieval on GRID corpus is inferior to that on LRW dataset, the reason being neither of the two feature extractors in our experiments were fine-tuned on GRID corpus. Still, the recognition-free retrieval showed an improvement over recognition-based. Quality of lip video is also important, as some words in Charlie Chaplin videos were not spotted, due to lower contrast and quality of the lip ROI, as shown in Figure 10.

## 7. Conclusion

We proposed a recognition-free retrieval pipeline and showed its precedence over recognition-based retrieval for the task of word-spotting. The base features from WAS and CMT lipreading models have been used to spot words in LRW dataset with an improvement of about 36% and 50% over the recognition-based counterpart. Pseudo-relevance feedback and re-ranking techniques, using spatio-temporal geometrical cues available in the lip videos, has been incorporated in the pipeline to further improve the retrieval results. We also showed domain invariance of our pipeline through out-of-vocabulary word spotting on GRID corpus dataset with an improvement of 106% and 195% over the baseline using WAS and CMT features respectively. Lastly, we presented the practical applicability of our proposed pipeline by spotting words in 11 out of 33 sentences in the “Charlie Chaplin, *The great dictator*” speech video.

## References

- [1] R. Arandjelović and A. Zisserman. Three things everyone should know to improve object retrieval. In *CVPR*, 2012.
- [2] Y. M. Assael, B. Shillingford, S. Whiteson, and N. de Freitas. Lipnet: Sentence-level lipreading. *arXiv preprint arXiv:1611.01599*, 2016.
- [3] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [4] S. Basu, N. Oliver, and A. Pentland. 3d modeling and tracking of human lip motions. In *ICCV*, 1998.
- [5] H. A. Bourlard and N. Morgan. *Connectionist speech recognition: a hybrid approach*, volume 247. Springer Science & Business Media, 2012.
- [6] G. Bradski. The opencv library. *Dr. Dobb's Journal: Software Tools for the Professional Programmer*, 25(11), 2000.
- [7] S. S. Brooke N.M. Pca image coding schemes and visual speech intelligibility. In *Proceedings of the Institute of Acoustics*, volume 16, 1994.
- [8] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *BMVC*, 2014.
- [9] F. Chollet et al. Keras, 2015.
- [10] J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman. Lip reading sentences in the wild. In *CVPR*, 2016.
- [11] J. S. Chung and A. Zisserman. Lip reading in the wild. In *ACCV*, 2016.
- [12] J. S. Chung and A. Zisserman. Out of time: automated lip sync in the wild. In *ACCV*, 2016.
- [13] M. Cooke, J. Barker, S. Cunningham, and X. Shao. An audio-visual corpus for speech perception and automatic speech recognition. *The Journal of the Acoustical Society of America*, 120(5), 2006.
- [14] P. Doetsch, M. Kozielski, and H. Ney. Fast and robust training of recurrent neural networks for offline handwriting recognition. In *ICFHR*, 2014.
- [15] R. Fergus, P. Perona, and A. Zisserman. A visual category filter for google images. *ECCV*, 2004.
- [16] S. Fernández, A. Graves, and J. Schmidhuber. An application of recurrent neural networks to discriminative keyword spotting. *ICANN*, 2007.
- [17] A. Fischer, A. Keller, V. Frinken, and H. Bunke. HMM-based word spotting in handwritten documents using subword models. In *ICMR*, 2010.
- [18] M. A. Fischler and R. C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6), 1981.
- [19] V. Frinken, A. Fischer, R. Manmatha, and H. Bunke. A novel word spotting method based on recurrent neural networks. *IEEE TPAMI*, 34(2), 2012.
- [20] A. P. Giotis, G. Sfikas, B. Gatos, and C. Nikou. A survey of document image word spotting techniques. *Pattern Recognition*, 68, 2017.
- [21] H. Gish and K. Ng. A segmental speech model with applications to word spotting. In *ICASSP*, volume 2, 1993.
- [22] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *ICML*, 2006.
- [23] A. Graves, S. Fernández, and J. Schmidhuber. Bidirectional LSTM networks for improved phoneme classification and recognition. *ICANN*, 2005.
- [24] A. B. Hassanat. Visual words for automatic lip-reading. *arXiv preprint arXiv:1409.6689*, 2014.
- [25] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [26] M. E. Hennecke. Audio-visual speech recognition: Preprocessing, learning and sensory integration. *PhD thesis, Stanford Univ.*, 1997.
- [27] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6), 2012.
- [28] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8), 1997.
- [29] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE TPAMI*, 35(1), 2013.
- [30] J. Keshet, D. Grangier, and S. Bengio. Discriminative keyword spotting. *Speech Communication*, 51(4), 2009.
- [31] D. E. King. Dlib-ml: A machine learning toolkit. *JMLR*, 2009.
- [32] P. Krishnan and C. V. Jawahar. Bringing semantics in word image retrieval. In *ICDAR*, 2013.
- [33] J.-S. Lee and C. H. Park. Robust audio-visual speech recognition based on late integration. *IEEE TMM*, 10(5), 2008.
- [34] R. Manmatha, C. Han, and E. M. Riseman. Word spotting: A new approach to indexing handwriting. In *CVPR*, 1996.
- [35] J. R. Rohlicek, W. Russell, S. Roukos, and H. Gish. Continuous hidden markov modeling for speaker-independent word spotting. In *ICASSP*, 1989.
- [36] T. Stafylakis and G. Tzimiropoulos. Combining residual networks with LSTMs for lipreading. *arXiv preprint arXiv:1703.04105*, 2017.
- [37] S. Sudholt and G. A. Fink. Phocnet: A deep convolutional neural network for word spotting in handwritten documents. In *ICFHR*, 2016.
- [38] S. S. Tsai, D. Chen, G. Takacs, V. Chandrasekhar, R. Vedantham, R. Grzeszczuk, and B. Girod. Fast geometric re-ranking for image-based retrieval. In *ICIP*, 2010.
- [39] M. Wand, J. Koutník, and J. Schmidhuber. Lipreading with long short-term memory. In *ICASSP*, 2016.
- [40] K. Wang and S. Belongie. Word spotting in the wild. In *ECCV*, 2010.
- [41] P. Wu, H. Liu, X. Li, T. Fan, and X. Zhang. A novel lip descriptor for audio-visual keyword spotting based on adaptive decision fusion. *IEEE TMM*, 18(3), 2016.
- [42] W. Xiong, J. Droppo, X. Huang, F. Seide, M. Seltzer, A. Stolcke, D. Yu, and G. Zweig. Achieving human parity in conversational speech recognition. *arXiv preprint arXiv:1610.05256*, 2016.

- [43] X.-Y. Zhang, F. Yin, Y.-M. Zhang, C.-L. Liu, and Y. Bengio. Drawing and recognizing chinese characters with recurrent neural network. *IEEE TPAMI*, 2017.
- [44] Z. Zhou, G. Zhao, X. Hong, and M. Pietikäinen. A review of recent advances in visual speech decoding. *Image and vision computing*, 32(9), 2014.