

## Article

# A Two-stage Approach to Segmentation-free Query-by-example Word Spotting

Leonard Rothacker, Marçal Rusiñol, Josep Lladós, and Gernot A. Fink | Dortmund – Barcelona

## Abstract

With the ongoing progress in digitization, huge document collections and archives have become available to a broad audience. Scanned document images can be transmitted electronically and studied simultaneously throughout the world. While this is very beneficial, it is often impossible to perform automated searches on these document collections. Optical character recognition usually fails when it comes to handwritten or historic documents. In order to address the need for exploring document collections rapidly, researchers are working on word spotting. In query-by-example word spotting scenarios, the user selects an exemplary occurrence of the query word in a document image. The word spotting system then retrieves all regions in the collection that are visually similar to the given example of the query word. The best matching regions are presented to the user and no actual transcription is required.

An important property of a word spotting system is the computational speed with which queries can be executed. In our previous work, we presented a relatively slow but high-precision method. In the present work, we will extend this baseline system to an integrated two-stage approach. In a coarse-grained first stage, we will filter document images efficiently in order to identify regions that are likely to contain the query word. In the fine-grained second stage, these regions will be analyzed with our previously presented high-precision method. Finally, we will report recognition results and query times for the well-known George Washington benchmark in our evaluation. We achieve state-of-the-art recognition results while the query times can be reduced to 50% in comparison with our baseline.

based representation like ASCII or UTF-8.<sup>1</sup> Even though researchers have been investigating this topic for decades, only the recognition of clear machine-printed texts can be considered solved. For handwritten documents and historical manuscripts, results are far from satisfactory. Difficulties with handwritten documents are caused by the extreme variability in human writing. Historic documents often show severe degradation, such as ink bleed-through, bad contrast due to changes in paper color, text line deviations, old fonts and other artifacts caused by old technical standards and storage. Nevertheless, algorithmic methods for human assistance exist in order to explore these kinds of documents. They usually work in a relatively constrained scenario but will save a significant amount of work if they are applicable.

One of the most prominent techniques for this purpose is word spotting.<sup>2</sup> The task is not to transcribe entire images of text but to automatically detect regions in document images where the query word is likely to be found. The results are presented to the user in the form of a ranked list of these regions. In comparison with a full transcription of the document images, word spotting is much more robust against recognition errors. As long as the relevant regions are among the top ranks of the list, the user can eventually decide what he actually wants to use from the given results.<sup>3</sup>

Errors in a text transcription are not as easy to recognize and even small mistakes can corrupt further processing such as a full-text search.

In this scenario, the query is an exemplary occurrence of the respective word in the image. Query-by-example word spotting works only as long as the variability in the text is

## 1. Introduction

Text recognition in scanned documents usually refers to transcribing images showing text into a machine-

<sup>1</sup> Cf. Doermann and Tombre 2014, chap. 8–14.

<sup>2</sup> Cf. Lladós et al. 2012.

<sup>3</sup> Rath and Manmatha 2007.

low, as in single-writer scenarios or documents printed in a single font. Apart from that, no prior knowledge of the problem domain is required.

A widely recognized approach to word spotting was presented by Rath and Manmatha. Their method follows a relatively classic pattern recognition pipeline.<sup>4</sup> A document image is first segmented into word regions. After skew and slant normalization, each region is represented by a sequence of feature vectors. Words that are similar to the query are then found by computing distances between feature vector sequences using Dynamic Time Warping.<sup>5</sup> An effect of a prior document segmentation is that subsequent processing steps can be specifically designed for working with word regions. Examples include specialized features, such as the upper and lower word contour or the number of ink background transitions along the word image's columns.<sup>6</sup>

However, while this may seem advantageous, a severe disadvantage is that such a system does not recover from segmentation errors. It is implicitly assumed that perfect segmentation is possible. If that assumption fails, these errors cannot be handled and all further steps will be based upon them. For instance, if the word segmentation fails, the information encoded in the feature representation will not be useful. In addition, the processing steps in such a pipeline are only able to work in a locally optimal manner and no information regarding the actual objective, i.e. the recognition, can be taken into account. When segmenting documents into lines or words, a recognition step is already incorporated because knowledge of their appearance is assumed. These steps are usually based on heuristics and, consequently, the recognition will also be based on heuristics. In challenging unconstrained handwriting recognition scenarios, it is impossible to rely on knowledge justifying those assumptions. Therefore, we propose to avoid early decisions and, rather, to integrate as much information as possible into the final recognition.

The first methods for segmentation-free word spotting were presented in Leydier et al.<sup>7</sup> and Gatos and Pratikakis<sup>8</sup>.

Leydier et al. used a keypoint-based approach. Local zones-of-interest from the query word are matched with the most similar local zones in the document image. Regions are retrieved where the spatial configuration of matching zones in the query and the document image fit. Gatos and Pratikakis use a patch-based framework for segmentation-free word spotting. After preprocessing and normalization, they obtain text regions in a filtering step. Within these salient image regions, patches are sampled that are finally matched with the query. By over-segmenting the text regions, they are able to analyze all possible word locations.

In our previous work on segmentation-free word spotting,<sup>9</sup> we demonstrated how very accurate results can be achieved using a statistical sequence model. The model captures the spatial sequential structure of the query word in a dynamic probabilistic way. Similarity scores are then computed between the model and patches that were densely sampled over the document image.

However, a drawback of this method is the high computational effort that is required when processing a complete document image without focusing the recognition on potentially relevant regions as, for example, presented by Gatos and Pratikakis. A word spotting system that is not able to respond instantly is not likely to be accepted by potential users. In order to improve computational efficiency, we propose a two-stage approach. In the first stage, we identify potentially interesting, i.e. salient regions, in a fully segmentation-free manner. In this step, we do not incorporate any heuristic decisions but only identify regions in the document image that are visually similar to the query image. We implement this efficiently using an index structure for looking up small image patches in the document image that are similar to small image patches in the word region. This index is independent of a particular query and can be precomputed for any document image. In the second step, we investigate the salient regions in more detail. We apply our statistical sequence model for obtaining highly accurate matches. A major advantage of the approach is that both stages are based upon small image patches that can be found in the query and the document image. No additional preprocessing or feature extraction is required.

<sup>4</sup> Cf. Duda and Hart 2001, chap. 1.

<sup>5</sup> Cf. Rabiner and Juang 1993, 221-226.

<sup>6</sup> Marti and Bunke 2000.

<sup>7</sup> Leydier et al. 2009.

<sup>8</sup> Gatos and Pratikakis 2009.

<sup>9</sup> Rothacker et al. 2013; Rothacker, Rusiñol, and Fink 2013.

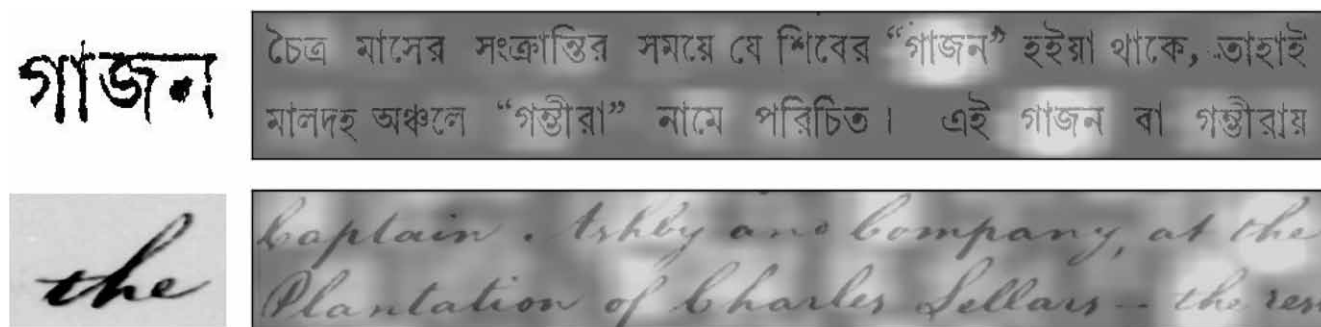


Fig. 1: Segmentation-free query-by-example word spotting. Top: Old printed Bangla book. Bottom: Historical George Washington dataset. Queries are shown on the left. On the right, the word occurrence probability is indicated by brightness.

## 2. Related work

A recent trend in document analysis is to adapt methods from computer vision. These algorithms usually have to work in completely unconstrained scenarios dealing with a huge variety of different objects and their appearances in natural scene images.<sup>10</sup> One of the best known approaches that has found its way into document analysis is Bag-of-Features (BoF).<sup>11</sup> These are statistically estimated feature representations that automatically integrate information from the problem domain, for example, document images. In contrast, a heuristically designed feature representation only captures the much more limited number of aspects that have been considered in the manual design process.

The basic idea behind BoF representations is to estimate feature representatives that are typical for the problem domain. When applied to images, gradient-based local image descriptors are very popular, for example SIFT descriptors,<sup>12</sup> which encode the local neighborhood at a given point in the image, i.e. a small image patch.

As the typical image descriptors are not directly available but have to be estimated statistically, three steps are required to compute a BoF. First, the typical image descriptors have to be estimated, usually by clustering a large set of descriptors from an image sample set.<sup>13</sup> Second, the descriptors in a given image have to be mapped to the most similar typical image descriptor.<sup>14</sup>

Finally, the BoF is obtained as the histogram of typical image descriptor frequencies. One of the first and most widely accepted methods for BoF in computer vision was presented by Sivic and Zisserman.<sup>15</sup> They demonstrated how objects can be efficiently retrieved from a large number of movie frames. For this purpose, objects as well as movie frames are represented by BoF representations. Large-scale retrieval is possible because features from the query can be efficiently localized in each movie frame with an inverted file structure.<sup>16</sup> This index saves all locations of each feature in a frame.

With respect to word spotting, the methods of Shekhar and Jawahar<sup>17</sup>, Rusiñol et al.,<sup>18</sup> and Almazán et al.<sup>19</sup> are especially relevant to our work. Shekhar and Jawahar use BoF representations for retrieving query words from large volumes of segmented word images. In analogy to Sivic and Zisserman, they use an inverted file structure for efficient indexing. The index stores pointers for each feature to all segmented word images that contain that specific feature. Given the features occurring in a query word image, the relevant word images in the database can be retrieved almost instantly. Rusiñol et al. were the first to apply BoF to segmentation-free word spotting. Almazán et al. showed how recognition accuracy and speed can be improved by using a Histogram-of-Oriented-Gradients in the same segmentation-free scenario. What both representations have in common is that no information about line and word locations is needed.

<sup>10</sup> Cf. Szeliski 2011, chap. 1.

<sup>11</sup> Cf. O'Hara and Draper 2011; Szeliski 2011, chap. 14.4.

<sup>12</sup> Lowe 2004.

<sup>13</sup> Cf. Gersho and Grey 1992, 362-370.

<sup>14</sup> Gersho and Grey 1992, chap. 10.

<sup>15</sup> Sivic and Zisserman 2003.

<sup>16</sup> Cf. Baeza-Jates and Ribeiro-Neto 1999, chap. 8.

<sup>17</sup> Shekhar and Jawahar 2012.

<sup>18</sup> Rusiñol et al. 2011.

<sup>19</sup> Almazán et al. 2012.

Instead, they are computed uniformly over the document image. While the BoF captures the simple occurrences of typical gradient-based local image features, the Histogram-of-Oriented-Gradients representation directly captures the image gradients in a grid of cells. For segmentation-free word spotting, the query image is transformed according to the respective feature description. This query feature description is then compared with patch feature descriptions that are densely extracted from the document image. This way, no prior assumptions about word positions have to be made but all possible locations are taken into account. An over-segmentation strategy of this kind is computationally very demanding. Rusiñol et al. reduced the search space by only considering a single patch size for all queries. Almazán et al. adapted the patch size to the query size. The huge amount of feature descriptions is efficiently compressed and stored in memory with product quantization.<sup>20</sup> Note, however, that limited variability of the script's visual appearance is assumed. Only a single patch size is used for spotting a certain query word. For that reason, it cannot be guaranteed that differently scaled instances of the query will be reliably found.

A drawback of the methods presented by Rusiñol et al. and Almazán et al. is limited flexibility with respect to modeling spatial gradient configurations. The more detailed the spatial information, the more specific becomes the feature representation of a query word image. While for single writer scenarios, a relatively explicit representation is advantageous, more abstraction will be needed once the variability increases. We, therefore, propose to integrate the BoF with a statistical sequence model, specifically Hidden Markov Models (HMM).<sup>21</sup> As shown in numerous examples, HMMs are able to model this spatial information in a dynamic probabilistic way.<sup>22</sup>

For a segmentation-free application, we adapted the patch-based framework presented in Rusiñol et al.<sup>23</sup> and Almazán et al.<sup>24</sup> but created a sequence of BoF representations from the query word image as well as from each patch. After the

BoF-HMM was estimated from the query, we obtained a probabilistic similarity score for each patch position with Viterbi decoding.<sup>25</sup>

Fig. 1 illustrates this for two datasets that we used for evaluating the method: old printed Bangla documents<sup>26</sup> and historic documents handwritten by George Washington and his associates.<sup>27</sup>

The probabilistic score maps show that the variability in the printed document scenario is much smaller than in the handwritten document scenario. But even though the detections in the handwritten word spotting case are less prominent, their scores are still good enough to be considered as most important in the given example. In terms of word spotting accuracy, we clearly outperformed the results reported by Rusiñol et al. and Almazán et al. on the same benchmark. In terms of computational speed, however, the massive generation of BoF sequences from patch representations and their evaluation with the Viterbi algorithm was considerably slower.

In the remainder of this paper, we will present a two-stage method for improving computational efficiency by only applying the BoF-HMM at document image locations that are salient with respect to the query word. The query-specific saliency map is based on an inverted file structure. This way, locations of typical features in the document image can be efficiently indexed. The major difference in computing the saliency map compared with Gatos and Pratikakis<sup>28</sup> is that we already integrate information from the query instead of just looking for arbitrary text areas. This makes the salient regions specific to the query and we can apply the BoF-HMM in a more focused manner.

In our experimental evaluation of the George Washington benchmark,<sup>29</sup> we will show that speed-ups of more than 50% are possible while the word spotting accuracy is only marginally affected.

<sup>20</sup> Jégou, Douze, and Schmid 2011.

<sup>21</sup> Rothacker, Vajda, and Fink 2012; Rothacker, Rusiñol, and Fink 2013.

<sup>22</sup> Cf. Fink 2014, chap. 5.

<sup>23</sup> Rusiñol et al. 2011.

<sup>24</sup> Almazán et al. 2012.

<sup>25</sup> Cf. Fink 2014, chap. 5.6.

<sup>26</sup> Rothacker et al. 2013.

<sup>27</sup> Rothacker, Rusiñol, and Fink 2013.

<sup>28</sup> Gatos and Pratikakis 2009.

<sup>29</sup> Rusiñol et al. 2011; George Washington Papers at the Library of Congress.

### 3. A two-stage approach for segmentation-free word spotting

The sole application of the BoF-HMMs for segmentation-free word spotting is very costly with respect to its computational efficiency. Running such a high precision method for all densely sampled patches on the document image also does not always seem to be appropriate. Most patches are visually very dissimilar to the query, which raises the question if these patches can be rejected with a more efficient approach.

In this section, we will present a method that performs a segmentation-free coarse-grained analysis of the document image followed by a fine-grained application of the BoF-HMM. The coarse-grained analysis produces a saliency map identifying regions in the document image that are similar to the query. The detailed analysis in the second stage is only applied in the local neighborhood of these salient locations. For the overall segmentation-free property of the method, it is very important that both stages rely on the same local image features. This means that if there is no indication of the query word in some region in the saliency map, the BoF-HMM also produces low scores in this particular area. The same features that did not match with the query in the first stage will not start matching in the second stage. For this reason, we do not incorporate any explicit prior segmentation step based on heuristic assumptions, such as distances between words or lines in the document image.

The two-stage method for segmentation-free word spotting consists of three processing steps: document image representation, model estimation and model decoding. As the computation of local image features is independent of any particular query and the features are shared among the two stages, they have to be computed once for each document image. In order to search for a query, it has to be modeled with the BoF-HMM. This incorporates a statistical model estimation procedure. Once the query model is available, it can be used for retrieving relevant regions in document images. The decoding step consists of the coarse-grained and the fine-grained analysis stages. After efficiently identifying regions in the document image that are similar to the query, these salient areas are analyzed in more detail using the BoF-HMM. Finally, regions are ranked according to their similarity with the query and presented to the user. The overall process is visualized in fig. 2 for an exemplary document image section of George Washington's letters.<sup>30</sup>

<sup>30</sup> George Washington Papers at the Library of Congress.

#### 3.1. Document image representation

Documents are represented by typical local image features that are extracted on a dense grid in the image. This is shown at the top of fig. 2. The descriptors that have been considered here (SIFT<sup>31</sup>) consist of histograms of oriented gradients and are intended for capturing the main directions of the pen stroke in the local neighborhood of a grid point. Each of the highly overlapping image features is an abstraction of the document's visual appearance representing mainly the information that is relevant for recognizing handwritten words. This has previously been demonstrated in many applications to document analysis.<sup>32</sup> The 'Dense Grid of Descriptors' in fig. 2 exemplarily shows the local neighborhood of a single image feature in the dense grid.

In the next step, typical image features are found with a cluster analysis using Lloyd's algorithm.<sup>33</sup> For this purpose, a sample set of local image features is required that is representative of the problem domain. In case of a word spotting system, all document images are known a priori and no unknown documents need to be considered. If new documents are added in the future, the cluster analysis can simply be repeated. The basic idea of the cluster analysis is to group similar image features into a given number of clusters. Each cluster has a representative that maximizes the average similarity to all its elements. Note that the representative is usually not among the set of local image features used in the cluster analysis. Finally, all descriptors in the dense grid are assigned to their most similar cluster representative, i.e., the typical image feature. This process is known as quantization. In fig. 2, the 'Descriptor Quantization' visualizes the dense grid with points. The points' colors indicate the typical image feature that they have been assigned to. Note how similar color patterns correspond to similar patterns of the pen stroke in the section of the document image. In the following, words will be spotted by the occurrence of typical image features in the document image that also appear in the query word region.

#### 3.2. Model estimation

When users want to query the word spotting system, they must select an exemplary occurrence of the word in a

<sup>31</sup> Lowe 2004.

<sup>32</sup> Rusiñol et al. 2011; Rothacker, Rusiñol, and Fink 2013; Shekhar and Jawahar 2012; Lladós et al. 2012.

<sup>33</sup> Cf. Gersho and Grey 1992, 362–370.



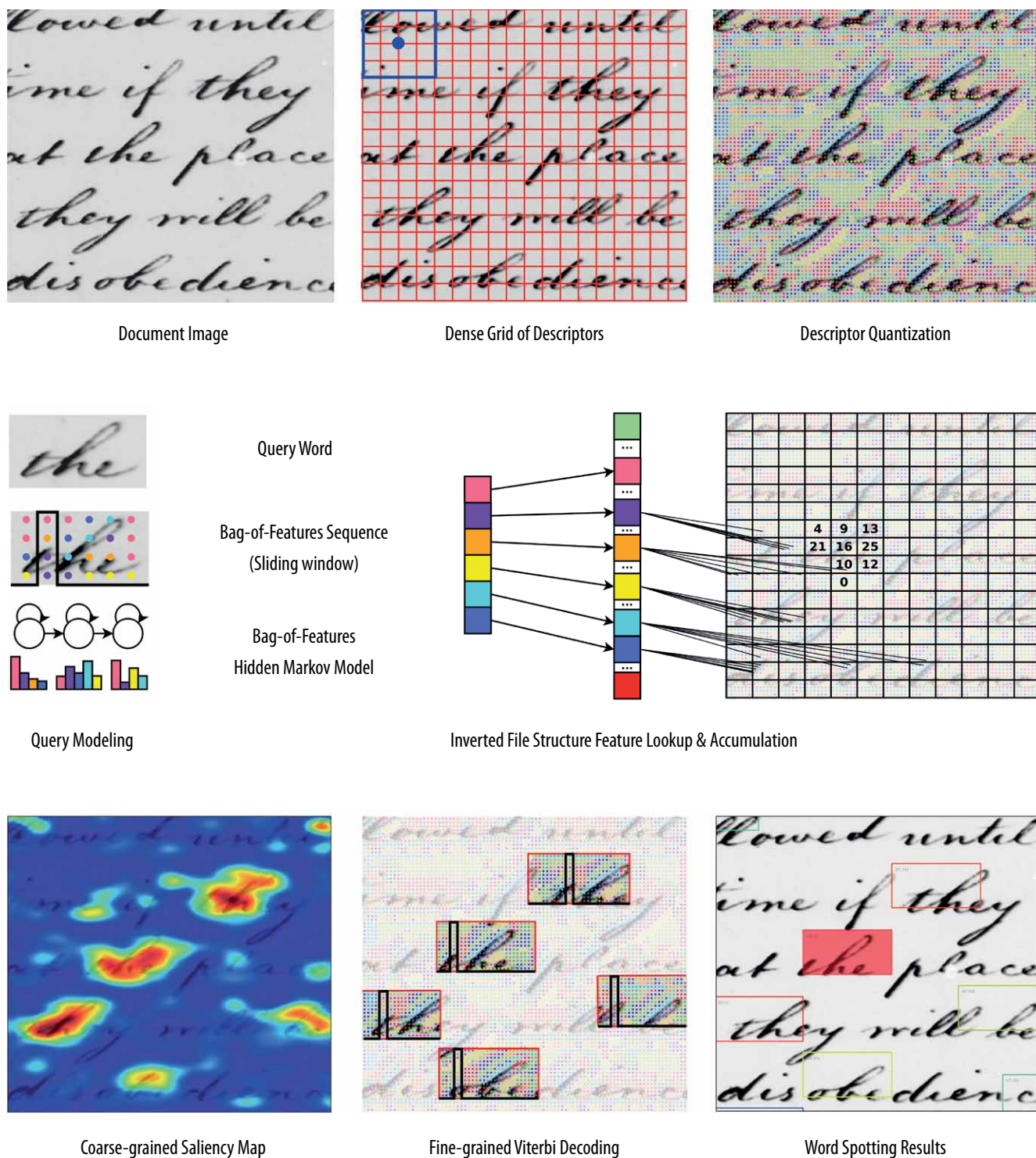


Fig. 2: Overview of the two-stage word spotting method. Top: Document image representation. Each document image is represented by typical local image features. In the dense grid these are indicated with different colors. Middle: Query word modeling and inverted file structure lookup & accumulation. Some accumulator cells show counts of how many features from the query model have been detected. Bottom: Coarse-grained saliency map and fine-grained analysis with the BoF-HMM. Saliency scores are indicated with blue to red colors. For the fine-grained analysis the BoF sequence extraction is visualized for all detected regions. The rank of the finally retrieved regions is indicated with blue to red colors.

document image. The respective region also specifies the typical image features that represent the query word. In order to capture the query word's spatial sequential structure, we extract a BoF representation for each column within the dense grid. The BoF only captures the relative frequency of each typical image feature, not their order.

This will be beneficial in the subsequent decoding step. The sequence of BoF representations is then modeled with a BoF-HMM. HMMs are generative finite state machines. At each point in time, one state is active and generates an output according to an underlying statistical process. In the estimation procedure, parameters are optimized with respect to the probability of generating this sequence of BoF representations from the query word. The process is referred to as the Baum-Welch algorithm.<sup>34</sup> In the query-by-example scenario, it is important that the BoF-HMM can be estimated with just a single sample.<sup>35</sup> This is usually impossible with continuous HMMs.<sup>36</sup> Fig. 2 exemplarily sketches the model estimation for the query word 'the'. Below the query word region, the extraction of the BoF sequence is visualized. A sliding window, shown in black, is moved over the dense grid in the direction of writing. At each window position, a histogram of typical image feature frequencies is created. For the window position, visualized in fig. 2, the purple feature has a higher frequency than the red and orange features. The BoF-HMM models feature probabilities for the query word directly within the states. In the given example, the states roughly represent the features occurring at the beginning, in the middle and at the end of the query word. For example, the cyan features have a higher probability in the middle, while the red features are rather probable at the beginning and the end.

### 3.3. Model decoding

After the model has been estimated, it can be used for retrieving regions in document images that are similar in terms of typical image feature occurrences. The overall process consists of two stages. In the first stage, a coarse-grained analysis is performed. It roughly identifies areas in a document image containing typical image features that also occur in the query word. In the second stage, these regions

are explored in more detail by measuring their similarity to the BoF-HMM.

#### 3.3.a First stage

The most important component in the first stage is an inverted file structure that we use for localizing typical features in the document image. A similar approach has been used for object detection and retrieving segmented word images. The inverted file structure index contains an entry for each typical image feature that has been determined in the cluster analysis. Each entry indexes all feature locations. The inverted file structure has to be computed only once and feature location lookups become very efficient afterwards. In our scenario, we want to look up features from the query word. Every feature with a non-zero probability in any of the states of the query model will be localized this way. For an actual occurrence of the query word in a document, we expect to find a larger number of detections in roughly the same area. We, therefore, accumulate detections in a cell structure over the document image. For each cell, the number of feature detections is counted. This is similar to the generalized Hough transform.<sup>37</sup>

In fig. 2, the inverted file structure lookup and accumulation is visualized for the exemplary query. All features occurring in the query model are localized in the section of the document image. Note that only a few features with respect to the total number of features will actually occur in a specific query word region. Also, not all features might reappear in a document. In fig. 2, this is indicated for the red feature. Inverted file structure lookups are very efficient for these reasons.

The resulting accumulator matrix can be interpreted as a coarse-grained saliency map. It is coarse-grained because no spatial relations between features have been taken into account yet. In fig. 2, saliency measures are visualized with blue to red colors. Due to the coarse-grained character, detections are usually not precisely located over the relevant regions.

Regions-of-interest that are positioned at the peaks in the saliency map are the final output of the first stage, i.e. a region is created for each locally optimal score. The regions' sizes are equal to the size of the query word. The locally optimal scores are determined in such a way that the regions do not overlap.

<sup>34</sup> Cf. Fink 2014, chap. 5.7.

<sup>35</sup> Rothacker, Rusiñol, and Fink 2013.

<sup>36</sup> Plötz and Fink 2011.

<sup>37</sup> Cf. Szeliski 2011, chap. 4.3.2.

Table 1: Evaluation of the two-stage approach for segmentation-free word spotting

Method	RoI dilation	Overlap threshold	mAP	mR	mQT	Speed-up
BoF	—	50%	30.4%	71.1%	340 ms	
HoG	—	50%	54.4%	—	15 ms	
VT	—	50%	69.7%	83.0%	4100 ms	Baseline
IFS	—	50%	50.1%	60.1%	380 ms	10x
IFS+VT	—	50%	60.4%	60.0%	540 ms	7x
IFS+VT	3x3	50%	64.9%	69.2%	980 ms	4x
IFS+VT	5x5	50%	69.6%	80.7%	1830 ms	2x
IFS+VT	7x7	50%	70.0%	82.5%	2750 ms	1.5x
VT	—	25%	71.4%	96.8%	4100 ms	Baseline
IFS	—	25%	48.0%	93.1%	380 ms	10x
IFS+VT	—	25%	53.9%	92.1%	540 ms	7x
IFS+VT	3x3	25%	64.7%	94.7%	990 ms	4x
IFS+VT	5x5	25%	70.9%	96.2%	1830 ms	2x
IFS+VT	7x7	25%	71.6%	96.3%	2760 ms	1.5x

### 3.3.b Second stage

The regions-of-interest from the first stage are analyzed in more detail in the second stage. A sequence of BoF representations is extracted from each region. This captures the spatial sequential structure in the writing direction. As no information about the feature order is encoded in a BoF, these representations are robust against smaller vertical displacements of the detected region with respect to the relevant region in the document. The horizontal displacements can be handled by the BoF-HMM, which models the sequential structure in a dynamic probabilistic way. The probability of generating the BoF sequence from a detected region can be computed with the Viterbi algorithm. Each region obtains a new similarity score with respect to the query model. However, even though the BoF-HMM is relatively robust against displacements of the detected regions, better scores can be obtained when these regions exactly fit the relevant regions in the document. For this reason, we dilate the search area locally in order to overcome the coarse detections of the first stage. Finally, we extract detected regions with locally optimal scores and rank them accordingly. Fig. 2 visualizes the fine-grained evaluation of the HMM and the retrieved regions that are presented to

the user. The regions' scores are indicated with blue to red colors.

In the following evaluation, we will focus on performance measures in terms of the two stages. We will evaluate their individual and joint accuracy as well as their speed. The region-of-interest dilation will be of special interest with respect to the joint evaluation of both stages.

## 4. Evaluation

The effect of the two-stage method is evaluated using papers written by George Washington and his associates.<sup>38</sup> The full collection consists of over 65,000 documents and covers many aspects 'of colonial and early American history'.<sup>39</sup> Document types include correspondences, diaries, journals, military records, notes etc. that were collected by George Washington from 1741 to 1799. At the Library of Congress, the collection is organized accordingly into nine series. For the word spotting benchmark, a small dataset of 20 pages that are in overall good condition has been compiled from 'Series

<sup>38</sup> Rusiñol et al. 2011.

<sup>39</sup> George Washington Papers at the Library of Congress.



2: Letterbooks 1754–99’ and consists of pages 270–279 and 300–309. The dataset was first used for evaluating a word spotting system by Rath and Manmatha.<sup>40</sup> In order to measure and compare segmentation-free word spotting performance, Rusiñol et al.<sup>39</sup> defined a benchmark containing 4,860 queries from these 20 pages. Ground truth annotations consisting of a bounding box in the document image and a word label are available for all words. The pages are written in an overall very similar style, thus the benchmark can be considered as a single writer scenario. Following this evaluation protocol, we use every word as a query without any further modification, like filtering short words or stemming words. For each query, we retrieve a ranked list of regions throughout all 20 pages. By comparing each of those regions with the ground truth annotations, we can decide if a single region is relevant with respect to the query or not. A region is considered relevant if it overlaps with a bounding box from the ground truth by more than a given threshold and the corresponding word label matches the query word. The choice of the overlap threshold is critical for performance measures and we will report results for two different values in the following. In table 1, we refer to this as the ‘Overlap threshold’.

Given the list of relevant and non-relevant regions, two aspects are of prime importance to users of word spotting systems. All relevant regions should be ranked first and the list should contain all relevant words.<sup>41</sup> The first requirement is measured by average precision, while the second requirement is measured by recall. An average precision of 100% refers to a list where all relevant regions are listed first. A recall of 100% refers to a list that is complete, i.e., no relevant regions are missing. In order to report results over all queries, we compute means over the individual results, thus mean average precision and mean recall. In table 1, we refer to these measures as ‘mAP’ and ‘mR’.

The motivation for the two-stage approach is to improve the computational speed. Mean runtimes for retrieving a single query on a single page are reported as well as the relative speed-ups with respect to our baseline method.<sup>42</sup> Results have been measured on a Xeon 3.0 GHz. In table 1, we refer to the mean query time per page as ‘mQT’ and to its relative improvement as ‘Speed-up’.

Our experiments focus on the results obtained individually and jointly with the stages. Additionally, we report results for different overlap thresholds. As the first stage is based on an inverted file structure, it is referred to as ‘IFS’ in table 1. Analogously, scores in the second stage are computed with the Viterbi algorithm, thus referring to it as ‘VT’. For the individual application of the first stage (IFS), we directly use the regions-of-interest as the output of the word spotting system. For the individual application of the second stage (VT), we use the same approach as presented in Rothacker, Rusiñol, and Fink,<sup>43</sup> i.e. we densely sample patches that are all decoded with the Viterbi algorithm for obtaining similarity scores. Please note that the results reported differ due to some parametric optimizations. In the joint evaluation of both stages, the region-of-interest dilation is important. In table 1, the different dilation masks can be found in the column ‘RoI dilation’. A dilation mask of 3x3 refers to extending the search area to all neighboring cells in the accumulator matrix from a detected region in the first stage. Larger masks extend the search area further.

Starting with the results obtained for our baseline system (VT),<sup>44</sup> we observe an upper boundary for retrieval accuracy. This is expected as we are applying a fine-grained analysis throughout entire document images. When, in contrast, only the coarse-grained analysis in first stage is applied, a significant drop in mean average precision can be observed. As the spatial sequential structure of the query word is not modeled in the coarse-grained analysis, the simple occurrence of features from the query words is already a strong indication for a relevant region. Problems occur, for example, when single or multiple characters from the query word appear within other, typically longer words. When putting both stages together, we first rerank the results with the Viterbi algorithm without dilating the search area. While the mean recall stays constant, there is a considerable increase in mean average precision. This nicely demonstrates the effect of adding spatial information to the query word modeling. The constant mean recall is due to the fact that the search space has not been extended in the second stage. The list of relevant regions stays the same.

When increasing the search area, improvements can be observed for both mean average precision and mean recall. Now, the regions-of-interest can be positioned better over

<sup>40</sup> Rath and Manmatha 2007.

<sup>41</sup> Cf. Baeza-Yates and Ribeiro-Neto 1999, chap. 3; Lladós et al. 2012, 13–15.

<sup>42</sup> Rothacker, Rusiñol, and Fink 2013.

<sup>43</sup> Ibid..

<sup>44</sup> Ibid.

the relevant regions in the document image and the BoF-HMM produces better matches when being evaluated with the Viterbi algorithm. Regions-of-interest that overlap only slightly with relevant regions in the document are discarded for better matches in the local neighborhood. Additionally, more regions will be regarded as relevant if the detections from the first stage were not precise enough to produce sufficient overlap percentages.

Regarding the overlap threshold, a reduction from 50% to 25% shows that over 90% of the relevant words can roughly be located in the first stage. The localizations are just too imprecise to produce over 50% overlaps. This nicely motivates the application of the regions-of-interest dilation.

Finally, we measured strong improvements in the mean query time with respect to our baseline system. The sole application of the inverted file structure is more than ten times faster. Using the two-stage approach, the size of the search area (RoI dilation) is strongly related to the mean query time and recognition accuracy. When matching the accuracy of the baseline system, a speed-up by a factor of two is still possible.

When comparing our results with BoF-based results reported in Rusiñol et al.<sup>45</sup> and the Histogram-of-Oriented-Gradient (HoG)-based results reported in Almazán et al.,<sup>46</sup> we clearly outperform their recognition accuracy with our integrated approach. Low recognition scores in Rusiñol et al. are due to the uniform patch size used for all queries in the retrieval stage. Small words that are hard to detect because patches contain a lot of context in the document image besides the relevant word. Results reported in Almazán et al. show good recognition accuracy and very fast retrieval times. The HoG features nicely encode handwritten words in the single-writer scenario. In comparison with our two-stage approach, we see two important aspects. Encoding full-page HoG representations in memory does not scale for very large collections of document images. By contrast, the inverted file structure in our coarse-grained stage is able to rapidly reduce the search space over all document pages in the collection. This capability has already been demonstrated for segmented word images.<sup>47</sup> Finally, the application of Hidden Markov Models offers more flexibility, especially when more than a single sample of a query is available.

Please note that the query times for neither Rusiñol et al. nor Almazán et al. are directly comparable because different machines have been used in the evaluation. However, they are suitable for showing the general order of magnitude in which the methods operate.

## 5. Conclusion

In this paper, we presented a two-stage approach for segmentation-free word spotting based on the George Washington benchmark. With respect to our baseline system, considerable improvements in computational speed have been observed. With the sole application of the coarse-grained first stage, speed-ups of more than ten times are possible. This comes at the cost of decreased recognition accuracy. When adding the fine-grained word spotting stage, the benefits of integrating a statistical sequence model can be demonstrated. The recognition accuracy increases substantially. Computational speed can still be reduced to 50% without losing any precision. In comparison with Rusiñol et al.<sup>48</sup> and Almazán et al.,<sup>49</sup> more accurate results can be achieved at the cost of higher query execution times.

Furthermore, we present the effect of different overlap thresholds for the evaluation. We can show that the fine-grained stage handles regions-of-interest having smaller overlaps with the ground truth very well. For larger dilation masks, very high mean recall scores can be achieved without any loss in mean average precision. We doubt that an overlap threshold of 25% or 50% will make any difference for users of word spotting systems as detected regions will usually be presented in a larger context in the document image.

In our future research, we will work on better region localizations in the coarse-grained stage. This could be achieved by taking some limited amount of the query word's spatial structure into account.

## ACKNOWLEDGMENTS

This work is partially supported by the German Academic Exchange Service on the basis of a DAAD-Doktorandenstipendium and by the Spanish project TIN2012-37475-C02-02.

<sup>45</sup> Rusiñol et al. 2011.

<sup>46</sup> Almazán et al. 2012.

<sup>47</sup> Shekhar and Jawahar 2012.

<sup>48</sup> Rusiñol et al. 2011.

<sup>49</sup> Almazán et al. 2012.

## REFERENCES

- Almazán, J., Gordo, A., Fornés, A., and Valveny, E. (2012), 'Efficient exemplar word spotting', *Proceedings of the British Machine Vision Conference*, 67.1–67.11.
- Baeza-Yates, R., and Ribeiro-Neto, B. (1999), *Modern Information Retrieval* (Addison Wesley).
- Doermann, D., and Tombre, K. (2014), *Handbook of Document Image Processing and Recognition* (Springer).
- Duda, R. O., Hart, P. E., and Stork, D. G. (2001), *Pattern Classification*, 2nd edition (Wiley).
- Fink, G. A. (2014), *Markov Models for Pattern Recognition, From Theory to Applications*, 2nd edition (Springer; Advances in Computer Vision and Pattern Recognition).
- Gatos, B., and Pratikakis, I. (2009), 'Segmentation-free word spotting in historical printed documents', *Proceedings of the International Conference on Document Analysis and Recognition*, 271–275.
- George Washington Papers at the Library of Congress, 1741–1799*, Manuscript Division, Library of Congress, Washington, D.C. <http://memory.loc.gov/ammem/gwhtml/gwhome.html>.
- Gersho, A., and Grey, R. M. (1992), *Vector Quantization and Signal Compression* (Kluwer Academic; Communications and Information Theory).
- Jégou, H., Douze, M., and Schmid, C. (2011), 'Product quantization for nearest neighbor search', *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33: 117–128.
- Marti, U. V., and Bunke, H. (2000), 'Handwritten Sentence Recognition', *Proceedings of the International Conference on Pattern Recognition*, 3: 463–466.
- Leydier, Y., Ouji, A., LeBourgeois, F., and Emptoz, H. (2009), 'Towards an omnilingual word retrieval system for ancient manuscripts', *Pattern Recognition*, 42.9: 2089–2105.
- Lladós, J., Rusiñol, M., Fornés, A., Mota, D. F., and Dutta, A. (2012), 'On the influence of word representations for handwritten word spotting in historical documents', *International Journal of Pattern Recognition and Artificial Intelligence*, 26.5.
- Lowe, D. (2004), 'Distinctive Image Features from Scale-Invariant Keypoints', *International Journal of Computer Vision*, 60.2: 91–110.
- O'Hara, S., and Draper, B. A. (2011), *Introduction to the Bag of Features Paradigm for Image Classification and Retrieval*, (Cornell University Library, <http://arxiv.org/abs/1101.3354>).
- Plötz, T., and Fink, G. A. (2011), *Markov Models for Handwriting Recognition* (SpringerBriefs in Computer Science, Springer).
- Rabiner, L. R., and Juang, B.-H. (1993), *Fundamentals of Speech Recognition* (Prentice-Hall, Englewood Cliffs).
- Rath, T., and Manmatha, R. (2007), 'Word spotting for historical documents', *International Journal on Document Analysis and Recognition*, 9.2–4: 139–152.
- Rothacker, L., Fink, G. A., Banerjee, P., Bhattacharya, U., and Chaudhuri, B. B. (2013), 'Bag-of-Features HMMs for Segmentation-Free Bangla Word Spotting', in *International Workshop on Multilingual OCR*, article no. 5.
- , Rusiñol, M., and Fink, G. A. (2013), 'Bag-of-Features HMMs for Segmentation-Free Word Spotting in Handwritten Documents', *Proceedings of the International Conference on Document Analysis and Recognition*, 1305–1309.
- , Vajda, S., and Fink, G. A. (2012), 'Bag-of-Features Representations for Offline Handwriting Recognition Applied to Arabic Script', *Proceedings of the International Conference on Frontiers in Handwriting Recognition*, 149–154.
- Rusiñol, M., Aldavert, D., Toledo, R., and Lladós, J. (2011), 'Browsing heterogeneous document collections by a segmentation-free word spotting method', *Proceedings of the International Conference on Document Analysis and Recognition*, 63–67.
- Shekhar, R. and Jawahar, C. (2012), 'Word image retrieval using bag of visual words', *International Workshop on Document Analysis Systems*, 297–301.
- Sivic, J., and Zisserman, A. (2003), 'Video Google: A text retrieval approach to object matching in videos', *Proceedings of the International Conference on Computer Vision*, 2: 1470–1477.
- Szeliski, R. (2011), *Computer Vision, Algorithms and Applications* (Springer).