# Fast Approximate Modelling of the Next Combination Result for Stopping the Text Recognition in a Video

Konstantin Bulatov
Smart Engines Service LLC,
FRC CSC RAS, Moscow, Russia
Email: kbulatov@smartengines.com

Nadezhda Fedotova
Smart Engines Service LLC,
Moscow, Russia
Email: nfedotova@smartengines.com

Vladimir V. Arlazarov
Smart Engines Service LLC
MIPT, Moscow, Russia
Email: vva@smartengines.com

*Abstract*—In this paper, we consider a task of stopping the video stream recognition process of a text field, in which each frame is recognized independently and the individual results are combined together. The video stream recognition stopping problem is an under-researched topic with regards to computer vision, but its relevance for building high-performance video recognition systems is clear.

Firstly, we describe an existing method of optimally stopping such a process based on a modelling of the next combined result. Then, we describe approximations and assumptions which allowed us to build an optimized computation scheme and thus obtain a method with reduced computational complexity.

The methods were evaluated for the tasks of document text field recognition and arbitrary text recognition in a video. The experimental comparison shows that the introduced approximations do not diminish the quality of the stopping method in terms of the achieved combined result precision, while dramatically reducing the time required to make the stopping decision. The results were consistent for both text recognition tasks.

## I. Introduction

Video processing has become a rich and dynamic branch of the research in the computer vision field. The problems of video stream analysis include object detection and segmentation [1], [2], object tracking [3], [4], super-resolution [5], text recognition [6] and many more. Modern computer vision applications employ methods of automated video processing not only for high-end industrial vision systems but for mobile devices such as smartphones or tablet computers as well.

One of the most relevant computer vision problems is automated data entry by means of text recognition. Text detection and recognition have applications in such fields as business processes automation, road traffic monitoring, government services, mobile payments [7]–[9], life augmentation for people with disabilities [10] and more. The research is targeted on improving the speed and accuracy of camera-based information extraction, given particular challenging conditions, such as poor illumination, low camera quality [11], optical distortions and noise, poor focus and motion blur [12]. Each year new and improved methods for arbitrary scene text detection recognition are published [13]–[15], including the works which focus on processing text in videos [16], [17]. An important requirement in many applications of text recognition



Fig. 1. Illustration of identity document recognition and data extraction on a mobile device

is for the system to be able to operate in real time, which is especially relevant for such use cases as recognition of road scene objects such as traffic signs [18], assisting the visually impaired [19]–[21], and others.

A special case of text recognition is represented within mobile identity document recognition systems. The recognition of identity documents is encumbered by specific features of such documents, e.g. textured background, holographic security elements which are obstructing the text recognition, reflective document surfaces which is prone to highlights, etc. [22]. At the same time, an important aspect of identity document recognition systems is their low error tolerance – the cost of recognition mistakes are high, as the recognized data is then used for personal identification, government services, financial transactions and in other sensitive fields. The scope of computer vision problems related to identity documents recognition includes document detection and location [23], [24], document layout analysis [25], face detection [26], and, of course, text fields recognition [27]–[29]. Fig. 1 illustrates the use case for information extraction from an identity document using a mobile device camera.

Using video input in recognition systems presents an opportunity to reduce the text recognition errors and thus increase the information extraction reliability, both in the context of arbitrary text recognition or in a more specialized case of document fields recognition. Combination of multiple recognition results of the same object obtained from different frames has been shown to be an effective way of improving the recognition accuracy in a video [27], [30]. This approach, however, gives rise to another problem – how to decide when there is enough accumulated information and the video stream recognition process should be stopped. Without access to the ground truth, the stopping methods should be able to make a decision whether the result could be improved or not and whether it is justified to spend more time to recognize additional frame in an effort to improve the combined result. The stopping rules are particularly crucial for real-time tracking and recognition of multiple objects, such as vehicle license plates or traffic signs [18].

However, after extensive research, we have found that the problem of optimal stopping is left almost unexplored in the field of computer vision, despite its importance for video processing. At the same time, optimal stopping is a known problem in the field of decision theory, mathematical statistics, and economics, and new theoretical results continue to be produced for its different variations [31], [32]. A few methods have been proposed for the problem of stopping the video stream recognition process [33], [34] in the context of document recognition. The method proposed in [33] is based on the clustering of the input results sequence and making the stopping decision based on the statistical characteristics of the obtained clusters. The stopping method described in [34] is based on the modelling of the next combined result and making the stopping decision based on the estimated expected distance between the current combined result and the next one. The latter method was tested on text recognition results using both Levenshtein-based and end-to-end text string distance metrics and exhibited higher effectiveness in comparison with the clusters analysis method [33]. The method was also tested for text recognition result model with per-character class membership estimations and also proved to be more effective [35] in comparison with the other methods.

All the related works considered only the combined result accuracy characteristics and the achieved mean number of processed frames, without any attention to the time required to compute the necessary estimations on which the stopping decision is based. In particular, the process of modelling of the combined result at the next stage, proposed in [34] has a high computational complexity which could diminish the positive effects of the stopping method, especially if executed in real time on a mobile device. The goal of this paper is to construct an approximate method of modelling of the next combined result required to make the stopping decision according to the method described in [34], which would have reduced computational complexity, and evaluate the constructed methods for the tasks of arbitrary text recognition as well as document fields recognition in videos.
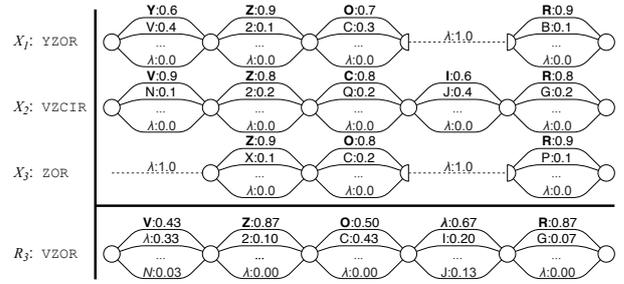


Fig. 2. An illustration of per-frame text string recognition results alignment and combination

Section II provides an overview of the studied stopping method. The method was originally proposed in [34] and further evaluated in [35]. In section III the proposed approximations and optimizations are described. The experimental evaluation of the base stopping method and the proposed optimizations is presented in the final section IV.

## II. EXISTING METHOD DESCRIPTION

In this section, we will provide a detailed description of the procedure of combining the individual per-frame text string recognition results using the algorithm described in [30], and stopping the video stream recognition process using the method described in [34], [35].

A text string recognition result $X$ with per-character alternatives can be represented as a matrix:

$$X = (x_{jk}) \in [0.0, 1.0]^{M \times K}, \quad \forall j \quad \sum_{k=1}^{K} x_{jk} = 1.0, \quad (1)$$

where $M$ is the length of the string (number of characters) and $K$ is the number of character classes. Each row $(x_{j1}, x_{j2}, \ldots, x_{jK})$ of the matrix represents the classification result for each individual character and contains membership estimations for each class. The value of each membership estimation is a real number from the range $[0.0, 1.0]$, and the sum of all membership estimations for each given character classification result equals to $1.0$. The combined recognition result of the text string in a video stream is represented with the same data structure.

During the recognition process, we observe a sequence of text string recognition results $X_1, X_2, \ldots$. The goal is to produce a combined result with the highest accuracy, that is, the closest possible to the correct text string value $X^*$ in terms of some predefined metric. In a more general case, there is also a non-negative weight $w_i$, associated with each $X_i$, which represents the desired contribution of the result $X_i$ in the combination.

On stage $n$, after obtaining the observation $X_n$, it is combined with the previously accumulated recognition results as follows. Firstly, each character classification result is expanded with an "empty" class label $\lambda$ with membership estimation $0.0$. In terms of the matrix representation, this corresponds to adding a zero-valued column at the beginning of the matrix (1). Then an alignment is calculated between $X_n$ and the

previously obtained combined result $R_{n-1}$ using a dynamic programming procedure to determine the optimal matching between rows of $X_n$ and $R_{n-1}$. After the character alignment is determined the corresponding character classification results are combined by calculating a weighted average of membership estimations for each class, and using an "empty" classification result (with class $\lambda$ having membership estimation 1.0) for pairing with unaligned characters.

Fig. 2 illustrates the alignment and the combination result for three text string recognition results. For ease of visual representation, each frame result is represented as a weighted identity transducer [36] with each character classification result corresponding to a set of labelled transitions between two consequent states.

After the combination is performed and the result $R_n$ is obtained, the stopping method is applied to make a decision whether $R_n$ should be returned as a final result or whether the process should continue (i.e. the observation $X_{n+1}$ should be acquired). Stopping method introduced in [34] operates under the assumption that the expected distances between two consecutive combined results do not increase from stage to stage. Under this assumption the problem can be viewed as a monotone stopping problem, at least starting from a certain stage, and an optimal stopping rule for it should behave in a "myopic" way, i.e. make the decision as if the next stage of the process will be the last one. To approximate an optimal stopping rule on stage $n$ the expected distance is calculated from the current combined result $R_n$ to the next $R_{n+1}$. The process is stopped if this expected distance is not higher than the predefined threshold which represents the cost of each observation.

The proposed method of calculation of the expected distance to the next combined result is to perform modelling of the next result by sampling already accumulated observations as candidates for the next one:

$$\hat{\Delta}_n = \frac{1}{n+1}\left(\delta + \sum_{i=1}^{n} \rho\left(R_n, R(X_1, \ldots, X_n, X_i)\right)\right), \quad (2)$$

where $n$ is the stage number, $\delta$ – external parameter, $\rho$ is a metric function defined on the set of text string recognition results, $R_n$ is a combined result obtained on the $n$-th stage, and $R(X_1, \ldots, X_n, X_i)$ is the combined result obtained by testing the individual frame recognition result $X_i$ as a candidate for the next observation.

For the validity of the stopping methods as it is described in [34] any metric function $\rho$ between the text string recognition results can be used, provided that it satisfies the triangle inequality. For experiments in [34] and [35] a normalized version of the Generalized Levenshtein Distance (GLD) [37] was used. For calculating GLD a metric $\rho_C$ must be defined on the set of individual character classification results (on individual rows of the matrices). For this purpose a scaled taxicab metric can be used:

$$\rho_C(a, b) = \frac{1}{2}\sum_{k=0}^{K}|a_k - b_k|, \quad (3)$$

where $a$ and $b$ are two matrix rows, $a_0$ and $b_0$ are the respective membership estimations for the "empty" class $\lambda$, and $a_k$ and $b_k$ are the respective membership estimations for character classes in the alphabet for all $k > 0$.

The described stopping method was tested on text string recognition results without membership estimations [34], using ROVER [38] as a combination algorithm. Later it was tested with another per-frame recognition method and with an extended text string recognition result model [35] using a combination algorithm [30] based on a ROVER approach. In both experiments, it was shown that using such an approach for a given average number of processed frames the lower error level could be achieved and vice versa.

However, the obvious downside of this stopping method is the complexity of the decision making algorithm, in particular the time required to compute the expected distance estimation (2). Given two text string recognition results $X$ and $Y$ with lengths $|X|$ and $|Y|$ the complexity of their combination using algorithm [30] is $O(|X| \cdot |Y| \cdot K)$, as the alignment procedure requires $O(|X| \cdot |Y|)$ calculations of the individual character classification metric function $\rho_C$ (3). For the same reason, the complexity of calculating GLD $\rho(X, Y)$ is also $O(|X| \cdot |Y| \cdot K)$. Consider $M$ as the maximal length of individual results $X_i$, and $S_n$ as the length of combined result $R_n$ (both are in terms of the number of rows). The complexity of performing each test combination required for computing the sum in (2) is $O(S_n M K)$ and the worst-case estimation for the length of the next combined result candidate is $O(S_n + M)$. Thus, to compute the distance sample from the current combined result to the next, another $O(S_n K(S_n + M))$ operations have to be performed. The total aggregate complexity of computing the expected distance estimation $\hat{\Delta}_n$ (2) on stage $n$ is $O(nS_n K(S_n + M))$.

## III. PROPOSED OPTIMIZATIONS

In this section we will introduce approximations which would help us to compute the approximate value of the estimation (2) more efficiently.

### A. General approximations

To obtain a more efficient method for computing the expected distance estimation (2) let us introduce the following approximations:

**Approximation 1:** naive alignment. During the test combination $R(X_1, \ldots, X_n, X_i)$ the candidate frame recognition result $X_i$ has to be aligned with $R_n$, which already contains $X_i$ (as it was aligned with $R_{i-1}$ on the $i$-th stage of the process). As an approximation of this alignment, we will assume that rows of $X_i$ will be aligned with the same rows of $R_n$ with which the corresponding components were combined on stage $i$. This coarse assumption will allow to skip the costly alignment for each test combination, as the row indices of $R_n$ with which the rows of $X_i$ should be aligned are known in advance.

**Approximation 2:** naive Levenshtein. Given approximation 1 the length of the combined result $R(X_1, \ldots, X_n, X_i)$

stays they same as the length of $R_n$. We will then assume that alignment of $R_n$ and $R(X_1, \ldots, X_n, X_i)$, which is required to compute the GLD between them, is direct, i.e. each $j$-th row of $R_n$ is aligned with $j$-th row of $R(X_1, \ldots, X_n, X_i)$. Thus, the distance between $R_n$ and $R(X_1, \ldots, X_n, X_i)$ is a sum of distances in terms of the scaled taxicab metric $\rho_C$ (3) between their rows with the same index.

To quickly compute the approximate value of an expected distance estimation (2), during the combination of input results $X_1, X_2, \ldots$ we will maintain and update a three-dimensional matrix $Y_n$:

$$Y_n = (y_{ijk}) \in [0.0, 1.0]^{n \times S_n \times (K+1)}, \qquad (4)$$

such that $i$ is the index of input per-frame result, $j$ is the index of a row of the current combined result, $k$ is the character class index, $y_{ijk}$ is a membership estimation of class $k$ from a row of input $X_i$ which was aligned and merged into the $j$-th row of the current combined result $R_n$.

Given the approximations 1 and 2, using a GLD as the metric function $\rho$, for $j$-th component of the combined result $R_n$ and for each individual character class index $k \in \{0, 1, \ldots, K\}$ ($k = 0$ being the index of an "empty" class label $\lambda$), considering $X_i$ as a candidate for the next frame recognition result adds the following contribution to the sum of distances in (2):

$$\Delta_{ijk} = \frac{1}{2} \left| \frac{A_{jk}}{W} - \frac{A_{jk} + y_{ijk} w_i}{W + w_i} \right|, \qquad (5)$$

where $w_i$ is the weight associated with an input per-frame result $X_i$; $A_{jk}$ is the weighted sum of membership estimations of class $k$ corresponding to $j$-th component of the combined result: $A_{jk} = \sum_{i=1}^{n} y_{ijk} w_i$; and $W$ is the sum of all weights: $W = \sum_{i=1}^{n} w_i$.

The approximation of the expected distance estimation (2) can now be computed as follows:

$$\hat{\Delta}_n \approx \frac{1}{n+1} \left( \delta + \sum_{i=1}^{n} \sum_{j=1}^{S_n} \sum_{k=0}^{K} \Delta_{ijk} \right). \qquad (6)$$

Since for all $j$ and $k$ the weighted sum $A_{jk}$ of membership estimations can be computed on-the-fly during combination of $X_n$ and $R_{n-1}$, the approximate estimation (6) can be computed with complexity $O(nS_nK)$.

### B. Unweighted case

Often there are no weights $w_i$ associated with input text string recognition results $X_i$, i.e. all input recognition results have an equivalent contribution to the final combined result. In this case the expression (5) can be further simplified:

$$\Delta_{ijk} = \frac{1}{2} \left| \frac{A_{jk}}{n} - \frac{A_{jk} + y_{ijk}}{n+1} \right| = \frac{|A_{jk} - n \cdot y_{ijk}|}{2n(n+1)}. \qquad (7)$$

If $A_{jk}$ are precalculated the three sums in the approximate estimation scheme (6) are independent, thus the higher-level summation across the input sequence can be

brought to the lowest level. The sum $\sum_{i=1}^{n} \Delta_{ijk}$ can then be computed with complexity lower than $O(n)$. Consider $L_{jk} \subset \{1, 2, \ldots, n\}$ as a subset of indices such that $\forall i \in L_{jk} : n \cdot y_{ijk} < A_{jk}$. Let $B_{jk}$ denote the sum of elements with such indices: $B_{jk} = \sum_{i \in L_{jk}} y_{ijk}$. By performing separate summations across indices in $L_{jk}$ we can remove the absolute value bars in the expression (7):

$$\sum_{i=1}^{n} \Delta_{ijk} = \frac{1}{2n(n+1)} \sum_{i=1}^{n} |A_{jk} - n \cdot y_{ijk}| =$$
$$= \frac{1}{2n(n+1)} \Big( |L_{jk}| \cdot A_{jk} - n \cdot B_{jk} +$$
$$+ n \cdot (A_{jk} - B_{jk}) - A_{jk} \cdot (n - |L_{jk}|) \Big) =$$
$$= \frac{1}{n(n+1)} (A_{jk} \cdot |L_{jk}| - n \cdot B_{jk}). \qquad (8)$$

Thus, to efficiently compute (8) we need to be able to quickly calculate the values of $|L_{jk}|$ and $B_{jk}$ for each $j$ and $k$. In order to do that we need to set up a data structure for $y_{1jk}, y_{2jk}, \ldots, y_{njk}$ which supports fast insertion (which is required when $Y_{n-1}$ is updated to incorporate $X_n$ and produce $Y_n$) and fast computation of the quantity $|L_{jk}|$ and sum $B_{jk}$ of elements lower than the average. For this purpose we can use balanced binary search trees such as treaps [39], with which both insertion and the queries would require $O(\log n)$ operations. The complexity of computing the approximate expected distance estimation (6) is thus reduced to $O(S_n K \log n)$.

For large $n$ this computation scheme is significantly more efficient than the direct summation (6), however for small values of $n$ it could be impractical to implement due to high computational overhead relative to the number of elements in each binary search tree. Thus in the experimental section IV we will evaluate both computational schemes.

### C. Using a normalized GLD

Optimizations proposed in subsections III-A and III-B consider a GLD as the distance function $\rho$. In this subsection we will consider a normalized version of the GLD [37], which was used for measuring a character-level recognition error rate in [34] and [35]. The normalized GLD always has a value in the range $[0.0, 1.0]$ and satisfies triangle inequality. It is defined as follows:

$$\text{nGLD}(X, Y) = \frac{2 \cdot \text{GLD}(X, Y)}{\text{GLD}(X, Y) + \alpha \cdot (|X| + |Y|)}, \qquad (9)$$

where $\alpha$ is the maximal component-wise distance between empty and non-empty characters. With a scaled taxicab metric $\rho_C$ (3) the value of $\alpha$ is 1.

If the approximation of the expected distance $\hat{\Delta}_n$ (2) is calculated directly using the computation scheme (6), the two internal sums in (6) compute the approximation of GLD between $R_n$ and $R(X_1, \ldots, X_n, X_i)$. Thus, each computed GLD can then be normalized under the higher-level sum sign if a normalized version of the GLD is used as a metric function $\rho$.

However, to apply the further optimization (8) for an unweighted case, an additional assumption needs to be allowed to convert the computed summation of GLD values to a sum of normalized GLD values.

Let us denote as $G_n$ the sum of GLD between the current result $R_n$ and the modelled candidates for the next combined result: $G_n = \sum_{i=1}^{n} \text{GLD}(R_n, R(X_1, \ldots, X_n, X_i))$. In order to be able to apply the optimization (8) we will introduce the following approximation:

**Approximation 3:** naive normalization. the GLD constituents of $G_n$ may be normalized after summation. Given approximation 1 (naive alignment) the lengths of $R_n$ and $R(X_1, \ldots, X_n, X_i)$ are both equal to $S_n$, thus approximation 3 can be expressed as:

$$\sum_{i=1}^{n} \text{nGLD}(R_n, R(X_1, \ldots, X_n, X_i)) \approx \frac{2G_n}{G_n + 2S_n}. \quad (10)$$

Since, given the approximation 2 (naive Levenshtein), the value of $\text{GLD}(R_n, R(X_1, \ldots, X_n, X_i))$ is not higher than $S_n$, the approximation 3 in its essence relies on an assumption of local linearity of function $f(t) = (2t)/(t + 2)$ in the range $t \in [0.0, 1.0]$. Using the optimization (8) and balanced binary search trees we can efficiently compute the approximation of $G_n$, and then use the approximate equation (10) to convert it to a sum of normalized GLD between the current result and the modelled candidates for the next one, required to compute the estimation $\hat{\Delta}_n$.

## IV. EXPERIMENTAL EVALUATION

In this section, we provide results of an experimental evaluation of the proposed optimization against the direct application of method evaluated in [35]. The source code and data necessary to reproduce the experiments are available at the following link: https://github.com/SmartEngines/stoppers_modelling.

### A. Experimental setting

The proposed approximations were tested on two different text recognition tasks: the recognition of text fields of identity documents in a video stream, and the recognition of arbitrary text in videos.

For the first task we used open datasets MIDV-500 [22] and MIDV-2019 [40] which contain video sequences of identity documents of various types. The provided ground truth contains coordinates of the document on each frame, for each unique document there are coordinates of each text field and their correct values. The experimental setting closely followed the one described in [35] (which only considered MIDV-500) in order to provide a just comparison of the methods. As in [35] we considered only frames which fully contain the document boundaries, and to avoid normalization effects each clip was repeated in a loop until the common length of 30 was reached. Four text field groups were analyzed: document numbers, numerical dates, names written in the Latin alphabet, and machine-readable zone lines. In total there were 2239 evaluated clips of MIDV-500 dataset and 992 clips

of MIDV-2019 dataset. Each text field image was cropped with the resolution of 300 DPI and with margins equal to 30% of the smallest text field bounding box side, then recognized using a text string recognition subsystem of Smart IDReader document recognition solution [29], thus obtaining string recognition results in form (1).

For evaluating the stopping rules on the task of arbitrary text recognition we used a training subset for the Text in Videos task of the ICDAR 2015 Competition on Robust Reading (IC15-Train) [41] and the YouTube Video Text dataset (YVT) [42]. From both datasets, we selected text objects with only alphanumeric characters and which were present on at least 30 frames, then split the video sequences for each object into clips of exactly 30 frames. Thus, a total of 851 clips were produced from IC15-Train dataset, and 409 clips from YVT dataset. Each text object was then cropped using the coordinates provided in the ground truth in the original resolution and recognized using the pre-trained model described in [15] with thin-plate spline (TPS) transformation, ResNet feature extraction, BiLSTM sequence modelling, and attention-based sequence prediction. The recognition results in form (1) were extracted by removing the softmax outputs corresponding to "start" and "end-of-sentence" tokens and normalization of the remaining character estimation values.

For both tasks the per-frame text recognition results were combined using the algorithm described in [30] without weights, and a normalized GLD [37] was used as a metric $\rho$ between recognition results. The value of $\delta$ parameter in the expected distance estimation (2) was taken to be $0.1$.

For each experiment, three methods were evaluated:

1) *Base method [35].* Estimation value $\hat{\Delta}_n$ is computed using direct modelling (2). The complexity of computing the estimation on stage $n$ is $O(nS_nK(S_n + M))$.
2) *Method A.* Estimation is computed using an approximate modelling using direct summation (6) with $\Delta_{ijk}$ computed using a simplified expression (7) for unweighted case. The complexity of computing the estimation on stage $n$ is $O(nS_nK)$.
3) *Method B.* Estimation is computed using an approximate modelling using a summation scheme optimization (8) with balanced binary search trees. Conversion to normalized GLD is performed using the approximation 3 (10) (naive normalization, see section III-C). Treaps with random priorities [39] were used as balanced binary search trees due to implementation simplicity. The complexity of computing the approximation of $\hat{\Delta}_n$ on stage $n$ is $O(S_nK \log n)$.

### B. Comparing the distance estimation value

For the first experiment we compared the value of the estimated expected distance $\hat{\Delta}_n$ from the current combined result to the potential combination result on the next stage. Fig. 3 illustrates the plot of the mean estimation value on each simulation stage on both datasets for both tasks.

If can be seen from Fig. 3 that the estimation value calculated with Method A equals on average to the estimation
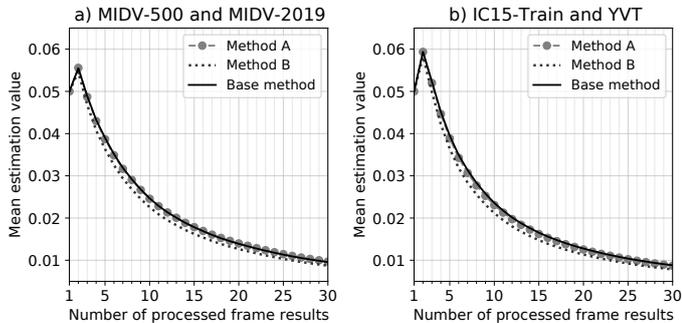
Fig. 3. Mean estimation value $\hat{\Delta}_n$ computed using the three evaluated stopping methods on each stage: (a) fields from MIDV-500 and MIDV-2019, (b) text objects from IC15-Train and YVT
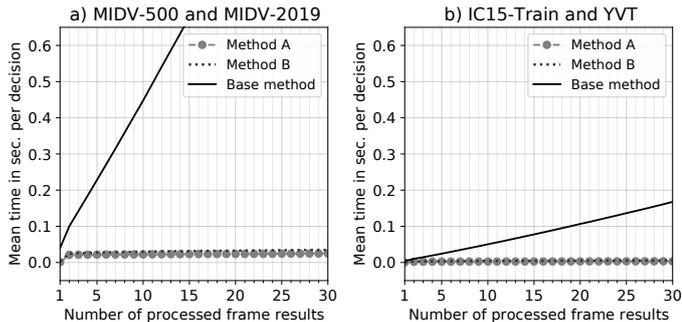


Fig. 4. Comparison of the time required to construct the combined result $R_n$ and calculate distance estimation $\hat{\Delta}_n$, for the evaluated methods: (a) fields from MIDV-500 and MIDV-2019, (b) text objects from IC15-Train and YVT

calculated using the full modelling (2). This signifies that the approximations 1 and 2 (naive alignment and naive Levenshtein) introduced in subsection III-A are justified for efficient approximate computation of the expected distance estimation. The approximation error of Method B is an effect of the naive normalization (10).

### C. Comparing the modelling time

To evaluate the computational performance of the proposed approximations we compared the combined time required to produce the result $R_n$ (as the proposed method involve on-the-fly modifications of some internal structures during the results combination process) and calculate on stage $n$ the estimation of the expected distance from the current combined result $R_n$ to the next result $R_{n+1}$. The combined time to perform such operations for the fields in MIDV-500 and MIDV-2019 are presented in Table I, and the same for the text objects in IC15-Train and YVT is presented in Table II. Fig. 4 shows the timing comparison for each stage $n$ for both tasks. Measurements were performed for a single-thread Python prototype implementation (Python 3.7.4 under Jupyter 6.0.1, AMD Ryzen 9 3950X, 64Gb RAM, GTX 1050 GPU).

It is evident from the data in Fig. 4 that approximate computation of the estimation $\hat{\Delta}_n$ allows to make the stopping decision much quicker in comparison with direct modelling, and that the increase of required computations from stage to stage becomes negligible. As mentioned in section III-B, the

TABLE I
TIME REQUIRED TO CONSTRUCT THE COMBINED RESULT AND CALCULATE THE EXPECTED DISTANCE ESTIMATION: FIELDS FROM MIDV-500 AND MIDV-2019

| Method | Time on stage $n$ (in seconds) | | | | |
|---|---|---|---|---|---|
| | $n = 5$ | $n = 10$ | $n = 15$ | $n = 20$ | $n = 25$ |
| Base [35] | 0.228 | 0.448 | 0.678 | 0.906 | 1.143 |
| Method A | **0.022** | **0.022** | **0.023** | **0.024** | **0.024** |
| Method B | 0.027 | 0.030 | 0.032 | 0.033 | 0.034 |

TABLE II
TIME REQUIRED TO CONSTRUCT THE COMBINED RESULT AND CALCULATE THE EXPECTED DISTANCE ESTIMATION: TEXT OBJECTS FROM IC15-TRAIN AND YVT

| Method | Time on stage $n$ (in seconds) | | | | |
|---|---|---|---|---|---|
| | $n = 5$ | $n = 10$ | $n = 15$ | $n = 20$ | $n = 25$ |
| Base [35] | 0.024 | 0.050 | 0.078 | 0.107 | 0.136 |
| Method A | **0.002** | **0.003** | **0.003** | **0.003** | **0.003** |
| Method B | 0.004 | 0.004 | 0.005 | 0.005 | 0.005 |

difference between the direct summation in Method A and the optimized summation scheme (8) with balanced search trees in Method B is insignificant, and the latter is even slightly less computationally efficient, due to the data structure overhead.

### D. Comparing the stopping method performance

The final experiment was aimed at comparing the performance of the resulting modified stopping methods. A "good" stopping method should be able to achieve lower error level given a fixed mean number of frames, or, respectively, the lower mean number of processed frames given a fixed mean error level. A convenient method of comparing stopping methods is to compare their expected performance profiles, which are used for anytime algorithms analysis [43]. Such performance profile can be obtained by plotting the mean error level (in terms of distance to the correct result) of the combined result at the stopping stage and the mean number of the stopping stage, while varying observation cost value. The lower position of the plotted curve indicates the greater performance of the stopping method. The observation cost in the analyzed methods corresponds to the threshold with which the estimation $\hat{\Delta}_n$ is compared when making the stopping decision.

Fig. 5 illustrates the expected performance profiles of the three evaluated methods, as well as a baseline stopping method which stops the process after a fixed stage, for text fields in the MIDV-500 and MIDV-2019 datasets. The varied parameter for the baseline stopping method is the number of the stage on which it should stop. The corresponding profiles for the arbitrary text recognition task are presented in Fig. 6.

It can be seen from Fig. 5 and 6 that the approximations introduced in sections III-A and III-C had almost no effect
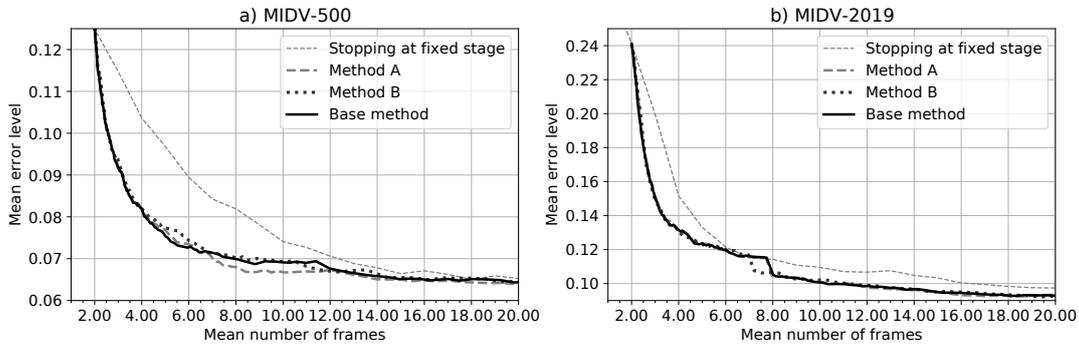
Fig. 5. Expected performance profiles of the three evaluated methods: mean distance from the obtained combination result to the correct value against the mean number of processed frames before stopping, with a varied stopping threshold: (a) fields from MIDV-500 dataset, (b) fields from MIDV-2019 dataset
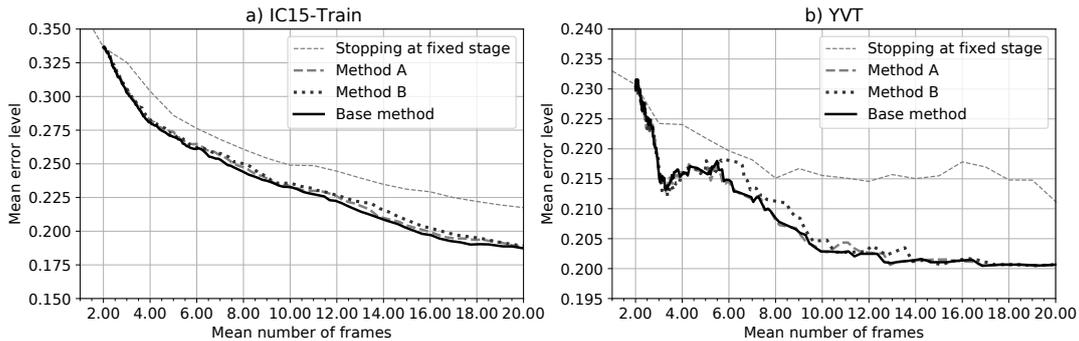


Fig. 6. Expected performance profiles of the three evaluated methods: mean distance from the obtained combination result to the correct value against the mean number of processed frames before stopping, with a varied stopping threshold: (a) text objects from IC15-Train, (b) text objects from YVT dataset

on the performance of the stopping method in terms of the achieved error level and required number of observations. There is a slight disadvantage of the Method B against the Method A, which could be an effect of utilizing the naive normalization approximation 3 (10) when using a normalized GLD as a metric function $\rho$. Given the significant advantage in computational performance, it can be concluded that the proposed approximate modelling method is more favourable for stopping the text recognition process in a video stream, applicable for both document analysis, and for recognition of an arbitrary text.

## V. CONCLUSION

In this paper, we described the problem of making the stopping decision efficiently for text string recognition in a video stream. An overview was given for an existing method based on modelling of the next combined result. The disadvantage of this method is its high computational complexity of the modelling, resulting in a slow decision making process. In the paper, we proposed two approximate modelling schemes, one which is applicable for a general weighted case, and one for unweighted case. Both optimizations were evaluated on two open datasets and in the same experimental setting as the previously introduced method.

It can be seen from the presented experimental evaluation that the assumptions and approximations introduced in the paper had almost no effect on the performance of the

stopping method in terms of the achieved mean error level and a mean number of consumed observations. At the same time, the proposed computational schemes have significantly higher computational performance. The obtained results were consistent for two different text recognition tasks: identity documents recognition, and arbitrary text recognition, each with a different text recognition algorithm.

Optimal stopping of object recognition in a video stream is an important pattern recognition problem. The role of stopping methods as components of modern video recognition systems is very valuable, not only from the perspective of the real-time interaction with the user but from the more general perspective of being able to achieve more accurate recognition results with adequate response time.

### REFERENCES

[1] Y. Chen, J. Pont-Tuset, A. Montes, and L. V. Gool, "Blazingly fast video object segmentation with pixel-wise metric learning," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1189–1198.

[2] H. Xiao, J. Feng, G. Lin, Y. Liu, and M. Zhang, "Monet: Deep motion exploitation for video object segmentation," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1140–1148.

[3] R. Girdhar, G. Gkioxari, L. Torresani, M. Paluri, and D. Tran, "Detect-and-track: Efficient pose estimation in videos," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 350–359.

[4] C. Lin and Y. Hung, "A prior-less method for multi-face tracking in unconstrained videos," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 538–547.

[5] Y. Jo, S. W. Oh, J. Kang, and S. J. Kim, "Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3224–3232.

[6] Z. Cheng, J. Lu, J. Xie, Y. Niu, S. Pu, and F. Wu, "Efficient video scene text spotting: Unifying detection, tracking, and recognition," *arXiv preprint:1903.03299*, 2019.

[7] B. A. Dangiwa and S. S. Kumar, "A business card reader application for iOS devices based on Tesseract," in *2018 International Conference on Signal Processing and Information Security (ICSPIS)*, 2018, pp. 1–4.

[8] N. Islam, Z. Islam, and N. Noor, "A survey on optical character recognition system," *arXiv preprint:1710.05703*, 2017.

[9] K. Ravneet, "Text recognition applications for mobile devices," *Journal of Global Research in Computer Science*, vol. 9, no. 4, pp. 20–24, 2018.

[10] H. Jabnoun, F. Benzarti, and H. Amiri, "A new method for text detection and recognition in indoor scene for assisting blind people," in *Proc. SPIE (ICMV 2016)*, vol. 10341, 2017.

[11] M. Ki, M. Hradi, and O. Kodym, "Brno Mobile OCR Dataset," in *International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2019, pp. 1352–1357.

[12] M. O. V. Ngoc, J. Fabrizio, and T. Graud, "Document detection in videos captured by smartphones using a saliency-based method," in *International Conference on Document Analysis and Recognition Workshops (ICDARW)*, 2019, pp. 19–24.

[13] S. Mohanty, T. Dutta, and H. P. Gupta, "An efficient system for hazy scene text detection using a deep cnn and patch-nms," in *2018 24th International Conference on Pattern Recognition (ICPR)*, 2018, pp. 2588–2593, doi:10.1109/ICPR.2018.8545198.

[14] W. Sui, Q. Zhang, J. Yang, and W. Chu, "A novel integrated framework for learning both text detection and recognition," in *2018 24th International Conference on Pattern Recognition (ICPR)*, 2018, pp. 2233–2238, doi:10.1109/ICPR.2018.8545047.

[15] J. Baek, G. Kim, J. Lee, S. Park, D. Han, S. Yun, S. J. Oh, and H. Lee, "What is wrong with scene text recognition model comparisons? dataset and model analysis," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 4714–4722, doi:10.1109/ICCV.2019.00481.

[16] Z. Cheng, J. Lu, Y. Niu, S. Pu, F. Wu, and S. Zhou, "You only recognize once: Towards fast video text spotting," in *Proceedings of the 27th ACM International Conference on Multimedia*, 2019, pp. 855–863, doi:10.1145/3343031.3351093.

[17] Y. Cai and W. Wang, "Robustly detect different types of text in videos," *Neural Computing and Applications*, 2020, doi:10.1007/s00521-020-04729-6.

[18] J. Greenhalgh and M. Mirmehdi, "Recognizing Text-Based Traffic Signs," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 3, pp. 1360–1369, 2015, doi:10.1109/TITS.2014.2363167.

[19] F. S. Bashiri, E. LaRose, J. C. Badger, R. M. D'Souza, Z. Yu, and P. Peissig, "Object detection to assist visually impaired people: A deep neural network adventure," in *Advances in Visual Computing*, 2018, pp. 500–510.

[20] M. Cutter and R. Manduchi, "Towards Mobile OCR: How to take a good picture of a document without sight," in *Proceedings of the ACM Symposium on Document Engineering*, 2015, pp. 75–84.

[21] E. Tekin, J. M. Coughlan, and H. Shen, "Real-time detection and reading of led/lcd displays for visually impaired persons," in *IEEE Workshop on Applications of Computer Vision (WACV)*, 2011, pp. 491–496.

[22] V. Arlazarov, K. Bulatov, T. Chernov, and V. Arlazarov, "MIDV-500: a dataset for identity document analysis and recognition on mobile devices in video stream," *Computer Optics*, vol. 43, pp. 818–824, October 2019, doi:10.18287/2412-6179-2019-43-5-818-824.

[23] M. O. V. Ngoc, J. Fabrizio, and T. Graud, "Saliency-based detection of identy documents captured by smartphones," in *13th IAPR International Workshop on Document Analysis Systems (DAS)*, 2018, pp. 387–392.

[24] N. Skoryukina, V. V. Arlazarov, and D. P. Nikolaev, "Fast method of ID documents location and type identification for mobile and server application," in *International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2019, pp. 850–857.

[25] M. A. Povolotskiy, D. V. Tropin, T. S. Chernov, and B. I. Savelyev, "Dynamic programming approach to textual structured objects segmentation in images," *Informatsionnye tekhnologii i vychislitelnye sistemy (Information technologies and computational systems)*, vol. 3, pp. 66–78, 2019, (In Russian).

[26] S. Bakkali, Z. Ming, M. M. Luqman, and J.-C. Burie, "Face detection in camera captured images of identity documents under challenging conditions," in *International Conference on Document Analysis and Recognition Workshops (ICDARW)*, 2019, pp. 55–60.

[27] K. Bulatov, V. V. Arlazarov, T. Chernov, O. Slavin, and D. Nikolaev, "Smart IDReader: Document recognition in video stream," in *14th International Conference on Document Analysis and Recognition (ICDAR)*, vol. 6. IEEE, 2017, pp. 39–44.

[28] M. Ryan and N. Hanafiah, "An examination of character recognition on ID card using template matching approach," *Procedia Computer Science*, vol. 59, pp. 520–529, 2015.

[29] Y. S. Chernyshova, A. V. Sheshkus, and V. V. Arlazarov, "Two-Step CNN Framework for Text Line Recognition in Camera-Captured Images," *IEEE Access*, vol. 8, pp. 32 587–32 600, 2020, doi:10.1109/ACCESS.2020.2974051.

[30] K. Bulatov, "A method to reduce errors of string recognition based on combination of several recognition results with per-character alternatives," *Bulletin of the South Ural State University. Ser. Mathematical Modelling, Programming & Computer Software*, vol. 12, no. 3, pp. 74–88, 2019.

[31] S. Christensen and A. Irle, "The monotone case approach for the solution of certain multidimensional optimal stopping problems," *arXiv preprint:1705.01763*, 2019.

[32] T. Ferguson and M. Klass, "House-hunting without second moments," *Sequential Analysis*, vol. 29, no. 3, pp. 236–244, 2010.

[33] V. V. Arlazarov, K. Bulatov, T. Manzhikov, O. Slavin, and I. Janiszewski, "Method of determining the necessary number of observations for video stream documents recognition," in *Proc. SPIE (ICMV 2017)*, vol. 10696, 2018.

[34] K. Bulatov, N. Razumnyi, and V. V. Arlazarov, "On optimal stopping strategies for text recognition in a video stream as an application of a monotone sequential decision model," *International Journal on Document Analysis and Recognition*, vol. 22, no. 3, pp. 303–314, 2019.

[35] K. Bulatov, B. Savelyev, and V. V. Arlazarov, "Next integrated result modelling for stopping the text field recognition process in a video using a result model with per-character alternatives," in *Proc. SPIE (ICMV 2019)*, vol. 11433, January 2020, pp. 710–716, doi:10.1117/12.2559447.

[36] R. Llobet, J. Cerdan-Navarro, J. Perez-Cortes, and J. Arlandis, "OCR post-processing using weighted finite-state transducers," in *20th International Conference on Pattern Recognition*, 2010, pp. 2021–2024.

[37] L. Yujian and L. Bo, "A normalized levenshtein distance metric," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 6, pp. 1091–1095, 2007.

[38] J. G. Fiscus, "A post-processing system to yield reduced word error rates: Recognizer Output Voting Error Reduction (ROVER)," in *IEEE Workshop on Automatic Speech Recognition and Understanding*, 1997, pp. 347–354.

[39] G. E. Blelloch and M. Reid-Miller, "Fast set operations using treaps," in *Proceedings of the Tenth Annual ACM Symposium on Parallel Algorithms and Architectures*, ser. SPAA '98. ACM, 1998, pp. 16–26. [Online]. Available: http://doi.acm.org/10.1145/277651.277660

[40] K. Bulatov, D. Matalov, and V. Arlazarov, "MIDV-2019: challenges of the modern mobile-based document OCR," in *Twelfth International Conference on Machine Vision (ICMV 2019)*, vol. 11433. SPIE, January 2020, pp. 717–722, doi:10.1117/12.2558438.

[41] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu, F. Shafait, S. Uchida, and E. Valveny, "Icdar 2015 competition on robust reading," in *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, 2015, pp. 1156–1160, doi:10.1109/ICDAR.2015.7333942.

[42] Phuc Xuan Nguyen, K. Wang, and S. Belongie, "Video text detection and recognition: Dataset and benchmark," in *IEEE Winter Conference on Applications of Computer Vision*, 2014, pp. 776–783, doi:10.1109/WACV.2014.6836024.

[43] S. Zilberstein, "Using anytime algorithms in intelligent systems," *AI Magazine*, vol. 17, no. 3, pp. 73–83, 1996.