**SN**

**ORIGINAL RESEARCH**

# Optimized Data Structuring and Preprocessing Techniques for Belfort Civil Registers of Birth Transcription

**Wissam AlKendi[1]** [ID] · **Franck Gechter[1,4]** · **Laurent Heyberger[2]** · **Christophe Guyeux[3]**

## Abstract

Historical documents offer invaluable insights into the past, shaping our understanding of the world's rich tapestry of stories. This paper presents methodologies designed to aid in the transcription of French Belfort birth registers from 1807 to 1919. The approach emphasizes preprocessing techniques, including binarization, skew correction, and text line segmentation, to effectively tackle the challenges arising from varied text styles, marginal annotations, and the combination of printed and handwritten text. We present this archive as a newly developed database, employing a structured methodology with XML tags to ensure accurate formatting and alignment of transcriptions with image elements at both the paragraph and text line levels. Our preprocessing phase demonstrates an accuracy rate of 95.9%, highlighting the effectiveness of our techniques in preserving and facilitating the study of this rich cultural heritage. This work contributes substantially to the field of handwritten text recognition and sets the stage for further advancements in the automated transcription of historical documents. This paper is an extended version of our previous work presented at IMPROVE2024, including additional experimental results and in-depth analysis.

Wissam AlKendi, Franck Gechter, Laurent Heyberger and Christophe Guyeux these author contributed to this work.

✉ Wissam AlKendi
wissam.al-kendi@utbm.fr

Franck Gechter
franck.gechter@utbm.fr

Laurent Heyberger
laurent.heyberger@utbm.fr

Christophe Guyeux
christophe.guyeux@univ-fcomte.fr

1 CIAD (UMR 7533), Université Marie et Louis Pasteur, UTBM, Belfort, France

2 FEMTO-ST Institute/RECITS (UMR 6174 CNRS), Université Marie et Louis Pasteur, UTBM, Belfort, France

3 FEMTO-ST Institute/DISC (UMR 6174 CNRS), Université Marie et Louis Pasteur, UFC, Belfort, France

4 LORIA (UMR 7503), SIMBIOT Team, Vandoeuvre-lès-Nancy F-54506, France

## Introduction

Handwritten text recognition (HTR) has become an essential element in the field of pattern recognition. Numerous researchers have devised methods to transcribe different kinds of documents, such as historical archives [1], books, letters, and general forms, employing either spatial (offline) [2] or temporal (online) [3] techniques.

Transcription involves automatically converting handwritten text in a digital image into a machine-readable format, making the text more accessible. This process includes several steps: preprocessing methods to prepare the image for analysis (e.g., noise reduction and normalization), deep learning techniques to recognize the handwriting, and postprocessing methods to fine-tune the output, correct errors, and enhance readability [4]. Figure 1 depicts the primary preprocessing steps necessary to prepare the image for analysis prior to implementing deep learning methods for text recognition.

Currently, systems possess the ability to analyze document layouts [5] and recognize text at various levels, such as characters, text lines, paragraphs, and complete documents.
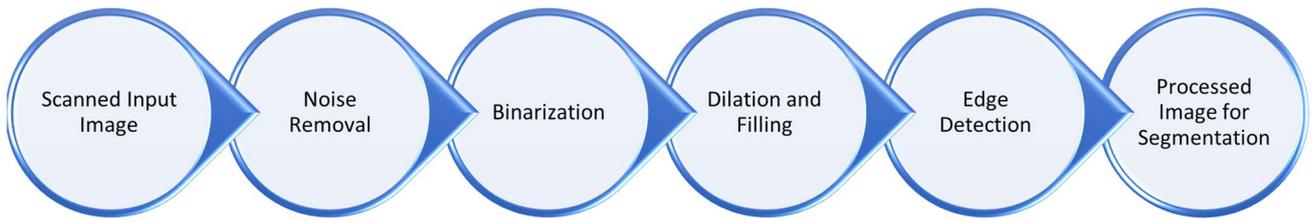
**Fig. 1** Preprocessing components: the essential steps involved in the transcription of handwritten text images

Moreover, these systems can distinguish between different handwriting styles and are applicable to a wide range of languages [6].

Despite significant advancements in recognizing modern handwritten text, many challenges still need to be addressed in transcribing historical documents. These documents frequently display distinctive features, including angular and spiky letters, elaborate flourishes, and overlapping words and text lines, all presented in a variety of handwriting styles [7]. Moreover, the quality of reproduction in these documents considerably extends the time needed for the preprocessing phase [8].
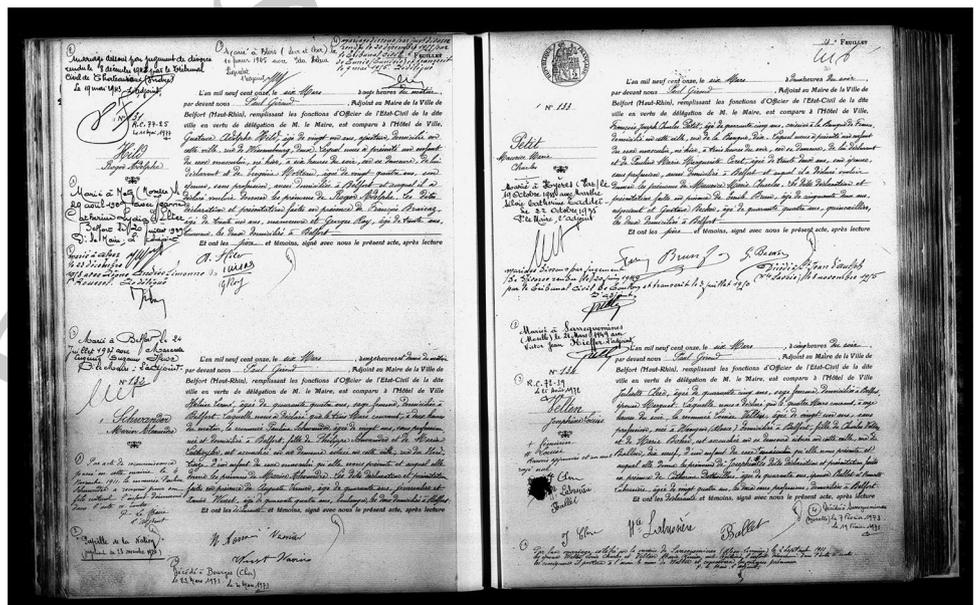
The historical significance of Belfort's records arises from the city's unique demographic history in the nineteenth century. After France's defeat by Prussia in 1871 and the subsequent annexation of Alsace-Lorraine by Germany, Belfort experienced a rapid population increase. This growth was primarily due to the influx of Alsatians, especially those from Mulhouse, who chose to stay in French territory or moved to Belfort to work in the burgeoning textile and mechanical industries established after 1871 [9]. Consequently, this metropolis functions as a distinctive

observatory for three pivotal transformations: urban growth, migration, and the evolving dynamics of social relations concerning sexuality.

The examination of birth records offers significant historical insights into various aspects, including the roles and practices of midwives, shifts in cohabitation patterns and the ages of parents at the birth of their first child, occurrences of out-of-wedlock births, trends in child recognition, and the selection of witnesses for official declarations. Additionally, these records offer insights into child naming conventions and the prevalence of home versus hospital births. Each of these aspects provides a distinct perspective on the societal and cultural dynamics of the time. Figure 2 provides a visual representation of a sample page from these civil birth registers.

To effectively investigate and address these historical concerns, it is crucial to create a comprehensive knowledge database. This process involves transcribing archival content through two primary methods. The first method is manual transcription, which, although accurate, is costly and time-consuming. The second method is automatic transcription, which employs deep learning techniques.

**Fig. 2** Sample of Belfort civil registers of births (Source: [10])

This paper tackles the challenges faced in the preprocessing phase of transcribing the Belfort birth records. The focus is on segmenting the documents into paragraph and text line levels, and creating structured data files to make these historical records more accessible and better preserved. Recognizing handwritten text in these documents is particularly difficult, especially with the mix of handwritten and printed text. Even though Optical Character Recognition (OCR) techniques have advanced significantly in recent years, they still fall short when it comes to these specific documents.

This paper is an extended version of our previous work presented at The 4th International Conference on Image Processing and Vision Engineering (IMPROVE2024) [10]. In this extended version, we provide additional experimental results and analyses. Specifically, we have conducted more experiments on the dataset, performed experiments on other benchmark datasets using additional evaluation metrics, and included a parameters sensitivity analysis study to determine the most suitable parameters and threshold values. Additionally, we compare the accuracy of our segmentation process with that of online handwritten text recognition commercial systems. These enhancements significantly contribute to the robustness and effectiveness of our methodology.

The rest of the paper is structured as follows: Sect. 2 reviews recent advancements in the field. In Sect. 3, we outline the characteristics and challenges associated with the Belfort civil birth registers. Section 4 presents a summary of the experimental findings. Lastly, Sect. 5 concludes the paper with recommendations for future research directions.

## State-of-the-Art Overview

Numerous studies have been conducted to enhance the field of image processing and address the challenges encountered during the preprocessing phase. This includes methods such as binarization/thresholding, noise removal, and contrast enhancement. These methods are widely used in various significant applications, as utilized in [11, 12].

In [1], authors examined the initial phase of historical document transcription, covering steps like binarization, skew correction, and segmentation. Their findings highlighted the crucial role of accurate transcription in ensuring effective information retrieval from archival texts. Furthermore, authors of [13] carried out an in-depth review of diverse document layout analysis (DLA) techniques designed to detect and annotate the physical structure of documents. Their study covered multiple stages of DLA algorithms, including preprocessing, layout analysis methods, postprocessing, and performance assessment.

The main steps of the preprocessing phase are outlined as follows.

## Noise Removal

This step entails removing extraneous pixels from the digital image that might obscure the original information. Such noise can originate from the image sensor and electronic components of a scanner or digital camera.

Authors of [14] introduced a hybrid binarization approach aimed at improving the quality of historical and ancient documents. This approach combines global and local thresholding techniques while maintaining low computational and temporal costs. The study also emphasizes the importance of identifying noise patterns to choose the most effective removal methods. Additionally, [15] conducted a survey focused on removing marginal noise from historical handwritten document images, such as those in the Australian archives, which feature unique layout complexities. This survey aims to guide researchers in selecting the most effective methods based on the specific type of marginal noise in their datasets.

Many effective algorithms have been designed to denoise images by improving sparse representation in the transform domain and clustering similar 2-D image patches into 3-D data arrays. A prominent example is the Block-Matching and 3D Filtering (BM3D) algorithm, which can be applied before the binarization step to substantially improve the results of the preprocessing phase [16, 17].

## Binarization

This step involves converting digital images into binary form, comprising two sets of pixels: black and white (0 and 1) [18]. Binarization is essential for distinguishing the foreground text from the background.

Authors of [19] concentrated on the binarization of historical document images, creating a standard benchmark that significantly advanced research in this field. The study explores a range of approaches and techniques, including statistical models and pixel classification with learning algorithms. Furthermore, the study addresses evaluation metrics and datasets, offering a thorough examination of the topic.

## Edges Detection

This step entails identifying the edges or boundaries of the text within an image by connecting continuous points that share the same color or intensity. Common techniques used for this purpose include Laplacian, Sobel, Canny, and Prewitt edge detection [20].

## Skew Detection and Correction

Skew refers to the misalignment of text within a digital image, representing the degree of rotation required to

align the text properly along the horizontal or vertical axis. Numerous methods for skew detection and correction have been developed. These approaches are detailed in [21], where authors discuss the detection of skew angles in text images. These include Projection Profile-based methods (PP), Nearest Neighbor clustering (NN), Hough Transform, and Cross Correlation.

## Segmentation

This stage involves segmenting the handwritten text image into individual letters, words, lines, and paragraphs, typically using the pixel properties of the image. Various methodologies have been applied to this process, including threshold methods, edge-based methods, region-based methods, watershed-based methods, and clustering methods. Effective segmentation is crucial as it significantly enhances the accuracy of Handwritten Text Recognition (HTR) models.

Authors of [22] explore a variety of classical and learning-based line segmentation techniques for Handwritten Text Recognition (HTR) in historical manuscripts. Their study highlights that connected components and graph-based techniques are particularly effective in overcoming challenges such as overlapping lines and touching characters in historical manuscripts. Additionally, authors of [23] utilized a binarization algorithm based on Gaussian filtering to tackle the challenge of non-uniform luminance in selected pages from the Latin Theologica Miscellanea documents. They also applied a Hough Transform Mapping technique to effectively handle the challenges posed by handwritten manuscripts, such as overlapping characters.

Furthermore, in [24], authors developed a modular analytic pipeline that leverages advanced image processing and machine learning algorithms to automate data extraction from French vital records. This pipeline includes document preprocessing and text segmentation stages, employing histograms and the EAST method to detect text boxes. The main objective is to simplify the manual annotation process for the training dataset. The study concluded that the EAST approach is highly effective in detecting text boxes within these documents.

Conversely, authors of [25] employed an Energy Map for both binary and grayscale images to minimize non-text areas by extracting connected components along text lines. They applied a seam carving technique to identify seams (paths of least energy) in the energy map, facilitating text line identification. This study also presented a new benchmark dataset consisting of historical documents in multiple languages. Similarly, [26–28] proposed important methodologies for text line segmentation. However, [15] point out that there is a limited amount of research on page segmentation for historical handwritten documents compared to non-historical documents, indicating a notable research gap in this area.

## Methodology

### Belfort Civil Registers of Births

The birth records in the Belfort commune's civil registries comprise 39,627 entries written in French, each digitized at 300 dpi. These records were chosen for their uniformity, spanning Gregorian birth dates from 1807 to 1919, and adhering to legal constraints. Originally, these registers consisted of handwritten entries but later transitioned to a partially printed version with specific areas left blank for recording precise information about the newborn.

In Belfort, this transition occurred in 1885, impacting about 57.5% of the 39,627 recorded declarations. The archive is accessible online up to 1902 at https://archives.belfort.fr/search/form/e5a0c07e-9607-42b0-9772-f19d7bfa180e (accessed on July 05, 2024). Additionally, we have received authorization from the municipal archives to study data extending to 1919.

### Records Structure

These records include essential details such as the child's name, parents' names, witnesses, and other pertinent information. Figure 3 illustrates a sample record from the civil registers. Table 1 presents a summary of the structure and content of a record within the archive.

### Automatic Transcription Challenges

The transcription of the Belfort archive records poses several challenges, as outlined below.

### Document Layout

The Belfort birth documents exhibit two distinct layouts. The first layout comprises double pages, with each page containing a single complete record. The second layout also consists of double pages but accommodates two complete records per page. However, some pages may have entries that begin on one page and continue onto the next page.

### Reading Order

Understanding the correct order for reading text areas, including primary paragraphs and marginal annotations, is essential. The reading sequence should be: first the record number, then the name, followed by the primary paragraph,

**Fig. 3** A record from the Belfort civil registers of birth with annotations indicating different sections. The primary paragraph containing detailed information, is marked in blue. Marginal annotations are marked in yellow. The record number, which uniquely identifies the record, is marked in green. The record name is marked in red. Both the record number and name are part of the record's header margin (Source: [10])
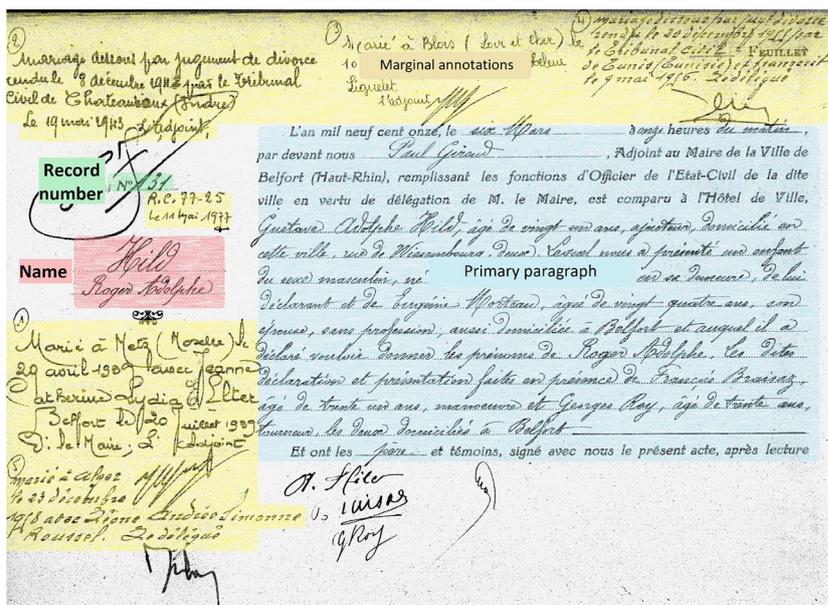


**Table 1** The structure and contents of a record in the Belfort civil registers of births as presented in [29].

| Structure | Content |
| --- | --- |
| Head margin | Registration number |
| | First and last name of the person born |
| Primary paragraph | Time and date of declaration |
| | Surname, first name and position of the official registering |
| | Surname, first name, age, profession and address of declarant |
| | Sex of the newborn |
| | Time and date of birth |
| | First and last name of the father (if different of the declarant) |
| | Surname, first name, status (married or other), profession (sometimes) and address (sometimes) of the mother. |
| | Surnames of the newborn |
| | Surnames, first names, ages, professions and addresses (city) of the 2 witnesses. |
| | Mention of absence of signature or illiteracy of the declarant (very rarely). |
| Marginal annotations | Mention of official recognition of paternity/maternity (by father or/and mother): surname, name of the declarant, date of recognition (by marriage or declaration). |
| | Mention of marriage: date of marriage, wedding location, surname and name of spouse. |
| | Mention of divorce: date of divorce, divorce location |
| | Mention of death: date and place of death, date of the declaration of death. |

and lastly, any marginal annotations. This order should be followed from left to right and top to bottom.

## Hybrid Format

Some registers feature entries with a mix of printed and handwritten texts, as illustrated in Fig. 2. The handwritten sections can greatly differ in style, legibility, and ink density, necessitating a smooth transition between OCR and handwriting recognition modes during transcription.

Additionally, the spatial arrangement of printed and handwritten text can be intricate, with handwritten annotations often appearing in the margins, between printed lines, or even overlapping the printed content.

## Marginal Annotations

These annotations, which provide additional information about the individual's birth, are typically added afterward. They often appear in different handwriting styles and are

written with various instruments, positioned differently from the primary paragraph of the declaration.

### Text Styles

The registers display various handwriting techniques, including angular and spiky letters, different character sizes, and intricate flourishes. Consequently, this diversity often leads to overlapping words and text lines within the script.

### Skewness

Many text lines in the primary paragraphs and marginal annotations display skew anomalies, including vertical text rotated 90 degrees. Efficient methods are necessary to adjust for image skewness, regardless of the rotation angle.

### Deterioration

The images exhibit text deterioration due to fading ink and page smearing, as well as ink spots and yellowing of the pages.

### Preprocessing Methods

Various methods and filters have been utilized to process the images, aimed at reducing noise and improving text visibility to facilitate effective edge detection and automated text extraction. The proposed methodology involves segmenting text images into two levels: text paragraphs and lines. This is achieved through a manual document layout analysis approach to identify key components of the records such as record number, name, primary paragraphs, and marginal annotations. The process includes determining the coordinates of these components using a custom interface tool, while the coordinates of text lines within them are automatically identified. The essential steps of the preprocessing phase are detailed below. Figure 4 depicts the complete pipeline of methods applied to the images.

- *Grayscale conversion*: the original images are converted into grayscale, streamlining further processing and reducing computational demands.
- *Gaussian blur*: grayscale images undergo a Gaussian blur (GB) to reduce noise and enhance feature recognition. Equation 1 illustrates the Gaussian blur operation applied to the grayscale images.

$$GB(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp\left( -\frac{x^2}{2\sigma_x^2} - \frac{y^2}{2\sigma_y^2} \right) * gray(x, y).$$

(1)

where $\sigma_x$ and $\sigma_y$ denote the kernel sizes, and $\sigma$ represents the standard deviation.

- *Adaptive Thresholding*: this technique is employed to produce a binary image that emphasizes edges and relevant features. Equation 2 illustrates the adaptive thresholding method applied, resulting in binary images. The Gaussian approach utilizes a fixed block size to compute the adaptive threshold value $T(x, y)$ based on local image characteristics, ensuring distinct separation between lines if overlapping occurs.

$$thresh(x, y) = \begin{cases} 255, & \text{if } GB(x, y) \leq T(x, y) \\ 0, & \text{if } GB(x, y) > T(x, y) \end{cases}$$

(2)

$$T(x, y) = mean(x, y) - 2 \times stddev(x, y).$$

where $mean(x, y)$ denotes the mean pixel value of the local neighborhood centered at $(x, y)$, and $stddev(x, y)$ represents the standard deviation of pixel values within that same local neighborhood.

- *Morphological operations*: dilation operations refine the detected edges by expanding the boundaries of the text, enhancing the brightness of bright regions and darkness of dark regions. These operations manipulate and improve the edge structures observed in the previous stage, thereby enhancing their significance and clarity.
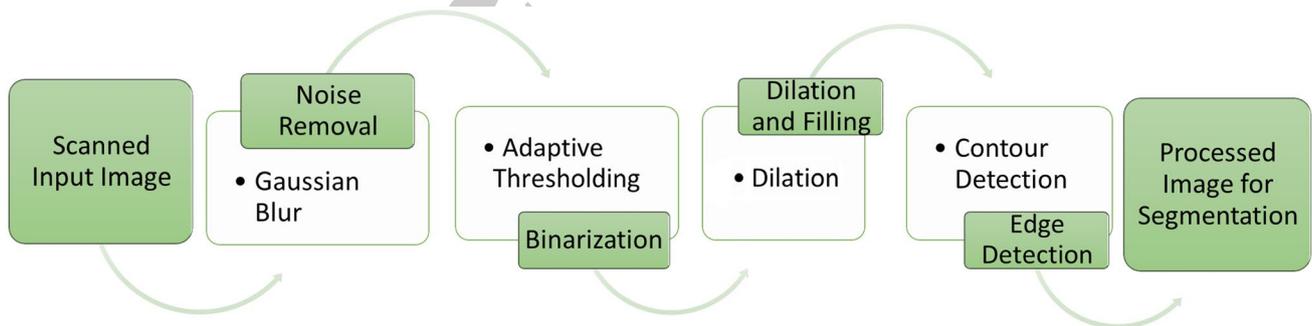


**Fig. 4** Preprocessing components: key steps in preparing Belfort civil registers of births for transcription (Source: [10])

Eq. 3 illustrates the dilation operation applied to images using a horizontal kernel of size 1x205. This process aids in locating the central axis of the text line by identifying the maximum value among neighboring pixels in the horizontal direction.

$$\text{Image dilation}(x, y) = \max_{i=0}^{204}(\text{thresh}(x, y + i)). \quad (3)$$

where $(x, y)$ denotes the position of the pixel within the images, and $i$ represents the horizontal offset within the kernel.

- *Contour detection*: this technique identifies and extracts contours from the filtered images, specifically isolating external outlines of text lines. Parameters are adjusted to simplify contours by excluding redundant points, thereby optimizing memory usage. This process yields a list of contours along with a hierarchical representation of these contours. Subsequently, these contours are scrutinized and filtered based on thresholds that indicate the width and height of the text lines within paragraphs. Table 2 presents the key parameters used during the preprocessing stage, such as those applied in blurring, thresholding, dilation, and line padding.

- *Skew correction*: the skew correction process for text line images involves several steps, including the application of filters, the Canny edge detection algorithm, and the Hough Line Transform. First, we used the NumPy function "np.arctan2" [30] to calculate the slopes of the detected text lines. These slopes provide information about the orientation of the lines, helping to identify the main horizontal line. Next, we determined the median angle from these slopes and converted it into degrees.

Secondly, we located the center of the text line image and computed a rotation matrix using "cv2.getRotationMatrix2D" [31], illustrated in Eq. 6. This matrix is essential for the subsequent rotation operation. Additionally, we determined the most frequent color in the border region of the text line image to fill empty spaces in the background post skew correction. This involved counting occurrences of each unique pixel value and selecting the

color with the highest count to accurately identify the predominant color.

Finally, this process utilizes an affine transformation on the text line image with the calculated matrix $M$. The interpolation method specified by the flags (cv2.INTER_CUBIC) ensures cubic interpolation, while cv2.BORDER_CONSTANT maintains a constant border value for smooth rotation around the image center. This approach enhances the accuracy of text line image representation post skew correction.

$$\text{Angle} = \text{median}(\arctan 2(\Delta y, \Delta x)) \times \frac{180}{\pi} \quad (4)$$

where $\Delta x$ and $\Delta y$ represent the differences in the x and y coordinates of the line endpoints, respectively.

$$\text{Center} = \left( \frac{N_x}{2}, \frac{N_y}{2} \right) \quad (5)$$

where $N_x$ and $N_y$ represent the dimensions of the image.

$$M = \text{cv2.getRotationMatrix2D(Center, Angle, 1)}. \quad (6)$$

where 1 denotes that no scaling is applied during the transformation process. Figure 5 demonstrates an example of the skew correction step applied to a text line image. Table 7 presents the parameters utilized in this process.

The parameter values and thresholds are dynamically adjusted based on the characteristics of the input images. Figure 6 provides visual examples illustrating the preprocessing phase.

## Structured Data Generation

Based on our prior study [29], there is a pressing demand for a novel deep learning approach tailored for transcribing the French Belfort civil registers of births, considering the manifold challenges involved. This phase encompasses manually transcribing a segment of the archival pages to establish a training dataset for the deep learning model. Additionally, it involves devising a precise methodology to correlate identified image regions with the transcribed text, thereby enhancing the labeling process for future attempts.

### Manual Transcription

Various tags were used throughout the manual transcription process to ensure accurate structuring and identification of the components within each record.

Table 3 displays the different types of tags employed in the manual transcription process. Currently, 319.txt files have been successfully transcribed, encompassing 1,010

**Table 2** Key Parameters and Their Values Utilized in Image Preprocessing and Text Line Segmentation Steps

| Parameter | Value |
|---|---|
| Gaussian Blur - Kernel Size | (101, 51) |
| Gaussian Blur - Standard Deviation | 61 |
| Adaptive Thresholding - Block Size | 71 |
| Adaptive Thresholding - Constant (C) | 2 |
| Dilation - Kernel Size | (1, 205) |
| Line Padding Above/Below | 5–10 pixels |

**Fig. 5** Examples of the skew correction process applied to a text line image from the Belfort civil birth registers: **a** Original text line image. **b** Dominant border color calculated to fill the rotated areas. **c** Image rotated to correct the skew. **d** Final corrected text line image
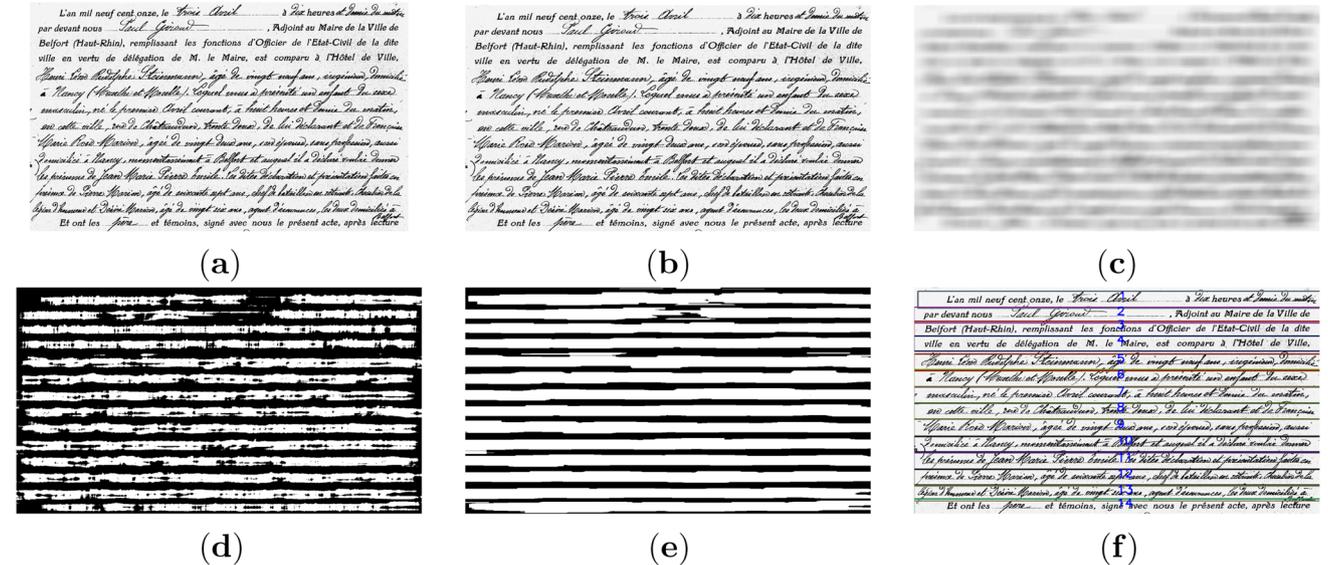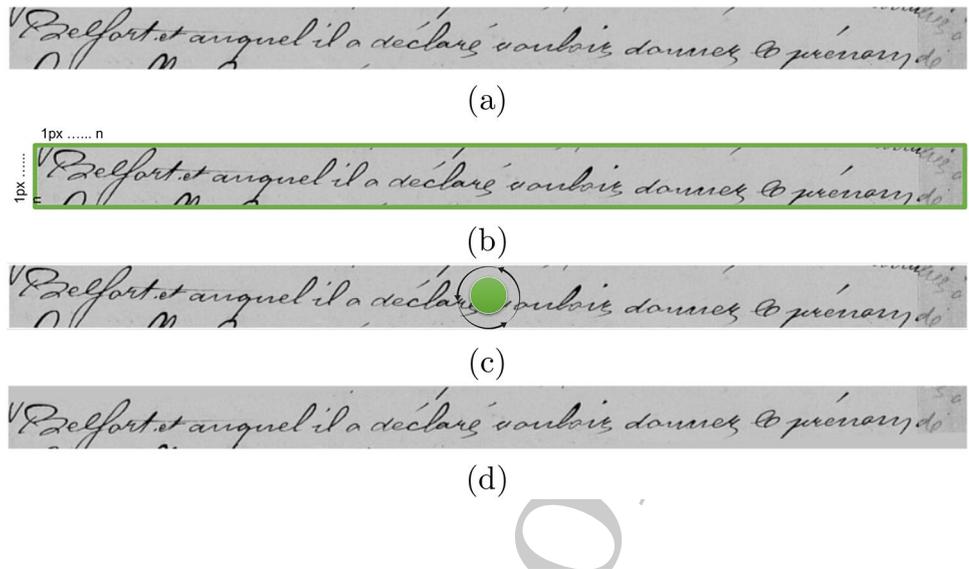


(a)

(b)

(c)

(d)



(a)

(b)

(c)

(d)

(e)

(f)

**Fig. 6** Examples illustrating the preprocessing steps applied to a record from Belfort civil registers of births: **a** Original record. **b** Grayscale conversion. **c** Gaussian blur application. **d** Adaptive thresh- olding. **e** Morphological operations. *f* Contour detection for text line extraction (Source: [10])

records. Each.txt file may contain up to 4 records. Additionally, these files include 984 margin texts. The total number of text lines across all record components amounts to 21,939, comprising 189,976 words and 1,177,354 characters in total. Figure 7 illustrates an example of transcriptions with tags.

### XML File Generator

A tool has been developed to generate XML files, streamlining the preparation of transcribed records for input into the deep learning model. This tool associates each component within the images with its corresponding transcription at both the paragraph and text line levels using added tags. The XML files are structured to include the following properties:

- Image Information: This section includes details such as the image name, type (single or double page), image height, and image width.
- Reading Order: As described in Sect. 3.3.2, this section establishes the sequence of reading by assigning a unique index number to each record. It specifies the types of components included in the reading order (number, name, paragraph, marginal annotation).
- Record number and Name: This section of the XML file includes details about the record's number and name,

**Table 3** The set of tags employed during the manual transcription process (Source: [10])

| Tag | Description |
|---|---|
| < begin> | Begin of the record |
| < text>...< \text> | Primary paragraph of the record |
| < margin>...< \margin> | Margins text |
| < ptext>...< \ptext> | Printed text |
| < striped>...< \striped> | Striped text |
| < unreadable>...< \unreadable> | Unreadable text |
| < added above>...< \added above> | Small text added above the text line |
| < added below>...< \added below> | Small text added below the text line |
| < page> | Start new page |

along with their coordinates within the image and the corresponding transcribed text.

- Paragraphs and Marginal annotations: This section details the primary paragraphs and marginal annotations within the record, including their coordinates in the image and the associated transcribed text. Additionally, it covers information about the text lines within these components, including unique IDs, coordinates in the image, and the corresponding transcribed text.

## Results

The proposed text line segmentation method has been applied to the Belfort civil registers of births, focusing on both primary paragraphs and marginal annotations. Images were chosen from various registers across different time periods. These images include a range of challenges such as hybrid text, text skewness, text overlapping, and diverse handwritten styles, making them ideal for evaluating the method's effectiveness.

Furthermore, we assessed the performance of our method on additional benchmark datasets, such as the IAM historical document database (Saint-Gall) and the Reconnaissance et Indexation de données Manuscrites et de fac-similés (RIMES) dataset. The ground truth for these samples was manually annotated using VGG Image Annotator (VIA) version 2.0.12 [32]. Figure 8 illustrates examples of the annotated datasets used in the evaluation process.

Parameters sensitivity analysis study has been carried out to assess the performance of the proposed method and to identify the optimal values for parameters and thresholds necessary for accurate segmentation. Additionally, a comparison of accuracy with online commercial systems has been conducted to validate the effectiveness of the proposed approach.

Finally, XML files have been created to structure the significant components of the Belfort civil registers of births. These files are crucial for streamlining the data labeling process during the training phase of text recognition models.

## Evaluation Metrics

In this study, various evaluation metrics have been utilized to assess the performance of the proposed method, which were employed in the ICDAR image segmentation contests [33]. These metrics were also discussed in our recently published work [34].

The Intersection over Union ($IoU$) quantifies the overlap between the predicted text line image regions and the

```
<begin>
N° 199
Steinmann.
Jean Marie Pierre Emile
<text>
<ptext>L'an mil neuf cent onze, le <\ptext>trois Avril<ptext>, à<\ptext> dix <ptext>heures <\ptext>et demie du matin<ptext>,
par devant nous <\ptext> Paul Giroud<ptext>, Adjoint au Maire de la Ville de
Belfort (Haut-Rhin), remplissant les fonctions d'Officier de l'Etat civil de la dite
ville en vertu de délégation de M. le Maire, est comparu à l'hôtel de ville,<\ptext>
Henri Léon Rudolphe Steinmann, âgé de vingt neuf ans, ingénieur, domicilié
à Nancy (Meurthe et Moselle). Lequel nous a présenté un enfant du sexe
masculin, né le premier Avril courant, à huit heures et demie du matin,
en cette ville, rue de Châteaudun, trente deux, de lui déclarant et de Françoise
Marie Rose Marion, âgée de vingt deux ans, son épouse, sans profession, aussi
domiciliée à Nancy, momentanément à Belfort et auquel il a déclaré vouloir donner
les prénoms de Jean Marie Pierre Emile. Les dites déclaration et présentation faites en
présence de Pierre Marion, âgé de soixante sept ans, chef de bataillon en retraite, Chevalier de la
légion d'honneur et Désiré Marion, âgé de vingt six ans, agent d'assurances, les deux domiciliés à <added below>Belfort<\added below>.
<ptext>Et ont les<\ptext> père <ptext>et témoins, signé avec nous le présent acte, après lecture<\ptext>
<\text>
<margin>
Décédé à PETRA (Jordanie)
le 8 avril 1963.
Le 20 janvier 1965
<\margin>
```
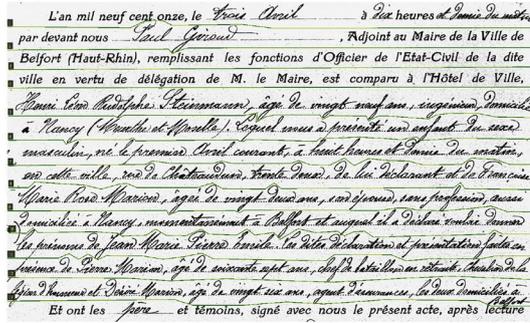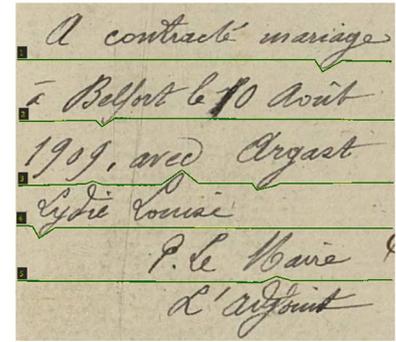
**Fig. 7** Examples of the tags employed in the manual transcription process of Belfort civil registers of births (Source: [10])
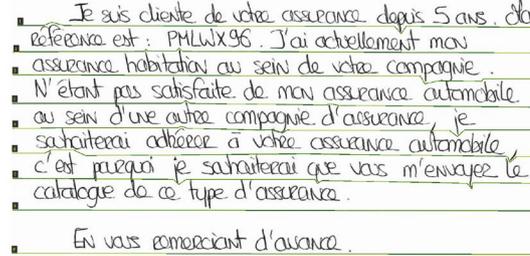
**Fig. 8** Samples of the ground truth annotations generated using the VGG Image Annotator tool are shown as follows: **a** Primary paragraph from the Belfort civil registers of birth. **b** Marginal annotation from the Belfort civil registers of birth. **c** Paragraph from the RIMES dataset. **d** Page from the Saint Gall dataset
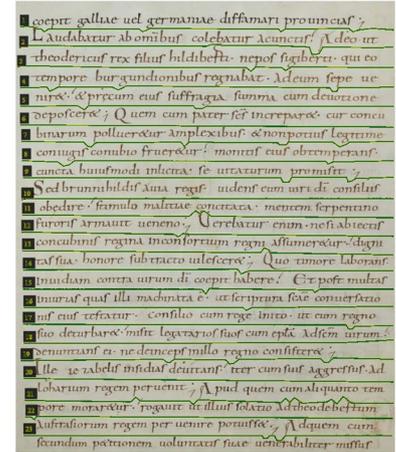
(**a**)

(**b**)

(**c**)

(**d**)

ground truth image regions. It is calculated using the following formula:

$$IoU = \frac{|\text{predicted mask} \cap \text{ground truth mask}|}{|\text{predicted mask} \cup \text{ground truth mask}|}, \quad (7)$$

where |predicted mask ∩ ground truth mask| is the common area shared by both the predicted and ground truth regions, and |predicted mask ∪ ground truth mask| is the total area encompassed by the predicted and ground truth regions.

Moreover, we employed the Detection Rate (DR) metric to measure the ratio of correctly detected text lines to the total number of ground truth text lines. An acceptance threshold of 0.9 is used to count correct text line detections. DR is calculated using the following formula:

$$DR = \frac{O2O}{N_{gt}}, \quad (8)$$

where one-to-one matching ($O2O$) represents the number of correctly matched text lines between the predicted and ground truth text lines, and $N_{gt}$ denotes the total number of ground truth text lines.

Additionally, the Recognition Accuracy (RA) metric is employed to determine the ratio of correctly detected text lines to the total number of predicted text lines. RA is calculated using the following formula:

$$RA = \frac{O2O}{N_{pred}}, \quad (9)$$

where $O2O$ represents the number of correctly matched text lines between the predicted and ground truth text lines, and $N_{pred}$ denotes the total number of predicted text lines.

Finally, the F-Measure (FM) is used to calculate the harmonic mean of the Detection Rate (DR) and the Recognition Accuracy (RA), serving as an overall metric for performance evaluation. It is computed using the following formula:

$$FM = 2 \times \frac{DR \times RA}{DR + RA}. \quad (10)$$

Using these metrics offers a robust and comprehensive evaluation of the proposed method's performance across various aspects. Detection Rate (DR) and Recognition Accuracy (RA) measure the method's ability to detect and recognize

text lines, providing practical performance insights. The F-Measure (FM), as a classical F1-score, provides a balanced assessment by considering both DR and RA.

## Text Line Segmentation

The proposed segmentation approach is applied on the primary paragraphs and marginal annotations, achieving competitive accuracy for the primary paragraphs. However, in cases of heavy text overlap, such as cursive writing, the tool may incorrectly detect multiple lines as a single line. To address this, a second phase refines the initial results by checking the height of the segmented lines. If the height exceeds a predetermined threshold, the lines are horizontally split. Additionally, a 5-pixel padding was added above and below the height of the segmented text lines when determining their coordinates. This padding improves both the visibility and accuracy of the segmentation, especially for various writing styles. Figure 9 depicts the accuracy obtained from the text line segmentation process applied to the Belfort civil birth registers.

We have also tested our method on other datasets, such as the IAM historical document database (Saint-Gall) and the RIMES dataset. These datasets were part of those used in the International Conference on Document Analysis and Recognition (ICDAR) and the International Conference on Frontiers in Handwriting Recognition (ICFHR).

The Saint Gall dataset, detailed in [35], consists of 60 pages from a handwritten Latin manuscript dating back to the 9th century. This dataset is accessible online in JPEG format (300 dpi) and includes binarized, normalized, and text line-level transcriptions. The Reconnaissance et Indexation de données Manuscrites et de fac-similés (RIMES) dataset [36] comprises French handwritten grayscale images at a resolution of 300 dpi. It includes 12,723 pages from 5605

letters written by 1300 volunteers, simulating mail sent to companies via fax or postal service. The dataset contains 99 unique characters and categorizes page regions into several classes: sender, recipient, subject, date, location, opening, body, and attachment. Additionally, segmentation and transcription are available at the paragraph, line, and word levels. Figure 10 shows result examples of the segmentation process applied to the datasets, and Table 4 summarizes the performance accuracy of the proposed method.

An accuracy comparison was performed to evaluate the effectiveness of the proposed segmentation method against online AI commercial systems like DOCSUMO, Ocelus, and Transkribus, which provide free document recognition trials. The experiment focused on text line segmentation using the paragraph image shown in Fig. 3. Table 5 presents the accuracy rates of these systems.

## Parameters Sensitivity Analysis

An analysis study was conducted to systematically test variations and determine the optimal configuration for the segmentation process. This study evaluated the method's performance with different thresholds and parameter values, including the kernel sizes for the Gaussian blur filter, morphological operations, and padding parameter. Table 6 summarizes the key configurations tested in this study and reports the performance accuracy for each configuration. Additionally, Table 7 shows the parameter values used in the text lines skew correction step.

## Results of XML File Generation

An automatic verification tool has been developed to ensure the proper formatting of our manual transcriptions. This tool checks the correct usage of tags within the transcriptions

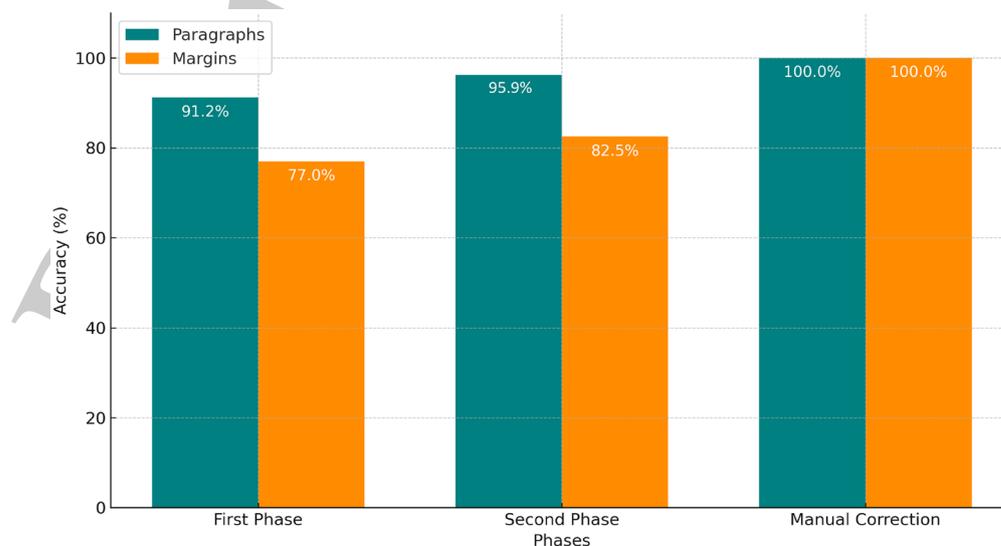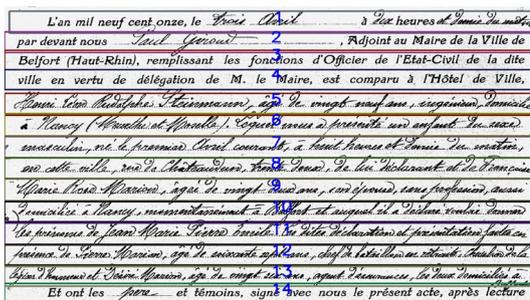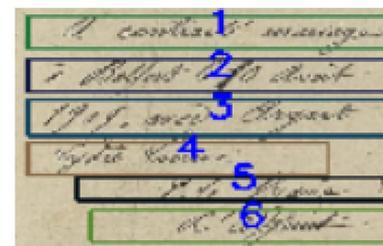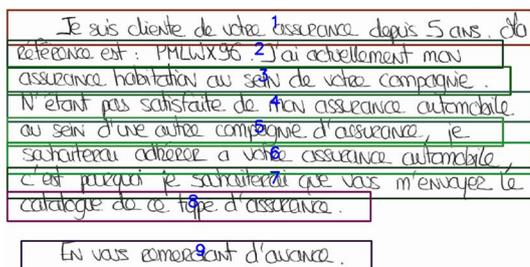**Fig. 9** Accuracy percentages of text line segmentation process

**Fig. 10** Examples of the segmentation results are shown as follows: **a** Segmentation result of the primary paragraph in the Belfort civil registers of birth. **b** Segmentation result of the marginal annotation in the Belfort civil registers of birth. **c** Segmentation result of the paragraph in the RIMES dataset. **d** Segmentation result of the page in the Saint Gall dataset
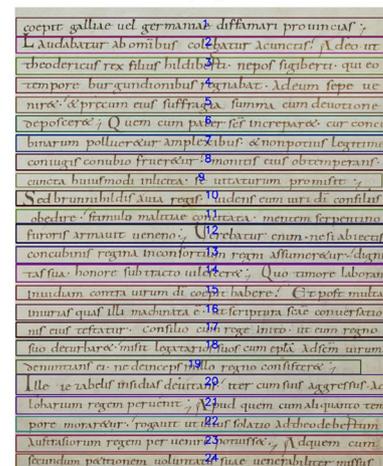
(**a**)

(**b**)

(**c**)

(**d**)

**Table 4** Accuracy results reported on Belfort civil registers of birth, Saint Gall, and RIMES datasets

| Component | Accuracy | | | |
|---|---|---|---|---|
| | IoU (%) | DR (%) | RA (%) | FM (%) |
| Primary paragraphs | 93.3 | 97 | 95.2 | 95.9 |
| Marginal annotations | 77 | 83 | 82 | 82.5 |
| Saint Gall | 98.3 | 98 | 98 | 98.04 |
| RIMES | 96.6 | 95 | 94 | 94.4 |

**Table 5** Accuracy comparison of Belfort civil registers of birth segmentation process with commercial systems. The links of these systems were accessed on July 3, 2024

| System | Accuracy (%) |
|---|---|
| DOCSUMO (URL) | 97 |
| Ocelus (URL) | 99 |
| Transkribus (URL) | 97 |
| Our method | 95.9 |

and maintains accurate alignment between the image components and their corresponding transcriptions during the generation of the .xml files.

We have generated a total of 306 .xml files, comprising 965 images of record numbers, 1935 images of names, 994 images of primary paragraphs, 982 images of marginal annotations, and 17,820 images of text lines. An example of one such XML file is presented in Fig. 11.

## Conclusions

This paper introduces a comprehensive methodology for facilitating the transcription of the Belfort civil birth registers. Our approach encompasses several preprocessing steps, including binarization, skew correction, and segmentation. Despite the inherent challenges of historical documents, such as varied handwriting styles, marginal annotations, and a hybrid of printed and handwritten text, we have

**Table 6** Configuration variations and performance accuracy

| Parameter/ threshold | Configuration | | | Accuracy | | | |
|---|---|---|---|---|---|---|---|
| | GB-KS | Morph-KS | padding | IoU (%) | DR (%) | RA (%) | FM (%) |
| GB-KS | (51, 1) | (1, 201) | 5 | 42 | 45 | 43.5 | 44.2 |
| | (101, 51) | (1, 201) | 5 | 93.3 | 97 | 95.2 | 95.9 |
| | (151, 101) | (1, 201) | 5 | 23.5 | 35 | 32.25 | 33.6 |
| | (201, 151) | (1, 201) | 5 | 10 | 12.4 | 8.5 | 10.1 |
| | (251, 201) | (1, 201) | 5 | 10 | 12.4 | 8.5 | 10.1 |
| Morph-KS | (101, 51) | (1, 51) | 5 | 55.1 | 80 | 67.55 | 73.2 |
| | (101, 51) | (1, 101) | 5 | 88.2 | 90 | 88.1 | 89.0 |
| | (101, 51) | (1, 151) | 5 | 90.14 | 92 | 90.07 | 91.0 |
| | (101, 51) | (1, 201) | 5 | 93.3 | 97 | 95.2 | 95.9 |
| | (101, 51) | (1, 251) | 5 | 87.3 | 90 | 87.15 | 88.6 |
| Padding | (101, 51) | (1, 201) | 1 | 91.7 | 92 | 90.35 | 91.1 |
| | (101, 51) | (1, 201) | 3 | 93.1 | 95 | 93.05 | 94.0 |
| | (101, 51) | (1, 201) | 5 | 93.3 | 97 | 95.2 | 95.9 |
| | (101, 51) | (1, 201) | 7 | 92.6 | 92 | 91.3 | 91.6 |
| | (101, 51) | (1, 201) | 9 | 90.3 | 90 | 90.15 | 90.1 |

*GB-KS: Gaussian blur kernel size, *Morph-KS: Morphological operation kernel size

**Table 7** Parameters values used in the text lines skew correction step (Source: [10])

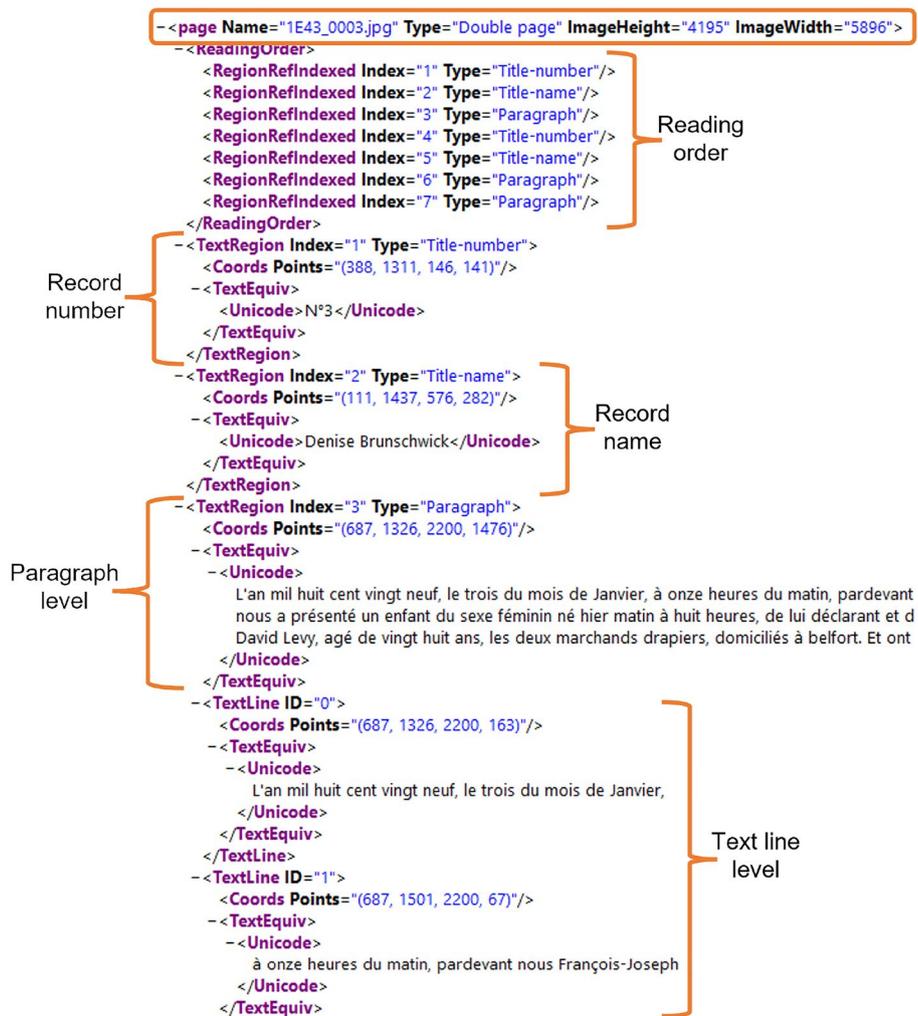| Parameter | Value |
|---|---|
| Gaussian Blur Kernel Size | (5, 5) |
| Canny Edge Threshold1 | 50 |
| Canny Edge Threshold2 | 150 |
| HoughLinesP Threshold | 100 |
| HoughLinesP MinLineLength | 100 |
| HoughLinesP MaxLineGap | 55 |

devised effective solutions to improve the accuracy of text line segmentation.

Our experimental results show a high level of accuracy in text line segmentation, validated against benchmark datasets such as the IAM historical document database (Saint-Gall) and the RIMES dataset, as well as compared to commercial systems. The development of an automatic verification tool and an XML file generator ensures that transcriptions are accurately formatted and aligned with their corresponding image components.

This work significantly contributes to the field of hand-written text recognition by introducing a new structured dataset tailored for historical documents. Future research will aim at developing an automatic document layout analysis tool and a deep learning model to automate the transcription of the remaining records, thus aiding in the recognition of this valuable cultural heritage.

**Fig. 11** Example of the.xml
files of Belfort civil registers of
births (Source: [10])

**Author Contributions** Wissam AlKendi: Conceptualization, Writing - review and editing, Writing - original draft, Visualization, Validation, Software, Methodology, Formal analysis. Franck Gechter: Conceptualization, Review and editing, Supervision. Laurent Heyberger: Conceptualization, Writing - original draft, Review and editing, Supervision. Christophe Guyeux: Conceptualization, Review and editing, Supervision.

**Data Availability** The datasets (Belfort civil registers of births) analyzed during the current study are available online as mentioned, and from the corresponding author on reasonable request.

## Declarations

**Conflict of Interest** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

1. Philips J, Tabrizi N. Historical document processing: Historical document processing: A survey of techniques, tools, and trends. ArXiv **abs/2002.06300** 2020. https://doi.org/10.5220/0010177403410349

2. Wang Y, Xiao W, Li S. Offline handwritten text recognition using deep learning: a review. J Phys: Conf Ser. 2021;1848(1):012015. https://doi.org/10.1088/1742-6596/1848/1/012015.

3. Gan J, Wang W, Lu K. In-air handwritten chinese text recognition with temporal convolutional recurrent network. Pattern Recognition. 2020;97: 107025. https://doi.org/10.1016/j.patcog.2019.107025.

4. Chacko BP, P, BA. Pre and post processing approaches in edge detection for character recognition. In: 2010 12th International Conference on Frontiers in Handwriting Recognition, 2010;676–681. https://doi.org/10.1109/ICFHR.2010.111

5. Diem M, Kleber F, Sablatnig R. Text classification and document layout analysis of paper fragments. In: 2011 International Conference on Document Analysis and Recognition, 2011;854–858. https://doi.org/10.1109/ICDAR.2011.175

6. Carbune V, Gonnet P, Deselaers T, Rowley HA, Daryin A, Calvo M, Wang L-L, Keysers D, Feuz S, Gervais P. Fast

multi-language lstm-based online handwriting recognition. IJDAR. 2020;23(2):89–102.

7. Bugeja M, Dingli A, Seychell D. An overview of handwritten character recognition systems for historical documents. Rediscovering Heritage Through Technology: A Collection of Innovative Research Case Studies That Are Reworking The Way We Experience Heritage, 2020;3–23. https://doi.org/10.1007/978-3-030-36107-5_1

8. Nikolaidou K, Seuret M, Mokayed H, Liwicki M. A survey of historical document image datasets. IJDAR. 2022;25(4):305–38.

9. Delsalle P. Histoires de Familles. Les Registres Paroissiaux et D'état Civil, du Moyen Âge à Nos Jours: Démographie et Généalogie (Family History: Parish and Civil Status Registers, from the Middle Ages to the Present Day: Demography and Genealogy). Presses universitaires de Franche-Comté, Besançon 2009.

10. AlKendi W, Gechter F, Heyberger L, Guyeux C. Belfort birth records transcription: Preprocessing, and structured data generation. In: IMPROVE, 2024;32–43.

11. Hussain SA-K, Al-Nayyef H, Al Kindy B, Qassir SA. Human earprint detection based on ant colony algorithm. Int J Intell Syst Appl Eng. 2023;11(2):513–7.

12. Al-Khalidi FQ, Alkindy B, Abbas T. Extract the breast cancer in mammogram images. Int J Civ Eng Technol. 2019;10(02):96–105.

13. Binmakhashen GM, Mahmoud SA. Document layout analysis: a comprehensive survey. ACM Comput Surv. 2019;52(6):1–36.

14. Ganchimeg G. History document image background noise and removal methods. Int J Knowl Content Dev Technol. 2015;5(2):11–24.

15. Chakraborty A, Blumenstein M. Marginal noise reduction in historical handwritten documents–a survey. In: 2016 12th IAPR Workshop on Document Analysis Systems (DAS). IEEE. 2016;323–328. https://doi.org/10.1109/DAS.2016.78

16. Montrésor S, Memmolo P, Bianco V, Ferraro P, Picart P. Comparative study of multi-look processing for phase map de-noising in digital fresnel holographic interferometry. J Opt Soc America A. 2019;36(2):59–66. https://doi.org/10.1364/JOSAA.36.000A59.

17. Dabov K, Foi A, Katkovnik V, Egiazarian K. Image denoising by sparse 3-D transform-domain collaborative filtering. IEEE Trans Image Proc. 2007;16(8):2080–95.

18. Mustafa WA, Kader MMMA. Binarization of document images: A comprehensive review. J Phys: Con Ser. 2018;1019:012023. IOP Publishing. https://doi.org/10.1088/1742-6596/1019/1/012023

19. Tensmeyer C, Martinez T. Historical document image binarization: a review. SN Comput Sci. 2020;1(3):173.

20. Ahmed AS. Comparative study among Sobel, Prewitt and Canny edge detection operators used in image processing. J Theor Appl Inf Technol. 2018;96(19):6517–25.

21. Biswas B, Bhattacharya U, Chaudhuri BB. Document image skew detection and correction: A survey 2023. Int J Innovative Res Technol. 2023;9(10):949–63.

22. Singh H, Kaur N, Kaur H. Analysis of line segmentation methods for historical manuscripts. In: 2022 2nd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE). IEEE. 2022;2043–2045. https://doi.org/10.1109/ICACITE53722.2022.9823610

23. Pach JL, Bilski P. A robust binarization and text line detection in historical handwritten documents analysis. Int J Comput. 2016;15(3):154–61.

24. Plateau-Holleville C, Bonnot E, Gechter F, Heyberger L (2021) French vital records data gathering and analysis through image processing and machine learning algorithms. J Data Min Digital Humanit. https://doi.org/10.46298/jdmdh.7327

25. Saabni R, Asi A, El-Sana J. Text line extraction for historical document images. Pattern Recognit Lett. 2014;35:23–33.

26. Louloudis G, Gatos B, Pratikakis I, Halatsis C. Text line and word segmentation of handwritten documents. Pattern Recognit. 2009;42(12):3169–83.

27. Papavassiliou V, Stafylakis T, Katsouros V, Carayannis G. Handwritten document image segmentation into text lines and words. Pattern Recognit. 2010;43(1):369–77.

28. Alaei A, Pal U, Nagabhushan P. A new scheme for unconstrained handwritten text-line segmentation. Pattern Recognit. 2011;44(4):917–28.

29. AlKendi W, Gechter F, Heyberger L, Guyeux C. Advancements and challenges in handwritten text recognition: a comprehensive survey. J Imaging. 2024;10(1):18.

30. Oliphant TE, et al. Guide to Numpy vol. 1. Trelgol Publishing USA, 2006.

31. Mohaideen Abdul Kadhar K, Anand G. Image processing using opencv. In: Industrial Vision Systems with Raspberry Pi: Build and Design Vision Products Using Python and OpenCV, pp. 87–140. Springer, 2024. https://doi.org/10.1007/979-8-8688-0097-9_5

32. Dutta A, Gupta A, Zissermann A. VGG Image Annotator (VIA) 2016.

33. Phillips IT, Chhabra AK. Empirical performance evaluation of graphics recognition systems. IEEE Trans Pattern Anal Mach Intell. 1999;21(9):849–70.

34. AlKendi W, Gechter F, Heyberger L, Guyeux C. Unsupervised approach to text line extraction in Belfort civil registers of births. IJDAR. 2024. https://doi.org/10.1007/s10032-024-00507-5.

35. Fischer A, Frinken V, Fornés A, Bunke H. Transcription alignment of latin manuscripts using hidden markov models. In: Proceedings of the 2011 Workshop on Historical Document Imaging and Processing, 2011;29–36. https://doi.org/10.1145/2037342.2037348

36. Augustin E, Carré M, Grosicki E, Brodin J-M, Geoffrois E, Prêteux F. Rimes evaluation campaign for handwritten mail processing. In: International Workshop on Frontiers in Handwriting Recognition (IWFHR'06),, 2006;231–235.