



# Article Exploiting Script Similarities to Compensate for the Large Amount of Data in Training Tesseract LSTM: Towards Kurdish OCR

Saman Idrees and Hossein Hassani \*D

\* Correspondence: hosseinh@ukh.edu.krd

Featured Application: This work helps in the preparation of OCR for the Kurdish language. In particular, its focus is on Kurdish texts written in Persian-Arabic script. Currently, Kurdish OCR is in its early stages. This work can assist in preparing the environment for a full-fledged OCR application for Kurdish.

**Abstract:** Applications based on Long-Short-Term Memory (LSTM) require large amounts of data for their training. Tesseract LSTM is a popular Optical Character Recognition (OCR) engine that has been trained and used in various languages. However, its training becomes obstructed when the target language is not resourceful. This research suggests a remedy for the problem of scant data in training Tesseract LSTM for a new language by exploiting a training dataset for a language with a similar script. The target of the experiment is Kurdish. It is a multi-dialect language and is considered less-resourced. We choose Sorani, one of the Kurdish dialects, that is mostly written in Persian-Arabic script. We train Tesseract using an Arabic dataset, and then we use a considerably small amount of texts in Persian-Arabic to train the engine to recognize Sorani texts. Our dataset is based on a series of court case documents in the Kurdistan Region of Iraq. We also fine-tune the engine using 10 Unikurd fonts. We use Lstmeval and Ocreval to evaluate the outputs. The result indicates the achievement of 95.45% accuracy. We also test the engine using texts outside the context of court cases. The accuracy of the system remains close to what was found earlier indicating that the script similarity could be used to overcome the lack of large-scale data.

**Keywords:** optical character recognition; tesseract; printed-document OCR; Kurdish-OCR system; offline character recognition system

# 1. Introduction

The less-resourced languages face various issues from the language technology perspective. The lack of resources often is one of the main obstacles in the forefront of the resolution of those challenges. The Optical Character Recognition (OCR) systems play a crucial role in resource preparation and language processing. However, it is not feasible to develop or adapt an appropriate OCR for those languages if the required data does not exist. Kurdish is considered a less-resourced language that is facing several issues [1] including lack of OCR systems that are accurate and widely available. Tesseract Long-Short-Term Memory (LSTM) is a popular OCR engine that has been adapted for various languages, but LSTM-based methods require large amounts of data for their training. In this research, we suggest a resolution for that issue by using a training dataset for a language with a similar script and then applying a considerably small amount of data from the target language. Our target language is Kurdish, and we choose Sorani, one of its dialects that is mostly written in Persian-Arabic script. Therefore, we use an Arabic dataset as the base of the training, and we use a Kurdish dataset based on a series of court case



Citation: Idrees, S.; Hassani, H. Exploiting Script Similarities to Compensate for the Large Amount of Data in Training Tesseract LSTM: Towards Kurdish OCR. *Appl. Sci.* 2021, *11*, 9752. https://doi.org/ 10.3390/app11209752

Academic Editor: Agnese Magnani

Received: 8 August 2021 Accepted: 27 September 2021 Published: 19 October 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/).

Department of Computer Science and Engineering, University of Kurdistan Hewlêr, 30 Meter, Kurdistan Region, Erbil 44001, Iraq; saman.idrees@ukh.edu.krd

documents in the Kurdistan Region of Iraq. We also investigate the fine-tuning phase of Tesseract using 10 Unikurd fonts and evaluate the approach by Lstmeval and Ocreval.

The LSTM-based methods are not only used in image recognition systems, but they are also applicable to address and solve a wide range of problems. For example, they are used in time-series forecasting [2], energy consumption forecasting [3], commercial vacancy prediction [4], and various other areas. Although this paper focuses on suggesting a solution for the requirement of large amounts of data in one aspect of the applications of LSTM-based approaches, the suggestion could be extended to similar cases that have similarities of some sort in their training data.

The rest of this paper is organized as follows. Section 2 provides a background about the Kurdish language focusing on its alphabets and script similarities with Arabic, and then briefly introduces the OCR history focusing on Tesseract. Section 3 reviews the related work and literature. In Section 4, we discuss our approach. Section 5 presents and discusses the results. Finally, Section 6 concludes the paper.

# 2. Background

In recent years, research on Kurdish language technology has gained more attraction. Kurdish is spoken by approximately 30 million people [1]. It is a multi-dialect language, and its two most widely spoken dialects, Kurmanji (Northern Kurdish) and Sorani (Southern Kurdish), are spoken by roughly 75% of native Kurds [5,6]. Sorani usually is written in an adapted Persian-Arabic script with a cursive style and from Right To Left (RTL). Kurmanji mostly uses Latin for writing, except in the Kurdistan Region of Iraq and Kurdish areas of Syria, where they use the same script that Sorani uses [1]. The alphabet of the mentioned script includes 34 letters that appear in different shapes according to their position in words (see Table 1). Kurdish language technology and its challenges have been addressed and categorized as: dialect heterogeneity, script variety, lack of standards, lack of resources, and lack of investment [1,7,8].

Contextual Form				
#	Isolated	Initial	Medial	Final
1	-	ءِ 1	-	
2	I		L	
3	ب	ب	<del></del>	ب
4	پ	÷	<del></del>	پ
5	ت	Ľ.	<u></u>	۔ ت
6	ج	<del>ج</del>	÷	_ج
7	ş	<del>ي</del>	<del></del>	-5
8	ζ	۔ ح	۔ ح	-ح
9	ź	خ	خ	۔ _خ
10		د		<u> </u>
11		ر		_ر
12		ړ		_ر
13		j		-ز
14		ژ		_ژ
15	س	س_		ے
16	ŵ	<u></u> ســ	<u></u>	ىش
17	ع	<u>عـ</u>	ے	ح
18	ż	غـ	ف	_غ
19	ف	<u> </u>	<u> </u>	ف
20	ڤ	<u>ف</u>	<u></u>	ف

Table 1. Kurdish alphabet forms, both Ligature and Isolated.

Contextual Form					
#	Isolated	Initial	Medial	Final	
21	ق	<u></u>	ä	ـق	
22	ك	ک	5	<u></u>	
23	گ	گ	<u> </u>	گ	
24	ل	ل_	1	ـل	
25	ť	-	Ť	ـٽ	
26	م	م_		ح	
27	ن	<u>ن</u>	<u></u>	_ن	
28		ھ	<u>_</u> &	-	
29		٥		4_	
30		و		_و	
31		ۆ		_ۆ	
32		وو		_وو	
33	ى	<u></u>	<del></del>	_ى	
34	ێ	<u></u>	<u> </u>	_ێ	

Table 1. Cont.

Table 2 shows the order of the alphabet of the Persian-Arabic script based on the Kurdish Academy proposal. The Arabic alphabet has 28 letters with a cursive style (see Table 3). The Persian writers augmented the Arabic alphabet to include graphemes for the phonemes that did not exist in Arabic. The Kurdish writers also did the same on the Persian-Arabic version [9]. The augmented scripts are known as Persian-Arabic scripts though the Kurdish one includes more new graphemes for the phonemes than the one that is used in Persian (Farsi) writing.

Table 2. Central Kurdish alphabet.

ئ <u>ہ</u>	י	ب	پ	ت	ج	ङ	с	ċ	د	ر	ر	ز	ژ	س	ش	٤
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
خ	ف	ڤ	ق	و	گ	ل	ڵ	م	ن	ھـ	ک	ہ	ۆ	وو	ى	ێ
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34

The Kurdish Persian-Arabic script uses 21 out of 28 Arabic letters, keeps all four extra letters added by Persians, and adds nine letters. Seven Arabic letters were found unsuitable for Kurdish because they have no corresponding sound (phoneme) in the language. Unsuitable letters for the Kurdish language are فر فر, فر, فر, فر, فر, فر, فر, خر, ث. Five new letters have been formed out of the residual letters, as shown in Tables 4 and 5. As a result, the current Kurdish (Sorani) alphabet consists of 34 characters.

Contextual Form					
#	Isolated Form	Initial	Medial	Final	
	¢		١	Ĺ	
1	(used mainly in the m	edial and final	ٳ	L_s	
	position, which is an o		ۇ	_ؤ	
		ئ	<u><u></u></u>	ئ ہئ	
2	1		L		
3	ب	÷	<u>.</u>	÷	
4	ت	ت	<u></u>	ت	
5	ث	ث	<u></u>	ؿ	
6	5	÷	÷	_ح	
7	2	<u>ح</u>	<u>ح</u>	_ح	
8	ż	خـ	خ		
9	-	د		<u>۲</u>	
10		ć		ż	
11		ر		_ر	
12		j		۔ _ز	
13	ىس	ىب_		_س	
14		صـ		ے	
15	_ ض	ضـ	ے	_ 	
16	ط	ط	لط	ط	
17	ظ	ظ	ظ	ظ	
18	ş	<u>عـ</u>	æ	۶_	
19	÷	<u> </u>	ف	د نر	
20	ف	<u> </u>	<u>.</u>	ف	
21		<u> </u>	<u></u>	, Ŭ	
22	ك	ک	2	ل ب	
23	(1	ل	1	, L	
24	ھ			-ت ب	
25		÷	<u>·</u>	\ 	
26	۵	 ه_		 د_	
27	-		-6-	_ _	
28	ي	و بـ	<del></del>	-ر -ى	

# Table 3. Arabic alphabet.

Table 4. Modified Arabic letters to Kurdish letters.

Modified Arabic Letters to Kurdish					
#	Arabic Letters	ТО	Kurdish Letters		
1	ب		ڀ		
2	5		چ		
3	į		ć		
4	ف		ف		
5	ک		گ		

As indicated by Hashemi [10], the Kurdish alphabet has seven vowels, which are presented with seven corresponding letters (see Table 6).

Arabic Letters with Kurdish Diacritics					
#	Arabic Letter	ТО	Kurdish Letter		
1	ر		ر		
2	ل		ť		
3	و		ۆ		
4	ى		ێ		

Table 5. Diacritics added to the Arabic letters to create Kurdish letters.

Table 6. Vowel letters in Kurdish alphabet.

Kurdish Alphabet Vowel				
#	Kurdish Vowel	API	Example	
1	1	ä	ب (air)	
2	٥	$\epsilon$	(head) سەر	
3	و	u	(Kurd) کورد	
4	ۅٚ	ö	(You) تۆ	
5	وو	ü	(far) دوور	
6	ى	ï	(blue) شين	
7	ێ	ë	(village) دئ	

The Kurdish alphabet also uses a specific character named Zero-Width Non-Joiner (ZWNJ) for digital writing. ZWNJ is a non-printing character placed between two characters to avoid inappropriate concatenation. For example, after the character "ه". ZWNJ is encoded in Unicode as U+200C. It is illustrated in the following words: هەزاير - هەزاير - هەزاير - هەزاير - هەزاير - هەزاير - هەزار حمهزار حمهزار عوار حمهزار على المحالية (CWNJ) was not used, those words appeared as

## **Optical Character Recognition**

Optical Character Recognition (OCR) is a research area that benefits from several computing fields, such as machine learning, computer vision, and natural language processing [11]. OCR essentially converts two types of documents into texts: handwritten and machine-typed [12]. Those documents may hold different kinds of data, such as passport documents, invoices, bank statements, digital receipts, business cards, mails, and newspapers, and OCR can make them ready for text processing.

Early character recognition systems followed the telegraphy applications technology. That allowed the industry to develop devices that could help blind and visually impaired to read texts [13]. In 1914, Emanuel Goldberg developed a machine that could read and convert characters into standard telegraph codes [14]. Concurrently, Edmund Edward Fournier developed what was a handheld scanner used to move across a printed page to produce tones that corresponded to a specific letter or character [15]. In the 1930s, Emanuel Goldberg built a statistical machine using optical code recognition for searching microfilm archives. Furthermore, in the 1950s, the US Department of Defense developed Geographic Information Systems and Mapping Operations (GISMO) that could read Morse Codes and words on printed pages, character by character [14].

In 1974, Rey Kurzweil developed the Omni-font OCR machine for blind users. The device was able to recognize printed text in any font. In the early 1990s, A. G. Ramakrishnan built a print-to-braille device that could convert scanned images of printed books to braille books. In the new millennium, OCR became widely available both as an online and offline service under different computing platforms [12].

# 3. Related Work

In this review, our focus is on OCR systems, but we also review some other LSTMbased applications that could benefit from our approach. For the OCR applications, the research focuses on machine-typed documents in the Kurdish language. However, because

6 of 20

the work on Kurdish OCR is in its early stages, we concentrate on OCRs for languages with Arabic-based scripts particularly, Arabic, Urdu, and Farsi (Persian). We also look into available platforms and technologies for OCR development.

#### 3.1. OCR Studies for Persian-Arabic Based Scripts

Yaseen and Hassani [16] worked on OCR for Kurdish (Sorani) texts in Persian-Arabic script. They discussed the features of the script in terms of its special diacritics and cursive style, and accordingly, they proposed a new set of rules for the segmentation process. They used optimized horizontal histogram projection for line segmentation and contour-based segmentation for character segmentation. The research used the Gamera [17] toolkit for classification and reported a 90.82% of accuracy for their approach.

While conducting this study, we also noticed work on Sorani (Central Kurdish) OCR using Tesseract OCR engine (more detail could be found on its GitHub page at https: //github.com/Shreeshrii/tesstrain-ckb [accessed on 7 August 2021]). The mentioned work created a language model for Sorani using synthetic data generated over 26 fonts (see Table 7). However, at the time of writing this paper, we could not find any papers, whether peer-reviewed or as work in progress, to explain more than what the GitHub page describes.

Table 7. Fonts used by Shreeshrii to fine-tune the Arabic language model for Sorani.

#	Font Name	#	Font Name	#	Font Name
1	Abd Akre	10	Arial	19	Segoe UI Bold
2	Abd Gare	11	Arial Bold	20	Speda
3	Abd Halabja OLD	12	NRT Bold,	21	Speda Bold
4	Abd Hewler	13	NRT Reg	22	Tahoma
5	Abd Metin Bold	14	Peshang Des 1,	23	Tahoma Bold
6	Abd Umed	15	RudawRegular	24	Ubuntu Kurdish 0.81 met
7	Adobe Arabic Bold	16	Sarchia_Banoka_1	25	Unikurd Chimen
8	Arabic Typesetting	17	Segoe UI	26	Unikurd Goran
9	Unikurd Web	18	Shasenem-kurd		

Radhiah et al. [18] used Artificial Neural Network (ANN) and Hidden Markov Model (HMM) models in the segmentation process of their Arabic OCR. Their focus was on both isolated and concatenated forms of Arabic characters. They reported a 100% recognition accuracy with ANN and 69% with HMM. For HMM, they reported an accuracy rate of 71% and 50% for isolated and concatenated letters, respectively. According to their results, the ANN method misidentified five letters, while HMM approach misidentified nine. However, it is not clear how they achieved a 100% accuracy with ANN while the method wrongly identified five letters.

Nashwan et al. [19] conducted a study for isolated Arabic letters using ANN with Backpropagation and Learning Vector Quantization (LVQ) method. The research showed a 98.81% accuracy. However, it reported an issue about the confusion between isolated and concatenated letter recognition.

The Urdu language is also written from right to left except for the numerals that follow the similar languages wherein they are written from left to right [20]. Urdu consists of 38 letters and ten numeric characters.

A study by Naz et al. [21] proposed preprocessing techniques for OCR systems using four steps. The first step includes document scanning, orientation detection, skew correction, noise removal, binarization, and then other standard preprocessing methods. The second step is segmentation which separates the paragraphs and text lines. The third step is to extract features such as the right, left, top, and bottom points, curves, loops, crosses, and the height of characters in the text. The last step is the classification that uses fuzzy logic rules in two phases to recognize the diacritics and visible features. They used Nasta'liq font of the size of 36 for the proposed method. They reported an accuracy rate of 100% and 94% for the baseline/single character identification and ligature/concatenated identification, respectively.

Developers of Persian OCR have adopted several methods, such as movement on edge to capture character features, HMM, and morphology of characters. For example, using a chain-code-based approach and nearest neighbor, they trained a classifier, font by font, and then compared the result with a different font as the test set. The method, on average, showed an accuracy rate of 97.4% [22]. Using morphological operators, Jelodar et al. [23] applied hit/miss operators to describe all sub-words combined with the template matching. The study used one font (Lotus) of two different sizes. The test set contained 3000 words (15,000 characters), and the experiments achieved an accuracy rate of 99.9%.

#### 3.2. OCR Technology and Platforms

A variety of commercial companies along with non-commercial communities provide OCR services and applications. The services are available online or offline. However, currently, a functional and reliable Kurdish OCR is not available. In this section, we review some available OCR services in which we concentrate on Tesseract as one of the widely used engines. We also address some other engines that are available as open-source.

## 3.2.1. Tesseract

Tesseract is one of the most widely used OCR engines that provides a high accuracy rate compared with other available engines [11]. Figure 1 shows a brief history of Tesseract. Tesseract has been adapted for many languages (up to 140 different languages) [24]. Since its version 4.0, it presented a new engine based on Long Short-Term Memory (LSTM). LSTM, as a specific form of Artificial Recurrent Neural Network (RNN), provides substantially higher accuracy on image recognition than Tesseract's earlier versions. Tesseract can be trained from scratch or be fine-tuned based on already trained languages. In 2005, it became open-source, and it is freely available on http://code.google.com/p/tesseract-ocr (accessed on 7 August 2021).



Figure 1. Tesseract Timeline [25].

## 3.2.2. Other OCR Technologies

Other OCR technologies also exist and are used for different purposes. Gamera is a toolkit for building document image recognition systems [26]. The images or symbols must be assigned to the classes manually, and the result of this step creates a database in the Extensible Markup Language (XML) format. Gamera uses the K Nearest Neighbors (KNN) algorithm as its classifier. Gamera is available on https://gamera.informatik.hsnr. de/index.html (accessed on 7 August 2021).

Google Drive OCR is a multilingual and online OCR service provided by Google Incorporation and supporting up to 200 languages in almost 25 writing systems. According to Genzel and Ashok Popat [27], Hidden Markov Model (HMM) has been employed to work potentially with all languages and writing systems. Although it is a powerful OCR system, it does not support the Kurdish language [28]. Google drive can partially recognize Kurdish documents based on its Arabic model, but it misses the special Kurdish characters that do not exist in the Arabic alphabet.

GOCR/JOCR is an open-source OCR tool under the GNU license [29] which can read and convert images files such as PNM, PGM, PBM, and PP into text format. GOCR does not require training data that results in a fairly low accuracy rate for different languages. The recognition process happens in two stages. In the first stage, GOCR takes the whole document as input, and in the second one, it processes misidentified characters from the same document [29].

Tafti et al. [11] provided a comparative analysis on various OCR engines, their capability, working platform, and other attributes. Recently, a new survey with a broader coverage was also published [30].

Finally, literature reports on research about handwritten Kurdish texts. For example, Zarro and Anwer [31] and Ahmed [32] have reported on handwritten character recognition in Kurdish, focusing on single and isolated characters. However, as we mentioned, our focus is on machine-type documents, and therefore, we are not reviewing those studies in detail.

To summarize, although OCR is a well-developed technology in language processing, most less-resourced languages have not benefited from the advances in the field. Furthermore, as most modern OCR engines use artificial intelligence methods such as LSTM, they require a large amount of data for their training. We are interested in showing that leveraging the script similarity between languages is a solution to that issue. Although the recent literature presents some work on using limited data in OCR adaption, particularly on LSTM-based engines [33–36], it does not show script-similarity-based approaches based on using well-resourced language to the advantage of less-resourced ones. That is what we intend to show in this research. In our method, the emphasis and concentration are on a specific engine (Tesseract).

#### 3.3. Using LSTM-Based Approaches in Other Fields or Study

Peng et al. [2] used LSTM-based methods to forecast time-series forecasting. They used LSTM with fruit FOA-Fly Optimization Algorithm with LSTM (LSTM) to solve time series problems. Peng et al. [3] applied LSTM-based approaches on energy consumption forecasting. They suggested that their model could be retrained with data collected after the COVID-19 pandemic. As these data might not be as large as enough for an LSTM-based system, the idea of using a small dataset based on a trained model over a similar one, as we suggest in this paper, could be advantageous in the mentioned study too. Commercial vacancy prediction is another example of using LSTM-based methods [4]. In this study, Lee et al. [4] applied LSTM-based methods to predict commercial vacancy. Again, a dataset from the after COVID-19 era that could be reasonably smaller than the original dataset could be used for retraining the model.

# 4. Proposed Method

In this research, we train the Tesseract engine using an Arabic dataset and fine-tune it with a small amount of Kurdish (Sorani) texts. The LSTM-based systems require a large amount of training data. That is typically expensive, labor-intensive, and time-consuming.

Tesseract provides various training methods such as fine-tuning the existing language model and training from scratch. It follows two different approaches for its training: training with synthetic data/image lines and providing ground truth transcriptions. Augmentation could also be applied, and it would be useful when a pre-trained dataset is available.

We train Tesseract and fine-tune it by augmenting the Arabic language model by providing image lines and ground truth transcriptions as input datasets. Using Kurdish documents, we introduce the Kurdish-specific characters. We also examine to what extent a post-processing phase could resolve the issues regarding ZWNJ. Afterward, we retrain the fine-tuned language model with synthetic data to add new Unikurd fonts to compare the accuracy with the post-processing approach.

# 4.1. Tesseract Training Modes and Available Language Models

Tesseract 4 provides various data models for Latin languages. The models are trained on approximately 400,000 text lines covering roughly 4500 fonts [37]. Although the same number of fonts have not been available for non-Latin scripts, they have also been trained with a similar amount of text lines. The training phase is a long process, and depending on the computation power, it might take a few days or weeks. Despite having different training options, the training steps are identical. Tesseract training options are as follows:

- 1. Fine-tune: This method adds new data to an already trained model close to the target language model.
- Cut off the top layer (or some arbitrary layers): If fine-tuning did not produce the desired outcome, this method cuts off the top layer and retrains it with additional data. It can train a new language or a language having script-similarity with a trained one.
- 3. Retrain from scratch: This is a method that starts from scratch. We do not recommend it if a large amount of data is not available.

According to Google [37], in Tesseract 4, old engines still exist, but they are deprecated and could be removed in the future. Tesseract has two OCR engines (Legacy and LSTM) that can be selected by OCR Engine Mode (–OEM) option, see Table 8. We used the LSTM engine that is faster and appears more promising to provide higher accuracy.

OEM	Туре
0	Legacy engine only
1	Neural nets LSTM engine only
2	Legacy + LSTM engines
3	Default, based on what is available

Table 8. Tesseract OCR engine modes.

4.1.1. Two Main Sources of Data for Training Tesseract LSTM Language Model

Below we describe the two main sources of data for training Tesseract LSTM. Other methods also exist but do not suit our approach.

- 1. Training with Synthetic data: Synthetic data is defined as artificially generated data from text files rather than being captured from actual events. It is used in data mining and other different systems such as fraud detection applications. Synthetic data is usually used in machine learning applications to create a model that could be applied to real data [38].
- 2. Training with Image lines: This is the process of generating data from actual documents. The documents are obtainable from printed sources, digital images (or scanned documents), or any other suitable and available formats. Ground truth descriptions for each image must be provided.

# 4.1.2. Data Collection and Preparation

We collect court cases from several courts in the Kurdistan Region of Iraq. We scan the documents with a resolution of 300 Dots Per Inch (DPI) according to Tesseract requirements. For each line of the text, we prepare an image line in PNG. Although various image formats such as Tagged Image File Format (TIFF) are acceptable for this purpose, we chose PNG that is more common for this process. Since no specific tool is recommended in the literature, any image editor with the capability of cropping and rotation that supports the mentioned image formats can be used for the process. For the synthetic data we use available online sources.

# 4.2. Evaluation Method

Tesseract has an evaluation program named *lstmeval*. Using a list of *lstmf* files for evaluation, it performs its evaluations. However, other tools are also available for evaluating the accuracy of OCR language models, for instance, *ocreval*, ISRI-OCR's evaluation tool [39]. *Ocreval* is a powerful tool for evaluating character level and word level accuracy that supports all characters represented in UTF-8 encoding and provides more details in comparison to the Tesseract *lstmeval* program. Tesseract *lstmeval* evaluates LSTM-based networks. In our case, it evaluates the performance of the network in an optical character recognition model, while *ocreval* evaluates the overall performance of an OCR system by comparing the output of the OCR-ed document with the original one.

We initially use *lstmeval*, and we compare the results between the original Arabic language model with the fine-tuned *kur\_ara* language model. We show the result after applying the post-processing method using Algorithm 1. Then we present the result after training/fine-tuning for adding new fonts. We also use the *ocreval* program to present more detail on results.

Alg	orithm 1 Post-processing method
1: :	if (ه) without ZWNJ exist <b>then</b>
2:	Replace all(ه) without ZWNJ to (ه) including ZWNJ ;
3:	end if
4: :	if (۱۵۵) exist then
5:	Replace all (۱۵۵۱);
6:	end if
7: 3	if ( ٥٥٥) exist then
8:	Replace all (۵۵۵) to (۵۵۵);
9:	end if
10: 3	if (٥٥) exist then
11:	Replace all (۵۵) to (هـ);
12:	end if
13: 3	if 10 exist then
14:	Replace all (۵۱) to (۵۱);
15:	end if
16: 3	if (هۆ) exist then
17:	Replace all (هۆ) to (هۆ);
18:	end if
19: 3	if (ەئ) exist <b>then</b>
20:	Replace all (هێ) to (هێ);
21:	end if

- exist then (دی) exist
- 23: Replace all (هى) to (هي);
- 24: end if

# 5. Results

We collected 110 documents in paper format from different town courts in three cities. They are court-registered complaints by civilians about various issues. Documents in the courts are prepared and typed on computer devices by court employees in the desired court format. There are different types of fonts and sizes among the collected documents. We only have limited data set, and therefore, we used 90%, equating to 99 papers, for training the system, and we used 10%, equating to 11 papers for testing and evaluation. A summary of the collected data is presented in Table 9. Table 10 shows the parameters we used in the training process.

#	Type of Data	Amount	Rate %	Total
1	Data utilized for training	99 printed documents	90 %	110 printed documents
2	Data utilized for evaluation	11 printed documents	10 %	
3	Image lines prepared from the documents	522 files	1044 files	
4	Ground Truth transcriptions prepared from documents	522 files		

Table 9. Summary of data utilized for the training and evaluation process.

Table 10. Training Parameters.

Param	eter	Value
Page s	egmentation	Set to 13 default for RTL languages
OCR e	ngine mode	1 (LSTM mode only)
Debug	-Interval	-1
Max-It	eration	10,000
Lang-]	Гуре	RTL
Norm-	Mode	1
Rando	m-Seed	0
OMP-	Thread-Limited	8

#### 5.1. Dataset Preparation

We removed all personally identifiable information (see Figure 2) and image lines with a naming protocol (e.g., image0.png) (see Figure 3). A summary of collected data is presented in Table 9.

Furthermore, we prepared the Ground Truth (GT) Transcription file, which is singleline plain text corresponding to the line images. We saved GT files with 8-bit encoding and named them with the identical name as the line image. We saved with the extension of .gt.txt,(e.g., image0.gt.txt), see Figure 4.

We used Notepad++ text editor for the preparation of the single-line plain text and created 522 image lines and 522 GT files, in total 1044 files.

The other preparation activities were performed as follows. Tesseract engine performed *binarization*. Basically, the process aims at converting color images into black and white equivalents. In our case, the documents were already in black and white. However, Tesseract still performed the process. *Dewarping* is another pre-process activity in data preparation. However, it is required when the images have been taken using different kinds of input devices, and therefore they might have distortions, skews, and other attributes that could negatively affect the result of the recognition. In our case, we scanned the documents manually and all with the same settings and a single scanner. Therefore, we did not find it necessary to apply *dewarping*. Because of a small set of data, we also applied *segmentation, cropping*, and image *rotation* manually.

In the first phase of training/fine-tuning the Arabic Trained data, the language model was not able to handle ZWNJ properly, and therefore, we replaced all "。" without ZWNJ character to "。" with ZWNJ character, to unbind wrong concatenated characters. After that, other replacement methods will be applied subsequently to correct common wrongly concatenated characters, see Algorithm 1.

بەريۆز/دادوەرى دادگاى بارى كەسى
داواكار : المحمد الم
داوالنِّكران : بين ما
رووی داوا
داوالنِک راوی سے درموہ هاوسے دری منے ہے پنے گرینبہ سے تی هاوسے درگیری
ژماره و و و و و و و و و و و و و و و و و و و
پَیْشُهکی (۱۹مسقال) و پاشهکی(نیه)وه لهسهر سهرینی هاوسهری مندالیّکمان ههیه به ناوی
که له محمد الله الله الله دایك بووه وه مالّی هاوسهری بۆ من دانهناوه تاوهكو ئيّستا
ســـهرهرای ئــهوهی چــهندین جــار داوام لێَکـردووه، بۆيــه داوا لــهدادگای بــهرێز دهکــهم کــه
داوالێِڪراو بانگبکرێِت بـۆ دادبينـى و بريـار دەربکرێِت بـه پابەنـد کردنـى بـه دانـانى مـالى
هاوسەرى سەربەخۆ بۆ من و تەواوى خەرجى داواو ماندوبوونى پاريزەرى بخريّتە ئەستۆ.
لەگەل ٚڕێزماندا
بەلگە سەلمێدەرەكان
رۆژى دادىيىتى
داواکار

Figure 2. Scanned court document image sample for data set preparation.

داوالێکراوی سهرهوههاوسهری بریکاردارمه و بریکارارمی گواستۆتهوه به پی ی گرییهستی مارهبرینی ژماره

Figure 3. Line Image sample.

ا داوالتِکراوی سەرەوەھاوسەری بریکاردارمە و بریکارارمی گواستۆتەوە بە پی ی گرتِبەستی مارەبرینی ژمارە

Figure 4. Ground Truth transcription sample.

# 5.2. Result after Training/Fine-Tuning Arabic Trained Data

We trained/fine-tuned the Arabic language model with the targeted dataset and created the *kur\_ara* language model. We used *lstmeval* to evaluate and showed the result after training. We achieved a 32 percentage points difference in character error rate and a 51 percentage points difference in word error rate. The accuracy error rate for *kur\_ara* language model is reported as follows, the character error rate is 2.93, and the word error rate is 13.53. The detailed results are presented in Table 11. As some Kurdish characters do not exist in the Arabic language model, that magnitude of difference was expected.

Language Models	Char Error Rate	Word Error Rate
Arabic.traineddata	35.10 %	65.03 %
kur_ara.traineddata	2.93 %	13.52 %
Improve rate	32.18 %	51.50 %

Table 11. The evaluation result of trained language model with the starter language model.

To evaluate the accuracy of the trained language model by actual court-printed document images, we used *ocreval* tool. A sample is shown in Figure 5.

		يني ژماره	، مارەبري	كريبهستى	به به پی ی	داواليكراوى سەرەوە زنى بريكاردارە
وه پاشهکی	مسقال زير	پێؚۺ۬؋کی( سی	مارەيى	ووه ئەسەر	دەرچ	که له دادگای باریی کهسی
دادگا طلاق	ئەدەرەوەى	)داواليّكراوى	2014	/12/7)	لەرىككەوتى	بیست مسقال زیّر) بریکاردارم
و داد بيني وه	نگ بکریْت ب	داوا ليْكراو بان	دكەم كە	له دادگا د	به بۆيە داوا	اوەئەسەر رەزامەندى داوائيكراو بۆ
هرجی داوا و	د د د د د د د د د د د د د د د د د د د	م داواليّكراو و	ريكاردار	ی نيوان ب	للاقى دەرەك	ېريار دەربكريّت بەپەسەندكرنى م
				اليكراو	ئەستۆى داوا	ماندوبوونی پارێزەريە تی بخرێته
				ماندا	لمكمان ف	

Figure 5. Actual court document image.

We OCRed the image with *kur\_ara* language model, see Figure 6, and compared this with the original text file, see Figure 7.

داوالتِکراوی سهرهوه ژنی بریکاردارمه به پێ ی کرزبهسنی ماردهبرینی ژ که له دادگای باربی کهسی دهرچووه لهسهر مارهیی پیشهکی(سی مسقال زبّر) وه پاشهکی (بیست مسقال زبّر) بریکاردارم لهریّککهوتی (2014/12/7) داوالیّکراوی لهدهرهوهی دادگا طلاق داوهلهسهر رهزامهندی داوالیّکراو بۆیه بۆیه داوا له دادگا دهکهم که داوا لیّکراو بانگ بکریّت بۆ داد پینی وه بریار دهربکریّت به پهسهندکرنی طلاقی دهرهکی نیّوان بریکاردارم داوالیّکراو وهتهواوی خهرجی داوا و ماندوبوونی پارتزهریه تی بخریّته ئهستۆی داوالیّکراو

لهگهڵ رێزماندا

Figure 6. OCRed text of actual Court's document image.

داوالتِکراوی سهرەوه ژنی بریکاردارمه به پێ ی گرێبهستی ماردەپرینی ژماره که له دادگای باریی کهسی دەرچووه لهسهر مارەیی پێشهکی(سی مسقال زێر) وه پاشهکی (بیست مسقال زێر) بریکاردارم لهریککهوتی (2014/12/ )داوالیَکراوی لهدەرەوەی دادگا طلاق داوەلهسهر رەزامەندی داوالیَکراو بۆیه بۆیه داوا له دادگا دەکەم که داوا لیکراو بانگ بکریّت بۆ داد پینی وه بپریار دەربکریّت به پەسەندکرنی طلاقی دەرەکی نیّوان بریکاردارم داوالیَکراو وەتەواوی خەرجی داوا و ماندوبوونی پاریّزەریه تی بخریّته ئەستۆی داوالیکراو لهگهڵ ریزماندا

Figure 7. Original text of actual Court's document image.

The result showed an accuracy rate of 87.5%. Most of the errors that occurred are related to ZWNJ. The detailed report is shown in Figure 8. The result for the same document image after applying the post-processing method for text correction showed 97.75% accuracy.

ocreval	Accurac	cy Report	Vers	sion 7.0
488	Chara	acters		
61	Errors			
87.50	% Acc	uracy		
0	Reject	Characte	ers	
0	Suspe	ct Marker	S	
0	False I	Marks		
0.00	% Char	acters M	arked	I
87.50	% Acc	uracy Afte	er Co	rrection
Ins	Subst	Del	Erro	rs
0	0	0	0	Marked
5	55	1	61	Unmarked
5	55	1	61	Total
Coun	t Misse	ed %Rig	ght	
72	1	98.61	ASC	II Spacing Characters
8	0	100.00	ASC	II Special Symbols
7	0	100.00	ASC	II Digits
250	4	98.40	Bas	ic Arabic
150	55	63.33	Ara	bic Extended
1	0	100.00	Gene	eral Punctuation
488	60	87.70	Tota	al
Errors	Marke	d Corre	ct-Ge	nerated
51	0	{a}-{a}		
4	0	{ک-{مارہ}		
1	0	{}-{}		
1	0	{}-{<\n>}		
1	0	{;;}-{¤}		
1	0	رر {گ}-{گ}		
1	0	{ {}-{}		
1	0	{j}-{ï}		
	-			

Figure 8. Detailed report on actual Court's document image.

5.3. Result after Training/Fine-Tuning for Adding New Font

We Trained/Fine-tuned the kur\_ara language model for adding new fonts. We used synthetic data, which includes AWN and AEN, and trained the model by utilizing 20 pages of approximately 1000 lines for each Unikurd font. We prepared those pages using a script that compiled the pages for each font based on the text at https://github.com/Shreeshrii/tesstrain-ckb/blob/master/langdata/ckb.training\_text (accessed on 7 August 2021). We selected the 20 pages data from beginning of the document onward. Figure 9 shows a sample of the data that we prepared for Unikurd\_Tishk font.

ئەچ مانگە ئەنجام دەدرىّت 15 : 57 2012 / 01 ھەۋلار 1 مارس/ئازار – يەكىّتى تۆپى پى ئاسيا رايگەياندوو
بۆ ئاو كۆنگرە 17 : 39 2011 / / 02 سلاّمانى 4تسرينى دووەم/ ئۆكتۈبەر( )— بزيارە ٤ى ئەم مانگە بە بەژدارى
فيستيڤاٽي فيلمي "مۆنتراڭ" بەدەست بێنێت. بە گوێرەي ھەواٽي ئاژانسى ھەواٽنێرى "فارس".
ھەولۆر _دوۆنى دووشەممە ھەڭبۋاردنەگانى يانەى وەرزشى ھەندرۆن ئەنمامدرا بە مەبەستى ديارى

Figure 9. A Sample of Synthetic Data for Unikurd\_Tishk font.

We evaluated and showed the result of the base *kur\_ara* language model versus the fine-tuned version in Table 12. We started by training one font at a time and then repeated the process of training over ten fonts. The average error rate for characters before training/fine-tuning over ten fonts is 10.34%, while it is 22.61% for words. It is 1.60% for characters after training/fine-tuning, while it is 4.83% for words. That shows an improvement of 8.74 and 17.78 percentage points for characters and words, respectively. Table 12 presents a more detailed result for each font.

Table 12. The evaluation result of training/fine-tuning for adding new fonts.

	Unikurd	Evaluation Before Train	ing	Evaluation After Trainin	ıg
*2 —- #	Font Name	Char Err Rate %	Word Err Rate %	Char Err Rate %	Word Err Rate %
1	Jino	30.76	58.01	2.51	6.81
2	Web	6.15	15.51	1.35	4.78
3	Chimen	12.04	21.11	3.02	6.41
4	Hiwa	2.37	6.79	1.35	4.08
5	Goran	2.10	6.17	1.31	4.21
6	Kawe	11.94	33.65	1.35	4.70
7	Hejar	2.39	8.91	0.91	3.62
8	Hemen	11.88	30.22	2.05	6.12
9	Nali	5.52	16.83	0.73	3.52
10	Tishk	18.25	28.96	1.38	4.09
Average		10.34	22.61	1.60	4.83

We used the final *kur\_ara* language model which was trained/fine tuned over ten fonts for recognizing actual court's documents images, see Figures 10 and 11.

The average accuracy rate for 11 documents reported as 95.45818182%. See details in Table 13.

#	<b>Court's Documents</b>	Accuracy Rate %
1	Image 1	95.88
2	Image 2	93.54
3	Image 3	93.20
4	Image 4	97.06
5	Image 5	96.34
6	Image 6	92.09
7	Image 7	97.24
8	Image 8	97.77
9	Image 9	92.83
10	Image 10	97.57
11	Image 11	96.52
Average ac	ccuracy rate	95.46

Table 13. Accuracy rate for court's document images using ocreval.

	باری دهسی	بەرير دادوەرى دادكاى	
	ناونیشان :	: دشيو	اواكار
	ناونیشان:	پیشه : ژنی مال	اوا ليْكراو :
			رووي داوا :
له ۲۰/۱۰/۳۰ که له دادگای	رینی ژماره که	منه به پیٰ ی گریٰیهستی ماره	اواليکراوی سەرەوە ژنی
ه پاشهکی ( ۱۹ مسقال زیّر ) وه	۳۰۰سیٰ همزار دینار ) و	ووه لهسهر مارهیی پیّشهکی (۰	باریی کهسی دەرچ
مرمر به من دمگهينيٽ به جنيو	وى سەرەوە بەردەوام زە	، سیٰ مندالیان هەیه داوالیّکرار	لەسەر سەرينى ھاوسەرى
و لهمندالهکانم دمدات و وه	دموام جگەرە دەكىْشىْت	برووم وه سهرمرای ئهومش بهره	دان و ئەكە داركردنى ئاب
یشی نەخۆشیەکی دەرونی بوه	ۆكوشتنى داوەو وە تو	چەندىن جارى تر ھەولى خ	جارێك خۆى سوتاندوەو
هشێوهیهکن که ناتوانم چیتر	بلينت وه ئهو زمرمرانه ب	دالهكانى بەتەنيا بەجى دەھ	ومشهو لممال رادمكات من
م که داوا ٹیکراو بانگ بکریّت	يه داوا له دادگا دمكه	وام بم لەژيانى ھاوسەرى بۆ	نەگەل داوا ئيّكراو بەردە
سهر بنهمای زمرمر ( التفریق	ان من و داواليّكراو له	ربكريّت به جيابوونهوه له نيّو	ېۆ داد بينى وه بريار ده
ی خەرجى داوا و ماندوبوونى	س عيْراقي وه تهواوي	.می ( ٤٠ ) له یاسای باری کهس	نلضرر) پالپشت به ماده
		ستۆى داوائىكراو	پارێزەريەتى بخرێتە ئە
		ئەگەڻ رێزماندا	
			بەلگە سەلمىنەر
			ۆژى دادېينى
	داماكا		
2			

Figure 10. Court document image for evaluation by final model (Image1).

		باری کهسی	بەرێز⁄دادومرى دادگاى	
	باريزمر	۔ انڈنی مال بریکارہکہی ہ	ييشه	اواکاران ۱_
			. 4	
		: فوتابی	پيشه	
		ال ناونیشان :	پیشه :ژنی م	اواليّكراو/
		زچکردو	سەرەراي ميراتى كۆ	
				ووی داوا
) ئەريككــەوتى	بووه بهناوی (	م ومباوکی داواکاری دوومم	مەرەوە مىردى بريكاردارە	کوری داوائیکراوی س
ززده مسقال ونيو زير)	، پێشـ ۵۸ ، ۱۹، نو	ئه نجام داوه لهسهرمارهيي	يكاردارم هاوسهرگيريان	۱۹۹۰/۳/۱)ئەگەل بر
هاوســهرىدوو منــداليان	روهو و ئەسەرســەرينى ه	۲۰۰۲ )كۆچىي دوايىي كىردو	ی برکاردەرم ئــه( ۲/۱۱/۹	باشهكى ( نيه )ومميّرد
۱۹۹۸که کچه )که کـوردن	ه ( )مواليد ٥/٨/	کهداواکاری دووممه وکچهو	مواليد ١٩٩٣/٧/١	مەيەبەناوەكانى (
انگبکرێت بــۆدادبينى و	، ميراتى كۆچىكردوو ب	اواليكراوى سەرەوەسەرەراي	ۆيەداوائەدگادەكەم كەدا	موسولان و عراقين ب
الى رەچــەلەكى ھــەردوو	راوى كۆچىكردوودانە پ	بريكادەرم و كورى داوائيّك	لماندنى هاوسەرى نێوان	ريار دەربكريّت بەسە
ــتانی بــه پِێِی تــاریخی	ۆمارەكانى بــارى شارس	سەرەوە وەتۆماركردنى ئەت	ارم و کوری داوالیّکراوی ،	ندالەكان بۆبريكارد
راوارى وەتەواوى خەرجى	بنار دەخەمليْنم داواليْكم	ردمرم به ( ۱۰۰۰۰ )دمههزاردی	وه وهبهشهمیراتی بریکار	اماژه پيکراو لهسهره
		اليُكراو.	زەريىتى بخرئەستۆى داوا	اواو ماندوبونی پارێ
		يْزماندا	ئەڭەن ر	
اکار (۱)	ٹہجیاتے داہ		داماکار (۲)	
	بریکارهکهی :پاریزمر			
	یکارنامهی گشتی ژماره	به پێی بر		

Figure 11. Court document image for evaluation by final model (Image2).

#### 5.4. Discussion

The research showed that using a training dataset that has a script similarity with a less-resourced language could increase the overall accuracy of the LSTM-based Tesseract engine. We also showed that two methods could be used to improve the output of the engine: using post-processing and fine-tuning the engine using different fonts. Although both provide similar improvement, the latter is more straightforward and preferable. However, using the latter approach depends on the fonts available for the target language. Highly accurate character recognition of the created model (*kur\_ara*) is currently limited to documents written in ten *Unikurd* fonts, see Table 12. The approach can be applied to other Kurdish Unicode fonts by using synthetic data as well. Currently, the developed language model is restricted in recognizing only machine-typed documents.

The advantage of our approach is in the size of the data. A small amount of data could be used to train a Tesseract model that has been trained over an appropriate dataset for a well-resourced language with a similar script. The result language model of that process provides a highly accurate output. However, this research did not intend to investigate documents that include tables and test older documents typed using mechanical type machines and other older typing tools.

## 6. Conclusions

The unavailability of OCR systems is an obstacle that hinders the text processing task in under-resourced languages. Tesseract is an engine that provides a proper basis to develop OCR models for various languages. However, its LSTM-based approach requires large amounts of data for their training. This research suggested a method to overcome this issue. For our experiment, we used Kurdish (Sorani), which is considered an underresourced language. The study concluded that a small dataset for a target language along with a larger dataset for a second language that has script similarity with the target language could compensate for the requirement of a large amount of data in LSTM-based applications. The study showed that the resulted model could provide a reasonable result even in the absence of a large dataset. We also showed that, in the absence of fine-tuning, the output could be improved by post-processing. However, if fine-tuning using different fonts is possible, it could eliminate the post-processing phase.

The conclusion suggests three main areas for future work. First, the expansion of the training data set to cover more complex documents, for instance, documents with tables, newspapers, multicolumn documents, to name a few. That could be combined with an appropriate set from different types of Tesseract segmentation to extend the coverage of the adapted engine. Second, working on historical documents from the pre-digital era is paramount that could help many other sectors of Kurdish processing, particularly corpora development. Third, to investigate the applicability and efficiency of the suggested approach for OCR systems in other fields that use LSTM-based methods such as energy consumption, commercial vacancy prediction, and time-series forecasting.

Author Contributions: Conceptualization, H.H. and S.I.; Methodology, H.H. and S.I.; Software, S.I.; Validation, S.I.; Formal analysis H.H. and S.I.; Investigation, S.I. and H.H.; Resources, S.I.; Data curation, S.I.; Writing—original draft preparation, S.I.; Writing—review and editing, H.H.; Visualization, H.H. and S.I.; Supervision, H.H.; Project administration, H.H. Revision: H.H. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

**Data Availability Statement:** The dataset will be publicly available on Kurdish-BLARK under the GPL upon the acceptance of the paper.

**Acknowledgments:** We appreciate the feedback we received from the anonymous reviewers. They were constructive, professional, valuable, and helpful. We are grateful to the reviewers for their opinion that significantly improved the paper. We are also grateful to Lesley A T Gaj for her generous assistance in proofreading the manuscript.

Conflicts of Interest: The authors declare no conflict of interest.

#### Abbreviations

The following abbreviations are used in this manuscript:

AEN	Arabic Eastern Number
ANN	Artificial Neural network
AWN	Arabic Western Number
BLARK	Basic Language Resource Kit
DPI	Dots Per Inch
FOA	Fly Optimization Algorithm with LSTM
FOA-LSTM	Fly Optimization Algorithm with LSTM

GB	GigaByte
GIF	Graphics Interchange Format
GISMO	Geographic Information Systems and Mapping Operations
GNU	GNU's Not Unix
GUI	Graphical User Interface
GT	Ground Truth
HMM	Hidden Markov Model
ISRI	Information Science Research Institute
JOCR	Jorg Optical Character Recognition
JPEG	Joint Photographic Experts Group
KNN	K Nearest Neighbors
KRI	Kurdistan Region of Iraq
KRG	Kurdistan Regional Government
LR	language resources
LSTM	Long Short-Term Memory
LTR	Left To Right
LVQ	Learning Vector Quantization
OCR	Optical Character Recognition
OEM	OCR Engine Mode
PBM	Portable Bitmap
PGM	Portable Gray Map
PNG	Portable Network Graphics
PNM	Portable Any Map
PP	Pocket Physics
PPP	Public-Private Partnership
PSM	Page Segmentation Mode
RNN	Recurrent Neural Network
RTL	Right To Left
SDK	Software Development Kit
TIFF	Tagged Image File Format
UTF-8	8-bit Unicode Transformation Format
XML	Extensible Markup Language
ZWNJ	Zero-Width Non-Joiner

## References

- 1. Hassani, H. BLARK for Multi-dialect Languages: Towards The Kurdish BLARK. Lang. Resour. Eval. 2018, 52, 625–644. [CrossRef]
- Peng, L.; Zhu, Q.; Lv, S.X.; Wang, L. Effective long short-term memory with fruit fly optimization algorithm for time series forecasting. *Soft Comput.* 2020, 24, 15059–15079. [CrossRef]
- 3. Peng, L.; Wang, L.; Xia, D.; Gao, Q. Effective energy consumption forecasting using empirical wavelet transform and long short-term memory. *Energy* **2021**, 238, 121756. [CrossRef]
- 4. Lee, J.; Kim, H.; Kim, H. Commercial Vacancy Prediction Using LSTM Neural Networks. Sustainability 2021, 13, 5400. [CrossRef]
- Esmaili, K.S.; Salavati, S. Sorani Kurdish Versus Kurmanji Kurdish: An Empirical Comparison. In Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Sofia, Bulgaria, 4–9 August 2013; pp. 300–305.
- 6. Hassani, H.; Medjedovic, D. Automatic Kurdish Dialects Identification. Comput. Sci. Inf. Technol. 2016, 6, 61–78.
- 7. Esmaili, K.S. Challenges In Kurdish Text Processing. arXiv 2012, arXiv:1212.0074.
- 8. Ahmadi, S. Why Does Kurdish Language Processing Matter? 2019. Available online: https://sinaahmadi.github.io/posts/why-kurdish-language-processing-matters.html (accessed on 2 September 2021).
- 9. Marouf, M. Kurdish Academia Journal NO. 16. 2015. Available online: https://govkrd.b-cdn.net/OtherEntities/Kurdish%20 Academy/Kurdish/%D8%A8%DA%B5%D8%A7%D9%88%DA%A9%D8%B1%D8%A7%D9%88%DB%95%DA%A9%D8%A7 %D9%86/%DA%AF%DB%86%DA%A4%D8%A7%D8%B1%DB%8C%20%D9%8A%D9%94%DB%95%DA%A9%D8%A7%D8 %AF%DB%8C%D9%85%DB%8C%D8%A7/Govari%20Ekadimi%2016.pdf (accessed on 5 August 2021)
- 10. Hashemi, D. Kurdish Orthography. 2016. Available online: http://yageyziman.com/Renusi\_Kurdi.htm (accessed on 2 September 2021).
- Tafti, A.P.; Baghaie, A.; Assefi, M.; Arabnia, H.R.; Yu, Z.; Peissig, P. OCR as a Service: An Experimental Evaluation of Google Docs OCR, Tesseract, ABBYY FineReader, Furthermore, Transym. In *Advances in Visual Computing*. *ISVC 2016*; Lecture Notes in Computer Science; Bebis, G., Ed.; Springer: Cham, Switzerland, 2016; Volume 10072.
- 12. Mithe, R.; Indalkar, S.; Divekar, N. Optical Character Recognition. Int. J. Recent Technol. Eng. 2013, 2, 72–75.

- 13. Herbert, H. *The History of OCR, Optical Character Recognition*; Recognition Technologies Users Association: Bennington County, VT, USA, 1982.
- 14. Dhavale, S.V. Advanced Image-Based Spam Detection and Filtering Techniques; IGI Global: Hershey, PA, USA, 2017.
- 15. d'Albe, E.F. On A Type-reading Optophone. Proc. R. Soc. Lond. Ser. A Contain. Pap. A Math. Phys. Character 1914, 90, 373–375.
- 16. Yaseen, R.; Hassani, H. Kurdish Optical Character Recognition. UKH J. Sci. Eng. 2018, 2, 18–27. [CrossRef]
- 17. Gamera. The Gamera Project. Available online: https://gamera.informatik.hsnr.de/ (accessed on 7 August 2021)
- Radhiah, A.; Machbub, C.; Hidayat, E.M.I.; Prihatmanto, A.S. Printed Arabic Letter Recognition Based On Image. In Proceedings of the 2018 International Conference on Signals and Systems (ICSigSys), Bali, Indonesia, 1–3 May 2018; pp. 86–91.
- Nashwan, F.; Rashwan, M.A.; Al-Barhamtoshy, H.M.; Abdou, S.M.; Moussa, A.M. A Holistic Technique For An Arabic OCR System. J. Imaging 2018, 4, 6. [CrossRef]
- 20. Husnain, M.; Saad Missen, M.M.; Mumtaz, S.; Jhanidr, M.Z.; Coustaty, M.; Muzzamil Luqman, M.; Ogier, J.M.; Sang Choi, G. Recognition of Urdu Handwritten Characters Using Convolutional Neural Network. *Appl. Sci.* **2019**, *9*, 2758. [CrossRef]
- 21. Naz, S.; Hayat, K.; Razzak, M.I.; Anwar, M.W.; Madani, S.A.; Khan, S.U. The Optical Character Recognition Of Urdu-like Cursive Scripts. *Pattern Recognit.* 2014, 47, 1229–1248. [CrossRef]
- Izakian, H.; Monadjemi, S.; Ladani, B.T.; Zamanifar, K. Multi-font Farsi/Arabic Isolated Character Recognition Using Chain Codes. World Acad. Sci. Eng. Technol. 2008, 43, 67–70.
- Jelodar, M.S.; Fadaeieslam, M.J.; Mozayani, N.; Fazeli, M. A Persian OCR System Using Morphological Operators. In Proceedings of the World Academy of Scienc, Engineering and Technology, Istanbul, Turkey, 25–27 February 2005; Volume 2, pp. 137–140.
- 24. Smith, R. An Overview Of The Tesseract OCR Engine. In Proceedings of the Ninth International Conference on Document Analysis and Recognition (ICDAR 2007), Curitiba, Brazil, 23–26 September 2007; Volume 2, pp. 629–633.
- 25. Smith, R. Motivation and History of the Tesseract OCR Engine; Google Inc.: Menlo Park, CA, USA, 2020.
- Dalitz, C. A Tutorial Introduction to the Gamera Framework. 2009. Available online: https://gamera.informatik.hsnr.de/docs/ gamera-tutorial.pdf (accessed on 7 August 2021).
- 27. Genzel, D.; Ashok Popat, D.N. Paper to Digital in 200 Languages. 2015. Available online: https://ai.googleblog.com/2015/05/ paper-to-digital-in-200-languages.html (accessed on 5 August 2021).
- 28. Google Drive Help. Convert PDF and Photo Files to Text. 2020. Available online: https://support.google.com/drive/answer/17 6692?hl=en&co=GENIE.Platform%3DDesktop (accessed on 5 August 2021).
- 29. Dhiman, S.; Singh, A. Tesseract vs. GOCR A Comparative Study. Int. J. Recent Technol. Eng. 2013, 2, 80.
- 30. Jain, P.; Taneja, K.; Taneja, H. Which OCR toolset is good and why: A comparative study. Kuwait J. Sci. 2021, 48. [CrossRef]
- Zarro, R.D.; Anwer, M.A. Recognition-based Online Kurdish Character Recognition Using Hidden Markov Model Furthermore, Harmony Search. Eng. Sci. Technol. Int. J. 2017, 20, 783–794.
- Ahmed, R.M. Kurdish Handwritten Character Recognition Using Deep Learning Techniques. Master's Thesis, University of Kurdistan Hewlêr, Erbil, Iraq, 2019.
- Sinha, A.; Jenckel, M.; Bukhari, S.S.; Dengel, A. Unsupervised OCR Model Evaluation Using GAN. In Proceedings of the 2019 International Conference on Document Analysis and Recognition (ICDAR), Sydney, NSW, Australia, 20–25 September 2019; pp. 1256–1261.
- 34. Martínek, J.; Lenc, L.; Král, P. Building an efficient OCR system for historical documents with little training data. *Neural Comput. Appl.* **2020**, *32*, 17209–17227. [CrossRef]
- Hula, J.; Mojžíšek, D.; Adamczyk, D.; Čech, R. Acquiring Custom OCR System with Minimal Manual Annotation. In Proceedings of the 2020 IEEE Third International Conference on Data Stream Mining & Processing (DSMP), Lviv, Ukraine, 21–25 August 2020; pp. 231–236.
- 36. Kiss, M.; Benes, K.; Hradis, M. AT-ST: Self-Training Adaptation Strategy for OCR in Domains with Limited Transcriptions. *arXiv* 2021, arXiv:2104.13037.
- 37. Google. Tesseract Documentation. 2020. Available online: https://tesseract-ocr.github.io/ (accessed on 5 August 2021).
- Patki, N.; Wedge, R.; Veeramachaneni, K. The Synthetic Data Vault. In Proceedings of the 2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA), Montreal, QC, Canada, 17–19 October 2016; pp. 399–410.
- Santos, E.A. OCR Evaluation Tools for the 21st Century. In Proceedings of the 3rd Workshop on the Use of Computational Methods in the Study of Endangered Languages Volume 1 (Papers), 3rd Workshop on Computational Methods for Endangered Languages, Honolulu, HI, USA, 26–27 February 2019; pp. 23–27.