

# An attention based method for offline handwritten Urdu text recognition

Tayaba Anjum and **Nazar Khan**

Computer Vision and Machine Learning Group,  
Punjab University College of Information Technology (PUCIT)  
University of the Punjab, Lahore, Pakistan

17th International Conference on Frontiers in Handwriting Recognition  
September 9, 2020



# **Problem: Offline Handwritten Text Recognition (OHTR) for Urdu**

# Offline Handwritten Text Recognition (OHTR)

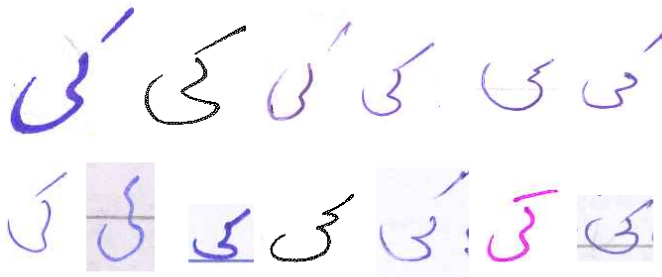
Automatic recognition of text from images of handwritten text.

کراچی پاکستان کا سب سے بڑا شہر اور صنعتی، تجارتی، تعلیمی، مواصلاتی و اقتصادی مرکز ہے۔ کراچی دنیا کا دوسرا بڑا شہر ہے۔ کراچی پاکستان کے صوبہ سندھ کا دارالحکومت ہے۔ شہر دریائے سندھ کے مغرب میں بحیرہ عرب کی شمالی ساحل پر واقع ہے۔ پاکستان کی سب سے بڑی بندرگاہ اور ہوائی اڈہ بھی کراچی میں قائم ہے۔ کراچی ۱۹۴۷ء سے ۱۹۶۵ء تک پاکستان کا دارالحکومت بھی رہا۔ موجودہ کراچی کی جگہ پر واقع قدیم ماہی گروں کی بستیوں میں سے ایک مانا کو لائی جو کوٹ تھا۔ انگریزوں نے انیسویں صدی میں اس شہر کی تعمیر و ترقی کی۔ بینظیر دہلوی ۱۹۴۷ء میں پاکستان کی آزادی

Urdu is a right-to-left language written in a mixture of mostly cursive and occasionally non-cursive form.

# Why is handwritten Urdu recognition difficult?

- Sayre's paradox: cursively written word cannot be recognized without being segmented and cannot be segmented without being recognized.
- Wide range of variations in pen type, writing style, writing size, page background.
- Almost every rule can be violated until text becomes illegible.



**Figure:** An Urdu word pronounced as “key” written by 13 different people.

# Handwritten vs. Typed Urdu vs. Characters

دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔

دیکھتا جاتا تھا اور **مسکراتا** اور خوش ہوتا تھا

دی کہت اجات اتھا اور **مسکراتا** اور خوش ہوتا تھا

- Within cursive words, characters to be recognized often appear with widely varying ligatures.
- While typed Urdu conforms to some rules, handwriting can be restricted to very few rules.

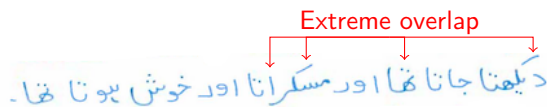
# Challenges of Urdu OHTR

Extreme overlap

دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔

# Challenges of Urdu OHTR

Extreme overlap



دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔

The diagram illustrates the challenge of extreme overlap in Urdu text. A red line with four arrows points to the overlapping characters 'ا' and 'ت' in the words 'تھا' and 'اور' of the sentence 'دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔'.

Context dependent shape and location

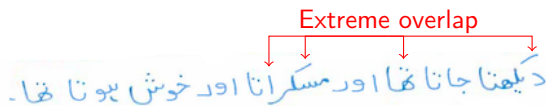


دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔

The diagram illustrates the challenge of context dependent shape and location in Urdu text. A blue line with two arrows points to the character 'ا' in the words 'تھا' and 'اور' of the sentence 'دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔', showing how its shape and position vary based on the surrounding context.

# Challenges of Urdu OHTR

Extreme overlap



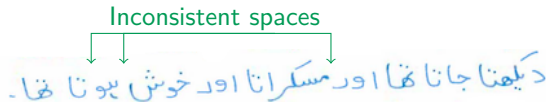
دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔

Context dependent shape and location



دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔

Inconsistent spaces



دیکھتا جاتا تھا اور مسکراتا اور خوش ہوتا تھا۔

# State of OHTR for Arabic-like scripts

## Techniques

- Urdu
  - Raw pixels + Bidirectional LSTM<sup>1</sup>
  - CNN + Bidirectional LSTM<sup>2</sup>
- Arabic
  - CNN + LSTM with context-window<sup>3</sup>
  - Raw pixels + Multidimensional LSTM<sup>4</sup>
- Similar situation for Farsi.
- The concept of **attention has not been explored**.
- Error reporting has been restricted to the level of characters. **World level accuracy is not reported.**

---

<sup>1</sup>Saad Bin Ahmed et al. "Handwritten Urdu character recognition using one-dimensional BLSTM classifier". In: *Neural Computing and Applications* 31.4 (2019), pp. 1143–1151.

<sup>2</sup>S. Hassan et al. "Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting". In: *International Conference on Deep Learning and Machine Learning in Emerging Applications, Deep-ML 2019*. 2019.

<sup>3</sup>Gui et al. "Adaptive context-aware reinforced agent for handwritten text recognition". In: *Proceedings of the British Machine Vision Conference (BMVC)*. 2018.

<sup>4</sup>Riaz Ahmad et al. "KHATT: A deep learning benchmark on Arabic script". In: *14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2017. doi: 10.1109/ICDAR.2017.358.

# State of OHTR for Arabic-like scripts

## Datasets

- Urdu

- CENIP-UCCP<sup>5</sup>
- UCOM/UNHD<sup>6</sup>
- Hassan et al.<sup>7</sup>

- Arabic

- IFN/ENIT<sup>8</sup>
- OpenHaRT<sup>9</sup>
- KHATT<sup>10</sup>

- Farsi

- Sadri et al.<sup>11</sup>

- Urdu datasets are **not entirely accessible**.

<sup>5</sup>Raza et al., "An Unconstrained Benchmark Urdu Handwritten Sentence Database with Automatic Line Segmentation".

<sup>6</sup>Ahmed et al., "Handwritten Urdu character recognition using one-dimensional BLSTM classifier".

<sup>7</sup>Hassan et al., "Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting".

<sup>8</sup>Pechwitz et al., "IFN/ENIT-database of handwritten Arabic words".

<sup>9</sup>Tong et al., "NIST 2013 open handwriting recognition and translation evaluation".

<sup>10</sup>Mahmoud et al., "KHATT: An open Arabic offline handwritten text database".

<sup>11</sup>Sadri, Yeganehzad, and Saghi, "A novel comprehensive database for offline Persian handwriting recognition".

## **Contributions**

- ① We present the first *attention-based*, encoder-decoder model for recognition of offline handwritten Urdu text.

# Paper Contributions

- ① We present the first *attention-based*, encoder-decoder model for recognition of offline handwritten Urdu text.
- ② We present a new dataset of offline handwritten Urdu text containing 7,309 unique text line images with ground-truth.

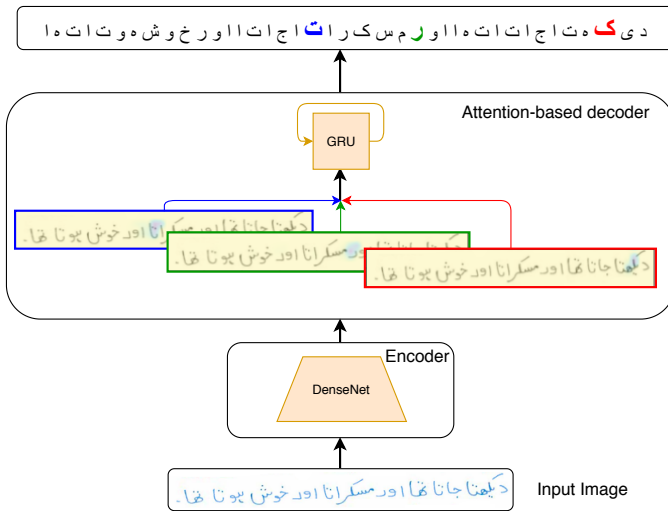
- ① We present the first *attention-based*, encoder-decoder model for recognition of offline handwritten Urdu text.
- ② We present a new dataset of offline handwritten Urdu text containing 7,309 unique text line images with ground-truth.
- ③ For the first time, we report word level accuracy instead of only character level accuracy.

- ① We present the first *attention-based*, encoder-decoder model for recognition of offline handwritten Urdu text.
- ② We present a new dataset of offline handwritten Urdu text containing 7,309 unique text line images with ground-truth.
- ③ For the first time, we report word level accuracy instead of only character level accuracy.
- ④ We show that attention can enable uni-directional decoders to out-perform bidirectional decoders.

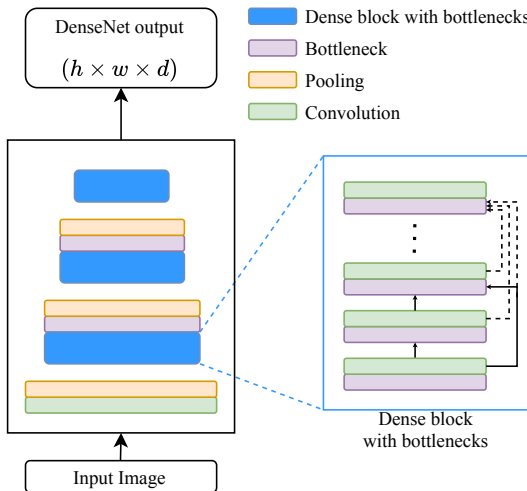
- ① We present the first *attention-based*, encoder-decoder model for recognition of offline handwritten Urdu text.
- ② We present a new dataset of offline handwritten Urdu text containing 7,309 unique text line images with ground-truth.
- ③ For the first time, we report word level accuracy instead of only character level accuracy.
- ④ We show that attention can enable uni-directional decoders to out-perform bidirectional decoders.
- ⑤ Compared to previous approaches, we report close to 2x accuracy improvement at character level and 37x at word level.

## **Proposed Solution**

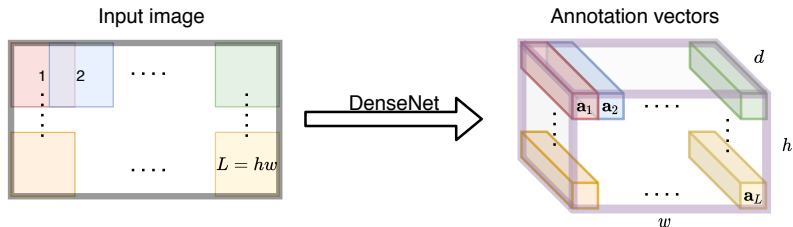
# Encoder/Decoder Model



# Encoder: DenseNet with bottlenecks



# Interpretation of DenseNet Output



Output volume of size  $h \times w \times d$  represents  $d$ -dimensional *annotation vectors* of  $h \times w$  overlapping blocks of the input image.

# Decoder: Gated Recurrent Unit (GRU)

Produces text sequence one character at a time conditioned on the previously generated character  $\mathbf{y}_{t-1}$  and current hidden state  $\mathbf{h}_t$ .

$$\mathbf{z}_t = \sigma(\mathbf{W}_{yz}\mathbf{E}\mathbf{y}_{t-1} + \mathbf{U}_{hz}\mathbf{h}_{t-1})$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{yr}\mathbf{E}\mathbf{y}_{t-1} + \mathbf{U}_{hr}\mathbf{h}_{t-1})$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{yh}\mathbf{E}\mathbf{y}_{t-1} + \mathbf{U}_{rh}(\mathbf{r}_t \otimes \mathbf{h}_{t-1}))$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \otimes \mathbf{h}_{t-1} + \mathbf{z}_t \otimes \tilde{\mathbf{h}}_t$$

where update gate  $\mathbf{z}_t$  and reset gate  $\mathbf{r}_t$  and candidate hidden state  $\tilde{\mathbf{h}}_t$  are used to compute the hidden state and GRU output  $\mathbf{h}_t$ .

# Decoder: Gated Recurrent Unit (GRU) with Attention

Produces text sequence one character at a time conditioned on the previously generated character  $\mathbf{y}_{t-1}$ , current hidden state  $\mathbf{h}_t$  and a context vector  $\mathbf{c}_t$ .

$$\mathbf{z}_t = \sigma(\mathbf{W}_{yz}\mathbf{E}\mathbf{y}_{t-1} + \mathbf{U}_{hz}\mathbf{h}_{t-1} + \mathbf{C}_{cz}\mathbf{c}_t)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{yr}\mathbf{E}\mathbf{y}_{t-1} + \mathbf{U}_{hr}\mathbf{h}_{t-1} + \mathbf{C}_{cr}\mathbf{c}_t)$$

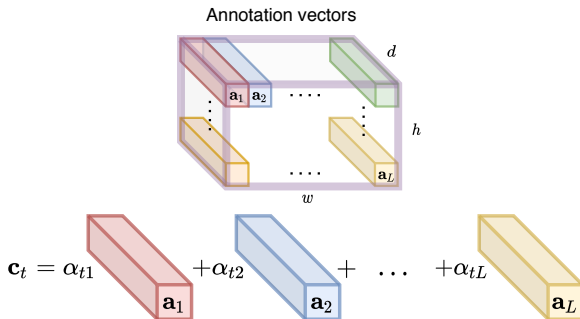
$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{yh}\mathbf{E}\mathbf{y}_{t-1} + \mathbf{U}_{rh}(\mathbf{r}_t \otimes \mathbf{h}_{t-1}) + \mathbf{C}_{ch}\mathbf{c}_t)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \otimes \mathbf{h}_{t-1} + \mathbf{z}_t \otimes \tilde{\mathbf{h}}_t$$

where update gate  $\mathbf{z}_t$  and reset gate  $\mathbf{r}_t$  and candidate hidden state  $\tilde{\mathbf{h}}_t$  are used to compute the hidden state and GRU output  $\mathbf{h}_t$ .

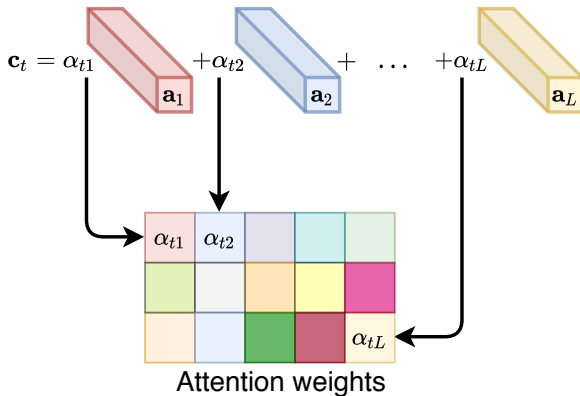
# Attention via the context vector

The context vector  $\mathbf{c}_t$  is a dynamic representation of the relevant part of the image at time  $t$ . Computed as a weighted sum of annotation vectors

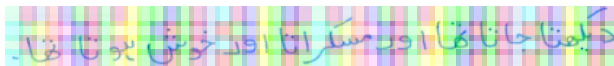


# Attention via the context vector

Attention weight  $\alpha_{ti}$  determines importance of location  $i$  in determining context vector  $\mathbf{c}_t$ .



# Computing Attention



## Standard model

Attention weights  $\alpha_{ti}$  computed via softmax over locations

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{j=1}^L \exp(e_{tj})}$$

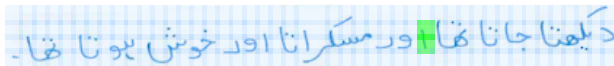
Per-location exponents computed via

$$e_{ti} = \mathbf{v}_a^T \tanh(\mathbf{W}_a \mathbf{h}_{t-1} + \mathbf{U}_a \mathbf{a}_i)$$

depend on

- the hidden state, and
- each location's content encoded in the annotation vector

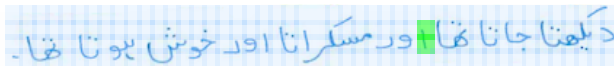
# Computing Attention



دیکھتا جاتا تھا

- Ideally, we would like our text recognizer to “read” like we do.
- It should focus on the relevant image region when recognizing each character.

# Computing Attention



دیکھتا جاتا تھا۔

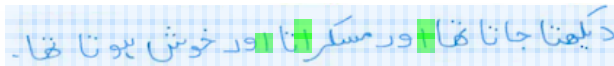
- Ideally, we would like our text recognizer to “read” like we do.
- It should focus on the relevant image region when recognizing each character.



دیکھتا جاتا تھا۔

- However, a character or word can appear at multiple locations.
- Nothing stops an attention model from re-attending a previously attended location.
- The decision for attention *in text* needs to depend on the history of attention.

# Computing Attention with Coverage



## Proposed model

Attention weights  $\alpha_{ti}$  computed via softmax over locations

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{j=1}^L \exp(e_{tj})}$$

Exponents computed via

$$e_{ti} = \mathbf{v}_a^T \tanh(\mathbf{W}_a \mathbf{h}_{t-1} + \mathbf{U}_a \mathbf{a}_i + \mathbf{U}_f \mathbf{f}_i)$$

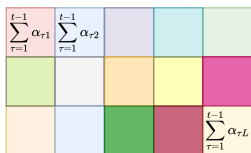
where *coverage vector*<sup>12</sup>  $\mathbf{f}_i$  represents history of attention already given to location  $i$ .

<sup>12</sup>Zhang, Du, and Dai, "Multi-scale attention with dense encoder for handwritten mathematical expression recognition".

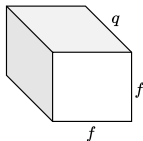
# Coverage vectors

2-step computation of coverage vectors  $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_L$  at time  $t$

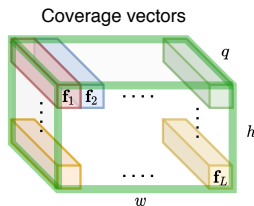
- 1 Compute aggregated attention  $\sum_{\tau=1}^{t-1} \alpha_{\tau i}$  at each location  $i$ .
- 2 Convolve attention aggregates with  $q$  filters of size  $f \times f$ .



Attention aggregates



Learnable convolution filter



Attention of location  $i$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{j=1}^L \exp(e_{tj})}$$
$$e_{ti} = \mathbf{v}_a^T \tanh(\mathbf{W}_a \mathbf{h}_{t-1} + \mathbf{U}_a \mathbf{a}_i + \mathbf{U}_f \mathbf{f}_i)$$

depends on

- hidden state,
- description of location  $i$ , and
- attention already given to location  $i$ .

# Output Character Probability

Conditional probabilities of next character computed via softmax

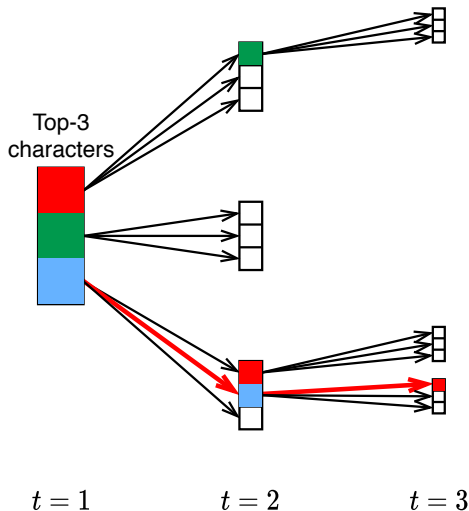
$$p(\mathbf{y}_t | \mathbf{y}_{t-1}; \mathbf{X}) = \text{softmax}(\mathbf{W}_o(\mathbf{E}\mathbf{y}_{t-1} + \mathbf{W}_h\mathbf{h}_t + \mathbf{W}_c\mathbf{c}_t))$$

depend on

- previous output character  $\mathbf{y}_{t-1}$ ,
- hidden state  $\mathbf{h}_t$ , and
- context  $\mathbf{c}_t$  which
  - encodes localized image region(s),
  - based on a history of localizations/coverage.

# Optimal character sequence

Optimal sequence of characters  $y_1, y_2, \dots, y_c$  is found via beam search.



**New Dataset**

## Statistics

- 7,309 offline handwritten Urdu text lines.
- 78,870 words written by 100 different writers.
- 98 unique Urdu characters.

## Collection

- 100 undergraduate students between the ages 20-24 years.
- Requested to submit
  - any handwritten Urdu text, and
  - a corresponding ground-truth text file.
- No restriction on pen type, page type, and ink colour.
- No restriction on what to write.
- Pages scanned at 200 DPI and text lines manually segmented.
- No deskewing.
- Submitted ground-truth was thoroughly checked and corrected/completed by a team of 3 persons.

# Comparison of datasets

Dataset	Total Lines	Total Words	Total characters	Total Writers	Vocabulary Size
UCOM/UNHD <sup>13</sup>	10,000	312,000	1,872,000	500	59
CENIP-UCCP <sup>14</sup>	2,051	23,833	-	200	-
Hassan et al. <sup>15</sup>	6,000	86,400	432,000	600	-
<b>Proposed</b>	<b>7,309</b>	<b>78,870</b>	<b>283,664</b>	<b>100</b>	<b>98</b>

---

<sup>13</sup>Saad Bin Ahmed et al. "Handwritten Urdu character recognition using one-dimensional BLSTM classifier". In: *Neural Computing and Applications* 31.4 (2019), pp. 1143–1151.

<sup>14</sup>A. Raza et al. "An Unconstrained Benchmark Urdu Handwritten Sentence Database with Automatic Line Segmentation". In: *International Conference on Frontiers in Handwriting Recognition*. 2012, pp. 491–496.

<sup>15</sup>S. Hassan et al. "Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting". In: *International Conference on Deep Learning and Machine Learning in Emerging Applications, Deep-ML 2019*. 2019.

# Comparison of datasets

Dataset	Total Lines	Total Words	Total characters	Total Writers	Vocabulary Size
UCOM/UNHD <sup>13</sup>	10,000	312,000	1,872,000	500	59
CENIP-UCCP <sup>14</sup>	2,051	23,833	-	200	-
Hassan et al. <sup>15</sup>	6,000	86,400	432,000	600	-
<b>Proposed</b>	<b>7,309</b>	<b>78,870</b>	<b>283,664</b>	<b>100</b>	<b>98</b>

Existing datasets:

- Unreliable statistics. Only 700 lines in UCOM/UNHD are unique.
- Either publicly unavailable or only partially available.

---

<sup>13</sup>Saad Bin Ahmed et al. "Handwritten Urdu character recognition using one-dimensional BLSTM classifier". In: *Neural Computing and Applications* 31.4 (2019), pp. 1143–1151.

<sup>14</sup>A. Raza et al. "An Unconstrained Benchmark Urdu Handwritten Sentence Database with Automatic Line Segmentation". In: *International Conference on Frontiers in Handwriting Recognition*. 2012, pp. 491–496.

<sup>15</sup>S. Hassan et al. "Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting". In: *International Conference on Deep Learning and Machine Learning in Emerging Applications, Deep-ML 2019*. 2019.

# Comparison of datasets

Dataset	Total Lines	Total Words	Total characters	Total Writers	Vocabulary Size
UCOM/UNHD <sup>13</sup>	10,000	312,000	1,872,000	500	59
CENIP-UCCP <sup>14</sup>	2,051	23,833	-	200	-
Hassan et al. <sup>15</sup>	6,000	86,400	432,000	600	-
<b>Proposed</b>	<b>7,309</b>	<b>78,870</b>	<b>283,664</b>	<b>100</b>	<b>98</b>

Existing datasets:

- Unreliable statistics. Only 700 lines in UCOM/UNHD are unique.
- Either publicly unavailable or only partially available.

Proposed dataset publicly available in full at

<http://faculty.pucit.edu.pk/nazarkhan/>

---

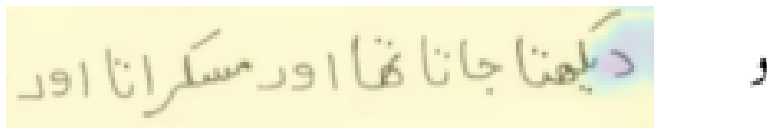
<sup>13</sup>Saad Bin Ahmed et al. "Handwritten Urdu character recognition using one-dimensional BLSTM classifier". In: *Neural Computing and Applications* 31.4 (2019), pp. 1143–1151.

<sup>14</sup>A. Raza et al. "An Unconstrained Benchmark Urdu Handwritten Sentence Database with Automatic Line Segmentation". In: *International Conference on Frontiers in Handwriting Recognition*. 2012, pp. 491–496.

<sup>15</sup>S. Hassan et al. "Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting". In: *International Conference on Deep Learning and Machine Learning in Emerging Applications, Deep-ML 2019*. 2019.

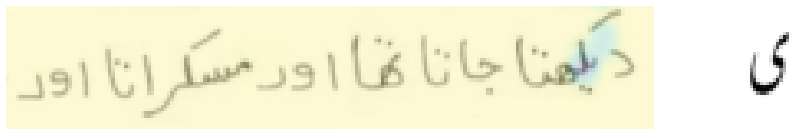
## **Experiments and Results**

# Attention Visualization



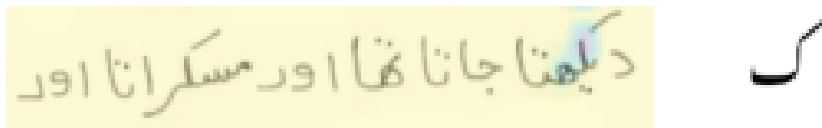
**Figure:** Our model implicitly segments and recognizes every character by focusing only on localized, relevant areas. It also learns to focus in context.

# Attention Visualization



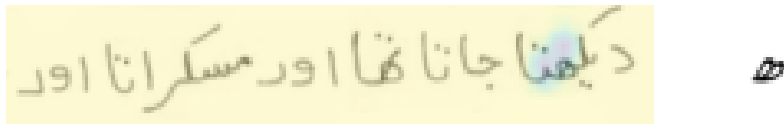
**Figure:** Our model implicitly segments and recognizes every character by focusing only on localized, relevant areas. It also learns to focus in context.

# Attention Visualization



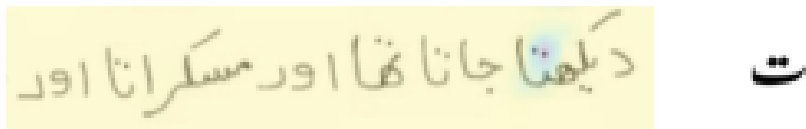
**Figure:** Our model implicitly segments and recognizes every character by focusing only on localized, relevant areas. It also learns to focus in context.

# Attention Visualization



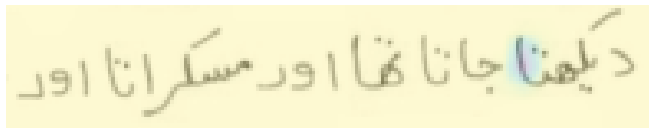
**Figure:** Our model implicitly segments and recognizes every character by focusing only on localized, relevant areas. It also learns to focus in context.

# Attention Visualization



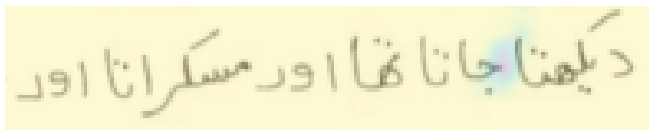
**Figure:** Our model implicitly segments and recognizes every character by focusing only on localized, relevant areas. It also learns to focus in context.

# Attention Visualization



**Figure:** Our model implicitly segments and recognizes every character by focusing only on localized, relevant areas. It also learns to focus in context.

# Attention Visualization



**Figure:** Our model implicitly segments and recognizes every character by focusing only on localized, relevant areas. It also learns to focus in context.

ج

دیکھتا جاتا تھا اور مسکراتا اور

ا

دیکھتا جاتا تھا اور مسکراتا اور

ت

دیکھتا جاتا تھا اور مسکراتا اور

ا

دیکھتا جاتا تھا اور مسکراتا اور

دیکھتا جاتا تھا اور مسکراتا اور

د	دیکھنا جانا تھا اور مسکراتا اور	۱
ی	دیکھنا جانا تھا اور مسکراتا اور	و
ک	دیکھنا جانا تھا اور مسکراتا اور	ر
ھ	دیکھنا جانا تھا اور مسکراتا اور	
ت	دیکھنا جانا تھا اور مسکراتا اور	م
ا	دیکھنا جانا تھا اور مسکراتا اور	س
	دیکھنا جانا تھا اور مسکراتا اور	ک
ج	دیکھنا جانا تھا اور مسکراتا اور	ر
ا	دیکھنا جانا تھا اور مسکراتا اور	ا
ت	دیکھنا جانا تھا اور مسکراتا اور	ت
ا	دیکھنا جانا تھا اور مسکراتا اور	ا
	دیکھنا جانا تھا اور مسکراتا اور	
ت	دیکھنا جانا تھا اور مسکراتا اور	ا
ھ	دیکھنا جانا تھا اور مسکراتا اور	و
ا	دیکھنا جانا تھا اور مسکراتا اور	ر
	دیکھنا جانا تھا اور مسکراتا اور	

## Character Level

$$\text{CER} = \frac{\text{ins} + \text{sub} + \text{del}}{n} \times 100$$

Percentage of insertions, substitutions and deletions required to transform target of length  $n$  into the output.

$$\text{CLA} = 100 - \text{CER}$$

## Word Level

Same formulae as above but with role of characters replaced by words.

# Comparison of different models

Models	CLA	WLA
CNN+GRU	39.79	1.24
CNN+LSTM	39.26	1.40
CNN+BGRU	41.45	1.41
CNN+BLSTM <sup>16</sup>	40.90	1.17
<b>DenseNet+GRU+ Attention</b>	<b>77.05</b>	<b>43.35</b>

---

<sup>16</sup>S. Hassan et al. "Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting". In: *International Conference on Deep Learning and Machine Learning in Emerging Applications, Deep-ML 2019*. 2019.

# Comparison of different models

Models	CLA	WLA
CNN+GRU	39.79	1.24
CNN+LSTM	39.26	1.40
CNN+BGRU	41.45	1.41
CNN+BLSTM <sup>16</sup>	40.90	1.17
<b>DenseNet+GRU+ Attention</b>	<b>77.05</b>	<b>43.35</b>

## Reasons for improvement

- More diverse features learned by DenseNet compared to CNN.
- Attention reduces the need for bidirectional decoding.
- Coverage captures right-to-left nature of Urdu.

---

<sup>16</sup>S. Hassan et al. "Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting". In: *International Conference on Deep Learning and Machine Learning in Emerging Applications, Deep-ML 2019*. 2019.

# Conclusion

- We have proposed an encoder/decoder framework for recognition of offline handwritten Urdu text.

# Conclusion

- We have proposed an encoder/decoder framework for recognition of offline handwritten Urdu text.
- Using attention in context learns the right-to-left nature of Urdu.

# Conclusion

- We have proposed an encoder/decoder framework for recognition of offline handwritten Urdu text.
- Using attention in context learns the right-to-left nature of Urdu.
- Close to 2x improvement in character level accuracy.

# Conclusion

- We have proposed an encoder/decoder framework for recognition of offline handwritten Urdu text.
- Using attention in context learns the right-to-left nature of Urdu.
- Close to 2x improvement in character level accuracy.
- Around 37x improvement in word level accuracy.

# Conclusion

- We have proposed an encoder/decoder framework for recognition of offline handwritten Urdu text.
- Using attention in context learns the right-to-left nature of Urdu.
- Close to 2x improvement in character level accuracy.
- Around 37x improvement in word level accuracy.
- New publicly accessible dataset of >7K text lines consisting of around 79K words written by 100 writers.

<http://faculty.pucit.edu.pk/nazarkhan/>

# Conclusion

- We have proposed an encoder/decoder framework for recognition of offline handwritten Urdu text.
- Using attention in context learns the right-to-left nature of Urdu.
- Close to 2x improvement in character level accuracy.
- Around 37x improvement in word level accuracy.
- New publicly accessible dataset of >7K text lines consisting of around 79K words written by 100 writers.

<http://faculty.pucit.edu.pk/nazarkhan/>

**Thank you for your attention.**



17<sup>th</sup> International Conference on Frontiers in Handwriting Recognition, September 7 – 10, Dortmund, Germany

# References I

- Ahmad, Riaz et al. “KHATT: A deep learning benchmark on Arabic script”. In: *14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2017. DOI: 10.1109/ICDAR.2017.358.
- Ahmed, Saad Bin et al. “Handwritten Urdu character recognition using one-dimensional BLSTM classifier”. In: *Neural Computing and Applications* 31.4 (2019), pp. 1143–1151.
- Gui et al. “Adaptive context-aware reinforced agent for handwritten text recognition”. In: *Proceedings of the British Machine Vision Conference (BMVC)*. 2018.
- Hassan, S. et al. “Cursive Handwritten Text Recognition using Bi-Directional LSTMs: A Case Study on Urdu Handwriting”. In: *International Conference on Deep Learning and Machine Learning in Emerging Applications, Deep-ML 2019*. 2019.
- Mahmoud, Sabri A. et al. “KHATT: An open Arabic offline handwritten text database”. In: *Pattern Recognition* 47.3 (2014), pp. 1096–1112.

# References II

- Pechwitz, Mario et al. "IFN/ENIT-database of handwritten Arabic words". In: *Proc. of CIFED*. Vol. 2. 2002, pp. 127–136.
- Raza, A. et al. "An Unconstrained Benchmark Urdu Handwritten Sentence Database with Automatic Line Segmentation". In: *International Conference on Frontiers in Handwriting Recognition*. 2012, pp. 491–496.
- Sadri, Javad, Mohammad Reza Yeganehzad, and Javad Saghi. "A novel comprehensive database for offline Persian handwriting recognition". In: *Pattern Recognition* 60 (2016), pp. 378–393.
- Tong, A. et al. "NIST 2013 open handwriting recognition and translation evaluation". In: *Proc. NIST*. 2013.
- Zhang, Jianshu, Jun Du, and Lirong Dai. "Multi-scale attention with dense encoder for handwritten mathematical expression recognition". In: *2018 24th international conference on pattern recognition (ICPR)*. IEEE. 2018, pp. 2245–2250.