# Attention is all you need to read

Denis Coquenet

2023/09/19
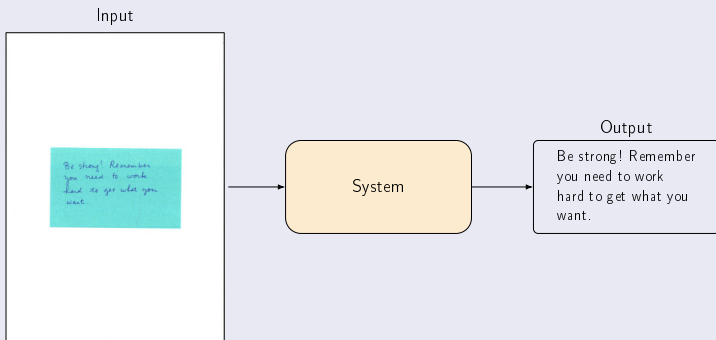
# Table of contents

# Handwritten Text Recognition (HTR)

## An image-to-sequence problem



Input: an image
Output: a sequence of characters

# Challenges

### A wide variety of documents

**Writing styles**, layout, size / resolution, background

## Challenges

### A wide variety of documents

Writing styles, **layout**, size / resolution, background

**Context**
○●○○○○○

Hybrid 1D attention
○○○○○○○

2D attention
○○○○○○○○○○○

Parallel 2D attention
○○○○○○○○○

# Challenges

## A wide variety of documents

Writing styles, layout, size / resolution, **background**

# Challenges

### A wide variety of documents

Writing styles, layout, size / resolution, background

### No a priori knowledge about the document

- Number of lines
- Number of characters per line
- Reading order

# The line-level sequential paradigm

- Segmentation
- Ordering
- **Recognition**

# Related works: Recognition stage

Challenges:

- going from a 2D input image to a 1D sequence of characters
- a variable, unknown number of ordered characters to predict

## Related works: Recognition stage

The Connectionist Temporal Classification (CTC) paradigm [1]



- A frame-by-frame decision process
- Special blank token $\varnothing$
- A left-to-right constrained alignment
- ➤ CTC loss
- ➤ Limited to 1d sequences

[1] Graves *et al.*, ICML 2006

# Related works: Recognition stage

## The attention paradigm (at character level) [2, 3]



t=1, "c"

t=2, "o"

⋮

t=8, "e"

t=9, <eot>

- Iterative decoding process
- Implicit character segmentation
- Special end-of-transcription token <eot>
- Unconstrained attention → reading order must be learned
- ➤ Cross-Entropy loss

$$c^t = \sum_i \alpha_i^t f_i$$

$$\sum_i \alpha_i^t = 1$$

[2] Michael et al., ICDAR 2019

## Conclusion

### The sequential paradigm: a mature approach... with some limitations

- Three steps treated independently
- A complex pipeline, hard to maintain
- Cumulative errors between steps
- Additional segmentation annotations
- Rule-based reading order

Goal: to overcome these limitations

Strategy: designing end-to-end HTR models step by step

➤ from line to document level

# Table of contents

# Related works: Paragraph recognition

## Challenges from line to paragraph recognition

- An additional vertical reading order
- Variable number of text lines
- Variable interline spacing, indent

Context
0000000

Hybrid 1D attention
0●00000

2D attention
00000000000

Parallel 2D attention
000000000

# Related works: Paragraph recognition

## CTC-only approaches

- OrigamiNet [4]



[4] Yousef *et al.*, CVPR 2020

# Related works: Paragraph recognition

## CTC-only approaches

- OrigamiNet [4]
- **Contribution:** Simple Predict & Align Network (SPAN) [5]



SPAN 2D prediction of shape $\frac{H}{32} \times \frac{W}{8}$, before reshaping.
Only the most probable character is represented for each 2D position.

Flatten over vertical axis (reshaping operation)

Remove successive identical tokens

Remove null symbols

[5] Coquenet *et al.*, ICDAR 2021

# Related works: Paragraph recognition

## CTC-only approaches

- OrigamiNet [4]
- **Contribution:** Simple Predict & Align Network (SPAN) [5]

## Attention-based approaches

- Line-level attention [6]

[6] Bluche *et al.*, NIPS 2016

# Related works: Paragraph recognition

## CTC-only approaches

- OrigamiNet [4]
- **Contribution:** Simple Predict & Align Network (SPAN) [5]

## Attention-based approaches

- Line-level attention [6]
- Character-level attention [7, 8]



[7] Bluche *et al.*, ICDAR 2017

Context
○○○○○○○

**Hybrid 1D attention**
○○○●○○○○

2D attention
○○○○○○○○○○○

Parallel 2D attention
○○○○○○○○○

# Contribution: Vertical Attention Network (VAN) [9]

## Overview



[9] Coquenet *et al.*, TPAMI 2023

## Line-level vertical hybrid attention

$$\alpha_i^t = \text{softmax}\big(W_a \tanh(W_f f_i' + W_j j_{t,i} + W_h h_{W_f(t-1)})\big)$$

Context
0000000

Hybrid 1D attention
0000●000

2D attention
00000000000

Parallel 2D attention
000000000

# Datasets



RIMES 2011



IAM



READ 2016

# Paragraph-level recognition results

Paragraph-level state-of-the-art approaches, without language model, external data, nor lexicon constraints.

| Architecture | IAM | | RIMES 2011 | | READ 2016 | | # Param. |
|---|---|---|---|---|---|---|---|
| | CER (%) test | WER (%) test | CER (%) test | WER (%) test | CER (%) test | WER (%) test | |
| Best line-level approach | 4.87[1] | | 2.3[2] | 9.6[2] | 4.66[1] | | |
| [7] CNN+MDLSTM[b] | 16.2 | | | | | | |
| [6] CNN+MDLSTM[a] | 7.9 | 24.6 | 2.9 | 12.6 | | | |
| [8] CNN+Transformer[b] | 6.7 | | | | | | 27.8 M |
| [5] SPAN (FCN) | 5.45 | 19.83 | 4.17 | 15.61 | 6.20 | 25.69 | 19.2 M |
| [4] OrigamiNet (GFCN) | 4.7 | | | | | | 16.4 M |
| [9] VAN (FCN+LSTM)[a] | **4.45** | **14.55** | **1.91** | **6.72** | **3.59** | **13.94** | 2.7 M |

[1] Results from [2] CNN+BLSTM[b].
[2] Results from [10] CNN+BLSTM.
[a] With line-level attention.
[b] With character-level attention.

Context
○○○○○○○

Hybrid 1D attention
○○○○○●○

2D attention
○○○○○○○○○○○

Parallel 2D attention
○○○○○○○○○

# VAN demonstration

https://youtu.be/OXi1birmbuw

# Conclusion

## Attention is powerful but:

- Attention mechanisms → slower convergence
  ➤ vertical attention (1D) + pre-training
- Hybrid attention
  ➤ recurrent training (OK for lines, KO for chars)

## Bridging the gap between line-level and paragraph-level approaches...

- State-of-the-art results on RIMES 2011, IAM and READ 2016
- Able to deal with slightly inclined lines

## ... but still the same limitations, inherent to the sequential paradigm

➤ Rethinking the paradigm

# Table of contents

## HTR at document level



### Challenges from paragraph to document

- Layout-dependent reading order
- Larger input images and output sequences
  - ➤ GPU constraints
  - ➤ More complex attention

# Handwritten Document Recognition (HDR)

Goal: joint recognition of both text and layout from whole documents



Handwritten Document Recognition

# How to encode both text and layout ?



```
<document>
  <page>
    <page_number>
      204
    </page_number>
    <section>
      <body>
        Schgrafer, [...] gehalt.
      </body>
    </section>
    <section>
      <annotation>
        Genneral [...] Raitüng
      </annotation>
      <body>
        Aüf den: [...] werden,
      </body>
    </section>
  </page>
  <page>
    <page_number>
      204
    </page_number>
    <section>
      <annotation>
        Schmalz. [...] bet:
      </annotation>
      <body>
        Verer [...] dar¬
      </body>
    </section>
  </page>
</document>
```

➤ XML paradigm

## How to evaluate the performance ?

### Evaluate the text recognition

- CER / WER
- ➤ Normalized edit distance between sequences of characters / words

Prediction: "<A><B>HTR</B>2<B>HDR</B></A>"
Metric computed on: "HTR2HDR"

# How to evaluate the performance ?

**Evaluate the text recognition**

- CER / WER

**Evaluate the layout recognition**

- LOER (Layout Ordering Error Rate)
- ➤ Normalized edit distance between graphs

Prediction: "`<A><B>HTR</B>2<B>HDR</B></A>`"
Metric computed on: "`<A><B></B><B></B></A>`"

# How to evaluate the performance ?

## Evaluate the text recognition

- CER / WER

## Evaluate the layout recognition

- LOER (Layout Ordering Error Rate)

⚠ **Not sufficient:**

Ground truth: "&lt;A&gt;&lt;B&gt;HTR&lt;/B&gt;2&lt;B&gt;HDR&lt;/B&gt;&lt;/A&gt;"
Prediction: "&lt;A&gt;&lt;B&gt;&lt;/B&gt;&lt;B&gt;&lt;/B&gt;&lt;/A&gt;HTR2HDR"

LOER = 0%    CER = 0%

# How to evaluate the performance ?

## Evaluate the text recognition

- CER / WER

## Evaluate the layout recognition

- LOER (Layout Ordering Error Rate)

## Evaluate text and layout recognition altogether

- $mAP_{CER}$
- ➤ Area under the precision / recall curve

Prediction: "`<A><B>`HTR`</B>`2`<B>`HDR`</B></A>`"
Metric computed on: "HTR2HDR", "HTR", "HDR"

Context
0000000

Hybrid 1D attention
0000000

2D attention
00000●000000

Parallel 2D attention
000000000

# Document Attention Network (DAN) [11]



$$\boldsymbol{c}^t = \underbrace{\text{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^T}{\sqrt{d_k}}\right)\boldsymbol{V}}_{\boldsymbol{\alpha}^t}$$

$\rightarrow$ Teacher forcing

[11] Coquenet *et al.*, TPAMI 2023

## DAN - Training strategy

- Pre-training on synthetic text line images.
- Curriculum learning with synthetic documents:



(a) $l = 3$.



(b) $l = 15$.



(c) $l = l_{\max} = 30$ (end of curriculum stage, no crop).

## Datasets





| Dataset | Level | Training | Validation | Test | # char tokens | # layout tokens |
|---------|-------|----------|------------|------|---------------|-----------------|
| RIMES 2009 [12] | Page | 1,050 | 100 | 100 | 108 | 14 |
| READ 2016 [13] | Page | 350 | 50 | 50 | 89 | 10 |
|  | Double page | 169 | 24 | 24 |  |  |

Context
0000000

Hybrid 1D attention
0000000

**2D attention**
0000000●0000

Parallel 2D attention
000000000

# DAN results on the RIMES dataset

⚠ Metrics do not take into account the segmentation step

| Dataset | Approach | CER (%) ↓ | WER (%) ↓ | LOER (%) ↓ | mAP$_{CER}$ (%) ↑ |
|---|---|---|---|---|---|
| RIMES 2011 | **Line level** | | | | |
| | [9] FCN | 3.04 | 8.32 | ✗ | ✗ |
| | [10] CNN+BLSTM[a] | **2.3** | 9.6 | ✗ | ✗ |
| | [11] DAN (FCN+transformer)[c] | 2.63 | **6.78** | ✗ | ✗ |
| | **Paragraph level** | | | | |
| | [5] SPAN (FCN) | 4.17 | 15.61 | ✗ | ✗ |
| | [6] CNN+MDLSTM[b] | 2.9 | 12.6 | ✗ | ✗ |
| | [9] VAN (FCN+LSTM)[b] | 1.91 | 6.72 | ✗ | ✗ |
| | [11] DAN (FCN+transformer)[c] | **1.82** | **5.03** | ✗ | ✗ |
| RIMES 2009 | **Paragraph level** | | | | |
| | [11] DAN (FCN+transformer)[c] | 5.46 | 13.04 | ✗ | ✗ |
| | **Page level** | | | | |
| | [11] DAN (FCN+transformer)[c] | 4.54 | 11.85 | 3.82 | 93.74 |

[a] This work uses a slightly different split (10,203 for training, 1,130 for validation and 778 for test).
[b] with line-level attention.
[c] with character-level attention.

Context
0000000

Hybrid 1D attention
0000000

2D attention
0000000000●00

Parallel 2D attention
000000000

# DAN results on the READ 2016 dataset

⚠ Metrics do not take into account the segmentation step

| Approach | CER (%) ↓ | WER (%) ↓ | LOER (%) ↓ | mAP$_{CER}$ (%) ↑ |
|---|---|---|---|---|
| **Line level** | | | | |
| [2] CNN+BLSTM[a] | 4.66 | ✗ | ✗ | ✗ |
| [13] CNN+RNN | 5.1 | 21.1 | ✗ | ✗ |
| [9] VAN (FCN+LSTM)[b] | **4.10** | **16.29** | ✗ | ✗ |
| [11] DAN (FCN+transformer)[a] | **4.10** | 17.64 | ✗ | ✗ |
| **Paragraph level** | | | | |
| [5] SPAN (FCN) | 6.20 | 25.69 | ✗ | ✗ |
| [9] VAN (FCN+LSTM)[b] | 3.59 | 13.94 | ✗ | ✗ |
| [11] DAN (FCN+transformer)[a] | **3.22** | **13.63** | ✗ | ✗ |
| **Single-page level** | | | | |
| [11] DAN (FCN+transformer)[a] | 3.53 | 13.33 | 5.94 | 92.57 |
| **Double-page level** | | | | |
| [11] DAN (FCN+transformer)[a] | 3.69 | 14.20 | 4.60 | 93.92 |

[a] with character-level attention.
[b] with line-level attention.

# DAN demonstration

https://youtu.be/HrrUsQfW66E

## Conclusion

DAN: the first end-to-end model for HDR

➤ Structured output sequence

➤ No need for any physical segmentation annotation

➤ Can follow the slant of the lines (character-level attention)

### Line-level / paragraph-level limitations

- ~~Three steps treated independently~~
- ~~A complex pipeline, hard to maintain~~
- ~~Cumulative errors between steps~~
- ~~Additional segmentation annotations~~
- ~~Rule-based reading order~~

Drawback: prediction times grow with the character sequence

# Table of contents

# Faster DAN: parallelizing text line recognition [14]



(a) DAN

(b) Faster DAN

[14] Coquenet et al., ICDAR 2023

Context
○○○○○○○

Hybrid 1D attention
○○○○○○○

2D attention
○○○○○○○○○○○

Parallel 2D attention
○●○○○○○○○

# Faster DAN - Multi-target queries

Context
0000000

Hybrid 1D attention
0000000

2D attention
00000000000

Parallel 2D attention
000●000000

# Faster DAN - Positional encoding

| <sot> | <page> | T | H | E | ↵ | F | A | S | T | E | R | ↵ | D | A | N | </page> | <eot> |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | |

(a) DAN single-pass prediction process

First pass: start-of-line character recognition

Second pass: completing lines in parallel

| | <sot> | <page> | T | F | D | </page> | <eot> |
|---|---|---|---|---|---|---|---|
| **line index** ⟶ | 0 | 0 | 0 | 0 | 0 | 0 | |
| **index in line** ⟶ | 0 | 1 | 2 | 3 | 4 | 5 | |

| T | H | E | <eot> |
|---|---|---|---|
| 1 | 1 | 1 | |
| 0 | 1 | 2 | |

| F | A | S | T | E | R | <eot> |
|---|---|---|---|---|---|---|
| 2 | 2 | 2 | 2 | 2 | 2 | |
| 0 | 1 | 2 | 3 | 4 | 5 | |

| D | A | N | <eot> |
|---|---|---|---|
| 3 | 3 | 3 | |
| 0 | 1 | 2 | |

(b) Faster DAN two-pass prediction process

# Faster DAN - Context



(a) Context used by the DAN

(b) Context used by the Faster DAN

# Results

| Architecture | READ 2016 (single-page) | | | | READ 2016 (double-page) | | | |
|---|---|---|---|---|---|---|---|---|
| | CER ↓ | WER ↓ | LOER ↓ | mAP$_{CER}$ ↑ | CER ↓ | WER ↓ | LOER ↓ | mAP$_{CER}$ ↑ |
| DAN [11] | **3.43** | **13.05** | 5.17 | 93.32 | **3.70** | **14.15** | 4.98 | 93.09 |
| Faster DAN [14] | 3.95 | 14.06 | **3.82** | **94.20** | 3.88 | 14.97 | **3.08** | **94.54** |

| Architecture | RIMES 2009 | | | |
|---|---|---|---|---|
| | CER ↓ | WER ↓ | LOER ↓ | mAP$_{CER}$ ↑ |
| DAN [11] | **4.54** | **11.85** | **3.82** | **93.74** |
| Faster DAN [14] | 6.38 | 13.69 | 4.48 | 91.00 |

## Prediction times

|  | RIMES 2009 | READ 2016 | | MAURDOR | | |
|---|---|---|---|---|---|---|
|  |  | single-page | double-page | C3 | C4 | C3 & C4 |
| Dataset details (averaged for a document on the test set) | | | | | | |
| width (px) | 1,235 | 1,190 | 2,380 | 1,336 | 1,240 | 1,297 |
| height (px) | 1,751 | 1,755 | 1,755 | 1,658 | 1,754 | 1,697 |
| # chars | 578 | 528 | 1,062 | 481 | 706 | 575 |
| # lines | 18 | 23 | 47 | 16 | 22 | 18 |
| # chars / line | 31 | 22 | 22 | 30 | 31 | 30 |
| # layout tokens | 11 | 15 | 30 | 0 | 0 | 0 |
| Prediction times (in seconds) | | | | | | |
| DAN [11] | 5.6 | 4.6 | 8.5 | 5.8 | 7.7 | 6.6 |
| Faster DAN [14] | **1.4** | **0.9** | **1.9** | **1.0** | **1.6** | **1.3** |
| Speed factor | x4 | x5.1 | x4.5 | x5.8 | x4.8 | x5.1 |

# Faster DAN demonstration

https://youtu.be/_pBsO2W8XRE

Context
0000000

Hybrid 1D attention
0000000

2D attention
00000000000

Parallel 2D attention
000000000

# VAN vs DAN vs Faster DAN



VAN: hybrid line-level vertical attention

$\alpha^t$ → $l^t$ → CTC

DAN: transformer, character-level 2D attention

$\alpha^t$ → $c^t$ → Cross Entropy

Faster DAN: transformer, character-level, parallel 2D attention

$\alpha^t$ → $c^t$ → Cross Entropy

Input image $X$

FCN Encoder

Features $f_{2D}$

$H_f$, $W_f$, $C_f$

One-shot encoding

Iterative decoding

# General conclusion

## Attention for reading systems

Line → Paragraph → Document
➤ From text recognition to reading

## Perspectives

Recognizing more:

- Heterogeneous documents (layout)
- Multilingual documents
- Combining HDR with other tasks: Named Entity Recognition, Mathematical Expression Recognition, Table Recognition

Context
0000000

Hybrid 1D attention
0000000

2D attention
00000000000

Parallel 2D attention
000000000

Thank you for your attention

# References I

[1]   Alex Graves, Santiago Fernández, Faustino J. Gomez, and Jürgen Schmidhuber. "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks". In: *International Conference on Machine Learning (ICML)*. Vol. 148. 2006, pp. 369–376.

[2]   Johannes Michael, Roger Labahn, Tobias Grüning, and Jochen Zöllner. "Evaluating Sequence-to-Sequence Models for Handwritten Text Recognition". In: *International Conference on Document Analysis and Recognition (ICDAR)*. 2019, pp. 1286–1293.

[3]   Christoph Wick, Jochen Zöllner, and Tobias Grüning. "Transformer for Handwritten Text Recognition Using Bidirectional Post-decoding". In: *16th International Conference on Document Analysis and Recognition (ICDAR)*. Vol. 12823. 2021, pp. 112–126.

[4]   Mohamed Yousef and Tom E. Bishop. "OrigamiNet: Weakly-Supervised, Segmentation-Free, One-Step, Full Page Text Recognition by learning to unfold". In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2020, pp. 14698–14707.

# References II

[5]   Denis Coquenet, Clément Chatelain, and Thierry Paquet. "SPAN: a Simple Predict & Align Network for Handwritten Paragraph Recognition". In: *International Conference on Document Analysis and Recognition (ICDAR)*. Vol. 12823. 2021, pp. 70–84.

[6]   Théodore Bluche. "Joint Line Segmentation and Transcription for End-to-End Handwritten Paragraph Recognition". In: *Advances in Neural Information Processing Systems 29 (NIPS)*. 2016, pp. 838–846.

[7]   Théodore Bluche, Jérôme Louradour, and Ronaldo O. Messina. "Scan, Attend and Read: End-to-End Handwritten Paragraph Recognition with MDLSTM Attention". In: *International Conference on Document Analysis and Recognition (ICDAR)*. 2017, pp. 1050–1055.

[8]   Sumeet S. Singh and Sergey Karayev. "Full Page Handwriting Recognition via Image to Sequence Extraction". In: *16th International Conference on Document Analysis and Recognition (ICDAR)*. Vol. 12823. 2021, pp. 55–69.

[9]   Denis Coquenet, Clément Chatelain, and Thierry Paquet. "End-to-end Handwritten Paragraph Text Recognition Using a Vertical Attention Network". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 45.1 (2023), pp. 508–524.

# References III

[10]    Joan Puigcerver. "Are Multidimensional Recurrent Layers Really Necessary for Handwritten Text Recognition?" In: *14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. 2017, pp. 67–72.

[11]    Denis Coquenet, Clément Chatelain, and Thierry Paquet. "DAN: a Segmentation-free Document Attention Network for Handwritten Document Recognition". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2023).

[12]    Emmanuele Grosicki, Matthieu Carré, Jean-Marie Brodin, and Edouard Geoffrois. "Results of the RIMES Evaluation Campaign for Handwritten Mail Processing". In: *10th International Conference on Document Analysis and Recognition (ICDAR)*. 2009, pp. 941–945.

[13]    Joan-Andreu Sánchez, Verónica Romero, Alejandro H. Toselli, and Enrique Vidal. "ICFHR2016 Competition on Handwritten Text Recognition on the READ Dataset". In: *15th International Conference on Frontiers in Handwriting Recognition (ICFHR)*. 2016, pp. 630–635.

[14]    Denis Coquenet, Clément Chatelain, and Thierry Paquet. "Faster DAN: Multi-target Queries with Document Positional Encoding for End-to-end Handwritten Document Recognition". In: *International Conference on Document Analysis and Recognition (ICDAR)*. Vol. 14190. 2023, pp. 182–199.