# Long-Range Transformer Architectures for Document Understanding

Thibault Douzon, Stefan Duffner, Christophe Garcia and Jérémy Espinas

August 25, 2023

ICDAR 2023 – VINALDO

## Motivations

■ Transformers provide best performance ;

## Motivations

■ Transformers provide best performance ;

■ Transformers don't scale well with long sequences ;

## Motivations

■ Transformers provide best performance ;

■ Transformers don't scale well with long sequences ;

■ How to process long documents?
 – Multi-page documents
 – Dense text (scientific, legal, . . . )
 – Both?

Multi-page customer order with annotations for information extraction.

Single-page document from Docbank [Li et al., 2020] with its layout annotation mask.

Transformer [Vaswani et al., 2017] and LayoutLM [Xu et al., 2020] architectures.

Default dot-product attention with $O(N^2)$ complexity.

Efficient transformer architectures Venn diagram from Tay et al. [2022].

Linformer [Wang et al., 2020] attention approximation.

Key's and Value's sequence length dimensions are projected onto a smaller space.

Cosformer [Qin et al., 2022] kernel attention approximation.
Where $\Phi$ is a kernel function, $\mathrm{ReLU}$ is used in this work.

## Efficiency

| Model Name | Time (s) / *Memory (GiB)* | | | | | |
|---|---|---|---|---|---|---|
| | Sequence Length | | | | | |
| | 512 | 1024 | 2048 | 4096 | 8192 | 16384 |
| LayoutLM | 1.41/*1.25* | 2.83/*2.50* | 7.39/*5.01* | 23.43/*13.69* | - | - |
| LayoutLinformer | 1.18/*1.35* | 1.92/*2.26* | 3.54/*3.28* | 6.90/*5.19* | 13.08/*8.96* | 25.65/*16.78* |
| LayoutCosformer | 2.03/*1.36* | 2.50/*2.37* | 4.68/*3.38* | 9.00/*5.38* | 17.23/*9.59* | 33.96/*17.59* |

Duration and memory use for various sequence lengths on a reference inference task.

Cumulated F1-scores per document length category.

| Item | Quantity | Description | Unit Price | Total Price |
|---|---|---|---|---|
| 1 | 10 | T5F.L, DN15<br>5.015.F.L.K | $ 101,31 | $ 1 013,10 |
| 2 | 4 | T5F.L, DN25<br>5.025.F.L.K | $ 153,04 | $ 612,16 |
| 3 | 2 | T5F.L, DN32<br>5.032.F.L.K | $ 172,97 | $ 345,94 |
| 4 | 10 | T6F.L, DN15<br>6.015.F.L.K | $ 96,55 | $ 965,50 |
| 5 | 4 | T6F.L, DN25<br>6.025.F.L.K | $ 142,87 | $ 571,48 |
| 6 | 2 | T6F.L, DN32<br>6.032.F.L.K | $ 157,46 | $ 314,92 |
| 7 | 3 | T38V.E, DN15<br>38.015.V.E | $ 210,80 | $ 632,40 |

A closer look at an order's line items table.

Squircle (left) and Cross-shaped (right) attention bias patterns.

# Results - Docbank

| Model Name | F1 Score | | | | | | | | | Macro Average |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Categories | | | | | | | | | |
| | Abst. | Auth. | Capt. | Equa. | Footer | List | Sect. | Table | Title | |
| LayoutLM | 97.8 | 87.5 | 94.9 | 87.2 | 90.5 | 84.0 | 92.8 | 85.7 | 88.6 | 91.6 |
| LayoutLM$_{SQUIRCLE}$ | **98.4** | 90.2 | **96.1** | 89.7 | 92.0 | **88.9** | **94.6** | 87.7 | 90.3 | **93.2** |
| LayoutLM$_{CROSS}$ | **98.4** | **90.3** | 96.0 | 89.6 | **92.1** | 88.7 | **94.6** | 87.5 | **90.7** | **93.2** |

## Results - Docbank

| Model Name | F1 Score | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Categories | | | | | | | | | Macro Average |
| | Abst. | Auth. | Capt. | Equa. | Footer | List | Sect. | Table | Title | |
| LayoutLM | 97.8 | 87.5 | 94.9 | 87.2 | 90.5 | 84.0 | 92.8 | 85.7 | 88.6 | 91.6 |
| LayoutLM$_{SQUIRCLE}$ | **98.4** | 90.2 | **96.1** | 89.7 | 92.0 | **88.9** | 94.6 | 87.7 | 90.3 | **93.2** |
| LayoutLM$_{CROSS}$ | **98.4** | **90.3** | 96.0 | 89.6 | **92.1** | 88.7 | **94.6** | 87.5 | **90.7** | **93.2** |
| LayoutLinformer | 97.9 | 88.9 | 93.7 | 90.0 | 91.1 | 87.9 | 91.3 | 87.6 | 88.7 | 92.3 |
| LayoutCosformer | 97.2 | 87.2 | 91.0 | 88.1 | 90.6 | 87.4 | 81.4 | 87.0 | 88.3 | 90.7 |
| LayoutCosformer$_{SQUIRCLE}$ | 97.0 | 85.4 | 92.4 | 89.2 | 90.7 | 84.2 | 85.6 | 87.9 | 86.8 | 90.7 |
| LayoutCosformer$_{CROSS}$ | 97.4 | 86.9 | 93.8 | **91.2** | 91.7 | 87.5 | 87.4 | **89.0** | 88.1 | 91.9 |

## Takeaway

- Efficient transformer architectures can be valuable for document understanding ;

## Takeaway

- Efficient transformer architectures can be valuable for document understanding ;

- But the cost of attention approximation negatively impact performance on short documents ;

## Takeaway

■ Efficient transformer architectures can be valuable for document understanding ;

■ But the cost of attention approximation negatively impact performance on short documents ;

■ 2D relative attention is hard to tune on efficient architectures and does not perform as anticipated.

Thanks everyone !

Any questions ?

## Paper & contact information

Any in-depth question about this work ? Please contact me !

**Thibault Douzon**, Stefan Duffner, Christophe Garcia and Jérémy Espinas

thibault.douzon@esker.com

Code and models will be available on
github.com/thibaultdouzon/long-range-document-transformer.git

Get the source of this theme and the demo presentation from

github.com/matze/mtheme

The theme *itself* is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License.

## References

Minghao Li, Yiheng Xu, Lei Cui, Shaohan Huang, Furu Wei, Zhoujun Li, and Ming Zhou. DocBank: A Benchmark Dataset for Document Layout Analysis, November 2020.

Zhen Qin, Weixuan Sun, Hui Deng, Dongxu Li, Yunshen Wei, Baohong Lv, Junjie Yan, Lingpeng Kong, and Yiran Zhong. cosFormer: Rethinking Softmax in Attention. *arXiv:2202.08791 [cs]*, February 2022.

Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient Transformers: A Survey. *arXiv:2009.06732 [cs]*, March 2022.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.*, June 2017.

Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-Attention with Linear Complexity. *arXiv:2006.04768 [cs, stat]*, June 2020.

Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. LayoutLM: Pre-training of Text and Layout for Document Image Understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1192–1200, August 2020. doi: 10.1145/3394486.3403172.

# Results - Business Documents - Relative Attention



Left chart:

| Encoding length categories | Label | | | | |
|---|---|---|---|---|---|
| | Number | Date | Total | ItemId | ItemQty |
| **Short** LayoutLM Split Page | 87.6 | 95.9 | 88.2 | 81.3 | 97.5 |
| **Short** LayoutLM Split Page Relative Squircle | 87.7 | 95.8 | 88.4 | 82.5 | 97.7 |
| **Short** LayoutLM Split Page Relative Cross | 87.8 | 96.1 | 87.2 | 81.7 | 97.5 |
| **Medium** LayoutLM Split Page | 85.8 | 89.5 | 70.1 | 75.6 | 90.1 |
| **Medium** Relative Squircle | 86.7 | 89.5 | 71.7 | 76.3 | 90.6 |
| **Medium** Relative Cross | 85.2 | 89.3 | 69.8 | 75.5 | 90.4 |
| **Long** LayoutLM Split Page | 81.0 | 89.2 | 62.1 | 70.6 | 82.8 |
| **Long** Relative Squircle | 83.7 | 87.1 | 64.2 | 71.5 | 82.2 |
| **Long** Relative Cross | 82.4 | 87.3 | 64.0 | 71.8 | 82.1 |

Right chart:

| Encoding length categories | Label | | | | |
|---|---|---|---|---|---|
| | Number | Date | Total | ItemId | ItemQty |
| **Short** LayoutLM Split Page | 87.6 | 95.9 | 88.2 | 81.3 | 97.5 |
| **Short** LayoutCosformer | 85.3 | 94.7 | 87.9 | 74.6 | 95.9 |
| **Short** LayoutCosformer Relative Squircle | 82.8 | 92.4 | 86.7 | 71.3 | 95.5 |
| **Short** LayoutCosformer Relative Cross | 85.6 | 94.5 | 88.1 | 73.7 | 96.0 |
| **Medium** LayoutLM Split Page | 85.8 | 89.5 | 70.1 | 75.6 | 90.1 |
| **Medium** LayoutCosformer | 84.0 | 89.9 | 79.8 | 76.1 | 94.7 |
| **Medium** Relative Squircle | 81.0 | 86.4 | 76.8 | 72.8 | 93.7 |
| **Medium** Relative Cross | 85.7 | 89.0 | 79.8 | 76.1 | 95.1 |
| **Long** LayoutLM Split Page | 81.0 | 89.2 | 62.1 | 70.6 | 82.8 |
| **Long** LayoutCosformer | 86.7 | 84.1 | 68.0 | 69.7 | 86.2 |
| **Long** Relative Squircle | 80.0 | 74.0 | 65.6 | 65.8 | 81.7 |
| **Long** Relative Cross | 84.3 | 82.8 | 66.0 | 67.2 | 85.4 |