

ÉCLAIR – Extracting Content and Layout with Integrated Reading Order for Documents

Ilia Karmanov*, Amala Sanjay Deshmukh*, Lukas Vögtle, Philipp Fischer, Kateryna Chumachenko,
Timo Roman, Jarno Seppänen, Jupinder Parmar, Joseph Jennings, Andrew Tao, Karan Sapra †

NVIDIA

{ikarmanov, amalasanjayd, lvoegtle, pfischer, kchumachenko,
troman, jseppanen, jupinderp, jjennings, atao, ksapra}@nvidia.com

Abstract

Optical Character Recognition (OCR) technology is widely used to extract text from images of documents, facilitating efficient digitization and data retrieval. However, merely extracting text is insufficient when dealing with complex documents. Fully comprehending such documents requires an understanding of their structure — including formatting, formulas, tables, and the reading order of multiple blocks and columns across multiple pages — as well as semantic information for detecting elements like footnotes and image captions. This comprehensive understanding is crucial for downstream tasks such as retrieval, document question answering, and data curation for training Large Language Models (LLMs) and Vision Language Models (VLMs). To address this, we introduce ÉCLAIR, a general-purpose text-extraction tool specifically designed to process a wide range of document types. Given an image, ÉCLAIR is able to extract formatted text in reading order, along with bounding boxes and their corresponding semantic classes. To thoroughly evaluate these novel capabilities, we introduce our diverse human-annotated benchmark DROBS for document-level OCR and semantic classification. ÉCLAIR achieves state-of-the-art accuracy on this benchmark, outperforming other methods across key metrics. Additionally, we evaluate ÉCLAIR on established benchmarks, demonstrating its versatility and strength across several evaluation standards.

1. Introduction

Optical Character Recognition (OCR) has allowed machines to extract text from images and transformed the way we interact with textual information. The recent success of Large Language Models (LLMs) is partly attributed to the availability of extremely large text datasets, placing an

increasing demand for high-quality tokens extracted from text-dense documents such as scientific textbooks and journals. This is a challenging task since it necessitates an understanding of reading order across complex layouts, which in turn requires identifying different semantic elements and their relationships on the page. Maintaining a seamless flow requires separating relevant elements (e.g., paragraphs or tables) from irrelevant ones (e.g., page headers, page footers and other floating text).

Traditional OCR systems operate on a word or line level and are unable to properly understand the spatial and semantic relationships and hierarchies present in text-dense documents. More complex systems that possess such capabilities are generally not end-to-end and combine several models into a brittle pipeline. This shortcoming has spawned an interest in developing end-to-end models [6, 31, 52] that can extract formatted and structured text from complex documents such as those shown in Figure 1. Such capabilities provide downstream benefits for a multitude of tasks, including retrieval, document question answering, and increasing the availability of text tokens for LLM training. However, recent models proposed in this area still have crucial drawbacks: Kosmos-2.5 [31] lacks the ability to extract formatted text that is at the same time spatially aware (since it has two mutually-exclusive prompts), while GOT [52] and Nougat [6] do not predict any spatial information at all. In addition, none of these models predict semantic classes of bounding boxes, which can be used as a conditioning element for retrieval, help filter out irrelevant information for LLM training, and assist when combining multiple pages within a document (e.g., placing footnotes only after a text section has ended).

To address these concerns, we present ÉCLAIR: a multi-modal LLM (MLLM) comprised of a ViT-like encoder and an auto-regressive decoder, architecturally similar to Donut [23]. ÉCLAIR is able to extract text (formatted as markdown/LaTeX), bounding boxes of text blocks with their

*Equal contribution.

†Project Lead.

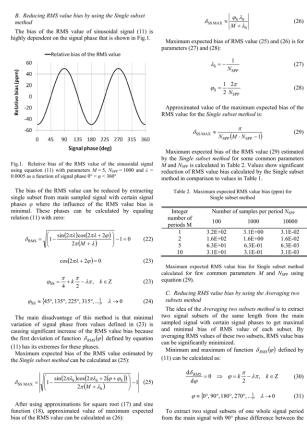


Figure 1. ÉCLAIR outperforms other methods on complex documents: (a) tables, formulas, figure, page header and multiple columns; (b) uneven columns, styling, figure; (c) non-obvious reading order and visual elements like background coloring.

semantic classes, and any combination of these simultaneously, while preserving the reading order.

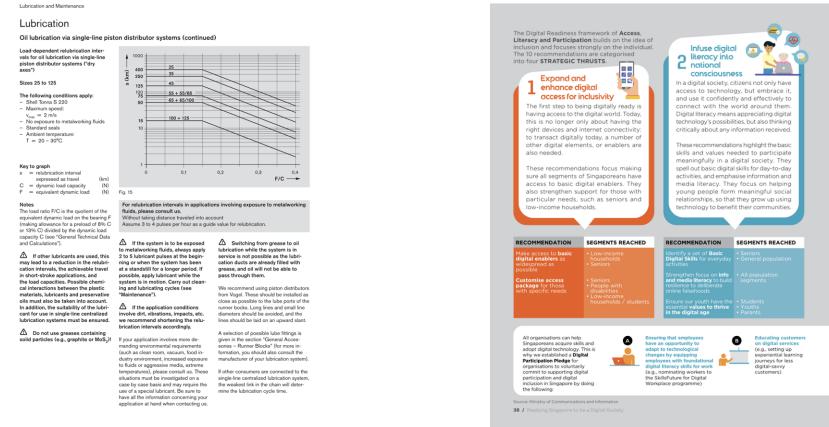
Training such a versatile model necessitates a dataset that encompasses all these annotation types. To address this problem, we generate arXiv-5M, a large-scale dataset that is sampled from arXiv papers, covers all desired annotation capabilities and serves as a link between existing datasets of varying annotation coverage.

Our proposed novel data generation pipeline includes a modified L^AT_EX compiler which generates ground truth labels directly from the L^AT_EX sources.

Furthermore, existing document-level OCR benchmarks are limited by partial annotations: GOT [52] allows for document level reading order but lacks block-level spatial and semantic labels, while DocLayNet [40] lacks reading order information. To address these shortcomings, we release a new benchmark DROBS consisting of 789 visually-diverse pages sampled from different sources, see Figure 3 for examples. The annotations come from human labeling and contain text in reading-order along with bounding boxes and semantic classes. ÉCLAIR achieves state-of-the-art (SOTA) accuracy on DROBS when compared to other recent models, as well as competitive metrics on several existing benchmarks spanning different tasks, including general OCR, document layout understanding, and extraction of high-quality data for LLM training.

To summarize, our contributions are as follows:

- We create an end-to-end document-level OCR model which is the first to extract formatted text with their respective bounding boxes and semantic classes.
 - We develop a novel data generation pipeline where we can control the rendering of L^AT_EX sources ourselves. This is a pre-requisite for bridging the gap between ex-



isting datasets with fewer label types.

- We release a new benchmark DROBS, with high-quality human-labeled data and show SOTA accuracy as well as competitive metrics on existing benchmarks.

2. ÉCLAIR

2.1. Architecture

ÉCLAIR uses a transformer encoder-decoder architecture. The **vision encoder**, denoted as \mathcal{E} , is initialized from RADIO [42] which follows a ViT-H /16 [10] architecture, and maps an image $\mathbf{I} \in \mathbb{R}^{3 \times H \times W}$ to a latent representation $\mathbf{Z} \in \mathbb{R}^{N \times d}$, where d is the hidden dimension and N is the sequence length. The **neck** \mathcal{N} then reduces the dimensionality of the latent space as well as the sequence length.

The **decoder**, denoted as \mathcal{D} , uses mBART [28] and predicts text tokens $\mathbf{T} = \{t_{P+1}, t_{P+2}, \dots, t_L\}$ by conditioning on the latent encoder representation, $\mathcal{N}(\mathbf{Z})$, and the context $t_{<i}$, $P(t_i | \mathcal{N}(\mathbf{Z}), t_{<i})$, where $\mathbf{Z} = \mathcal{E}(\mathbf{I})$ and $\{t_1, t_2, \dots, t_P\}$ are the prompt tokens and where L is the prompt-augmented sequence length.

Since autoregressive models scale poorly with the decoder size and sequence length at inference time, we adopt a larger vision encoder (657M parameters) and combine it with a lightweight decoder (279M parameters). This follows from the observation that OCR is not fundamentally a generative task but rather depends on the content in the input image. We describe further modifications to improve inference time in Section 3.5.

2.2. Prompt

We use the input prompt to specify the desired format of the model outputs. Each prompt is a tuple of three options,

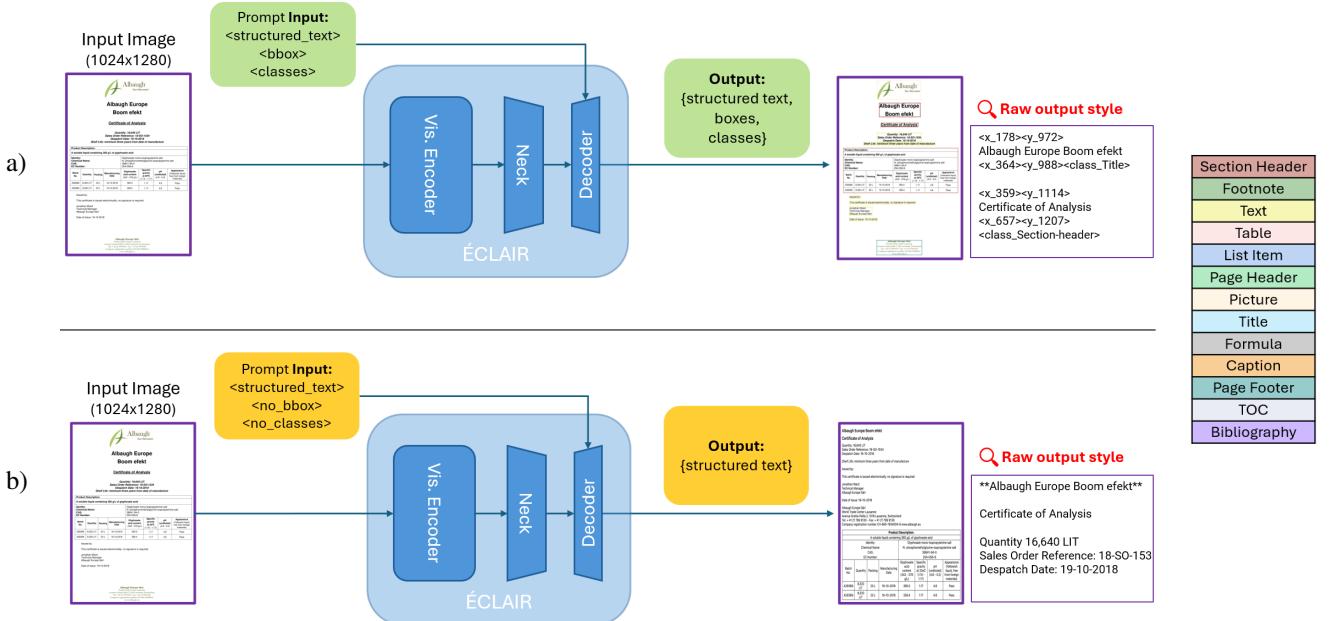


Figure 2. Meta architecture for ÉCLAIR showcasing the usage with two different (out of eight valid) prompts: Example a) uses the maximal information prompt to return bounding boxes along with their semantic class, markdown text, and tables and formulae. In b) we ask the model to return only markdown text without boxes or classes. All supported semantic classes are listed on the right.

with 8 possible valid combinations (ignoring the trivial case of no output, and the cases where semantic classes are requested without the corresponding bounding boxes):

- <structured_text> or <plain_text>
or <no_text>
- <bbox> or <no_bbox>
- <classes> or <no_classes>

For each of the three groups, the first option specifies the most information, while the last option suppresses this output type. With the <structured_text> prompt, the text is predicted in markdown format and inline formulae are formatted as \LaTeX , whereas with the <plain_text> prompt both are formatted as plain text. Tables and block formulae are formatted as \LaTeX for both modes. We define the maximal-information prompt (MIP) as:

<structured_text><bbox><classes>

The novelty of ÉCLAIR compared to existing methods lies in its ability to handle any of the 8 valid prompt combinations. This is achieved by pre-training on a custom dataset that has labels for the maximal-information setting and then decreasing the information density for each group with some dataset-dependent probability during the fine-tuning stage. This allows the model to leverage visually diverse datasets with partial annotations for training. A schematic structure of ÉCLAIR along with possible prompts and corresponding output is presented in Figure 2.

2.3. Output Format and Tokenization

ÉCLAIR predicts bounding boxes of the semantic blocks in the form of discrete coordinates, similar to Kosmos [31]. These bounding boxes are predicted in a canonical reading order, which is described further in the supplementary material. The following regular expression shows the output format for each box in the maximal-information setting:

`<x_(\d+)><y_(\d+)>(.*)<x_(\d+)><y_(\d+)>
<class_([^\>]+)>`

where the the **first group** denotes the coordinates of the top-left corner, the **second group** denotes the text contained within the bounding box, the **third group** denotes the coordinates of the bottom-right corner, and the **final group** represents the semantic class.

Note that each of these groups is optional and their presence in the model's output for a given image would depend on the prompt combination specified for that sample.

We adopt the tokenizer used by Taylor et al. [47], as their model is also specialized for the scientific text domain. The coordinates of the bounding boxes, the semantic classes and the seven prompt components are all added as dedicated special tokens. This adds $H + W + C + 7$ tokens in total to the tokenizer vocabulary, where C is the number of semantic classes.

2.4. Datasets

Compared to existing methods, such as Kosmos-2.5 [31] and Nougat [6], ÉCLAIR is trained on a relatively smaller dataset as summarized in Table 1.

| Dataset | Size | Modality |
|---|--------|----------------------------|
| arXiv-5M | 5M | Structured, Boxes, Classes |
| SynthTabNet [33] | 480K | Structured, Boxes, Classes |
| README | 302K | Structured |
| DocLayNet [40] | 56K | Plain, Boxes, Classes |
| G1000 [49] | 324K | Plain |
| Human-labeled Common Crawl samples | 14K | Plain, Boxes, Classes |
| Total | 6.176M | |

Table 1. Summary of the datasets used to train ÉCLAIR, including a description of the maximum information available in the annotations of each dataset.

The arXiv-5M dataset makes up a large portion of our training data and it supports the maximum-information prompt (MIP) described in Section 2.2. The generation pipeline used to create this dataset is discussed further in Section 2.5. We pre-train ÉCLAIR on this dataset.

We find that recent models such as Nougat [6], that are only trained on academic documents, do not handle visually-diverse documents very well, often either degenerating into hallucinations or repetition loops or simply terminating early by predicting the end-of-sequence token. We hypothesize that this is because the training data lacks the heterogeneity needed to handle more complex layouts such as magazines, leaflets, and picture-books. To address this, we fine-tune ÉCLAIR further on the arXiv-5M along with several publicly available datasets with diverse layouts and domains, such as DocLayNet [40], SynthTabNet [33] and G1000 [49]. We also create a high-quality human-annotated dataset consisting of documents sampled from the Common Crawl corpus [12]. Additionally, we create a README dataset by sampling README documents from the Stack [24] and rendering them using Pandoc [37]. Most of these datasets contain only partial annotations and the maximum information available in each is summarized in Table 1. The pre-processing steps for these datasets are described in more detail in the supplementary material.

2.5. The arXiv-5M Dataset

In the introduction, we briefly discussed the need for a dataset that provides labels for our maximum information setting, i.e. bounding boxes, semantic classes and formatted text with formulas and tables, all in reading order. Since no such dataset exists, we created a new one.

Our approach is inspired by Nougat [6], where the authors create ground truth image/markdown pairs from arXiv papers. Their pipeline relies on LatexML, a tool to convert L^AT_EX source code to HTML, which they convert to markdown subsequently. We follow a different approach, which handles both the L^AT_EX compilation and the conversion to structured output at the same time (instead of using separate processing pipelines for each) and hence retains the re-

lationship between text and image down to character-level bounding boxes and allows us to extract semantic classes for each box. Our representation for the structured text output

- consists of rectangular boxes
- has a semantic class assigned to each box
- represents normal text and formatted text as markdown
- represents tables and formulas as L^AT_EX

The box and class information can be used to re-arrange the order of content (e.g. footnotes at the end) and to filter unwanted content, for example page headers and footers.

We modify the open-source T_EX Live distribution by adding hooks inside the T_EX compiler itself and embedding a Python interpreter for further processing on-the-fly. We hook the internal T_EX methods for node, character and hbox/vbox allocations, token reading and output generation and forward these to a custom Python class that keeps track of the elements from allocation to output on the PDF page. Multiple stacks are used to keep track of how the elements are nested in the input and output, and a rule-based system generates a nested hierarchy with the elements of interest.

With this method we generated a high-quality ground-truth dataset consisting of around roughly 5 million pages which we call arXiv-5M.

3. Results

The details about our experimental setup and training strategy can be found in the supplementary material.

3.1. Reading Order benchmark

DROBS Evaluation. We evaluate the reading order accuracy of ÉCLAIR against known SOTA methods like Kosmos-2.5 [31] and GOT [52]. Both of these methods have two output modalities - a plain OCR mode and a markdown mode, and we compare ÉCLAIR with both modes. For this evaluation, we utilize an internally curated and human-labeled diverse set of PDFs, comprising a total of 789 pages sampled from various sources such as magazines, books, and the Common Crawl corpus [45] which we call DROBS (Document Reading Order, Bounding boxes, and Semantic classes), see Figure 3. This approach aims to cover a diversity of layouts similar to those found in DocLayNet [40]. We instructed human annotators to provide annotations on this dataset following the same labeling system as DocLayNet due to its comprehensive human annotation guidelines. However, we added additional requirements, the most significant being the inclusion of reading order. We will make the selected pages and associated annotations available to the research community to serve as an additional and complementary benchmark for document understanding and OCR.

| Method | Mask out | Counting F1** ↑ | WER ↓ | Edit distance ↓ | F1 ↑ | Precision ↑ | Recall ↑ | BLEU ↑ | METEOR ↑ |
|-----------------------|----------|-----------------|--------------|-----------------|--------------|--------------|--------------|--------------|--------------|
| ÉCLAIR-MIP | ✗ | 0.934 | 0.142 | 0.109 | 0.942 | 0.960 | 0.942 | 0.886 | 0.930 |
| ÉCLAIR-MIP | ✓ | 0.937 | 0.146 | 0.108 | 0.941 | 0.966 | 0.936 | 0.885 | 0.927 |
| Kosmos-2.5 (ocr-mode) | ✓ | 0.919 | 0.195 | 0.114 | 0.937 | 0.932 | 0.950 | 0.862 | 0.927 |
| Kosmos-2.5 (md-mode) | ✓ | 0.843 | 0.249 | 0.184 | 0.890 | 0.941 | 0.876 | 0.805 | 0.851 |
| GOT (ocr-mode) | ✓ | 0.776 | 0.302 | 0.216 | 0.818 | 0.863 | 0.825 | 0.713 | 0.795 |
| GOT (md-mode) | ✓ | 0.825 | 0.259 | 0.157 | 0.879 | 0.908 | 0.875 | 0.760 | 0.852 |

Table 2. Evaluation results on DROBS. Reported standard NLTK metrics [41] are character level (Edit-distance) or word level (F1, Precision, Recall, BLEU, METEOR) metrics typically used by the OCR and natural language processing (NLP) communities. We also report Counting F1 and word error rate/word edit distance metrics.

*MIP-maximal-information prompt

**Counting F1 score is computed over the set { he₁, said₁, that₁, she₁, said₂, that₂, they₁, said₃, that₃, he₂, said₄, something₁ }. This allows to track and penalize words that missed but has more than one occurrence in the document.



Figure 3. Example pages from DROBS, our visually diverse document benchmark.

Prior to evaluation, we perform three preprocessing steps on the predictions and corresponding ground truth labels. First, we apply string normalization [35] to remove all non-alphanumeric characters, convert sequences of whitespaces to a single space, and strip any leading or trailing whitespaces to ensure fair comparison. Second, to address variations in the output formats for tables and equations among different methods, and given that our current evaluation benchmark does not include labels for equations and tables, we mask out these elements in the images used for inference across all considered methods. Additionally, for GOT (md), we also mask-out headers and footers from the images, as the model seems to ignore these elements. Lastly, we also filter out TeX commands present in GOT prediction e.g. for title, section and sub-section headers, since those would otherwise penalize the model in text-only metrics.

The results in Table 2 show how ÉCLAIR outperforms both Kosmos and GOT on most metrics. Since these models are trained with multiple prompt modes similar to ÉCLAIR, we compare the output of ÉCLAIR in MIP mode against Kosmos-2.5 and GOT in both OCR and MD modes. We observe that Kosmos-2.5 performs better in OCR mode, where it produces bounding boxes; however, GOT exhibits the opposite behavior, performing better in MD mode as opposed to OCR mode. We believe this discrepancy is due to differences in data blending during training. ÉCLAIR in MIP mode produces both MD for text inside the bounding boxes and shows superior accuracy compared to both other methods.

Since there is no common validation set available for comparing equation and table extraction across methods, we evaluate our formula and table extraction accuracy on a validation set derived from arXiv (See Section 3.2).

Training & Inference Ablation. All the results for ÉCLAIR presented in Tables 2 and 3 were obtained with a repetition penalty [22] of 1.1 applied during inference. Table 5 demonstrates the value of this inference-time hyperparameter and of the additional datasets added in the fine-tuning stage.

GOT Benchmark. Along with DROBS, we evaluate ÉCLAIR on the GOT benchmark proposed in Fox [27], with the results shown in Table 3. ÉCLAIR with a size of 936M parameters and in MIP mode demonstrates competitive or superior accuracy across most metrics, particularly

| Method | Size | Edit Distance↓ | | F1-score↑ | | Precision↑ | | Recall↑ | | BLEU↑ | | METEOR↑ | |
|------------------|------|----------------|-------|--------------|-------|--------------|-------|--------------|-------|--------------|-------|--------------|-------|
| | | en | zh | en | zh | en | zh | en | zh | en | zh | en | zh |
| Nougat [6] | 250M | 0.255 | - | 0.745 | - | 0.720 | - | 0.809 | - | 0.665 | - | 0.761 | - |
| TextMonkey [29] | 7B | 0.265 | - | 0.821 | - | 0.778 | - | 0.906 | - | 0.671 | - | 0.762 | - |
| DocOwl1.5 [16] | 7B | 0.258 | - | 0.862 | - | 0.835 | - | 0.962 | - | 0.788 | - | 0.858 | - |
| Vary [50] | 7B | 0.092 | 0.113 | 0.918 | 0.952 | 0.906 | 0.961 | 0.956 | 0.944 | 0.885 | 0.754 | 0.926 | 0.873 |
| Vary-toy [51] | 1.8B | 0.082 | 0.142 | 0.924 | 0.914 | 0.919 | 0.928 | 0.938 | 0.907 | 0.889 | 0.718 | 0.929 | 0.832 |
| Qwen-VL-Plus [3] | - | 0.096 | 0.121 | 0.931 | 0.895 | 0.921 | 0.903 | 0.950 | 0.890 | 0.893 | 0.684 | 0.936 | 0.828 |
| Qwen-VL-Max [3] | 72B+ | 0.057 | 0.091 | 0.964 | 0.931 | 0.955 | 0.917 | 0.977 | 0.946 | 0.942 | 0.756 | 0.971 | 0.885 |
| Fox [26] | 1.8B | 0.046 | 0.061 | 0.952 | 0.954 | 0.957 | 0.964 | 0.948 | 0.946 | 0.930 | 0.842 | 0.954 | 0.908 |
| GOT [52] | 580M | 0.035 | 0.038 | 0.972 | 0.980 | 0.971 | 0.982 | 0.973 | 0.978 | 0.947 | 0.878 | 0.958 | 0.939 |
| ÉCLAIR | 936M | 0.032 | - | 0.968 | - | 0.962 | - | 0.974 | - | 0.950 | - | 0.980 | - |

Table 3. Accuracy comparison of various methods across different metrics in both English and Chinese (zh). Currently ÉCLAIR doesn’t train with additional chinese data or other form of multi-lingual data. The numbers in top row are obtained from GOT [52].

| Method | Size | Modality | Edit Distance↓ | BLEU↑ | METEOR↑ | Precision↑ | Recall↑ | F1-score↑ |
|-----------------|------|----------|----------------|-------|---------|------------|---------|-----------|
| Nougat-Base [6] | 350M | All | 0.071 | 0.891 | 0.930 | 0.935 | 0.928 | 0.931 |
| | | Text | 0.058 | 0.912 | 0.946 | 0.962 | 0.953 | 0.957 |
| | | Math | 0.128 | 0.569 | 0.754 | 0.765 | 0.766 | 0.765 |
| | | Tables | 0.211 | 0.697 | 0.791 | 0.754 | 0.807 | 0.780 |
| ÉCLAIR | 963M | All | 0.026 | 0.952 | 0.998 | 0.970 | 0.970 | 0.970 |
| | | Text | 0.015 | 0.979 | 0.996 | 0.992 | 0.990 | 0.990 |
| | | Math | 0.123 | 0.679 | 0.934 | 0.858 | 0.860 | 0.853 |
| | | Tables | 0.064 | 0.871 | 0.992 | 0.918 | 0.916 | 0.916 |

Table 4. Evaluation of the Nougat-Base model on the Nougat validation set (as reported in [6]), and of ÉCLAIR pre-trained on the arXiv-5M dataset and validated on the corresponding validation set. Note: We do not aim to provide a direct comparison between Nougat and ÉCLAIR here due to the concerns discussed in Section 3.2.

| Pre Training | Fine Tuning | Repetition Penalty= 1.1 | Counting F1** ↑ | WER ↓ |
|--------------|-------------|-------------------------|-----------------|-------|
| ✓ | ✗ | ✗ | 0.663 | 0.264 |
| ✓ | ✓ | ✗ | 0.925 | 0.151 |
| ✓ | ✓ | ✓ | 0.934 | 0.142 |

Table 5. Comparison of ÉCLAIR on DROBS before and after the fine-tuning stage, and also with and without a repetition-penalty (after the fine-tuning stage).

excelling in Edit Distance and Recall, where it achieves the best scores among the models compared. Notably, ÉCLAIR outperforms several larger models, such as Qwen-VL-Max (72B) [3], Fox (1.8B) [27], despite having a significantly smaller parameter size. While GOT (580M) [52] and ÉCLAIR exhibit similar accuracy, ÉCLAIR employs a decoder that is approximately half the size of GOT’s decoder.

3.2. Extraction of Formulas and Tables

In this section, we evaluate the extraction quality of ÉCLAIR on some important semantic-classes: formulae, tables and text. The latter consists of all semantic classes excluding formulae and tables. We report our findings on the validation set (10,000 samples) associated with the arXiv-5M dataset in Table 4. Our results demonstrate good extraction quality overall, with table and math elements being harder for ÉCLAIR to transcribe compared to text elements. Note that these metrics are reported for ÉCLAIR pre-trained on arXiv-5M (i.e., prior to fine-tuning).

Since other methods such as Nougat [6] cannot be directly compared to ÉCLAIR on our validation set owing to non-trivial differences in their output formatting styles, we cannot provide a direct comparison here. However, since Nougat and the pre-trained ÉCLAIR model are both trained on academic documents and evaluated on data sampled from arXiv, we find it useful to present the extraction quality of Nougat on these categories on their own validation set as a point of reference. These results are also summarized in Table 4. We observe similar trends in Nougat’s extraction quality for math and table elements, as discussed

above.

3.3. Document Object Detection

We evaluate the accuracy of detection of semantic text blocks of ÉCLAIR on the DocLayNet benchmark. Following [4], we fine-tune ÉCLAIR solely on DocLayNet for 50k steps to ensure that the bounding box class labels are not biased by labeling styles of other datasets (such as merging of several header and footer boxes). In order to compare to SOTA methods that report coco-mAP, we report the same metric using class token logits for ranking the predicted bounding boxes. We note, however, that being an autoregressive generator, ÉCLAIR remains in inherent disadvantage on coco-mAP metric compared to standard detectors due to it predicting bounding boxes and classes inline with the text in reading order, leading to inability of over-prediction to 100 bounding boxes assumed by coco-mAP. On the other hand, this results in ÉCLAIR not requiring non-maximum-suppression or threshold selection. For this reason, previous autoregressive object detectors adopt various tricks to improve the recall at low precision [8]. We also follow this approach and adopt sequence augmentation [8] with noisy and duplicate bounding boxes as well as sampling of top-k class labels from each predicted bounding box during inference for reporting coco-mAP. The comparison with SOTA methods is presented in Table 6, where we compare to reported Mask R-CNN [14] metrics and reproduced SwinDocSegmenter [4]. As can be seen, ÉCLAIR is competitive even compared to specialized object detectors.

Nevertheless, in agreement with previous works on autoregressive detection [2], we find mAP to be a suboptimal metric for such scenario. We provide further discussion on the evaluation metrics in the supplementary material, with more detailed evaluation of ÉCLAIR and competing methods.

| Classes | Mask-RCNN [14] | SwinDoc Segmenter[4] | ÉCLAIR |
|-------------|-------------------|-------------------------|-------------|
| Caption | 71.5 | 83.5 | 83.5 |
| Footnote | 71.8 | 67.8 | 66.9 |
| Formula | 63.4 | 64.2 | 65.7 |
| List-item | 80.8 | 84.1 | 79.0 |
| Page-footer | 59.3 | 65.1 | 62.0 |
| Page-header | 70.0 | 71.3 | 70.7 |
| Picture | 72.7 | 85.6 | 76.9 |
| Sec-header | 69.3 | 68.0 | 67.0 |
| Table | 82.9 | 86.0 | 77.6 |
| Text | 85.8 | 84.5 | 82.0 |
| Title | 80.4 | 66.8 | 82.0 |
| All | 73.5 | 75.2 | 73.9 |

Table 6. COCO-mAP (with defaults IoU=0.5:0.95, area=all, maxDets=100) on DocLayNet Benchmark.

3.4. LLM Benchmark

ÉCLAIR enables content extraction from PDFs, PPTs, and other scanned documents to meet the growing demands for high-quality data to train large language models (LLMs) [11, 34, 38, 48]. Unlike conventional extraction tools, e.g., PyMuPDF4LLM [19], ÉCLAIR is engineered to preserve semantic integrity and textual coherence. In this section, we compare the effectiveness of ÉCLAIR and PyMuPDF4LLM [19] for this task. We do this by training the Nemotron-8B LLM model from scratch on the text extracted by both of these methods from a common set of PDF documents, and compare the trained models on the Massive Multitask Language Understanding (MMLU) [15] benchmark, an average of multiple other benchmark scores including: ARC-Easy and ARC-Challenge [9], HellaSwag [53], OpenBooxQA [32], PIQA [5], RACE [25], WinoGrande [44], TriviaQA [21]. The results of this experiment, summarized in Table 7, highlight ÉCLAIR’s effectiveness in extracting high quality training data for improved LLM accuracy. Details about the training setup and post-processing steps for ÉCLAIR can be found in the supplementary material.

| Method | Tokens Extracted (B) | Other Bench | |
|------------------|----------------------------|----------------|-------------|
| | MMLU ↑ | Avg ↑ | |
| PyMuPDF4LLM [19] | 43.6 | 37.2 | 55.72 |
| ÉCLAIR | 55.1 | 39.1 | 56.7 |

Table 7. Comparison of the Nemotron-8B accuracy when trained on data extracted with ÉCLAIR or PyMuPDF4LLM [19].

3.5. Multi-token Inference

An important shortcoming of autoregressive models, including those targeted at OCR applications, is the large number of decoding steps necessary for text extraction, resulting in slow inference speed. In a standard autoregressive decoding formulation, each subsequent l^{th} token in the sequence is decoded incrementally, based on the context of $t_0 : t_{l-1}$ tokens. For text-dense images, such as documents, this results in a large amount of decoding steps, at least equal to the number of tokens in the sequence.

To mitigate this, we investigate multi-token generation as an alternative inference method. Instead of next-token prediction, we train ÉCLAIR to predict n subsequent tokens at a single step, and therefore reduce the number of necessary decoding steps by a factor of n .

Specifically, for predicting n tokens simultaneously, during training we introduce $n - 1$ new linear layers on top of the final hidden state of the decoder. The output of each of these linear layers is subsequently input into the shared decoder head. During training, standard teacher forcing is applied, with next n tokens representing the groundtruth for

each corresponding context. During inference, we greedily decode the sequence n tokens at a time. We do not perform token verification [7] but rely on purely greedy decoding in the interest of maximal throughput at batch-inference.

Following prior work [13], we additionally perform experiments where we only consider the first predicted token during inference while the rest are discarded, while during training next- n tokens are predicted. In other words, we evaluate n -token trained models at next-token prediction.

We experimented with 2-, 3-, or 4-token prediction and the results are reported in Table 8, where $\frac{tkn}{step}$ corresponds to the number of tokens kept from multi-token prediction at a single decoding step during inference. As can be seen, multi-token ÉCLAIR is matching or outperforming the baseline at 2 or 3 tokens, which is equivalent to around 2x inference speed increase. At 4 tokens, the accuracy degrades. We nevertheless find that keeping only the first token during inference is a valid strategy for improving OCR metrics for any of the variants. We additionally report the inference speed of multitone ÉCLAIR as well as competing methods, with average time per image on DROBS test set. As per-image speed of each method is related to its text extraction capabilities (due to the possibility of hallucination loops), we also report the average time per 100 tokens. As can be seen, multi-token approach allows ÉCLAIR to outperform competing methods speed-wise, despite it being a bigger model parameter-wise. The details on evaluation protocol can be found in the supplementary material.

| Method | $\frac{tkn}{step}$ | WER ↓ | F1 ↑ | $\frac{sec}{img} \downarrow$ | $\frac{sec}{100} \downarrow$ |
|-------------|--------------------|-------------|-------------|------------------------------|------------------------------|
| Nougat [6] | 1 | - | - | 4.7 | 0.41 |
| GOT [52] | 1 | 0.25 | 0.82 | 9.8 | 0.90 |
| ÉCLAIR | 1 | 0.14 | 0.93 | 3.8 | 0.42 |
| ÉCLAIR-2tkn | 2 | 0.13 | 0.94 | 2.5 | 0.31 |
| | 1 | 0.12 | 0.95 | 3.8 | 0.42 |
| ÉCLAIR-3tkn | 3 | 0.15 | 0.92 | 1.77 | 0.23 |
| | 1 | 0.13 | 0.94 | 3.8 | 0.42 |
| ÉCLAIR-4tkn | 4 | 0.17 | 0.90 | 1.32 | 0.20 |
| | 1 | 0.14 | 0.94 | 3.8 | 0.42 |

Table 8. Results and speed of multi-token models and competing methods. We report the average speed per image on DROBS test set ($\frac{sec}{img}$), and speed per 100 tokens ($\frac{sec}{100}$). These values are obtained from a PyTorch-based inference pipeline on an NVIDIA H-100 GPU.

4. Related Work

Document Understanding Models Models like LayoutLMv3 [18] excel in parsing complex documents for tasks such as layout analysis and visual question answering. However, they rely heavily on pre-extracted text, im-

ages, and bounding boxes, forming a brittle pipeline that can be error-prone due to its dependence on external systems. SwinDocSegmenter [4] and specialized variants of YOLO [20] have been trained for document-specific detection tasks without requiring additional inputs. While they effectively detect objects within documents, they generally do not output any text associated with these objects, lacking integrated OCR capabilities.

Object Detection in Documents is crucial for identifying and localizing elements within documents, aiding tasks like OCR and determining reading order. Traditional models such as Faster R-CNN [43] and Mask R-CNN [14] have been adapted for document analysis, effectively detecting and segmenting components like text blocks, images, and tables. Despite their success, these models typically do not provide textual content alongside the detected objects, limiting their usefulness for comprehensive document understanding.

End-to-End OCR-Free Models that do not depend on external OCR systems have gained attention. Donut [23] introduced a transformer-based encoder-decoder architecture pre-trained on general text documents. Building on this, Nougat [6] extended training to scientific documents, outputting structured markdown with L^AT_EX tables and equations. GOT [52] focused on enhancing the recognition of specialized documents containing molecular formulas, sheet music, and charts. Kosmos-2.5 [31] incorporated both markdown and plain text data with bounding boxes, introducing a prompt structure that allows users to choose between different output formats. However, these models may require compromises in prompt structures or may not handle a wide variety of document layouts effectively. Our proposed model, ÉCLAIR, is specifically trained to handle a greater variety of document layouts without requiring compromises in the prompt structure.

Multimodal Large Language Models like QwenVL [3], GPT-4O [36] and Claude [1] have demonstrated impressive OCR and document understanding capabilities, including the extraction of complex equations and tables in structured formats. While powerful, these models are large and computationally expensive, making them impractical for scaling to millions of pages. In contrast, ÉCLAIR is a sub-1B parameter model optimized for inference speed with multi-token decoding.

5. Conclusion

In this work, we have presented ÉCLAIR, a general-purpose end-to-end text-extraction model. ÉCLAIR is the first model that extracts structured text, bounding boxes and semantic classes all at the same time. We are releasing a new benchmark dataset DROBS to capture the variety of layouts of various online documents and have shown that ÉCLAIR outperforms all current competitors on this benchmark. Ad-

ditionally, we investigate and provide a technique to improve the inference time for ÉCLAIR. We hope that this will aid the OCR community in improving document-based text extraction, and benefit the LLM community by increasing the availability of previously unseen text data for training.

6. Acknowledgments

We would like to thank Shrimai Prabhumoye, John Kamalu, Brandon Norick, Mostafa Patwary, Thomas Breuel, Osvald Nitski, Ushnish De, Mehrzad Samadi, Guilin Liu, Zhiding Yu, Mike Ranninger, Teo Ene, Mohammad Shoeybi and Bryan Catanzaro, for their feedback and valuable discussions.

References

- [1] Anthropic. Claude 3: A new era in ai with advanced reasoning and intelligence. <https://www.anthropic.com/clause>, March 2024. Claude 3 is a family of large language models including Haiku, Sonnet, and Opus variants, designed for advanced reasoning, analysis, and complex problem-solving. 8
- [2] Armen Avetisyan, Christopher Xie, Henry Howard-Jenkins, Tsun-Yi Yang, Samir Aroudj, Suvam Patra, Fuyang Zhang, Duncan Frost, Luke Holland, Campbell Orme, Jakob Engel, Edward Miller, Richard Newcombe, and Vasileios Balntas. Scenescrit: Reconstructing scenes with an autoregressive structured language model, 2024. 7, 14
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 6, 8
- [4] Ayan Banerjee, Sanket Biswas, Josep Lladós, and Umapada Pal. SwinDocSegmenter: An end-to-end unified domain adaptive transformer for document instance segmentation. In *International Conference on Document Analysis and Recognition*, pages 307–325. Springer, 2023. 7, 8, 16
- [5] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020. 7
- [6] Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. Nougat: Neural optical understanding for academic documents. In *International Conference on Learning Representations*, 2024. 1, 3, 4, 6, 8, 13
- [7] Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D Lee, Deming Chen, and Tri Dao. Medusa: Simple llm inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv:2401.10774*, 2024. 8
- [8] Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. *arXiv preprint arXiv:2109.10852*, 2021. 7, 16
- [9] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018. 7
- [10] Alexey Dosovitskiy, Lukas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 2
- [11] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Kordovaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasudevan Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidor, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonnia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparth, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra,

Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Asperegn, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madiyan Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev,

Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghatham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wencheng Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. 7

- [12] The Common Crawl Foundation. <https://commoncrawl.org/>. 4
- [13] Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Roziere, David Lopez-Paz, and Gabriel Synnaeve. Better and faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*, 2024. 8
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 7, 8
- [15] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020. 7
- [16] Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, et al. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding. *arXiv preprint arXiv:2403.12895*, 2024. 6

- [17] Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, and Jingen Zhou. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding, 2024. 13
- [18] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4083–4091, 2022. 8
- [19] Artifex Software Inc. Pymupdf4llm: A python package for extracting pdf content in markdown format, 2024. Version 0.0.17. 7
- [20] Glenn Jocher. Ultralytics yolov5, 2020. 8
- [21] Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, Vancouver, Canada, July 2017. Association for Computational Linguistics. 7
- [22] Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*, 2019. 5, 13
- [23] Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeong Yeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *European Conference on Computer Vision*, pages 498–517. Springer, 2022. 1, 8
- [24] Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von Werra, and Harm de Vries. The stack: 3 tb of permissively licensed source code. *Preprint*, 2022. 4, 13
- [25] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017. 7
- [26] Chenglong Liu, Haoran Wei, Jinyue Chen, Lingyu Kong, Zheng Ge, Zining Zhu, Liang Zhao, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Focus anywhere for fine-grained multi-page document understanding. *arXiv preprint arXiv:2405.14295*, 2024. 6
- [27] Chenglong Liu, Haoran Wei, Jinyue Chen, Lingyu Kong, Zheng Ge, Zining Zhu, Liang Zhao, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Focus anywhere for fine-grained multi-page document understanding, 2024. 5, 6
- [28] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742, 2020. 2, 13
- [29] Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. Textmonkey: An ocr-free large multimodal model for understanding document. *arXiv preprint arXiv:2403.04473*, 2024. 6
- [30] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 13
- [31] Tengchao Lv, Yupan Huang, Jingye Chen, Lei Cui, Shuming Ma, Yaoyao Chang, Shaohan Huang, Wenhui Wang, Li Dong, Weiyao Luo, et al. Kosmos-2.5: A multimodal literate model. *arXiv preprint arXiv:2309.11419*, 2023. 1, 3, 4, 8
- [32] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018. 7
- [33] Ahmed Nassar, Nikolaos Livathinos, Maksym Lysak, and Peter Staar. Tableformer: Table structure understanding with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4614–4623, 2022. 4, 13
- [34] Nvidia, : Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H. Anh, Pallab Bhattacharya, Annika Brundyn, Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, Sirshak Das, Ayush Dattagupta, Olivier Delalleau, Leon Derczynski, Yi Dong, Daniel Egert, Ellie Evans, Aleksander Ficek, Denys Fridman, Shaona Ghosh, Boris Ginsburg, Igor Gitman, Tomasz Grzegorzek, Robert Hero, Jining Huang, Vibhu Jawa, Joseph Jennings, Aastha Jhunjhunwala, John Kamalu, Sadaf Khan, Oleksii Kuchaiev, Patrick LeGresley, Hui Li, Jiwei Liu, Zihan Liu, Eileen Long, Ameya Sunil Mahabaleshwarkar, Somshubra Majumdar, James Maki, Miguel Martinez, Maer Rodrigues de Melo, Ivan Moshkov, Deepak Narayanan, Sean Narendiran, Jesus Navarro, Phong Nguyen, Osvald Nitski, Vahid Noroozi, Guruprasad Nutheti, Christopher Parisien, Jupinder Parmar, Mostofa Patwary, Krzysztof Pawelec, Wei Ping, Shrimai Prabhumoye, Rajarshi Roy, Trisha Saar, Vasanth Rao Naik Sabavat, Sanjeev Satheesh, Jane Polak Scowcroft, Jason Sewall, Pavel Shamis, Gerald Shen, Mohammad Shoeybi, Dave Sizer, Misha Smelyanskiy, Felipe Soares, Makesh Narsimhan Sreedhar, Dan Su, Sandeep Subramanian, Shengyang Sun, Shubham Toshniwal, Hao Wang, Zhilin Wang, Jiaxuan You, Jiaqi Zeng, Jimmy Zhang, Jing Zhang, Vivienne Zhang, Yian Zhang, and Chen Zhu. Nemotron-4 340b technical report, 2024. 7
- [35] Ocropus. String normalization. <https://github.com/ocropus/ocropus4-eval/blob/main/ocroeval/eval.py>. 5
- [36] OpenAI. GPT-4o: Large language model. <https://openai.com/index/hello-gpt-4o>, 2024. Accessed: 2024-11-12. 8
- [37] Pandoc. A universal document converter. <https://pandoc.org/index.html>. 4
- [38] Jupinder Parmar, Shrimai Prabhumoye, Joseph Jennings, Mostofa Patwary, Sandeep Subramanian, Dan Su, Chen Zhu, Deepak Narayanan, Aastha Jhunjhunwala, Ayush Dattagupta, Vibhu Jawa, Jiwei Liu, Ameya Mahabaleshwarkar, Osvald Nitski, Annika Brundyn, James Maki, Miguel Martinez, Jiaxuan You, John Kamalu, Patrick LeGresley, Denys Fridman, Jared Casper, Ashwath Aithal, Oleksii Kuchaiev, Mohammad Shoeybi, Jonathan Cohen, and Bryan Catanzaro. Nemotron-4 15b technical report, 2024. 7
- [39] Keiran Paster, Marco Dos Santos, Zhangir Azerbayev, and Jimmy Ba. Openwebmath: An open dataset of high-quality mathematical web text. *arXiv preprint arXiv:2310.06786*,

2023. 17
- [40] Birgit Pfitzmann, Christoph Auer, Michele Dolfi, Ahmed S Nassar, and Peter Staar. DocLayout: A large human-annotated dataset for document-layout segmentation. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 3743–3751, 2022. 2, 4
- [41] NLTK Project. Source code for nltk.metrics.scores. https://www.nltk.org/_modules/nltk/metrics/scores.html. 5
- [42] Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. AM-RADIO: Agglomerative vision foundation model reduce all domains into one. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12490–12500, 2024. 2, 13
- [43] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6):1137–1149, 2016. 8
- [44] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021. 7
- [45] Sebastian Spiegler. Statistics of the Common Crawl Corpus 2012, 2013. 4
- [46] Ray Smith. An overview of the tesseract ocr engine. In *Ninth international conference on document analysis and recognition (ICDAR 2007)*, volume 2, pages 629–633. IEEE, 2007. 13
- [47] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poultton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science, 2022. 3
- [48] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin Mc-Donell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kociský, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotrata, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical size, 2024. 7
- [49] L. Vincent. Google book search: Document understanding on a massive scale. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, volume 2, pages 819–823, 2007. 4
- [50] Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Vary: Scaling up the vision vocabulary for large vision-language models, 2023. 6
- [51] Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, En Yu, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Small language model meets with reinforced vision vocabulary, 2024. 6
- [52] Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, et al. General OCR theory: Towards OCR-2.0 via a unified end-to-end model. *arXiv preprint arXiv:2409.01704*, 2024. 1, 2, 4, 6, 8
- [53] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019. 7

Supplementary Material

S1. Architecture Details

The entire architecture has a total of 937M parameters.

Vision Encoder Our vision encoder is initialized from the RADIO [42] ViT-H /16 model (657M parameters). The inputs to the encoder are images of resolution 1280×1024 . We resize and pad the original image to this size, while maintaining the aspect ratio. This results in 80×64 patches.

Vision Neck Following the observation in [17] that text is more correlated within a line rather than across lines, we compress the sequence with a horizontal convolutional kernel of size 1×4 and stride 1×4 . Given an input image of size 1280×1024 which produces 5120 patches, this reduces the sequence length to 1280.

Decoder We use an MBart [28] decoder with 10 layers (279M parameters) trained from scratch. The maximum sequence length for the decoder is set to 3584.

S2. Training & Inference

We follow a two-step training strategy, similar to Nougat [6]. We pre-train ÉCLAIR on the arXiv-5M dataset and fine-tune it further on all the training datasets. In both stages, we use the AdamW optimizer [30] and train for 130000 iterations with an effective batch size of 128. During pre-training, we use a constant learning rate of 2×10^{-5} with 5000 linear warmup steps. For fine-tuning, we employ a constant learning rate of 8×10^{-6} with 500 linear warmup steps. During inference, we use a greedy decoding strategy with a repetition penalty [22] of 1.1.

Inference Speed Comparison To measure the speed of competing methods compared to ÉCLAIR and its multi-token variants, we utilize their publicly shared huggingface pipelines and pre-trained models. We report the speed of exclusively the model excluding any pre-processing, i.e., `model.generate()` step. All models are evaluated with the provided pre-processing pipelines and the recommended parameters are used. We perform evaluation in bf16 for fair comparison. Additionally, since the speed of a model is affected by potential hallucinations that can lead to infinite loops bound by maximal sequence length, we also report the speed per 100 tokens. We can see that in this case speed of ÉCLAIR is almost identical to that of Nougat, as their decoder architecture is the same. At the same time, while ÉCLAIR has the most parameters out of competing methods, being front-heavy, it enjoys faster inference speed which is primarily dictated by the size of the decoder. For multi-token inference, we adopt a simple greedy decoding strategy (as opposed to more complex speculative strategies) for ensuring maximal throughput in batch mode during inference. Each linear layer has 1024 nodes.

S3. Datasets

S3.1. Reading Order

The reading order for the bounding boxes in a page is the order in which the contents of these boxes would be visited if a person was reading the page out aloud. We only include relevant text-like semantic classes in this order, i.e., Text, Section-Header, List-Item, Title and Formula. See Figure S1 for an illustration. Additionally, we reorder all Page-Header elements to be at the start of the page, and all Footnote, Page-Footer, Picture, Table and Caption elements to be at the end.

S3.2. Pre-processing steps for the training datasets

SynthTabNet The SynthTabNet dataset [33] consists of images of tabular data along with HTML annotations. We convert these HTML annotations to L^AT_EX.

G1000 We obtain the annotations for G1000 using Tesseract [46].

README We create a dataset of project documents by sampling README files from the Stack [24]. To normalize the markdown source code, we convert it to HTML using Pandoc ¹ and render the converted HTML as PDF using wkhtmltopdf ². To obtain the ground truth for each page within the rendered PDFs, we first convert the HTML source back into our markup format and adopt Nougat’s [6] data processing pipeline³ to split and align the markdown with each page.

Table Auto-Labeling The DocLayNet dataset and the human-annotated CommonCrawl samples do not contain L^AT_EX annotations for tables and equation blocks. When we mask these elements out (both in the input image and the targets) in the fine-tuning stage owing to lack of ground truth, we observe that the model tends to mis-classify tables as pictures during inference. We hypothesize that this is caused by the lack of tables in the visually diverse subset of the training data, rendering such samples out-of-distribution (OOD) during inference. We address this by fine-tuning ÉCLAIR on table crops from SynthTabNet and arXiv-5m and using it to auto-label table crops from DocLayNet and the human-annotated samples.

S4. Post-processing: Hallucinations and Bad-box detection.

Similar to Nougat [6], we observe that ÉCLAIR sometimes degenerates into repetition loops wherein it repeats the same phrase, sentence or paragraph over and over again at the end of the prediction. Nougat [6] detect hallucinations by tracking a moving average of logits and flagging the outputs when a certain threshold is reached.

¹<https://pandoc.org/>

²<https://wkhtmltopdf.org/>

³<https://github.com/facebookresearch/nougat>

landscaping restoration plan will be available at www.naperville.il.us/uverse.aspx in early 2009 prior to the spring planting season.

Additional information regarding the upcoming construction work can be obtained at www.naperville.il.us/uverse.aspx. All interested parties are invited to visit the city's Web site, which will continuously track the progress of AT&T's installation efforts. For more information on AT&T's U-verse service, visit <http://uverse.att.com>.

November NAHC Meeting to Focus on Schools

We're finalizing the details, but mark your calendar for the Confederation November meeting where we will focus on School Districts 203 and 204. With construction projects ongoing in both districts, a new superintendent on tap for District 203, and School Board elections in the Spring, there will be a lot to talk about. As the largest line item on our Property Tax Bills, schools are of interest to homeowners with and without children in the districts!

Right now we're planning on taking this meeting on the road to a site in one of the School Districts - keep an eye on the Confederation website, www.napervillehomeowners.com, for more details!

NAHC Candidate Forums – Coming In February

We were all set to hold Candidate Forums for the November State House and County Board races impacting Naperville, but it is kind of hard to hold a forum without the candidates. While we had some responses, we did not get the type of bipartisan area wide reaction needed to hold a fair and informative forum.

So, stay tuned – we are planning forums for the late Winter, early Spring, focused on the City Council, Park District, and School District April elections. Hopefully these candidates will be willing to spend the time to hear what is on the minds of the homeowners.

Our format will include a brief opportunity for the candidates to introduce themselves followed by our moderator posing questions prepared by the NAHC and submitted by our audience. Our focus will be on issues important to Homeowners. The forums will be held in the City Council Chambers of the Naperville Municipal Center and will be broadcast live on Naperville cable channels 6 (WOW) and 10 (Comcast). If AT&T U-Verse is a reality by then, we'll try to figure out how they are handling public access broadcasting and provide you with information on where to find it on their system as well (from what we hear, it won't be quite as easy as clicking on a channel number).

Watch the website and future newsletters for more details.

Neither Snow nor Rain nor Gloom of Night?

Are you having problems with mail delivery? Some subdivisions with curbside individual boxes are reporting not receiving their daily delivery when there are cars parked (legally) on the street blocking access to the box. We're trying to figure out if this is a widespread problem or indicative of a couple of letter carriers who don't want to get out of their trucks. Send an e-mail to nahc@napervillehomeowners@wownway.com and let us know.

10
The NAHC Monthly Newsletter is a tool for communicating to and within our member associations. Please feel free to reproduce any newsletter content in your own Association newsletters.

Medicaid Formula Narrows Differences in Some States' Funding Ability and Widens Differences in Others

chosen as a proxy for a state's ability to fund public services. Consistent with the purpose described in the formula's legislative history, PCI is used as a proxy for both state resources and the low-income population. As a state's PCI increases, relative to the national average, the formula provides for a decreasing federal matching rate, meaning the federal government shares a smaller portion of a state's costs. By statute, the federal matching rate may range from 50 percent to 83 percent.¹³ The formula's multiplier, currently 0.45, represents the state's share of its total Medicaid costs for a state with PCI equal to the national average, and the federal government thus pays a 55 percent share of total costs.

Formula Reduces Overall Differences in States' Funding Ability by 20 Percent

The Medicaid formula reduces by 20 percent the differences among states in their ability to fund program services, compared with the national average funding ability. While the formula narrows differences for 30 states, making the average difference in funding ability smaller, it moves 21 states farther away from the national average, making the average difference wider. These 21 states include 3 that are among those with the largest populations in poverty—California, Florida, and New York. Because of the formula's current structure, in many instances, two states devoting the same proportion of their own resources toward funding Medicaid services are unable, after receiving federal matching aid, to spend the same amounts per person in poverty, adjusted for cost differences related to age and geographic location.

Because state resources, numbers of people in poverty, and the cost of serving this population vary widely across the states, there also are wide differences in states' funding ability to fund health care services. Considering these indicators of state funding ability, Alaska has the highest funding ability—exceeding the national average by 119 percent—and Mississippi has the lowest funding ability—46 percent below the national average, as measured using states' TTR and the number of people in poverty, adjusting the poverty count for age and geographic cost differences (see fig. 1). Nationwide, the average difference between a state's funding ability and

¹³In fiscal year 2003, Mississippi had the highest federal matching rate of any state—76.6 percent.

Text**Section-header****Picture****Footnote****Page-footer****Page 6****GAO-03-620 Medicaid Formula**

Figure S1. Illustrations of reading order over relevant text-like elements, i.e. Text, Section-header, List-item, Title and Formula. Other semantic classes (such as Picture, Footnote and Page-footer in the examples here) are not included in the reading order of the main body. (Note: We are not showing all the classes)

For ÉCLAIR, we adopt a simple hallucination mitigation strategy to filter out such occurrences: the inference-time prompt is always set to the Maximal Informative Prompt (MIP) and we do a strict syntax check on the resulting predictions to reject non-compliant boxes. Some examples of hallucinations detected and filtered out using this strategy are shown in Figure S2. We also enforce the spatial and categorical validity of the remaining boxes by verifying that the bottom-right corner of each bounding box exceeds the top-left corner and that classes conform to a validated schema. By implementing this layered filtering strategy, we observe a substantial reduction in model hallucinations.

S5. Object detection

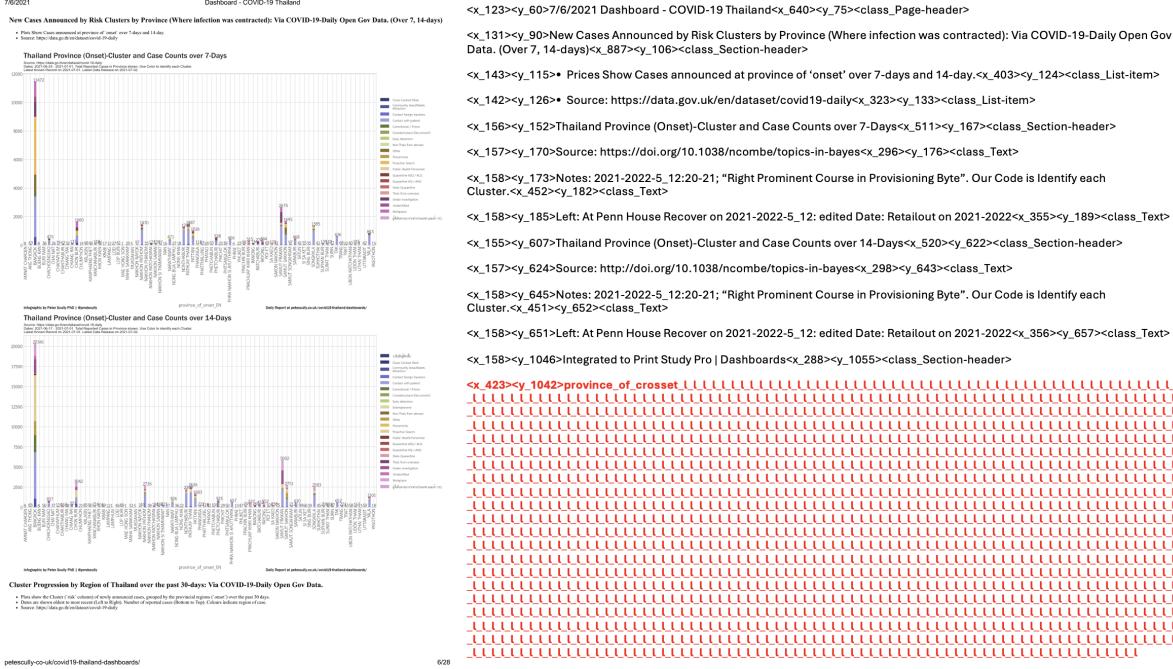
S5.1. mAP - Back to Simple Metrics

In line with previous works on autoregressive detection [2], we find mAP to be inherently unfavorable to end-to-end models like ÉCLAIR. Dedicated object detectors generally predict a set of bounding boxes of a fixed size as a raw output, where each bounding box is associated with a confidence score. Consequently, it is possible to control the recall-precision trade-off of the model by adjusting the confidence threshold by which the raw predictions are filtered.

Naturally, this results in a possibility of achieving a high recall for the model, albeit at low precision. Instead, ÉCLAIR predicts a set of bounding boxes in line in the output stream of tokens, requiring no filtering or thresholding.

An inherent problem with our end-to-end detector is the absence of a score for detected boxes that could be used to rank them. Using the likelihood/logits of the initial coordinate tokens is not ideal, as they indicate the distribution over potential starting points rather than an independent probability. Similarly, class-token logits only provide a distribution over class choices, not the probability of the box's existence. Considering text tokens is also impractical, as they represent the actual text rather than the existence of the surrounding box. Therefore, our predictor does not generate a box score. On one hand, as there is no over-prediction, no subsequent filtering or post-processing (such as non-maximum suppression) as well as no score is necessary. On the other hand, this makes comparison on the average precision metric challenging, as when considering all of the predicted bounding boxes jointly, only a single recall level exists for ÉCLAIR, making area calculation not meaningful.

Therefore, it can be seen that comparing ÉCLAIR against other works on the *mAP* metric poses challenges:



Program Status Report - CAPRA 2012

Printed: 2/16/2012 4:20:35 PM

Appendix-10

Figure S2. Examples of hallucinations in the ÉCLAIR predictions. The hallucinations (in this case, repetition loops), marked in red, are detected and filtered out by our hallucination-mitigation strategy.

1. AP is the area under the PR-Curve, which degenerates to a single point without the possibility to rank predictions, making the calculation of the area not meaningful.
 2. The COCO implementation assumes scores are unique. Identical scores (as in our case) lead to incorrect PR-Curves and inconsistent results.⁴

3. COCO mAP is computed per class independently (i.e. first separate classes, then match boxes). We propose to first match boxes over all classes and then compute the per-class precision/recall, which allows us to plot a confusion matrix, to better visualize problematic cases. Pre-

⁴See https://github.com/MiXaiLL76/faster_coco_

cision/recall and confusion matrix can still be averaged over multiple IoUs (0.5 to 0.95) like the COCO framework does. For an example of this on the DocLayNet evaluation dataset, see Fig. S3.

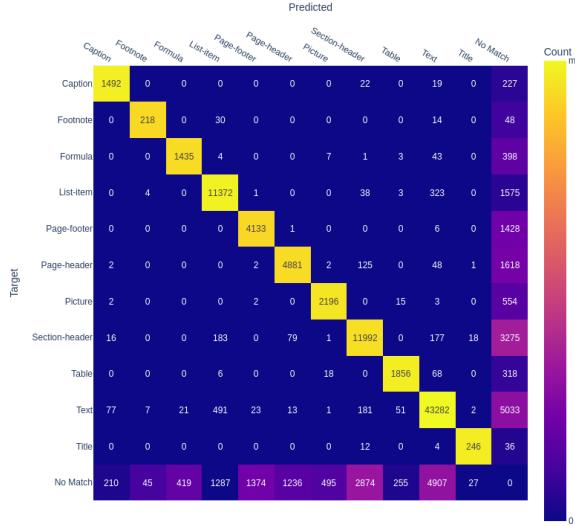


Figure S3. Confusion matrix for ÉCLAIR boxes matched with ground truth on the DocLayNet evaluation dataset, averaged over thresholds of $IoU \geq \{0.5, 0.55, \dots, 0.9, 0.95\}$.

Figure S4 shows PR-Curves of individual classes for $IoU \geq 0.5$ and 1001 recall-bins vs predictions from SwinDocSegmenter [4]. For ÉCLAIR, scores are taken from class-token logits, which apparently are not a good separator for true positives vs false positives, compared to SwinDocSegmenter’s scores. ÉCLAIR does not predict additional low-score boxes, causing curves to drop to zero when all boxes are included. SwinDocSegmenter curves are smoother due to over-predicted boxes, allowing proper score thresholds for each class.

The second part of Fig. S4 shows averaged PR-Curve over all classes for $IoU \geq 0.5$ and 1001 recall-bins. Steps in the ÉCLAIR curve come from averaging over classes, cutting off at the mean precision (mP)/mean recall (mR) point.

In the main paper, we have shown ÉCLAIR to be competitive on mAP metric when using methods from the earlier autoregressive object detection literature [8] such as sequence augmentation and top-k selection that are able to increase the recall of the model at the cost of the low precision. Still, we show ÉCLAIR to be a strong predictor without the necessity to introduce such augmentations in the next subsection.

S5.2. Comparisons at corresponding precision/recall levels

We additionally present a comparison of ÉCLAIR to SwinDocSegmenter in a point-to-point manner. Given ÉCLAIR’s precision and recall value of each class when considering all the predicted bounding boxes, we compare the recall and precision of SwinDocSegmenter at corresponding precision and recall levels (that is, considering equal recall, we want to evaluate which model achieves higher precision, and vice versa). We evaluate each class separately, and report the mean precision and mean recall for both methods. Here, we train a standard ÉCLAIR model without sequence augmentation and perform no filtering or post-processing on the predicted bounding boxes, reporting the mean precision and mean recall of the full prediction set. The results are summarized in Tab. S1. As can be seen, ÉCLAIR achieves higher precision at the same recall as SwinDocSegmenter, as well as higher recall at the same precision point, for the vast majority of the classes, as well as on average.

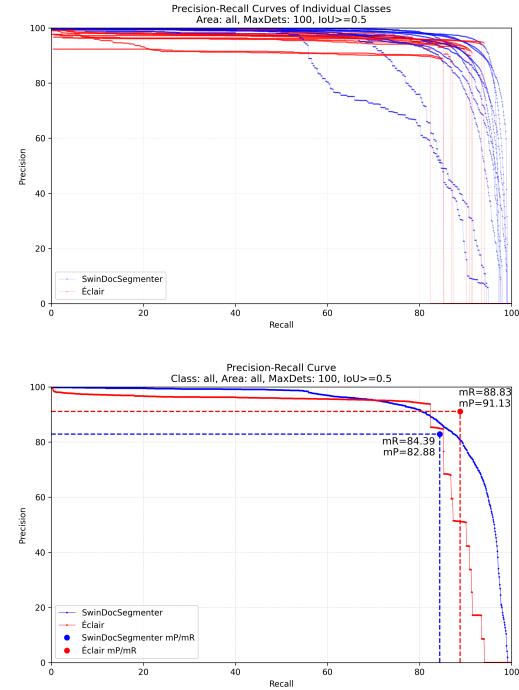


Figure S4. PR-Curves for individual classes and averaged over all classes for $IoU \geq 0.5$ and 1001 recall-bins evaluated on the DocLayNet evaluation dataset. For ÉCLAIR, the scores are taken from the class-token logits. The mean precision and recall are taken from Tab. S1.

| Class | $IoU \geq 0.5$ | | | |
|----------------|-----------------------|-------------|-------------|-------------|
| | mP | | mR | |
| | ÉCLAIR | SDS* | ÉCLAIR | SDS* |
| Caption | 89.2 | 83.2 | 91.2 | 89.1 |
| Footnote | 92.4 | 65.6 | 82.1 | 70.8 |
| Formula | 91.7 | 78.7 | 90.9 | 83.6 |
| List-item | 91.0 | 87.8 | 91.4 | 89.6 |
| Page-footer | 94.6 | 93.3 | 94.0 | 93.2 |
| Page-header | 93.4 | 96.7 | 86.8 | 91.1 |
| Picture | 86.9 | 92.8 | 85.2 | 90.9 |
| Section-header | 93.0 | 90.5 | 90.1 | 87.8 |
| Table | 88.4 | 87.8 | 85.2 | 84.8 |
| Text | 94.0 | 91.0 | 93.4 | 90.9 |
| Title | 87.8 | 44.2 | 87.0 | 56.6 |
| All | 91.1 | 82.9 | 88.8 | 84.4 |
| Class | $IoU \geq 0.5 : 0.95$ | | | |
| | mP | | mR | |
| | ÉCLAIR | SDS* | ÉCLAIR | SDS* |
| Caption | 82.8 | 72.1 | 84.7 | 79.5 |
| Footnote | 79.0 | 59.8 | 70.1 | 62.8 |
| Formula | 76.8 | 56.5 | 75.8 | 56.1 |
| List-item | 85.0 | 68.2 | 85.4 | 82.1 |
| Page-footer | 74.7 | 45.3 | 74.2 | 59.8 |
| Page-header | 78.6 | 64.0 | 73.0 | 68.4 |
| Picture | 80.7 | 88.6 | 79.1 | 85.3 |
| Section-header | 78.7 | 53.4 | 76.2 | 59.8 |
| Table | 84.9 | 84.3 | 81.8 | 81.5 |
| Text | 88.5 | 60.6 | 88.0 | 79.6 |
| Title | 83.2 | 20.6 | 82.3 | 49.7 |
| All | 81.2 | 61.2 | 79.2 | 69.5 |

Table S1. The mean precision and mean recall of ÉCLAIR for each class and the corresponding mean recall and mean precision of SwinDocSegmenter for the respective recall/precision on the PR-curve evaluated on the DocLayNet evaluation dataset. Computed for $IoU \geq 0.5$ (corresponding to Fig. S4) and for averaged thresholds of $IoU \geq \{0.5, 0.55, \dots, 0.9, 0.95\}$ (default for COCO metrics). *SDS: SwinDocSegmenter.

S6. LLM Training

We train both models with a total of 300B tokens obtained from a combination of sources. We ensure that the models are trained for 3.3 epochs of the tokens extracted using ÉCLAIR and PyMuPDF4LLM. These tokens are extracted from a common set of PDFs for both methods. The rest of the 300B training tokens come from various sources including CommonCrawl snapshots, Stack Exchange, OpenWebMath [39], PubMed Abstracts, PubMed Central, bioRxiv, SEC filings, Wikipedia and ArXiv data.

S6.1. Postprocessing and Joining Pages

To join pages and handle the positioning of floating objects, we follow these steps:

1. **Process Pages Individually:** Each page is processed separately. To manage paragraphs that span across pages, we need to carry open paragraphs over to next pages.
2. **Reassign Floating Objects:** Floating objects (e.g., im-

ages, tables, captions) are removed and captions are reassigned to their respective objects using Hungarian matching based on the Manhattan distance of the bounding boxes.

3. **Concatenate Pages:** Pages are concatenated while skipping sections like Table of Contents, Bibliography, and Indexes by detecting typical headings. Floating text blocks (e.g., Text and List-item) are merged based on specific rules, such as not ending with punctuation.
4. **Remove Markdown Formatting:** All markdown formatting is removed from the inner text to ensure consistency.
5. **Flush Floating Objects:** After processing each page, floating objects that are not part of the floating text are flushed to the output blocks.

S7. Examples of predictions

In this section, we present examples of predictions from ÉCLAIR on samples from the Common Crawl dataset. Figure S5 contains samples with tables, formulae, pictures and a variety of other elements.

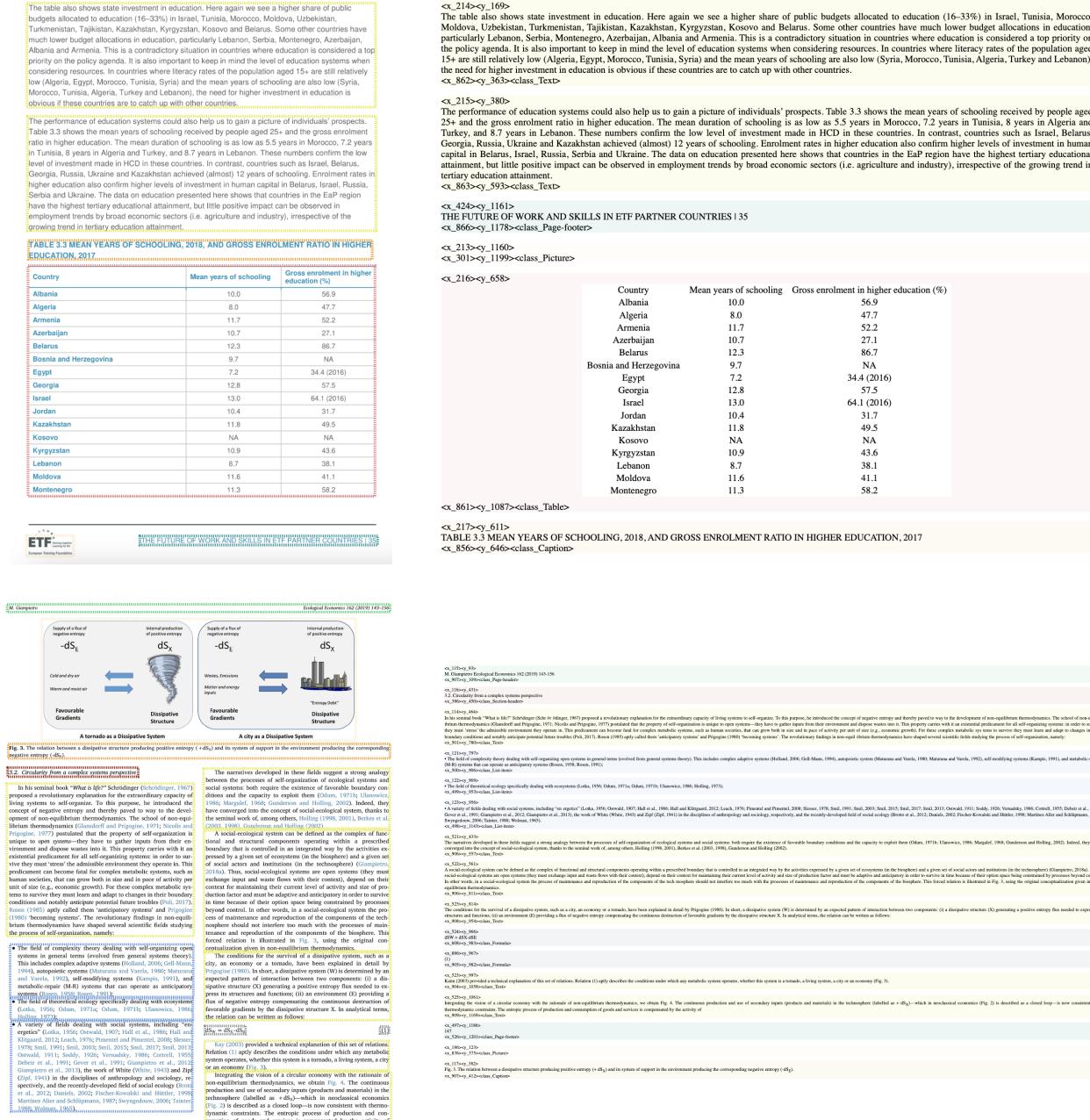


Figure S5. Examples of pages with tables, formulae and pictures. On the left, predicted bounding boxes superimposed on the original sample image. On the right, the corresponding full predictions.