

# LMDX: LANGUAGE MODEL-BASED DOCUMENT INFORMATION EXTRACTION AND LOCALIZATION

Vincent Perot<sup>1\*</sup> Kai Kang<sup>2</sup> Florian Luisier<sup>2</sup> Guolong Su<sup>1</sup> Xiaoyu Sun<sup>2</sup>  
 Ramya Sree Boppana<sup>2</sup> Zilong Wang<sup>3†</sup> Jiaqi Mu<sup>1</sup> Hao Zhang<sup>1</sup> Nan Hua<sup>1</sup>  
<sup>1</sup>Google Research <sup>2</sup>Google Cloud <sup>3</sup>University of California, San Diego  
 {vperot, kangkai, fluisier, gsu, sunse}@google.com  
 {ramyasreeb, jqmu, haozhangthu, nhua}@google.com  
 zlwang@ucsd.edu

## ABSTRACT

Large Language Models (LLM) have revolutionized Natural Language Processing (NLP), improving state-of-the-art on many existing tasks and exhibiting emergent capabilities. However, LLMs have not yet been successfully applied on semi-structured document information extraction, which is at the core of many document processing workflows and consists of extracting key entities from a visually rich document (VRD) given a predefined target schema. The main obstacles to LLM adoption in that task have been the absence of layout encoding within LLMs, critical for a high quality extraction, and the lack of a grounding mechanism ensuring the answer is not hallucinated. In this paper, we introduce *Language Model-based Document Information EXtraction and Localization* (LMDX), a methodology to adapt arbitrary LLMs for document information extraction. LMDX can do extraction of singular, repeated, and hierarchical entities, both with and without training data, while providing grounding guarantees and localizing the entities within the document. In particular, we apply LMDX to the PaLM 2-S LLM and evaluate it on VRDU and CORD benchmarks, setting a new state-of-the-art and showing how LMDX enables the creation of high quality, data-efficient parsers.

## 1 INTRODUCTION

The recent advent of transformers (Vaswani et al., 2017) and self-supervised pretraining procedures has led to significant progress in Visually Rich Document (VRD) Understanding. Within that field, the task of document information extraction (IE), which consists of extracting key entities within a semi-structured document (e.g. invoice, tax form, paystub, receipt, etc) given a predefined schema, has received a lot of attention from industry and academia due to its importance and wide applicability to intelligent document processing workflows. However, document information extraction still remains challenging for current generation systems. In particular, information in semi-structured forms is organized in complex layout across many possible templates, which requires understanding of the document context, spatial alignment among the different segments of text, and tabular arrangement of structured entities (e.g. line items on an invoice, deduction items on a paystub, etc.). Content on the document can be printed or handwritten, with scanning artefacts like rotation and contrast issues. Moreover, since some business automation workflows require certain level of accuracy, they are often integrated with human-in-the-loop interactions for auditing and correction of predictions, requiring knowing the precise location of extracted entities to make it a tractable task for a human rater. Finally, since a quasi-infinite number of document types exist, and that organizations have limited annotation resources, most parsers are built with very small amount of data.

From those complexities emerge the following desiderata of document information extraction systems: they should support high-quality extraction of singular, repeated, and hierarchical entities, while localizing those entities in the document, and doing so with very low or no amount of training data. So far, no publicly disclosed system has been able to address all of those desiderata.

\*Corresponding author

†Work done during Google internship

Many current approaches divide the problem in two stages: a text recognition/serialization step, usually achieved by an off-the-shelf Optical Character Recognition (OCR) service, followed by a parsing step, which finds the relevant entity values from the recognized text. Since the text serialization is imperfect, much attention has been given to fusing the text and layout together in the parsing step (Majumder et al., 2020; Garncarek et al., 2021; Hwang et al., 2021; Katti et al., 2018; Denk & Reisswig, 2019). Hong et al. (2021) proposes to encode the relative 2D distances of text blocks in the attention of the transformer, and learning from unlabeled documents with an area-masking strategy. Lee et al. (2022) proposes encoding the relative token positions with a graph neural network with edges constructed from a beta-skeleton algorithm. It further frames information extraction as a NER sequence tagging task with an IOB scheme (Ramshaw & Marcus, 1995; Palm et al., 2017) which allows them to localize the entities. However, IOB does not support extracting hierarchical entities, and is not robust to text serialization errors, where an entity is broken in disjoint segments.

Since text and layout do not contain all the information in the document (e.g. table boundaries, logos), leveraging the image modality has also been extensively explored (Xu et al., 2021; Lee et al., 2023; Appalaraju et al., 2021; 2023; Zhang et al., 2022). Xu et al. (2020) uses a separate image encoder before adding the output as feature to the token encodings, while Huang et al. (2022) jointly models the page image patches alongside the tokens, using a word-patch alignment self-supervised pretraining task to learn an alignment between the modalities.

Other approaches treat extraction as a sequence generation problem. Powalski et al. (2021) adds an auto-regressive decoder on top of a text-layout-image encoder, all initialized from T5 (Raffel et al., 2020). Kim et al. (2022) foregoes the text recognition step completely, using a Vision Transformer encoder with an auto-regressive decoder pretrained on a pseudo-OCR task on a large document image corpora, and finetuned on the final extraction parse tree with XML tags for the target extraction schema. While this approach allows to predict hierarchical entities, it does not allow localizing entities in the document.

None of the previously discussed approaches attempt to understand the semantics of the schema and its entity types, and instead opt to encode the schema in the model weights through training, hence requiring training data for unseen schemas and document types. QueryForm (Wang et al., 2023b) utilizes a prompt encoding both the schema and entity types, allowing the model to do zero-shot extraction. Likewise, PPN (Wei et al., 2023) inputs the raw entity types in the encoder itself, and uses a scoring matrix to predict the link classes between document tokens and types, with great few-shot performance. However, both approaches are not able to predict hierarchical entities.

In parallel, Large Language Models (OpenAI, 2023; Google et al., 2023; Hoffmann et al., 2022) have revolutionized Natural Language Processing, showing the capabilities to solve a task with simply an instruction (Wei et al., 2022) or a few examples added to the prompt (Brown et al., 2020). This paradigm opens the possibility of extracting entities with very little to no training data. Wang et al. (2023a) transforms the NER task to a sequence generation task suitable for LLMs by incorporating special tokens in the sequence, marking the entity boundaries, and proposes a self-verification strategy limiting the LLM hallucinations. However, this is applicable to text-only scenarios, with hallucinations still a possibility.

This motivates us to introduce *Language Model-based Document Information EXtraction and Localization* (LMDX), a methodology for leveraging off-the-shelf LLMs for information extraction on semi-structured documents. Our contributions can be summarized as follows:

- We propose a **prompt** that enables LLMs to perform the document IE task on leaf and hierarchical entities with precise localization, including without any training data.
- We also propose a **layout encoding scheme** that communicate spatial information to the LLM without any change to its architecture.
- We introduce a **decoding algorithm transforming the LLM responses into extracted entities** and their bounding boxes on the document, while discarding all hallucination.
- We systematically evaluate the data efficiency of LMDX on multiple public benchmarks and establish a new state-of-the-art on those by a wide margin, especially at low-data regimes.

A comparison of LMDX characteristics and other popular document information extraction systems can be found at Table 1.

Table 1: Comparison of document information extraction systems.

Document Information Extraction Systems	Hierarchical entity	Entity localization	Zero-shot support
FormNet(v2), LayoutLM(v2), Docformer, Glean, ...	✗	✓	✗
QueryForm, PPN	✗	✓	✓
Donut	✓	✗	✗
<b>LMDX (Ours)</b>	✓	✓	✓

## 2 METHODOLOGY

### 2.1 OVERVIEW

Overall, our pipeline is divided into five stages: OCR, chunking, prompt generation, LLM inference and decoding, detailed in the following sections. An overview with a simple example can be found in Figure 1, with the input and output of each stage showcased.

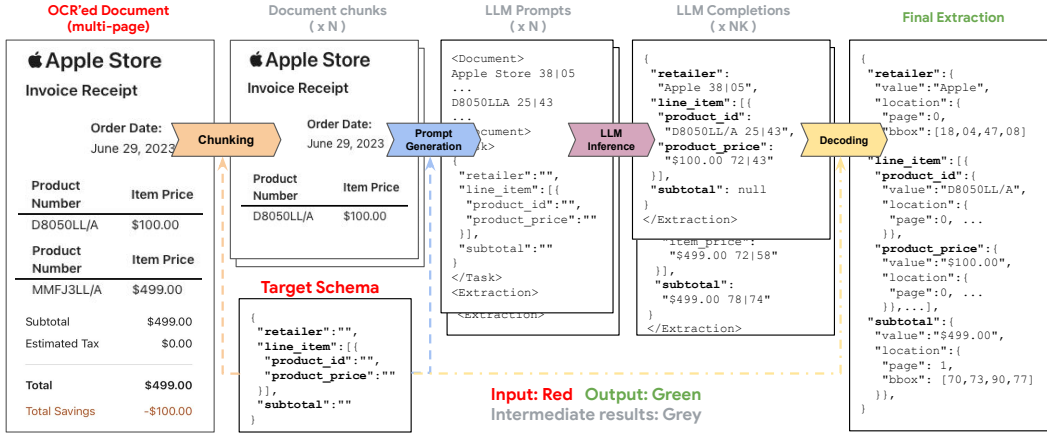


Figure 1: Overview of the LMDX methodology.

### 2.2 OPTICAL CHARACTER RECOGNITION

We first use an off-the-shelf OCR service on the document image to obtain words and lines segments, along with their corresponding spatial position (bounding box) on the document. An example of output from that stage on a sample document is given in Appendix A.6.

### 2.3 CHUNKING

Since a document can be arbitrarily long and that LLMs have limited input token length, the document is divided into document chunks so that each is small enough to be processed by the LLM. To achieve this, we first divide the document into individual pages, then we iteratively remove the last line segments (coming from OCR) until the prompt containing this chunk is below the maximum input token length of the LLM. Lastly, we group those removed lines as a new document page, and repeat the same logic until all chunks are below the input token limit of the LLM. At the end of this stage, we have  $N$  chunks. The decision to first divide the document by page stems from the observation that entities rarely cross page boundaries, and as such this chunking scheme will have minimal impact on the final extraction quality. The algorithm is described in pseudo-code in Appendix A.1.

### 2.4 PROMPT GENERATION

The prompt generation stage takes in the  $N$  document chunks and creates a LLM prompt for each of them. As seen in Figure 2, our prompt design contains the document representation, a description

of the task, and the target schema representation containing the entities to extract. XML-like tags are used to define the start and end of each component.

```
<Document>
{DOCUMENT_REPRESENTATION}
</Document>
<Task>
{TASK_DESCRIPTION}
{SCHEMA_REPRESENTATION}
</Task>
<Extraction>
```

Figure 2: Structure of the LLM prompts.

**Document Representation.** The chunk content is represented in the prompt as the concatenation of all its segment texts, suffixed with the coordinates of those segments in the following format:  $\langle \text{segment text} \rangle XX|YY_{\text{segment}}$ . Coordinate tokens are built by normalizing the segment’s X and Y coordinate, and quantizing them in  $B$  buckets, assigning the index of that bucket as the token for a coordinate.

This coordinate-as-tokens scheme allows us to communicate the layout modality to the LLM, without any change to its architecture. There are many variation to that scheme: using OCR line versus OCR words as segment, the granularity of the quantization, and the number of coordinates to use per segment (e.g.  $[x_{\text{center}}, y_{\text{center}}]$  versus  $[x_{\text{min}}, y_{\text{min}}, x_{\text{max}}, y_{\text{max}}]$ ). Appendix A.4 shows how those variations affect the prompt token length. In practice, since LLM context length is still limited, we use line-level segments with 2 coordinates and  $B = 100$  quantization buckets in all our experiments.

**Task Description.** The task description is simply a short explanation of the task to accomplish. In our experiments, we hard code it to the following: *From the document, extract the text values and tags of the following entities:.*

**Schema Representation.** The schema is represented as a structured JSON object, where the keys are the entity types to be extracted, and the values correspond to their occurrence (single or multiple), and sub-entities (for hierarchical entities). For instance,  $\{ "foo": "", "bar": [{"baz": []}] \}$  means that the LLM should extract only a single entity of type *foo* and multiple hierarchical entities of type *bar*, that could each hold multiple entities of type *baz*.

After this step, we have  $N$  prompts, one for each document chunk. A full example of a prompt on a document can be found in Appendix A.6.

## 2.5 COMPLETION TARGETS

In this section, we describe the expected LLM completion format, which can be observed in Figure 1. Like the schema, the completion is a JSON structured object with the keys being the entity types, and values being the extracted information from the document chunk. JSON was chosen as a format for the completion and schema since it supports hierarchical objects, is very token-efficient, and usually present in LLMs training data mixtures. Note that the keys in the completion have the same ordering, occurrence and class (hierarchical or leaf) as the entity types in the schema. The values of leaf entities must follow a specific format:

$$\langle \text{text on segment}_1 \rangle XX|YY_{\text{segment}_1} \backslash n \langle \text{text on segment}_2 \rangle XX|YY_{\text{segment}_2} \backslash n \dots$$

An entity can span multiple (potentially disjoint) OCR segments (lines or words). For each segment of the entity, the value contains the entity text on that segment, along with the coordinate tokens of that segment, which act as a *segment identifier*, allowing us to localize the entities and ground the model prediction (e.g. making sure the extracted value is not a hallucination), as will be detailed in Section 2.7.

---

Missing entity types are completed by the model with *null* for singular types, and  $\square$  for repeated types. Samples of completions can be found in Appendix A.6.

## 2.6 LLM INFERENCE

In this stage of the pipeline, we run inference on the LLM with the  $N$  prompts. For each prompt, we sample  $K$  completions from the LLM (for a total of  $NK$  completions for the entire document) using  $\text{Top}_K$  sampling. This randomness in the sampling allows to do error correction (e.g. if a response is not valid JSON, have hallucinated segment coordinate identifier, etc), and increase the extraction quality as will be shown in further sections. Note that we still want the inference to be fully deterministic so that LMDX’s extractions are the same across two identical documents. To do so, we rely on pseudo-random sampling using a fixed seed.

## 2.7 DECODING

In this stage, we parse the raw LLM completions into structured entities and their locations.

**Conversion to structured entities.** We begin by parsing each model completion as a JSON object. Completions that fail to parse are discarded. For each key-value pair in the JSON object, we interpret the key as the entity type and parse the value to get the entity text and bounding box (as detailed in the next paragraph). Predicted entity types that are not in the schema are discarded. If the model unexpectedly predicts multiple values for single-occurrence entity types, we use the most frequent value as the final predicted value. Hierarchical JSON object are recursively parsed as hierarchical entities in a similar manner. This algorithm is described in pseudo-code in Appendix A.3.

**Entity Value Parsing.** We expect the JSON value to include both text extractions and segment identifiers for each predicted entity, as described in Section 2.5. We first parse the value into its (*segment text*, *segment identifier*) pairs. For each pair, we look up the corresponding segment in the original document using the segment identifier and verify that the extracted text is exactly included on that segment. Finally, once we have the entity location on all its segments, we get the entity bounding box by computing the smallest bounding box encompassing all the OCR words included in the entity. Entity values with any segments that fail to ground (invalid entity value format, non-existent segment identifier, or non-matching segment text) in the original document are discarded. The entity value parsing algorithm is described in pseudo-code in Appendix A.2.

**Prediction Merging.** We first merge the predicted entities for the same document chunk from the  $K$  LLM completions through majority voting (Wang et al., 2022). For each entity type, we gather the predicted entities, including empty predictions, across the  $K$  completions. The most common prediction(s) are selected as the predicted value for that entity type. We then merge the predictions among the  $N$  document chunks by concatenating them to obtain the document level predictions.

**Prediction Merging for hierarchical entities.** For hierarchical entities, we use the entire predicted tree value from a single LLM completion, as this method best preserves the parent-child relationship predicted by the model. For each top-level hierarchical entity type, we perform majority voting on all affiliated leaf, intermediate and top-level entity types among  $K$  completions as if they are flattened. We then tally the votes with equal weight to determine which completion to use for the prediction, and select the most common one for that hierarchical entity.

## 3 EVALUATION

We evaluate the methodology explained in section 2 on public benchmarks using the PaLM 2-S LLM, which we call  $\text{LMDX}_{\text{PaLM 2-S}}$ . Note that we use the small version of this LLM due to limited accelerator resources, but larger versions could be used, likely leading to higher extraction quality.

Our training process is composed of two phases as shown in Figure 3. In the first phase we finetune PaLM 2-S on a data mixture containing a variety of (*document*, *schema*, *extraction*) tuples. In particular, this data mixture contains the *Payment* dataset (Majumder et al., 2020), along with a diverse set of publicly available PDF form templates obtained from government websites that we filled with

synthetic data using an internal tool, and annotated for schema and entities to extract. The goal of this phase is to train the model to interpret the semantics of the entity types and extraction hierarchy specified in the schema, and find them within the document, along with learning the extraction syntax. Hence, the variety of schemas and documents in this phase is of utmost importance.

During the second phase, starting from the base entity extractor checkpoint from the previous phase, we finetune the LLM on the target to specialize it to do high quality extraction on the target benchmark. At this stage, only the target benchmark data is included in the training mixture. Note that, for zero-shot experiments, this second phase is skipped. Furthermore, no document or schema contained in the base extraction training phase overlap with the documents and schemas used in the specialization training phase. For all training phases, we follow the input and target syntax described in section 2.4 and 2.5.

### 3.1 PARAMETERS

For training, we finetune PaLM 2-S using a batch size of 8, a dropout probability of 0.1 and a learning rate of  $10^{-6}$  with a standard cross-entropy loss. Once training is done, we select the checkpoint with the lowest loss on the dev set, and report performance on the test set. For LLM inference, we use a temperature of 0.5 and a  $\text{Top}_K$  of 40, sampling 16 responses for each chunk processed by the LLM, as described in section 2.6. Finally, for both training and inference, we use an input token length of 6144 and output token length of 2048. We use line-level segments and only two coordinates  $[x_{\text{center}}, y_{\text{center}}]$  with 100 quantization buckets to save on the number of input and output tokens consumed by the coordinate-as-tokens scheme.

### 3.2 DATASETS

**Visually Rich Document Understanding (VRDU).** Wang et al. (2023c) introduces a public benchmark for entity extraction from visually-rich documents that includes two datasets: Registration Form, containing 6 semantically rich entity types, and Ad-buy Form, containing 14 entity types with one hierarchical *line\_item* entity. For each dataset, VRDU proposes samples of 10, 50, 100 and 200 train documents to evaluate the data efficiency of models. It also offers different tasks to evaluate the generalization powers of extraction systems: Single Template Learning (STL) where train/test share the same single template, Mixed Template Learning (MTL) where train/test contains overlapping sets of templates, and Unseen Template Learning (UTL) where train/test contains disjoint sets of templates. For our experiments, we finetune  $\text{LMDX}_{\text{PaLM 2-S}}$  for 4000 steps on each dataset, training data size, and task setup independently and report Micro-F1 through the provided evaluation tool. We then compare  $\text{LMDX}_{\text{PaLM 2-S}}$  to the published state-of-the-art baselines.

**Consolidated Receipt Dataset (CORD).** Park et al. (2019) introduces a benchmark of Indonesian receipts from shops and restaurants, with a target schema of 30 fine-grained entities, grouped into *menu*, *total* and *subtotal* hierarchical entities. CORD<sup>1</sup> does not provide a standard evaluation toolkit, so we adopt the normalized Tree Edit Distance accuracy metric (Zhang & Shasha, 1989), previously introduced by Kim et al. (2022) on that benchmark, since it is agnostic to the output scheme used and considers the hierarchical entities as part of the metric. For our experiments, we use the official 800*train*/100*dev*/100*test* split, but also sample the first  $D = 10/50/100/200$  documents from the train split to assess the data efficiency of LMDX on this benchmark. For each data setup, we fine-

<sup>1</sup><https://huggingface.co/datasets/naver-clova-ix/cord-v1>

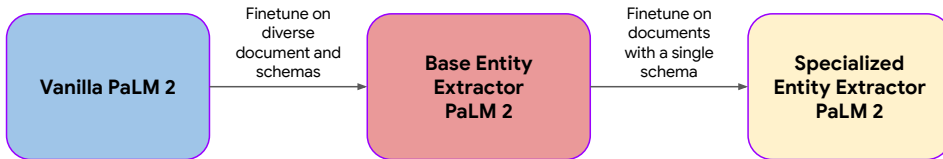


Figure 3: LMDX training phases.

tune LMDX for 12000 steps. For comparison, we also train and evaluate state-of-the-art baselines LayoutLMv3<sub>LARGE</sub> and *Donut*. Those baselines are detailed in Appendix A.7.

### 3.3 RESULTS

Table 2: Results of LMDX<sub>PaLM 2-S</sub> on the different tasks and train data size setups  $|\mathcal{D}|$  of VRDU, with best performing model results in bold. Unlike the baselines, LMDX can do zero-shot extraction.

$ \mathcal{D} $	Model	Registration Form			Ad-buy Form		
		Single	Mixed	Unseen	Mixed Template		Unseen
		Micro-F1	Micro-F1	Micro-F1	Micro-F1	Line Item F1	Micro-F1
<b>0</b>	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>73.81</b>	<b>71.65</b>	<b>74.94</b>	<b>39.74</b>	<b>21.21</b>	<b>39.33</b>
<b>10</b>	FormNet	74.22	63.61	50.53	20.47	5.72	20.28
	LayoutLM	65.91	36.41	25.54	20.20	6.95	19.92
	LayoutLMv2	80.05	69.44	54.21	25.36	9.96	25.17
	LayoutLMv3	72.51	60.72	21.17	10.16	5.92	10.01
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>90.88</b>	<b>87.72</b>	<b>86.87</b>	<b>54.35</b>	<b>39.35</b>	<b>54.82</b>
<b>50</b>	FormNet	89.38	85.38	68.29	40.68	19.06	39.52
	LayoutLM	86.21	80.15	55.86	39.76	19.50	38.42
	LayoutLMv2	88.68	84.13	61.36	42.23	20.98	41.59
	LayoutLMv3	87.24	81.36	47.85	39.49	19.53	38.43
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>93.06</b>	<b>91.42</b>	<b>88.43</b>	<b>75.08</b>	<b>65.42</b>	<b>75.70</b>
<b>100</b>	FormNet	90.91	88.13	72.58	40.38	18.80	39.88
	LayoutLM	88.70	86.02	63.68	42.38	21.26	41.46
	LayoutLMv2	90.45	88.36	65.96	44.97	23.52	44.35
	LayoutLMv3	89.23	87.32	57.69	42.63	22.08	41.54
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>93.97</b>	<b>92.41</b>	<b>89.70</b>	<b>78.05</b>	<b>69.77</b>	<b>75.99</b>
<b>200</b>	FormNet	92.12	90.51	77.29	43.23	21.86	42.87
	LayoutLM	90.47	87.94	70.47	44.66	23.90	44.18
	LayoutLMv2	91.41	89.19	72.03	46.54	25.46	46.31
	LayoutLMv3	90.89	89.77	62.58	45.16	24.51	44.43
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>93.97</b>	<b>92.78</b>	<b>90.22</b>	<b>79.82</b>	<b>72.09</b>	<b>78.42</b>

Results for VRDU are presented in Table 2. For all data regimes and tasks, LMDX<sub>PaLM 2-S</sub> sets a new state-of-the-art by a wide margin. In particular, we find that LMDX<sub>PaLM 2-S</sub> can extract decently with no training data: it exhibits similar extraction quality at zero-shot than baselines at 10-100 train dataset size (for instance 39.74% Micro-F1 on Ad-Buy Form Mixed Template vs 40.68% for FormNet at 50 train documents, or 73.81% Micro-F1 on Registration Single Template vs 74.22% for FormNet at 10 train documents). LMDX<sub>PaLM 2-S</sub> is also much more data efficient than the baselines: it is at 5.06% Micro-F1 of its peak performance at 10 training documents for Registration Form Mixed Template (87.72% vs 92.78% Micro-F1) while LayoutLMv2, the strongest baseline, is within 19.75% of its peak performance (69.44% vs 89.19% Micro-F1). Finally, we notice that LMDX<sub>PaLM 2-S</sub> generalizes much better to unseen templates than baselines: on Registration Form, LMDX<sub>PaLM 2-S</sub> has a drop lesser than 5% Micro-F1 on Unseen Template compared to Single Template across all data regimes, while baselines like LayoutLMv2 sees a drop between 19.38% and 27.32%.

On CORD, with results in Table 3, we observe similar trends, highlighting the generalization of the results. At 10 documents, LMDX<sub>PaLM 2-S</sub> is 4.03% from its peak performance attained at 800 documents, versus 22.34% for the strongest baseline LayoutLMv3<sub>LARGE</sub>, showcasing the data efficiency of the LMDX methodology.

**Performance on Hierarchical Entities.** As seen on Ad-Buy Form tasks, LMDX<sub>PaLM 2-S</sub> is capable of grouping line items much better than the baselines (which are using heuristics) for all data regimes. In particular, LMDX<sub>PaLM 2-S</sub> has similar line\_item grouping performance at zero-shot than the best baseline at 200 train documents (21.21% versus 25.46% F1 respectively). With all the training data, LMDX<sub>PaLM 2-S</sub> scores a 72.09% F1 on line\_item, an absolute improvement of 46.63% over the best baseline LayoutLMv2.

Table 3: LMDX<sub>PaLM 2-S</sub> results on CORD. Normalized Tree Edit Distance Accuracy is reported.

$ \mathcal{D} $	Model	n-TED Accuracy
<b>0</b>	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>67.47</b>
<b>10</b>	Donut	33.01
	LayoutLMv3 <sub>LARGE</sub>	73.87
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>92.27</b>
<b>50</b>	Donut	75.44
	LayoutLMv3 <sub>LARGE</sub>	87.29
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>93.80</b>
<b>100</b>	Donut	82.17
	LayoutLMv3 <sub>LARGE</sub>	91.83
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>93.64</b>
<b>200</b>	Donut	84.49
	LayoutLMv3 <sub>LARGE</sub>	94.44
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>94.73</b>
<b>800</b>	Donut	90.23
	LayoutLMv3 <sub>LARGE</sub>	96.21
	<b>LMDX<sub>PaLM 2-S</sub></b>	<b>96.30</b>

### 3.4 ABLATIONS

In this section, we ablate different facets of the LMDX methodology to highlight their relative importance. The results can be found in Table 4 and are discussed below. For all ablations, we evaluate on the VRDU Ad-Buy Form Mixed Template task, only changing the ablated facet.

Table 4: Ablations of Base Entity Extraction Training, Coordinate Tokens, and Sampling and their relative effects on extraction quality. All ablations are done on VRDU Ad-Buy Mixed Template.

$ \mathcal{D} $	LMDX <sub>PaLM 2-S</sub>	Without Base EE Training		Without Coordinate Tokens		Without Sampling Strategy	
	Micro-F1	Micro-F1	$\Delta$ (%)	Micro-F1	$\Delta$ (%)	Micro-F1	$\Delta$ (%)
<b>0</b>	39.74	0.00	-39.74	27.59	-12.15	39.53	-0.21
<b>10</b>	54.35	42.91	-11.44	39.37	-14.98	52.85	-1.50
<b>50</b>	75.08	66.51	-8.57	62.35	-12.73	73.88	-1.20
<b>100</b>	78.05	68.87	-9.18	65.14	-12.91	77.30	-0.75
<b>200</b>	79.82	72.25	-7.57	65.70	-14.12	78.43	-1.39

**Effects of Base Entity Extraction Training.** In this ablation, we remove the first stage training on the varied data mixture from Figure 3 and directly finetune on the VRDU target task. As seen on columns 3-4 of Table 4, ablating that training stage leads to significant drop in extraction quality in few-shot scenarios and the complete loss of zero-shot extraction ability due to the model not respecting the extraction format, hence failing decoding. As the train set size increases, the degraded performance lessens, from -11.44% to -7.57%, as the model learns the extraction task and the desired completion format.

**Effects of Coordinate Tokens.** In this ablation, we replace the coordinate tokens, which communicate the position of each line within the document, by the index of that line. This index still acts as a unique identifier for the line segment (required for entity localization and grounding) but does not communicate any position information. An example of a prompt with line index can be found in Appendix A.6. As can be seen on columns 5-6 of Table 4, the coordinate tokens are substantially important to the extraction quality, ranging from 12.15% to 14.98% absolute micro-F1 improvement across the data regimes.

**Effects of Sampling Strategy.** In this ablation, we discard our sampling strategy, and instead sample a single response from the model. As seen in columns 7-8 of Table 4, this leads to a 0.21% to 1.5% drop in micro-F1. While overall minor for quality, the sampling strategy also allows to correct extraction format mistakes.



### 3.5 ERROR ANALYSIS AND LIMITATIONS

In this section, we perform an error analysis on the test set to identify common error patterns of LMDX. A very common error type we observe is caused by OCR lines grouping multiple semantically different segments. We show two instance of those cases observed in LMDX<sub>PaLM 2-S</sub> on the VRDU Ad-Buy Form in Figure 4. In the first example, prediction for the entity *line\_item/program\_desc* includes text from the previous column "Channel" along with the value in the column "Description". From the OCR line bounding boxes, we can see that these two columns are grouped as the same OCR line. In the second example, the model confuses between the adjacent keys "Invoice Period" and "Flight Dates" and extracts invoice dates as flight dates. Similar to the first example, OCR line bounding boxes show that the invoice dates and the key "Flight Dates" are grouped together in the same line although they are semantically different. As LMDX<sub>PaLM 2-S</sub> uses only coarse line layout information ( $[x_{center}, y_{center}]$  with 100 quantization buckets), the model fails in these cases. This is a current limitation of the LMDX system. We believe that incorporating the image modality will help make LMDX more performant and robust to those OCR errors, addressing the current limitations of the system.

Line	Channel	Description
1	WJZ	Local News 6a-630a
	Level 3 - Immediately Preemptable	
	WJZ	

**Example 1:**  
*line\_item/program\_desc*  
**Groundtruth:** Local News 6a-630a  
**Prediction:** WJZ Local News 6a-630a

Property	KXTV
Invoice #	1903525-1
Order #	1903525
Invoice Date	12/29/19
Alt Order #	WOC12331150
Invoice Month	December 2019
Deal #	
Invoice Period	11/25/19 - 12/29/19
Flight Dates	12/24/19 - 12/30/19
Advertiser	POL/ Tom Steyer / D / PRES / US

**Example 2**  
**Groundtruth:** flight\_from: 12/24/19  
flight\_to: 12/30/19  
**Prediction:** flight\_from: 11/25/19  
flight\_to: 12/29/19

Figure 4: Typical error pattern of LMDX<sub>PaLM 2-S</sub>. In both examples, the detected OCR lines are shown in red, the model predicted entities are shown in blue, and the groundtruth entities are shown in green. In both cases, the detected OCR lines merge two semantically distinct segments, causing the model to wrongly associate them in its predictions.

## 4 CONCLUSION

In this paper, we have introduced LMDX, a methodology that enables using LLMs for information extraction on visually rich documents, setting a new state-of-the-art on public benchmarks VRDU and CORD. LMDX is the first methodology to allow the extraction of singular, repeated and hierarchical entities, while grounding its predictions and localizing the entities in the document. LMDX is extremely data efficient, and even allows high quality extraction at zero-shot on entirely new document types and schemas. Nonetheless, since it relies on a LLM, LMDX is more resource-intensive than previous approaches, and its coordinate-as-tokens scheme requires long inputs and outputs. As future work, we will explore applying the methodology to open-source LLMs and adding the image modality to the system using Large Vision-Language Models.

## 5 REPRODUCIBILITY STATEMENT

In order to increase reproducibility, we’ve provided all details of the LMDX methodology. We’ve included our LLM prompts and completions in Appendix A.6, along with all our algorithms for chunking and decoding in Appendix A.1, A.2 and A.3. Furthermore, we’ve provided the exact target schemas used in our experiments in Appendix A.5. For CORD specifically, we’ve used a metric with a public implementation (<https://github.com/clovaai/donut/blob/master/donut/util.py>) and an easy to reproduce sampling strategy for the data-efficiency splits (first  $D$  train documents). Finally, our baselines are publicly available (<https://github.com/microsoft/unilm/tree/master/layoutlmv3>, <https://github.com/clovaai/donut>) and thoroughly detailed in Appendix A.7.

---

## REFERENCES

- Srikar Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R. Manmatha. Doc-former: End-to-end transformer for document understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 993–1003, October 2021.
- Srikar Appalaraju, Peng Tang, Qi Dong, Nishant Sankaran, Yichu Zhou, and R. Manmatha. Doc-formerv2: Local features for document understanding, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- Timo I. Denk and Christian Reisswig. Bertgrid: Contextualized embedding for 2d document representation and understanding, 2019.
- Łukasz Garncarek, Rafał Powalski, Tomasz Stanisławek, Bartosz Topolski, Piotr Halama, Michał Turski, and Filip Graliński. LAMBERT: Layout-aware language modeling for information extraction. In *Document Analysis and Recognition – ICDAR 2021*, pp. 532–547. Springer International Publishing, 2021. doi: 10.1007/978-3-030-86549-8\_34. URL [https://doi.org/10.1007/978-3-030-86549-8\\_34](https://doi.org/10.1007/978-3-030-86549-8_34).
- Rohan Anil Google, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Keane, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Mousaleem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. Palm 2 technical report, 2023.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Henighan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022.
- Teakgyu Hong, Donghyun Kim, Mingi Ji, Wonseok Hwang, Daehyun Nam, and Sungrae Park. BROS: A layout-aware pre-trained language model for understanding documents. *CoRR*, abs/2108.04539, 2021. URL <https://arxiv.org/abs/2108.04539>.
- 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*, 2022.

- 
- Wonseok Hwang, Jinyeong Yim, Seunghyun Park, Sohee Yang, and Minjoon Seo. Spatial dependency parsing for semi-structured document information extraction. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 330–343, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.28. URL <https://aclanthology.org/2021.findings-acl.28>.
- Anoop Raveendra Katti, Christian Reisswig, Cordula Guder, Sebastian Brarda, Steffen Bickel, Johannes Höhne, and Jean Baptiste Faddoul. Chargrid: Towards understanding 2d documents, 2018.
- Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon 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 (ECCV)*, 2022.
- Chen-Yu Lee, Chun-Liang Li, Timothy Dozat, Vincent Perot, Guolong Su, Nan Hua, Joshua Ainslie, Renshen Wang, Yasuhisa Fujii, and Tomas Pfister. FormNet: Structural encoding beyond sequential modeling in form document information extraction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3735–3754, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.260. URL <https://aclanthology.org/2022.acl-long.260>.
- Chen-Yu Lee, Chun-Liang Li, Hao Zhang, Timothy Dozat, Vincent Perot, Guolong Su, Xiang Zhang, Kihyuk Sohn, Nikolay Glushnev, Renshen Wang, Joshua Ainslie, Shangbang Long, Siyang Qin, Yasuhisa Fujii, Nan Hua, and Tomas Pfister. FormNetV2: Multimodal graph contrastive learning for form document information extraction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9011–9026, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.501. URL <https://aclanthology.org/2023.acl-long.501>.
- Bodhisattwa Prasad Majumder, Navneet Potti, Sandeep Tata, James Bradley Wendt, Qi Zhao, and Marc Najork. Representation learning for information extraction from form-like documents. In *ACL*, 2020.
- OpenAI. Gpt-4 technical report, 2023.
- Rasmus Berg Palm, Ole Winther, and Florian Laws. Cloudscan - a configuration-free invoice analysis system using recurrent neural networks. In *Proceedings of 2017 14th IAPR International Conference on Document Analysis and Recognition*, pp. 406–413, United States, 2017. IEEE. ISBN 9781538635858. doi: 10.1109/ICDAR.2017.74.
- Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwal-suk Lee. Cord: A consolidated receipt dataset for post-ocr parsing. In *Workshop on Document Intelligence at NeurIPS 2019*, 2019.
- Rafał Powalski, Łukasz Borchmann, Dawid Jurkiewicz, Tomasz Dwojak, Michał Pietruszka, and Gabriela Pałka. Going full-tilt boogie on document understanding with text-image-layout transformer. In Josep Lladós, Daniel Lopresti, and Seiichi Uchida (eds.), *Document Analysis and Recognition – ICDAR 2021*, pp. 732–747, Cham, 2021. Springer International Publishing. ISBN 978-3-030-86331-9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Lance Ramshaw and Mitch Marcus. Text chunking using transformation-based learning. In *Third Workshop on Very Large Corpora*, 1995. URL <https://aclanthology.org/W95-0107>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).

- 
- Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. Gpt-ner: Named entity recognition via large language models, 2023a.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Zifeng Wang, Zizhao Zhang, Jacob Devlin, Chen-Yu Lee, Guolong Su, Hao Zhang, Jennifer Dy, Vincent Perot, and Tomas Pfister. QueryForm: A simple zero-shot form entity query framework. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4146–4159, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.255. URL <https://aclanthology.org/2023.findings-acl.255>.
- Zilong Wang, Yichao Zhou, Wei Wei, Chen-Yu Lee, and Sandeep Tata. Vrdu: A benchmark for visually-rich document understanding. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23*, pp. 5184–5193, New York, NY, USA, 2023c. Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599929. URL <https://doi.org/10.1145/3580305.3599929>.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022.
- Kaiwen Wei, Jie Yao, Jingyuan Zhang, Yangyang Kang, Fubang Zhao, Yating Zhang, Changlong Sun, Xin Jin, and Xin Zhang. Ppn: Parallel pointer-based network for key information extraction with complex layouts, 2023.
- Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, Min Zhang, and Lidong Zhou. Layoutlmv2: Multi-modal pre-training for visually-rich document understanding. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL) 2021*, 2021.
- 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*, pp. 1192–1200, 2020.
- Kaizhong Zhang and Dennis Shasha. Simple fast algorithms for the editing distance between trees and related problems. *SIAM Journal on Computing*, 18(6):1245–1262, 1989. doi: 10.1137/0218082. URL <https://doi.org/10.1137/0218082>.
- Zhenrong Zhang, Jiefeng Ma, Jun Du, Licheng Wang, and Jianshu Zhang. Multimodal pre-training based on graph attention network for document understanding, 2022.

## A APPENDIX

### A.1 CHUNKING ALGORITHM

---

**Algorithm 1** Document Chunking

---

```

1: function CHUNK( $D, L, F$ ) ▷  $D$  is a document containing multiple pages.  $L$  is token limit.
2: ▷  $F$  is a function that outputs prompt token length given some segments (e.g. lines).
3:    $C = \phi$  ▷  $C$  is to record all produced chunks.
4:   for  $i = 1$  to  $|D.pages|$  do
5:      $S = D.pages[i].segments$ 
6:     while  $S \neq \phi$  do
7:       for  $j = |S|$  to 1 do ▷ Start pruning from the end of the page.
8:         if  $F(S[1:j]) \leq L$  then
9:            $C = C \cup \{S[1:j]\}$ 
10:           $S = S[j+1 : |S|]$  ▷ Continue to work on the rest of the segments.
11:          Exit for loop
12:        end if
13:      end for
14:    end while
15:  end for
16:  return  $C$ 
17: end function

```

---

### A.2 ENTITY VALUE PARSING ALGORITHM

---

**Algorithm 2** Entity Value Parsing

---

```

1: function PARSEENTITYVALUE( $D, E$ ) ▷  $D$  is a document chunk.
2: ▷  $E$  is raw extraction results for one entity type parsed from one LLM sample.
3:    $G = \phi$  ▷  $G$  is to record all parsed entity values.
4:    $R = \text{Regex}("\\d\\d\\d\\d\\d")$  ▷  $R$  is a regex that captures the segment identifiers.
5:    $M = \{ "s.x|s.y" \mapsto s \mid s \in D.segments \}$  ▷  $M$  holds a mapping between segment id and segment.
6:   for  $i = 1$  to  $|E|$  do
7:      $W = \phi$  ▷  $W$  is to hold all words for this entity.
8:      $P = R.split(E[i])$  ▷  $P$  is expected to be interleaved text values and segment ids.
9:     for  $j = 1$  to  $|P|/2$  do
10:      if  $P[j*2] \notin M$  then
11:        Go to next  $i$  ▷ Segment ID is hallucinated. Grounding failure.
12:      end if
13:       $S = M[P[j*2]]$  ▷ Retrieve the stored segment from  $M$  with parsed segment ID.
14:       $T = P[j*2-1]$  ▷  $T$  is to hold the predicted text.
15:      if  $T$  not substring of  $S$  then
16:        Go to next  $i$  ▷ Grounding failure, skip the current entity.
17:      end if
18:       $W = W \cup (S \cap T)$ 
19:    end for
20:     $G'.value = \bigcup_{w \in W} w.text\_value$  ▷  $G'$  is to hold the entity to return.
21:     $G'.bounding\_box = \{ \min(b.x), \min(b.y), \max(b.x), \max(b.y) \}_{w \in W, b=w.bounding\_box}$ 
22:     $G = G \cup \{G'\}$ 
23:  end for
24:  return  $G$ 
25: end function

```

---

### A.3 DECODING ALGORITHM

---

**Algorithm 3** Responses Decoding
 

---

```

1: function DECODEFORTYPE( $J, T, D$ )                                 $\triangleright J$  is one or more JSON objects.
2:                                                                 $\triangleright T$  is an entity type.
3:                                                                 $\triangleright D$  is a document chunk.
4:    $E = \phi$                                                         $\triangleright E$  is to record all parsed and grounded entities.
5:   for  $j = 1$  to  $|J|$  do
6:      $J' = J[j][T.type]$                                            $\triangleright J'$  is to hold entities for  $T$ 's type before grounding.
7:     if  $T.subtypes = \phi$  then                                     $\triangleright T$  is leaf entity type.
8:        $E = E \cup ParseEntityValue(D, J')$ 
9:     else                                                          $\triangleright T$  is hierarchical entity type.
10:       $E'.subtypes = \bigcup_{T' \in T.subtypes} DecodeForType(J', T', D)$   $\triangleright E'$  is hierarchical entity.
11:       $E = E \cup \{E'\}$ 
12:    end if
13:  end for
14:  return  $E$ 
15: end function
16:
17: function MAJORITYVOTING( $T, E$ )                                 $\triangleright T$  is an entity type.
18:                                                                 $\triangleright E$  is a 2D vector of entities of type  $T$  from all LLM responses.
19:    $V = [0, 0, \dots, 0] \in \mathbb{R}^{|E|}$                              $\triangleright V$  is to record all votes.
20:    $L = \{T\}$ 
21:   while  $L \neq \phi$  do
22:      $T' = L[0]$ 
23:      $E' = \phi$ 
24:     for  $j = 1$  to  $|E|$  do
25:        $E' = E' \cup \{e | e \in E[j], e.type = T'\}$                $\triangleright E'[j]$  holds entities with type  $T'$  from  $E[j]$ .
26:     end for
27:     for  $i = 1$  to  $|E'| - 1$  do
28:       for  $j = i + 1$  to  $|E'|$  do
29:         if  $E'[i] = E'[j]$  then
30:            $V[i] = V[i] + 1$ 
31:            $V[j] = V[j] + 1$ 
32:         end if
33:       end for
34:     end for
35:      $L = L[1 : |L|]$                                                $\triangleright$  Remove  $T'$  and inject its sub-types for recursion.
36:      $L = L \cup T'.subtypes$ 
37:   end while
38:   return  $E[argmax(V)]$                                            $\triangleright$  Return the entity values with the highest votes.
39: end function
40:
41: function DECODEALLSAMPLES( $S, T, D$ )                             $\triangleright S$  is all LLM response samples on  $D$ .
42:                                                                 $\triangleright T$  is a list of entity types.
43:                                                                 $\triangleright D$  is a document chunk.
44:   return  $\bigcup_{T' \in T} MajorityVoting(\bigcup_{S' \in S} DecodeForType(ParseJson(S'), T', D))$ 
45: end function

```

---

#### A.4 TOKEN LENGTH STATISTICS

Table 5 details the token length (50<sup>th</sup> and 99<sup>th</sup> percentiles) of the prompt and completion targets for the train split of datasets used in our experiments. We select the line level segment, 2 coordinate scheme, no JSON indentation so that all datasets fit within our 6144 prompt token length and 2048 output token length.

Table 5: Prompt and target token length of different coordinate-as-tokens schemes on VRDU and CORD benchmarks, using the vocabulary of PaLM 2-S. We vary the number of coordinates and their quantization buckets in the localization tags, the segment level (e.g. line versus word), chunking style (e.g. page versus max input tokens) and JSON indentation in the schema and completion targets.

VRDU Ad-Buy Form								
# Coord.	# Quant.	Segment	Chunking	JSON Indent	Input		Target	
					50 <sup>th</sup>	99 <sup>th</sup>	50 <sup>th</sup>	99 <sup>th</sup>
2	100	Line	Page	None	2377	3920	602	1916
2	100	Word	Page	None	3865	13978	718	2328
4	100	Line	Page	None	3329	5284	777	2473
2	1000	Line	Page	None	2687	4322	660	2095
2	100	Line	Page	4	2417	3328	689	2234
2	100	Line	6144 tokens	None	2377	3920	602	1916

VRDU Registration Form								
# Coord.	# Quant.	Segment	Chunking	JSON Indent	Input		Target	
					50 <sup>th</sup>	99 <sup>th</sup>	50 <sup>th</sup>	99 <sup>th</sup>
2	100	Line	Page	None	963	1578	79	147
2	100	Word	Page	None	3083	5196	101	349
4	100	Line	Page	None	1232	2017	91	177
2	1000	Line	Page	None	1052	1723	83	155
2	100	Line	Page	4	977	1592	92	160
2	100	Line	6144 tokens	None	963	1578	79	147

CORD								
# Coord.	# Quant.	Segment	Chunking	JSON Indent	Input		Target	
					50 <sup>th</sup>	99 <sup>th</sup>	50 <sup>th</sup>	99 <sup>th</sup>
2	100	Line	Page	None	342	869	355	1495
2	100	Word	Page	None	396	1067	375	1638
4	100	Line	Page	None	408	1139	422	1801
2	1000	Line	Page	None	364	959	376	1957
2	100	Line	Page	4	411	938	474	1997
2	100	Line	6144 tokens	None	342	869	355	1495

## A.5 SCHEMAS

In this section, we present the schemas used for the experiments of this paper. The schema for VRDU Ad-Buy Form, VRDU Registration Form, and CORD can be found in Figure 5, Figure 6 and Figure 7 respectively.

```
{
  "advertiser": "",
  "agency": "",
  "contract_num": "",
  "flight_from": "",
  "flight_to": "",
  "gross_amount": "",
  "line_item": [
    {
      "channel": "",
      "program_desc": "",
      "program_end_date": "",
      "program_start_date": "",
      "sub_amount": ""
    }
  ],
  "product": "",
  "property": "",
  "tv_address": ""
}
```

Figure 5: VRDU Ad-Buy Form Schema.

```
{
  "file_date": "",
  "foreign_principle_name": "",
  "registrant_name": "",
  "registration_num": "",
  "signer_name": "",
  "signer_title": ""
}
```

Figure 6: VRDU Registration Form Schema.



---

```

{
  "line_item": [
    {
      "discount_price": "",
      "identifier": "",
      "name": "",
      "other": "",
      "quantity": "",
      "sub_name": [],
      "sub_price": [],
      "sub_quantity": [],
      "subtotal_price": "",
      "total_price": "",
      "unit_price": ""
    }
  ],
  "subtotal": {
    "discount_price": "",
    "other": [],
    "service_price": "",
    "subtotal_price": [],
    "tax_price": []
  },
  "total": {
    "cash_price": [],
    "change_price": "",
    "credit_card_price": "",
    "emoney_price": "",
    "line_item_quantity_count": "",
    "line_item_type_count": "",
    "other": "",
    "total_price": []
  }
}

```

Figure 7: CORD Schema. Note that the original entity types have been renamed to more semantically meaningful names.

## A.6 SAMPLE PROMPTS AND COMPLETIONS

In this section, we present example of LMDX prompts and completions from the LLM on the VRDU Ad-Buy dataset to better showcase the format used. Figure 8 shows the original document with the line bounding boxes from OCR, Figure 9 shows the corresponding prompt and completion on that document with coordinate grounding tags, and Figure 10 shows the same prompt and completion, but with line index tags (used in ablation studies to showcase how the LLM can interpret the layout).

Print Date: 02/28/20 14:21:20 Page 1 of 1

### ORDER

**WSIL**  
**HARRISBURG**  
**PADUCAH**  
**CAPE GIRARDEAU**

<b>Order#</b>	<b>Order / Rev:</b> 14086		<b>WSIL-TV</b>
	<b>Alt Order #:</b>		<b>David Cisco</b>
	<b>Product Desc:</b> Mike Carr for Jackson Co States Atty	<b>Primary AE:</b>	<b>Local</b>
	<b>Estimate:</b>	<b>Sales Office:</b>	<b>Local</b>
	<b>Flight Dates:</b> 03/03/20 - 03/03/20	<b>Sales Region:</b>	
	<b>Original Date / Rev:</b> 02/28/20 / 02/28/20		
	<b>Order Type:</b> GENERAL		

<b>Agency</b>	<b>Name:</b> Committee to Elect Mike Carr	<b>Billing Type:</b> Cash	
	<b>Buying Contact:</b>	<b>Billing Calendar:</b> Calendar	
	<b>Billing Contact:</b>	<b>Billing Cycle:</b> EOM/EOC	
	101 S TOWER RD	<b>Agency Commission:</b> 0%	
	CARBONDALE, IL 62901-1930		

<b>Advertiser</b>	<b>Name:</b> Committee to Elect Mike Carr	<b>New Business Thru:</b>	
	<b>Demographic:</b> HB	<b>Order Separation:</b> 00:15:00	
	<b>Product Codes:</b> Candidates	<b>Advertiser External ID:</b> 43917	
	<b>Priority:</b> P-01	<b>Agency External ID:</b> 43917	
	<b>Revenue Codes:</b> DIR_POL, POL-CAND		

Bill Plan					Totals				
Start Date	End Date	# Spots	Gross Amount	Net Amount	Month	# Spots	Gross Amount	Net Amount	Rating
03/01/20	03/03/20	\$	\$600.00	\$600.00	March 2020	\$	\$600.00	\$600.00	0.00
					<b>Totals</b>	\$	\$600.00	\$600.00	0.00

Account Executives				
Account Executive	Sales Office	Sales Region	Start Date / End Date	Order %
David Cisco			Start Of Order - End Of Order	100%

Ln	Ch	Start	End	Inventory Code	Break	Start/End Time	Days	Len	Spots	Rate	Pri	Rtg	Type	Spots	Amount
N 1	3	03/03/20	03/03/20	Tuesday Prime Other	CM	7:00 PM-10:00 PM	-3-----	:30		\$200.00	P-03	0.00	NM		\$600.00
				Your Voice Your Vote		(7:00 PM-10:00 PM)									
		Start Date	End Date	Weekdays	Spots/Week	Rate	Rating								
		Week: 03/02/20	03/08/20	-3-----	3	\$200.00	0.00								
<b>Totals</b>														\$	\$600.00

Figure 8: Example document from VRDU Ad-Buy Form. The line bounding boxes detected by OCR are shown in red.

```

===== Prompt =====

<Document>
Print Date 02/28/20 14:21:20 75|02
Page 1 of 1 91|02
ORDER 49|04
WSIL 77|07
Orders 04|09
Order/Rev: 17|09
14086 29|09
HARRISBURG 78|09
PADUCAH 78|10
CAPE GIRARDEAU 78|12
Alt Order #: 17|11
Product Desc: 18|13
Mike Carr for Jackson Co States Atty 38|13
Estimate: 16|15
WSIL-TV 73|15
Flight Dates: 17|17
03/03/20 - 03/03/20 33|17
Primary AE: 57|17
David Cisco 73|17
Original Date/ Rev: 19|19
02/28/20/02/28/20 33|19
...
Start Date End Date Weekdays Spots/Week Rate Rating 31|66
Week: 03/02/20 03/08/20 -3-- 15|67
3 $200.00 0.00 47|67
Totals 79|69
3 87|69
$600.00 94|69
</Document>
<Task>
From the document, extract the text values and tags of the following
entities:
{"advertiser": "", "agency": "", "contract_num": "", "flight_from": "",
"flight_to": "", "gross_amount": "", "line_item": [{"channel": "",
"program_desc": "", "program_end_date": "", "program_start_date": "",
"sub_amount": ""}], "product": "", "property": "", "tv_address": ""}
</Task>
<Extraction>

===== Completion =====

{"advertiser": "Committee to Elect Mike Carr 37|36", "agency":
"Committee to Elect Mike Carr 37|24", "contract_num": "14086 29|09",
"flight_from": "03/03/20 33|17", "flight_to": "03/03/20 33|17",
"gross_amount": "$600.00 77|51", "line_item": [{"channel": "3 31|63",
"program_desc": "Tuesday Prime 31|63\nYour Voice Your Vote 29|64",
"program_end_date": "03/03/20 31|63", "program_start_date":
"03/03/20 31|63", "sub_amount": "$600.00 94|63"}], "product":
"Mike Carr for Jackson Co States Atty 38|13", "property": "WSIL 77|07",
"tv_address": null}
</Extraction>

```

Figure 9: VRDU Ad-Buy Form sample prompt and completion with 2 Coordinates for grounding tags. The document is truncated for easier visualization.

```

===== Prompt =====

<Document>
Print Date 02/28/20 14:21:20 <0>
Page 1 of 1 <1>
ORDER <2>
WSIL <3>
Orders <4>
Order/Rev: <5>
14086 <6>
HARRISBURG <7>
PADUCAH <8>
CAPE GIRARDEAU <9>
Alt Order #: <10>
Product Desc: <11>
Mike Carr for Jackson Co States Atty <12>
Estimate: <13>
WSIL-TV <14>
Flight Dates: <15>
03/03/20 - 03/03/20 <16>
Primary AE: <17>
David Cisco <18>
Original Date/ Rev: <19>
02/28/20/02/28/20 <20>
...
Start Date End Date Weekdays Spots/Week Rate Rating <105>
Week: 03/02/20 03/08/20 -3-- <106>
3 $200.00 0.00 <107>
Totals <108>
3 <109>
$600.00 <110>
</Document>
<Task>
From the document, extract the text values and tags of the following
entities:
{"advertiser": "", "agency": "", "contract_num": "", "flight_from": "",
"flight_to": "", "gross_amount": "", "line_item": [{"channel": "",
"program_desc": "", "program_end_date": "", "program_start_date": "",
"sub_amount": ""}], "product": "", "property": "", "tv_address": ""}
</Task>
<Extraction>

===== Completion =====

{"advertiser": "Committee to Elect Mike Carr <44>", "agency":
"Committee to Elect Mike Carr <29>", "contract_num": "14086 <6>",
"flight_from": "03/03/20 <16>", "flight_to": "03/03/20 <16>",
"gross_amount": "$600.00 <84>", "line_item": [{"channel": "3 <102>",
"program_desc": "Tuesday Prime Other CM <102>", "program_end_date":
"03/03/20 <102>", "program_start_date": "03/03/20 <102>",
"sub_amount": "$600.00 <100>"}], "product": "Mike Carr for Jackson Co
States Atty <12>", "property": "WSIL-TV <14>", "tv_address": null}
</Extraction>

```

Figure 10: VRDU Ad-Buy Form sample prompt and completion with line index for grounding tags, which does not communicate layout information. The document is truncated for easier visualization.

---

## A.7 CORD BASELINES DETAILS

**LayoutLMv3<sub>LARGE</sub> Baseline.** We follow the released implementation<sup>2</sup> for the LayoutLMv3<sub>LARGE</sub> model and the training protocol described in Huang et al. (2022) as closely as possible. In particular, we train the model for 80 epochs for each experiment on CORD (namely, 10, 50, 100, 200, and 800-document training sets), on the IOB tags of the leaf entities. One difference in our training is that, due to computational resource constraints, we use  $batch\_size = 8$  and  $learning\_rate = 2 \cdot 10^{-5}$ .

As the LayoutLMv3 model can only extract leaf entities, we design and heavily optimize a heuristic algorithm to group the leaf entities into hierarchical entities *menu*, *subtotal* and *total*. The best heuristics we could find are as follows:

- For the *subtotal* and *total* hierarchical entity types, since they appear only once per document, we group all their extracted sub-entities under a single *subtotal* and *total* entity, respectively.
- For *menu* hierarchical entity type, we observe that those entities usually occur multiple times on a document, and each *menu* has at most one *nm*, *num*, *unitprice*, *cnt*, *discountprice*, *price*, *itemsubtotal*, *etc* sub-entities and potentially multiple *sub\_nm*, *sub\_price* and *sub\_cnt* sub-entities. We also notice that the sub-entities aligned horizontally overwhelmingly belong to the same *menu* entity, and a *menu* entity can sometimes span over two or more consecutive horizontal lines. To leverage those observations, we perform a two-step grouping process for *menu* entities. First, we merge the extracted leaf sub-entities into horizontal groups, where a threshold of 0.5 on the intersection-over-union of the Y-axis was used for the determination of horizontal alignment. Second, we further merge the *consecutive* horizontal groups into *menu* entities, if and only if the horizontal groups do not have type duplication in any of the *nm*, *num*, *unitprice*, *cnt*, *discountprice*, *price*, *itemsubtotal*, and *etc* sub-entities (namely, those sub-entities only show up in at most one of the consecutive horizontal groups to be merged). We allow duplication of *sub\_nm*, *sub\_price* and *sub\_cnt* sub-entity types. After those two steps, we obtain the final *menu* entities.

**Donut Baseline.** We follow Donut released implementation<sup>3</sup> for the Donut benchmarking results on CORD. We use the default training configuration for all experiments on CORD (namely, 10, 50, 100, 200, and 800-document training sets), with the following difference: we reduce batch size from 8 to 4 due to computational resource constraints, and increase the number of train epochs from 30 to 60. For each experiment, checkpoint with the lowest loss on the dev set is selected and we report performance on test set. Normalized Tree Edit Distance accuracy scores produced by Donut evaluation code are reported (similar to all our other models).

---

<sup>2</sup><https://github.com/microsoft/unilm/tree/master/layoutlmv3>

<sup>3</sup><https://github.com/clovaai/donut>