

LongFin: Multimodal Document Understanding Model for Long Financial Domain Documents

Ahmed Masry, Amir Hajian

Arteria AI

ahmed.masry@arteria.ai, amir.hajian@arteria.ai

Abstract

Document AI is a growing research field that focuses on the comprehension and extraction of information from scanned and digital documents to make everyday business operations more efficient. Numerous downstream tasks and datasets have been introduced to facilitate the training of AI models capable of parsing and extracting information from various document types such as receipts and scanned forms. Despite these advancements, both existing datasets and models fail to address critical challenges that arise in industrial contexts. Existing datasets primarily comprise short documents consisting of a single page, while existing models are constrained by a limited maximum length, often set at 512 tokens. Consequently, the practical application of these methods in financial services, where documents can span multiple pages, is severely impeded. To overcome these challenges, we introduce LongFin, a multimodal document AI model capable of encoding up to 4K tokens. We also propose the LongForms dataset, a comprehensive financial dataset that encapsulates several industrial challenges in financial documents. Through an extensive evaluation, we demonstrate the effectiveness of the LongFin model on the LongForms dataset, surpassing the performance of existing public models while maintaining comparable results on existing single-page benchmarks.

Introduction

There has been a noticeable industrial interest surrounding the automation of data extraction from various documents, including receipts, reports, and forms to minimize manual efforts and enable seamless downstream analysis of the extracted data (Zhang et al. 2020; Xu et al. 2020). However, the process of parsing documents poses several challenges, including obscure information within scanned documents that may result in Optical Character Recognition (OCR) errors, complex layouts (such as tables), and intricate content structures.

To investigate and address these challenges, several datasets have been made available. These datasets encompass a wide range of tasks, such as classification (Pramanik, Mujumdar, and Patel 2022), semantic entity recognition (Park et al. 2019; Guillaume Jaume 2019), relation extraction (Guillaume Jaume 2019), question answering (Mathew,

Figure 1: First page from a 4-page example financial form in the LongForms dataset. The information in these documents is spread over a mix of tables and text spanning multiple pages which makes it challenging for short-context models.

Figure 1: First page from a 4-page example financial form in the LongForms dataset. The information in these documents is spread over a mix of tables and text spanning multiple pages which makes it challenging for short-context models.

Karatzas, and Jawahar 2021), and key information extraction (Huang et al. 2019). Nonetheless, a significant limitation shared by these datasets is that they mostly consist of single-page documents with a limited amount of content. As a consequence, these datasets fail to capture various challenges inherent in parsing lengthy documents spanning multiple pages, which are commonly encountered in the financial industry. Financial reports and documents can become exceedingly lengthy, necessitating a comprehensive understanding of the entire context to effectively analyze and extract pertinent information.

The limitations inherent in existing datasets have a direct impact on the capabilities of the proposed models. In the literature, two primary lines of work have emerged: (i) OCR-dependent architectures (Wang, Jin, and Ding 2022; Xu et al. 2020, 2021; Huang et al. 2022; Tang et al. 2023) (ii) OCR-free models (Kim et al. 2022; Lee et al. 2023). OCR-dependent models typically employ transformer-based text encoders and incorporate spatial information by leveraging the words' coordinates in the documents as additional embeddings. One notable exception is UDOP (Tang et al. 2023) which consists of an encoder-decoder architecture.

Conversely, OCR-free models typically employ a vision encoder to process the scanned document image and a text decoder to generate the desired information. Nevertheless, a common limitation shared by most of these models is their design and pretraining to handle a maximum of 512 tokens or process a single input image.

In this work, we introduce two main contributions. Firstly, we present the LongForms dataset, a comprehensive financial dataset primarily comprising 140 long forms where the task is formulated as named entity recognition. Due to privacy concerns and proprietary limitations, we were unable to utilize our internal resources to construct this dataset. Consequently, we obtained financial statements from the SEC website¹, aligning our tasks to encompass the significant challenges encountered in the financial documents which require a deep understanding of lengthy contexts. Secondly, we propose LongFin, a multimodal document understanding model capable of processing up to 4K tokens. Our approach builds upon LiLT (Wang, Jin, and Ding 2022), one of the state-of-the-art multimodal document understanding models. Additionally, we incorporate techniques that effectively extend the capabilities of text-only models, such as RoBERTa (Liu et al. 2019), to handle longer sequences, as demonstrated by Longformer (Beltagy, Peters, and Cohan 2020). By leveraging these techniques, our proposed model exhibits enhanced performance in processing lengthy financial forms. The efficacy of our approach is extensively evaluated, showcasing its effectiveness and paving the way for numerous commercial applications in this domain.

Related Work

Document Datasets

Several recently released datasets in the field of document understanding have contributed significantly to advancing research in this area. The RVL-CDIP dataset (Pramanik, Mujumdar, and Patel 2022) introduced a classification task, encompassing 400K scanned documents categorized into 16 classes, such as forms and emails. Another notable dataset, DocVQ (Mathew, Karatzas, and Jawahar 2021), focuses on document question answering and comprises 50K question-answer pairs aligned with 12K scanned images. In addition, the CORD dataset (Park et al. 2019) consists of 11K scanned receipts, challenging models to extract 54 different data elements (e.g., phone numbers and prices). Furthermore, the FUNSD dataset (Guillaume Jaume 2019) was proposed, featuring 200 scanned forms. This dataset primarily revolves around two key tasks: semantic entity recognition (e.g., header, question, answer) and relation extraction (question-answer pairs). FUNSD is particularly relevant to our dataset, LongForms, as it also mainly consists of forms. However, FUNSD and all the above-mentioned datasets mainly focus on short contexts, as they typically consist of single-page documents. In contrast, our LongForms dataset primarily consists of multi-page documents, presenting unique challenges that demand a comprehensive understanding of lengthy contexts which is common in the financial industry.

¹<https://www.sec.gov/edgar/>

Document Understanding Models

Numerous document understanding models have been developed to tackle the challenges posed by the aforementioned benchmark datasets. These models can be broadly categorized into two main groups: OCR-free and OCR-dependent models. OCR-free models, exemplified by Donut (Kim et al. 2022) and Pix2Struct (Lee et al. 2023), typically employ vision transformer-based encoders to process input images and text decoders to handle output generation. These models are often pretrained on OCR-related tasks, enabling them to comprehend the text embedded within scanned documents effectively. On the other hand, OCR-dependent models, including LayoutLM (Xu et al. 2020), LayoutLMv2 (Xu et al. 2021), LayoutLMv3 (Huang et al. 2022), LiLT (Wang, Jin, and Ding 2022), DocFormer (Pappalari et al. 2021) and UDOP (Tang et al. 2023), rely on external OCR tools to initially extract underlying text from scanned documents. To incorporate layout information, these models utilize specialized positional embeddings, encoding the coordinates of each word in the document. Additionally, some models, such as LayoutLMv2, LayoutLMv3, DocFormer, and UDOP, employ visual embeddings created by splitting the image into patches. These visual embeddings, along with the text and layout embeddings, are fed into the models. While LayoutLM, LayoutLMv2, LayoutLMv3, DocFormer, and LiLT adopt an encoder-only architecture, UDOP is based on the T5 model (Raffel et al. 2020), which follows an encoder-decoder architecture. Despite the impressive achievements of these models, they share a common limitation: they are typically designed to process a single page or a maximum of 512 tokens, thereby restricting their applicability to multi-page documents. (Pham et al. 2022) proposed a multimodal document understanding model that can process up to 4096 tokens, however their code is not publicly available and their model performance deteriorates on the short-context datasets such as FUNSD (Guillaume Jaume 2019). In contrast, our proposed model, LongFin, works efficiently on both short and long contexts (to up 4096 tokens), making it particularly well-suited for a variety of real-world industrial applications.

LongForms Dataset

Due to privacy constraints, we are unable to utilize internal documents for dataset construction. Instead, we turn to publicly available financial reports and tailor our dataset, LongForms, to emulate the challenges encountered in our proprietary datasets. This approach ensures the task’s alignment with real-world financial contexts without violating privacy.

Dataset Collection & Preparation

To construct LongForms, we leverage the EDGAR database², a comprehensive repository of financial filings and reports submitted by US companies. These filings are based on different financial form types (e.g., 10-K, 10-Q) which vary in structure and content. Our dataset primarily centers around the SEC Form 10-Q, which provides a detailed quarterly report on a company’s finances. This specific form is chosen

²<https://www.sec.gov/edgar/>

Dataset Split	#Forms	#Pages	#Words	#Entities
Train	105	514	125094	843
Test	35	171	43364	285
Overall	140	685	168458	1128

Table 1: LongForms dataset statistics.

due to its similarity in both structure and content to the documents we frequently encounter in the financial services industry.

We download 140 10-Q forms that were published between 2018 and 2023. This deliberate decision to keep the dataset relatively small is intended to mirror the limited data challenges commonly encountered in real-world scenarios, particularly in the finance domain, where strict data confidentiality prevents access to large-scale datasets. Consequently, it is common practice to construct smaller datasets that mimic the proprietary datasets (Madl et al. 2023). Furthermore, our dataset size aligns with recently published datasets, such as the FUNSD dataset (Guillaume Jaume 2019) which primarily consists of single-page forms. Inspired by the FUNSD dataset, we perform a random split of the LongForms dataset and divide the dataset into 105 training documents, which account for 75% of the total dataset, and 35 testing documents, representing the remaining 25%.

Dataset Description & Setup

Our dataset, LongForms, is formulated as a Named Entity Recognition (NER) task. The dataset consists of N examples, denoted as $D = \{d_i \mid w_i, b_i, n_i\}_{i=1}^N$, where d_i represents a PDF document, w_i represents the list of words, b_i represents the list of bounding boxes, and n_i represents a list of entities present in the document. To obtain the words (w_i) and their bounding boxes (b_i), each PDF document is processed using the `pdftotext`³ tool. Moreover, we define six entity types: (i) Total assets, (ii) Cash at the beginning of the period (Beginning Cash), (iii) Cash at the end of the period (End Cash), (iv) Cash provided by financial activities (Financial Cash), (v) Net change in cash (Change in Cash), and (vi) Quarter Keys. As shown in Table 1, our LongForms dataset contains 140 forms that consist of 685 pages, 168458 words, and 1128 entities in total. The models are trained to predict n_i given both w_i and b_i .

LongFin Model

Architecture

Figure 2 illustrates the overall architecture of our proposed model, LongFin, which builds upon recently published models: LiLT (Wang, Jin, and Ding 2022) and Longformer (Beltagy, Peters, and Cohan 2020). Similar to LiLT (Wang, Jin, and Ding 2022), LongFin comprises three primary components: a text encoder, a layout encoder, and the Bi-CM (bidirectional attention complementation mechanism) layer (Wang, Jin, and Ding 2022). However, LongFin introduces additional mechanisms, namely sliding window local attention and interval-based global attention, to effectively handle

long contexts within both the text and layout encoders. One key advantage of LongFin is its ability to scale linearly with the input sequence length, in contrast to the quadratic scaling ($O(n^2)$) observed in the original transformers’ (Vaswani et al. 2017) attention mechanism. This linear scaling, inspired by the Longformer model (Beltagy, Peters, and Cohan 2020), allows LongFin to efficiently handle long contexts up to 4K tokens.

Text Encoder For the text encoder in LongFin, we adopt the Longformer (Beltagy, Peters, and Cohan 2020) model, which has been pretrained to handle long textual contexts of up to 4096 tokens. As depicted in Figure 2a, the input to the text encoder consists of two types of embeddings: text embeddings (E_T) and absolute position embeddings (E_P). These embeddings are added together to produce the final embeddings (E_{final}). Subsequently, a layer normalization (Ba, Kiros, and Hinton 2016) operation is applied, and the resulting output is fed into the encoder.

The attention mechanism in LongFin incorporates two types of attention: local attention and global attention. The local attention employs a sliding window approach, where each token attends to the 512 local tokens surrounding it. On the other hand, the global attention involves a set of global tokens, selected at intervals of 100. While other approaches (Beltagy, Peters, and Cohan 2020; Pham et al. 2022) may employ different methods for selecting global tokens, such as random selection or task-specific strategies, we limit our experimentation to interval-based selection for simplicity and due to limited computational resources. Each token in the input sequence attends to these global tokens, in addition to its local context as shown in Figure 2b. This combination of local and global attention mechanisms enhances the model’s ability to capture both local context and broader global dependencies within the long input sequences.

Layout Encoder For the layout encoder in LongFin, we adopt the layout encoder utilized in the LiLT model (Wang, Jin, and Ding 2022). Similar to the text encoder, the input for the layout encoder comprises two types of embeddings: absolute position embeddings and layout embeddings. Each word in the input document is associated with a bounding box that defines its location within the document layout. This bounding box is represented by four numbers: x, y, x_1 , and y_1 , which correspond to the coordinates of the top-left and bottom-right points of the bounding box. To normalize these coordinates within the range [0,1000], we use the page’s height and width.

To generate the layout embedding for each word, each coordinate in the normalized bounding box is used to obtain an embedding vector. The different coordinates’ embedding vectors are then concatenated and projected using a linear layer. The resulting layout embeddings are added to the absolute position embeddings to obtain the final embeddings. These final embeddings are then fed into the layout encoder. Similar to the text encoder, we also employ the local & global attention mechanisms in the layout encoder to process long sequences.

Bi-CM To facilitate communication between the text encoder and layout encoder, we incorporate the Bi-CM layer

³<https://pypi.org/project/pdftotext/>

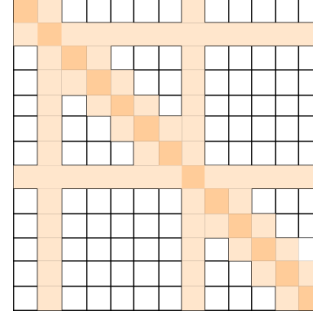
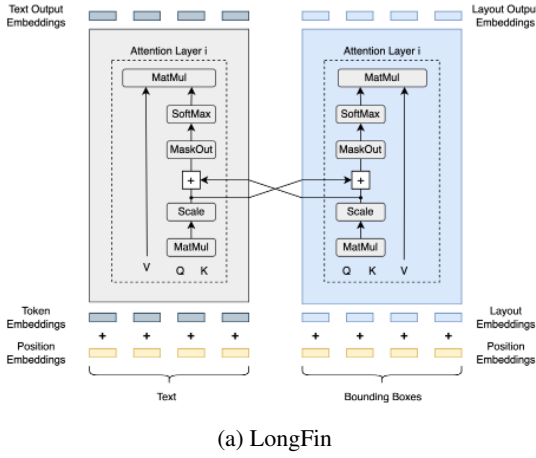


Figure 2: (a) The architecture of the LongFin model. It mainly consists of two encoders: text encoder and layout encoder which are connected through the Bi-CM layer. (b) visualization of the employed local (sliding window) and global attention mechanisms to process long sequences.

from the LiLT model (Wang, Jin, and Ding 2022). As depicted in Figure 2a, the Bi-CM layer adds the scores resulting from the multiplication of keys and queries from both encoders. In LiLT, a detach operation is applied to the scores generated by the text encoder before passing them to the layout encoder. This detachment prevents the layout encoder from backpropagating into the text encoder during pretraining, promoting better generalization when fine-tuning the model with different language text encoders. However, since our focus is primarily on the English language for our applications, we have chosen to remove the detach operation to expedite pretraining, given our limited computational resources.

Pretraining

To pretrain LongFin, we utilize the IIT-CDIP (Lewis et al. 2006) dataset which contains 11M scanned images that make up 6M documents. We obtain the OCR annotations (words and their bounding boxes) from OCR-IDL (Biten et al. 2022) which used the WS-Texttract-PI⁴. We initialize our text encoder from Longformer (Beltagy, Peters, and Cohan 2020) and our layout encoder from LiLT (Wang, Jin, and Ding 2022) layout encoder. Since the LiLT layout encoder was pretrained on inputs with a maximum length of 512 tokens, we copy LiLT’s pretrained positional embeddings eight times to initialize our layout encoder positional embeddings, which consist of 4096 embedding vectors. This enables the layout encoder to handle longer sequences while leveraging the pretrained positional information from the LiLT model.

For the pretraining of LongFin, we employ the Masked Visual-Language Modeling task (Devlin et al. 2019; Wang, Jin, and Ding 2022). In this task, 15% of the tokens in the input to the text encoder are masked. In 80% of the cases, we replace the masked tokens with the [MASK] token. In 10% of the cases, we replace the masked tokens with random tokens. In the remaining 10%, we keep the original token unchanged. Inspired by Longformer (Beltagy, Peters,

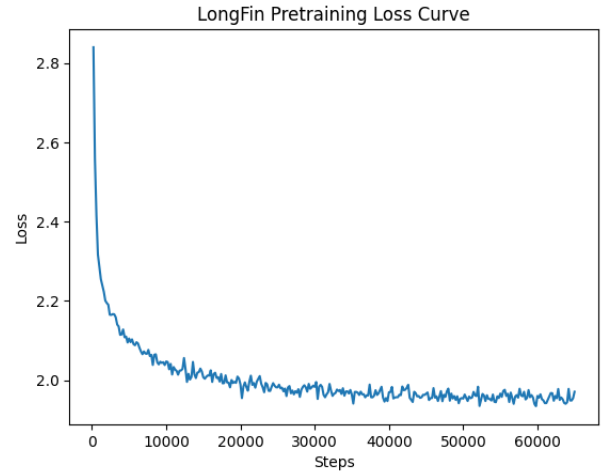


Figure 3: LongFin pretraining loss curve. The loss starts at 2.84 and oscillated between 1.97 and 1.94 near convergence.

and Cohan 2020), we pretrain the model for 65K steps with a learning rate of $3e-5$ and batch size of 12 on one 100 GPU. We set the warmup steps to 500 and use the AdaFactor optimizer (Shazeer and Stern 2018). Also, we utilize gradient checkpointing (Chen et al. 2016) to enable using a large batch size. The pretraining loss curve is shown in Figure 3

Experiments & Evaluation

Tasks & Datasets

To assess the generalizability of LongFin on both short and long contexts, we evaluate LongFin on two existing short (single-page) datasets: FUNSD (Guillaume Jaume 2019) and CORD (Park et al. 2019) to show the generalizability of our model on short contexts as well as our newly created LongForms dataset.

FUNSD: This dataset comprises 200 scanned forms and requires models to extract four main entities: headers, questions, answers, and other relevant information. Additionally, it involves linking questions with their corresponding an-

⁴<https://aws.amazon.com/texttract/>

Model	Modalities	FUNSD ()	CORD ()
Short Context Models (512 tokens)			
BERT _{B SE}	T	60.26	89.68
RoBERTa _{B SE}	T	66.48	93.54
LayoutLM _{B SE}	T+L	79.27	94.72
LayoutLMv2 _{B SE}	T+L+V	82.76	94.95
LayoutLMv3 _{B SE}	T+L+V	<u>90.29</u>	<u>96.56</u>
DocFormer _{B SE}	T+L+V	83.34	96.33
LiLT _{B SE}	T+L	88.41	96.07
Long Context Models (4096 tokens)			
Longformer _{B SE}	T	71.4	90.41
(Pham et al. 2022)	T+L	77.1	— ⁶
LongFin _{B SE} (ours)	T+L	<u>87.03</u>	<u>94.81</u>

Table 2: Accuracy of the different models on FUNSD and CORD datasets. The second column shows the modalities used by each model where T refers to Text, L refers to Layout, and V refers to Vision.

swers, thereby encompassing named entity recognition and relation extraction tasks. We mainly focus on the named entity recognition task and use the entity-level F1 score as our evaluation metric.

CORD: With over 11,000 receipts, this dataset focuses on extracting 54 different data elements (e.g., phone numbers) from receipts. The task can be formulated as named entity recognition or token classification. For evaluation, we use the entity-level F1 score.

Baselines

To demonstrate the effectiveness of LongFin on our LongForms dataset, we compare it against a set of publicly available text and text+layout baselines that are capable of handling both short and long input sequences. For the text baselines, we select the following models: (i) BERT (Devlin et al. 2019) which is a widely used text-based model known for its strong performance on short context tasks (512 tokens), (ii) Longformer (Beltagy, Peters, and Cohan 2020) which is specifically designed to handle text long texts (up to 4096 tokens). For the text+layout baseline, we utilize LiLT (Wang, Jin, and Ding 2022), which is one of the state-of-the-art models for document understanding⁵. For the short context models, we split the LongForms documents into chunks that can fit within 512 tokens. Table 5 shows the hyperparameters of the different models when finetuning on the LongForms dataset. It also presents the hyperparameters we used when finetuning LongFin on the previous single-page datasets. All the finetuning experiments were performed on one A100 and one T4 GPUs.

Results

Previous (Single-Page) Datasets

As shown in Table 2, LongFin outperforms other long-context models such as Longformer (Beltagy, Peters, and Cohan 2020) and (Pham et al. 2022) on the previous datasets that mainly consist of single-page documents. The performance disparity is particularly pronounced on the FUNSD dataset (Guillaume Jaume 2019), where all documents have

⁵LayoutLMv3 (Huang et al. 2022) is another state-of-the-art document understanding model, but its usage is limited to non-commercial applications

⁶The code of (Pham et al. 2022) is not publicly available.

Model	Modalities	Precision	Recall	F1
Short Context Models				
BERT _{B SE}	T	35.82	30.18	32.76
LiLT _{B SE}	T+L	35.55	43.24	39.02
Long Context Models				
Longformer _{B SE}	T	49.50	45.20	47.25
LongFin _{B SE} (ours)	T+L	47.67	56.16	51.57

Table 3: Results of the different models on LongForms. We evaluate using the entity-level F1 score.

Model	Beginning Cash	Ending Cash	Financial Cash	Change in Cash	Quarter Keys	Total Assets
LiLT _{B SE}	27.39	35.89	7.40	30.55	64.93	47.05
LongFin _{B SE} (ours)	47.61	56.16	15.15	45.94	75.17	45.71

Table 4: Ablation results of LiLT and LongFin on the LongForms dataset by entity.

very short textual content (less than 512 tokens). Notably, LongFin also achieves comparable performance to the short-context models on these datasets. This comparison highlights the superior generalization ability of our model, LongFin, which performs well on both short and long contexts. In contrast, the performance of (Pham et al. 2022) model deteriorates on short-context documents.

LongForms Dataset

As presented in Table 3, the performance results on our LongForms dataset highlight the advantage of our model, LongFin, compared to the short-context models. This observation emphasizes the significance of long-context understanding when working with financial documents. There is also a noticeable difference in performance between the text models (BERT (Devlin et al. 2019) and Longformer (Beltagy, Peters, and Cohan 2020)) and text+layout models (LiLT (Wang, Jin, and Ding 2022) and LongFin). This is mainly because the documents in LongForms contain diverse layouts that might be challenging for text-only models.

To provide a deeper analysis of the results on the LongForms dataset, we conduct ablations and report metrics by entity for both LiLT (Wang, Jin, and Ding 2022) and LongFin, as shown in Table 4. We notice that the gap in performance is more significant in the entities that are typically found in long tables such as Beginning Cash, Ending Cash, Financial Cash, and Change in Cash. To illustrate the challenges posed by long tables, we present an example from our test set in Figure 4. In the example, the table header indicates "Nine Months," implying that the table includes information for a nine-month period that should not be extracted as we are only interested in the financial information per quarter "Three Months". Due to the large number of rows and content in the table, the short-context models may not be able to include all the table information in a single forward pass of 512 tokens. Consequently, when the long documents are split into chunks, such tables might be divided as well, leading to the short-context models losing important context when making predictions.

Limitations

Despite the effectiveness of our model, LongFin, on both short and long context document understanding datasets, it has a few limitations. First, LongFin was trained and evaluated on the English language only. In future, we plan to

7/7/23, 4:26 PM mlb_04.htm
[Table of Contents](#)

MAJOR LEAGUE FOOTBALL, INC.
STATEMENTS OF CASH FLOWS
(UNAUDITED)

	For the Nine Months Ended, January 31,	
	2021	2020
CASH FLOWS FROM OPERATING ACTIVITIES		
Net loss	\$ (3,024,627)	\$ (990,021)
Adjustments to reconcile net loss to net cash used in operating activities:		
Amortization of debt discount on convertible secured promissory note	-	10,083
Amortization of debt discount on convertible unsecured promissory notes	4,010	132,072
Conversion fees on convertible unsecured promissory notes	3,750	1,500
Accretion of put premium liability	-	193,971
Initial fair value of conversion option liability	-	350,072
Loss (gain) from change in fair value of conversion option liability	2,630,554	(255,529)
Changes in operating assets and liabilities:		
Prepaid consulting	-	(2,500)
Accounts payable	74,695	104,736
Accounts payable - related parties	95,000	(10,947)
Accrued expenses	16,052	15,120
Accrued interest	68,857	61,147
Accrued interest - related party	3,304	-
Net cash used in operating activities	(128,405)	(390,296)
CASH FLOWS FROM INVESTING ACTIVITIES		
Trademark legal and filing fees	-	(500)
Football equipment	-	(46,223)
Office equipment	-	(11,000)
Net cash used in investing activities	-	(57,723)
CASH FLOWS FROM FINANCING ACTIVITIES		
Proceeds from issuance of convertible unsecured promissory notes, net of issue costs	-	302,200
Proceeds from issuance of notes payable	-	55,284
Proceeds from issuance of notes payable - related party	30,000	-
Proceeds from issuance of common stock	97,400	100,000
Net cash provided by financing activities	127,400	457,484
NET INCREASE (DECREASE) IN CASH	(1,005)	9,465
CASH - BEGINNING OF PERIOD	3,796	5,417
CASH - END OF PERIOD	\$ 2,791	\$ 14,882
SUPPLEMENTAL DISCLOSURE OF CASH FLOWS		
CASH PAID FOR INCOME TAXES	\$ -	\$ -
CASH PAID FOR INTEREST	\$ -	\$ -

See accompanying condensed notes to these unaudited financial statements.

F-6

Figure 4: Page 6 from an example document from the LongForms test set. Since the original document has 6 pages which can not fit in a single forward pass of 512 tokens, the document is split into several chunks, leading to a loss of important content. For example, in this table from the sixth page, the context from the top is crucial to decide whether to pick the net change in cash entity or not, since we are only interested to extract quarter information "Three months" periods only.

expand it to support multiple languages. Second, although LongFin maximum input length (4096 tokens) can accommodate the multi-page documents in the LongForms dataset as well as most of our proprietary datasets, it might not accommodate certain financial documents that contain tens of pages. To overcome this limitation, we may consider further expanding the positional embeddings to accommodate 16K tokens similar to the LED model (Beltagy, Peters, and Cohan 2020) or explore utilizing a model architecture that uses relative position embeddings (Shaw, Uszkoreit, and Vaswani 2018) such as T5 (Raffel et al. 2020) instead of the absolute position embeddings. Third, due to limited computational resources, we have not explored many different hyperparameters setup. Hence, there might be room for improvement in our model performance. Finally, while our LongForms shed the light on long context understanding challenges which are frequent in the financial industry, it is still limited in size. We encourage the research community to explore this undercharted area of research since it has various commercial applications in many industries such as finance and legal.

Conclusion

We introduce LongFin, a multimodal document I model designed to handle lengthy documents. Additionally, we present the LongForms dataset, which aims to replicate real-world challenges in understanding long contexts, specifically in the financial industry. Through our evaluation, we demonstrate the superior performance of LongFin on the LongForms dataset, which comprises multi-page doc-

Experiment	Steps	Learning Rate	Batch Size
Finetuning on the LongForms dataset			
BERT (Devlin et al. 2019)	10000	2e-5	4
LiLT (Wang, Jin, and Ding 2022)	8000	2e-5	4
Longformer (Beltagy, Peters, and Cohan 2020)	6000	2e-5	4
LongFin (Ours)	6000	2e-5	4
Finetuning LongFin on the previous datasets			
LongFin _{FUNSD}	6000	2e-5	4
LongFin _{CORD}	6000	5e-5	8

Table 5: Training details for finetuning the different models on the LongForms dataset. The lower section also shows the hyperparameters used in finetuning LongFin on the previous single-page datasets.

uments, while achieving comparable results on previous datasets consisting of single-page documents. Moving forward, our plan is to deploy LongFin after training it on our proprietary datasets in the finance domain. Furthermore, we are working on extending LongFin to support different languages.

Ethical Statement

If the documents used in our LongForms dataset is collected from the EDG R database which grants the right to use and distribute their data without permissions⁷. The dataset annotation process were accomplished by data anno-

⁷<https://www.sec.gov/privacy#dissemination>

tators who are fairly compensated. We provide the hyperparameters and experimental setups of our experiments to ensure the reproducibility of our work. Moreover, the models, LiLT (Wang, Jin, and Ding 2022) and Longformer (Beltagy, Peters, and Cohan 2020), on which our LongFin model is built are published under permissive licenses⁸⁹ that allow commercial use.

References

- ppalaraju, S.; Jasani, B.; Kota, B. U.; Xie, Y.; and Manmatha, R. 2021. DocFormer: End-to-End Transformer for Document Understanding. *arXiv:2106.11539*.
- Ba, J. L.; Kiros, J. R.; and Hinton, G. E. 2016. Layer Normalization. *arXiv:1607.06450*.
- Beltagy, I.; Peters, M. E.; and Cohan, . 2020. Longformer: The Long-Document Transformer. *arXiv:2004.05150*.
- Biten, . F.; Tito, R.; Gomez, L.; Valveny, E.; and Karatzas, D. 2022. Ocr-idl: Ocr annotations for industry document library dataset. *arXiv preprint arXiv:2202.12985*.
- Chen, T.; Xu, B.; Zhang, C.; and Guestrin, C. 2016. Training Deep Nets with Sublinear Memory Cost. *arXiv:1604.06174*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805*.
- Guillaume Jaume, J.-P. T., Hazim Kemal Ekenel. 2019. FUNSD: Dataset for Form Understanding in Noisy Scanned Documents. In *accepted to ICD R-OST*.
- Huang, Y.; Lv, T.; Cui, L.; Lu, Y.; and Wei, F. 2022. LayoutLMv3: Pre-Training for Document I with Unified Text and Image Masking. In *Proceedings of the 30th ACM International Conference on Multimedia, MM '22*, 4083–4091. New York, NY, US : ssociation for Computing Machinery. ISBN 9781450392037.
- Huang, Z.; Chen, K.; He, J.; Bai, X.; Karatzas, D.; Lu, S.; and Jawahar, C. V. 2019. ICD R2019 Competition on Scanned Receipt OCR and Information Extraction. In *2019 International Conference on Document nalysis and Recognition (ICD R)*. IEEE.
- Kim, G.; Hong, T.; Yim, M.; Nam, J.; Park, J.; Yim, J.; Hwang, W.; Yun, S.; Han, D.; and Park, S. 2022. OCR-free Document Understanding Transformer. *arXiv:2111.15664*.
- Lee, K.; Joshi, M.; Turc, I.; Hu, H.; Liu, F.; Eisenschlos, J.; Khandelwal, U.; Shaw, P.; Chang, M.-W.; and Toutanova, K. 2023. Pix2Struct: Screenshot Parsing as Pretraining for Visual Language Understanding. *arXiv:2210.03347*.
- Lewis, D.; gam, G.; rgamon, S.; Frieder, O.; Grossman, D.; and Heard, J. 2006. Building a Test Collection for Complex Document Information Processing. In *Proceedings of the 29th nual International CM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '06*, 665–666. New York, NY, US : ssociation for Computing Machinery. ISBN 1595933697.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: Robustly Optimized BERT Pretraining pproach. *arXiv:1907.11692*.
- Madl, T.; Xu, W.; Choudhury, O.; and Howard, M. 2023. pproximate, dapt, nonymize (3): a Framework for Privacy Preserving Training Data Release for Machine Learning. *arXiv:2307.01875*.
- Mathew, M.; Karatzas, D.; and Jawahar, C. V. 2021. DocVQ : Dataset for VQ on Document Images. *arXiv:2007.00398*.
- Park, S.; Shin, S.; Lee, B.; Lee, J.; Surh, J.; Seo, M.; and Lee, H. 2019. COD: Consolidated Receipt Dataset for Post-OCR Parsing.
- Pham, H.; Wang, G.; Lu, Y.; Florêncio, D. . F.; and Zhang, C. 2022. Understanding Long Documents with Different Position- ware ttentions. *rXiv, abs/2208.08201*.
- Pramanik, S.; Mujumdar, S.; and Patel, H. 2022. Towards a Multi-modal, Multi-task Learning based Pre-training Framework for Document Representation Learning. *arXiv:2009.14457*.
- Raffel, C.; Shazeer, N.; Roberts, .; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, 21(140): 1–67.
- Shaw, P.; Uszkoreit, J.; and Vaswani, . 2018. Self-ttention with Relative Position Representations. In *Proceedings of the 2018 Conference of the North merican Chapter of the ssociation for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 464–468. New Orleans, Louisiana: ssociation for Computational Linguistics.
- Shazeer, N.; and Stern, M. 2018. dafactor: daptive Learning Rates with Sublinear Memory Cost. *arXiv:1804.04235*.
- Tang, Z.; Yang, Z.; Wang, G.; Fang, Y.; Liu, Y.; Zhu, C.; Zeng, M.; Zhang, C.; and Bansal, M. 2023. Unifying Vision, Text, and Layout for Universal Document Processing. *arXiv:2212.02623*.
- Vaswani, .; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, . N.; Kaiser, L.; and Polosukhin, I. 2017. ttention Is ll You Need. *arXiv:1706.03762*.
- Wang, J.; Jin, L.; and Ding, K. 2022. LiLT: Simple yet Effective Language-Independent Layout Transformer for Structured Document Understanding. In *Proceedings of the 60th nual Meeting of the ssociation for Computational Linguistics (Volume 1: Long Papers)*, 7747–7757. Dublin, Ireland: ssociation for Computational Linguistics.
- Xu, Y.; Li, M.; Cui, L.; Huang, S.; Wei, F.; and Zhou, M. 2020. LayoutLM: Pre-Training of Text and Layout for Document Image Understanding. In *Proceedings of the 26th CM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20*, 1192–1200. New York, NY, US : ssociation for Computing Machinery. ISBN 9781450379984.

⁸<https://github.com/allenai/longformer>

⁹<https://github.com/jpWang/LiLT>

Xu, Y.; Xu, Y.; Lv, T.; Cui, L.; Wei, F.; Wang, G.; Lu, Y.; Florencio, D.; Zhang, C.; Che, W.; Zhang, M.; and Zhou, L. 2021. LayoutLMv2: Multi-modal Pre-training for Visually-rich Document Understanding. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2579–2591. Online: Association for Computational Linguistics.

Zhang, R.; Yang, W.; Lin, L.; Tu, Z.; Xie, Y.; Fu, Z.; Xie, Y.; Tan, L.; Xiong, K.; and Lin, J. 2020. Rapid adaptation of BERT for Information Extraction on Domain-Specific Business Documents. arXiv:2002.01861.