CHULO: CHUNK-LEVEL KEY INFORMATION REPRE SENTATION FOR LONG DOCUMENT PROCESSING

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformer-based models have achieved remarkable success in various Natural Language Processing (NLP) tasks, yet their ability to handle long documents is constrained by computational limitations. Traditional approaches, such as truncating inputs, sparse self-attention, and chunking, attempt to mitigate these issues, but they often lead to information loss and hinder the model's ability to capture long-range dependencies. In this paper, we introduce ChuLo, a novel chunk representation method for long document classification that addresses these limitations. Our ChuLo groups input tokens using unsupervised keyphrase extraction, emphasizing semantically important keyphrase based chunk to retain core document content while reducing input length. This approach minimizes information loss and improves the efficiency of Transformer-based models. Preserving all tokens in long document understanding, especially token classification tasks, is especially important to ensure that fine-grained annotations, which depend on the entire sequence context, are not lost. We evaluate our method on multiple long document classification tasks and long document token classification tasks, demonstrating its effectiveness through comprehensive qualitative and quantitative analyses.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 029

Transformer-based models (Vaswani et al., 2017), including LLMs (Radford, 2018; Radford et al., 2019; Brown et al., 2020; Ouyang et al., 2022; Touvron et al., 2023a;b; Chowdhery et al., 2023; 031 Anil et al., 2023; Dubey et al., 2024), have achieved remarkable success across a wide range of Natural Language Processing (NLP) tasks, including Machine Translation, Text Summarization, Text 033 Generation, and Text Classification. A key factor behind their success is the self-attention mech-034 anism, which allows the model to capture long-range dependencies by computing the similarity between any two tokens and aggregating information accordingly. However, this mechanism incurs 036 a quadratic computational cost in terms of both time and space, relative to input length. This compu-037 tational burden makes it difficult for Transformer-based models to scale to long documents, limiting 038 their application to real-world data with unrestricted document lengths.

To address this challenge, several approaches have been proposed for applying Transformer-based 040 models to long documents while managing computational resources. One of them is truncating, 041 where the model discards content exceeding a predefined input length. For instance, BERT (Kenton 042 & Toutanova, 2019) processes up to 512 tokens, and LLaMa (Touvron et al., 2023a) handles up 043 to 2048 tokens, with any additional content being ignored. Another one is sparse self-attention, 044 which reduces computational complexity by restricting each query token to attend only to a subset of key tokens (Child et al., 2019; Beltagy et al., 2020; Zaheer et al., 2020; Wei et al., 2021; Li et al., 2023a). Lastly, chunking divides long documents into smaller, manageable segments that are 046 processed independently by the model (Zhao et al., 2021; Zhang et al., 2022). 047

 While these methods enable Transformer-based models to process long documents, they have limitations. Truncation risks discarding important information that falls beyond the maximum input length. Although more efficient, Sparse attention reduces each token's receptive field, leading to potential information loss from the neglected tokens. Similarly, chunking breaks the input into isolated segments, which can disrupt long-range dependencies critical for a comprehensive understanding of the document. Preserving all tokens is particularly important in tasks that require fine-grained tokenlevel understanding, such as token classification. In such tasks, dropping tokens can severely impact

the accuracy of fine-grained annotations, which often depend on the full context of the document. 055 Therefore, there is a need for methods that can handle long documents efficiently while retaining all 056 key information from the input.

057 In this paper, we introduce ChuLo, a novel chunk-level key information representation method that addresses these challenges in long document classification and token classification. Our method 059 reduces input length while minimizing information loss by strategically grouping tokens using un-060 supervised keyphrase extraction. By identifying and emphasizing semantically important tokens, 061 ChuLo ensures that each chunk retains the core content of the document. The resulting chunk repre-062 sentation is used for training Transformer models, with more weight assigned to keyphrases to make 063 them more salient in each chunk. We evaluate ChuLo on various long document classification tasks 064 and long document token classification tasks, demonstrating its effectiveness through competitive results and thorough analysis. 065

066 The key contributions of this paper are as follows: 1) Novel Chunk Representation Method: We 067 introduce ChuLo, a chunk representation method for long document understanding that leverages 068 unsupervised keyphrase extraction to prioritize semantically important information, effectively re-069 ducing input length while preserving core content. 2) Enhanced Document and Token Classi-070 fication: Our proposed method is designed to handle both document-level and token-level tasks, 071 addressing the limitations of existing models in retaining fine-grained annotations and global context in long documents. 3) Scalable and Efficient Solution: ChuLo offers a scalable and efficient 072 approach for long document processing, making it suitable for various NLP applications where han-073 dling long-range dependencies and context preservation are critical. 074

075 076

077 078

079

2 **RELATED WORK**

2.1 LONG DOCUMENT UNDERSTANDING

Document understanding includes two directions: global understanding (e.g., document classification) and token-level understanding (e.g., named entity recognition). With Transformer-based 081 models, document length impacts classification performance. Approaches, shown in Appendix A.1, to applying Transformer-based models on long document classification can be divided into input 083 processing and Transformer architecture optimization. Input processing involves truncating and 084 chunking. Truncating drops tokens exceeding the model's input length (Park et al., 2022), while 085 chunking segments documents into smaller parts processed separately. For example, Hierarchical Transformer (Pappagari et al., 2019) splits documents into non-overlapping chunks and computes 087 chunk representations. RoR (Zhao et al., 2021) generates regional answers from chunks, which 088 are combined for the final answer. However, neither considers the entire document context when 089 chunking. Transformer architecture optimization includes two strategies: improving self-attention efficiency and incorporating RNN concepts. Longformer (Beltagy et al., 2020), BigBird (Zaheer et al., 2020), and others (Roy et al., 2021) use sparse attention, restricting queries to focus on spe-091 cific keys. Other methods (Peng et al., 2021; Wang et al., 2020; Choromanski et al., 2020) approx-092 imate self-attention with reduced complexity. Meanwhile, (Dai et al., 2019; Hutchins et al., 2022; 093 Li et al., 2023b) integrate cache memory to update history information. These approaches involve 094 performance-efficiency trade-offs, making it valuable to explore improving performance by reduc-095 ing input length. In document NER, text length has received less attention. Recent research focuses 096 on low-resource languages (Çetindağ et al., 2023; Mengliev et al., 2024), complex domain-specific texts (Park et al., 2023; Bhattacharya et al., 2023), prompt-based large model methods (Wang et al., 098 2023; Dagdelen et al., 2024; Hu et al., 2024), and multimodal NER (Yu et al., 2023; Zhang et al., 099 2023). Our work addresses these existing challenges by introducing a novel chunk representation 100 that reduces input length while retaining key information, improving both classification and token-101 level tasks through semantic grouping of important phrases. Our method will preserve global and 102 local features, achieving better efficiency and performance compared to existing models.

103

104 2.2 UNSUPERVISED KEYPHRASE EXTRACTION

105

Unsupervised keyphrase extraction automatically identifies representative phrases from a document 106 to summarize its content without requiring labeled data (Hasan & Ng, 2014). Based on the fea-107 tures used, unsupervised methods can be categorized into three types: statistics-based, graph-based,

108 and embedding-based (Kong et al., 2023). Statistics-based methods rank candidate phrases using 109 features like TfIdf (El-Beltagy & Rafea, 2009), co-occurrence (Liu et al., 2009), and context statis-110 tics (Campos et al., 2020; Won et al., 2019). Graph-based methods construct a graph where nodes 111 are candidate phrases and edges represent their relationships, such as TextRank (Mihalcea & Ta-112 rau, 2004) and its variants (Wan & Xiao, 2008; Bougouin et al., 2013; Florescu & Caragea, 2017; Yu & Ng, 2018). Embedding-based methods leverage distributional representations of words or 113 sentences to achieve state-of-the-art performance, as seen in EmbedRank (Bennani-Smires et al., 114 2018), SIFRank (Sun et al., 2020), and PromptRank (Kong et al., 2023). Although these methods 115 have shown effectiveness in capturing keyphrases, they focus on optimizing phrase extraction and 116 ranking independently, rather than enhancing downstream tasks like long document representation. 117 Our work bridges this gap by integrating unsupervised keyphrase extraction with chunk represen-118 tation to improve long document understanding. We highlight the limitations of existing methods, 119 which fail to preserve key content when reducing document length, and propose a novel approach 120 that groups semantically important phrases to maintain critical information for downstream tasks. 121



Figure 1: The Overall ChuLo Framework proposed in this paper. Each chunk is surrounded by a pink box. $C_1 \dots C_n$ represents the chunk representation.

3 CHULO

135

136

137 138 139

140

141 We introduce a novel chunk representation method, ChuLo, tailored to enhance Transformer-based 142 long document classification by effectively reducing input length while preserving semantic content. 143 Our approach addresses the challenges posed by existing techniques like truncation and standard 144 chunking, which often result in information loss and disrupted contextual dependencies. Specif-145 ically, we first segment the document into non-overlapping, fixed-length chunks to manage long input sequences. Next, we employ unsupervised keyphrase extraction to capture the semantically 146 critical information across the document. By integrating these keyphrases into the chunk representa-147 tions, we strategically assign higher weights to the keyphrase tokens, emphasizing essential content 148 in each chunk. The resulting chunk representations are then used to train a Transformer-based chunk 149 attention module, which ensures that the model retains the global context while focusing on impor-150 tant local information. This methodology not only mitigates the issue of information loss but also 151 enables the model to handle long documents with improved efficiency and performance. The details 152 of each component are in the following subsections, and the overall framework is in Figure 1. 153

154 3.1 DOCUMENT INPUT CHUNKING

To effectively manage long document inputs, we employ a chunking strategy that reduces input
length while preserving all relevant information. Transformer models, despite their success in capturing long-range dependencies through self-attention, suffer from quadratic computational complexity as input length increases (Ivgi et al., 2023; Beltagy et al., 2020). This complexity imposes
limitations on the maximum input length and, consequently, on the amount of text the model can
process at once. Common solutions such as truncating and sparse attention either disregard parts of
the document (Lewis et al., 2020; Park et al., 2022) or limit the receptive field of individual tokens

162 (Beltagy et al., 2020; Zaheer et al., 2020; Brown et al., 2020), leading to information loss. Our 163 approach mitigates these issues by segmenting the document into non-overlapping chunks before 164 feeding them into the model. This strategy enables complete self-attention among chunk, ensuring 165 that all information is retained and enabling the model to process larger portions of the document 166 context. Specifically, we first tokenize the document $D = (t_0, t_1, \ldots, t_{l_D-1})$ and divide it into fixed-length chunks $\mathbb{C}_D = (C_0, C_1, \dots, C_{m-1})$, where l_D is the number of the tokens, each chunk 167 *C* consists of *n* tokens, and $m = \lceil \frac{l_D}{n} \rceil$ is the number of chunks. The incomplete chunks will be padded with the [PAD] tokens. The chunk size *n* is a hyper-parameter controlling the degree of 168 input length reduction. By grouping tokens this way, we maintain the integrity of the input content 170 while alleviating the computational burden associated with processing long sequences. 171

172 173

3.2 SEMANTIC KEY INFORMATION EXTRACTION

174 The fundamental reason for extracting keyphrases from the chunks, as defined in the document 175 chunking step, is to maintain the integrity of the document's semantic content while reducing input 176 length. During chunking, the document is divided into smaller segments, which can inadvertently 177 distribute important semantic information unevenly across chunks or even cause it to be diluted. 178 Simply treating each chunk equally may lead to overlooking critical context essential for accurate 179 document classification and token-level understanding. Identifying and highlighting critical phrases within these chunks ensures that the most relevant information is preserved and emphasized, allow-181 ing the model to focus on the core content even within a limited input space. This compensates for the information fragmentation caused by chunking and guides the Transformer's attention mecha-182 nism to prioritize the most informative parts of the text, enhancing the model's ability to capture the 183 document's overall meaning and relationships. Thus, extracting keyphrases from chunks is crucial for bridging the gap between document segmentation and semantic coherence, ultimately improving 185 the effectiveness of the chunk-based representation for long document understanding. To achieve 186 this, we extract semantically important keyphrases to identify the core content of the entire docu-187 ment. Since document understanding, such as document classification or token classification, in-188 herently involves semantic understanding, it is crucial to highlight the most informative parts of the 189 text to create meaningful chunk representations. By making the extracted keyphrases more salient, 190 we can effectively emphasize the content that contributes most to the document's overall mean-191 ing. Hence, we employ unsupervised keyphrase extraction methods, ensuring our approach remains adaptable across diverse domains without requiring annotated data. Building on the principles of 192 PromptRank (Kong et al., 2023), we adapt and enhance its template-based approach to prioritize 193 keyphrases that are contextually significant across the entire document. Our modified strategy, the 194 Semantic Keyphrase Prioritization (SKP) Algorithm, leverages prompts to assess the importance of 195 each candidate keyphrase, ensuring that semantically crucial information is highlighted for down-196 stream document understanding. The details of this process are provided in Algorithm 1: 197

Algorithm 1 Semantic Keyphrase Prioritization (SKP) Algorithm

200 201 202 203	Input : A tokenized document <i>D</i> , an encoder-decoder pretrained model represented by \mathcal{F}_E and \mathcal{F}_D , a POS tagger \mathcal{F}_{POS} , a regular expression $\mathcal{F}_{REG} = \langle \text{NN}. * \text{JJ} \rangle * \langle \text{NN}. * \rangle$ Parameter : Experimentally determined α, γ Output : Sorted keyphrases set \mathbb{K}_s	 6: Construct the prompt P "The * mainly discusses k_i" and to-kenize, * is the category of the document. 7: for j < m do 8: Calculate the probability p_{ij} of the phrase k_i:
204 205 206 207 208 209 210	 Let S = Ø, K_s = Ø, i = 0, j = 0. Get the candidate phrases set: K = F_{REG}(F_{POS}(D)) = {k₀, k₁,, k_{n-1}} Split D into segments S = {D₀, D₁,, D_{m-1}} to meet the input requirement of F_E Galculate the position penalty r_i = L_c/l_d + γ (l_d)³ where L_c is the first occurrence position of k_i in the docu- ment, l_d is the length of the document 	$p_{ij} = \frac{1}{(l_P)^{\alpha}} \sum_{g=h}^{h+l_k-1} \log p(t_g \mid t_{\leq g}),$ $p(t_g \mid t_{\leq g}) = \mathcal{F}_D(\mathcal{F}_E(D_j), t_{\leq g})$ where l_P is the length of the tokenized P , h is the start index of k_i in the prompt, l_k is the length of k_i , t is the token of the prompt. 9: end for 10: Calculate the final score of k_i : $s_i = r_i \times \sum_{j=0}^{j \leq m} p_{ij}$ 11: end for 12: return $\mathbb{K}_s = Sort(\mathbb{K}, s)$

211

199

212 While PromptRank uses prompts to rank keyphrases across the first segment of the document de-213 termined by its encoder model, our SKP applies this concept at the entire document level to ensure 214 that each chunk can preserve the most informative content for the entire document. After obtaining 215 the sorted phrases set \mathbb{K}_s , we select top-n phrases as the keyphrases of the document, which can be 216 regarded as ranked phrases according to their contextual significance within the entire document. 216 This approach emphasizes keyphrases during chunk representation, making critical semantic infor-217 mation more salient. Consequently, our method bridges the gap between document segmentation 218 and semantic coherence by ensuring that key content is preserved and highlighted within the entire 219 document, despite input length constraints. By integrating these keyphrase extraction techniques, 220 our framework effectively identifies and emphasizes the most informative parts of a long document, resulting in improved chunk-based representations. This ensures that the Transformer retains critical 221 context and relationships, ultimately enhancing its performance in long document classification and 222 token-level understanding tasks. 223

224 225

3.3 CHUNK REPRESENTATION WITH SEMANTIC EMPHASIS

After extracting the semantically significant keyphrases, we construct a chunk representation that preserves and highlights this key information, ensuring that the chunk retains the core semantic content of the document. While chunking helps reduce the input length, it may also result in an uneven distribution of meaningful content across chunks. Thus, it is crucial to re-emphasize the importance of these keyphrases within the chunk to maintain semantic integrity. Our approach dynamically adjusts the representation of each chunk by assigning greater importance to keyphrase tokens, enabling the model to focus on the most relevant content during downstream processing.

To achieve this, we label the tokens corresponding to the extracted keyphrases in the original text as keyphrase tokens T_k , while other tokens are labeled as non-keyphrase tokens T_{nk} . Then, we feed these chunked tokens t into the embedding layer to obtain their embeddings. The chunk embedding e_C is then computed using a weighted average of these token embeddings, as defined in Formula 1:

$$\begin{cases} w_t = \begin{cases} a, t \text{ is } T_k \\ b, t \text{ is } T_{nk} \\ e_c = \frac{\sum w_t * e_t}{\sum w_t} \end{cases}$$
(1)

Here, w_t represents the weight assigned to each token t in the chunk, where a and b are hyperparameters with a > b. e_t denotes the embedding of token t, and e_c is the resulting chunk embedding that captures the weighted importance of keyphrase and non-keyphrase tokens. By assigning a higher weight a to keyphrase tokens, we ensure that the resulting chunk representation emphasizes the most critical information while maintaining a compact input length.

Finally, the chunk embeddings are fed into the Transformer-based model, allowing it to effectively leverage the enhanced chunk representations during long document classification or token-level understanding tasks. This method not only preserves the semantic coherence of the document but also allows the model to retain meaningful context and relationships, ultimately enhancing its performance on long document tasks.

252 253 254

3.4 CHUNK REPRESENTATION TRAINING

255 In this final step, we train a Transformer-based model using our keyphrase-enhanced chunk representations to effectively incorporate the core semantic content of the document. We selected BERT-base 256 according to Table 8. By emphasizing key information in the chunk embeddings, we ensure that the 257 model can focus on the most relevant aspects of the text, thereby improving its ability to handle long 258 document inputs without losing critical context. To achieve this, we leverage a Transformer-based 259 backbone model, which is used to initialize the weights of the chunk attention module, as illustrated 260 in Figure 1. This chunk attention module is designed to capture the intricate contextual relationships 261 among chunks while retaining the influence of keyphrases. By doing so, the module can better un-262 derstand local and global semantic patterns, enabling the model to perform robustly across various 263 long document tasks. The chunk embeddings are processed through the chunk attention module to 264 produce refined contextual representations, which are then fed into a classification head to generate 265 the final predictions. Our framework, ChuLo, is adaptable to any transformer-based architecture, 266 from pretrained to large language models, making it versatile for tasks involving long document 267 understanding. Through integrating keyphrase-enhanced chunk representations, the model achieves superior performance in both document classification and token-level tasks, highlighting the effec-268 tiveness of our approach in leveraging semantic information to tackle the challenges associated with 269 long document processing.

270 4 EXPERIMENTS SET-UP

271 272

We evaluate ChuLo on document and token classification tasks, highlighting our motivation for
including both types. While document classification primarily requires global contextual understanding, token classification tasks test the model's ability to retain and utilize detailed token-level
information within long documents. We compare it with BERT (Kenton & Toutanova, 2019) and
BERT variants (Park et al., 2022), Longformer (Beltagy et al., 2020), ToBERT (Pappagari et al.,
2019), CogLTX (Ding et al., 2020), ChunkBERT (Jaiswal & Milios, 2023), and instructions with
LLMs, GPT4o¹ and Gemini1.5pro². Baselines Details are listed in Appendix A.2.

279 **Datasets:** We conduct experiments using three(3) datasets for long document classification and 280 two(2) for long document token classification. For document classification, we use HP(Kiesel et al., 281 2019), LUN (Rashkin et al., 2017), and Eurlex57k (Chalkidis et al., 2019). These datasets vary in 282 average document length from 707 to 1147 tokens, enabling us to assess performance on documents 283 of different lengths and complexities. Further details on dataset statistics and splits are available in Appendix A.4. 1) HP (Hyperpartisan News Detection): evaluates the classification of news 284 articles based on hyperpartisan argumentation (Kiesel et al., 2019). We use the 'byarticle' version, 285 which contains 238 hyperpartisan and 407 non-hyperpartisan articles. The same train-test split as 286 in (Beltagy et al., 2020) is adopted. 2) LUN Used for unreliable news source classification, this 287 dataset includes 17,250 articles from satire, propaganda, and hoaxes (Rashkin et al., 2017). Our 288 goal is to predict the source type for each article. 3) Eurlex57k A multi-label classification dataset 289 consisting of 47,000 English EU legislative documents with 4,271 EUROVOC concepts. We also 290 include a simulated Inverted-Eurlex57k version, where the header and recitals are moved to the 291 end, compelling the model to read the entire document for key information. For token classification, 292 we use GUM and CoNLL-2012 for Named Entity Recognition (NER) tasks: 1) GUM (Georgetown 293 University Multilayer) is a richly annotated collection of 235 documents across various genres such as academic texts, news, fiction, and interviews (Lin & Zeldes, 2021). GUM's various linguistic 294 styles and structures make it an excellent benchmark for assessing token-level understanding in 295 lengthy documents, ensuring that the model captures complex entity relationships over extended 296 contexts. 2) CoNLL-2012 dataset Originally designed for coreference resolution, this dataset is 297 adapted for NER in our work (Pradhan et al., 2013). We convert the data to a document-level format 298 and select the top-k longest documents in each split, emphasizing the model's ability to understand 299 and process extended text sequences for token classification tasks. 300

Metrics and Implementation: For the HP and LUN datasets, we use **accuracy** as the evaluation 301 metric, while for Eurlex57k, Inverted Eurlex57k, GUM, and CoNLL-2012, we adopt micro F1. 302 These metrics are chosen to maintain consistency with prior work and facilitate direct comparison. 303 Regarding implementation, we provide key details here, with the complete setup in Appendix A.5. 304 We use CrossEntropy loss for training on the Hyperpartisan, LUN, CoNLL and GUM datasets, 305 and Binary CrossEntropy loss for the Eurlex57k and Inverted Eurlex57k datasets. All models are 306 optimized using the AdamW optimizer, and training employs early stopping based on the respective 307 validation metric, with a patience threshold set to 10 epochs. A learning rate search is conducted for 308 each experiment to ensure optimal model performance for comparison. Top-n value is set to 15^3 .

309 310

5 RESULTS

311 312 313

5.1 DOCUMENT CLASSIFICATION PERFORMANCE

We evaluate the effectiveness of our ChuLo by comparing it with fine-tuned PLMs and previously published baselines (Park et al., 2022; Jaiswal & Milios, 2023) on several benchmark datasets: HP, LUN, EURLEX57K, and Inverted EURLEX57K. The comparative results are summarized in Table 1, with input configurations provided in Table 2 and detailed descriptions available in Appendix A.3. Our method demonstrates clear superiority on three of the four datasets, achieving a significant improvement of 6.43% accuracy on the LUN dataset compared to the second-best model, BERT. This marked improvement presents ChuLo's ability to capture comprehensive document context through

²https://deepmind.google/technologies/gemini/pro/

³²¹

³²² 323

¹https://openai.com/index/hello-gpt-40/

³We tested with different n values, but 15 was generally better in most datasets

324	Model	HP	LUN	EURLEX57K	I-EURLEX57K
325	BERT (Kenton & Toutanova, 2019)	0.9200	0.5797	0.7309	0.7053
000	ToBERT (Pappagari et al., 2019)	0.8954	0.3697	0.6757	0.6731
326	CogLTX (Ding et al., 2020)	0.9477	-	0.7013	0.7080
327	Longformer (Beltagy et al., 2020)	0.9569	0.5552	0.5453	0.5647
	BERT+TextRank (Park et al., 2022)	0.9115	0.4880	0.7287	0.7130
328	BERT+Random (Park et al., 2022)	0.8923	0.3015	0.7322	0.7147
329	ChunkBERT (Jaiswal & Milios, 2023)	0.9300	-	0.6494	0.6294
220	Ours	0.9538	0.6440	0.7332	0.7244

Table 1: Document classification Result. Following previous work, we use accuracy for HP and
LUN, and micro F1 for other datasets. Results for LUN are obtained by our own experiment based on
provided baseline codes and methods, while baseline results for the other datasets are from previous
work(Park et al., 2022; Jaiswal & Milios, 2023). The best performance for each dataset is bolded
while the second best is underscored, and we can see that our final model, a BERT-based backbone,
generally outperforms other baselines across all datasets by achieving the best or second-best.

337 its keyphrase-based chunk representation, despite using only 512 input embeddings. The results 338 suggest that our method effectively mitigates the limitations of traditional truncation and chunk-339 ing strategies by preserving critical semantic information, which contributes to higher classification 340 accuracy. On the EURLEX57K and Inverted EURLEX57K datasets, ChuLo achieves consistent per-341 formance gains over baselines, further validating its capability to handle long documents efficiently. 342 In these datasets, which have hierarchical labels and require understanding complex semantic struc-343 tures, our model benefits from enhanced chunk representations that emphasize key content. This 344 allows ChuLo to capture document semantics better, even when compared to models that can process larger input lengths. While our model delivers competitive results on the HP dataset, it trails 345 behind Longformer by a slight margin of 0.0031 in accuracy. This marginal difference corresponds 346 to only one additional correctly classified instance out of a total of 65 test samples. 347

Model	The Usage of Input
BERT (Kenton & Toutanova, 2019)	F-512 tokens
ToBERT (Pappagari et al., 2019)	All
CogLTX (Ding et al., 2020)	S-512 tokens
Longformer (Beltagy et al., 2020)	F-4096 tokens
BERT+TextRank (Park et al., 2022)	F-512 + S-512 tokens
BERT+Random (Park et al., 2022)	F-512 + S-512 tokens
ChunkBERT (Jaiswal & Milios, 2023)	All
Ours	All (512*Chunk Size)

Table 2: The usage of the input content in the experiments. "F-512" and "F-4096" means the first 512 tokens and the first 4096 tokens, "S-512" means the selected 512 tokens.

Interestingly, for the other datasets, Longformer underperforms compared to models like BERT variants or CogLTX, which use the first 512 tokens and focus on selecting key sentences. This observation indicates that unfiltered additional content can introduce noise, negatively impacting classification accuracy. In contrast, ChuLo expands the input content and strategically emphasizes key semantic elements during chunk representation. This approach mitigates noise interference, ensuring that only the most relevant information is retained and highlighted. ChuLo achieves superior perfor-

mance by balancing content and semantic emphasis, outperforming other models. Overall, the re sults confirm that ChuLo consistently outperforms standard PLM baselines and existing chunking
 methods in long document classification tasks. Its ability to retain and emphasize key semantic con tent, while efficiently handling long inputs, makes it a robust solution for various document classifi cation challenges. The subsequent sections delve deeper into the impact of our chunk representation
 strategy and discuss its contributions to improving document classification performance.

366

348

349

350

351

352

353

354

355

356

357

358

359

367

5.2 DOCUMENT CLASSIFICATION PERFORMANCE IN LONGER DOCUMENTS

368 To further validate the robustness of our model, we evaluate its classification performance across 369 various document length ranges, with a particular focus on longer documents. For this analysis, 370 we consider the documents with more than 1024 tokens and more than 2048 tokens in the test set. 371 To provide a fair comparison, we use Longformer based on its original code and hyperparameters 372 as described in Table 1 and off-the-shelf LLMs, GPT40 and Gemini1.5 pro. As shown in Table 3, 373 our model consistently outperforms others on longer documents in the LUN dataset. Specifically, 374 for documents exceeding 2,048 tokens, ChuLo maintains a higher accuracy compared to all base-375 lines, demonstrating its capacity to handle lengthy inputs effectively. This performance gain can be attributed to our chunk representation's emphasis on keyphrases, which preserves crucial semantic 376 content even when document length increases. On the HP dataset, ChuLo and Longformer achieve 377 perfect accuracy (1.0) for documents longer than 2,048 tokens. However, for shorter documents (more than 1,024 tokens), ChuLo surpasses Longformer. This improvement is likely due to our chunk representation strategy, which selectively highlights key content rather than averaging information across the entire document. As a result, ChuLo maintains high semantic fidelity, leading to better overall performance even with condensed text inputs.

382 383 384

385

386

387

388

389

390

391

392

404 405

406

407

408

409

410

411

412

413

414

415

LUN	All(2250)	1024(243)	2048(49)
Longformer	0.5552	0.4062	0.5306
GPT40	-	-	0.7143
Gemini1.5pro	-	-	0.6531
Ours	0.6741	0.5911	0.7959
	(a) LUN da	taset	
HP	(a) LUN da All(65)	taset 1024(28)	2048(9)
HP Longformer	(a) LUN da All(65) 0.9538	taset 1024(28) 0.8929	2048(9) 1.000
HP Longformer GPT40	(a) LUN da All(65) 0.9538 -	taset 1024(28) 0.8929	2048(9) 1.000 0.8889
HP Longformer GPT40 Gemini1.5pro	(a) LUN da All(65) 0.9538 - -	taset 1024(28) 0.8929 - -	2048(9) 1.000 0.8889 0.7778

(b) HP dataset

Table 3: Document classification results for com-394 parison on documents of different lengths: all 395 documents in the test set, the subset of documents 396 longer than 1024 tokens, and longer than 2048 to-397 kens respectively. Values in brackets indicate the 398 number of documents for each specific document 399 set. The best performance (Accuracy) for each 400 document set is bolded. 401

We also benchmarked against newly released LLMs, GPT-40 and Gemini 1.5 Pro, using longer document inputs for both the LUN and HP datasets. On LUN, GPT-40 achieved an accuracy of 0.7143 and Gemini 1.5 Pro scored 0.6531, both surpassing Longformer. However, ChuLo achieved the highest accuracy of 0.7959, showcasing its superiority in handling long documents with diverse content. On the HP dataset, GPT-40 (0.8889) and Gemini 1.5 Pro (0.7778) performed worse than Longformer and ChuLo, both of which achieved a perfect accuracy of 1.0 on the longer documents. This highlights ChuLo's robustness and consistency in classifying documents with varying length, even compared to advanced language models. The prompt and response samples are in Appendix A.6 and A.7. Overall, these results demonstrate that ChuLo not only outperforms standard PLM baselines and chunking methods on long documents but also remains competi-

tive against the latest large language models. By prioritizing key semantic elements and effectively
 managing document length, ChuLo maintains stable performance across varying input lengths.

5.3 TOKEN CLASSIFICATION PERFORMANCE

Model	CoNLL	GUM
Longformer (4096)	0.5560	0.9427
BigBird (4096)	0.5553	0.9418
GPT40	0.2290	0.3231
Gemini1.5	0.3036	0.3262
Ours (All)	0.9334	0.9555

Table 4: Results on token classification tasks. The best performance for each dataset is bolded, and our model achieves the best on both datasets.

To further demonstrate the effectiveness of our chunk representation method, we evaluated it on a token-level classification task—specifically, Named Entity Recognition (NER) using long documents. We compared our model against two state-of-the-art longdocument pre-trained models, Longformer (Beltagy et al., 2020) and BigBird (Zaheer et al., 2020), as well as newly released large language models, GPT-40 and Gemini 1.5 Pro.

416 As shown in Table 4, our model consistently outperforms both Longformer and BigBird and LLM 417 models on the NER tasks, particularly on the CoNLL, where document lengths often exceed the input limitations of these baseline models. To leverage the broader context captured by our chunk 418 representation, we integrate a BERT-decoder module that utilizes the enhanced chunk embeddings 419 to predict token labels more accurately. This configuration allows the model to maintain a global 420 understanding of the document while focusing on the local dependencies necessary for precise 421 token classification. The superior performance of our ChuLo on the CoNLL demonstrates the 422 significance of retaining global contextual information when dealing with long documents. All 423 baselines struggle with these longer inputs due to their limited capacity for handling extensive 424 sequences. In contrast, our method's ability to encode the entire document's context through 425 keyphrase-based chunk representations enables it to achieve higher accuracy in recognizing and 426 classifying named entities. This is particularly evident in cases where long-distance dependencies 427 and contextual nuances play a critical role in determining the correct labels. Overall, the results 428 indicate that our model's chunk representation not only enhances performance on document-level 429 classification tasks but also proves highly effective for token-level tasks such as NER. By retaining global context while emphasizing key semantic content, our method enables more accurate token 430 classification, validating its application in downstream tasks that require detailed and comprehensive 431 understanding of long document tokens.

432 5.4 TOKEN CLASSIFICATION PERFORMANCE IN LONGER DOCUMENTS

We further analyze the NER performance across different document length ranges. As presented in Table 5a and Table 5b, we report the number of documents exceeding specific length thresholds and their corresponding performance metrics. On the CoNLL, as document lengths exceed the maximum input capacities of Longformer and BigBird, both models exhibit significant performance drops to 31.56% and 31.45%, respectively. In contrast, our model experiences a minimal decrease of 1.28%, showcasing its resilience and effectiveness in handling long sequences. For the GUM, where all document lengths are within the acceptable range for these models, performance remains stable across all models, with our approach consistently achieving the best results.

CoNLL	Entire dataset (20)	Longer than 2048 (17)	Longer than 4096(6)	Longer than 8192 (2
Longformer	0.5560	0.5268	0.3156	0.3116
BigBird	0.5553	0.5261	0.3145	0.3106
GPT40	0.2290	0.2217	0.1252	0.0282
Gemini 1.5	0.3036	0.2633	0.1652	0.0584
Ours	0.9334	0.9325	0.9287	0.9206
GUM	Entire dataset (2	6) - Longer than 512 L	onger than 1000(8) I	(1 1040 (6)
GUM	Entre uataset (2	au - Louger man 514 L		ABBON THAN IN/1/1/K
Longform	ar 0	0427	0.9427	$\frac{\text{onger than 1042 (6)}}{0.0439}$
Longform	er <u>0</u>	<u>1.9427</u> 0.9418	<u>0.9427</u> 0.9417	$\frac{0.9439}{0.9426}$
Longform BigBird GPT40	er <u>0</u>	<u>.9427</u> .9418	<u>0.9427</u> 0.9417 0.3018	$\frac{0.9439}{0.9426}$
Longform BigBird GPT40 Comini 1	er <u>0</u> 0	<u>.9427</u> .9418 .3231	0.9427 0.9417 0.3018 0.3002	$\frac{0.9439}{0.9426}$ 0.2808 0.2215

(b) Results on GUM dataset.

Table 5: NER results for comparison on documents of different lengths. Values of brackets indicate the # of documents for each document. The best performance (Micro F1) is bolded and the second best is underscored, and our model consistently outperforms all the baselines for each document set.



Figure 2: Comparison of performance in different length ranges for CoNLL and GUM datasets. Values of brackets includes the min and max length of each dataset

Figures 2a and 2b visualize the performance breakdown across varying length ranges. For the CoNLL, our model maintains high performance in all length intervals, while Longformer and Big-Bird exhibit comparable performance within the [1k-2k) range but degrade significantly for longer texts, even for documents that do not exceed their maximum input length. This discrepancy suggests that the uneven distribution of document lengths in their pretraining corpora may lead to inconsistent performance on longer sequences. In contrast, our model's ability to compress the entire document into 512-length chunks for the decoder enables it to leverage complete contextual information, re-sulting in better stability and accuracy even on longer documents. For the GUM, where document lengths are shorter (up to 1.3k tokens), our model consistently outperforms Longformer and Big-Bird in all intervals. The stable performance of all models on GUM aligns with the results on CoNLL, further confirming that our approach's chunk representation is particularly effective when documents reach lengths that exceed the standard input capacities of the baselines. These results underscore the effectiveness of our chunk representation, which emphasizes keyphrase informa-tion, for coarse-grained document classification and fine-grained token-level classification tasks like NER. The ability to maintain performance across varying document lengths highlights the impor-tance of incorporating global contextual information in NER tasks—a largely underexplored aspect. Additionally, off-the-shelf LLMs such as GPT-40 and Gemini 1.5 Pro show suboptimal performance on NER tasks without fine-tuning, and their performance deteriorates further as document length in-creases. This indicates that, despite their advancements, LLMs still require substantial optimization for token classification tasks and effective long document understanding.

486 5.5 ABLATION STUDIES 487

512 513

519

530

531 532

533

488	Keyphrase method	HP	LUN	Sentence Embedding H	IP	LUN	Backbone	HP	LUN
489	Average	0.9538	0.5951	w/o sentence emb. 0.	.9538	0.6440	BERT (Ours)	0.9538	0.6440
400	YAKE	0.8769	0.5951	sentence emb. 0.	.9076	0.5537	RoBERTa	0.8615	0.5906
490	PromptRank	0.9538	0.6440	I			Longformer	0.8923	0.5600
491				Table 7: Effect of	conto	n 00 0 m			

Table 6: Effect of keyphrase ex-492 traction methods; Average: Av-493 erage Chunk Representations 494

Table 7: Effect of sentence embedding, adding the sentencelevel information to the chunk representations.

Table 8: Effect of different backbone models for the chunk attention.

495 We performed a few ablation studies on the HP and LUN to assess the impact of various components 496 within our model. First, we analyzed the effectiveness of different keyphrase extraction methods 497 and the effect of using average chunk representations. As shown in Table 6, the PromptRank-498 based method yields the highest performance across both datasets, outperforming alternatives like 499 YAKE-based. This improvement can be attributed to PromptRank's ability to extract higher-quality 500 keyphrases by considering semantic relationships within the document, whereas YAKE relies primarily on statistical features such as phrase frequency, resulting in less semantically rich keyphrases. 501 Then, we explored the effect of incorporating sentence embeddings into the chunk representations to 502 introduce global sentence-level context. Surprisingly, as shown in Table 7, the results indicate a per-503 formance drop when sentence embeddings are included. We hypothesize that adding sentence-level 504 information at the initial representation stage may cause chunk embeddings within the same sentence 505 to become too similar, hindering the model's ability to learn distinctive patterns and reducing overall 506 classification performance. We also evaluated the performance of different backbone models for the 507 chunk attention module while keeping the keyphrase extraction and chunk representation settings 508 consistent. Table 8 shows that BERT outperforms Longformer as the backbone. This result suggests 509 that, after document chunking, the input sequences become relatively short, making it difficult for 510 Longformer to leverage its long-range attention capabilities fully. Consequently, Longformer may suffer underutilization during training, resulting in suboptimal performance compared to BERT. 511

5.6 QUALITATIVE ANALYSIS

514 We performed a qualitative analysis by visualizing a sample document from the GUM, comparing 515 the outputs of Longformer, GPT-40, Gemini 1.5 Pro, and our ChuLo. ChuLo captures the context 516 and semantic patterns of the document, providing accurate predictions, whereas the other models 517 struggle to maintain coherence and consistency. We have more examples in Appendix A.7. 518

Case 5, length: 895 - Document and Prompt	Case 5 – Our Model
"In the task of Named Entity Recognition, the B-, I-, and O- prefixes are commonly used to annotate slot	[B-event', 'I-event', 'I-event', 'I-event', 'I-event', 'I-event', 'B-abstract', 'I-abstract',]
types, indicating the boundaries and types of slots. These labels typically represent: B- (Begin):	
Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior	
of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not	
part of any slot. For instance, in a sentence where we want to label a """date" slot, words	Case 5 - GPT-4o Response
followed by consecutive words carrying date information tagged as """"I-date""" (indicating the	[I-place', 'B-substance', 'I-person', 'B-animal', 'I-object', 'B-time', 'I-organization', 'I-place',
continuation of the date slot), while words not containing date information would be tagged as	
""""O"""" (indicating they are outside any slot). Definition: In this task, you are given a conversation,	
where the words spoken by a person are shown as a list. Your job is to classify the words in the	
following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B-	
place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal',	Case 5 – Gemini1.5pro Response
"B-organization", "I-event", "B-person", "B-event", "I-plant", "I-organization", "O', "I-object", "B-animal". Only	[B-person', '0', '0', '0', '0', '0', '0', '0',]
[Sam' 'bas' 'been' '' 'bas' 'taken' such' 'an' interest' 'in' 'this' 'retirement' 'hit' 'that' 'it' ''	<u></u>
'it'. 'really', 'surprises', 'me', '.', 'Well', ""she's"", 'she', ""'s"", 'begun', 'to', 'listen', '.', 'Yes', 'she', 'has', '.',	
'You', 'know', '.', 'She', 'has', '.', 'Uh', ',', 'she', 'used', 'to', 'go', 'over', 'and', 'read', 'a', 'book', 'or',	
'something', '.', 'Yeah', ',', 'or', 'turn', 'a', 'deaf', 'ear', '.', 'That', 'was', 'for', 'sure', '.', 'But', 'some', 'basil',	Case 5 – Longformer Output
',', 'yes', ',', 'Yeah', ',', 'I', ""don't"", 'do', ""n't"", 'have', 'any', 'this', 'year', ',', 'I', 'forgot',]."	
	[[<mark>'O', </mark> I-event, 'I-event', 'I-event', 'I-event', 'I-event', 'O', 'O', 'O', 'O', 'O', ']

Figure 3: Prompt and output for a sample document of length 895 in **GUM** dataset for NER task, where correct predictions are highlighted in green and wrong predictions are highlighted in red.

6 CONCLUSION

534 We introduced ChuLo, a novel chunk representation method that enhances the performance of Transformer-based models on long document classification and token-level tasks. By utilizing un-536 supervised keyphrase extraction, ChuLo effectively reduces input length while preserving critical 537 information, addressing the limitations of truncation and sparse attention. Extensive experiments demonstrate that ChuLo outperforms existing methods by maintaining both global context and high 538 accuracy, even for lengthy inputs. Our results highlight the effectiveness of ChuLo as a robust solution for long document understanding, enabling processing of complex texts in NLP applications.

540 REFERENCES

547

554

565

566

567

573

- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
 arXiv preprint arXiv:2004.05150, 2020.
- Kamil Bennani-Smires, Claudiu Musat, Andreea Hossmann, Michael Baeriswyl, and Martin Jaggi.
 Simple unsupervised keyphrase extraction using sentence embeddings. In *Proceedings of the* 22nd Conference on Computational Natural Language Learning, pp. 221–229, 2018.
- Medha Bhattacharya, Swati Bhat, Sirshasree Tripathy, Anvita Bansal, and Monika Choudhary. Improving biomedical named entity recognition through transfer learning and asymmetric tritraining. *Procedia Computer Science*, 218:2723–2733, 2023.
- Adrien Bougouin, Florian Boudin, and Béatrice Daille. Topicrank: Graph-based topic ranking for
 keyphrase extraction. In *International joint conference on natural language processing (IJCNLP)*,
 pp. 543–551, 2013.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 561
 562
 563
 564
 564
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 765
 764
 766
 766
 766
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 767
 - Can Çetindağ, Berkay Yazıcıoğlu, and Aykut Koç. Named-entity recognition in turkish legal texts. *Natural Language Engineering*, 29(3):615–642, 2023.
- Ilias Chalkidis, Emmanouil Fergadiotis, Prodromos Malakasiotis, and Ion Androutsopoulos. Large scale multi-label text classification on eu legislation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6314–6322, 2019.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Krzysztof Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane,
 Tamas Sarlos, Peter Hawkins, Jared Davis, David Belanger, Lucy Colwell, et al. Masked language modeling for proteins via linearly scalable long-context transformers. *arXiv preprint arXiv:2006.03555*, 2020.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- John Dagdelen, Alexander Dunn, Sanghoon Lee, Nicholas Walker, Andrew S Rosen, Gerbrand Ceder, Kristin A Persson, and Anubhav Jain. Structured information extraction from scientific text with large language models. *Nature Communications*, 15(1):1418, 2024.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov.
 Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the* 57th Annual Meeting of the Association for Computational Linguistics, pp. 2978–2988, 2019.
- Ming Ding, Chang Zhou, Hongxia Yang, and Jie Tang. Cogltx: Applying bert to long texts. Advances in Neural Information Processing Systems, 33:12792–12804, 2020.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
 arXiv preprint arXiv:2407.21783, 2024.

- Samhaa R El-Beltagy and Ahmed Rafea. Kp-miner: A keyphrase extraction system for english and arabic documents. *Information systems*, 34(1):132–144, 2009.
- Corina Florescu and Cornelia Caragea. Positionrank: An unsupervised approach to keyphrase extraction from scholarly documents. In *Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: long papers)*, pp. 1105–1115, 2017.
- Kazi Saidul Hasan and Vincent Ng. Automatic keyphrase extraction: A survey of the state of the art. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1262–1273, Baltimore, Maryland, June 2014. Association for Computational Linguistics. doi: 10.3115/v1/P14-1119. URL https://aclanthology.org/P14-1119.
- Yan Hu, Qingyu Chen, Jingcheng Du, Xueqing Peng, Vipina Kuttichi Keloth, Xu Zuo, Yujia Zhou,
 Zehan Li, Xiaoqian Jiang, Zhiyong Lu, et al. Improving large language models for clinical named
 entity recognition via prompt engineering. *Journal of the American Medical Informatics Association*, pp. ocad259, 2024.
- DeLesley Hutchins, Imanol Schlag, Yuhuai Wu, Ethan Dyer, and Behnam Neyshabur. Block recurrent transformers. *Advances in Neural Information Processing Systems*, 35:33248–33261, 2022.
- Maor Ivgi, Uri Shaham, and Jonathan Berant. Efficient long-text understanding with short-text models. *Transactions of the Association for Computational Linguistics*, 11:284–299, 2023. doi: 10.1162/tacl_a_00547. URL https://aclanthology.org/2023.tacl-1.17.
- Aman Jaiswal and Evangelos Milios. Breaking the token barrier: Chunking and convolution for
 efficient long text classification with bert. *arXiv preprint arXiv:2310.20558*, 2023.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pp. 4171–4186, 2019.
- Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David Corney, Benno Stein, and Martin Potthast. Semeval-2019 task 4: Hyperpartisan news detection. In
 Proceedings of the 13th International Workshop on Semantic Evaluation, pp. 829–839, 2019.
- Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong Qin, Ruiqi Sun, and Xiaoyan Bai.
 PromptRank: Unsupervised keyphrase extraction using prompt. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9788– 9801, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/ 2023.acl-long.545. URL https://aclanthology.org/2023.acl-long.545.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer
 Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pretraining for natural language generation, translation, and comprehension. In *Proceedings of the*58th Annual Meeting of the Association for Computational Linguistics, pp. 7871–7880, 2020.
- Miao Li, Eduard Hovy, and Jey Han Lau. Towards summarizing multiple documents with hierarchical relationships. *arXiv preprint arXiv:2305.01498*, 2023a.
- Kianming Li, Zongxi Li, Xiaotian Luo, Haoran Xie, Xing Lee, Yingbin Zhao, Fu Lee Wang, and Qing Li. Recurrent attention networks for long-text modeling. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 3006–3019, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.188. URL https://aclanthology.org/2023. findings-acl.188.
- Jessica Lin and Amir Zeldes. WikiGUM: Exhaustive entity linking for wikification in 12 genres. In
 Proceedings of The Joint 15th Linguistic Annotation Workshop (LAW) and 3rd Designing Mean- ing Representations (DMR) Workshop (LAW-DMR 2021), pp. 170–175, Punta Cana, Dominican
 Republic, 2021. URL https://aclanthology.org/2021.law-1.18.

657

658

659

660

661

668

669

670

673

689

691

- 648 Zhiyuan Liu, Peng Li, Yabin Zheng, and Maosong Sun. Clustering to find exemplar terms for 649 keyphrase extraction. In Proceedings of the 2009 conference on empirical methods in natural 650 language processing, pp. 257-266, 2009. 651
- Davlatyor Mengliev, Vladimir Barakhnin, Nilufar Abdurakhmonova, and Mukhriddin Eshkulov. 652 Developing named entity recognition algorithms for uzbek: Dataset insights and implementation. 653 Data in Brief, 54:110413, 2024. 654
- 655 Rada Mihalcea and Paul Tarau. Textrank: Bringing order into text. In Proceedings of the 2004 656 conference on empirical methods in natural language processing, pp. 404–411, 2004.
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35: 27730-27744, 2022.
- Raghavendra Pappagari, Piotr Zelasko, Jesús Villalba, Yishay Carmiel, and Najim Dehak. Hierar-662 chical transformers for long document classification. In 2019 IEEE automatic speech recognition 663 and understanding workshop (ASRU), pp. 838-844. IEEE, 2019. 664
- 665 Hyunji Park, Yogarshi Vyas, and Kashif Shah. Efficient classification of long documents using 666 transformers. In Proceedings of the 60th Annual Meeting of the Association for Computational 667 Linguistics (Volume 2: Short Papers), pp. 702–709, 2022.
 - Yeon-Ji Park, Min-a Lee, Geun-Je Yang, Soo Jun Park, and Chae-Bong Sohn. Web interface of ner and re with bert for biomedical text mining. Applied Sciences, 13(8):5163, 2023.
- 671 H Peng, N Pappas, D Yogatama, R Schwartz, N Smith, and L Kong. Random feature attention. In 672 International Conference on Learning Representations (ICLR 2021), 2021.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga 674 Uryupina, Yuchen Zhang, and Zhi Zhong. Towards robust linguistic analysis using OntoNotes. 675 In Proceedings of the Seventeenth Conference on Computational Natural Language Learning, 676 pp. 143-152, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL 677 https://aclanthology.org/W13-3516. 678
- Alec Radford. Improving language understanding by generative pre-training. 2018. 679
- 680 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language 681 models are unsupervised multitask learners. OpenAI blog, 1(8):9, 2019. 682
- Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. Truth of varying 683 shades: Analyzing language in fake news and political fact-checking. In Proceedings of the 2017 684 conference on empirical methods in natural language processing, pp. 2931–2937, 2017. 685
- 686 Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse 687 attention with routing transformers. Transactions of the Association for Computational Linguis-688 tics, 9:53-68, 2021.
- Yi Sun, Hangping Qiu, Yu Zheng, Zhongwei Wang, and Chaoran Zhang. Sifrank: a new baseline 690 for unsupervised keyphrase extraction based on pre-trained language model. *IEEE Access*, 8: 10896-10906, 2020. 692
- 693 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 694 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023a. 695
- 696 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-697 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-698 tion and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023b. 699
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 700 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural informa-701 tion processing systems, 30, 2017.

Xiaojun Wan and Jianguo Xiao. Single document keyphrase extraction using neighborhood knowl- edge. 2008.
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. arXiv preprint
arXiv:2304.10428, 2023.
Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention
with linear complexity. arXiv preprint arXiv:2006.04768, 2020.
Jason Wai Maartan Rosma Vincant V Zhao, Kalvin Guu, Adams Wai Vu, Brian Lester, Nan Du
Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. <i>arXiv preprint arXiv:2109.01652</i> , 2021.
Miguel Won, Brune Marting and Filing Deimunde, Automatic extraction of relevant keynhroses for
the study of issue competition. In <i>International Conference on Computational Linguistics and Intelligent Text Processing</i> , pp. 648–669. Springer, 2019.
Jianfei Yu, Ziyan Li, Jieming Wang, and Rui Xia, Grounded multimodal named entity recognition
on social media. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 9141–9154, 2023.
Yang Yu and Vincent Ng. Wikirank: Improving keyphrase extraction based on background knowl-
edge. arXiv preprint arXiv:1803.09000, 2018.
Man-il Zahara Cum Cumanada Kuman Asiana Dubay Jashua Airalia Chais Alberti Cartiana
Ontanon Philip Pham Anirudh Ravula Oifan Wang Li Yang et al. Rig bird: Transformers for
longer sequences. Advances in neural information processing systems, 33:17283–17297, 2020.
Xin Zhang, Jingling Yuan, Lin Li, and Jianguan Liu, Reducing the bias of visual objects in multi-
modal named entity recognition. In Proceedings of the Sixteenth ACM international conference
on web search and data mining, pp. 958–966, 2023.
Yusen Zhang, Ansong Ni, Ziming Mao, Chen Henry Wu, Chenguang Zhu, Budhaditya Deb, Ahmed
Awadallah, Dragomir Radev, and Rui Zhang. Summn: A multi-stage summarization framework
for long input dialogues and documents: A multi-stage summarization framework for long input
Computational Linguistics (Volume 1: Long Papers), pp. 1592–1604, 2022.
Jing Zhao, Junwei Bao, Yifan Wang, Yongwei Zhou, Youzheng Wu, Xiaodong He, and Bowen
Zhou. Ror: Read-over-read for long document machine reading comprehension. In <i>Findings of</i>
the Association for Computational Linguistics: EMNLP 2021, pp. 1862–1872, 2021.

756 A APPENDIX

758 A.1 RELATED WORKS 759

760 As shown in Table 9, most of the previous works addressing the problem of processing long docu-761 ments cannot fully utilize all the content. Those methods either reduce input length via truncation 762 or focus on local context learning to improve efficiency by applying sparse attention, approximated attention or RNN integration. Such approaches will lead to a certain level of information loss, unlike 763 our chunking approach which can take all the content into consideration. Hierarchical Transformer 764 (Pappagari et al., 2019) splits documents into non-overlapping chunks and computes chunk repre-765 sentations. RoR (Zhao et al., 2021) generates regional answers from chunks, which are combined 766 for the final answer. However, neither considers the entire document context when chunking. In 767 addition, previous works applying the chunking method for processing long document context only 768 focus on a single task, either document classification or token classification, while our framework 769 can be applied to both tasks to guarantee both document-level and token-level understanding. 770

771	Model	Year	Task	Lengthy Document Solution	Core Architecture
772	Efficient Classification (Park et al., 2022)	2022	D	Truncating	Transformer
770	Hierarchical transformer (Pappagari et al., 2019)	2019	D	Chunking (Partial, Phrase)	Transformer
113	RoR (Zhao et al., 2021)	2021	Т	Chunking (Partial, Voting)	Transformer
774	Longformer (Beltagy et al., 2020)	2020	D, T	Sparse Attention	Transformer
775	BigBird (Zaheer et al., 2020)	2020	D, T	Sparse Attention	Transformer
(1)	Routing Transformer (Roy et al., 2021)	2021	D, T, G	Sparse Attention	Transformer
776	Macformer (Peng et al., 2021)	2021	D, T	Approximated Attention	Transformer
	Linformer (Wang et al., 2020)	2020	D, T, G	Approximated Attention	Transformer
///	Performer (Choromanski et al., 2020)	2020	D, T, G	Approximated Attention	Transformer
778	Transformer-xl (Dai et al., 2019)	2019	G	RNN Integration	Transformer
770	Block-Recurrent Transformer (Hutchins et al., 2022)	2022	G	RNN Integration	Transformer
//9	RAN (Li et al., 2023b)	2023	D, T	RNN Integration	Attention
780	(Çetindağ et al., 2023)	2023	Т	N/A	LSTM
704	(Mengliev et al., 2024)	2024	Т	N/A	Neural Network
781	(Park et al., 2023)	2023	Т	N/A	Transformer
782	(Bhattacharya et al., 2023)	2023	Т	N/A	LSTM
700	Gpt-NER (Wang et al., 2023)	2023	Т	N/A	Transformer
783	(Dagdelen et al., 2024)	2024	Т	N/A	Transformer
784	(Hu et al., 2024)	2024	Т	N/A	Transformer
705	(Yu et al., 2023)	2023	Т	N/A	Transformer
(00	(Zhang et al., 2023)	2023	Т	N/A	Transformer
786	Ours	2024	D, T	Chunking (Entire)	Transformer

Table 9: Summary of Related Works. D, T, G represent tasks of document classification, token classification, and text generation, respectively.

A.2 BASELINES

787

788

789 790 791

792

799

800

801

802

804

805

We use BERT (Kenton & Toutanova, 2019) as our backbone model, comparing it with ToBERT (Pappagari et al., 2019), CogLTX (Ding et al., 2020), Longformer (Beltagy et al., 2020), various BERT variants (Park et al., 2022) and ChunkBERT (Jaiswal & Milios, 2023) for the document classification task. For the NER task, we compare against Longformer, BigBird (Zaheer et al., 2020), and two large language models, GPT40 and Gemini1.5pro. Below are brief descriptions of the baseline models:

- **BERT**: A transformer model pre-trained on masked language modeling (MLM) and nextsentence prediction (NSP). We fine-tuned the BERT-base variant on each dataset.
 - **ToBERT**: A BERT variant designed for long document classification, utilizing an additional transformer layer to learn inter-chunk relationships.
- **CogLTX**: A framework for applying BERT to long documents by training a key sentence identification model to assist in document understanding
- **Longformer**: Optimized for long sequences using sparse attention, combining dilated sliding window and global attention patterns
- **BigBird**: Utilizes block sparse attention, integrating sliding window, global, and random attention patterns across token blocks.

810	Model	Input	
811	BERT (Kenton & Toutanova, 2019)	The first 512 tokens	
812	ToBERT (Pappagari et al., 2019)	Segmented all input tokens	
813	Longformer (Beltagy et al., 2020)	The first 4096 tokens	
91/	BigBird (Zaheer et al., 2020)	The first 4096 tokens	
014	BERT+TextRank (Park et al., 2022) BERT+Random (Park et al., 2022)	The first 512 tokens with the selected 512 tokens The first 512 tokens with the selected 512 tokens	
815	ChunkBERT (Jaiswal & Milios, 2023)	The first 4096 tokens	
816	GPT40	All input tokens with instruction	
817	Gemini1.5pro	All input tokens with instruction	
818	Table 10: The inr	out of the baseline models	
819			
820			
821	REPT_TaytBank and REPT_Pand	lom: Proposed to select other tokens ran	domly or with
822	the help of TextRank (Mihalcea & Ta	rau 2004) to feed into the BERT model	as the supple-
823	mentation for long document underst	tanding	as the supple-
824	mentation for long document underst	anding.	
825	 ChunkBERT: A BERT variant fo 	r long document classification that pr	rocesses self-
826	attention within document chunks a	nd adds a TextCNN module for classif	fication using
827	the chunk representation.		
828	• CPT-40 : A transformer-based multi	-modal large language model develope	d by OpenAI
829	which leverages large-scale pretraini	ng data to process diverse language task	s via instruc-
830	tion prompts	ing data to process arrense language task	
831			
832	• Gemini 1.5 Pro: an advanced multi	i-modal AI model from Google, leverage	ging a Sparse
833	Mixture-of-Experts (MoE) Transform	mer architecture, with a context windo	w of up to 2
924	million tokens. This architecture allo	ws for the efficient handling of long doc	uments.
005			
836	A.3 BASELINE INPUT		
837	We calcuted these boseling models because the	an nonnacht divinger mathada far measag	ing long dog
838	uments As summarized in Table 10 REPT	truncates the input to 512 tokens. Lo	ngformer and
839	BigBird utilize sparse attention mechanisms	allowing them to process up to 4006 toke	ingiornier and
840	serving computational resources ToBERT di	vides the input into 200-token chunks	feeds them to
841	BFRT for chunk representations and uses a	transformer layer for downstream tasks	However it
842	cannot capture dependencies across the entire i	input sequence. CogI TX selects key sent	ences to form
843	a 512-token input limiting its input size to t	hat constraint BERT variants like BE	RT+TextRank
844	and BERT+Random select up to 512 tokens	using TextRank or random selection	They concate-
044	nate the [CLS] representation of the initial 5	12 tokens with the selected tokens, cre	ating an aug-
045	mented input for a fully connected classification	on layer, with a maximum input length of	1024 tokens
840	ChunkBERT splits the input into chunks, con	nputes self-attention, and feeds chunk re	presentations
847	into a TextCNN module for classification. The	e original implementation processes up t	o 4096 tokens
848	per document. It has the same limitation as the	e ToBERT. For GPT40 and Gemini1.5pro	o, we input all
849	tokens together with our instruction in the pror	npt due to the large input size supported	by these large
850	language model. In contrast to these baseline n	nodels, our approach flexibly segments th	ne entire input
851	into chunks of varying sizes, using semantic k	eyphrases to minimize information loss.	Additionally,
852	we compute chunk-level attention to capture l	ong-range dependencies more effectivel	у.
853			

A.4 DETAILS OF DATASETS

Datasets	Train/Dev/Test	#Classes	Avg. Length
HP	516/64/65	2	705
LUN	12003/2992/2250	3	480
EURLEX57k	45000/6000/6000	4271	707
-INVERTED	45000/6000/6000	4271	707
GUM	179/26/26	21	972
CoNLL	120/20/20	37	5065

Table 11: The split and statistics of the datasets, including document classification task (HP, LUN, EURLEX57K, and Inverted EURLEX57K) and token classification task (GUM, CoNLL)

We analyzed the data distribution across the datasets used in this paper. Of these, the CoNLL dataset has the highest average number of tokens per document at 5,065. In contrast, LUN has the shortest average length, with 480 tokens per document. Both HP and EURLEX57K have similar average document lengths, measuring 705 and 707 tokens respectively. GUM presents a relatively higher token count, averaging 972 tokens per document.

Regarding the number of classes, EURLEX57K is the most complex dataset, containing 4,271 unique labels. In comparison, HP and LUN are more limited, with only 2 and 3 classes respec-tively. GUM and CoNLL are more diverse, with 21 and 37 different classes. Beyond label diversity, EURLEX57K also has the largest sample size, comprising 45,000 training samples, 6,000 devel-opment samples, and 6,000 test samples. LUN is the second-largest dataset, with 12,003 training samples, 2,992 development samples, and 2,250 test samples. Due to our selection strategy, CoNLL has the longest average document length, with the fewest samples. It has a total of 160 documents split into 120/20/20 for training, development, and test sets. GUM follows a similar distribution with 179/26/26 samples. The HP dataset includes 516 training samples, 64 development samples, and 65 test samples.

A.5 IMPLEMENTATION DETAILS

A.5.1 EXPERIMENT HYPERPARAMETERS

We performed extensive experiments to select the hyperparameters, including chunk size, token weights, learning rates, and warm-up strategies and steps. The optimal hyperparameters for each dataset for our proposed ChuLo model are presented in Table 12.

Hyperparameter	HP	LUN	EURLEX57K	I-EURLEX57K	CoNLL	GUM
Number of top-n phrases	15	15	15	15	15	15
Chunk size n	10	50	5	5	20	50
Weight for T_k	0.8	0.5	0.8	0.8	0.8	0.8
Weight for T_{nk}	0.1	0.1	0.1	0.1	0.1	0.1
Learning Rate	5e-5	5e-5	5e-5	5e-5	5e-5	5e-5
Batch Size	16	32	16	16	2	8
Warm-up Strategy	Linear	Linear	Cosine	Cosine	Linear	Linear
Warm-up Steps	10%	10%	5%	5%	10%	10%
Mex epoch	100	100	100	100	100	100
Stop Patience	10	10	10	10	10	10
Optimizer	AdamW	AdamW	AdamW	AdamW	AdamW	AdamW
Optimizer Weight Decay	1e-2	1e-2	1e-2	1e-2	1e-2	1e-2
Optimizer Betas	0.9, 0.999	0.9, 0.999	0.9, 0.999	0.9, 0.999	0.9, 0.999	0.9, 0.999

Table 12: The optimal hyperparameters used in our experiments.

A.5.2 HARDWARE INFORMATION

Our experiments are run on the Linux platform with an A6000 Nvidia graphic card and an AMD Ryzen Threadripper PRO 5955WX 16-core CPU, and the RAM is 128G.

A.6 PROMPT METHOD

- 917 We employed zero-shot prompting with large language models (LLMs), specifically Gemini 1.5 Pro and GPT40, in our experiments. The prompts used for each dataset are detailed in Table 13 and 14:

Dataset	Prompt
LUN	Task Definition: You are provided with a news article. Your task is to classify the article into one of the following categories: "Satire" "Hoax" or "Propaganda" Respond
HP	Task Definition: You are provided with a news article. Your task is to classify whether
	"False" if it is not. The news is: [{input}].
	·

Table 13: The prompt we used for each dataset in our experiments.

Dataset	Prompt
CoNLL	In the task of Named Entity Recognition, the B-, I-, and O- prefixes are commonly
	used to annotate slot types, indicating the boundaries and types of slots. These labels
	typically represent: B- (Begin): Signifies the beginning of a slot, marking the start
	of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation
	of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For
	instance, in a sentence where we want to label a "date" slot, words containing date
	information might be tagged as "B-date" (indicating the beginning of a date slot), fol-
	lowed by consecutive words carrying date information tagged as "I-date" (indicating
	the continuation of the date slot), while words not containing date information would
	be tagged as "O" (indicating they are outside any slot).
	Definition: In this task, you are given a conversation, where the words spoken by a
	versation into one of the 37 different entities. The entities are: "O" "R DEPCON"
	"I-PERSON", "B-NORP", "I-NORP", "B-FAC", "I-FAC", "B-ORG", "I-ORG", "R-
	GPE", "I-GPE", "B-LOC", "I-LOC", "B-PRODUCT", "I-PRODUCT", "B-DATE",
	"I-DATE", "B-TIME", "I-TIME", "B-PERCENT", "I-PERCENT", "B-MONEY", "I-
	MONEY", "B-QUANTITY", "I-QUANTITY", "B-ORDINAL", "I-ORDINAL", "B-
	CARDINAL", "I-CARDINAL", "B-EVENT", "I-EVENT", "B-WORK_OF_ART",
	"I-WORK_OF_ART", "B-LAW", "I-LAW", "B-LANGUAGE", "I-LANGUAGE".
	Only output entities. And the entity types should be output as a list without any ex-
	planation. The input is [{input}].
GUM	In the task of Named Entity Recognition, the B-, I-, and O- prefixes are commonly
	used to annotate slot types, indicating the boundaries and types of slots. These labels
	typically represent: B- (Begin): Signifies the beginning of a slot, marking the start
	of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation
	of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For
	instance, in a sentence where we want to label a "date" slot, words containing date
	information might be tagged as "B-date" (indicating the beginning of a date slot), fol-
	the continuation of the date slot) while words not containing date information would
	be tagged as " Ω " (indicating they are outside any slot)
	Definition: In this task, you are given a conversation, where the words spoken by a
	person are shown as a list. Your job is to classify the words in the following conversa-
	tion into one of the 37 different entities. The entities are: "I-abstract", "B-object".
	"B-place", "I-substance", "I-time", "I-place", "B-time", "B-abstract", "I-person",
	"B-plant", "B-substance", "I-animal", "B-organization", "I-event", "B-person", "B-
	event", "I-plant", "I-organization", "O", "I-object", "B-animal". Only output entities.
	And the entity types should be output as a list without any explanation. The input is
	[{input}].
	,

Table 14: The prompt we used for each dataset in our experiments.

Table 3 shows that LLMs outperform Longformer in the document classification task with zero-shot prompt tuning. However, their performance drops significantly in the NER task in Table 5a and Table 5b. For instance, in Figure 12, both GPT4o and Gemini1.5pro only predicted a single correct label, "O". Moreover, the models often fail to predict a sufficient number of token labels for longer inputs, or they repeatedly predict all "O" labels or redundant label sequences. These inconsistencies suggest that LLMs struggle to generate outputs matching the input length in token classification, highlighting substantial room for improvement in this area.

A.7 MORE CASE STUDIES

In this section, we will present several prompt and output samples for the long documents from the LUN (Figures 4) and 5) and Hyperpartisan (Figures 6 and 7) datasets for document classification, as well as GUM (Figures 9 and 10) and CoNLL (Figures 11, 12, 3 and 14) datasets for NER task. Documents with various lengths are randomly selected to see the comparison of our model against GPT-4, Gemini1.5pro and Longformer. While there is always at least one baseline which predicts wrongly for the difficult cases presented for the document classification task, we can observe that our model consistently classifies those documents well. For the token classification task, our model can also correctly classify more tokens than each baseline across the shown cases.

1004	
1005	Case 1, length: 3928 - Document and Prompt
1006	Task Definition: You are provided with a news article. Your task is to classify the article into one
1007	of the following categories: "Satire," "Hoax," or "Propaganda." Respond only with the
1008	appropriate category. The news is: [is obama a liar? or just loyal to his faith! the holy guran instructs its followors to lie to strongthon islam, guran 3:26, 3:54, 9:3, 40:28, and 16:106 are
1009	where you can find these instructions. obama has been a loyal follower and i am sure satan
1010	must be proud. any book or prophet that instructs its followers to chop the heads off non
1011	believers is a prophet of satan]
1012	Case 1 – Our Model
1013	Hoax
1014	
1015	Case 1 - GPI-40 Response
1016	Propaganda
1017	Case 1 – Gemini1.5pro Response
1018	Propaganda
1019	
1020	Case 1 – Longformer Output
1021	Propaganda

Figure 4: Prompt and output for a sample document of length 3928 in LUN dataset, where the correct prediction is highlighted in green and wrong predictions are highlighted in red. Compared to GPT40, Gemini1.5pro and Longformer, our model can **correctly** classify the given document as Hoax.

1026	
1027	Case 2, length: 2410 - Document and Prompt
1028	Task Definition: You are provided with a news article. Your task is to classify the article into one of the following categories: "Satire," "Hoax," or "Propaganda," Respond only with the
1029	appropriate category. The news is: [the irony is inescapable: in reaction to the historic drought
1030	that has transformed the california dream into california dust, the state is now embarking on the construction of a wave of desalination plants that will turn ocean water into fresh water.
1031	tragically, these power-hungry desalination plants will be running primarily on fossil fuel-
1032	generated electricity , meaning that california residents will have to commit]
1033	Case 2 – Our Model
1034	Propaganda
1035	
1036	Case 2 - GPT-40 Response
1037	Satire
1038	Case 2 – Gemini1 5nro Response
1039	
1040	Saure
1041	Case 2 – Longformer Output
1042	Satire
1043	

Figure 5: Prompt and output for a sample document of length 2410 in LUN dataset, where the correct prediction is highlighted in green and wrong predictions are highlighted in red. Compared to GPT40, Gemini1.5pro and Longformer, our model can **correctly** classify the given document as **Propaganda**.

1057	
1058	
1059	Case 3, length: 6800 - Document and Prompt
1060	Task Definition: You are provided with a news article. Your task is to classify whether the article
1061	is hyperpartisan. Respond only with "True" if the news is hyperpartisan or "False" if it is not.
1062	The news is: [<a "<="" -trump-s-missing-signature-force-him-to-be-deposed-="" 17799862="" td="">
1063	deal-shows-value-of-cfo-in-prosecution-1535448600" type="external">immunity deal shows
1064	value of cfo in prosecution. <a donalds-farewell-message-n904171"<="" td="">
1065	type="external">mccain's long goodbye.]
1066	Case 3 – Our Model
1067	False
1068	
1069	Case 3 - GPI-4o Response
1070	False
1071	Case 3 – Gemini1.5pro Response
1072	True
1073	
1074	Case 3 – Longformer Output
1075	False
1076	

Figure 6: Prompt and output for a sample document of length 6800 in Hyperpartisan dataset, where correct predictions are highlighted in green and the wrong prediction is highlighted in red. Compared to Gemini1.5pro, our model, GPT40 and Longformer can correctly classify the given document as False.

	Case 4, length: 2445 - Document and Prompt Task Definition: You are provided with a news article. Your task is to classify whether the article is hyperpartisan. Respond only with "True" if the news is hyperpartisan or "False" if it is not. The news is: [xp>alt-left refers to a loosely defined term to describe left-wing principles, organizations, politicians and activists, encompassing almost everything outside the norm of mainstream democratic liberal politics. used as both a pejorative and an affirmative, alike, alt-left has been used as a catchall term for far-left political ideologies such as socialism, anarchism, communism as well as antifa groups. however, because this term remains so undefined] Case 4 – Our Model False Case 4 - GPT-40 Response</a 	
	Task Definition: You are provided with a news article. Your task is to classify whether the article is hyperpartisan. Respond only with "True" if the news is hyperpartisan or "False" if it is not. The news is: [alt-left refers to a loosely defined term to describe left-wing principles, organizations, politicians and activists, encompassing almost everything outside the norm of mainstream democratic liberal politics. used as both a pejorative and an affirmative, alike, alt-left has been used as a catchall term for far-left political ideologies such as socialism, anarchism, communism as well as <a "="" type="">type="""TempCase 4 – Our ModelFalseCase 4 - GPT-40 ResponseTrue	
	hyperpartisan. Respond only with "True" if the news is hyperpartisan or "False" if it is not. The news is: [alt-left refers to a loosely defined term to describe left-wing principles, organizations, politicians and activists, encompassing almost everything outside the norm of mainstream democratic liberal politics. used as both a pejorative and an affirmative, alike, alt-left has been used as a catchall term for far-left political ideologies such as socialism, anarchism, communism as well as antifa groups. however, because this term remains so undefined] Case 4 – Our Model False Case 4 - GPT-40 Response True</a 	
	[alt-left refers to a loosely defined term to describe left-wing principles, organizations, politicians and activists, encompassing almost everything outside the norm of mainstream democratic liberal politics. used as both a pejorative and an affirmative, alike, alt-left has been used as a catchall term for far-left political ideologies such as socialism, anarchism, communism as well as <a <="" a="" type="internal"> Case 4 – Our Model False Case 4 - GPT-40 Response	
	politics. used as both a pejorative and an affirmative, alike, alt-left has been used as a catchall term for far-left political ideologies such as socialism, anarchism, communism as well as antifa groups. however, because this term remains so undefined] Case 4 – Our Model False Case 4 - GPT-40 Response</a 	
	far-left political ideologies such as socialism, anarchism, communism as well as antifa type="internal">Case 4 - Our Model False Case 4 - GPT-40 Response	
	Case 4 - Our Model False Case 4 - GPT-40 Response	
	Case 4 - Our Model False Case 4 - GPT-40 Response True	
	False Case 4 - GPT-40 Response True	
	Case 4 - GPT-4o Response	
	Case 4 – Gemini1.5pro Response	
	True	
	Case 4 – Longformer Output	
	False	
7	Description of successful and successful states of the set of 1445 in Hernemantices and states to a	.1
gure /	redictions are highlighted in green and wrong predictions are highlighted in red. Com	nere
GPT4	o and Gemini1 5pro, our model and Longformer can correctly classify the given docu	ment
False	and Commission, our model and hongronnel can correctly classify the given doca	mem
	Case 5, length: 895 - Document and Prompt	
	"In the task of Named Entity Recognition, the R-L, and O-prefixes are commonly used to apportate slot	
	types indicating the boundaries and types of clots. These labels typically represent: R. (Berin):	
	types, malating the boundaries and types of slots. These labels typically represent. b- (begin).	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, more the slot of the slot of a lot. (Inside): Represents the interior	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date""" slot, words	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """Indate"" date """ slot, words containing date information might be tagged as """B-date""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date""" slot, words containing date information might be tagged as """B-date""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """-I-date""" (indicating the continuation of the date slot), while words not containing date information would be tagged as	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """"B-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (indicating the continuation of the date slot), while words not containing date information would be tagged as """"O""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your iob is to classify the words in the	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """"B-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date""" (indicating the continuation of the date slot), while words not containing date information would be tagged as """"O""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B-	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """"date"""" (Indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"date""" (Indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"date""" (Indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"date information would be tagged as """"0"""" (Indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B- place', 'I-substance', 'I-time', 'I-place', 'B-bstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-point', 'D-object', 'B-animal', Only	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """"B-date"""" (Indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (indicating the beginning of a date slot), where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: I-abstract', 'B-object', 'B- place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-person', 'B-plant', 'D-siget', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """"B-date"""" (Indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (Indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (Indicating the continuation of the date slot), while words not containing date information would be tagged as """"O"""" (Indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: I-abstract', 'B-object', 'B- place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'been', 'L', 'Well', ""she's', "she', "mis', 'thet', 'I', 'berun', 'D', 'I', 'that', 'it', '-', 'it', 'really', 'surprises', 'me', 'L', 'Well', ""she's', "she', 'mis', ''', 'berun', 'to', 'I', 'that', 'it', '-',	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date""" slot, words containing date information might be tagged as """"B-date""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date""" (indicating the continuation of the date slot), while words not containing date information would be tagged as """"O"""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B- place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plat', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'been', 'I', 'Well', ""she', 'an', 'interest', 'in', 'this', 'retirement', 'bit', 'that', 'it', '', 'it', 'really', 'surprises', 'me', ',', 'Well', "she', 'used', 'to', 'go', 'over', 'and', 'read', 'a', 'book', 'or',	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date""" slot, words containing date information might be tagged as """B-date""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """I-date""" (indicating the continuation of the date slot), while words not containing date information would be tagged as """"O""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B- place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'L-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'been', '!, 'has', 'taken', 'such', 'an', 'interest', 'in', 'this', 'retirement', 'i', 'tat', 'it', 'a-', 'it', 'really', 'surprises', 'me', '!, 'Well', "'she', 'used', 'to', 'go', 'over,' and', 'read', 'a', 'book', 'or', 'something', '!, Yeah', '!, 'or', 'turn', 'a', 'deaf', 'ear', '!', 'That', 'was', 'for', 'sure', '!', 'But', 'some', 'basi', ',', 'yes', '', 'Yeah', '', '', ''', '''', ''''', '''', '''', ''''', ''''', ''''', ''''', ''''', ''''', ''''', ''''''	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """B-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (indicating the continuation of the date slot), while words not containing date information would be tagged as """"O"""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B- place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', ''B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'been', '!, 'has', 'taken', 'such', 'an', 'interest', 'in', 'this', 'retirement', 'bit', 'that', 'tr', '', 'it', 'really', 'surprises', 'me', '!, 'Well', "'she', 'used', 'to', 'go', 'over,' and', 'read', 'a', 'book', 'or', 'something', '!, Yeah', '!, 'or', 'turn', 'a', 'deaf', 'ear', '!, 'That', 'was', 'for', 'sure', '!, 'But', 'some', 'basil', '!, 'yes', '!, 'Yeah', '!, 'm', ''do', ""n't'", 'have', 'any', 'this', 'year', ',', '', 'forgot',]."	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represent: In the following of a slot, antiking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" date slot), followed by consecutive words carrying date information tagged as """I-date""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """I-date""" (indicating the generation of the date slot). While words not containing date information would be tagged as """"O""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: I-labstract', 'B-object', 'B-place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', has', 'been', 'I', 'aken', 'such', 'an', 'interest', 'in', 'this', 'retirement', 'bit', 'that', 'it', '', 'it', 'really', 'surprises', 'me', 'I', 'Well', ""she', 'used', 'to', 'go', 'over', 'and', 'read', 'a', 'book', 'or', 'something', 'Yeah', 'I', 'm', 'a', 'deaf', 'ar', 'I', That', 'was', 'for', 'surprise', 'mas', 'I', '', '', '', '', '', '', '', '',	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """B-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """I-date"""" (indicating the continuation of the date slot), while words not containing date information would be tagged as """"O"""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: Ii-abstract', 'B-object', 'B- place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', has', 'been', ',', 'has', 'taken', 'such', 'an', 'interest', 'in', 'this', 'retirement', 'bit', 'that', 'it', '', 'it', 'really', 'surprises', 'me', ',', 'Well', ""she', 'used', 'to', 'go', 'over', 'and', 'read', 'a', 'book', 'or', 'something', '.', 'Yeah', '.', 'r', ''un', 'a', 'deaf', 'ear', '.', 'That', 'was', 'for', 'sure', '.', 'but', 'somet', 'basil', ',', 'yes', ', 'Yeah', ',', 'r, '', ''event', 'I-event', 'I-event', 'B-abstract', 'I-abstract', 'I-abstrac	
	Signifies the beginning of a slot, marking the start of a new slot. I- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" date""" slot, words containing date information might be tagged as """B-date""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """I-date""" (indicating the continuation of the date slot), while words not containing date information would be tagged as """"O"""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: I-abstract', 'B-object', 'B- place', I-substance', I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'been', ',', 'has', 'taken', 'such', 'an', 'interest', 'in', 'this', 'retirement', 'bit', 'that', 'it', '', 'it', 'really', 'surprises', 'me', ',', 'Well', ""she', 'west', 'to', 'isot', 'to', 'isot', 'a', 'book', 'or', 'something', '.', Yeah', ',', 'or', 'turn', 'a', 'deaf', 'ear', ',', 'That', 'was', 'for', 'sure', '.', 'bas', 'a', 'book', 'or', 'something', '.', Yeah', ',', 'n', 'u', 'deaf', 'ear', ',', 'That', 'was', 'for', 'sure', ',', 'l', 'somet', 'basil', ',', 'yes', ', 'Yeah', ',', 'n', 'mon't'", 'have', 'any', 'this', 'year', ',', 'l', 'forgot',]." Case 5 - Our Model	
	Signifies the beginning of a slot, marking the start of a new slot. I - (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" date """ date slot, followed by consecutive words carrying date information tagged as """I-date""" (Indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """I-date""" (Indicating the tagged as """O""" (Indicating the y are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: I-abstract', 'B-object', 'B-place', 'I-substance', 'I-time', 'I-place', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-atima', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'been', '.', 'Well', "mshe's", 'she', 'ms', 'begun', 'to', 'listen', ', 'wes', 'she', 'has', '.', 'U', ', 'she', 'as', 'in', 'this', 'retirement', 'bit', 'that', 't', '', 'it', really', 'surprises', 'me', '.', 'Well', "she', 'use', 'to', 'go', 'over', 'and', 'read', 'a', 'book', 'or', 'something', '.', 'Yeah', '.', 'r', 'deaf', 'ear', '.', 'That', 'was', 'for', 'sure', '.', 'But', 'some', 'basi', ', 'yes', ', 'yes', '.', 'Fennt', 'I-event', 'I-event', 'I-event', 'B-abstract', 'I-abstract', 'I-abstract', 'I-abstract', 'I-abstract', '] Case 5 - Our Model ['B-event', 'I-event', 'I-event', 'I-event', 'I-event', 'B-abstract', 'I-abstract', 'I-abstract',]	
	Signifies the beginning of a slot, marking the start of a new slot. I - (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """B-date""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """-I-(Inside): Represents the interior of the date slot), while words not containing date information would be tagged as """"O""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B-place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-ospanization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list throut any explanation. The input is ['Sam', 'has', 'been', 'I, 'wel', 'wel', 'wel', 'wel', 'wel', 'wel', ''she', 'has', 't', '', '', '', '', '', '', '', '',	
	Signifies the beginning of a slot, marking the start of a new slot. I - (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" date """ slot, words containing date information might be tagged as """B-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"-date""" (indicating the y consecutive words carrying date information tagged as """"-date""" (indicating the y are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B-place', 'I-substance', 'I-time', 'I-place', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'taken', 'such', 'an', 'interest', 'in', 'this', 'retirement', 'bot', 'that', 'tr', '', 'it', 'really', 'surprises', 'me', '.', 'Well', '"she', 'used', 'to', 'go', 'over', 'and', 'read', 'a', 'book', 'or', 'something', '.', Yeah', '.', 'or', 'well', ''she', 'as', 'is', ''nter', 'as', 'go', 'over', 'and', 'read', 'a', 'book', 'or', 'something', '.', Yeah', '.', 'or', 'turn', 'a', 'deaf', 'ear', '.', 'That', 'was', 'for', 'sure', '.', 'But', 'somet', 'basi', '.', 'yes', '.', 'Yeah', '.', 'or', ''don', ''nter', 'l-event', 'l-event', 'l-event', 'l-event', 'l-event', 'l-event', 'l-event', 'l-event', 'l-event', 'l-gates', 'mor', '', 'l', 'forgot',]. Case 5 - Our Model [Palace', 'B-substance', 'l-preson', 'B-animal', 'l-object', 'B-lime', 'l-organization', 'l-place',] Case	
	Signifies the beginning of a slot, marking the start of a new slot. II- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"deter""" slot, words containing date information might be tagged as """"I-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (indicating the vare outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'I-abstract', 'B-object', 'B-place', 'I-substance', 'I-lipace', 'B-time', 'B-abstract', 'I-person', 'B-plant', 'B-substance', 'I-animal', 'B-organization', 'I-event', 'B-event', 'I-plant', 'Horganization', 'O', 'I-object', 'B-animal', 'B-organization', 'I-event', 'B-event', 'I-plant', 'Horganization', 'O', 'I-object', 'B-animal', 'B-organization', 'I-event', 'B-event', 'I-plant', 'Horganization', 'O', 'I-object', 'B-animal', Only output entities. And the entity types should be output as a list without any explanation. The input is ['Sam', 'has', 'been', 'J', 'wel', 'm', 'a', 'sea', 'm', 'sea', 'm', 'sea', 'm', 'sea', 'sea', 'm', 'some', 'something', 'J', 'Yeah', 'J', 'u', ', 'she', 'she', 'sae', 'sea', 'so', 'go', 'over', 'and', 'read', 'a', 'book', 'or', 'something', 'J', 'Yeah', 'J', 'event', 'I-event', 'I-event', 'I-event', 'I-event', 'I-event', 'B-abstract', 'I-abstract', 'I-abstract',] Case 5 - Our Model [B-person', O', O', O', O', O', O', O', O', O', 'L', 'D', 'L', 'B-place', 'B-lime', 'I-organization', 'I-place', '] Case 5 - Gemini1.5pro Response [B-person', O', O', O', O', O', O', O', O', O', O	
	Signifies the beginning of a slot, marking the start of a new slot. II- (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be tagged as """"I-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"I-date"""" (indicating the y are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are: 'l-abstract', 'B-object', 'B-object', 'B-object', 'B-object', 'B-object', 'B-object', 'B-object', 'B-object', 'B-animal', 'B-organization', 'I-event', 'I-planc', 'I-pganization', 'O', I-object', 'B-animal', 'B-organization', 'I-event', 'B-event', 'I-plant', 'I-organization', 'O', I-object', 'B-animal', 'B-organization', 'I-event', 'B-event', 'I-plant', 'B-object', 'B-animal', 'B-object', 'B-object', 'B-animal', 'B-organization', 'I-event', 'I-planc', 'I-organization', 'O', I-object', 'B-animal', 'B-organization', 'I-event', 'I-event', 'I-plant', 'I-organization', 'I-y, 'I-s', 'I-a', 'I', 'r-a', 'I', 'r-a', 'I', 'really', 'surprises', 'may, 'taken', 'such', 'an', 'Irt', ''-', '', '', '', '', '', '', '', '',	
	<pre>Signifies the beginning of a slot, marking the start of a new slot. In (Inside): Represents the interior of a slot, indicating a continuation of the slot. O (Outside): Denotes parts of the input that are not part of any slot. For instance, in a sentence where we want to label a """"date"""" slot, words containing date information might be targed as """B-date"""" (indicating the beginning of a date slot), followed by consecutive words carrying date information tagged as """"l-date"""" (indicating the continuition of the date slot), while words not containing date information would be tagged as """"O""" (indicating they are outside any slot). Definition: In this task, you are given a conversation, where the words spoken by a person are shown as a list. Your job is to classify the words in the following conversation into one of the 37 different entities. The entities are 'l-abstract', 'B-object', 'B- place', 'L-substance', 'L-ime', 'B-bestract', 'I-person', 'B-plant', 'B-substance', 'L-animal', 'B-organization', 'I-event', 'B-person', 'B-event', 'I-plant', 'I-organization', 'O', 'I-object', 'B-animal'. Only output entities. And the entity types should be output as a list without any explanation. The input is [I'Sam', 'has', 'been', 'J', 'Nat,' 'was', 'this', 'retirement', 'I', 'tet', '', 'it', 'really', 'surprises', 'me', '.', 'Well', ""she's''', 'she', 'used', 'to', 'go', 'over', 'and', read', 'a', 'book', 'or', 'something', '.', Yeah', '.', 'I', ''', 'U', '', 'abe', 'used', 'to', 'go', 'over', 'and', read', 'a', 'book', 'or', 'something', '.', Yeah', '.', 'I', '''', ''', ''', '''', '''', '''', '''', ''''', ''''', ''''', ''''''</pre>	

Figure 8: Prompt and output for a sample document of length 895 in **GUM** dataset for NER task, where correct predictions are highlighted in green and wrong predictions are highlighted in red.



Figure 10: Prompt and output for a sample document of length 1281 in **GUM** dataset for NER task, where correct predictions are highlighted in green and wrong predictions are highlighted in red.



Figure 12: Prompt and output for a sample document of length 3038 in **CoNLL** dataset for NER task, where correct predictions are highlighted in green and wrong predictions are highlighted in red.

