

DELTA: AN ONLINE DOCUMENT-LEVEL TRANSLATION AGENT BASED ON MULTI-LEVEL MEMORY

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) have achieved reasonable quality improvements in machine translation (MT). However, most current research on MT-LLMs still faces significant challenges in maintaining translation consistency and accuracy when processing entire documents. In this paper, we introduce DELTA, a **Document-level Translation Agent** designed to overcome these limitations. DELTA features a multi-level memory structure that stores information across various granularities and spans, including Proper Noun Records, Bilingual Summary, Long-Term Memory, and Short-Term Memory, which are continuously retrieved and updated by auxiliary LLM-based components. Experimental results indicate that DELTA significantly outperforms strong baselines in terms of translation consistency and quality across four open/closed-source LLMs and two representative document translation datasets, achieving an increase in consistency scores by up to 4.58 percentage points and in COMET scores by up to 3.16 points on average. DELTA employs a sentence-by-sentence translation strategy, ensuring no sentence omissions and offering a memory-efficient solution compared to the mainstream method. Furthermore, DELTA improves pronoun and [context-dependent](#) translation accuracy, and the summary component of the agent also shows promise as a tool for query-based summarization tasks. [The code and data is anonymously available at https://anonymous.4open.science/r/DelTA_Agent-7716](https://anonymous.4open.science/r/DelTA_Agent-7716).

1 INTRODUCTION

Large language models (LLMs) such as GPT-4 (OpenAI, 2023) have recently demonstrated reasonable performance on the machine translation (MT) task within the natural language processing domain (Garcia & Firat, 2022; Hendy et al., 2023; Zhang et al., 2023; Siu, 2023; Jiao et al., 2023). Numerous studies have been carried out to further unleash LLMs’ potential for MT (Ghazvininejad et al., 2023; Peng et al., 2023; Zeng et al., 2023; He et al., 2024; Wang et al., 2024c). However, the majority of these researches mainly focus on sentence-level translation, operating under the strong assumption that source sentences are independent of one another. This isolated approach may fail to model the discourse structure and overlook the coherence in continuous document texts (Scarton & Specia, 2015; Bawden et al., 2018).

Document-level machine translation (DocMT) systems have been receiving growing focus in recent years, which involves the whole document or some part of it to capture more context information to guide the translation process (Kim et al., 2019; Maruf et al., 2021). Researchers find that modeling discourse phenomena (Bawden et al., 2018) while translating the whole document helps increase the coherence and consistency in the generated translation (Maruf & Haffari, 2018; Wang et al., 2017; Zhang et al., 2018; Tan et al., 2019). Recently, a few studies have proposed to introduce LLMs to the DocMT task, utilizing their inherent context information modeling and long text processing capabilities (Wang et al., 2023b; Wu & Hu, 2023; Wu et al., 2024a). However, existing DocMT-LLMs still suffer from critical issues such as occasional content omissions and terminology translation inconsistency (Karpinska & Iyyer, 2023). These issues seriously affect the reliability of the developed system, especially when accurate document translations are required.

LLM-based autonomous agents equipped with specially designed memory components can efficiently store and retrieve key information embedded in the environment. These data assist the inference process of LLMs, facilitating the handling of complex tasks and environments through

self-directed planning and actions (Wang et al., 2024a; 2023a; Park et al., 2023; Lee et al., 2024). Inspired by this, we propose **DELTA**, an online **Document-levEL Translation Agent** based on multi-level memory components. Specifically, we store information in four memory components: Proper Noun Records, Bilingual Summary, Long-Term Memory, and Short-term Memory, and utilize LLMs to update and retrieve them. Proper Noun Records maintain a repository of previously encountered proper nouns and their initial translations within the document, ensuring consistency by reusing the same translation for each subsequent occurrence of the same proper noun. The Bilingual Summary contains summaries of both the source and target texts, capturing the core meanings and genre characteristics of the documents to enhance translation coherence. Long-Term Memory and Short-Term Memory store contextual sentences over broader and narrower scopes, respectively. Long-Term Memory is accessed by LLMs to retrieve sentences most relevant to the current source sentence, while Short-Term Memory provides instant context to support the translation process. During translation, sentence pairs are drawn from the Long-Term and Short-Term memory as the exemplars for few-shot learning demonstration, and the proper noun translation records and bilingual summaries are also integrated into the prompt for DocMT-LLMs as auxiliary information.

Experimental results indicate that DELTA achieves improvements in both translation consistency and quality. For translation consistency, DELTA achieves an average improvement of 4.36 percentage points across four translation directions from English and 4.58 percentage points across four directions into English. In terms of translation quality, DELTA yields an average improvement of 3.14 COMET points for four translation directions from English and 3.16 COMET points for four directions into English. Moreover, DELTA translate documents in a sentence-by-sentence manner (following an online approach) to avoid content omissions, ensuring sentence-level alignment of target documents with source documents. This manner also prevents memory bloat caused by data accumulation, making it more suitable for practical application scenarios.

Our main contributions are summarized as follows:

- We develop **DELTA**, an online DocMT agent employing a multi-level memory structure, which stores information across different granularities and spans.
- We demonstrate that DELTA substantially improves the consistency and quality of document translations. Additionally, the summary component of DELTA can function as an independent tool for query-based summarization tasks.
- We certificate that the sentence-wise translation approach employed by DELTA incurs a lower memory cost compared to existing document translation methods.
- We observe that DELTA is particularly effective in maintaining translation consistency over expended spans. Moreover, it enhances the pronoun translation accuracy in the document.

2 RELATED WORK

Document-Level Machine Translation In recent years, studies on DocMT have achieved rich results (Kim et al., 2019; Maruf et al., 2021; Fernandes et al., 2021). These studies can be separated into two categories. Studies of the first group employ a document-to-sentence (Doc2Sent) approach, where the source-side context sentences are encoded to generate the current target sentence (Wang et al., 2017; Tan et al., 2021; Lyu et al., 2021). However, these approaches suffer from limitations caused by separated encoding modules of the current sentences and their context (Sun et al., 2022; Bao et al., 2021), as well as the failure to utilize target-side context (Li et al., 2023b). Studies of the second group employ a document-to-document (Doc2Doc) approach, where the translation unit is extended from a single sentence to multiple sentences (Zhang et al., 2020; Liu et al., 2020; Lupo et al., 2022; Bao et al., 2021; Li et al., 2023b).

Autonomous Agents LLM-based autonomous agents have recently achieved remarkable performance in various NLP tasks. Park et al. (2023); Wang et al. (2023a); Lee et al. (2024) deal with long-context understanding and processing tasks by introducing carefully designed memory and retrieval workflows. Xu et al. (2024); Wang et al. (2024c); Feng et al. (2024) prompt LLMs to evaluate their own outputs and conduct refinement accordingly to improve the quality of the outputs. Li et al. (2023a); Liang et al. (2023); Li et al. (2024); Wu et al. (2024b) enhance the performance of LLMs on specific tasks by multi-agent interaction.

Window	LTCR-1	LTCR-1 _f	#Missing Sents	sCOMET	dCOMET
1	75.09	88.24	0	84.04	6.62
5	80.49	88.15	0	84.30	6.70
10	79.65	90.81	2	84.27	6.65
30	83.08	95.83	8	83.88	6.69
50	86.94	95.90	10	83.70	6.66

Table 1: Translation results with different translation window sizes. “#Missing Sents” represents the number of missing target sentences in the translated document.

Our method in this paper represents a Doc2Sent approach implemented through an LLM-based automatic agent. Instead of simply encoding source-side context to generate target sentences, our method directs LLMs to retrieve key information across varying granularities and spans, and store this data in memory components. During document translation, relevant information is incorporated into the prompts for DocMT-LLMs to assist the translation process.

3 MOTIVATION

3.1 MAIN CHALLENGES FOR DOCMT-LLMS

Due to the maximum context limitation inherent in LLMs, translating a lengthy document in a single pass becomes unfeasible. A conventional strategy involves segmenting the document into smaller translation windows and translating them sequentially. In our study, we initially leverage the GPT-3.5-turbo-0125 model to translate the IWSLT2017 En \Rightarrow Zh test set, comprising 12 documents sourced from TED talks. We employ a window of size l to facilitate document translation, where l source sentences are simultaneously processed to generate l hypothesis sentences. Once all source sentences are translated, they are concatenated to form the complete target document. The primary challenges associated with DocMT-LLMs arise from the following two aspects.

Translation Inconsistency Given a source document $D_s = (s_1, s_2, \dots, s_N)$ and its corresponding target document $D_t = (t_1, t_2, \dots, t_N)$, if there exists a proper noun $p \in P$ (P denotes the set of all proper nouns in D_s , including names of people, locations, and organizations), and p appears multiple times in D_s , we expect that all occurrences of its translation in D_t should be consistent.

Lyu et al. (2021) propose the Lexical Translation Consistency Ratio (LTCR), a metric that quantifies the proportion of consistent translation pairs among all proper noun translation pairs in the target document. However, we argue that the translations of the proper nouns are supposed to not only maintain consistency throughout the document but also align their first appearance. This consideration is particularly important for enhancing the reading experience of audiences. Therefore, we introduce the **LTCR-1** metric for the DocMT-LLMs, which calculates the proportion of proper noun translations that are consistent with the initial translation within the document:

$$\text{LTCR-1}(D_s, D_t) = \frac{\sum_{p \in P} \sum_{i=2}^{k_p} \mathbb{1}(\mathcal{T}_i(p) = \mathcal{T}_1(p))}{\sum_{p \in P} (k_p - 1)} \quad (1)$$

$\mathcal{T}_i(p)$ represents the i -th translation of p in D_t , and k_p denotes the number of occurrences of p in the document. The indicator function $\mathbb{1}(\mathcal{T}_i(p) = \mathcal{T}_1(p))$ returns 1 if the translations $\mathcal{T}_i(p)$ and $\mathcal{T}_1(p)$ are identical, and 0 otherwise. The numerator is the number of times the proper nouns appear again and their translation remains the same as their first appearance, and the denominator represents the sum of all occurrences except the first one of all proper nouns. To compute this metric, we initially annotate all proper nouns in the source document using spaCy¹. Subsequently, we utilize the token align tool awesome-align (Dou & Neubig, 2021)² to determine the translations of these proper nouns in the target document. To mitigate the impact of errors from the alignment tool, we introduce a fuzzy match version of this metric, where two proper noun translations are considered

¹<https://spacy.io/>

²<https://github.com/neulab/awesome-align/>

consistent when one is a substring of the other:

$$\text{LTCR-1}_f(\mathbf{D}_s, \mathbf{D}_t) = \frac{\sum_{p \in P} \sum_{i=2}^{k_p} \mathbb{1}(\mathcal{T}_i(p) \subseteq \mathcal{T}_1(p) \vee \mathcal{T}_1(p) \subseteq \mathcal{T}_i(p))}{\sum_{p \in P} (k_p - 1)} \quad (2)$$

As shown in Table 1, translating every sentence separately (window size = 1) causes poor translation consistency. An example is illustrated in Appendix A. Increasing the window size consistently leads to higher scores across all three consistency metrics. This suggests that when more sentences are processed within a single translation pass, the LLM is better able to model discourse phenomena and maintain consistent translation of proper nouns throughout the document. However, due to the inherent limitations in the context length of LLMs, resolving translation inconsistencies cannot be achieved solely by indefinitely expanding the window size.

Translation Inaccuracy When employing a large window size for document translation, LLMs tend to process the input source sentences as cohesive documents rather than as individual sentences. As a result, the model prioritizes maintaining the general meaning of the text and loses track of the detailed information in each sentence. This can lead to undertranslation issues and a decline in translation quality (Karpinska & Iyyer, 2023; Wu et al., 2024a). We utilize two neural metrics to assess the quality of document translation. The first is the sentence-level COMET (sCOMET) score³, for which we utilize the model Unbabel/wmt22-comet-da to obtain the scores. The second metric is the document-level COMET (dCOMET) score⁴ proposed by Vernikos et al. (2022), for which we use wmt21-comet-qe-mqm⁵ to derive reference-free scores. In calculating this document-level metric, the model encodes previous sentences as context rather than encoding only the hypothesis, making this approach more accurate for evaluating document translations.

As illustrated in Table 1, an increase in window size correlates with a higher tendency for the LLM to omit sentences from the source document, resulting in missing translations in the target document. An example of this undertranslation issue is presented in Appendix A. Additionally, quality metrics such as sCOMET and dCOMET do not demonstrate a consistent improvement with larger translation windows. Therefore, we conclude that translating documents in batches of several sentences at a time may introduce translation inaccuracy issues. These concerns are particularly significant in contexts where precise translations are essential, such as in technical manuals or official documents.

3.2 WHY USING A DOC2SENT APPROACH?

Previous experiments indicate that translating a document by processing multiple sentences at once may occasionally result in sentence omissions. Although human translators often translate entire paragraphs simultaneously, which can also lead to occasional omissions, their underlying translation mechanism is fundamentally distinct from that of DocMT-LLMs. DocMT-LLMs are prone to omitting source sentences due to hallucination issues or limited capabilities in handling long texts effectively. Therefore, we argue that, at this moment, a Doc2Sent approach offers a more promising alternative for DocMT-LLMs to produce precise and high-quality document translations. In our study, we provide LLMs with contextual information from the document while asking them to translate each source sentence separately. Once all sentences are translated, we concatenate the target sentences to form the final target document.

4 DELTA: DOCMT AGENT BASED ON MULTI-LEVEL MEMORY

Considering the multi-granularity and multi-scale of key information in the document during translation, we introduce **DELTA**, an online DocMT agent. DELTA employs a multi-level memory stream that captures and preserves critical information encountered throughout the translation process. This memory stream accommodates a wide range of perspectives, spanning from recent to historical, concrete to abstract, and coarse-grained to fine-grained details. DELTA translate the source document in a sentence-by-sentence manner while updating its memory in real-time. This approach addresses the

³<https://github.com/Unbabel/COMET/>

⁴<https://github.com/amazon-science/doc-mt-metrics/>

⁵<https://unbabel-experimental-models.s3.amazonaws.com/comet/wmt21/wmt21-comet-qe-mqm.tar.gz>

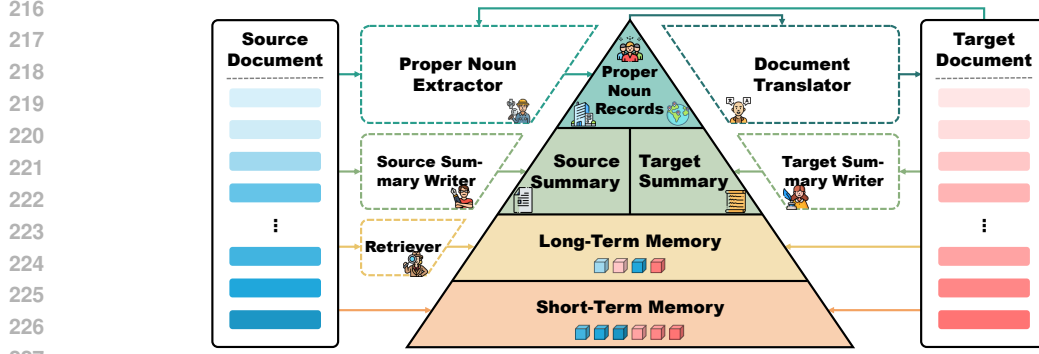


Figure 1: Framework of DELTA. The modules outlined with dashed lines represent the multi-level memory components, while those outlined with solid lines denote the LLM-based components. Memories closer to the top are more global, abstract, and densely packed with information. During translation, memory information is retrieved and incorporated into the translator LLM’s prompt. After the translation of each sentence, the LLM-based components extract key information from both the source and target documents and update the multi-level memory components.

context limitations of large language models and ensures the generation of a sentence-level aligned target document, thereby preserving both the quality and rigorousness of the translation. The main framework of DELTA is illustrated in Figure 1, the algorithm of DELTA is detailed in Algorithm 1, and the prompts used for each module are given in Appendix C.

Proper Noun Records The first level of the agent’s memory component we introduce is a dictionary called the Proper Noun Records \mathcal{R} to store proper nouns p in the document along with their translations upon first encounter $\mathcal{T}_1(p)$ within the document: $\mathcal{R}^{(i)} = \{(p, \mathcal{T}_1(p)) \mid p \in s_j, 1 \leq j < i\}$, where $\mathcal{R}^{(i)}$ represents the state of \mathcal{R} before translating the i -th sentence, and the same applies to other components. When translating the subsequent sentence s_i , the agent consults $\mathcal{R}^{(i)}$ to obtain all recorded proper nouns that are also contained in s_i : $\hat{\mathcal{R}}^{(i)} = \{(p, \mathcal{T}_1(p)) \mid p \in s_i, (p, \mathcal{T}_1(p)) \in \mathcal{R}^{(i)}\}$.

The Proper Noun Records are continuously updated by an LLM-based component known as the Proper Noun Extractor $\mathcal{L}_{\text{Extract}}$. After each sentence is translated, it extracts newly encountered proper nouns from the source sentence and their translations from the target sentence $\mathcal{L}_{\text{Extract}}(s_i, t_i) = \{(p, \mathcal{T}_j(p)) \mid p \in s_i, \mathcal{T}_j(p) \in t_i, \forall (p', \mathcal{T}_k(p)) \in \mathcal{R}^{(i)}, p \neq p'\}$ and add them to \mathcal{R} .

Bilingual Summary Unlike previous studies (Wang et al., 2023a; Lee et al., 2024), our research implements a bilingual summary approach as the second level of the agent’s memory component to address the challenges of extensive context on both the source and target sides. We maintain a pair of summaries throughout the translation process to enhance accuracy and fluency. The Source-Side Summary \mathcal{A}_s encapsulates the main content, domain, style, and tone of the previously translated sections of the document. This summary serves to preserve a coherent understanding of the text’s overall context, thereby aiding the LLMs in producing more accurate translations. Conversely, the Target-Side Summary \mathcal{A}_t focuses solely on the main content of the previously translated target text.

The pair of summaries are generated by two LLM-based components of the agent: the Source Summary Writer $\mathcal{L}_{\text{WriteS}}$ and the Target Summary Writer $\mathcal{L}_{\text{WriteT}}$. These summaries are updated every m sentences through a two-step process. Initially, the writers generate segment summaries for the last m sentences from both the source and target texts: $\tilde{\mathcal{A}}_s^{(i+1)} = \mathcal{L}_{\text{WriteS}}(s_{i-m+1}, \dots, s_i)$, $\tilde{\mathcal{A}}_t^{(i+1)} = \mathcal{L}_{\text{WriteT}}(t_{i-m+1}, \dots, t_i)$. Subsequently, these segment summaries are merged with the previous overall summaries for both sides to summary mergers to obtain new overall summaries: $\mathcal{A}_s^{(i+1)} = \mathcal{L}_{\text{MergeS}}(\mathcal{A}_s^{(i)}, \tilde{\mathcal{A}}_s^{(i+1)})$, $\mathcal{A}_t^{(i+1)} = \mathcal{L}_{\text{MergeT}}(\mathcal{A}_t^{(i)}, \tilde{\mathcal{A}}_t^{(i+1)})$. This process is repeated iteratively until all sentences in the source document have been read.

Long-Term & Short-Term Memory The last two levels of the agent’s memory component are the Long-Term Memory and the Short-Term Memory, respectively. These two components are

Algorithm 1: The Overall Framework of DELTA

```

270 input : Source document  $D_s = \{s_1, \dots, s_N\}$ , Large Language Model  $\mathcal{L}$ , Proper Noun
271         Records  $\mathcal{R} = \emptyset$ , Source-Side Summary  $\mathcal{A}_s = \emptyset$ , Target-Side Summary  $\mathcal{A}_t = \emptyset$ ,
272         Short-Term Memory  $\mathcal{M} = \emptyset$ , Long-Term Memory  $\mathcal{N} = \emptyset$ 
273 output: Target document  $D_t = \{t_1, \dots, t_N\}$ 
274  $D_t \leftarrow \emptyset$ 
275 for  $i = 1$  to  $N$  do
276     /* Retrieve memory */
277      $\hat{\mathcal{R}} \leftarrow \{(p, \mathcal{T}_1(p)) \mid p \in s_i, (p, \mathcal{T}_1(p)) \in \mathcal{R}\}$  /* Search Proper Noun Records */
278      $\hat{\mathcal{N}} \leftarrow \mathcal{L}_{\text{Retrieve}}(s_i, \mathcal{N})$  /* Match  $n$  relative sentences from Long-Term Memory */
279     /* Translate with hybrid memory information */
280      $t_i \leftarrow \mathcal{L}_{\text{Translate}}(s_i, \hat{\mathcal{R}}, \hat{\mathcal{N}}, \mathcal{A}_s, \mathcal{A}_t, \mathcal{M})$ 
281      $D_t \leftarrow D_t \cup \{t_i\}$  /* Add hypothesis to target document */
282     /* Update memory */
283      $\mathcal{R} \leftarrow \mathcal{R} \cup \mathcal{L}_{\text{Extract}}(s_i, t_i)$  /* Extract new proper nouns and add to records */
284      $\mathcal{N} \leftarrow \mathcal{N}[-l+1:] \cup \{(s_i, t_i)\}$  /* Last  $l$  sentences as Long-Term Memory */
285      $\mathcal{M} \leftarrow \mathcal{M}[-k+1:] \cup \{(s_i, t_i)\}$  /* Last  $k$  sentences as Short-Term Memory */
286     if  $i \bmod m = 0$  /* Update Bilingual Summary every  $m$  sentences */
287     then
288         /* Generate source and target segment summaries */
289          $\hat{\mathcal{A}}_s \leftarrow \mathcal{L}_{\text{WriteS}}(s_{i-m+1}, \dots, s_i)$   $\hat{\mathcal{A}}_t \leftarrow \mathcal{L}_{\text{WriteT}}(t_{i-m+1}, \dots, t_i)$ 
290         /* Merge segment summaries into document summaries */
291          $\mathcal{A}_s \leftarrow \mathcal{L}_{\text{MergeS}}(\mathcal{A}_s, \hat{\mathcal{A}}_s)$   $\mathcal{A}_t \leftarrow \mathcal{L}_{\text{MergeT}}(\mathcal{A}_t, \hat{\mathcal{A}}_t)$ 
292     end
293 end

```

designed to address the requisite coherence across document-level translations. The Short-Term Memory \mathcal{M} retains the last k source sentences along with their corresponding translations, where k represents a relatively small number: $\mathcal{M}^{(i)} = \{(s_{i-k}, t_{i-k}), \dots, (s_{i-1}, t_{i-1})\}$. This component is specifically designed to capture immediate contextual information in adjacent sentences, which is then seamlessly integrated into the translation prompt, serving as the context for the current sentence.

Similarly, the Long-Term Memory \mathcal{N} component preserves a broader range of context by maintaining a window of the last l sentences from the source document, with l being significantly greater than k , storing extended coherent information throughout the document. Before translating a given source sentence, an LLM-based component called the Memory Retriever $\mathcal{L}_{\text{Retrieve}}$ chooses n source sentences that are most relevant to the current translation query: $\hat{\mathcal{N}}^{(i)} = \mathcal{L}_{\text{Retrieve}}(s_i, \mathcal{N}^{(i)})$. These n sentences with their translations are subsequently employed as demonstration exemplars.

Document Translator We utilize an LLM-based component called Document Translator $\mathcal{L}_{\text{Translate}}$ to perform the final translation process. Information from the multi-level memory is integrated into the prompt to support the translator in producing high-quality and consistent translations: $t_i = \mathcal{L}_{\text{Translate}}(s_i, \hat{\mathcal{R}}^{(i)}, \hat{\mathcal{N}}^{(i)}, \mathcal{A}_s^{(i)}, \mathcal{A}_t^{(i)}, \mathcal{M}^{(i)})$. The sentence-by-sentence approach ensures that the resulting target document is consistently aligned with the source document at the sentence level, effectively minimizing the risk of missing target sentences. Additionally, this method allows for straightforward evaluation of translation quality using sentence-level metrics such as sCOMET.

5 EXPERIMENTS

5.1 SETTINGS

Datasets & Metrics We conduct our experiments on the two test sets. The first is the tst2017 test sets from the IWSLT2017 translation task⁶ (Akiba et al., 2004), which consists of parallel documents sourced from TED talks, covering 12 language pairs. Our experiments are conducted on eight language pairs: En \leftrightarrow Zh, De, Fr, and Ja. There are 10 to 12 sentence-level aligned parallel documents with approximately 1.5K sentences for each language pair. The second is Guofeng Webnovel⁷ (Wang et al., 2023c; 2024b), a high-quality and discourse-level corpus of web fiction. We conduct our experiments on the Guofeng V1 TEST.2 set in the Zh \Rightarrow En direction. The detailed dataset statistics are demonstrated in Appendix B. We employ LTCR-1 and LTCR-1_f for proper noun translation consistency evaluation and adopt sCOMET and dCOMET as translation quality metrics, which are all introduced in §3.1.

⁶<https://wit3.fbk.eu/2017-01-d/>

⁷<https://github.com/longyuewangdcu/GuoFeng-Webnovel/>

System	En \Rightarrow Xx				Xx \Rightarrow En			
	sCOMET	dCOMET	LTCR-1	LTCR-1 _f	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
NLLB	82.11	6.36	74.56	81.87	84.10	6.98	79.03	90.76
GOOGLE	80.41	5.83	81.38	84.72	80.17	5.96	81.43	90.81
GPT-3.5-Turbo								
Sentence	84.80	6.58	77.06	82.81	84.47	7.05	81.98	91.86
Context	85.40	6.70	77.34	83.12	84.97	7.15	85.03	95.27
Doc2Doc	-	6.62	79.12	86.39	-	6.96	85.17	92.98
DELTA	85.58	6.73	82.96	88.83	84.95	7.15	86.53	96.26
GPT-4o-mini								
Sentence	81.51	6.35	78.59	85.07	84.01	6.99	81.42	91.34
Context	84.78	6.65	80.01	86.99	84.95	7.15	84.40	94.34
Doc2Doc	-	6.75	80.54	85.39	-	7.01	83.50	93.39
DELTA	85.85	6.80	81.80	86.33	85.26	7.24	85.25	95.89
Qwen2-7B-Instruct								
Sentence	80.03	5.96	73.91	79.54	77.10	6.48	76.39	87.94
Context	80.84	6.08	79.59	85.35	83.09	6.84	81.48	92.56
Doc2Doc	-	5.83	77.32	84.59	-	6.59	85.03	93.68
DELTA	81.02	6.07	80.09	87.78	83.36	6.84	82.05	93.30
Qwen2-72B-Instruct								
Sentence	78.53	5.97	79.54	85.09	80.53	6.73	82.25	92.05
Context	80.79	6.22	79.14	85.40	83.27	6.99	82.86	92.21
Doc2Doc	-	6.45	73.58	78.64	-	6.87	83.00	90.74
DELTA	84.99	6.66	81.66	88.34	85.19	7.21	86.53	96.48
Average								
Sentence	81.22	6.21	77.27	83.13	81.53	6.81	80.51	90.80
Context	82.95	6.41	79.02	85.21	84.07	7.03	83.44	93.59
Doc2Doc	-	6.41	77.64	83.75	-	6.86	84.18	92.70
DELTA	84.36	6.57	81.63	87.82	84.69	7.11	85.09	95.48

Table 2: Test results on the IWSLT2017 dataset. Since the translations produced by the Doc2Doc method are not aligned at the sentence level with the source text, we do not report the sCOMET scores for this method. The highest score in each block is highlighted in **bold font**. The results in the “Average” block represent the mean scores across the four backbone models.

Models and Hyperparameters In this work, we utilize two versions of GPT models, GPT-3.5-Turbo-0125 and GPT-4o-mini, as our base models. We get access to these models through the official API provided by OpenAI⁸. We also introduce the open-source Qwen2-7B-Instruct⁹ and Qwen2-72B-Instruct¹⁰ in our experiments. The max_new_tokens is set to 2048 and other hyper-parameters remain default. The updating window of Bilingual summary m and length of Long-Term Memory l are set to 20. The number of retrieved relative sentences from Long-Term Memory n is set to 2. The length of Short-Term Memory k is set to 3.

Baseline Methods We include the following three approaches as our baselines. a) **Sentence**: We employ the same LLMs but conduct a sentence-level translation process to obtain the baseline results. b) **Context**: We follow Wu et al. (2024a) to provide the LLMs with three previously obtained source-target sentence pairs as the context of the current sentence. This method integrates more contextual information and helps adapt the LLMs for the document-level translation task. c) **Doc2Doc**: We reproduce the approach proposed by Wang et al. (2023b), translating 10 sentences in a single conversation turn and processing the entire document within a single chat box, thereby leveraging the long-term modeling ability of the LLMs. In computing the metric scores for the Doc2Doc results, we first perform sentence alignment using Bleualign¹¹ to obtain aligned source and target documents, after which we calculate the involved metrics. Furthermore, we also introduce the results of NLLB-3.3B (Costa-jussà et al., 2022) and GoogleTrans¹² for comparison.

⁸<https://platform.openai.com/docs/guides/text-generation/>

⁹<https://huggingface.co/Qwen/Qwen2-7B-Instruct/>

¹⁰<https://huggingface.co/Qwen/Qwen2-72B-Instruct/>

¹¹<https://github.com/rsennrich/Bleualign/>

¹²<https://py-googletrans.readthedocs.io/>

System	sCOMET	dCOMET	LTCR-1	LTCR-1 _f	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
GPT-3.5-Turbo								
Sentence	77.62	3.07	61.58	78.82	77.87	3.10	58.82	70.59
Context	78.57	3.19	70.10	81.37	78.56	3.19	64.32	74.37
Doc2Doc	–	2.82	77.46	89.02	–	2.96	82.04	91.62
DELTA	78.45	3.17	85.57	96.52	78.77	3.34	88.94	96.48
Qwen2-7B-Instruct								
Sentence	73.65	2.62	37.00	50.00	75.15	2.98	58.00	71.50
Context	76.54	3.01	52.82	61.54	77.87	3.20	58.21	70.15
Doc2Doc	–	2.69	73.25	84.08	–	2.77	80.79	90.07
DELTA	76.95	3.10	85.50	94.00	78.32	3.31	86.93	95.98

Table 3: Test results on the Guofeng dataset.

5.2 RESULTS

Test Results on IWSLT2017 The main experiment results on the IWSLT2017 test set are demonstrated in Table 2. For more detailed scores, please refer to Appendix D. It is evident that DELTA outperforms baseline approaches on LTCR-1 and LTCR-1_f metric scores across nearly all translation directions and models. This indicates that our approach yields significant enhancements in proper noun translation consistency for document-level translation. Furthermore, DELTA significantly improves the overall quality of document translation, as evidenced by consistently higher sCOMET and dCOMET scores. The superior dCOMET scores indicate that DELTA effectively captures contextual information to support the translation process. Translation consistency is improved significantly in the En \Rightarrow Zh direction, while gains in directions like En \Rightarrow De are modest. For instance, with GPT-3.5-Turbo, LTCR-1 improves by 6.17 percentage points (86.44 vs. 80.27) for En \Rightarrow Zh, but only 1.40 points (93.46 vs. 92.06) for En \Rightarrow De (see Table 13 of Appendix D). This disparity stems from linguistic differences: English proper nouns require conversion into Chinese characters, posing challenges for maintaining consistency, whereas in German, they can be directly copied due to the shared alphabet. Despite this, a reasonable LTCR-1 improvement in En \Rightarrow De still demonstrates our method’s effectiveness. [The p-values of t-tests for DELTA vs Sentence/Context in translation quality are less than 0.05 for En \$\Leftrightarrow\$ Xx, whereas that in translation consistency are less than 0.05 in En \$\Leftrightarrow\$ Zh. We also test DELTA on the low-resource language pair, as demonstrated in Appendix F.](#)

Test Results on Guofeng The test results on Guofeng are illustrated in Table 3. Our approach achieves superior results across almost all metrics and backbone models, demonstrating its robustness to data of the novel domain. Notably, stronger models, such as GPT-4o-mini and Qwen2-72B-Instruct, achieve greater improvements in translation consistency and quality metrics. This suggests that the stronger the backbone models are, the more substantial the gains achieved by DELTA. The Guofeng test set poses particular challenges for maintaining translation consistency due to the prevalence of proper nouns (mainly names) in the source text. Nevertheless, DELTA demonstrates significant improvements in the relevant metrics, with an increase in LTCR-1 of up to 48.50 percentage points (85.50 vs. 37.00). This indicates that DELTA represents a strong tool for addressing translation inconsistency issues and holds great potential for novel translation, as it effectively reduces inconsistent noun translations, thereby minimizing potential confusion for readers.

6 ANALYSIS

Ablation Study Table 4 presents an ablation study in the En \Rightarrow Zh direction using GPT-3.5-Turbo-0125 as the backbone model. [For more detailed ablation studies of each memory component and the effects of hyper-parameters, please refer to Appendix E.](#) When provided with context sentences (Model 2), the model exhibits improved translation quality scores, but no significant enhancement in translation consistency is observed. The introduction of long-term memory contributes to more consistent translations (Model 3). Incorporating bilingual summaries (Model 6) led to an increase in COMET scores as well as consistency metrics, indicating that this component not only enhances translation quality but also reinforces translation consistency. When proper noun records are introduced (Model 7), a slight decrease in sCOMET and dCOMET scores is observed, likely due to the perturbation introduced by incorporating additional information. However,

Id	Setting	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
1	Sentence-level	83.78	6.55	80.27	88.78
2	1 + Short-Term Memory	84.50	6.68	77.89	87.41
3	2 + Long-Term Memory	84.54	6.67	79.23	89.44
4	3 + Source Summary	84.61	6.68	76.09	91.25
5	3 + Target Summary	84.70	6.72	82.14	92.86
6	3 + Bilingual Summary	84.72	6.74	82.49	93.60
7	6 + Record (DELTA)	84.70	6.72	86.44	95.25

Table 4: Ablation Study.

translation consistency improves significantly, with LTCR-1 increasing by 3.95 points and LTCR-1_f increasing by 1.65 points compared to Model 6. Moreover, it is evident that the bilingual summary has a superior impact on both translation quality and consistency compared to using a summary on either the source side or the target side alone. Among Models 4, 5, and 6, Bilingual Summary achieves the highest scores across all four metrics.

Consistency Distance One significant challenge in document-level translation is maintaining long-term consistency. To evaluate whether our approach addresses this challenge, we divide the sentence-wise distance between each proper noun’s translation and its first occurrence into several intervals. We then report the proportion of consistent translations in each interval in $En \Rightarrow Xx$ (upper) and $Xx \Rightarrow En$ (lower) in Figure 2. We observe that our approach almost outperforms the Sentence and Context methods in achieving proper noun translation consistency across all distance intervals. Notably, our approach excels when the distances exceed 50 sentences, yielding a larger proportion of consistent translations than the baseline methods. This demonstrates the effectiveness of our approach in enhancing long-context translation consistency.

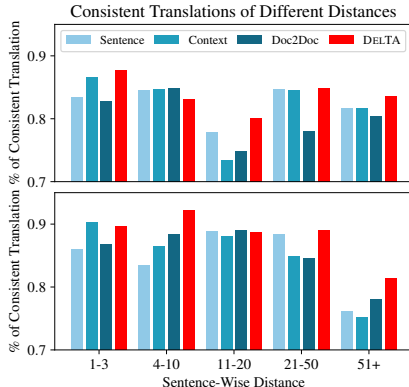


Figure 2: Proportions of consistent translations in different sentence-wise distances.

Pronoun & Context-Dependent Translation We follow Miculicich et al. (2018); Tan et al. (2019); Lyu et al. (2021) to evaluate the accuracy of pronoun translation (APT) of our system in $En \Rightarrow Zh$ using the reference-based metric proposed by Miculicich Werlen & Popescu-Belis (2017). We also evaluate our system on the first 1000 instances in the $En \Rightarrow De$ subset “mini.gender.opensubtitles” of a context-dependent translation benchmark called CTX-PRO Wicks & Post (2023). The results achieved by GPT-3.5-Turbo-0125 are demonstrated in Table 5 and Table 6. DELTA improves the performance of pronoun translation compared to the Sentence and Context baselines, and enhances translations where context information is explicitly needed. These results indicate that the multi-level memory in DELTA is beneficial to resolving coreference and discourse issues in the document.

Metric	Sentence	Context	Doc2Doc	DELTA
APT	59.96	60.84	56.11	61.07

Table 5: Evaluation results of pronoun translation accuracy (APT).

Metric	Sentence	DELTA
Generative Accuracy (%)	29.7	51.0

Table 6: Evaluation results of context-dependent translation.

DELTA as a Summarize Writer Our system employs an iterative approach to document summarization, where a summary writer generates partial summaries every 20 sentences, sequentially merging them with previous summaries to produce an updated version. To assess the effectiveness of this component, we conduct an experiment using the QMSum (Zhong et al., 2021) test set. QMSum is a benchmark for query-based multi-domain meeting summarization, where systems are required to generate summaries of the meeting transcripts in response to a specified query. In our experiment, the query is incorporated into the prompts for both summary generation and the merging process, enabling a query-based summarization approach.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

We compare our results with those of Lee et al. (2024), who paginate the document, generate summaries for each page, and then perform a lookup process on these summaries according to the query. The results are shown in Table 7. Our system’s segment-by-segment, query-oriented summary generation approach enables us to effectively locate relevant portions of the document and synthesize the information in these segments through the summary merging process. This demonstrates that the summary component of DELTA is also well-suited for general summarization tasks.

Memory Costs Our agent system consumes less memory than LLM-based Doc2Doc methods, such as those described by Wang et al. (2023b). In their research, multiple continuous sentences are translated in a single conversational turn, and the whole document is translated within a single chat box. However, this approach suffers from significant memory costs. As shown in Figure 3, we compared the memory usage of the Doc2Doc method with our agent-based online approach by utilizing Qwen2-72B-Instruction to translate a document in En \Rightarrow Zh on a device with 2 NVIDIA A800 80GB GPUs.

Our method demonstrates relatively slow memory growth as the number of processed sentences increases, primarily due to the increasing length of the summaries. In contrast, while the Doc2Doc method starts with lower memory consumption, its usage increases rapidly with the number of processed sentences. Memory consumption surpasses that of our method when the document length reaches 70 sentences, and it runs out of memory when the length reaches 490 sentences. This indicates that our approach is more memory-efficient and cost-effective for deployment on local devices.

7 CONCLUSION

In this paper, we start with analyzing two critical challenges for DocMT-LLMs, namely translation inconsistency and inaccuracy. To tackle these issues, we design an online document-level translation agent equipped with a multi-level memory component. This memory structure retrieves and stores key information to assist the document translation process, significantly enhancing translation consistency and quality. Notably, the effectiveness in maintaining proper noun translation consistency is particularly pronounced in novel translation, and our approach is still able to maintain consistency even when there is a large distance between the occurrences of a proper noun pair. The sentence-by-sentence online translation method avoids sentence omissions and reduces GPU memory consumption, in contrast to mainstream Doc2Doc approaches. Further analysis indicates that our framework is able to model discourse structures in the documents to improve pronoun translation and context-dependent translation accuracy, and the built-in summarizer component in our agent is also capable of the query-based summarization task.

LIMITATIONS

In this work, we present a framework for the DocMT agent, without prioritizing its inference efficiency. Given the complexity of DELTA’s inference process, LLMs are frequently invoked during document translation, leading to prolonged runtime. To address this issue, we identify several potential directions for improvement. First, by explicitly marking sentence boundaries with special boundary tags, we can enforce sentence-level alignment within generated paragraphs, allowing LLMs to process multiple sentences concurrently and thereby reduce invocation frequency. Second, employing more precise alignment tools and scripts to extract proper nouns and their translations, rather than relying on LLMs, can further enhance the efficiency of DELTA. Other components, such as Long-Term Retriever, could be implemented using a dense retriever rather than employing LLMs to identify related sentences. Finally, in the summary component, reducing the summary generation to a single step by directly merging sentences within the window into an overall summary can also decrease runtime. We consider these optimizations as future directions of our work.

System	ROUGE-L	Length
READAGENT	21.50	67.86
DELTA	23.60	82.28

Table 7: QMSum test results of ReadAgent and DELTA. “Length” denotes the word-wise length of the response.

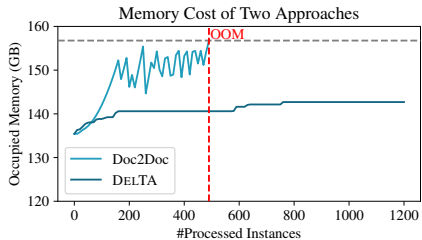


Figure 3: Memory cost of the Doc2Doc approach and our approach.

REFERENCES

- 540
541
542 Yasuhiro Akiba, Marcello Federico, Noriko Kando, Hiromi Nakaiwa, Michael Paul, and Jun'ichi
543 Tsujii. Overview of the IWSLT evaluation campaign. In *Proceedings of the First International*
544 *Workshop on Spoken Language Translation: Evaluation Campaign*, Kyoto, Japan, 2004. URL
545 <https://aclanthology.org/2004.iwslt-evaluation.1>.
- 546
547 Guangsheng Bao, Yue Zhang, Zhiyang Teng, Boxing Chen, and Weihua Luo. G-transformer for
548 document-level machine translation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli
549 (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*
550 *and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long*
551 *Papers)*, pp. 3442–3455, Online, 2021. Association for Computational Linguistics. doi: 10.18653/
552 v1/2021.acl-long.267. URL <https://aclanthology.org/2021.acl-long.267>.
- 553
554 Rachel Bawden, Rico Sennrich, Alexandra Birch, and Barry Haddow. Evaluating discourse phenom-
555 ena in neural machine translation. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceed-*
556 *ings of the 2018 Conference of the North American Chapter of the Association for Computational*
557 *Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1304–1313, New Or-
558 leans, Louisiana, 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1118.
URL <https://aclanthology.org/N18-1118>.
- 559
560 Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan,
561 Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. No language left behind: Scaling
562 human-centered machine translation. *ArXiv preprint*, abs/2207.04672, 2022. URL [https://](https://arxiv.org/abs/2207.04672)
563 arxiv.org/abs/2207.04672.
- 564
565 Zi-Yi Dou and Graham Neubig. Word alignment by fine-tuning embeddings on parallel corpora.
566 In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (eds.), *Proceedings of the 16th Conference*
567 *of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp.
568 2112–2128, Online, 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
eacl-main.181. URL <https://aclanthology.org/2021.eacl-main.181>.
- 569
570 Zhaopeng Feng, Yan Zhang, Hao Li, Wenqiang Liu, Jun Lang, Yang Feng, Jian Wu, and Zuozhu
571 Liu. Improving llm-based machine translation with systematic self-correction. *ArXiv preprint*,
572 abs/2402.16379, 2024. URL <https://arxiv.org/abs/2402.16379>.
- 573
574 Patrick Fernandes, Kayo Yin, Graham Neubig, and André F. T. Martins. Measuring and increasing
575 context usage in context-aware machine translation. In Chengqing Zong, Fei Xia, Wenjie Li, and
576 Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computa-*
577 *tional Linguistics and the 11th International Joint Conference on Natural Language Processing*
578 *(Volume 1: Long Papers)*, pp. 6467–6478, Online, 2021. Association for Computational Linguis-
579 tics. doi: 10.18653/v1/2021.acl-long.505. URL [https://aclanthology.org/2021.acl-long.](https://aclanthology.org/2021.acl-long.505)
505.
- 580
581 Xavier Garcia and Orhan Firat. Using natural language prompts for machine translation. *ArXiv*
582 *preprint*, abs/2202.11822, 2022. URL <https://arxiv.org/abs/2202.11822>.
- 583
584 Marjan Ghazvininejad, Hila Gonen, and Luke Zettlemoyer. Dictionary-based phrase-level prompt-
585 ing of large language models for machine translation. *ArXiv preprint*, abs/2302.07856, 2023.
URL <https://arxiv.org/abs/2302.07856>.
- 586
587 Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu,
588 Shuming Shi, and Xing Wang. Exploring human-like translation strategy with large language
589 models. *Transactions of the Association for Computational Linguistics*, 12:229–246, 2024. doi:
590 10.1162/tacl.a.00642. URL <https://aclanthology.org/2024.tacl-1.13>.
- 591
592 Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Mat-
593 sushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. How good are gpt models
at machine translation? a comprehensive evaluation. *ArXiv preprint*, abs/2302.09210, 2023. URL
<https://arxiv.org/abs/2302.09210>.

- 594 Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. Is
595 chatgpt a good translator? yes with gpt-4 as the engine. *ArXiv preprint*, abs/2301.08745, 2023.
596 URL <https://arxiv.org/abs/2301.08745>.
- 597 Marzena Karpinska and Mohit Iyyer. Large language models effectively leverage document-level
598 context for literary translation, but critical errors persist. In Philipp Koehn, Barry Haddow, Tom
599 Kocmi, and Christof Monz (eds.), *Proceedings of the Eighth Conference on Machine Translation*,
600 pp. 419–451, Singapore, 2023. Association for Computational Linguistics. doi: 10.18653/v1/
601 2023.wmt-1.41. URL <https://aclanthology.org/2023.wmt-1.41>.
- 602 Yunsu Kim, Duc Thanh Tran, and Hermann Ney. When and why is document-level context useful
603 in neural machine translation? In Andrei Popescu-Belis, Sharid Loáiciga, Christian Hardmeier,
604 and Deyi Xiong (eds.), *Proceedings of the Fourth Workshop on Discourse in Machine Translation*
605 (*DiscoMT 2019*), pp. 24–34, Hong Kong, China, 2019. Association for Computational Linguistics.
606 doi: 10.18653/v1/D19-6503. URL <https://aclanthology.org/D19-6503>.
- 607 Kuang-Huei Lee, Xinyun Chen, Hiroki Furuta, John Canny, and Ian Fischer. A human-inspired
608 reading agent with gist memory of very long contexts. *ArXiv preprint*, abs/2402.09727, 2024.
609 URL <https://arxiv.org/abs/2402.09727>.
- 610 Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. CAMEL:
611 communicative agents for "mind" exploration of large language model society. In Alice
612 Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine
613 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*
614 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December*
615 *10 - 16, 2023*, 2023a. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
616 a3621ee907def47c1b952ade25c67698-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/a3621ee907def47c1b952ade25c67698-Abstract-Conference.html).
- 617 Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang
618 Liu. Agent hospital: A simulacrum of hospital with evolvable medical agents. *ArXiv preprint*,
619 abs/2405.02957, 2024. URL <https://arxiv.org/abs/2405.02957>.
- 620 Yachao Li, Junhui Li, Jing Jiang, Shimin Tao, Hao Yang, and Min Zhang. P-transformer: Towards
621 better document-to-document neural machine translation. *IEEE/ACM Transactions on Audio,*
622 *Speech, and Language Processing*, 2023b. URL [https://ieeexplore.ieee.org/abstract/
623 document/10255243/](https://ieeexplore.ieee.org/abstract/document/10255243/).
- 624 Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng
625 Tu, and Shuming Shi. Encouraging divergent thinking in large language models through multi-
626 agent debate. *ArXiv preprint*, abs/2305.19118, 2023. URL [https://arxiv.org/abs/2305.
627 19118](https://arxiv.org/abs/2305.19118).
- 628 Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike
629 Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine trans-
630 lation. *Transactions of the Association for Computational Linguistics*, 8:726–742, 2020. doi:
631 10.1162/tacl.a.00343. URL <https://aclanthology.org/2020.tacl-1.47>.
- 632 Lorenzo Lupo, Marco Dinarelli, and Laurent Besacier. Focused concatenation for context-aware
633 neural machine translation. In Philipp Koehn, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen
634 Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus
635 Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck,
636 Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki
637 Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Marco
638 Turchi, and Marcos Zampieri (eds.), *Proceedings of the Seventh Conference on Machine Trans-*
639 *lation (WMT)*, pp. 830–842, Abu Dhabi, United Arab Emirates (Hybrid), 2022. Association for
640 Computational Linguistics. URL <https://aclanthology.org/2022.wmt-1.77>.
- 641 Xinglin Lyu, Junhui Li, Zhengxian Gong, and Min Zhang. Encouraging lexical translation consis-
642 tency for document-level neural machine translation. In Marie-Francine Moens, Xuanjing Huang,
643 Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical*
644 *Methods in Natural Language Processing*, pp. 3265–3277, Online and Punta Cana, Dominican
645 Republic, 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.
646 262. URL <https://aclanthology.org/2021.emnlp-main.262>.

- 648 Sameen Maruf and Gholamreza Haffari. Document context neural machine translation with memory
649 networks. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting*
650 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1275–1284, Mel-
651 bourne, Australia, 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1118.
652 URL <https://aclanthology.org/P18-1118>.
- 653 Sameen Maruf, Fahimeh Saleh, and Gholamreza Haffari. A survey on document-level neural ma-
654 chine translation: Methods and evaluation. *ACM Computing Surveys (CSUR)*, 54(2):1–36, 2021.
655 URL <https://dl.acm.org/doi/abs/10.1145/3441691>.
- 656
657 Lesly Miculicich, Dhananjay Ram, Nikolaos Pappas, and James Henderson. Document-level neural
658 machine translation with hierarchical attention networks. In Ellen Riloff, David Chiang, Julia
659 Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Meth-*
660 *ods in Natural Language Processing*, pp. 2947–2954, Brussels, Belgium, 2018. Association for
661 Computational Linguistics. doi: 10.18653/v1/D18-1325. URL <https://aclanthology.org/D18-1325>.
- 662
663 Lesly Miculicich Werlen and Andrei Popescu-Belis. Validation of an automatic metric for the
664 accuracy of pronoun translation (APT). In Bonnie Webber, Andrei Popescu-Belis, and Jörg
665 Tiedemann (eds.), *Proceedings of the Third Workshop on Discourse in Machine Translation*,
666 pp. 17–25, Copenhagen, Denmark, 2017. Association for Computational Linguistics. doi:
667 10.18653/v1/W17-4802. URL <https://aclanthology.org/W17-4802>.
- 668
669 OpenAI. Gpt-4 technical report, 2023. URL <https://arxiv.org/pdf/2303.08774>.
- 670
671 Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and
672 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings*
673 *of the 36th annual acm symposium on user interface software and technology*, pp. 1–22, 2023.
674 URL <https://dl.acm.org/doi/abs/10.1145/3586183.3606763>.
- 675
676 Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang,
677 and Dacheng Tao. Towards making the most of ChatGPT for machine translation. In Houda
678 Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Lin-*
679 *guistics: EMNLP 2023*, pp. 5622–5633, Singapore, 2023. Association for Computational Lin-
680 guistics. doi: 10.18653/v1/2023.findings-emnlp.373. URL <https://aclanthology.org/2023.findings-emnlp.373>.
- 681
682 Carolina Scarton and Lucia Specia. A quantitative analysis of discourse phenomena in ma-
683 chine translation. *Discours. Revue de linguistique, psycholinguistique et informatique. A jour-*
684 *nal of linguistics, psycholinguistics and computational linguistics*, (16), 2015. URL <https://journals.openedition.org/discours/9047>.
- 685
686 Sai Cheong Siu. Chatgpt and gpt-4 for professional translators: Exploring the potential of large
687 language models in translation. *Available at SSRN 4448091*, 2023. URL https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4448091.
- 688
689 Zewei Sun, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Lei Li.
690 Rethinking document-level neural machine translation. In Smaranda Muresan, Preslav Nakov,
691 and Aline Villavicencio (eds.), *Findings of the Association for Computational Linguistics: ACL*
692 *2022*, pp. 3537–3548, Dublin, Ireland, 2022. Association for Computational Linguistics. doi:
693 10.18653/v1/2022.findings-acl.279. URL <https://aclanthology.org/2022.findings-acl.279>.
- 694
695 Xin Tan, Longyin Zhang, Deyi Xiong, and Guodong Zhou. Hierarchical modeling of global con-
696 text for document-level neural machine translation. In Kentaro Inui, Jing Jiang, Vincent Ng,
697 and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural*
698 *Language Processing and the 9th International Joint Conference on Natural Language Process-*
699 *ing (EMNLP-IJCNLP)*, pp. 1576–1585, Hong Kong, China, 2019. Association for Computational
700 Linguistics. doi: 10.18653/v1/D19-1168. URL <https://aclanthology.org/D19-1168>.
- 701
702 Xin Tan, Longyin Zhang, and Guodong Zhou. Coupling context modeling with zero pronoun re-
covering for document-level natural language generation. In Marie-Francine Moens, Xuanjing

- 702 Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on*
703 *Empirical Methods in Natural Language Processing*, pp. 2530–2540, Online and Punta Cana,
704 Dominican Republic, 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
705 emnlp-main.197. URL <https://aclanthology.org/2021.emnlp-main.197>.
706
- 707 Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. Embarrassingly easy
708 document-level MT metrics: How to convert any pretrained metric into a document-level metric.
709 In Philipp Koehn, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-
710 jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Ro-
711 man Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom
712 Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa,
713 Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri
714 (eds.), *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pp. 118–128, Abu
715 Dhabi, United Arab Emirates (Hybrid), 2022. Association for Computational Linguistics. URL
716 <https://aclanthology.org/2022.wmt-1.6>.
- 717 Bing Wang, Xinnian Liang, Jian Yang, Hui Huang, Shuangzhi Wu, Peihao Wu, Lu Lu, Zejun Ma,
718 and Zhoujun Li. Enhancing large language model with self-controlled memory framework. *ArXiv*
719 *preprint*, abs/2304.13343, 2023a. URL <https://arxiv.org/abs/2304.13343>.
- 720 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Ji-
721 akai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous
722 agents. *Frontiers of Computer Science*, 18(6):186345, 2024a. URL <https://link.springer.com/article/10.1007/s11704-024-40231-1>.
- 723
724 Longyue Wang, Zhaopeng Tu, Andy Way, and Qun Liu. Exploiting cross-sentence context for neural
725 machine translation. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel (eds.), *Proceedings*
726 *of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2826–2831,
727 Copenhagen, Denmark, 2017. Association for Computational Linguistics. doi: 10.18653/v1/
728 D17-1301. URL <https://aclanthology.org/D17-1301>.
- 729
730 Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu.
731 Document-level machine translation with large language models. In Houda Bouamor, Juan Pino,
732 and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural*
733 *Language Processing*, pp. 16646–16661, Singapore, 2023b. Association for Computational Lin-
734 guistics. doi: 10.18653/v1/2023.emnlp-main.1036. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.emnlp-main.1036)
735 emnlp-main.1036.
- 736
737 Longyue Wang, Zhaopeng Tu, Yan Gu, Siyou Liu, Dian Yu, Qingsong Ma, Chenyang Lyu, Lit-
738 ing Zhou, Chao-Hong Liu, Yufeng Ma, Weiyu Chen, Yvette Graham, Bonnie Webber, Philipp
739 Koehn, Andy Way, Yulin Yuan, and Shuming Shi. Findings of the WMT 2023 shared task on
740 discourse-level literary translation: A fresh orb in the cosmos of LLMs. In Philipp Koehn, Barry
741 Haddow, Tom Kocmi, and Christof Monz (eds.), *Proceedings of the Eighth Conference on Ma-*
742 *chine Translation*, pp. 55–67, Singapore, 2023c. Association for Computational Linguistics. doi:
743 10.18653/v1/2023.wmt-1.3. URL <https://aclanthology.org/2023.wmt-1.3>.
- 744
745 Longyue Wang, Siyou Liu, Minghao Wu, Wenxiang Jiao, Xing Wang, Jiahao Xu, Zhaopeng Tu,
746 Liting Zhou, Yan Gu, Weiyu Chen, Philipp Koehn, Andy Way, and Yulin Yuan. Findings of
747 the wmt 2024 shared task on discourse-level literary translation. In *Proceedings of the Ninth*
748 *Conference on Machine Translation*, 2024b.
- 749
750 Yutong Wang, Jiali Zeng, Xuebo Liu, Fandong Meng, Jie Zhou, and Min Zhang. TasTe: Teach-
751 ing large language models to translate through self-reflection. In Lun-Wei Ku, Andre Martins,
752 and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-*
753 *putational Linguistics (Volume 1: Long Papers)*, pp. 6144–6158, Bangkok, Thailand, August
754 2024c. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.333. URL
755 <https://aclanthology.org/2024.acl-long.333>.
- 756
757 Rachel Wicks and Matt Post. Identifying context-dependent translations for evaluation set pro-
758 duction. In Philipp Koehn, Barry Haddow, Tom Kocmi, and Christof Monz (eds.), *Proceed-*
759 *ings of the Eighth Conference on Machine Translation*, pp. 452–467, Singapore, December

- 756 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.wmt-1.42. URL
757 <https://aclanthology.org/2023.wmt-1.42>.
758
- 759 Minghao Wu, Thuy-Trang Vu, Lizhen Qu, George Foster, and Gholamreza Haffari. Adapting large
760 language models for document-level machine translation. *ArXiv preprint*, abs/2401.06468, 2024a.
761 URL <https://arxiv.org/abs/2401.06468>.
- 762 Minghao Wu, Yulin Yuan, Gholamreza Haffari, and Longyue Wang. (perhaps) beyond human
763 translation: Harnessing multi-agent collaboration for translating ultra-long literary texts. *ArXiv*
764 *preprint*, abs/2405.11804, 2024b. URL <https://arxiv.org/abs/2405.11804>.
765
- 766 Yangjian Wu and Gang Hu. Exploring prompt engineering with GPT language models for
767 document-level machine translation: Insights and findings. In Philipp Koehn, Barry Haddow,
768 Tom Kocmi, and Christof Monz (eds.), *Proceedings of the Eighth Conference on Machine*
769 *Translation*, pp. 166–169, Singapore, 2023. Association for Computational Linguistics. doi:
770 10.18653/v1/2023.wmt-1.15. URL <https://aclanthology.org/2023.wmt-1.15>.
- 771 Wenda Xu, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, Biao Zhang, Zhongtao Liu,
772 William Yang Wang, Lei Li, and Markus Freitag. LLMRefine: Pinpointing and refining large
773 language models via fine-grained actionable feedback. In Kevin Duh, Helena Gomez, and
774 Steven Bethard (eds.), *Findings of the Association for Computational Linguistics: NAACL 2024*,
775 pp. 1429–1445, Mexico City, Mexico, 2024. Association for Computational Linguistics. URL
776 <https://aclanthology.org/2024.findings-naacl.92>.
- 777 Jiali Zeng, Fandong Meng, Yongjing Yin, and Jie Zhou. Tim: Teaching large language models to
778 translate with comparison. *ArXiv preprint*, abs/2307.04408, 2023. URL <https://arxiv.org/abs/2307.04408>.
779
- 780 Biao Zhang, Barry Haddow, and Alexandra Birch. Prompting large language model for machine
781 translation: A case study. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engel-
782 hardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning,*
783 *ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine*
784 *Learning Research*, pp. 41092–41110. PMLR, 2023. URL <https://proceedings.mlr.press/v202/zhang23m.html>.
785
- 786 Jiacheng Zhang, Huanbo Luan, Maosong Sun, Feifei Zhai, Jingfang Xu, Min Zhang, and Yang
787 Liu. Improving the transformer translation model with document-level context. In Ellen Riloff,
788 David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of the 2018 Con-*
789 *ference on Empirical Methods in Natural Language Processing*, pp. 533–542, Brussels, Bel-
790 gium, 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1049. URL
791 <https://aclanthology.org/D18-1049>.
792
- 793 Pei Zhang, Boxing Chen, Niyu Ge, and Kai Fan. Long-short term masking transformer: A sim-
794 ple but effective baseline for document-level neural machine translation. In Bonnie Webber,
795 Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empir-*
796 *ical Methods in Natural Language Processing (EMNLP)*, pp. 1081–1087, Online, 2020. Asso-
797 ciation for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.81. URL <https://aclanthology.org/2020.emnlp-main.81>.
798
- 799 Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadal-
800 lah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir Radev. QMSum: A new benchmark
801 for query-based multi-domain meeting summarization. In Kristina Toutanova, Anna Rumshisky,
802 Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy
803 Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North Amer-*
804 *ican Chapter of the Association for Computational Linguistics: Human Language Technologies*,
805 pp. 5905–5921, Online, 2021. Association for Computational Linguistics. doi: 10.18653/v1/
806 2021.naacl-main.472. URL <https://aclanthology.org/2021.naacl-main.472>.
807
808
809

810	Proper Noun Translation Inconsistency	
811	SRC	It's a story about this woman, Natalia Rybczynski .
812	HYP	这是关于这个女人 娜塔莉亚·雷布琴斯基 的故事。
813	SRC	Natalia Rybczynski : Yeah, I had someone call me "Dr. Dead Things."
814	HYP	娜塔莉亚·丽琴斯基 : 是的, 有人叫我“死物博士”。
816	Undertranslation	
817	SRC	But here's the truth.//Here's the epiphany that I had that changed my thinking.//From 1970 until today, the percentage of the world's population living in starvation levels...
818	HYP	但事实是, // *Missing translation* //自1970年至今, 生活在饥饿水平、每天生活在一美元以下 (当然要根据通货膨胀调整) 的全球人口比例下降了80%。
820	Low Translation Quality	
821	SRC	And we make decisions about where to live, who to marry and even who our friends are going to be, based on what we already believe.
822	REF	我们做的各种决定, 选择生活在何处, 与谁结婚甚至和谁交朋友, 都只基于我们已有的信念。
823	HYP ₁	我们根据自己已有的信念来做决定, 包括选择居住的地方, 结婚对象, 甚至决定谁会成为我们的朋友。
824	HYP ₅₀	我们决定居住地、婚姻对象, 甚至我们的朋友根据我们已经相信的事情。

Table 8: Demonstrations of proper noun translation inconsistency, undertranslation, and low translation quality issues encountered during document-level translation. In the first part, texts highlighted in red represent the same source proper noun with two different translations in the target document. In the second part, “//” indicates sentence boundaries in the document, illustrating that the translation of the second source sentence is absent from the target document. In the third part, “HYP_n” represents the hypothesis generated by the LLM using a translation window of size n .

A EXAMPLES OF TRANSLATION INCONSISTENCY AND INACCURACY

Examples of translation inconsistency, undertranslation, and low translation quality issues are presented in Table 8. In the first instance, the name “Natalia Rybczynski”, which appears twice in the source document, is translated into two different forms: “娜塔莉亚·雷布琴斯基” and “娜塔莉亚·丽琴斯基”. This variation leads to a notable inconsistency in the translation. In the second instance, the second sentence in the source document is omitted. Its corresponding translation is absent in the target document. This exemplifies an undertranslation issue in the document translation task. In the third instance, the translation produced by the LLM with a translation window of 50 sentences deviates significantly from the customary word order in Chinese compared to the sentence-level approach (using a window size of 1).

B DETAILED STATISTICS OF THE DATASETS IN OUR EXPERIMENTS

We conduct our experiments on two test sets: IWSLT2017 and Guofeng. The IWSLT2017 test set consists of parallel documents sourced from TED talks, covering 12 language pairs. Our experiments are conducted on eight of these pairs, En \leftrightarrow Zh, De, Fr, and Ja. Each language pair contains 10 to 12 sentence-level aligned parallel documents, totaling approximately 1.5K sentences, with an average of around 120 sentences per document. Detailed statistics for each language pair are shown in the first block of Table 9.

The Guofeng Webnovel corpus is a high-quality and discourse-level corpus of web fiction, and we conduct our experiments on the Guofeng V1 TEST_2 set, which is designed in the Zh \Rightarrow En language pair. The second block of Table 9 illustrated the statistics of the test sets.

C PROMPT TEMPLATES FOR LLM-BASED COMPONENTS IN DELTA

This part details the prompts used for each module of DELTA. The prompt template for the Proper Noun Extractor is depicted in Figure 4. We use a prompt in the few-shot style to ensure accurate and formatted outputs. The prompt templates of the source and target summary writers are shown in

Dataset	Language	S	D	S / D
IWSLT2017	Zh \Leftrightarrow En	1459	12	122
	De \Leftrightarrow En	1138	10	114
	Fr \Leftrightarrow En	1455	12	121
	Ja \Leftrightarrow En	1452	12	121
Guofeng VI TEST 2	Zh \Rightarrow En	857	12	71

Table 9: Statistics of the test sets used in our experiments. “|S|” represents the number of sentences in each test set, “|D|” represents the number of documents, and |S|/|D| represents the average number of sentences per document.

Id	Setting	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
1	Sentence-level	83.78	6.55	80.27	88.78
2	1 + Short-Term Memory	84.50	6.68	77.89	87.41
3	1 + Long-Term Memory	84.48	6.69	78.77	88.01
4	1 + Record	84.11	6.60	81.33	89.33
5	1 + Summary	84.51	6.73	79.73	90.70
6	2 + Long-Term Memory	84.54	6.67	79.23	89.44
7	2 + Record	84.45	6.70	82.37	92.54
8	3 + Source Summary	84.61	6.68	76.09	91.25
9	3 + Target Summary	84.70	6.72	82.14	92.86
10	3 + Bilingual Summary	84.72	6.74	82.49	93.60
11	10 + Record (DELTA)	84.70	6.72	86.44	95.25

Table 10: More detailed results of the ablation study.

Figure 5 and Figure 6, respectively. Note that the prompts for the summary components are written in their respective source or target language to avoid off-target issues. The prompt template for the Memory Retriever is shown in Figure 7. The prompt template for the Document Translator is illustrated in Figure 8. All the retrieved information from the multi-level memory components is formatted into this prompt to assist the translation process.

D DETAILED RESULTS OF THE MAIN EXPERIMENT

The scores for the En \Rightarrow Zh, De, Fr, Ja translation directions are presented in Table 13, while the scores for the Zh, De, Fr, Ja \Rightarrow En are shown in Table 14.

DELTA achieves improvements in both translation consistency, as indicated by the LTCR-1 and LTCR-1_f metrics, and translation quality, as indicated by the sCOMET and dCOMET metrics, across most translation directions compared to several baselines. The Qwen models show significant enhancements in sCOMET and dCOMET scores, demonstrating that our approach provides strong reinforcement for the document translation quality of these open-source models.

The quality improvements are most pronounced in the Ja \Leftrightarrow En directions across all backbone models, likely due to their modest baseline capabilities for these language pairs, which our approach effectively enhances. Translation consistency in the Zh \Leftrightarrow En directions benefits most from our approach, as the distinct character sets of Chinese and English pose challenges in maintaining proper noun consistency, highlighting the effectiveness of our method.

When applying GPT-3.5-Turbo as the backbone model, DELTA outperforms translation-specialized baselines, such as NLLB-3.3B and GoogleTrans, across most languages, demonstrating the promising capabilities of LLM-based autonomous agents in document translation.

E DETAILED RESULTS OF THE ABLATION STUDY

Short. Window Size	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
1	84.51	6.71	87.42	95.36
3	84.70	6.72	86.44	95.25
5	84.65	6.74	87.42	95.03
Long. Window Size	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
10	84.61	6.74	86.67	95.67
20	84.70	6.72	86.44	95.25
30	84.71	6.73	84.25	94.86
Growing	84.68	6.71	85.57	93.96

Table 11: Effect of the Short-Term and Long-Term window size.

System	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
Sentence	87.33	7.04	78.65	96.07
Context	87.87	7.19	78.09	96.07
DELTA	87.96	7.20	82.68	96.65

Table 12: Evaluation results on the Lt \Rightarrow En low-resource test set.

Effect of Each Memory Component The ablation results for each individual memory component are presented in Table 10. The backbone model used in our experiments is GPT-3.5-Turbo, evaluated on the IWSLT2017 En \Rightarrow Zh test set. We can observe that 1) Short-Term Memory enhances translation quality but negatively impacts consistency due to disruptions caused by the short-span context it introduces (Model 2). 2) Similarly, Long-Term Memory improves translation quality but also leads to a decline in consistency (Model 3). However, these two modules complement each other, as their integration leverages information from different spans, resulting in further quality enhancements while mitigating consistency issues (Model 6). 3) Proper Noun Records contribute to both translation quality and consistency, as improved accuracy in proper noun translation directly enhances overall quality (Model 4). 4) Bilingual Summary also positively impacts translation quality by introducing the document’s main idea, which guides the generation of target sentences more effectively (Model 5).

Each module independently contributes to performance improvements. However, none surpass the combination of all modules working together. We also conducted an additional comparative test between DelTA and the context method augmented with Proper Noun Records. The results, presented as Model 7 in Table 10, reveal a large performance gap favoring DelTA. We attribute this to the DelTA architecture’s ability to incorporate information across varying granularities and scales, enabling a synergistic enhancement of both translation quality and consistency within our framework.

Effect of Hyper-Parameters We perform experiments to evaluate the impact of different window sizes for Short-Term and Long-Term Memory (i.e. k and l introduced in §4). The results are presented in Table 11. Increasing the window size of Short-Term Memory results in a slight improvement in consistency but incurs a substantial additional computational cost, as two sentence pairs extend each sentence’s translation prompt. Given the trade-off between computational cost, translation quality, and consistency, we opted to employ a Short-Term window size of 3. Further increasing the Long-Term Memory window does not yield significant benefits. To address the idea that the model might require longer context windows in the later stages of document translation, we also experimented with a dynamic window setting, where the window size increases by one sentence pair for every eight sentences translated. However, this approach does not lead to performance improvements. We attribute this to the Bilingual Summary component, which effectively captures key information from more distant contexts. The iterative process of summary generation and merging inherently functions as a form of context window growth, rendering additional adjustments to the Long-Term Memory window unnecessary.

F PERFORMANCE ON LOW-RESOURCE LANGUAGES

We evaluate how well DELTA scales to low-resource languages. Considering that the spaCy NLP tool used for evaluation does not support many low-resource languages, we select Lt \Rightarrow En as the language pair to test. We randomly sample a document (with 556 sentences) from Europarl v9 training set¹³ as our evaluation set, and employ GPT-3.5-Turbo as the backbone model. The results are shown in Table 12, indicating that DELTA also performs well on the low-resource language pair.

Prompt for Proper Noun Extractor

You are an English-Chinese bilingual expert. Given an English source sentence with its Chinese translation, you need to annotate all the proper nouns in the English source sentence and their corresponding translations in the Chinese translation sentence. Here are some examples for you:

Example 1:

<English source> NASA’s Kepler mission has discovered thousands of potential planets around other stars, indicating that Earth is but one of billions of planets in our galaxy.

<Chinese translation> 美国国家航空航天局的开普勒任务已经发现了围绕着其他恒星的数千颗潜在的行星，这也表明了地球只是银河系中数十亿行星中的一颗。

<Proper nouns> “NASA” - “美国国家航空航天局”, “Kepler” - “开普勒”, “Earth” - “地球”

Example 2:

<English source> I had just driven home, it was around midnight in the dead of Montreal winter, I had been visiting my friend, Jeff, across town, and the thermometer on the front porch read minus 40 degrees – and don’t bother asking if that’s Celsius or Fahrenheit, minus 40 is where the two scales meet – it was very cold.

<Chinese translation> 我开车回到家，在Montreal的寒冬，大约午夜时分，我开车从城镇一边到另一边，去看望我的朋友杰夫，门廊上的温度计显示零下40度——不需要知道是摄氏度还是华氏度，到了零下40度，两个温度显示都一样——天气非常冷。

<Proper nouns> “Montreal” - “N/A”, “Jeff” - “杰夫”, “Celsius” - “摄氏度”, “Fahrenheit” - “华氏度”

Example 3:

<English source> To make the case to the National Health Service that more resources were needed for autistic children and their families, Lorna and her colleague Judith Gould decided to do something that should have been done 30 years earlier.

<Chinese translation> 为了向国家医疗保健系统证明，自闭症儿童和他们的家庭需要更多的资源，Lorna和她的同事朱迪思·古尔德决定去做一些三十年前就应该被完成的事情。

<Proper nouns> “National Health Service” - “国家医疗保健系统”, “Lorna” - “N/A”, “Judith Gould” - “朱迪思·古尔德”

If there isn’t any proper noun in the sentence, just answer with “N/A”. Now annotate all the proper nouns in the following sentence pair:

<English source> {SOURCE_SENTENCE}

<Chinese translation> {TARGET_SENTENCE}

<Proper nouns>

Figure 4: Prompt template for the Proper Noun Extractor. We provide several few-shot exemplars preceding the current input. This template is designed for the En \Rightarrow Zh translation direction. For other translation directions, adjust the corresponding content to match the specific languages.

¹³<https://www.statmt.org/europarl/v9/training/>

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Prompt for Source Summary Writer (Segment Summary Generation)

Below is a paragraph. Please provide a summary of this paragraph, including the main contents of these sentences, the overall domain, style and tone of it, while preserving key information as much as possible.
Paragraph: {SOURCE_SEGMENT}
Summary:

Prompt Template for Source Summary Writer (Summary Merging)

Below are the summaries of two adjacent paragraphs. Please merge them into a single summary, retaining as much key information as possible and ensuring that information about the domain, style, and tone are preserved.
Summary 1: {SUMMARY_1}
Summary 2: {SUMMARY_2}
Merged summary:

Figure 5: Prompt template for Source Summary Writer.

Prompt Template for Target Summary Writer (Segment Summary Generation)

Below is a paragraph. Please provide a summary of this paragraph, while preserving key information as much as possible.
Paragraph: {SOURCE_SEGMENT}
Summary:

Prompt Template for Target Summary Writer (Summary Merging)

Below are the summaries of two adjacent paragraphs. Please merge them into a single summary, retaining as much key information as possible.
Summary 1: {SUMMARY_1}
Summary 2: {SUMMARY_2}
Merged summary:

Figure 6: Prompt template for Target Summary Writer. We write the prompt in the target language to better align with the agent profile of the monolingual summary writer and reduce the off-target issues. For demonstration purposes, the prompts provided here are written in English.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

Prompt Template for Long-Term Memory Retriever

You are a linguistic expert. Given a list of sentences and a query, your task is to find the {TOP_NUM} sentences in the list that are most relevant to the request.

Sentence list:
 {SRC1}
 {SRC2}
 ...

Query:
 {QUERY}

Note that you should only respond with a list containing the numbers of these {TOP_NUM} sentences. For example, if you choose sentences 15, 16, and 19 as your answer, your response should be “[15, 16, 19]”.

Figure 7: Prompt template for Long-Term Memory Retriever.

Prompt Template for Document Translator

You are an {SRC_LANG}-{TGT_LANG} bilingual expert, translating a very long {SRC_LANG} document. Given the summary of the preceding text in both {SRC_LANG} and {TGT_LANG}, the historical translation of some proper nouns, source and translation texts preceding the current sentence, as well as some relevant translation instances from the preceding text, translate the current {SRC_LANG} source sentence into {TGT_LANG}. Please ensure that the translations of proper nouns in the source sentence are consistent with their historical translation, and the translation style remains consistent as well.

Summaries:
 <{SRC_LANG} summary> {SRC_SUMMARY}
 <{TGT_LANG} summary> {TGT_SUMMARY}

Historical translations of proper nouns:
 {HISTORY}

Preceding texts:
 <{SRC_LANG} text> {SRC_CONTEXT}
 <{TGT_LANG} text> {TGT_CONTEXT}

Relevant instances:
 {RELEVANT_INSTANCES}

Now translate the following {SRC_LANG} source sentence to {TGT_LANG}.

<{SRC_LANG} source> {SOURCE}
 <{TGT_LANG} translation>

Figure 8: Prompt template for Document Translator. Involved proper nouns and their corresponding translations are formatted into the “HISTORY” field. Source and target sentences from Short-Term Memory are concatenated and formatted into the “SRC_CONTEXT” and “TGT_CONTEXT” fields, respectively. Retrieved instances from Long-Term Memory are formatted into the “RELEVANT_INSTANCES” field as source-target pairs.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

System	sCOMET	dCOMET	LTCR-1	LTCR-1 _f	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
En ⇒ Zh				En ⇒ De				
NLLB	76.81	6.20	71.68	84.96	84.15	6.64	91.76	99.61
GOOGLE	78.46	5.76	89.45	92.36	80.23	5.78	93.55	99.19

GPT-3.5-Turbo								
Sentence	83.78	6.55	80.27	88.78	84.97	6.71	92.06	98.81
Context	84.50	6.68	77.89	87.41	85.12	6.74	93.70	99.21
Doc2Doc	–	6.29	82.04	94.29	–	6.81	88.89	97.94
Ours	84.70	6.72	86.44	95.25	85.37	6.78	93.46	99.23

GPT-4o-mini								
Sentence	82.13	6.43	78.04	91.89	81.41	6.39	90.70	98.84
Context	84.36	6.68	78.95	93.42	84.83	6.70	92.66	99.61
Doc2Doc	–	6.60	82.33	88.35	–	6.90	91.05	99.22
Ours	84.94	6.81	85.52	91.72	85.47	6.79	92.19	100.0

Qwen-7B-Instruct								
Sentence	83.05	6.51	77.78	83.84	76.24	5.54	82.11	90.65
Context	83.51	6.67	77.29	87.12	76.67	5.56	85.66	92.45
Doc2Doc	–	6.16	81.85	91.11	–	5.25	88.60	97.93
Ours	83.98	6.70	80.13	90.55	76.84	5.56	87.40	96.75

Qwen-72B-Instruct								
Sentence	77.76	5.88	73.13	83.96	78.98	6.13	92.18	98.35
Context	81.69	6.24	78.01	86.17	80.96	6.43	89.64	96.81
Doc2Doc	–	6.22	74.59	78.38	–	6.66	88.28	98.44
Ours	84.76	6.70	81.56	89.72	84.29	6.70	90.48	98.41

En ⇒ Fr				En ⇒ Ja				
NLLB	85.35	6.23	87.84	89.86	82.12	6.37	46.94	53.06
GOOGLE	82.21	5.53	88.61	91.46	80.74	6.23	53.92	55.88

GPT-3.5-Turbo								
Sentence	85.84	6.18	83.55	88.49	84.61	6.89	52.34	55.14
Context	86.49	6.27	83.06	89.25	85.50	7.09	54.72	56.60
Doc2Doc	–	6.28	92.28	94.63	–	7.10	53.26	58.70
Ours	86.48	6.30	88.96	94.16	85.76	7.13	62.96	66.67

GPT-4o-mini								
Sentence	80.79	5.82	88.89	90.91	81.72	6.74	56.73	58.65
Context	85.10	6.14	90.52	92.81	84.84	7.09	57.89	62.11
Doc2Doc	–	6.24	91.00	94.00	–	7.25	57.78	60.00
Ours	86.38	6.28	90.94	93.85	86.61	7.32	58.54	59.76

Qwen-7B-Instruct								
Sentence	80.61	5.47	82.53	84.93	80.21	6.30	53.21	58.72
Context	81.31	5.54	89.00	90.72	81.85	6.53	66.41	71.09
Doc2Doc	–	5.29	87.88	93.56	–	6.63	50.96	55.77
Ours	81.02	5.46	88.66	91.41	82.23	6.57	64.18	72.39

Qwen-72B-Instruct								
Sentence	81.02	5.75	90.00	92.33	76.34	6.11	62.86	65.71
Context	84.03	6.06	87.22	89.46	76.48	6.14	61.68	69.16
Doc2Doc	–	6.25	89.25	91.86	–	6.66	42.20	45.87
Ours	85.76	6.28	92.08	94.39	85.13	6.94	62.50	70.83

Table 13: Detailed results of our experiments in En ⇒ Xx directions.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

System	sCOMET	dCOMET	LTCR-1	LTCR-1 _f	sCOMET	dCOMET	LTCR-1	LTCR-1 _f
Zh ⇒ En				De ⇒ En				
NLLB	82.14	7.01	75.31	88.27	85.63	7.23	95.98	98.85
GOOGLE	78.06	5.68	72.89	86.14	82.31	6.47	96.53	98.27

GPT-3.5-Turbo								
Sentence	83.34	7.17	73.99	86.71	85.92	7.26	98.88	100.0
Context	83.88	7.29	76.92	90.53	86.10	7.30	98.30	100.0
Doc2Doc	–	7.08	76.77	88.39	–	7.16	98.24	98.82
Ours	83.88	7.30	80.00	93.53	86.14	7.30	98.33	100.0

GPT-4o-mini								
Sentence	83.55	7.24	71.93	84.80	85.11	7.17	98.20	100.0
Context	83.96	7.35	78.24	91.18	86.12	7.27	98.88	100.0
Doc2Doc	–	7.15	79.62	92.36	–	7.19	95.24	97.62
Ours	84.10	7.47	79.41	94.71	86.61	7.31	98.32	100.0

Qwen-7B-Instruct								
Sentence	79.44	6.76	71.17	87.73	77.20	6.75	93.25	95.09
Context	82.62	7.04	69.70	83.64	84.52	7.08	98.80	99.40
Doc2Doc	–	6.53	82.79	92.62	–	6.88	98.67	99.33
Ours	82.83	7.09	76.47	92.35	84.61	7.05	98.25	100.0

Qwen-72B-Instruct								
Sentence	80.35	6.95	73.49	84.34	81.17	6.94	97.19	100.0
Context	80.11	6.93	74.25	86.23	84.81	7.24	98.31	99.44
Doc2Doc	–	6.94	68.21	80.13	–	7.08	97.59	98.19
Ours	84.51	7.40	83.93	94.05	86.17	7.34	98.29	100.0
Fr ⇒ En				Ja ⇒ En				
NLLB	87.59	6.79	93.56	97.42	81.02	6.90	51.27	78.48
GOOGLE	84.64	6.19	95.63	96.83	75.67	5.50	60.67	82.00

GPT-3.5-Turbo								
Sentence	87.60	6.78	94.94	97.89	81.00	6.98	60.12	82.82
Context	88.03	6.84	94.96	97.90	81.85	7.17	69.94	92.64
Doc2Doc	–	6.78	94.78	97.39	–	6.80	70.90	87.31
Ours	88.02	6.86	96.17	98.30	81.76	7.13	71.60	93.21

GPT-4o-mini								
Sentence	87.32	6.77	94.42	97.85	80.04	6.76	61.11	82.72
Context	87.72	6.81	94.85	97.42	82.00	7.17	65.62	88.75
Doc2Doc	–	6.83	94.42	98.28	–	6.86	64.71	85.29
Ours	88.13	6.90	96.58	98.72	82.20	7.29	66.67	90.12

Qwen-7B-Instruct								
Sentence	81.21	6.27	87.90	95.56	70.56	6.15	53.25	73.38
Context	86.55	6.63	94.19	97.51	78.68	6.59	63.23	89.68
Doc2Doc	–	6.57	87.00	97.76	–	6.39	71.67	85.00
Ours	86.40	6.56	91.20	92.80	79.60	6.66	62.26	88.05

Qwen-72B-Instruct								
Sentence	82.67	6.38	95.44	98.34	77.94	6.63	62.89	85.53
Context	87.02	6.76	93.25	97.89	81.13	7.02	65.64	85.28
Doc2Doc	–	6.73	94.67	97.78	–	6.74	71.53	86.86
Ours	88.01	6.87	95.78	98.73	82.06	7.24	68.12	93.12

Table 14: Detailed results of our experiments in Xx ⇒ En directions.