

TRANSLLLAMA: LLM-BASED SIMULTANEOUS TRANSLATION SYSTEM

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Decoder-only large language models (LLMs) have recently demonstrated impressive capabilities in text generation and reasoning. Nonetheless, they have limited
 2 applications in simultaneous machine translation (SiMT), currently dominated by
 3 encoder-decoder transformers. This study demonstrates that, after fine-tuning on a
 4 small dataset comprising causally aligned source and target sentence pairs, a pre-
 5 trained open-source LLM can control input segmentation directly by generating a
 6 special "wait" token. This obviates the need for a separate policy and enables the
 7 LLM to perform English-German and English-Russian SiMT tasks with BLEU
 8 scores that are comparable to those of specific state-of-the-art baselines. We also
 9 evaluated closed-source models such as GPT-4, which displayed encouraging re-
 10 sults in performing the SiMT task without prior training (zero-shot), indicating a
 11 promising avenue for enhancing future SiMT systems.
 12

13 1 INTRODUCTION

14 Unlike conventional sequential translation, in which the target text is produced after the end of
 15 the corresponding source sentence (or long phrase), in simultaneous machine translation (SiMT)
 16 the target text is produced with minimal delay, aiming for the best listener experience expected
 17 from professional conference interpreters. While recent years have seen tremendous progress in
 18 sentence-based machine translation, mainstream adoption of SiMT systems requires solving a range
 19 of technical problems. Perhaps the most important of them is that, much like human conference in-
 20 terpreters, SiMT systems must make optimal decisions about *when* (rather than *how*) to translate. In
 21 particular, naively translating each source word immediately results in compromised target quality,
 22 given that the meaning of a source word often makes sense only in the context of later words. And
 23 while waiting until the end of a sentence might seem a viable solution, in practice it would introduce
 24 unacceptable delays between the source and target message. Consequently, the development of an
 25 effective SiMT system necessitates striking a balance between these two opposite scenarios.

26 Existing approaches to maintaining an optimal quality-latency tradeoff in SiMT, conventionally
 27 called *policies*, fall into two broad categories: fixed and adaptive. The policy's role is to signal
 28 to a separately trained translation model *when* to produce a partial translation (aka WRITE action
 29 (Gu et al., 2017)) based of the partial input; at other times the input, which represents either text
 30 chunks from an upstream ASR system (in cascade SiMT systems) or speech embeddings (in end-
 31 to-end systems), is just read in (READ action). While with a fixed policy (Dalvi et al., 2018; Ma
 32 et al., 2019a; Elbayad et al., 2020; Zhang & Feng, 2021), the decision to output translation is based
 33 on a simple heuristic, an adaptive policy (Arivazhagan et al., 2019; Ma et al., 2019b; Zhang &
 34 Feng, 2022) can be implemented as a separately trained model, for example an agent trained with
 35 reinforcement learning (RL) (Gu et al., 2017; Satija & Pineau, 2016).

36 To the best of our knowledge, state-of-the-art SiMT systems use encoder-decoder transformer ar-
 37 chitectures in a sequence-to-sequence paradigm. However, as of writing this paper the largest – and
 38 generally most expressive – language models are causal decoder-only architectures. We wanted to
 39 explore the utility of such models for SiMT tasks, focusing on the English-German and English-
 40 Russian language pairs, and specifically if they can be harnessed with minimal engineering effort.

41 Inspired by the recent success of LLMs – in particular their agential capabilities (Nascimento et al.,
 42 2023; Wang et al., 2023a;c) – here we propose TRANSLLLAMA, a policy-free SiMT system, in which
 43 an off-the-shelf pre-trained decoder-only LLM is fine-tuned on a dataset of causally aligned source

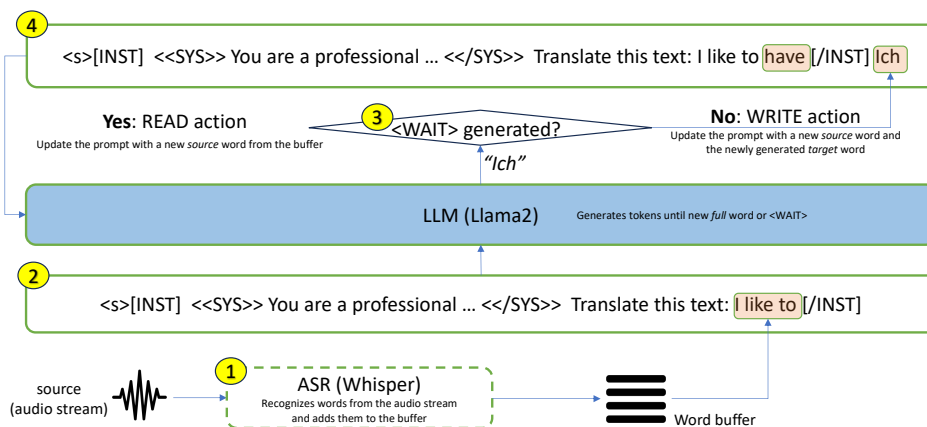


Figure 1: Model overview. The source audio stream is processed with an ASR model (1), which saves each recognized word into the buffer. The initial prompt (2) is built with k source words ($k = 3$ in this example). When the buffer has 3 words, the initial prompt is fed into the LLM, which generates output tokens until either a `<WAIT>` token or a full word is generated ("Ich" in this example). Then the prompt is updated with a new input ("have") and target ("Ich") word (WRITE action). Finally, the updated prompt (4) is fed back into the LLM. If `<WAIT>` is generated, the prompt is only updated with a new source word from the buffer (READ action).

44 and target sentences. The causality of the source is guaranteed by inserting one or more `<WAIT>`
 45 tokens into the target sentence to ensure that target content words never appear earlier than their clos-
 46 est equivalents in the source. We call our model policy-free, because as a result of fine-tuning on a
 47 causally aligned dataset the LLM becomes capable of deciding when to output translation and when
 48 to read in more of the source, without requiring a separate policy. At inference, the fine-tuned LLM
 49 is prompted with *part* of a source sentence concatenated with its corresponding (partial) translation
 50 and outputs one or more target tokens until either a full new word or a `<WAIT>` token is generated,
 51 which signals for more words to be read in. When extended with a off-the-shelf ASR model, in
 52 addition to text-to-text translation (T2TT), our system handles speech-to-speech translation (S2TT)
 53 tasks with quality (as measured by BLEU score (Papineni et al., 2002)) approaching that of some of
 54 the recently published baselines at comparable latencies.

55 Our main contributions are as follows:

- 56 1. We present the first system that leverages a decoder-only causal LLM for the SiMT task;
- 57 2. We propose a way to fine-tune a pre-trained LLM with direct supervision on a dataset of
58 causally aligned source-target sentence pairs;
- 59 3. We demonstrate that an LLM can perform both simultaneous translation and input segmen-
60 tation without a separate policy, with performance approaching or exceeding state of the
61 art.

62 The rest of the paper is structured as follows. Section 2 offers a brief overview of most recent SiMT
63 literature. In Section 3 we detail our system’s architecture, fine-tuning data preparation and training
64 procedure. In Section 4 we showcase its performance on *en-de* and *en-ru* language directions.
65 We conclude with Section 5 in which we discuss the limitations and directions for future work.

66 2 RELATED WORK

67 SiMT systems aim to deliver the best translation quality, usually measured with BLEU score (Pap-
68 ineni et al., 2002), while keeping its latency at an acceptable level. This quality-latency trade-off is
69 controlled by the "policy", which decides *when* to translate (i.e. perform a WRITE action) and when
70 to receive more input (i.e. perform a READ action). The various policies described in the literature
71 can be broadly categorized into fixed and adaptive (Zhang et al., 2020). Fixed policies (e.g. *wait-k*

(Ma et al., 2019a) are simple rules that determine the timing and order of WRITE and READ actions irrespective of the context. Early SiMT systems used *chunk-based* approaches (Fügen et al., 2007; Bangalore et al., 2012; Yarmohammadi et al., 2013; Sridhar et al., 2013), in which the input is split into sub-sentence phrases and translated independently of the previous chunk’s context, which compromised translation quality. Attempting to overcome this limitation, Dalvi et al. (2018) proposed an *incremental decoding* approach, in which chunk translations incorporate previous context encapsulated by an RNN’s hidden states. They showed that coupled with a simple segmentation strategy, their approach outperformed existing state of the art. On the other hand, adaptive policies (e.g. *wait-if* rules (Cho & Esipova, 2016)) make READ/WRITE actions more flexibly by taking account of the partial source and/or target. Adaptive policies can be implemented as separately trained agents (e.g. with reinforcement learning) (Grissom II et al., 2014; Gu et al., 2017; Satija & Pineau, 2016; Alinejad et al., 2018). In such policies, READ/WRITE actions can be taken based on attention (Raffel et al. (2017); Chiu & Raffel (2018); Arivazhagan et al. (2019); Ma et al. (2020b), or stability of the model’s outputs over n steps (so-called *local agreement* (Liu et al., 2020a; Ko et al., 2023; Polák et al., 2022)). More recent studies have also explored training the policy with binary search (Guo et al. (2023) aiming to maximize the gain in translation quality per each token read, or cast the problem of deciding when to translate as a hidden Markov transformer (Zhang & Feng (2023), in which hidden events correspond to the times at which to output translation.

Another promising line of work, related to the present study, aims to fine-tune encoder-decoder transformers, such as mBART (Liu et al., 2020b), originally pre-trained for sentence-level translation, for the SiMT task. For example, Fukuda et al. (2023); Kano et al. (2022) utilized fine-tuning on prefix-alignment data and Zhang et al. (2020) on meaningful units, achieving compelling performance on some language pairs.

Distinct from these approaches, we propose to fine-tune a large language model for the SiMT task on a dataset of causally aligned source-target sentence pairs, which we describe below.

97 3 METHOD

98 Although the LLMs we consider in this paper are designed to process only text input, we add an
99 ASR stage to enable it to also perform S2TT mode. Thus, we follow a cascaded approach shown in
100 Fig. 1.

101 **Causal alignment.** Training SiMT models, including optimal segmentation policies, with direct
102 supervision has remained a challenge (Guo et al., 2023) due to at least three reasons: (1) word
103 order inconsistencies between the source and target, (2) omissions of words from the target that
104 were present in the source, and/or (3) additions of words to the target not explicitly present in
105 the source, making it difficult to establish unambiguous correspondences between each source and
106 target words. This is less of a problem for offline translation models, because they are trained with
107 direct supervision on pairs of *complete* source and target sentences, and both during training and
108 inference the entire source context is revealed. However, it is not immediately clear how to use
109 direct supervision for the SiMT task, in which the model must begin translation based on *partial*
110 context. Nevertheless, we believe that direct supervision for the SiMT task is possible and propose a
111 way to accomplish that with a *causally aligned* dataset. In such a dataset, a target word never appears
112 before its corresponding (when such correspondence can be established) source word in time, which
113 is defined as the number of words from the sentence start. In other words, in a causally aligned
114 source-target sentence pair, source words are guaranteed to be causal relative to their corresponding
115 target words. We illustrate this in Fig. 2.

116 Note that the causal alignment is not always perfect: due to the word length mismatch between the
117 source and target, not all source words will have a corresponding target word, and vice versa, not
118 every target word will have a corresponding word in the source. However, as we demonstrate below,
119 fine-tuning an LLM on such a causally aligned dataset enables it to achieve results comparable to
120 some state-of-the-art baselines.

121 In order to causally align the source and target, we split each sentence using the `word_tokenize`
122 function from the `nltk` package (Bird et al., 2009), treating punctuation marks as "words", then find
123 the best correspondences between the source and target words with `SimAlign` (Jalili Sabet et al.,
124 2020), and finally insert as many `<WAIT>` tokens into the target as appropriate. If after alignment

original			causally aligned			original			causally aligned		
en	ru		en	ru		en	de		en	de	
1	They	→	Они	→	Они	1	He	→	Er	→	Er
2	live	→	живут	→	живут	2	took	→	befreite	→	@
3	in	→	глубоко	→	@	3	one	→	uns	→	@
4	the	→	в	→	@	4	of	→	von	→	@
5	depths	→	конголезских	→	глубоко	5	the	→	einer	→	@
6	of	→	джунглях	→	в	6	worst	→	der	→	@
7	the	→	.	→	@	7	scourges	→	schlimmsten	→	@
8	Congolese	→	где	→	конголезских	8	of	→	Geißeln	→	@
9	jungle	→	сложно	→	джунглях	9	mankind	→	der	→	@
10	and	→	проводить	→	.	10	away	→	Menschheit	→	befreite
11	it	→	исследования	→	где	11	from	→	.	→	@
12	has	→	.	→	@	12	us	→	uns	→	@
13	been	→	.	→	@	13	.	→	von	→	@
14	very	→	.	→	@	14	---	→	einer	→	@
15	difficult	→	.	→	сложно	15	---	→	der	→	@
16	to	→	.	→	проводить	16	---	→	schlimmsten	→	@
17	study	→	исследования	→	исследования	17	---	→	Geißeln	→	@
18	them	→	.	→	@	18	---	→	der	→	@
19	.	→	.	→	.	19	---	→	Menschheit	→	@
						20	---	→	.	→	@

Figure 2: **Causality-preserving alignment.** Two examples are shown: for en-ru (left) and en-de (right). If time is defined as the number of words from the beginning of the sentence, before alignment, some target words appear earlier than their corresponding English equivalents in the source. By inserting <WAIT> tokens (shown as "@"), we can shift those target words into the future, thereby achieving causality for every content word. " _ _ " are fillers added at the end of the source sentence if necessary to match its length with that of the target.

125 the target becomes longer than the source due to added <WAIT> tokens, we pad the source at the
 126 end with filler strings ensuring that the aligned source and target sentences have the same num-
 127 ber of "words". These filler strings are only used for convenient batching and are dropped before
 128 tokenization.

129 **Supervised Fine-Tuning.** We fine-tune the LLAMA-2 13B and and 70B models (Touvron et al.,
 130 2023)¹ to optimize the following objective:

$$\mathcal{L}_{\text{T2TT}} = - \sum_{t=1}^{|y|} \log p(y_t | y_{<t}, x_{\leq t}) \quad (1)$$

131 where y_t is the next target token, $y_{<t}$ are previously generated (and committed) tokens and $x_{\leq t}$ and
 132 the partial source tokens revealed up to the time step t . Following (Touvron et al., 2023), we zero
 133 out the loss on tokens corresponding to the system message and source, only backpropagating on
 134 the target tokens.

135 We use batches of prompt-response pairs collated in the following way. Before tokenization, each
 136 aligned sentence-target pair selected from the causally aligned dataset is trimmed from the right to
 137 leave first l words, where $l \sim U(1, L)$ and L is the full length of the causally aligned source-target
 138 pair. After trimming, all the <WAIT> tokens except the last one (if present) are dropped, because
 139 they are never plugged back into the input and only serve the purpose of signaling for more words

¹We found that the LLAMA-2-CHAT variants (both 13B and 70B), when fine-tuned on our causally aligned dataset performed slightly, but consistently, worse than LLAMA-2, and we report the results for the latter model only.

140 to be read in. Likewise, we drop all the fillers (if present) from the source. Finally, the system
 141 message, trimmed source and trimmed target are joined into the prompt (as shown in Fig. 4) and
 142 tokenized. Because there is no `<WAIT>` token in the LLAMA 2 tokenizer, we use 0 (which originally
 143 corresponds to the `<UNK>` token). Thus, the model is fine-tuned to either output the next token of
 144 a word or `<WAIT>`, if the partial source does not contain sufficient information needed to predict
 145 translation.

146 To save memory, we loaded the the base model in 4-bit precision. This allowed us to fine-tune
 147 LLAMA 2 70B on one NVIDIA A100 80GB device. We fine-tune the base model with LoRA (Hu
 148 et al., 2022) with $r = 16$ and $\alpha = 32$ for 3 epochs with a batch size of 25 and gradient accu-
 149 mulation of 4 steps. We save model checkpoints every 10 steps and select the one with the lowest
 150 validation loss for inference. For optimization, we used the `paged_adamw_32bit` optimizer with
 151 default parameters, and a learning rate schedule with a linear warm-up of 10 steps up to 0.00005,
 152 followed by a cosine decay. For parameter-efficient training, as well as for inference, we used the
 153 `transformers`² library.

154 **Inference.** At inference, given a prompt (Fig. 4) comprised of a system message, partial source and
 155 previously committed partial target, the LLM greedily generates one or more next tokens. We use
 156 modified `wait-k` (Ma et al., 2019a), in which `WRITE` actions are only allowed when the length of
 157 the `PARTIAL_SOURCE` is equal or greater than k . Since k controls the tradeoff between quality
 158 and latency, we report results for different values of k . After a full new word – which may consist of
 159 several tokens – is generated, the prompt is updated by appending a new source word to the partial
 160 source and the newly generated word to the partial target. This process is repeated until the LLM
 161 generates the `<EOS>` token. All the generation parameters were at default, except `top_p` which we
 162 set to 0.7. We did not use beam search during generation.

PARTIAL_SOURCE	PARTIAL_TARGET	Prediction
I		<code><WAIT></code>
I like		Я
I like to	Я	люблю
I like to have	Я люблю	<code><WAIT></code>
I like to have tea	Я люблю	пить
I like to have tea	Я люблю пить	чай
I like to have tea in the	Я люблю пить чай	<code><WAIT></code>
I like to have tea in the morning.	Я люблю пить чай	по
I like to have tea in the morning.	Я люблю пить чай по	утрам.
I like to have tea in the morning.	Я люблю пить чай по утрам.	<code><EOS></code>

Figure 3: An illustration of the inference process for the `en-ru` language pair. Assuming $k = 1$, given the prompt with one source and zero target words, the model first outputs `<WAIT>`, which signals for the next source word to be read in. At the next step, the model generates the first target word (Я), which is plugged into the prompt at the next step. This process continues until `<EOS>` is generated.

163 After all the source words have been revealed, the input is no longer partial and no new words are
 164 added to it, but the generation process continues until `<EOS>`. Importantly, if the model generates
 165 the `<WAIT>` token, a new source word is read in, but the `<WAIT>` token itself is not appended to
 166 the partial target. We illustrate the inference process in Fig. 3 and Algorithm 1.

167 **Prompt structure.** We follow a similar prompt structure as in Touvron et al. (2023) (Fig. 4). For
 168 the `SYSTEM_MESSAGE` we used the following text: "You are a professional conference interpreter.
 169 Given an English text you translate it into {TARGET_LANGUAGE} as accurately and as concisely
 170 as possible, NEVER adding comments of your own. You output translation when the information
 171 available in the source is unambiguous, otherwise you output the wait token ({WAIT_TOKEN}),
 172 not flanked by anything else. It's important that you get this right.". We note that while the system
 173 message is only necessary in zero-shot SiMT scenarios – which we discuss below – for consistency
 174 we still kept it in all the experiments reported here, including those involving supervised fine-tuning.

²<https://huggingface.co/docs/transformers/installation>

Algorithm 1 Inference process

```

partial_target = []
k = WAIT_K

while True:
    partial_source = SOURCE[:k]
    prompt = " ".join([SYS_MSG, partial_source, partial_target])

    # generate until next full word, or <EOS> or <WAIT>
    next_word = model.generate(prompt)

    if next_word == "<EOS>":
        break # finish sentence
    elif next_word == "<WAIT>":
        k += 1 # READ action
    else:
        partial_target.append(next_word) # WRITE action
        k += 1

```

```

<s>[INST]
<<SYS>>

[SYSTEM_MESSAGE]
<</SYS>>
Translate this text: PARTIAL_SOURCE [/INST] PARTIAL_TARGET

```

Figure 4: Prompt structure.

175 **Automatic speech recognition.** Given that the LLMs are designed to process text input, to enable
176 S2TT we first need to extract text from input audio, for which we use Whisper ³ (Radford et al.,
177 2023). Specifically, for each READ action, a new segment of audio, lasting 200 ms, is added to any
178 previously read audio chunks and then processed by Whisper. This method of fixed audio windowing
179 often results in partially clipped words. To address this, we discard the last word predicted by
180 Whisper during each READ action unless the entire source audio has been read in. We note that
181 this approach to online ASR is somewhat naive and has room for improvement – as indicated by a
182 roughly 1 BLEU point decrease due to ASR-related errors (Fig. 9). Since our main objective was
183 to assess the capability of LLMs to perform SiMT tasks, we leave exploring ways to decrease ASR
184 errors to future work.

185 **4 RESULTS**

186 **Data.** For supervised fine-tuning (SFT), validation and testing, we used MuST-C v2.0 (Di Gangi
187 et al., 2019) for English-to-German (en-de) and English-to-Russian (en-ru) translation direction.
188 We randomly selected 4000 sentences for training and 100 sentences for validation. However, since
189 it is possible that the dataset that LLAMA2 was pre-trained on and MuST-C v2.0 (including its vali-
190 dation and test set) might have overlapping content, we also compiled another test set, which we call
191 NEW-TED-2023. This test set has a similar content type (TED talks) and follows the same format
192 as the original MuST-C v2.0, but only includes talks posted after February 2023. The dataset has
193 two parts: 102 source-target pairs for en-de and 102 for en-ru language pair. Unless indicated
194 otherwise, we report the results obtained on this test set.

195 **T2TT.** We first analyzed the T2TT performance of our approach on the MuST-C dataset v2.0
196 (Di Gangi et al., 2019). To get a sense for the quality-latency tradeoff, we plot BLEU scores against
197 several different values of k (because k is the only way to control the translation latency). The

³We used whisper-large-v2.

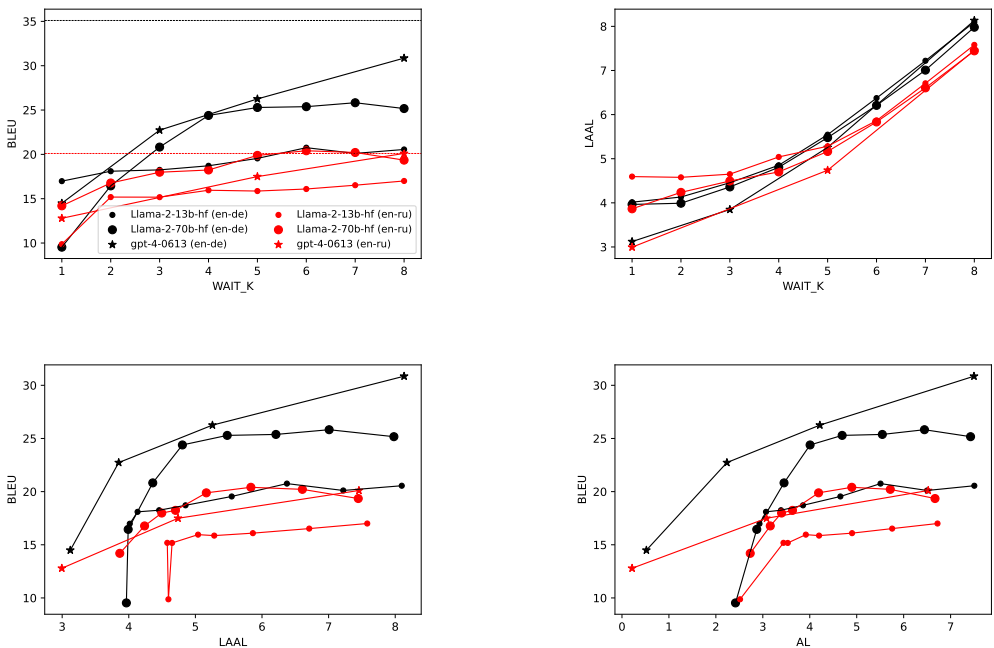


Figure 5: Dependence of latency and quality on k (top panels) and quality-latency tradeoff curves (bottom panels) for the T2TT mode on the MuST-C v2.0 dataset. For reference, dashed lines indicated GPT-4’s sentence-level (i.e. with k set to the sentence length) BLEU scores: black for `en-de` and red for `en-ru`.

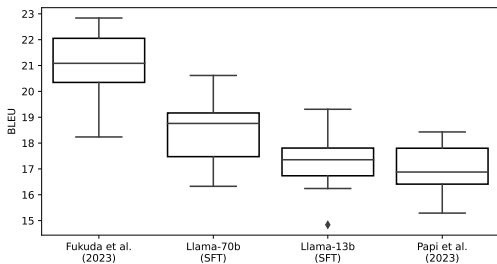


Figure 6: S2TT performance of SFT LLAMA-2 and two recently published models on the `en-de` language pair on TED-TST-2023. See also Appendix C.1.

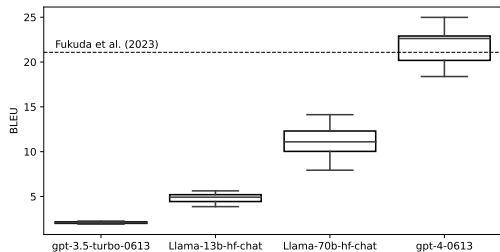


Figure 7: Zero-shot S2TT performance of our approach compared with GPT-3.5 and GPT-4 on the `en-de` language pair on TED-TST-2023.

198 results, shown in Fig. 5, suggest that the LLM’s size is a major factor determining the translation
 199 quality.

200 **S2TT.** We next test fine-tuned LLMs and compare them with two recently published S2TT baselines
 201 (Fukuda et al., 2023; Papi et al., 2023) as well as to OpenAI’s GPT-3.5 and GPT-4 (in zero-shot
 202 mode). To ensure as fair a comparison as possible, we ensured that average lagging (AL) of all of
 203 the models below approximately 2000 ms. For Llama-2 models we set $k = 5$ (the other models’
 204 settings are listed in Appendix D). The boxplots in Figs. 6, 7 and throughout are drawn based on
 205 data from 10 evaluation runs of the same model with the same parameters on sentence pairs sampled
 206 with replacement from TED-TST-2023. The results show a degradation of translation quality by
 207 approximately 1 BLEU score point compared to T2TT mode, which is to be expected due to ASR
 208 errors (Fig. 9).

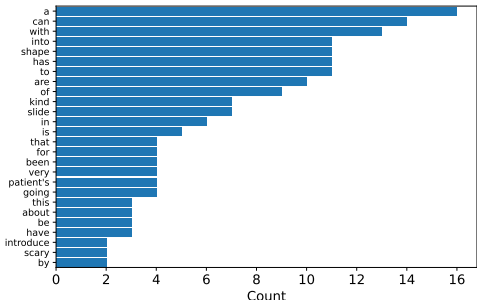


Figure 8: After fine-tuning, LLAMA-2 generates `<WAIT>` tokens predominantly after function words (especially articles and prepositions).

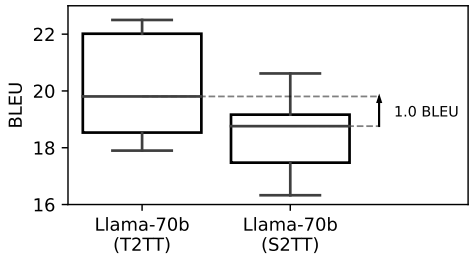


Figure 9: Performance decrease due to ASR-related errors. In T2TT mode, Llama2-70b performs by about 1 BLEU score point better than the the same model on the same data in S2TT mode.

k	w/ <code><WAIT></code>	w/o <code><WAIT></code>	k	w/ <code><WAIT></code>	w/o <code><WAIT></code>	k	w/ <code><WAIT></code>	w/o <code><WAIT></code>
1	15.23	14.88	1	14.76	10.80	1	17.17	4.64
2	17.17	15.66	2	14.97	11.94	2	16.83	7.84
			4	17.42	15.67	4	19.24	14.80

Table 1: Removing the instruction to generate or suppressing the `<WAIT>` token degrades performance. The numbers indicate BLEU scores on TED-TST-2023 (en-de) in T2TT mode for GPT-4 (a), supervised fine-tuned Llama-2-13b-hf (b) and Llama-2-70b-hf (c).

209 **Zero-shot T2TT.** Can the LLMs perform the SiMT task zero-shot, that is without any prior fine-
 210 tuning? To answer this question, we used LLMs that have been fine-tuned with RLHF for instruc-
 211 tion following: open-source LLAMA2-CHAT, as well as GPT-3.5 (gpt-3.5-turbo-0613) and
 212 GPT-4 (gpt-4-0613), which were among the strongest closed-source LLMs available at the time
 213 of writing this paper. In general, with the notable exception of GPT-4, zero-shot performance was
 214 poor. Inspection of the translations revealed that the models consistently failed to follow the prompt
 215 instruction, specifically, (1) generating output in English rather than the target language, (2) adding
 216 expressly prohibited explanatory comments, (3) restating or summarizing the task, or (4) explain-
 217 ing the reason for adding `<WAIT>` tokens). GPT-4 was surprisingly good, performing better than
 218 the supervised fine-tuned LLAMA2-70B, and we speculate that the performance of GPT-3.5 and
 219 GPT-4 could be further improved with SFT⁴, more sophisticated generation strategies and prompt
 220 engineering.

221 **Importance of wait tokens.** To evaluate the utility of `<WAIT>` tokens, we conduct two ablation
 222 experiments. In the first experiment we consider a zero-shot translation scenario in which GPT-4
 223 was not instructed to use `<WAIT>` tokens. In the second experiment, we suppress the generation
 224 of `<WAIT>` tokens in supervised fine-tuned LLMs. The results, as indicated in Table 1, reveal that
 225 GPT-4 demonstrates marginally inferior performance when $k \in \{1, 2\}$ ⁵ when not instructed about
 226 `<WAIT>` tokens. However, it is important to note that in a zero-shot context, the GPT-3.5 and GPT-
 227 4 seldom generated `<WAIT>` tokens (almost never for $k > 2$). Therefore, the directive to employ
 228 these tokens only exhibited a discernible impact for smaller values of k . By contrast, in the SFT
 229 scenario, suppressing `<WAIT>` tokens led to significantly decreased performance for both the 13
 230 and 70B versions of LLAMA-2 (Table 1 (b, c)).

231 To gain insight into where LLAMA-2 tended to insert the `<WAIT>` token, we plot the distribution
 232 of words after which the SFT models generated this token. Fig. 8 shows that most of the time the
 233 model generated `<WAIT>` after function words – which makes sense – rather than content words,
 234 indicating that it had learned to choose appropriately between READ and WRITE actions.

⁴SFT was not available for GPT-3.5 and GPT-4 at the time of writing this paper.

⁵We did not investigate the role of `<WAIT>` tokens for $k > 2$, because GPT-4 almost never generates them for those values of k .

235 5 CONCLUSION AND FUTURE DIRECTIONS

236 We have shown that with minimal fine-tuning and without resorting to sophisticated training tech-
 237 niques (e.g. checkpoint averaging (Fukuda et al., 2023)), an off-the-shelf pre-trained LLM can per-
 238 form simultaneous translation and achieve encouraging results that rival some of the recent SiMT
 239 models. This opens interesting directions to be explored in future work, such as multilingual fine-
 240 tuning, self-instruct (Wang et al., 2023b) and human preference tuning (Ouyang et al., 2022).

241 There are several reasons to believe that we are far from unlocking the full potential of LLMs
 242 for SiMT. First, we followed the practice – standard in the SiMT literature – of evaluating the
 243 model on individual sentences randomly sampled from continuous prose. However, many (if not
 244 the majority of) short sentences are ambiguous when taken out of context. Even human conference
 245 interpreters routinely prepare for an upcoming translation job, studying relevant materials, which
 246 means that they do not have to translate sentences taken out of context. For this reason, we believe
 247 that the most straightforward way to boost the performance of future LLM-based SiMT systems is to
 248 insert background information into the prompt. Second, the big difference in zero-shot performance
 249 between GPT-3.5 and GPT-4 suggests that size is likely the biggest factor determining the model’s
 250 translation quality, and that further gains can be achieved once SFT becomes available for these
 251 closed-source models.

252 In conclusion, we note that there are several performance bottlenecks that must be addressed be-
 253 fore our approach can be deployed for simultaneous translation in the real world. As we show in
 254 Appendix E, these bottlenecks result from a long system message, which is often longer than the
 255 source sentence itself, as well as delays introduced by the ASR subsystem and weight quantization.
 256 We believe that these issues are not prohibitive. Specifically, instead of using a separate ASR model,
 257 future work might follow an end-to-end approach similar to Fathullah et al. (2023), in which in-
 258 stead of being converted into text with an separate ASR model, the audio is directly mapped into
 259 the LLM’s embedding space, reducing the system’s overall latency. Efficient quantization schemes,
 260 faster algorithms and hardware support for low bit-width arithmetic are also promising directions.
 261 Finally, because LLAMA-2 was trained predominantly on English text, its tokenizer represents En-
 262 glish more efficiently than other languages. That is, fewer tokens on average are needed to encode a
 263 text in English than a text of the same length (in characters) in another, less represented, language.
 264 Thus, future LLMs pre-trained on a linguistically more balanced dataset, might be slightly faster at
 265 inference.

266 REFERENCES

- 267 Ashkan Alinejad, Maryam Siahbani, and Anoop Sarkar. Prediction improves simultaneous neu-
 268 ral machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in*
 269 *Natural Language Processing*, pp. 3022–3027, Brussels, Belgium, October–November 2018.
 270 Association for Computational Linguistics. doi: 10.18653/v1/D18-1337. URL [https://](https://aclanthology.org/D18-1337)
 271 aclanthology.org/D18-1337.
- 272 Naveen Arivazhagan, Colin Cherry, Wolfgang Macherey, Chung-Cheng Chiu, Semih Yavuz, Ruom-
 273 ing Pang, Wei Li, and Colin Raffel. Monotonic infinite lookback attention for simultaneous ma-
 274 chine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational*
 275 *Linguistics*, pp. 1313–1323, Florence, Italy, July 2019. Association for Computational Linguis-
 276 tics. doi: 10.18653/v1/P19-1126. URL <https://aclanthology.org/P19-1126>.
- 277 Srinivas Bangalore, Vivek Kumar Rangarajan Sridhar, Prakash Kolan, Ladan Golipour, and Aura
 278 Jimenez. Real-time incremental speech-to-speech translation of dialogs. In *Proceedings of the*
 279 *2012 Conference of the North American Chapter of the Association for Computational Linguis-*
 280 *tics: Human Language Technologies*, pp. 437–445, 2012.
- 281 Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing*
 282 *text with the natural language toolkit*. " O’Reilly Media, Inc.", 2009.
- 283 Colin Cherry and George Foster. Thinking slow about latency evaluation for simultaneous machine
 284 translation. *arXiv preprint arXiv:1906.00048*, 2019.

- 285 Chung-Cheng Chiu and Colin Raffel. Monotonic chunkwise attention. In *6th International Confer-*
286 *ence on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018,*
287 *Conference Track Proceedings*. OpenReview.net, 2018. URL [https://openreview.net/](https://openreview.net/forum?id=Hko85plCW)
288 [forum?id=Hko85plCW](https://openreview.net/forum?id=Hko85plCW).
- 289 Kyunghyun Cho and Masha Esipova. Can neural machine translation do simultaneous translation?
290 *arXiv preprint arXiv:1606.02012*, 2016.
- 291 Fahim Dalvi, Nadir Durrani, Hassan Sajjad, and Stephan Vogel. Incremental decoding and train-
292 ing methods for simultaneous translation in neural machine translation. In *Proceedings of the*
293 *2018 Conference of the North American Chapter of the Association for Computational Lin-*
294 *guistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 493–499, New Orleans,
295 Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2079.
296 URL <https://aclanthology.org/N18-2079>.
- 297 Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. MuST-
298 C: a Multilingual Speech Translation Corpus. In *Proceedings of the 2019 Conference of the*
299 *North American Chapter of the Association for Computational Linguistics: Human Language*
300 *Technologies, Volume 1 (Long and Short Papers)*, pp. 2012–2017, Minneapolis, Minnesota, June
301 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1202. URL [https://](https://aclanthology.org/N19-1202)
302 aclanthology.org/N19-1202.
- 303 Maha Elbayad, Laurent Besacier, and Jakob Verbeek. Efficient Wait-k Models for Simultaneous Ma-
304 chine Translation. In *Proc. Interspeech 2020*, pp. 1461–1465, 2020. doi: 10.21437/Interspeech.
305 2020-1241.
- 306 Yassir Fathullah, Chunyang Wu, Egor Lakomkin, Junteng Jia, Yuan Shangguan, Ke Li, Jinxi Guo,
307 Wenhan Xiong, Jay Mahadeokar, Ozlem Kalinli, Christian Fuegen, and Mike Seltzer. Prompting
308 large language models with speech recognition abilities, 2023.
- 309 Christian Fügen, Alex Waibel, and Muntsin Kolss. Simultaneous translation of lectures and
310 speeches. *Machine translation*, 21:209–252, 2007.
- 311 Ryo Fukuda, Yuta Nishikawa, Yasumasa Kano, Yuka Ko, Tomoya Yanagita, Kosuke Doi, Mana
312 Makinae, Sakriani Sakti, Katsuhito Sudoh, and Satoshi Nakamura. NAIST simultaneous speech-
313 to-speech translation system for IWSLT 2023. In *Proceedings of the 20th International Con-*
314 *ference on Spoken Language Translation (IWSLT 2023)*, pp. 330–340, Toronto, Canada (in-
315 person and online), July 2023. Association for Computational Linguistics. doi: 10.18653/v1/
316 2023.iwslt-1.31. URL <https://aclanthology.org/2023.iwslt-1.31>.
- 317 Alvin Grissom II, He He, Jordan Boyd-Graber, John Morgan, and Hal Daumé III. Don’t until the
318 final verb wait: Reinforcement learning for simultaneous machine translation. In *Proceedings*
319 *of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp.
320 1342–1352, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.
321 3115/v1/D14-1140. URL <https://aclanthology.org/D14-1140>.
- 322 Jiatao Gu, Graham Neubig, Kyunghyun Cho, and Victor O.K. Li. Learning to translate in real-
323 time with neural machine translation. In *Proceedings of the 15th Conference of the European*
324 *Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pp. 1053–
325 1062, Valencia, Spain, April 2017. Association for Computational Linguistics. URL [https://](https://aclanthology.org/E17-1099)
326 aclanthology.org/E17-1099.
- 327 Shoutao Guo, Shaolei Zhang, and Yang Feng. Learning optimal policy for simultaneous machine
328 translation via binary search. In *Proceedings of the 61st Annual Meeting of the Association for*
329 *Computational Linguistics (Volume 1: Long Papers)*, pp. 2318–2333, Toronto, Canada, July 2023.
330 Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.130. URL [https://](https://aclanthology.org/2023.acl-long.130)
331 aclanthology.org/2023.acl-long.130.
- 332 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
333 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-*
334 *ference on Learning Representations*, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=nZeVKeeFYf9)
335 [forum?id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).

- 336 Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. SimAlign: High quality
337 word alignments without parallel training data using static and contextualized embeddings. In
338 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing:
339 Findings*, pp. 1627–1643, Online, November 2020. Association for Computational Linguistics.
340 URL <https://www.aclweb.org/anthology/2020.findings-emnlp.147>.
- 341 Yasumasa Kano, Katsuhito Sudoh, and Satoshi Nakamura. Simultaneous neural machine trans-
342 lation with prefix alignment. In *Proceedings of the 19th International Conference on Spo-
343 ken Language Translation (IWSLT 2022)*, pp. 22–31, Dublin, Ireland (in-person and online),
344 May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.iwslt-1.3. URL
345 <https://aclanthology.org/2022.iwslt-1.3>.
- 346 Yuka Ko, Ryo Fukuda, Yuta Nishikawa, Yasumasa Kano, Katsuhito Sudoh, and Satoshi Naka-
347 mura. Tagged end-to-end simultaneous speech translation training using simultaneous interpre-
348 tation data. In *Proceedings of the 20th International Conference on Spoken Language Trans-
349 lation (IWSLT 2023)*, pp. 363–375, Toronto, Canada (in-person and online), July 2023. As-
350 sociation for Computational Linguistics. doi: 10.18653/v1/2023.iwslt-1.34. URL <https://aclanthology.org/2023.iwslt-1.34>.
- 352 Danni Liu, Gerasimos Spanakis, and Jan Niehues. Low-Latency Sequence-to-Sequence Speech
353 Recognition and Translation by Partial Hypothesis Selection. In *Proc. Interspeech 2020*, pp.
354 3620–3624, 2020a. doi: 10.21437/Interspeech.2020-2897.
- 355 Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike
356 Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine transla-
357 tion. *Transactions of the Association for Computational Linguistics*, 8:726–742, 2020b. doi:
358 10.1162/tacl_a_00343. URL <https://aclanthology.org/2020.tacl-1.47>.
- 359 Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng, Kaibo Liu, Baigong Zheng, Chuanqiang
360 Zhang, Zhongjun He, Hairong Liu, Xing Li, Hua Wu, and Haifeng Wang. STACL: Simultaneous
361 translation with implicit anticipation and controllable latency using prefix-to-prefix framework.
362 In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*,
363 pp. 3025–3036, Florence, Italy, July 2019a. Association for Computational Linguistics. doi:
364 10.18653/v1/P19-1289. URL <https://aclanthology.org/P19-1289>.
- 365 Xutai Ma, Juan Miguel Pino, James Cross, Liezl Puzon, and Jiatao Gu. Monotonic multihead
366 attention. *CoRR*, abs/1909.12406, 2019b. URL <http://arxiv.org/abs/1909.12406>.
- 367 Xutai Ma, Mohammad Javad Dousti, Changhan Wang, Jiatao Gu, and Juan Pino. SIMULEVAL: An
368 evaluation toolkit for simultaneous translation. In *Proceedings of the 2020 Conference on Empiri-
369 cal Methods in Natural Language Processing: System Demonstrations*, pp. 144–150, Online, Oc-
370 tober 2020a. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.
371 19. URL <https://aclanthology.org/2020.emnlp-demos.19>.
- 372 Xutai Ma, Juan Miguel Pino, James Cross, Liezl Puzon, and Jiatao Gu. Monotonic multihead
373 attention. In *International Conference on Learning Representations*, 2020b. URL <https://openreview.net/forum?id=Hyg96gBKPS>.
- 375 Nathalia Nascimento, Paulo Alencar, and Donald Cowan. Gpt-in-the-loop: Adaptive decision-
376 making for multiagent systems, 2023.
- 377 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
378 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-
379 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,
380 and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- 381 Sara Papi, Marco Gaido, Matteo Negri, and Marco Turchi. Over-generation cannot be rewarded:
382 Length-adaptive average lagging for simultaneous speech translation. In Julia Ive and Ruiqing
383 Zhang (eds.), *Proceedings of the Third Workshop on Automatic Simultaneous Translation*, pp.
384 12–17, Online, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
385 autosimtrans-1.2. URL <https://aclanthology.org/2022.autosimtrans-1.2>.

- 386 Sara Papi, Matteo Negri, and Marco Turchi. Attention as a guide for simultaneous speech transla-
387 tion. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*
388 (*Volume 1: Long Papers*), pp. 13340–13356, Toronto, Canada, July 2023. Association for Compu-
389 tational Linguistics. doi: 10.18653/v1/2023.acl-long.745. URL <https://aclanthology.org/2023.acl-long.745>.
- 391 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
392 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Associa-*
393 *tion for Computational Linguistics*, pp. 311–318, Philadelphia, Pennsylvania, USA, July 2002.
394 Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- 396 Peter Polák, Ngoc-Quan Pham, Tuan Nam Nguyen, Danni Liu, Carlos Mullov, Jan Niehues,
397 Ondřej Bojar, and Alexander Waibel. CUNI-KIT system for simultaneous speech translation
398 task at IWSLT 2022. In *Proceedings of the 19th International Conference on Spoken Lan-*
399 *guage Translation (IWSLT 2022)*, pp. 277–285, Dublin, Ireland (in-person and online), May
400 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.iwslt-1.24. URL
401 <https://aclanthology.org/2022.iwslt-1.24>.
- 402 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.
403 Robust speech recognition via large-scale weak supervision. In *International Conference on Ma-*
404 *chine Learning*, pp. 28492–28518. PMLR, 2023.
- 405 Colin Raffel, Minh-Thang Luong, Peter J. Liu, Ron J. Weiss, and Douglas Eck. Online and linear-
406 time attention by enforcing monotonic alignments. In *Proceedings of the 34th International Con-*
407 *ference on Machine Learning - Volume 70, ICML’17*, pp. 2837–2846. JMLR.org, 2017.
- 408 Harsh Satija and Joelle Pineau. Simultaneous machine translation using deep reinforcement learn-
409 ing. 2016. URL <https://api.semanticscholar.org/CorpusID:201718412>.
- 410 Vivek Kumar Rangarajan Sridhar, John Chen, Srinivas Bangalore, Andrej Ljolje, and Rathinavelu
411 Chengalvarayan. Segmentation strategies for streaming speech translation. In *Proceedings of the*
412 *2013 Conference of the North American Chapter of the Association for Computational Linguis-*
413 *tics: Human Language Technologies*, pp. 230–238, 2013.
- 414 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
415 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher,
416 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy
417 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
418 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
419 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
420 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
421 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
422 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
423 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
424 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
425 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,
426 2023.
- 427 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai
428 Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. A survey on large
429 language model based autonomous agents, 2023a.
- 430 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and
431 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions,
432 2023b.
- 433 Zekun Wang, Ge Zhang, Kexin Yang, Ning Shi, Wangchunshu Zhou, Shaochun Hao, Guangzheng
434 Xiong, Yizhi Li, Mong Yuan Sim, Xiuying Chen, Qingqing Zhu, Zhenzhu Yang, Adam Nik,
435 Qi Liu, Chenghua Lin, Shi Wang, Ruibo Liu, Wenhu Chen, Ke Xu, Dayiheng Liu, Yike Guo, and
436 Jie Fu. Interactive natural language processing, 2023c.

- 437 Mahsa Yarmohammadi, Vivek Kumar Rangarajan Sridhar, Srinivas Bangalore, and Baskaran
438 Sankaran. Incremental segmentation and decoding strategies for simultaneous translation. In
439 *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pp.
440 1032–1036, 2013.
- 441 Ruiqing Zhang, Chuanqiang Zhang, Zhongjun He, Hua Wu, and Haifeng Wang. Learning adaptive
442 segmentation policy for simultaneous translation. In *Proceedings of the 2020 Conference on Em-
443 pirical Methods in Natural Language Processing (EMNLP)*, pp. 2280–2289, Online, November
444 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.178. URL
445 <https://aclanthology.org/2020.emnlp-main.178>.
- 446 Shaolei Zhang and Yang Feng. Universal simultaneous machine translation with mixture-of-experts
447 wait-k policy. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Lan-
448 guage Processing*, pp. 7306–7317, Online and Punta Cana, Dominican Republic, November
449 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.581. URL
450 <https://aclanthology.org/2021.emnlp-main.581>.
- 451 Shaolei Zhang and Yang Feng. Information-transport-based policy for simultaneous translation. In
452 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.
453 992–1013, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
454 Linguistics. doi: 10.18653/v1/2022.emnlp-main.65. URL [https://aclanthology.org/
455 2022.emnlp-main.65](https://aclanthology.org/2022.emnlp-main.65).
- 456 Shaolei Zhang and Yang Feng. Hidden markov transformer for simultaneous machine translation,
457 2023.

458 A SUPPLEMENTARY RESULTS FOR THE S2TT TASK

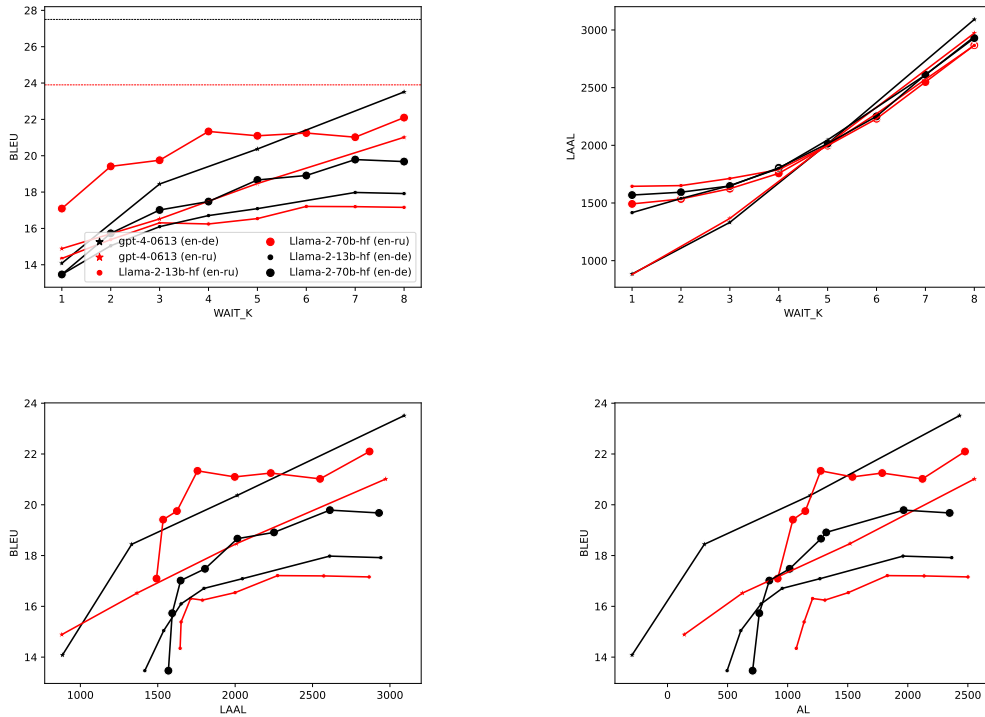


Figure 10: Dependence of latency and quality on k (top panels) and quality-latency tradeoff curves (bottom panels) for the S2TT mode on the NEW-TED-2023 dataset. For reference, dashed lines indicated GPT-4’s sentence-level (i.e. with k set to the sentence length) BLEU scores: black for en-de and red for en-ru.

459 B ENGLISH-RUSSIAN S2TT TASK ($k = 5$)

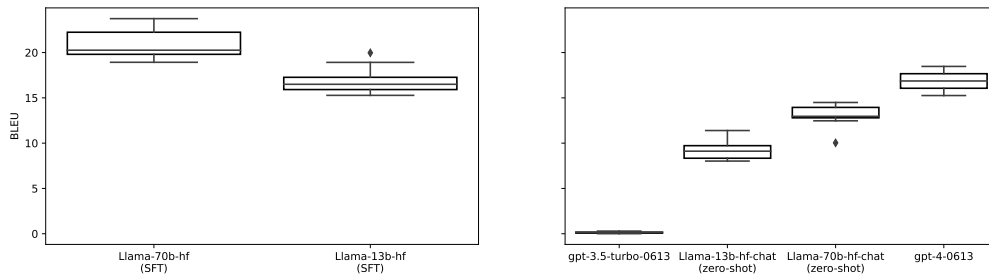


Figure 11: S2TT en-ru performance of our method on TED-TST-2023. Left panel: supervised fine-tuned LLAMA-2. Right panel: zero-shot S2TT performance of LLAMA-2-CHAT. All the runs were on TED-TST-2023, with $k = 5$ to ensure AL around 2000 ms. Each of the boxplots is drawn based on 10 evaluation runs on sentences randomly sampled with replacement from the test set.

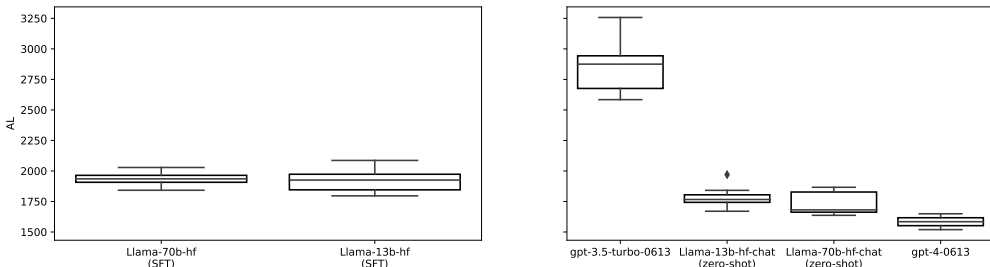


Figure 12: Average lagging in S2TT mode for the English-Russian language pair. Left panel: supervised fine-tuned LLAMA-2. Right panel: zero-shot S2TT performance of LLAMA-2-CHAT. All the runs were on TED-TST-2023, with $k = 5$ to ensure AL around 2000 ms. Each of the boxplots is drawn based on data from 10 evaluation runs on sentences randomly sampled with replacement from the test set.

460 C ADDITIONAL PERFORMANCE DATA FOR THE S2TT TASK

461 C.1 ENGLISH-GERMAN

462 Here we report additional comparisons including latency performance measured using several differ-
 463 ent metrics, including Average Lagging (AL) (Ma et al., 2019a), Length Adaptive Average Lagging
 464 (LAAL) (Papi et al., 2022), Average Proportion (AP) (Cho & Esipova, 2016) and Differentiable
 465 Average Lagging (DAL) (Cherry & Foster, 2019).

System	BLEU	LAAL	AL	AP	DAL
gpt-3.5-turbo-0613 (zero-shot)	2.08 (0.24)	2637.11 (252.79)	2574.98 (230.95)	0.35 (0.0)	2477.55 (146.26)
gpt-4-0613 (zero-shot)	21.82 (2.81)	2448.86 (74.74)	1998.63 (110.91)	0.94 (0.03)	2813.47 (69.48)
Llama-70b-hf (SFT)	18.41 (1.4)	2107.57 (59.68)	1619.64 (76.47)	0.84 (0.02)	2454.72 (67.84)
Llama-13b-hf (SFT)	17.07 (0.68)	2358.89 (34.11)	1880.76 (61.77)	0.88 (0.02)	2735.34 (40.88)
Papi et al. (2023)	17.01 (1.0)	2295.72 (41.54)	1867.1 (148.69)	0.77 (0.01)	3251.38 (139.12)
Fukuda et al. (2023)	21.08 (1.41)	2005.39 (71.04)	1397.33 (85.74)	0.9 (0.01)	3066.15 (122.01)

Table 2: Mean performance metrics of Llama-2 (SFT) compared to some recent S2TT systems and GPT-3.5 and GPT-4 (zero-shot). Then mean and standard deviation (in brackets) are computed over 10 runs of the same model on 102 source-target pairs sampled with replacement from TED-TST-2023.

466 C.2 ENGLISH-RUSSIAN

System	BLEU	LAAL	AL	AP	DAL
gpt-3.5-turbo-0613 (zero-shot)	0.14 (0.1)	2876.85 (240.03)	2861.22 (245.91)	0.28 (0.04)	2661.22 (231.0)
gpt-4-0613 (zero-shot)	16.86 (2.27)	2022.81 (20.3)	1584.38 (91.81)	0.82 (0.04)	2390.11 (23.65)
Llama-70b-hf (SFT)	20.96 (1.71)	2252.75 (49.77)	1937.76 (62.75)	0.9 (0.08)	2676.56 (62.11)
Llama-13b-hf (SFT)	16.9 (1.52)	2238.6 (48.38)	1917.46 (90.38)	0.87 (0.03)	2641.01 (45.73)

Table 3: Mean performance metrics of Llama-2 (SFT) compared to some recent S2TT systems and GPT-3.5 and GPT-4 (zero-shot). Then mean and standard deviation (in brackets) are computed over 10 runs of the same model on 102 source-target pairs sampled with replacement from TED-TST-2023.

467 **D** PARAMETERS USED FOR COMPARISONS WITH BASELINES ON THE S2ST
468 EN-DE TASK

469 [Papi et al. \(2023\)](#)

470 We used the open-source implementation of the model⁶. The evaluations were run in *SimulEval*⁷
471 ([Ma et al., 2020a](#)) with the following parameters:

472
473 `extract-attn-from-layer 5`
474 `frame-num 2`
475 `attn-threshold 0.25`
476 `speech-segment-factor 8`

477 [Fukuda et al. \(2023\)](#)

478 The source code for the model and weights were obtained on request from the authors. The
479 evaluations were run in *SimulEval* with the following parameters:

480
481 `source-segment-size 950`
482 `la-n 2`
483 `beam 5`
484 `sacrebleu-tokenizer 13a`

485 We chose these parameters aiming to maximize the BLEU score while keeping AL approximately
486 below 2000 ms.

⁶https://github.com/hlt-mt/FBK-fairseq/tree/master/fbk_works

⁷<https://github.com/facebookresearch/SimulEval>

487 E INFERENCE WALL TIME COMPARISONS

488 Here we compare real-time factors of our model in different sizes and compare them with those
 489 of the selected baselines and GPT-4. Real-time factor is the ratio of the amount of time taken to
 490 process source audio to the length of the source audio itself ⁸. We note that removing the system
 491 message from the prompt speeds up inference with no noticeable drop in quality for supervised
 492 fine-tuned models. Loading our model’s weights with 16-bit (instead of 4-bit) quantization further
 493 accelerates inference. Finally, we observe that the use of ASR in S2TT mode substantially reduces
 494 system speed. An end-to-end implementation, directly converting raw source audio into the LLM’s
 495 embedding space, could potentially alleviate this performance bottleneck.

model	mode	quantization	system message	size, bn param.	RTF
Ours	T2TT	16-bit	no	13	1.7
Ours	T2TT	4-bit	no	13	2.2
Ours	T2TT	16-bit	yes	13	2.9
Ours	T2TT	4-bit	yes	13	4.2
Ours	S2TT	16-bit	no	13	5.9
Ours	S2TT	4-bit	no	13	7.6
Ours	S2TT	16-bit	yes	13	8.0
Ours	S2TT	4-bit	yes	13	9.3
Ours	T2TT	4-bit	no	70	14.6
Ours	T2TT	4-bit	yes	70	20.2
Ours	S2TT	4-bit	no	70	15.3
Ours	S2TT	4-bit	yes	70	23.9
GPT-4	T2TT	unknown	yes	unknown	1.5
GPT-4	S2TT	unknown	yes	unknown	4.8
Fukuda et al. (2023)	S2TT	16-bit	N/A	1.04	0.7
Papi et al. (2023)	S2TT	16-bit	N/A	0.176	1.4

Table 4: Inference wall time comparisons for our model and selected baselines on en-de. Real-time factor (RTF) less than one means the model is faster than real time.

⁸For T2TT mode, the source audio transcripts were fed directly in the LLM. The RTF was calculated based on the known length of the audio corresponding to the source transcripts and the time to complete translation of that text.

496 F EXAMPLES OF TRANSLATION

497 en-de

498 Source: But in Antarctica, some scientists who are normally pretty levelheaded are getting a little
499 bit freaked out, I would say, is a fair definition, about the lowest level ever, at this point in the year,
500 of sea ice.

501 TransLLaMa: Aber in der Antarktis, einige Wissenschaftler, die normalerweise recht nüchtern sind,
502 werden ein wenig verängstigt, würde ich sagen, ist eine gerechte Definition über den niedrigsten
503 Stand, den wir jemals zu diesem Zeitpunkt im Jahr von Meereis hatten.

504 Back translation: But in Antarctica, some scientists who are normally quite sober are getting a little
505 frightened, I would say, is a fair definition of the lowest level we've ever had of sea ice at this point
506 in the year.

507 Source: I remember when there were so many fossil fuel delegates in Madrid, but by the time we
508 got to Glasgow a year and a half ago, the delegates from the fossil fuel companies made up a larger
509 group than the largest national delegation.

510 TransLLaMa: Ich erinnere mich an die Zeit, als es so viele Delegierte der fossilen Brennstoffind-
511 ustrie in Madrid gab, aber als wir ein Jahr und ein halbes später in Glasgow waren, machten die
512 Delegierten der fossilen Brennstoffunternehmen eine größere Gruppe aus als die größte nationale
513 Delegation.

514 Back translation: I remember when there were so many fossil fuel industry delegates in Madrid, but
515 when we were in Glasgow a year and a half later, the fossil fuel company delegates made up a larger
516 group than the largest national delegation.

517 Source: And I'm a really positive person, so it's actually difficult to say this, but the way we do
518 philanthropy right now, the way we've done it for decades, is broken.

519 TransLLaMa: Ich bin eine sehr positive Person, also ist es für mich eigentlich schwierig, das zu
520 sagen, aber die Art und Weise, wie wir Philanthropie betreiben, die Art und Weise, wie wir es seit
521 Jahrzehnten betreiben, ist kaputt.

522 Back translation: I'm a very positive person, so it's actually hard for me to say this, but the way we
523 do philanthropy, the way we've done it for decades, is broken.

524 Source: Work norms, of course, differ around the world, but we replicated our original American
525 result in South Korea, which is known by the numbers to be one of the hardest-working countries in
526 the OECD, and in France, which is known for other strengths.

527 TransLLaMa: Natürlich unterscheiden sich die Arbeitsnormen in der Welt, aber wir konnten unser
528 ursprüngliches Ergebnis in den USA in Südkorea wiederholen, das bekannt ist als eines der härtesten
529 arbeitenden Länder in der OECD, und in Frankreich, das für andere Stärken bekannt ist.

530 Back translation: Of course, labor standards vary around the world, but we were able to replicate
531 our original result in the US in South Korea, known as one of the hardest working countries in the
532 OECD, and in France, known for other strengths.

533 en-ru

534 Source: A few years ago, I called for toothpaste from my hotel room, and it was delivered to me by
535 a robot.

536 TransLLaMa: Несколько лет назад я попросила зубную пасту из моего номера в
537 отеле и она была доставлена мне роботом.

538 Back translation: A few years ago I asked for toothpaste from my hotel room and it was delivered
539 to me by robot.

540 Source: It conducted the first quantum video call, and it launched the first drone that can carry a
541 passenger in air.

542 TransLLaMa: Она провела первую видео-конференцию с использованием квантовых
543 технологий и запустила первый беспилотный самолет, способный перевозить пас-
544 сажира в воздухе.

545 Back translation: She held the first video conference using quantum technology and launched the
546 first unmanned aircraft capable of carrying a passenger in the air.

547 Source: And it's really gratifying to hear from people like Raymond who are using the technology
548 we are building, and others, for so many wonderful things.

549 TransLLaMa: И это действительно приятно слышать от людей, как Раймонд, кото-
550 рые используют технологии, которые мы создаём, и другие, для таких замеча-
551 тельных вещей.