# From Start to Finish: Latency Reduction Strategies for Incremental Speech Synthesis in Simultaneous Speech-to-Speech Translation

*Danni Liu[1], Changhan Wang[2], Hongyu Gong[2], Xutai Ma[2,3], Yun Tang[2], Juan Pino[2]*

[1]Maastricht University, The Netherlands [2]Meta AI, USA [3]Johns Hopkins University, USA

[1]`danni.liu@maastrichtuniversity.nl`,[2]`{changhan,hygong,yuntang,juancarabina}@fb.com`,[3]`xutai_ma@jhu.edu`

## Abstract

Speech-to-speech translation (S2ST) converts input speech to speech in another language. A challenge of delivering S2ST in real time is the accumulated delay between the translation and speech synthesis modules. While recently incremental text-to-speech (iTTS) models have shown large quality improvements, they typically require additional future text inputs to reach optimal performance. In this work, we minimize the initial waiting time of iTTS by adapting the upstream speech translator to generate high-quality pseudo lookahead for the speech synthesizer. After mitigating the initial delay, we demonstrate that the duration of synthesized speech also plays a crucial role on latency. We formalize this as a latency metric and then present a simple yet effective duration-scaling approach for latency reduction. Our approaches consistently reduce latency by $0.2-0.5$ second without sacrificing speech translation quality.[1]

**Index Terms**: speech translation, text-to-speech, low-latency

## 1. Introduction

Speech-to-speech translation (S2ST) [1, 2, 3] is the task of translating input speech utterances into speech in another language. Compared to translating into text alone, delivering translation into speech offers users with increased accessibility. While recent end-to-end S2ST methods [4, 5, 6, 7] have shown encouraging results, the pipeline approach based on intermediate text representations remains a strong baseline, which currently has not yet been surpassed by its end-to-end counterparts [5, 6, 7]. This motivates the relevance of bringing this approach under the more challenging condition of *simultaneous* S2ST. To enable real-time communication, all components of the S2ST pipeline must be optimized for incremental inference. Here we consider a direct speech-to-text (ST) translation and a text-to-speech (TTS) component. For the ST module, handling the complex mappings from acoustic to textual representations often requires large model size, which comes with high inference computation. Moreover, to account for word reordering in translation, the ST module often requires additional input context [8, 9] before outputting the translated text. Given the well-justified computational and algorithmic delays of the upstream simultaneous ST task, the incremental TTS (iTTS) module must be *lightweight* and *minimal-latency*. For latency reduction in iTTS, current approaches [10, 11, 12, 13] improve the system response speed by minimizing the amount of future input context (or *lookahead*) required by stand-alone TTS systems. In this work, we show that S2ST presents a unique opportunity to circumvent the initial delay from lookahead: we leverage the access to input speech for latency reduction. Moreover, we show that faster system response, or in other words latency reduction at the *start*, does not guarantee low latency when *finish-*
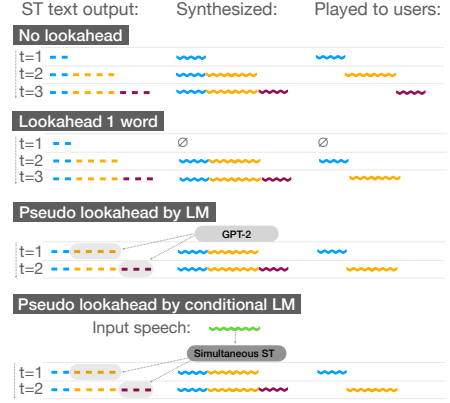
---

[1]`https://github.com/pytorch/fairseq/pull/4184`



Figure 1: *Example of synthesizing $\{w_1, w_2, w_3\}$. The initial delay caused by the lookahead could be mitigated by pseudo lookahead generated by monolingual LMs [11, 12]. We instead adapt the upstream ST system to generate pseudo lookahead.*

*ing*. Due to the non-overlapping constraints of incremental TTS (i.e., no upcoming audio can be played when the current word is still playing), when output speech has stretched duration, the final latency is substantially increased. We explicitly account for this constraint in our framework. Our contributions are:

- We propose a lightweight input-speech-guided pseudo lookahead mechanism that reduces the *starting* latency of iTTS.
- We show that the *duration* of synthesized speech plays a crucial role in *final* latency, and provide a simple yet effective duration-scaling approach that leverages recent advances in non-autoregressive TTS.
- We present an improved iTTS model that differentiates partial and full inputs, which, in combination with the above techniques, consistently reduces latency by $0.2-0.5s$ while retaining speech translation quality.

## 2. Related Work

Initial works on iTTS [14, 10, 15, 16] build upon autoregressive models, notably Tacotron [17, 18], and recently [11] extends the study to non-autoregressive models. An often-discussed topic in iTTS is the role of *lookahead*: [10] use a lookahead [8] for both the text-to-spectrogram model and vocoder, and achieve performance on par with offline systems. [15] show that shorter words need more lookahead. [11, 12, 13] utilize GPT-2 [19] to generate pseudo lookahead for iTTS. Since the lookahead causes initial delays for the iTTS module, current works mainly focus on minimizing its amount. Meanwhile, the relation between the lookahead and latency when *finishing* is still under-explored. Moreover, to the best of our knowledge, no previous work has leveraged source speech to generate lookahead for iTTS.

# 3. Approach

## 3.1. Input-Speech-Guided Pseudo Lookahead

Incremental TTS quality is often improved when allowing a lookahead to some future input words [10, 15, 11], which on the other hand causes delay. Several works [11, 12, 13] use pretrained language models (LMs) to generate pseudo lookahead. While this improves speech quality to some extent, the LM-predicted pseudo lookahead tends to disagree with the ground-truth lookahead [11, 12]. It has further been shown that only the *correct* pseudo lookahead improves output prosodic features [11]. Our initial experiments on controlled replacements of ground-truth lookaheads also confirms this finding.

While pretrained LMs are powerful at many generation tasks, generating the exact next words for TTS is challenging, especially at the beginning of the sentence, where the model can only predict likely words in the *average case*. In S2ST, however, we can additionally leverage the contextual information from the source speech. With a separate LM, assuming greedy decoding for illustrative purposes, the pseudo lookahead $w_t$ at time step $t$ would be $\text{argmax}_w P(w_t|w_{1,...t-1})$, which is only conditioned on the currently available words $w_{1,...,t-1}$. In contrast, the pseudo lookahead generated from the upstream ST system is additionally conditioned on the currently-available source speech $\mathbf{s}$: $\text{argmax}_w P(w_t|w_{1,...t-1}, \mathbf{s})$. We hypothesize the source speech can help us generate more accurate pseudo lookahead, and therefore minimize the performance gap to the case of using ground-truth lookahead. Figure 1 illustrates our idea in comparison to existing methods.

To generate the pseudo lookahead, for every token produced by the upstream ST system, we trigger an additional decoding step. At the next timestep, the decoder continues at the hidden state cached before the pseudo lookahead is generated. This mechanism ensures that the actual ST outputs remain unchanged regardless of pseudo lookahead generation. We note the procedure also applies to multiple pseudo lookahead steps.

Compared to separate LMs, the advantage of our approach is two-fold: First, it is easier to deploy as no extra component is introduced into the translation pipeline. Moreover, from a modeling point of view, this approach has access to the input speech, which likely contains more semantic information on the text to synthesize. We provide more analysis on this in §4.5.

## 3.2. Speech-Duration-Aware Latency Metric

When evaluating iTTS latency, existing works [11, 14, 10, 16] mainly focus on the amount of input text the model requires before starting to produce output audio, e.g. waiting for $X$ words. While this does correspond to the initial system response time, it does not fully account for the real-time constraints in S2ST applications that often translate *consecutive* utterances. As illustrated in Figure 2a, although the TTS module starts with minimal delay immediately after receiving the first text token, the full utterance still finishes comparatively late because the output speech has prolonged duration. Since multiple output words cannot play concurrently, output audios with constantly long durations with can congest the system and cause an ever-increasing latency. In comparison, in Figure 2b, although the model initially waits for an additional word, its final latency is lower thanks to the time saved from shorter output durations (i.e. speaking faster). While admittedly the speaking speed of synthesized speech cannot be considered alone without naturalness and intelligibility, this example showcases the crucial impact of duration on latency.

Motivated by this observation, we adopt a latency metric that account for the time elapsed between end points of the input and output illustrated in Figure 2. Algorithm 1 shows how we measure the latency for an utterance. It depends on three factors: 1) when the input for synthesizing a given word arrived (EMIT_TIME($\cdot$)), 2) computation time for synthesizing the word (COMPUTE_TIME($\cdot$)), 3) when the previous word finished playing (dependent on DURATION($\cdot$)). As the texts from a simultaneous ST system to the TTS module are time-stamped, EMIT_TIME($\cdot$) can be easily derived.



(a) *Start early, finish late*     (b) *Start late, finish early*

Figure 2: *Example of incrementally synthesizing words $w_1$, $w_2$, $w_3$ and the impact of output duration on latency. Lengths of wavy lines indicate duration of synthesized speech.*

---

**ALGORITHM 1** Utterance-level latency for synthesizing a sentence with tokens $\{w_1, \cdots, w_n\}$

---

1: **procedure** CALCULATELATENCY($\{w_1, \cdots, w_n\}$)
2:      $t \leftarrow 0$
3:      **for** $w_i \in \{w_1, \cdots, w_n\}$ **do**
4:          $s \leftarrow$ EMIT_TIME($w_i$)
5:          $p \leftarrow$ MAX($t, s +$ COMPUTE_TIME($w_i$))    ▷ Synthesized audio can only play if: 1) previous word finished playing; 2) current word has been synthesized
6:          $t \leftarrow p +$ DURATION($w_i$)
7:      **end for**
8:      **return** $t -$ EMIT_TIMESTAMP($w_n$)
9: **end procedure**

---

## 3.3. Improving Incremental TTS Module

Having demonstrated the impact of output speech duration on latency, we need a TTS system that enables direct duration control. Among current TTS models, non-autoregressive (NAR) ones are suitable candidates, as the source-target alignment (therefore the duration of phonetic units) come from a dedicated duration predictor. In contrast, with autoregressive models (e.g. [17, 18]), the output speech duration is a by-product of source-target attention and is less straightforward to control. Therefore, we build upon a widely-used NAR TTS model, FastSpeech 2 [20]. Given $j$ current phonemes and $k$ (pseudo) lookahead phonemes, we synthesize up till $X_{j+k}$ but only output those frames up to $X_j$. We use the duration predictor to derive word boundaries in the synthesized output, as also done in [11].

As the TTS module consumes partial inputs at inference time, we augment the training set by prefixes. The prefix augmentation is similar to the approach by [14], except that they use shorter subsets that are not necessarily prefixes[2]. Moreover, as prosodic features tend to differ during and at the end of sentences [21, 22], we use the presence of the end-of-sentence (EOS) token to distinguish between partial and full sentences: Full input sequences are followed by the EOS token while partial sequences are not. We expect the presence of EOS token to signal the model to synthesize the suitable prosody at test time.

---

[2]Our preliminary experiments with this approach did not show consistent improvements as prefix augmentation.

# 4. Experiments

## 4.1. Data

We build S2ST systems for Spanish to English (es→en) and vice versa (en→es). Table 1 shows the dataset statistics.

**ST Data** For es→en, the ST system is trained on the Fisher Spanish-English [23] dataset. For en→es we use MuST-C v1.0 [24] and report performance on `tst-COMMON`. For audio inputs, we extract 80-dimensional filterbank with global ceptral mean variance normalization. For text, we learn byte pair encoding [25] using SentencePiece [26] with size 500 and $10k$ for Fisher and MuST-C respectively. For Fisher, the vocabulary size is intentionally kept small due to the limited corpus size.

**TTS Data** We train incremental TTS models for English and Spanish on LJSpeech [27] and CSS10 Spanish [28] respectively. The phoneme durations are extracted by Montreal Forced Aligner (MFA) [29] using pretrained acoustic models in the MFA toolkit. The phoneme representations are stressed-ARPAbet and GlobalPhone [30] respectively.

Table 1: *Dataset statistics.*

| Task | Corpus | Lang. | Train | | Test | |
|------|--------|-------|-------|-----|------|-----|
| | | | # sent. | hrs | # sent. | hrs |
| ST | Fisher | es→en | $126k$ | 162 | 3,641 | 4.5 |
| | MuST-C | en→es | $270k$ | 504 | 2,052 | 4.2 |
| TTS | LJSpeech | en | $12k$ | 22 | 523 | 1.0 |
| | CSS10 | es | $11k$ | 23 | 107 | 0.2 |

## 4.2. Models and Training

The simultaneous ST system is a Transformer model with a wait-$k$ [8] policy. Following [9], we use a fixed pre-decision of 7 encoder states where each state is an audio frame of $40ms$. As a result, a text token is generated every $0.28s$. The model architecture is the `s2t_transformer_s` model from the FAIRSEQ S2T toolkit [31]. For TTS, we follow the hyperparameters from [20] and train for $200k$ updates. We train HiFi-GAN with configuration $V3$ [32]. When applying prefix augmentation, we keep an equal ratio between prefixes and full sentences. The prefixes are 1/3 or 2/3 of the full sentence lengths.

## 4.3. Evaluation

**Automatic quality evaluation.** We concatenate the incremental speech chunks into full utterances for evaluation. Following [33, 34], we use mean Mel ceptral distortion (MCD) [35] and character error rates (CER) to evaluate TTS quality. MCD is computed between the HiFi-GAN vocoded utterances from the ground-truth and synthesized Mel spectrograms. The English ASR system is a wav2vec 2.0 [36] model. The Spanish ASR system is an XLSR-53 [37] model fine-tuned on Common Voice [38] Spanish data[3]. For speech-to-speech experiments, we transcribe the synthesized utterances with the corresponding ASR system and evaluate BLEU[4] against the reference translation. As the ASR outputs are uncased and unpunctuated, we normalize the reference text accordingly before evaluating BLEU.

---

[3]https://huggingface.co/jonatasgrosman/wav2ve c2-large-xlsr-53-spanish

[4]sacrebleu [39] ID: BLEU+case.mixed+numrefs.1 +smooth.exp+tok.13a+version.1.5.1

Table 2: *iTTS performance. We report 95% confidence interval (CI) over 3 runs for CER, and over valid ratings for MOS. CI for MCD and latency is <0.1 and omitted for brevity.*

| Condition | Quality | | | Latency |
|-----------|---------|---------|------|---------|
| | MCD↓ | CER(%)↓ | MOS↑ | (seconds) |
| **LJSpeech (English)** | | | | |
| (1) Original audio | - | 3.3±0.0 | 4.16±0.15 | - |
| (2) Ground truth Mel | - | 3.2±0.0 | 3.73±0.13 | - |
| (3) Offline | 3.5 | 4.8±0.1 | 3.72±0.14 | 6.2 |
| (4) No lookahead | 4.7 | 8.2±0.2 | 3.02±0.18 | 2.4 |
| (5) Lookahead 1 word | 3.9 | 5.6±0.2 | 3.42±0.16 | 1.4 |
| **After augmentation:** | | | | |
| (6) Offline | 3.6 | 5.0±0.2 | 3.85±0.15 | 6.5 |
| (7) No lookahead | 3.8 | 5.9±0.3 | 3.13±0.18 | 1.1 |
| (8) Lookahead 1 word | 3.7 | 5.1±0.1 | 3.77±0.13 | 1.8 |
| **CSS10-es (Spanish)** | | | | |
| (9) Original audio | - | 4.9±0.0 | 4.78±0.05 | - |
| (10) Ground truth Mel | - | 5.2±0.0 | 4.56±0.07 | - |
| (11) Offline | 3.6 | 5.0±0.3 | 4.34±0.08 | 6.4 |
| (12) No lookahead | 5.1 | 13.2±0.5 | 1.48±0.08 | 2.1 |
| (13) Lookahead 1 word | 3.9 | 8.5±0.9 | 3.04±0.12 | 1.5 |
| **After augmentation:** | | | | |
| (14) Offline | 3.6 | 5.0±0.3 | 4.45±0.07 | 6.4 |
| (15) No lookahead | 3.8 | 7.1±0.5 | 3.11±0.13 | 0.9 |
| (16) Lookahead 1 word | 3.6 | 5.2±0.2 | 3.95±0.10 | 1.5 |

**Subjective quality evaluation.** As human evaluation of S2ST currently still lacks standardized procedures, we only conduct subjective evaluation of TTS quality. We collect mean opinion scores (MOS) on MTurk, where annotators rate speech naturalness on a Likert scale. We evaluate the first 100 utterances of the LJSpeech test set, and the full test set of CSS10. Each utterance receives 15 ratings from different annotators.

**Latency evaluation.** We measure latency using the metric from §3.2. The utterances are synthesized one at a time on an Nvidia V100 GPU. Computation time is included in the latency evaluation. When evaluating the TTS module alone, as we do not directly have timestamps of the input words, we use timestamps derived by forced alignment using MFA [29].

## 4.4. Incremental TTS Results

We first evaluate the incremental TTS module alone (Table 2). Results of the improved model are in lines (6-8) and (14-16).

**Impact of output duration on latency.** In lines (4) and (12), we see that naively decoding substantially degrades quality. Interestingly, although this approach does not wait for any future inputs, its latency is still much higher than other strategies that involve more waiting (lines (5) and (13)). By manual inspection, we confirm the increased latency is caused by prolonged durations of the synthesized audios[5]. In particular, the average length of the synthesized audios is 5.9 and $7.7s$ for lines (5) and (4) respectively, while the ground truth audio length is $6.0s$ on average. As incorporating lookahead leads to more suitable durations, we achieve better latency despite the initial waiting. This phenomenon corresponds to the example in Figure 2, where less lookahead does not necessarily reduce latency.

---

[5]The long-duration phenomenon is also reported in [11].

Table 3: *S2ST quality for es→en (Fisher) and en→es (MuST-C) in upper and lower sections. When using pseudo lookahead from the ST system, we achieve output speech quality on par with ground-truth lookahead without waiting for an extra word.*

| | **iTTS Decoding Strategy** | | | | |
| **Wait-$k$** | Lookahead 1 word | | No lookahead | | +Pseudo |
| | latency($s$)↓ | BLEU↑ | latency($s$)↓ | BLEU↑ | BLEU↑ |
| 3 | 1.3 | 27.5 | 0.9 | 26.0 | 27.1 |
| 5 | 1.4 | 31.5 | 1.1 | 29.8 | 31.3 |
| 7 | 1.6 | 32.7 | 1.3 | 31.1 | 32.5 |
| 9 | 1.7 | 34.5 | 1.5 | 33.0 | 34.3 |
| Offline [6] | | | 39.5 | | |
| 3 | 2.3 | 14.4 | 1.8 | 13.5 | 14.4 |
| 5 | 2.5 | 17.1 | 2.1 | 16.1 | 17.1 |
| 7 | 2.8 | 18.7 | 2.4 | 18.5 | 18.8 |
| 9 | 2.9 | 19.4 | 2.5 | 18.5 | 19.5 |
| Offline [40] | | | 24.4 | | |

**Improved iTTS module.** Table 2 shows the prefix-augmented model outperforms the baseline under all experimented decoding strategies. One reason for the gain is the distinguishing of partial and full inputs (§3.3). With the baseline model, every partial sentence would be synthesized as if it were a full sentence, therefore carrying end-of-sentence prosodic features such as duration stretching [21] and dropping intonations (confirmed by manual inspections). By explicitly indicating full inputs with the EOS tags, we ease the prosody prediction task. Given the favorable results in terms of both quality and latency, we use the prefix-augmented iTTS system in the next S2ST experiments.

### 4.5. Simultaneous Speech-to-Speech Translation Results

In Table 3, we report the performance of S2ST systems of different latency regimes. We include the BLEU score of recent offline systems [6, 40] as quality upper bounds. For each upstream ST model, we compare the performance of the following iTTS decoding strategies: 1) **Lookahead 1 word**: wait for the next word from the ST system; 2) **No lookahead**: directly synthesize currently available words without waiting; 3) **+Pseudo**: directly synthesize without waiting, additionally use pseudo lookahead generated by upstream ST system.

**S2ST quality-latency tradeoff.** In Table 3, by contrasting the first two decoding strategies ("lookahead 1 word" vs "no lookahead"), we clearly observe the quality-latency tradeoff: Without the lookahead, we consistently achieve lower latency of 0.2 to 0.5$s$. This comes with lower output quality, as shown by the loss of ~1.5 BLEU on en→es and ~1.0 BLEU on es→en. By incorporating our proposed ST pseudo lookahead, we are able to match the output quality when using the actual lookahead ("lookahead 1 word"), while avoiding the delay. These findings suggest that our ST-generated pseudo lookahead achieves more efficient quality-latency tradeoff than existing approaches. Moreover, to the best of our knowledge, this is the first evidence that pseudo lookahead matches the performance of ground-truth lookahead in term of output speech quality.

**Importance of accurate pseudo lookahead.** To confirm the gains indeed come from high-quality pseudo lookahead, we evaluate the accuracy of the ST-generated pseudo lookahead, i.e. how many pseudo lookahead tokens agree with the word

generated in the next time step. The accuracy is 70%, 76%, 81%, 83% for Wait-$k = 3, 5, 7, 9$ on Fisher es→en. Our initial experiments with separate LMs yielded low accuracy under 20% and did not improve output speech quality. In Table 4, we report BLEU scores of random pseudo lookahead sampled from the ST model's vocabulary. The consistent degradation (up to 3.1 BLEU vs "ST pseudo lookahead" and 1.9 BLEU vs "no lookahead") indicates that the pseudo lookahead must be accurate to improve synthesis quality. Low-accuracy pseudo lookahead is more harmful than not using any lookahead.

Table 4: *BLEU(↑) scores of different lookahead strategies for speech-to-speech translation experiments on es→en.*

| **Wait-$k$** | $k = 3$ | $k = 5$ | $k = 7$ | $k = 9$ |
|---|---|---|---|---|
| Lookahead 1 word | 27.5 | 31.5 | 32.7 | 34.5 |
| No lookahead | 26.0 | 29.8 | 31.1 | 33.0 |
| ST pseudo lookahead | 27.1 | 31.3 | 32.5 | 34.3 |
| Random lookahead | 24.4 | 28.4 | 28.8 | 31.1 |

### 4.6. Output Duration Control for Latency Reduction

Our upstream wait-$k$ ST system generates tokens at a fixed rate of 0.28s (§4.2). We now study the impact on latency when the tokens come in at a faster rate. As shown in upmost line in Figure 3, a minor difference of 0.06$s$ per token (0.28→0.22$s$) increases latency by 1.3 seconds, as a result of accumulated latency (c.f. Figure 2a). Therefore, we speed up the output utterance and investigate its tradeoff with quality. Recall from §3.3 that the duration prediction module enables us to scale the output duration without altering other prosodic features. As shown by the bottom two lines in Figure 3, when moderately speeding up the output speech by scaling the durations by 0.95 or 0.9, we can counteract the latency induced by faster incoming text tokens. From the quality metrics, this comes with a minimal impact on audio intelligibility and naturalness, with neither CER and MOS differing noticeably from the original values of 5.1% and 3.77. These findings showcase the effectiveness of our iTTS module in controlling duration for latency reduction.
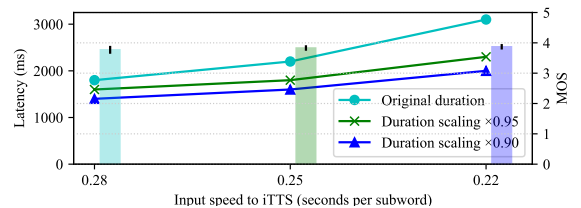


Figure 3: *S2ST latency is sensitive to the speed of incoming text tokens to iTTS. By duration scaling (§3.3), we reduce latency by ~ 0.7s without harming output naturalness measured by MOS.*

## 5. Conclusion

In this work, we improve speech-to-speech translation pipelines for simultaneous speech-to-speech translation. We adapt the upstream speech-to-text translation system to generate high-quality pseudo lookahead for the TTS module, allowing the latter to deliver output translation with minimal delay. After demonstrating the crucial role of duration on latency, we build on recent non-autoregressive models for duration-controllable TTS. We further show our prefix augmentation procedure substantially improves the incremental speech quality.

# 6. References

[1] A. Lavie, A. Waibel, L. S. Levin, M. Finke, D. Gates, M. Gavaldà, T. Zeppenfeld, and P. Zhan, "Janus-iii: speech-to-speech translation in multiple languages," in *Proc. ICASSP*, 1997.

[2] S. Nakamura, K. Markov, H. Nakaiwa, G. Kikui, H. Kawai, T. Jitsuhiro, J. Zhang, H. Yamamoto, E. Sumita, and S. Yamamoto, "The ATR multilingual speech-to-speech translation system," *IEEE Trans. Speech Audio Process.*, 2006.

[3] R. Fukuda, S. Novitasari, Y. Oka, Y. Kano, Y. Yano, Y. Ko, H. Tokuyama, K. Doi, T. Yanagita, S. Sakti, K. Sudoh, and S. Nakamura, "Simultaneous speech-to-speech translation system with transformer-based incremental asr, mt, and tts," in *Proc. O-COCOSDA*, 2021.

[4] A. Tjandra, S. Sakti, and S. Nakamura, "Speech-to-speech translation between untranscribed unknown languages," in *Proc. ASRU*, 2019.

[5] Y. Jia, R. J. Weiss, F. Biadsy, W. Macherey, M. Johnson, Z. Chen, and Y. Wu, "Direct speech-to-speech translation with a sequence-to-sequence model," in *Proc. Interspeech*, 2019.

[6] A. Lee, P.-J. Chen, C. Wang, J. Gu, S. Popuri, X. Ma, A. Polyak, Y. Adi, Q. He, Y. Tang, J. Pino, and W.-N. Hsu, "Direct speech-to-speech translation with discrete units," in *Proc. ACL*, 2022.

[7] Y. Jia, M. T. Ramanovich, T. Remez, and R. Pomerantz, "Translatotron 2: Robust direct speech-to-speech translation," in *Proc. ICML*, 2022.

[8] M. Ma, L. Huang, H. Xiong, R. Zheng, K. Liu, B. Zheng, C. Zhang, Z. He, H. Liu, X. Li, H. Wu, and H. Wang, "STACL: Simultaneous translation with implicit anticipation and controllable latency using prefix-to-prefix framework," in *Proc. ACL*, 2019.

[9] X. Ma, J. Pino, and P. Koehn, "SimulMT to SimulST: Adapting simultaneous text translation to end-to-end simultaneous speech translation," in *Proc. AACL*, 2020.

[10] M. Ma, B. Zheng, K. Liu, R. Zheng, H. Liu, K. Peng, K. Church, and L. Huang, "Incremental text-to-speech synthesis with prefix-to-prefix framework," in *Proc. EMNLP*, 2020.

[11] B. Stephenson, T. Hueber, L. Girin, and L. Besacier, "Alternate Endings: Improving Prosody for Incremental Neural TTS with Predicted Future Text Input," in *Proc. Interspeech*, 2021.

[12] T. Saeki, S. Takamichi, and H. Saruwatari, "Incremental text-to-speech synthesis using pseudo lookahead with large pretrained language model," *IEEE Signal Process. Lett.*, vol. 28, 2021.

[13] T. Saeki, S. Takamichi, and Saruwatari, "Low-latency incremental text-to-speech synthesis with distilled context prediction network," in *Proc. ASRU*, 2021.

[14] T. Yanagita, S. Sakti, and S. Nakamura, "Neural iTTS: Toward Synthesizing Speech in Real-time with End-to-end Neural Text-to-Speech Framework," in *Proc. ISCA SSW*, 2019.

[15] B. Stephenson, L. Besacier, L. Girin, and T. Hueber, "What the future brings: Investigating the impact of lookahead for incremental neural TTS," in *Proc. Interspeech*, 2020.

[16] D. S. R. Mohan, R. Lenain, L. Foglianti, T. H. Teh, M. Staib, A. Torresquintero, and J. Gao, "Incremental text to speech for neural sequence-to-sequence models using reinforcement learning," in *Proc. Interspeech*, 2020.

[17] Y. Wang, R. J. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. V. Le, Y. Agiomyrgiannakis, R. Clark, and R. A. Saurous, "Tacotron: Towards end-to-end speech synthesis," in *Proc. Interspeech*, 2017.

[18] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Ryan, R. A. Saurous, Y. Agiomyrgiannakis, and Y. Wu, "Natural TTS synthesis by conditioning wavenet on Mel spectrogram predictions," in *Proc. ICASSP*, 2018.

[19] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

[20] Y. Ren, C. Hu, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T. Liu, "Fastspeech 2: Fast and high-quality end-to-end text to speech," in *Proc. ICLR*, 2021.

[21] W. E. Cooper and M. Danly, "Segmental and temporal aspects of utterance-final lengthening," *Phonetica*, vol. 38, no. 1-3, pp. 106–115, 1981.

[22] R. Berkovits, "Utterance-final lengthening and the duration of final-stop closures," *Journal of Phonetics*, vol. 21, no. 4, pp. 479–489, 1993.

[23] M. Post, G. Kumar, A. Lopez, D. G. Karakos, C. Callison-Burch, and S. Khudanpur, "Improved speech-to-text translation with the fisher and callhome spanish-english speech translation corpus," in *Proc. IWSLT*, 2013.

[24] M. A. Di Gangi, R. Cattoni, L. Bentivogli, M. Negri, and M. Turchi, "MuST-C: a Multilingual Speech Translation Corpus," in *Proc. NAACL*, 2019.

[25] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," in *Proc. ACL*, 2016.

[26] T. Kudo and J. Richardson, "SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing," in *Proc. EMNLP*, 2018.

[27] K. Ito and L. Johnson, "The LJ Speech dataset," https://keithito.com/LJ-Speech-Dataset/, 2017.

[28] K. Park and T. Mulc, "CSS10: A collection of single speaker speech datasets for 10 languages," in *Proc. Interspeech*, 2019.

[29] M. McAuliffe, M. Socolof, S. Mihuc, M. Wagner, and M. Sonderegger, "Montreal forced aligner: Trainable text-speech alignment using kaldi," in *Proc. Interspeech*, 2017.

[30] T. Schultz and T. Schlippe, "GlobalPhone: Pronunciation dictionaries in 20 languages," in *Proc LREC*, 2014.

[31] C. Wang, Y. Tang, X. Ma, A. Wu, D. Okhonko, and J. Pino, "Fairseq S2T: Fast speech-to-text modeling with Fairseq," in *Proc. AACL*, 2020.

[32] J. Kong, J. Kim, and J. Bae, "HiFi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis," in *Proc. NeurIPS*, 2020.

[33] R. J. Weiss, R. J. Skerry-Ryan, E. Battenberg, S. Mariooryad, and D. P. Kingma, "Wave-Tacotron: Spectrogram-free end-to-end text-to-speech synthesis," in *Proc. ICASSP*, 2021.

[34] C. Wang, W.-N. Hsu, Y. Adi, A. Polyak, A. Lee, P.-J. Chen, J. Gu, and J. Pino, "fairseq s^2: A scalable and integrable speech synthesis toolkit," in *Proc. EMNLP*, 2021.

[35] J. Kominek, T. Schultz, and A. W. Black, "Synthesizer voice quality of new languages calibrated with mean mel cepstral distortion," in *Proc. SLTU*, 2008.

[36] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. NeurIPS*, 2020.

[37] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, "Unsupervised cross-lingual representation learning for speech recognition," in *Proc. Interspeech*, 2021.

[38] R. Ardila, M. Branson, K. Davis, M. Kohler, J. Meyer, M. Henretty, R. Morais, L. Saunders, F. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," in *Proc. LREC*, Marseille, France, 2020.

[39] M. Post, "A call for clarity in reporting BLEU scores," in *Proc. WMT*, 2018.

[40] X. Ma, H. Gong, D. Liu, A. Lee, Y. Tang, P.-J. Chen, W.-N. Hsu, P. Koehn, and J. Pino, "Direct simultaneous speech-to-speech translation with variational monotonic multihead attention," *CoRR*, vol. abs/2110.08250, 2022.