

LOW LATENCY ASR FOR SIMULTANEOUS SPEECH TRANSLATION

*Thai Son Nguyen, Jan Niehues, Eunah Cho, Thanh-Le Ha
Kevin Kilgour, Markus Müller, Matthias Sperber, Sebastian Stüker, Alex Waibel*

Institute for Anthropomatics and Robotics Karlsruhe Institute of Technology

firstname.lastname@kit.edu

ABSTRACT

User studies have shown that reducing the latency of our simultaneous lecture translation system should be the most important goal. We therefore have worked on several techniques for reducing the latency for both components, the automatic speech recognition and the speech translation module. Since the commonly used commitment latency is not appropriate in our case of continuous stream decoding, we focused on word latency. We used it to analyze the performance of our current system and to identify opportunities for improvements. In order to minimize the latency we combined run-on decoding with a technique for identifying stable partial hypotheses when stream decoding and a protocol for dynamic output update that allows to revise the most recent parts of the transcription. This combination reduces the latency at word level, where the words are final and will never be updated again in the future, from 18.1s to 1.1s without sacrificing performance in terms of word error rate.

Index Terms— ASR, Low Latency, Decoding

1. INTRODUCTION

In order for students to be able to follow a lecture by using our system’s automatic lecture transcription and translation, the system’s output needs to be as much in sync with the lecturer’s speech and presentation as possible. Thus, the speech and translation components of the systems do not only need to run in real-time, but must produce output with as low a latency as possible. The high importance of a low latency is also the result of a user study and test that we conducted during real-world operation [1]. This paper addresses the problems of latency measurement and latency reduction for our speech transcription component.

For applications that are turn based and operate on shorter queries, such as Google Voice and Apple Siri, the latency can be measured at utterance level, i.e. the response time after an utterance is finished. In these traditional applications users usually stop and wait for the results. But for our system, that is acting as an interpreter of a continuous, unsegmented stream of speech, the situation is different. In our scenario there is no clear notion of utterance breaks in the speech, and

thus measuring the latency becomes more difficult. The traditional approach for latency measurements in real-time speech recognition systems found in literature uses either real-time factor or delay between the ending time of a segment, e.g. an utterance, and when the recognition result is available. Such measurements might capture the overall speed of the underlying recognition system but not the real latency as perceived by the users. To better capture the user experience, latency measures need to measure the delay as the word sequence is incrementally constructed—thus our focus on the word latency. Further, it is not sufficient to look at an average of the overall latency of a whole test set. Instead we need to look at the variability of the latency and especially at the peak values in the per-word latency that are caused by problematic passages in the input audio. These peaks can lead to occasionally very high local latencies and need to be either avoided or dealt with appropriately.

In this paper we present two approaches to improve the latency of our lecture translator’s speech recognition systems while maintaining its accuracy. The first approach uses a real-time recognition system, utilising an incremental decoding framework to decode continuous audio streams, in combination with a trace back of stable partial hypotheses. This approach is used to output whole portions, i.e. several words in sequence, of the final hypothesis as soon as possible.

The second approach enhances the first one in combination with the display components by allowing to output partial hypotheses not only when they are stable, i.e., when it can be guaranteed that they will not change anymore in the future, but at any time as soon as they are available. In unison with the display and translation components the recognition system is then allowed to correct itself later on, i.e. to revise the most recent history of its output, when a different word sequence has become more probable. Since very often the system will not need to correct itself, but the early output turns out to be the stable one, even though this could not have been guaranteed at the time it was passed on to the translation and display components, the latency of the system is reduced further this way.

Both approaches are shown to reduce the latency of the speech transcription component significantly from 18.1s to 1.1s without trade-off between latency and accuracy.

2. RELATED WORK

On-line speech recognition differs from off-line recognition in that latency is a crucial issue. Although latency in general refers to the response time of the recognition system, it has been defined in different ways in the past. For instance, in dialogue systems such as Google Voice [2], latency is the time from when the user finishes speaking until the search results appear. In other related work on speech recognition for broadcast news [3], latency measurement has included the time for the input to be completed. Note that carefully defining latency is important because only then we can optimize our recognition system in a systematic fashion.

A related but different concept, the *real-time factor* (RTF), is calculated as the ratio between the utterance duration and its required decoding time. RTF is a common measure to evaluate the speed of a speech recognition system. Although distinct from the concept of latency, reducing the RTF can lead to a reduced latency in recognition systems, especially when the decoding starts after the input is completed. Recently, work to improve the latency of Apple’s digital assistant Siri by boosting the pruning behaviour of a deep neural network (DNN) acoustic model [4] resulted in a RTF reduction of 23%. In these systems, there is usually a trade-off between the RTF and accuracy. Larger sizes of pruning beams, acoustic models or language models can lead to a better recognition accuracy at the cost of computing time.

Some recent papers on incremental speech recognition [5, 6] have addressed the latency problem in dialogue systems. The authors conducted a study about the stability which defines how much a word or hypothesis portion remains unchanged over the incremental decoding. Low latency is achieved by early putting out hypotheses when they reached a certain level of stability. However, no results have been published for latency measurements in real applications.

Unlike in dialogue systems, there are no markers for utterance boundaries in continuous recognition systems. To the best of our knowledge, there is not yet a standard approach for latency measurement for continuous recognition system in the literature. A study on broadcast news [3] used the delay between the ending time of a partial output and when the partial recognition is available. While a study on speech recognition in meetings [7] measured the latency as the difference between the end time of a spoken word and the time when the word was output by the speech recogniser.

3. LATENCY DEFINITION

In order to more precisely discuss different types of delays and compare them to the previous studies we define and distinguish: **Word latency** as the difference between the time a spoken word and the time when its transcription is available at the display component. **Commitment latency** is the difference between the end time of an audio segment (Section 4.1)



Fig. 1: An example of hypothesis update.

or portion (Section 4.2) and when its transcription is available at the display component. **Word or commitment peak latency** is the highest word or commitment latency measured on a test set.

4. SPEECH RECOGNITION SYSTEM

4.1. Run-on recognition

Run-on recognition overlaps the decoding process with the audio recording process in order to reduce the latency. Our system used an adapted version of the run-on decoding described in [8]. As in [8] our system uses an audio segmenter in a separate process for pre-processing which writes the incoming audio stream into a shared memory with the decoder, while at the same time filtering out stretches of silence. We refer to everything between two stretches of silence as a *segment*, which are usually several seconds long but could be as long as a lecture.

Our recognition system’s search is re-initialised before processing a new segment and reads the segment’s audio data from the shared memory while the audio segmenter continues to write to it. The audio data read in *chunks* consisting of a fixed number of frames is incrementally decoded. The system therefore only has to wait until a chunk of audio is available in contrast to batch processing which requires complete segments before decoding.

4.2. Stable hypothesis portion

The decoder tries to find the most probable hypothesis. Normally, only at the end of a segment the most probable hypothesis is obtained. However, because waiting for an end-of-segment detected by the segmenter leads to a high latency, we use a partial trace-back [8] for finding stable portions of the hypothesis early. In our design, we detect partial hypotheses right after a chunk has been processed. Whenever a partial hypothesis is detected, its stable portion is extracted and delivered. The end of the portion will be tracked for the next detection.

4.3. Adaptive pruning

Although our recognition system runs, on average, significantly faster than real time, we frequently encounter individual chunks which are processed much slower than real-time. This happens when encountering chunks that are diffi-

cult to decode, e.g. speech with background music or noises in which case the beam might fail to prune away competing paths effectively. This problem results in an unstable response time and introduces latency peaks. To overcome the problem, we use an adaptive pruning scheme. When an audio chunk is processed slower than real-time, we will narrow the beam to reduce the processing time of the following chunks. Once the recognition system has caught up again with the live audio the beam size is set back to its normal size.

4.4. Hypothesis update

Next, we introduce another method that dramatically reduces the latency. We output probable parts of the unstable hypothesis and present them to the user. Later, the recognizer can revise its decision and overwrite the previous output if necessary. In this way, the recognition component does not need to wait until a stable portion or end-of-segment, instead it finds the most probable hypothesis every iteration of the incremental decoding, detects and sends the update portions to the display component.

Figure 1 illustrates how this works in detail. In the example, the incremental decoding was performed 7 times and each time the most probable hypothesis was generated. The updated parts (italic text) were detected each time and sent to the display component. At T4, T6 and T7 the system detected the stable portions (underlined text). These had however already appeared as part of the unstable hypothesis at much earlier times. At T5 and T7, the hypotheses had new start times as described in Section 4.2.

Ignoring the words that are later replaced this algorithm can be seen as inducing a partition of stable hypothesis resulting in an improved latency without any accuracy loss. For example, only at T6 the system was sure about the stable hypothesis portion “wear a flickering cheap”, but the parts of it were already sent, “wear” at T1, “a” at T2, “flickering” at T3 and “cheap” at T4. So the latency is again improved.

4.5. Limiting the length of the partial hypothesis

Since neither the segmenter nor the detection of partial hypothesis provide any guarantee regarding a maximal length of stable portions, we may still encounter situations in which the system does not output anything for a longer period of time. To deal with this issue, a threshold can be applied after which we force an output. If the waiting time exceeds this threshold, we simply output the most probable hypothesis at this point and track the end of the hypothesis.

5. EXPERIMENTS

5.1. System description

We evaluate two baseline systems and three variants using the techniques described above for reducing latency. The first

System	Run-on	AP	PH	Update
<i>Baseline-1</i>	✗	✗	✗	✗
<i>Baseline-2</i>	✓	✓	✗	✗
<i>Portion</i>	✓	✓	✓	✗
<i>Update</i>	✓	✓	✓	✓
<i>Update-NA</i>	✓	✗	✓	✓

Table 1: System summary (AP = Adaptive Pruning, PH = Partial Hypothesis).

System	WER	RTF	Commit. Latency	Word Latency
Baseline-1	18.6	0.51	7.02	18.1
Baseline-2	18.4	0.68	0.92	10.2
Portion	18.5	0.68	1.72	2.10
Update	18.5	0.68	0.83	1.09
Update-NA	18.5	0.71	1.03	1.23

Table 2: Overall performance.

baseline which demonstrates batch processing, waits for completed segments before performing the whole decoding. The segments are generated by our integrated energy based segmenter. In the second baseline, we replace the batch processing with the run-on decoding described in Section 4.1. The decoded results are still produced for whole segments. Run-on decoding is employed in all three of the examined experimental systems.

The first experimental system, labeled *Portion*, uses the algorithm from Section 4.2 for finding stable hypothesis portions. The second variant, called *Update*, applies the update protocol explained in Section 4.4. Both of these utilise adaptive pruning. The third variant, named *Update-NA*, applies the update protocol but without adaptive pruning. A chunk size of 40 frames is used in all run-on systems. Table 1 shows the summary of the applied techniques.

All systems share the same basic setup. It uses a hybrid DNN/HMM acoustic model with log-Mel features. The acoustic model uses a context dependent phoneme setup with three states per polyphone. The DNN has an input window of ± 7 frames, followed by 4 layers of 1,600 neurons and a classification output layer containing just over 8,000 neurons. We used a 4-gram language model with more than 150 thousand words for the decoding system.

5.2. Performance evaluation & Latency measurement

The test data used for the evaluation includes 8 TED talks from the development set of the IWSLT 2015 evaluation campaign. First we evaluate the overall performance by averaging the measurements of WER, RTF, commitment latency and word latency over all talks. RTF is measured by the ratio between the processing time and the length of the processed audio segment. Commitment latency and word latency are measured as defined in Section 3. Secondly, we measure and analyze the peaks in latency using both commitment and word latency.

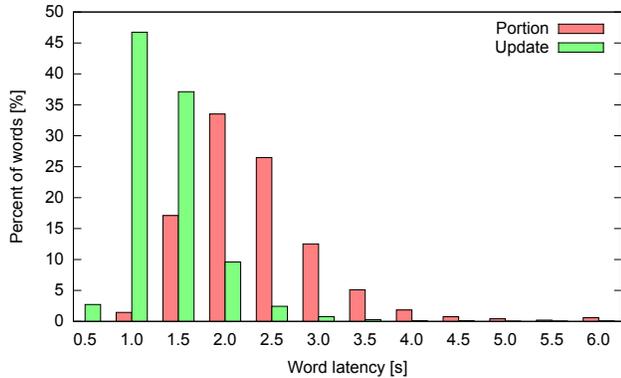


Fig. 2: Word latency distribution.

All audio is fed into the ASR in real-time, as if it were being recorded by a microphone. The display component receives the sequence of uttered words and measures at when the words arrive. The RTF is measured directly by the ASR component, while latency is measured on the display side.

For the systems that send updated parts, we do not consider the latency of the intermediate words which are later replaced. Instead, we measure the latency to the time when a word has been updated the last time. We can then directly compare these measurements to the other systems.

6. RESULTS

6.1. Overall performance

Table 2 shows the performance of all systems on all test talks in terms of overall WER as well as average RTF, commitment latency and word latency.

All the systems have similar WER performance. This confirms that our implemented algorithms did not change the accuracy. The batch processing *Baseline-1* achieves a lower RTF than the other systems that employ run-on processing. This is because it is less efficient for the DNN acoustic model to process multiple smaller chunks than a few large chunks.

Despite its low RTF *Baseline-1* has a large commitment latency since in the batch processing this latency mostly reflects the processing time of the segments. *Portion* has a larger commitment latency than *Baseline-2* and *Update* since it needs to wait until the output can be guaranteed to be stable. *Baseline-2* demonstrates that we can significantly reduce the commitment latency by following the run-on design. Note also that commitment latency, word latency, and RTF are only loosely correlated, indicating that commitment latency and RTF are not sufficient for evaluating the latency of the continuous recognition systems, and justifying our introduction of word latency.

Baseline-2 is especially interesting in this regard, because it has a low commitment latency but a very high word latency. This demonstrates the need for committing recognition results as quickly as possible in order to achieve a low latency. In this sense, *Portion* and *Update* perform better than the others.

System	Max Commit. Latency	Max Word Latency
Baseline-1	93.4	230
Baseline-2	2.7	145
Portion	21.8	23.0
Update	2.7	9.0
Update-NA	19.3	23.5

Table 3: Peak latency.

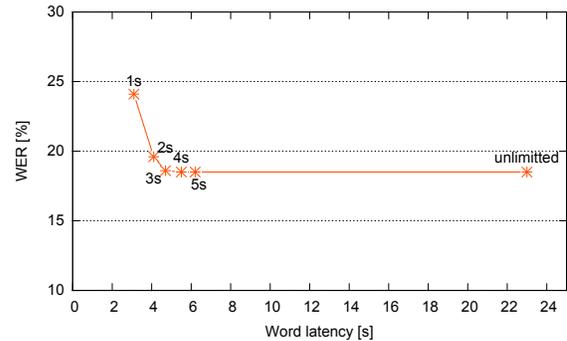


Fig. 3: WER vs. Peak Latency.

As a more detailed analysis, we provide the statistics in Figure 2. It shows the latency distribution of all uttered words in the test set. We only focus on *Portion* and *Update*. According to the diagram, most spoken words are recognised within 2 seconds in *Update*, and 3.5 seconds in *Portion*.

6.2. Peak latency

Looking only at the overall latency, the difference in latency between *Update* and *Update-NA* appears small. However, in practice we noticed that *Update-NA* has a much higher larger peak latency. Table 3 shows the peak latency of all systems. *Portion* required up to 23 seconds to identify a stable hypothesis while the worst word for *Update* was displayed after only 9 seconds. *Update-NA* has a similar overall latency as *Update*, but its peak latency is much worse. The results of *Baseline-2* emphasises the need for word latency measurements.

To explore whether the peak latency can be further improved, we performed an experiment by applying length limitations techniques from Section 4.5 to *Portion*. Figure 3 presents the accuracy and the peak latency of the system at different maximal threshold settings from unlimited over 5 s down to 1 s. We can see that the system imposing a 3s threshold hardly impacts its accuracy.

7. CONCLUSION

We have presented an evaluation for exploring and analysing different problems of latency and latency measurement in our continuous speech recognition system. We have also presented several techniques to deal with these latency problems. The latency improvement not only enhances the usability of our lecture translation system, but also enables the transcriptions to be in sync with the slides and gestures of the lecturer.

8. REFERENCES

- [1] Markus Müller, Sarah Fünfer, Sebastian Stüker, and Alex Waibel, “Evaluation of the kit lecture translation system,” in *Proc. of LREC*, 2016.
- [2] Johan Schalkwyk, Doug Beeferman, Françoise Beaufays, Bill Byrne, Ciprian Chelba, Mike Cohen, Maryam Garret, and Brian Strope, “Your Word is my Command: Google search by voice: A case study,” in *Advances in Speech Recognition*, pp. 61–90. Springer, 2010.
- [3] Murat Saraclar, Michael Riley, Enrico Bocchieri, and Vincent Goffin, “Towards automatic closed captioning: low latency real time broadcast news transcription,” in *Proc. of INTERSPEECH*, 2002.
- [4] Matthias Paulik, “Improvements to the pruning behavior of dnn acoustic models,” in *Proc. of INTERSPEECH*, 2015.
- [5] Ian McGraw and Alexander Gruenstein, “Estimating word-stability during incremental speech recognition,” *Training*, vol. 17, no. 27,327, pp. 6–4, 2011.
- [6] Ethan Selfridge, Iker Arizmendi, Peter A Heeman, and Jason D Williams, “Stability and accuracy in incremental speech recognition,” in *Proc. of SIGDIAL*, 2011.
- [7] Takaaki Hori, Shoko Araki, Takuya Yoshioka, Masakiyo Fujimoto, Shinji Watanabe, Takanobu Oba, Atsunori Ogawa, Kazuhiro Otsuka, Dan Mikami, Keisuke Kinoshita, et al., “Low-latency real-time meeting recognition and understanding using distant microphones and omni-directional camera,” *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 20, no. 2, pp. 499–513, 2012.
- [8] Christian Fügen, *A system for simultaneous translation of lectures and speeches*, Ph.D. thesis, Universität Karlsruhe (TH), 2008.