

# Transformer-based Automatic Speech Recognition of Formal and Colloquial Czech in MALACH Project

Jan Lehečka<sup>1</sup>[0000-0002-3889-8069], Josef V. Psutka<sup>1</sup>[0000-0003-4761-1645], and  
Josef Psutka<sup>1</sup>[0000-0002-0764-3207]

Department of Cybernetics, University of West Bohemia in Pilsen, Czech Republic  
{jlehecka,psutka\_j,psutka}@kky.zcu.cz

**Abstract.** Czech is a very specific language due to its large differences between the formal and the colloquial form of speech. While the formal (written) form is used mainly in official documents, literature, and public speeches, the colloquial (spoken) form is used widely among people in casual speeches. This gap introduces serious problems for ASR systems, especially when training or evaluating ASR models on datasets containing a lot of colloquial speech, such as the MALACH project. In this paper, we are addressing this problem in the light of a new paradigm in end-to-end ASR systems – recently introduced self-supervised audio Transformers. Specifically, we are investigating the influence of colloquial speech on the performance of Wav2Vec 2.0 models and their ability to transcribe colloquial speech directly into formal transcripts. We are presenting results with both formal and colloquial forms in the training transcripts, language models, and evaluation transcripts.

**Keywords:** Wav2Vec 2.0 · Colloquial speech · ASR.

## 1 Introduction

Formal Czech differs a lot from the colloquial Czech. Almost 20% of Czech words have different transcription in both varieties [17, p. 250]. This gap between the everyday, colloquial language, and the official codified formal language emerged during the Czech National Revival back in the 1830s when a group of Czech writers, poets, translators, editors, and teachers established new grammar rules and vocabularies independent of German influence. They took inspiration from other Slavic languages and outdated Czech Bible texts. However, common people did not adopt these new rules and words into their spoken language creating a very specific widely-spoken vernacular that persists to this day [7].

The gap between formal and colloquial Czech constitutes a serious problem for Automatic Speech Recognition (ASR) systems which automatically transcribe – possibly colloquial – spoken utterances into formal text [4]. The usual way how to deal with this phenomenon in a common Large-Vocabulary Continuous Speech Recognition (LVCSR) system is to train the acoustic model with

colloquial phonetic transcripts, define alternative (colloquial) pronunciations for formal words in the lexicon and finally use a formal language model to decode the speech into a formal transcript [14,16].

In the recent few years, self-supervised neural networks became a very popular alternative to LVCSR systems in speech recognition tasks. A significant milestone was the introduction of the Transformer architecture [18] into ASR systems [3,2,11,12,5]. The most studied transformer-based ASR model architecture is Wav2Vec 2.0 [3]. It is an end-to-end speech recognizer that alleviates the need for word pronunciation modeling and does not require any alignment of data. It is a single model converting the raw audio signal from the input into the sequence of tokens on the output, no matter whether these tokens are graphemes, phonemes, word pieces, or other speech units. Thus, the model has a very interesting ability: when the input audio data during fine-tuning contain colloquial speech and the target transcripts are in the formal Czech, it could internally learn the mapping between the two forms without any engineering or manual effort. In this paper, we are investigating the extent of this ability of Wav2Vec models.

## 2 MALACH Project

The whole story of the MALACH project began in 1994, when after the premiere of the film “Schindler’s List”, many survivors turned to Steven Spielberg to tell him their stories about the Holocaust. Inspired by these requests, Spielberg decided to establish the Shoah Visual History Foundation (VHF) so that as many survivors as possible could record their stories and save them for future generations. Nowadays are these video interviews located in the Shoah Foundation Institute at the University of Southern California (USC-SFI) along with another 54,000 video interviews with witnesses of the history of the entire 20th century.

The Shoah part of the archive contains testimonies in 32 languages of personal memories of survivors of the World War II Holocaust, in total it is 116,000 hours of video. Interviews (in all languages) contain natural, unrestricted speech, full of disfluencies, emotional excitements, heavy accents, and are often influenced by the high age of speakers (problems with keeping ideas). More than 550 testimonies are in the Czech (almost 1000 hours of video).

In 2014, the Linguistic Data Consortium (LDC) released the Czech part of the MALACH project [15]. There were published 420 testimonies along with their transcripts. The release contains 400 randomly selected testimonies for the purpose of acoustic model training. As only 15-minute segments were transcribed for each testimony, the acoustic training part, therefore, consists of 100 hours of Czech speech from theoretically up to 800 speakers (interviewer and interviewee). The rest of the Czech MALACH corpus consists of 20 testimonies, which have been completely transcribed and are intended for development (10 testimonies, i.e. 20 speakers) and testing (10 testimonies, i.e. 20 speakers) purposes. (see Tab. 1 for details).

**Table 1.** Statistics of training and test data-sets of the Czech part of the MALACH project.

	Train	Test
# of speakers	776	20
# of words	49k	10.3k
# of tokens	715k	63k
dataset length [hours]	87.5	8.9

### 3 Formal vs. Colloquial Czech

During the annotation process of the Czech Malach corpus, the transcribers were instructed to use the orthographic transcription of colloquial words (i.e., not to “formalize” them artificially) to bring the transcripts as close as possible to what was actually said. There were several reasons for this decision. Firstly, this procedure was very beneficial for classical acoustic modeling, because the resulting transcription is very close to the actual phonetic realization of the word. Secondly, transcribing colloquial sentences using formal words is not an easy task, especially for transcribers without a solid linguistic background. Another problem solved by the colloquial method of transcription was no need to unify the transcription of foreign words.

On the other hand, the effect of the abundance of colloquial words on the language model is rather negative. The orthographic transcription of colloquial words causes an unnecessary growth of the lexicon. There are often several different colloquial variants corresponding to one formal word form. Consequently, the already sparse language model training data became even sparser. To take advantage of formal word forms in language modeling, we decided to “formalize” the lexicon. We went through a lexicon built from the original (orthographic) transcriptions and added a corresponding standard form to each colloquial word form, but only in cases where it was unambiguous. The normalization of manual transcripts not only made the parameters of the estimated language model more robust but also brought this main and most useful source for language modeling much closer to other potential formal text sources. More details on this process can be found in [13].

A good example of the ambiguity of such a formalization is the word *sem*. While in formal Czech this word means *sem (here)*, in colloquial Czech it is also used instead of the correct form *jsem ((I) am)* which naturally occurs quite frequently (the fourth most frequent word in the corpus). To distinguish which formal variant of a word *sem* is the correct one, we would have to use a larger word context or better use a sophisticated method of text understanding. Nevertheless, by formalizing the lexicon, we found more than 13k unambiguous rules that reduced the number of colloquial words by almost 85%.

In order to illustrate that the number of colloquial forms for a single formal word form can be really high, we present a fragment from the “formalized” lexicon in Tab. 2. The new “formalized” text corpus was created by automatically replacing colloquial words in the original transcripts with their formal counter-

**Table 2.** Example of formalization rules.

formal	colloquial	<i>in English</i>
odjet	odejet odjec odject vodjet vodejet vodject vodeject	<i>to leave</i>
odtamtuđ	odtamtađ odtamtud vodtamtađ vodtamtuđ vodtamtuđ votamtađ votamtud	<i>from there</i>
bývalý	bejvalej bejvalý bývalej	<i>former</i>

parts using the above-mentioned 2-column lexicon. Note that such a procedure does not take into account the word context, and therefore the formalization process is far from perfect.

## 4 Wav2Vec 2.0

Wav2Vec 2.0 model [3] is one of the current state-of-the-art models for ASR. It is a deep neural network pretrained to reconstruct the corrupted signals. The input raw audio signal is processed by a multi-layer convolutional neural network into a sequence of latent-speech representations which are fed into a multi-layer Transformer [18]. Only the encoder part of the full encoder-decoder Transformer architecture is used. The output of the Transformer is a sequence of frame-level contextualized speech representations encoding both the frame itself and its context in the signal. This approach is motivated by very successful self-supervised text-based Transformers solving Natural Language Processing (NLP) and Natural Language Understanding (NLU) tasks [8].

The training of Wav2Vec models consists of two phases: pretraining and fine-tuning. During the first self-supervised pretraining phase, the model learns contextualized speech representations from large-scale unlabeled audio datasets. This approach is motivated by the learning skills of infants, who do not learn to understand speech by reading its transcripts, but rather by listening to adults around them and trying to catch the meaning from the context. By masking latent representations of the raw waveform and solving a contrastive task over quantized speech representations, the model learns contextualized representations jointly with discrete speech units without the need for any annotations or labels.

Since labeled data could be very expensive and precious, the pretraining phase equips the model with deep knowledge about the speech signals mined out from tens of thousands of hours of unlabeled speech. This knowledge constitutes a great advantage over models trained from scratch using labeled data only. From this point of view, the pretrained weights of the Wav2Vec model could be seen as very clever initializations of the model weights for supervised training.

During the second supervised fine-tuning phase, the model transfers the pre-trained knowledge into the ASR task. For input speech signals, the speech representations are fed into Connectionist Temporal Classification (CTC) layer [9]

and the most probable sequences of graphemes are decoded. The model is fine-tuned with frozen feature-encoder weights from labeled data optimizing the CTC loss.

CTC is an alignment-free method for grouping audio frames belonging to the same output token in order to convert a sequence of speech representations (one per audio frame) into a much shorter sequence of output tokens. The CTC classification process can be described – in a simplified way – in 3 steps: (1) assign the most probable output token to each audio frame, (2) group sequences with the same tokens into a single token, and (3) remove blank tokens. Tokens are usually graphemes (i.e. characters including also a word delimiter) but could be any speech units.

## 5 Experimental Setup

### 5.1 Pretraining

Public monolingual Wav2Vec models for non-English languages are very rare. For the Czech language, there are none. However, there are several public multilingual pretrained models of sizes from large [6] to extremely large [1]. These models included also Czech in the pretraining datasets. The common practice with these models is to select the most suitable pretrained model and fine-tune it on the labeled ASR data from the target language. Since we were not satisfied with results from multilingual models and, at the same time, we had access to large unlabeled datasets and a high-performance GPU cluster, we decided to pretrain our own base-sized monolingual Wav2Vec model from scratch and released it to the public.

Self-supervised audio transformers are known to scale well with the size of pretraining data, even with extremely huge datasets [1]. Hence, we tried to gather as much public and in-house unlabeled audio data as possible. Together, we were able to collect more than 80 thousand hours of Czech speech. The collection includes recordings from radio (22k hours), unlabeled data from VoxPopuli dataset [19] (18.7k hours), TV shows (15k hours), shadow speakers (12k hours), sports (5k hours), telephone data (2k hours), and a smaller amount of data from several other domains. We included also raw unlabeled audio data from the MALACH project (1k hours).

Since the feature extraction of the input signal is limited by the memory of GPUs in use, we sliced all records not to exceed 30s, which we found to be a reasonable input size for batching.

We followed the same pretraining setup as for the base Wav2Vec 2.0 model in [3]. We pretrained the model for 400 thousand steps with a batch size not exceeding 1.6 hours, corresponding to more than 11 epochs over the dataset. The pretraining took about two weeks on a machine with four NVIDIA A100 GPUs. We released our pretrained model under the nickname *CITRUS* (abbreviation for **C**zech language **TR**ansformer from **U**nabeled **S**peech) for public

non-commercial use<sup>1</sup>. We are not aware of any similar model for Czech mentioned in the literature so far.

## 5.2 Fine-tuning

When fine-tuning models, we used the same setup as in [3], i.e. we trained the pretrained model for 80 thousand update steps with the peak learning rate of  $2 \times 10^{-5}$  and the batch size about 27 minutes of audio, resulting in 270 training epochs over the dataset. We removed non-speech events and punctuation from the transcripts and mapped texts into lowercase. We used implementation from the `Fairseq` tool<sup>2</sup> to fine-tune models.

First, we trained the colloquial model, denoted as  $W2V_{\text{colloq}}$ , from the original transcripts. Since annotators were instructed to transcribe the speech in the spoken form, i.e. exactly as it was spoken in the underlying speech, these transcripts are mainly in colloquial Czech. However, it is in fact a mix of both forms, because some people tend to speak more formally when giving an interview, and sometimes annotators were not able to distinguish between the two forms, especially in the strong emotional and heavily accented speeches. We left the formal words untouched as the rules from formal to colloquial form would be ambiguous.

After that, we transformed the original transcripts into formal Czech using the prepared set of rules (see Sec. 3) and fine-tuned the second model, denoted as  $W2V_{\text{formal}}$ . The whole fine-tuning process is depicted in the upper part of Fig. 1. The fine-tuning of each model took about 14 hours on a machine with four NVIDIA A100 GPUs.

## 5.3 Decoding

We studied two different decoding setups: (1) Connectionist Temporal Classification (CTC) [9], which is the training loss we used during fine-tuning of the models, and (2) CTC beam search decoder with a Language Model (LM). Decoding setup (1) is a grapheme-based lexicon-free speech recognition without any language constraints. The only orthography-related knowledge the model could learn is the training transcripts fed in during the fine-tuning. Wav2Vec with the CTC decoding setup (1) decodes also word delimiters, so it is an end-to-end ASR system, which can be evaluated using standard word-based metrics like word error rate.

Decoding setup (2) incorporates an LM into the CTC beam search decoder which usually improves the speech recognition accuracy by bringing useful language information into the decoding process and penalizing improbable outputs. For our experiments, we prepared 2 different word-based n-gram LMs: (a)  $LM_{\text{colloq}}$  trained from all colloquial transcripts, and (b)  $LM_{\text{formal}}$  trained from the

<sup>1</sup> Available at <https://huggingface.co/fav-kky/wav2vec2-base-cs-80k-CITRUS>

<sup>2</sup> <https://github.com/pytorch/fairseq>

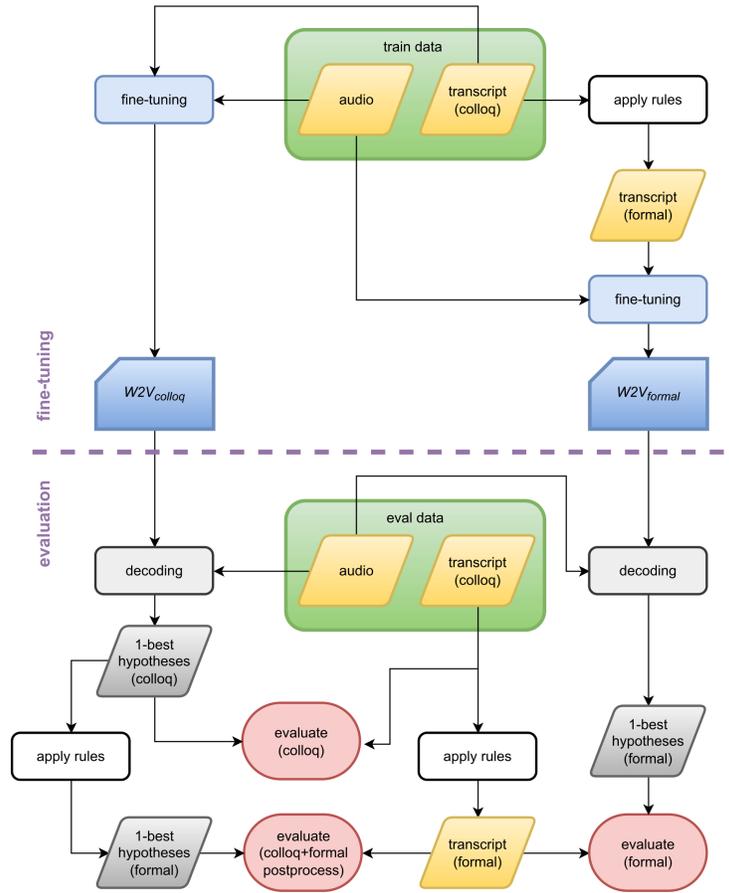


Fig. 1. Scheme of fine-tuning and evaluation.

formalized training transcripts, i.e. from the training data of  $W2V_{\text{formal}}$  model. We limited the maximum order of models to 4-grams for both LMs.

We used implementation from **Transformers** [20] for CTC decoding and `pyctcdecode`<sup>3</sup> decoder for CTC beam search decoder with n-gram LM. To train LMs, we used `KenLM` [10] and mapped all texts into lowercase.

#### 5.4 Evaluation

Both decoding setups described in Sec. 5.3 generated a 1-best hypothesis for each input signal. We aligned decoded hypotheses and reference transcripts using the

<sup>3</sup> <https://github.com/kensho-technologies/pyctcdecode>

minimum edit distance and evaluated the standard Word Error Rate (WER) and Character Error Rate (CER).

To evaluate colloquial models, we used the original reference transcripts processed in the same way as the training transcripts, i.e. we removed non-speech events and punctuation and mapped texts into lowercase. We denote test dataset with this colloquial reference as  $\text{TEST}_{\text{colloq}}$ . To evaluate formal models, we further converted the colloquial reference texts into formal texts using the prepared set of rules (see Sec. 3) and thus generated the test dataset with formal reference transcripts, denoted as  $\text{TEST}_{\text{formal}}$ .

We evaluated all combinations of formal and colloquial models and LMs against both formal and colloquial reference transcripts. From these combinations, we are particularly interested in three real-world scenarios:

1. Evaluation of the colloquial model, i.e. how well  $\text{W2V}_{\text{colloq}}$  model with  $\text{LM}_{\text{colloq}}$  transcribes the colloquial speech (thus evaluated against  $\text{TEST}_{\text{colloq}}$  dataset).
2. Evaluation of the formal model, i.e. how well  $\text{W2V}_{\text{formal}}$  model with  $\text{LM}_{\text{formal}}$  transcribes the colloquial speech into formal Czech (thus evaluated against  $\text{TEST}_{\text{formal}}$  dataset). This scenario is particularly interesting as it evaluates how well the Wav2Vec model internally learns the mapping between the two forms without any engineering or manual effort.
3. Transcripts generated from  $\text{W2V}_{\text{colloq}}$  with  $\text{LM}_{\text{colloq}}$  post-processed by rule-based formalization of texts evaluated against  $\text{TEST}_{\text{formal}}$  dataset. This scenario shows how the Wav2Vec model can use data prepared with a great manual effort for a standard LVCSR system in order to generate formal transcripts. We denote this colloquial model with Formalization Post-processing as  $\text{W2V}_{\text{colloq}}+\text{FP}$

These three scenarios are depicted in a flowchart diagram in the bottom part of Fig. 1 and corresponding error rates will be underlined in the results table.

Note, that the numbers of reference words in  $\text{TEST}_{\text{colloq}}$  and  $\text{TEST}_{\text{formal}}$  differ due to multi-word replacements in the rules. While the formal transcripts consist of 62 690 words, the colloquial has 62 918 words, so results evaluated against  $\text{TEST}_{\text{formal}}$  and  $\text{TEST}_{\text{colloq}}$  are not exactly comparable.

## 6 Results

Results of our experiments are tabulated in Tab. 3. First, we evaluated the existing LVCSR system developed specifically for MALACH dataset [16]. The system was a CNN-TDNN LF-MMI with iVectors, sMBR criterion, and system-specific 3-gram LM denoted as  $\text{LM}_{\text{LVCSR}}$ . The system was trained to transcribe colloquial speech into formal form, so we report only results evaluated against  $\text{TEST}_{\text{formal}}$ . A comparison of this system with the formal Wav2Vec model clearly reveals the superiority of transformer-based ASR systems.

As for the Wav2Vec models, the best results evaluated against  $\text{TEST}_{\text{colloq}}$  are significantly higher (i.e. worse) than best results evaluated against  $\text{TEST}_{\text{formal}}$ . It is mainly because the colloquial Czech does not have codified rules and one

**Table 3.** WER [%] and CER [%] of colloquial and formal models evaluated against colloquial and formal evaluation datasets ( $\text{TEST}_{\text{colloq}}$  and  $\text{TEST}_{\text{formal}}$ ). Each Wav2Vec model was decoded using three different decoding setups: as an end-to-end ASR with no LM and with the beam search CTC decoder with  $\text{LM}_{\text{formal}}$  and  $\text{LM}_{\text{colloq}}$  (see Sec. 5.3). Underlined values correspond to scenarios we are particularly interested in (see Sec. 5.4). Bold values are the best error rates for each model.

model	LM	TEST <sub>colloq</sub>		TEST <sub>formal</sub>	
		WER	CER	WER	CER
LVCSR	LM <sub>LVCSR</sub>	-	-	14.71	5.25
W2V <sub>colloq</sub>	-	12.24	<b>3.58</b>	19.73	5.28
	LM <sub>formal</sub>	13.85	4.05	15.96	4.68
	LM <sub>colloq</sub>	<u>11.55</u>	<u>3.64</u>	18.99	5.27
W2V <sub>formal</sub>	-	19.17	5.07	11.52	3.32
	LM <sub>formal</sub>	18.60	5.19	<u>10.48</u>	<b>3.31</b>
	LM <sub>colloq</sub>	18.60	5.16	10.85	3.37
W2V <sub>colloq</sub> +FP	-	19.02	5.05	11.18	3.33
	LM <sub>formal</sub>	18.47	5.07	11.09	3.53
	LM <sub>colloq</sub>	18.47	5.12	<u>10.43</u>	<b>3.30</b>

formal word could have many possible colloquial forms. Each speaker can use – based on his or her geographical background – a different set of colloquial words in the speech. Moreover, each annotator can perceive the spoken colloquial forms differently, especially in the strong emotional and heavily accented speeches. This ambiguity of transcribed speech leads to confusion when training and evaluating the colloquial models.

If we compare the underlined results of the last two Wav2Vec models in Tab. 3 (corresponding to scenarios 2. and 3. from Sec. 5.4), we see very similar error rates. The W2V<sub>colloq</sub>+FP is slightly better, which we found to be caused by an occasional incorrect exact match of formalized hypotheses with the formalized reference, as both were generated using the same rules. After analyzing errors from W2V<sub>formal</sub> model, we found that many recognition errors were actually errors in the reference as the rules were not covering all occurrences of colloquial form in the reference. For example, the formal reference contained (incorrectly) the word “německýho” (colloquial inflected form meaning “German”), because it was not covered by mapping rules due to its non-existence in training transcripts. Formalized output from W2V<sub>colloq</sub>+FP exactly matched the reference for the same reason, so there was no error counted. W2V<sub>formal</sub> predicted the correct formal form “německého”, which was, however, wrongly counted as a recognition error due to an error in the reference. We didn’t make more effort to clean the reference transcripts and fix these errors as they were infrequent and it would cost a lot of manual work with only a little effect on the error rates. Nevertheless, observing these types of errors was a clear sign of the generalization ability of the W2V<sub>formal</sub> model and we can conclude that W2V<sub>formal</sub> is – despite slightly higher error rates

– a more useful model than rule-based  $W2V_{\text{colloq}}+FP$  because of its generalization ability.

To sum up the results, Wav2Vec models are significantly better ASR systems for the MALACH project than LVCSR systems. They are able to learn the mapping from colloquial speech into a formal transcript and generalize this skill also to words not observed in training data, which is a more beneficial solution than limited rule-based formalization post-processing of the colloquial model. Moreover, the Wav2Vec’s internal mapping from colloquial speech to formal transcripts could make the acquisition of training transcripts much simpler as the annotators could be instructed to transcribe the speech directly into formal Czech alleviating the problems with ambiguous colloquial transcripts and manual listing of rules.

## 7 Conclusion

In this paper, we showed that the new paradigm models in ASR – Transformer-based models with CTC decoder (specifically Wav2Vec 2.0) – have a very interesting ability to learn how to transcribe Czech colloquial speech directly into formal transcripts. Such models not only perform better than common LVCSR systems, but also alleviate the need for complicated and ambiguous colloquial annotations, data alignments, phonetic transcriptions, and pronunciation lexicons. When collecting training transcripts for a new ASR dataset, we can instruct annotators just to transcribe the speech directly into formal Czech sentences, which is codified and unambiguous form, and that’s all that is needed for the Wav2Vec model to be fine-tuned. From the formal transcript and raw audio signal, the model is able to learn the alignment between the speech signal frames and graphemes, and also how to generalize the conversion between the colloquial speech and formal text. We believe our findings will simplify and accelerate the acquisition of training data for new challenging datasets containing a lot of colloquial speech.

**Acknowledgments.** This research was supported by the ITI project of the Ministry of Education of the Czech Republic CZ.02.1.01/0.0/0.0/17 048/0007267 In-teCom. Computational resources were supplied by the project ”e-Infrastruktura CZ” (e-INFRA CZ LM2018140 ) supported by the Ministry of Education, Youth and Sports of the Czech Republic.

## References

1. Babu, A., Wang, C., Tjandra, A., Lakhotia, K., Xu, Q., Goyal, N., Singh, K., von Platen, P., Saraf, Y., Pino, J., et al.: Xls-r: Self-supervised cross-lingual speech representation learning at scale. arXiv preprint arXiv:2111.09296 (2021)
2. Baevski, A., Mohamed, A.: Effectiveness of self-supervised pre-training for ASR. ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) pp. 7694–7698 (2020)

3. Baevski, A., Zhou, Y., Mohamed, A., Auli, M.: Wav2Vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems* **33**, 12449–12460 (2020)
4. Byrne, W., Doerman, D., Franz, M., Gustman, S., Hajič, J., Oard, D., Picheny, M., Psutka, J., Ramabhadran, B., Soergel, D., Ward, T., Zhu, W.: Automatic recognition of spontaneous speech for access to multilingual oral history archives. *IEEE Transactions on Speech and Audio Processing* **12**(4), 420–435 (2004). <https://doi.org/10.1109/TSA.2004.828702>
5. Chen, S., Wang, C., Chen, Z., Wu, Y., Liu, S., Chen, Z., Li, J., Kanda, N., Yoshioka, T., Xiao, X., et al.: WavLM: Large-scale self-supervised pre-training for full stack speech processing. *arXiv preprint arXiv:2110.13900* (2021)
6. Conneau, A., Baevski, A., Collobert, R., Mohamed, A., Auli, M.: Unsupervised cross-lingual representation learning for speech recognition. In: Hermansky, H., Cernocký, H., Burget, L., Lamel, L., Scharenborg, O., Motlíček, P. (eds.) *Interspeech 2021, 22nd Annual Conference of the International Speech Communication Association, Brno, Czechia, 30 August - 3 September 2021*. pp. 2426–2430. ISCA (2021). <https://doi.org/10.21437/Interspeech.2021-329>, <https://doi.org/10.21437/Interspeech.2021-329>
7. Cummins, G.M.: Literary Czech, common Czech, and the instrumental plural. *Journal of Slavic Linguistics* **13**(2), 271–297 (2005), <http://www.jstor.org/stable/24599659>
8. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018)
9. Graves, A., Fernández, S., Gomez, F., Schmidhuber, J.: Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In: *Proceedings of the 23rd international conference on Machine learning*. pp. 369–376 (2006)
10. Heafield, K.: KenLM: Faster and smaller language model queries. In: *Proceedings of the Sixth Workshop on Statistical Machine Translation*. pp. 187–197. Association for Computational Linguistics, Edinburgh, Scotland (Jul 2011), <https://aclanthology.org/W11-2123>
11. Hsu, W.N., Bolte, B., Tsai, Y.H.H., Lakhotia, K., Salakhutdinov, R., Mohamed, A.: Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **29**, 3451–3460 (2021)
12. Liu, A.T., Li, S.W., Lee, H.y.: TERA: Self-supervised learning of transformer encoder representation for speech. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **29**, 2351–2366 (2021)
13. Psutka, J., Ircing, P., Hajič, J., Radová, V., Psutka, J.V., Byrne, W., Gustman, S.: Issues in annotation of the Czech spontaneous speech corpus in the MALACH project. In: *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*. pp. 607–610. European Language Resources Association, Lisbon (2004)
14. Psutka, J., Ircing, P., Psutka, J.V., Hajič, J., Byrne, W., Mírovský, J.: Automatic transcription of Czech, Russian and Slovak spontaneous speech in the MALACH project. In: *Eurospeech 2005*. pp. 1349–1352. ISCA (2005)
15. Psutka, J., Radová, V., Ircing, P., Matoušek, J., Müller, L.: USC-SFI MALACH Interviews and Transcripts Czech LDC2014S04 (2014), <https://catalog.ldc.upenn.edu/LDC2014S04>

16. Psutka, J.V., Pražák, A., Vaněk, J.: Recognition of heavily accented and emotional speech of English and Czech Holocaust survivors using various DNN architectures. In: Karpov, A., Potapova, R. (eds.) *Speech and Computer*. pp. 553–564. Springer International Publishing, Cham (2021)
17. Tahal, K.: *A Grammar of Czech as a foreign language*. FACTUM CZ, s.r.o. (2010)
18. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. p. 6000–6010. NIPS'17, Curran Associates Inc., Red Hook, NY, USA (2017)
19. Wang, C., Riviere, M., Lee, A., Wu, A., Talnikar, C., Haziza, D., Williamson, M., Pino, J., Dupoux, E.: VoxPopuli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation. In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. pp. 993–1003. Association for Computational Linguistics, Online (Aug 2021), <https://aclanthology.org/2021.acl-long.80>
20. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T.L., Gugger, S., Drame, M., Lhoest, Q., Rush, A.M.: Transformers: State-of-the-art natural language processing. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. pp. 38–45. Association for Computational Linguistics, Online (Oct 2020), <https://www.aclweb.org/anthology/2020.emnlp-demos.6>