

# AVOID OVERTHINKING IN SELF-SUPERVISED MODELS FOR SPEECH RECOGNITION

Dan Berrebbi, Brian Yan, Shinji Watanabe

Carnegie Mellon University

## ABSTRACT

Self-supervised learning (SSL) models reshaped our approach to speech, language and vision. However their huge size and the opaque relations between their layers and tasks result in slow inference and *network overthinking*, where predictions made from the last layer of large models is worse than those made from intermediate layers. Early exit (EE) strategies can solve both issues by dynamically reducing computations at inference time for certain samples. Although popular for classification tasks in vision and language, EE has seen less use for sequence-to-sequence speech recognition (ASR) tasks where outputs from early layers are often degenerate. This challenge is further compounded when speech SSL models are applied on out-of-distribution (OOD) data. This paper first shows that SSL models do overthinking in ASR. We then motivate further research in EE by computing an optimal bound for performance versus speed trade-offs. To approach this bound we propose two new strategies for ASR: (1) we adapt the recently proposed patience strategy to ASR; and (2) we design a new EE strategy specific to ASR that performs better than all strategies previously introduced.

*Index Terms*— Self-Supervised Learning, Overthinking

## 1. INTRODUCTION

Even though SSL models significantly improved state of the art performances on most speech tasks [1–5], major drawbacks weaken their impact. First their huge stack of transformer layers [6] make **inference time quite slow**, which is non desirable for on-device problems for instance. Second as their training is task agnostic [7], each layer role and performance for ASR task remains unclear [8]. In particular such large models are prone to **overthinking** [9] which occurs when an accurate prediction is reached at an intermediate layer  $i$  and computations of layers  $j > i$  are wasteful and potentially destructive in terms of prediction quality.

Orthogonal to static model compression strategies such as knowledge distillation [10, 11] or weight pruning [12], early exit methods [13] address the overthinking problem by **dynamically reducing computations** for certain samples at **inference time**. Typical approaches add lightweight exit heads on top of some layers. Those heads, called branches, compute exit scores usually based on confidence metrics such as entropy of softmax prediction distribution. During inference if an intermediate layer’s exit head outputs an exit score bigger than a fixed threshold, the model outputs this layer’s prediction and stops forward computation, saving precious time.

Despite their huge success in vision [14] and language [15–20] and their portfolio of potential applications, EE have not yet met a huge enthusiasm in speech [21–23] and in particular in ASR [24]. Early exiting is very challenging for ASR. First ASR is a complex sequence-to-sequence task in which models fail to get 100% accuracy for most samples. On the contrary most EE approaches were

designed for classification problems [13, 15] in which early predictions are very likely to be correct and constant across layers. This enabled [15], inspired by early stopping strategies [25], to introduce a patience criteria which outperformed confidence based EE techniques on BERT [26]. Secondly, speech SSL models are commonly fine-tuned or used out of the box in **OOD setups** in which audio differs in language, noise level, type of speech and accent from the pre-training (and fine-tuning). This frequent scenario in ASR research remains very challenging [27] and it is legit to wonder whether EE strategies would be effective in such setup.

The contributions of this work are summarized as follows:

- We show that ASR models do overthinking and analyze such phenomenon for in-domain and OOD scenario;
- We compute theoretical lower bounds of speed/quality trade-offs for EE strategies through dynamic programming;
- We adapt patience based EE strategies to suit ASR task;
- We introduce *overlang*: a vocabulary based new EE strategy designed from our findings on overthinking in ASR.

## 2. MOTIVATIONS : OVERTHINKING IN ASR

For an audio sample  $x$  and a model composed of  $N$  layers, we will note  $\hat{y}_i(x)$  the output prediction of layer  $i \in \{i_{\min}, \dots, N\}$ . We assume that the model has exit heads on top of each layer except the  $i_{\min}$  first layers which predictions are often degenerated.  $\hat{y}_i(x)$  is a character sequence that can be compared to the reference transcription  $y(x)$ <sup>1</sup>. Following [9] a model overthinks for an input  $x$  if

$$\exists i < N \text{ such that } \text{error}(\hat{y}_i(x)) \leq \text{error}(\hat{y}_N(x)) \quad (1)$$

If Equation 1 is an equality then only time is wasted but performance is not impacted (we will call this scenario overthinking without degradation), but it could also be that quality of output is degraded. EE strategies aim to find trade-offs between saving computations, so outputting  $\hat{y}_i$  with a small  $i$ , and having  $\hat{y}_i$  quality better than  $\hat{y}_N$  (overthinking scenario) or at least not too much degraded.

### 2.1. Experimental details

**Model** - We used Hubert Large model [1] pre-trained (PT) on Librilight [28] 60k hours and fine-tuned (FT) on Librispeech [29] 960 hours. We used English characters as token set. We added exit heads over all layers from layer 10. Our exit heads are made of a feed-forward module and a CTC [30] head<sup>2</sup>. Heads training (HT) was done on Librispeech 100h clean split, using the intermediate CTC loss [31]. More details and hyperparameters are released in ESPnet [32].

**Data** - We used three datasets for inference: dev-clean and dev-other sets from Librispeech [29], and the English dev set from Commonvoice 5.1 [33]<sup>3</sup>. Table 1 exhibits the distribution shifts for those sets,

<sup>1</sup>For not overloading formulas we may write  $\hat{y}_i$  and  $y$ .

<sup>2</sup>Thus significantly smaller than [24] thus increasing computation savings.

<sup>3</sup>We call the sets respectively *Libri\_dev\_clean*, *Libri\_dev\_other* and *CV\_en*.

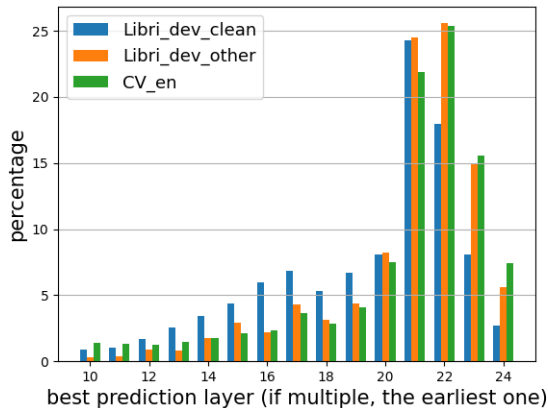
as accents, noise level and distributions (audiobooks/human generated sentences) differ. We see that *Libri\_dev\_clean* is fully in-distribution while *Libri\_dev\_other* is partly OOD and *CV\_en* is totally OOD.

**Table 1:** Distribution shifts for the inference sets.

	<i>Libri_dev_clean</i>	<i>Libri_dev_other</i>	<i>CV_en</i>
PT & FT	in distribution	in distribution	OOD
HT	in distribution	OOD	OOD

## 2.2. Do speech SSL model overthink?

Large networks such as SSL models are very prone to overthinking as shown in [9, 14, 15]. However, those studies mostly targeted classifications tasks such as sentiment analysis [15]. For such task, it seems intuitive that for some  $i$ , layers  $j \geq i$  will output the same class [15]. On the other hand ASR is a complex sequence to sequence task in which the model’s layers are refining a predicted sentence. In most cases  $\hat{y}_i(x)$  is different from  $y(x)$  for all  $i \in \llbracket 1, N \rrbracket$ . Thus we can think that overthinking does not occur since refining the prediction until the last layer could be more beneficial than harmful due to the complexity of the task.



**Fig. 1:** Overthinking in ASR for in-domain and OOD sets.

In Figure 1 we show that SSL models do overthinking for ASR: bar  $i$  is the percentage of samples per set for which the best prediction (for word error rate, WER) is reached **first** at layer  $i$ . We remark that :

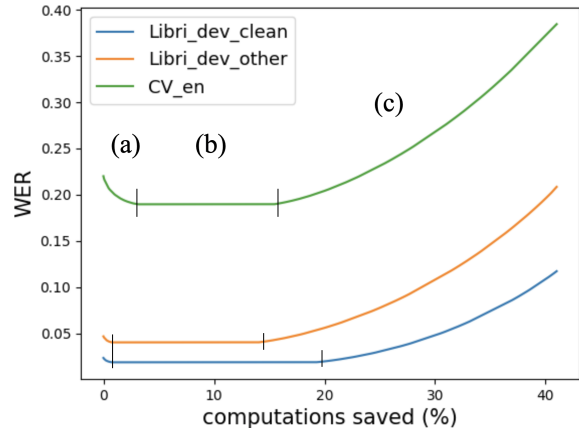
- For all sets less than 10% of the samples need to go to the last layer to reach best prediction, which is a clear case of overthinking.
- For *Libri\_dev\_clean* a non negligible portion of the samples reach their best prediction very early (i.e. before layer 20) whereas for the two OOD sets (*Libri\_dev\_other* and *CV\_en*) more than 70% of the samples need to wait the last 4 layers to reach best prediction.

Additionally to Figure 1, our aim is to differentiate if the model is overthinking without degradation (so only wasting time) or is wasting time **and** degrading performances. In other terms, among samples that reach best prediction before the last layer what is the share for which the last layer remains the best one<sup>4</sup>? We found that it is more than 90% for the Librispeech sets (so mainly overthinking without degradation) whereas for *CV\_en* more than 22% of the samples predictions are degraded when reaching the last layer.

<sup>4</sup>This cannot be inferred from Figure 1, we conducted additional analysis.

## 2.3. Lower bound for early exit strategies

Using every layer predictions for each sample we computed the best performances for any reachable number of saved computations through dynamic programming. Optimum are **convex** (see Figure 2) for the 3 datasets showing that (up to some point) computation savings increase relatively faster than quality degradation which is desirable.



**Fig. 2:** Theoretical bounds to early exits trade-offs

We remark that all curves are composed of (a) a small decreasing portion followed by (b) a non negligible plateau and finally (c) an increasing part. Portion (a) corresponds to a part where we can save computing time **and** increase quality of outputs. This is the evidence of quality degradation due to overthinking. As mentioned in Section 2.3 *CV\_en* suffers more from this phenomenon and so phase (a) is longer for this set in Figure 2. The plateau portion (b) corresponds to overthinking without degradation : wasteful but not harmful computations. Figure 2 optimum reflects Figure 1 findings : the in-distribution set (*Libri\_dev\_clean*) is more prone to overthinking so this plateau is longer for this set. For this set, an optimal strategy can save up to 20% of computations with no degradation (on average). Finally (c) corresponds to the trade-off phases where no more computations can be saved with no degradation cost. In that part we see that the slope is slightly lower for the in-domain set so that computations will be cheaper to save for *Libri\_dev\_clean* than for the OOD sets.

**Note on computation saved** - In network compression studies speed is often calculated through wall clock time or floating point operations saved. However as computations for transformer blocks are quadratic in the input length while the total number of errors is more likely to be linear in the length of the sentence we argue that reporting those time measurements is not desirable for EE. Indeed a strategy exiting long utterances very early and small ones late would be better than the reverse strategy because it would suffer a linear degradation of performances but would get a quadratic gain of time. As we do not want the strategy to favor short utterances, we rather compute relative gain of time for each sample that is simply reported as the proportion of layers that were skipped<sup>5</sup>.

## 3. PROPOSED METHODS

We formally note  $\zeta$  the criterion (taking only boolean values) such that for a sample  $x$  the model exits at layer  $i^* = \min_i \{i | \zeta(\hat{y}_i(x)) = 1\}$ .

<sup>5</sup>As this could have the opposite effect, i.e. favor very early exit for short sentences, we will report results for sentences of more than 10 words only.

We note  $T(x)$  (simplified as  $T$ ) the number of frames of sample  $x$  and  $C$  the number of tokens.  $f_{i,t,c}(x)$  is the output of the softmax of exit branch at layer  $i$ , for frame  $t \in \llbracket 1, T \rrbracket$  and token  $c \in \llbracket 1, C \rrbracket$ <sup>6</sup>.

### 3.1. Confidence based methods

We first tried the commonly used entropy based confidence score [13], computed from the softmax output following Eq. 2.

$$\text{score}(\hat{y}_i(x)) = -\frac{1}{T \cdot C} \sum_{t=1}^T \sum_{c=1}^C f_{i,t,c}(x) \cdot \log(f_{i,t,c}(x)) \quad (2)$$

Following [24] we also tried a variant where the confidence score is given by the average over frames of the maximum softmax(ed) probability instead of the entropy of the distribution. As a well calibrated model would be confident about accurate prediction, our criteria is is given by Equation 3,

$$\zeta(\hat{y}_i(x)) = \mathbb{1}_{\text{score}(\hat{y}_i(x)) < \tau} \quad (3)$$

where  $\tau$  is a fixed threshold and  $\mathbb{1}$  is the indicator function<sup>7</sup>.

### 3.2. Patience based methods

Strategies introduced in Section 3.1 can however be weak when a model is overconfident or poorly calibrated and cannot handle regression cases [15]. Inspired by early stopping, PABEE [15] solved those issues by introducing patience based strategies. The principle is that criteria  $\zeta$  allows exit at layer  $i$  if the prediction at layer  $i$  is consistent with the predictions of the few layers before  $i$ . It enables to (1) avoid unstable predictions; (2) exit when the model somehow converged to a prediction; and (3) indirectly predict from an ensemble of layers. Formally with a distance  $d$ , a threshold  $\tau$  and a patience threshold  $\rho \in \mathbb{N}$ , the patience strategy can be written as in Equation 4.

$$\zeta(\hat{y}_i(x)) = 1 \Leftrightarrow \forall j \in \llbracket i - \rho, i \rrbracket, d(\hat{y}_j(x), \hat{y}_{j-1}(x)) < \tau \quad (4)$$

Distances proposed in [15] (e.g. L2 distance of two predictions) cannot be applied to sequential outputs. We used two distances for patience strategies : **cross-entropy** at the logits level and **Levenshtein distance** at the output sentence level.

### 3.3. Overlang : a vocabulary based strategy

Finally, we propose a novel criterion that directly targets the EE aim: **output high quality predictions** while **avoiding overthinking**. As we stick to Hubert configuration [1] we adopt characters as our token set. Predicted words (sequences of letters delimited by spaces) may thus be any combination of characters which are potentially out-of-vocabulary (e.g. due to misspelling)<sup>8</sup>. A first observation is that when an early layer prediction is of poor quality, most of its **incorrect words are out-of-vocabulary**. Table 2 provides examples of predicted sentences from early layers with several of such out-of-vocabulary errors. These predictions are phonetically proximate to the reference, but consist of many invalid words – we can reasonably detect these by checking against a known vocabulary.

Overthinking happens only when accurate predicted words (which are in-vocabulary) are modified by later layers and so cannot

<sup>6</sup>Our models are CTC-based only however the strategies can be applied to encoder-decoder or any other architecture.

<sup>7</sup>For the maximum probability variant, the inequality sign is changed.

<sup>8</sup>We note that this findings and method also apply for any type of subword.

**Table 2:** Example ASR errors flagged by an out-of-vocabulary check.

Reference Sentence	Predicted sentence
now active exploitation was required	now <b>actiev</b> <b>explotation</b> was <b>requie</b>
he asked on seeing the prisoners	he <b>ased</b> on seeing the <b>prisiners</b>

occur on words that are out-of-vocabulary (under assumption that out-of-vocabulary words are incorrect)<sup>9</sup>. As our model was never guided by any language model in its training, we assume that when early layers predict a word that is in vocabulary this word is very likely to be the accurate prediction. Overthinking may then happen if this correct word is changed to another in or out-of-vocabulary word.

This leads us to the following criteria : exit is allowed at layer  $i$  if the prediction of this layer does not contain too many out of vocabulary words. We note  $\mathcal{V}$  our English vocabulary<sup>10</sup>. We define  $W(\hat{y}_i(x))$  in Equation 5 as the proportion of in-vocabulary words in layer  $i$ 's prediction, where  $\text{card}(\cdot)$  is the cardinal of a set and  $\text{length}(\cdot)$  is the number of words of a sentence.

$$W(\hat{y}_i(x)) = \frac{\text{card}(\{w \in \hat{y}_i(x) \text{ such that } w \in \mathcal{V}\})}{\text{length}(\hat{y}_i(x))} \quad (5)$$

For a threshold  $\tau$ , an exit criterion can be derived as in Equation 6:

$$\zeta(\hat{y}_i(x)) = \mathbb{1}_{W(\hat{y}_i(x)) \geq \tau} \quad (6)$$

This criterion is orthogonal to the ones defined in Section 3.1 and 3.2 so that they also can be combined together. We combined  $\zeta$  defined in Equation 6 with a patience criterion on the  $W(\hat{y}_j(x))$  values. For a chosen integer  $\rho$ , if  $W(\hat{y}_j(x))$  is constant for  $j \in \llbracket i - \rho, i \rrbracket$  then the model exits at layer  $i$  even if  $W(\hat{y}_j(x)) < \tau$ . For a given  $\rho$  and  $\tau$ , Equation 7 defines *overlang*, our vocabulary based EE criteria,

$$\zeta(\hat{y}_i) = \max(\mathbb{1}_{W(\hat{y}_i) \geq \tau}, \mathbb{1}_{W(\hat{y}_i) = W(\hat{y}_{i-1}) = \dots = W(\hat{y}_{i-\rho})}) \quad (7)$$

Where max is equal to one if at least one of the indicators is satisfied.

## 4. RESULTS AND ANALYSIS

### 4.1. Additional details on experiments

Table 3 gathers the values we used for thresholds in our EE strategies.

**Table 3:** Experimental details on EE strategies

Strategy	Measure	range of hyperparameters
Confidence	Entropy	$\tau \in \{0.002, 0.0025, 0.003, \dots, 0.006\}$
Confidence	Maximum Probability	$\tau \in \{0.93, 0.935, 0.94, \dots, 0.97\}$
Patience	Cross-Entropy	$\tau \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ $\rho \in \{1, 2, 3, 4, 5\}$
Patience	Levenshtein distance	$\tau \in \{0.05, 0.07, 0.1, 0.15, 0.2\}$ $\rho \in \{1, 2, 3, 4, 5\}$
Overlang	-	$\tau \in \{0.6, 0.625, 0.65, \dots, 0.95\}$ $\rho = 2$

### 4.2. Efficiency of the strategies

Figure 3 presents the speed/performance trade-offs of the EE strategies introduced in Section 3 for *Libri\_dev\_other* set<sup>11</sup>. We add in green the optimal bound curve computed in Section 2.3 and in red the

<sup>9</sup>This assumption is acceptable given that proportion of out-of-vocabulary words in reference transcription (e.g. named entities) is quite small.

<sup>10</sup>Given by any English dictionary, we used english-words library.

<sup>11</sup>Similar behaviors are observed on the two other inference sets.

curve linking the trade-offs reached by the naive strategies: *always exit at a fixed layer i* (i.e.  $\zeta_i(\hat{y}_j(x)) = \mathbb{1}_{j=i}$ ). Those  $\zeta_i$  are naive strategies as they exit at the same layer for every input sample. Thus we use them as an informal upper bound: any EE strategy should reach better trade-offs than those naive methods. We color in green an 'acceptable trade-offs zone' which is in between the optimal bound and the naive trade-offs. Under 10% of computations saved, the gap between optimal bound, proposed methods and constant layer strategies is very small and this range of time savings is of limited interest for potential applications. On the other hand, performances with large time savings are too much degraded to be applicable in the real world. We thus plot the 10% to 40% computations saved zone.

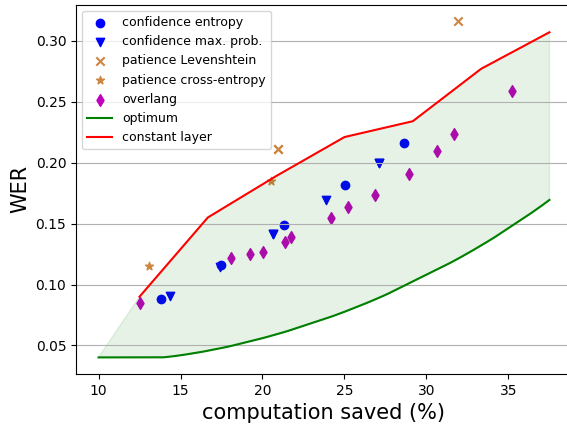


Fig. 3: Trade-off comparison for different EE strategies

We first remark that the patience based strategies (both with cross-entropy and Levenshtein distance) give bad trade-offs compared to confidence, *overlang*, and even naive strategies. Confidence and *overlang* strategies give better results than the naive strategies and *overlang* outperforms the trade-offs obtained by confidence-based strategies quite consistently over the considered zone. Table 4 presents a comparison of Word Error Rate (WER) obtained by confidence and *overlang* for a same (or very similar) amount of computations saved (vertically aligned points on Figure 3). We remark that *overlang* indeed reaches better trade-offs than confidence based approaches for this range of savings with WER reduction<sup>12</sup> of about 10%.

Table 4: Trade-offs comparisons for confidence-based strategies and *overlang*. We also compute relative reductions of the WER.

Method	COMPUTATION SAVED				
	~ 21.3%	~ 24.0%	~ 25.1%	~ 27.0%	~ 28.7%
Confidence (WER)	14.9	16.9	18.2	20.0	21.6
<i>Overlang</i> (WER)	13.5	15.5	16.4	17.4	19.0
<b>Relative Reduction</b>	<b>9.4%</b>	<b>8.3%</b>	<b>9.9%</b>	<b>13.0%</b>	<b>12.0%</b>

#### 4.3. Overthinking for Patience, Confidence and Overlang

We confirm intuitions from Figure 3 by computing overthinking percentage for three strategies saving approximately 21% computa-

<sup>12</sup>WER reduction is given by :  $100 \cdot \frac{\text{WER}(\text{confidence}) - \text{WER}(\text{overlang})}{\text{WER}(\text{confidence})}$

tions<sup>13</sup>. For a sample  $x$  we note  $l(x)$  the exit layer chosen by the strategy. We report overthinking by computing the proportion of  $x$  for which  $\exists i < l(x)$  such that  $\text{WER}(\hat{y}_i(x)) \leq \text{WER}(\hat{y}_{l(x)}(x))$ . Overthinking happens with the patience strategy at 68% whereas confidence based strategy only gets 53% and *overlang* 49%.

Utterance 1630-96099-0016 from *Libri\_dev\_other* provides an example of how the strategies can behave (see Table 5).

Table 5: Early exit layers for utterance 1630-96099-0016.

Strategy	Exit Layer	Sentence
Reference	-	he left everything behind
Patience (cross-entropy)	24	he left everything behind
Confidence (entropy)	13	he left <b>evrthing</b> behind
Overlang	15	he left everything behind

First we see that the language strategy benefits over the confidence one as *evrthing* is not in-vocabulary: *overlang* exits two layers later only but with a correct sentence. On the contrary, the patience strategy using cross-entropy exits at the last layer as the output probability distribution is not stable (in terms of entropy), thus overthinking a lot.

#### 4.4. Additional remarks on experiments

We add empirical findings of practical interest for the readers:

- Combining patience and confidence strategies does not improve over confidence ones contrary to combinations of patience and *overlang*. A plausible explanation is that when prediction is constant across layers, *overlang* ratio  $W$  of in-vocabulary words is constant and so if  $W < \tau$  the model do not exit early. On the other hand, confidence based methods are not impacted by this scenario because even if predicting the same sentence across different layers, the model usually becomes more confident about its prediction and thus quickly reaches the confidence threshold and exits.
- No behavioral difference was found between in-distribution and OOD. It demonstrates the feasibility of EE strategies for OOD.

## 5. CONCLUSION AND FUTURE DIRECTIONS

We demonstrate that overthinking happens in ASR for both in-domain and OOD scenario. Our introduced *overlang* strategy as well as confidence-based methods enable reaching good speed/quality trade-offs thus avoiding overthinking for some samples. We believe that research in this topic is very promising for a wide range of applications such as on-device systems or semi-supervised ASR where pseudo-labels could be generated (more rapidly) by early layers for some samples as [34] showed that bootstrapping from poor quality pseudo-labels is viable.

## 6. ACKNOWLEDGMENTS

We would like to thank Navdeep Jaitly, Tatiana Likhomanenko and Ronan Collobert for helpful discussions on the topic. This work used the Extreme Science and Engineering Discovery Environment (XSEDE) [35], which is supported by National Science Foundation grant number ACI-1548562. Specifically, it used the Bridges system [36], which is supported by NSF award number ACI-1445606, at the Pittsburgh Supercomputing Center (PSC).

<sup>13</sup>cross-entropy patience with  $\tau = 0.05$  and  $\rho = 3$ , entropy confidence with  $\tau = 0.055$  and *overlang* with  $\tau = 0.8$ .

## 7. REFERENCES

- [1] Wei-Ning Hsu et al., “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 2021.
- [2] Alexei and others Baevski, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” in *NeurIPS 2020*.
- [3] Sanyuan Chen et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” *ArXiv*, vol. abs/2110.13900, 2022.
- [4] Shaoshi Ling et al., “Deep contextualized acoustic representations for semi-supervised speech recognition,” in *ICASSP 2020*.
- [5] Shu-Wen Yang et al., “SUPERB: speech processing universal performance benchmark,” *CoRR*, 2021.
- [6] Ashish Vaswani et al., “Attention is all you need,” in *NeurIPS 2017*.
- [7] Lucio M. Dery et al., “Should we be pre-training? an argument for end-task aware training as an alternative,” in *ICLR 2022*.
- [8] Ankita Pasad et al., “Layer-wise analysis of a self-supervised speech representation model,” *ASRU 2021*.
- [9] Yigitcan Kaya et al., “Shallow-deep networks: Understanding and mitigating network overthinking,” in *ICML 2019*.
- [10] Jimmy Ba et al., “Do deep nets really need to be deep?,” in *NIPS 2014*.
- [11] Geoffrey Hinton et al., “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [12] Steven A Janowsky, “Pruning versus clipping in neural networks,” *Physical Review A*, 1989.
- [13] Surat Teerapittayanon et al., “Branchynet: Fast inference via early exiting from deep neural networks,” *ICPR 2016*.
- [14] Zhengcong Fei et al., “Deecap: Dynamic early exiting for efficient image captioning,” in *CVPR 2022*.
- [15] Wangchunshu Zhou et al., “Bert loses patience: Fast and robust inference with early exit,” *NIPS’20*.
- [16] Ji Xin et al., “DeeBERT: Dynamic early exiting for accelerating BERT inference,” in *ACL 2020*.
- [17] Ji Xin et al., “BERxiT: Early exiting for BERT with better fine-tuning and extension to regression,” in *EACL 2021*.
- [18] Kaiyuan Liao et al., “A global past-future early exit method for accelerating inference of pre-trained language models,” in *NAACL 2021*.
- [19] Keli Xie et al., “Elbert: Fast albert with confidence-window based early exit,” in *ICASSP 2021*.
- [20] Tianxiang Sun et al., “Early exiting with ensemble internal classifiers,” *ArXiv*, 2021.
- [21] Sanyuan Chen and others, “Don’t shoot butterfly with rifles: Multi-channel continuous speech separation with early exit transformer,” *ICASSP 2021*.
- [22] Andong Li et al., “Learning to inference with early exit in the progressive speech enhancement,” *EUSIPCO 2021*.
- [23] Raphael Tang et al., “Temporal early exiting for streaming speech commands recognition,” in *ICASSP 2022*.
- [24] Ji Won Yoon et al., “Hubert-ee: Early exiting hubert for efficient speech recognition,” 2022.
- [25] Lutz Prechelt, “Early stopping-but when?,” in *Neural Networks: Tricks of the trade*. Springer, 1998.
- [26] Jacob Devlin et al., “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *NAACL 2019*.
- [27] Wei-Ning Hsu et al., “Robust wav2vec 2.0: Analyzing domain shift in self-supervised pre-training,” in *Interspeech 2021*.
- [28] J. Kahn et al., “Libri-light: A benchmark for asr with limited or no supervision,” in *ICASSP 2020*.
- [29] Panayotov et al., “Librispeech: an asr corpus based on public domain audio books,” in *ICASSP 2015*.
- [30] Alex Graves et al., “Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks,” in *ICML 2006*.
- [31] Jaesong Lee et al., “Intermediate loss regularization for ctc-based speech recognition,” *ICASSP 2021*.
- [32] Shinji Watanabe et al., “ESPnet: End-to-end speech processing toolkit,” in *Interspeech 2018*.
- [33] Rosana Ardila et al., “Common voice: A massively-multilingual speech corpus,” in *Language Resources and Evaluation Conference*, 2020.
- [34] Dan Berrebbi et al., “Continuous pseudo-labeling from the start,” *arXiv preprint arXiv:2210.08711*, 2022.
- [35] J. Towns and other, “Xsede: Accelerating scientific discovery,” *Computing in Science & Engineering*, 2014.
- [36] Nicholas A Nystrom et al., “Bridges: a uniquely flexible hpc resource for new communities and data analytics,” in *Proceedings of the 2015 XSEDE Conference*.