# ALTERNATE INTERMEDIATE CONDITIONING WITH SYLLABLE-LEVEL AND CHARACTER-LEVEL TARGETS FOR JAPANESE ASR

*Yusuke Fujita, Tatsuya Komatsu, Yusuke Kida*

LINE Corporation, Tokyo, Japan

## ABSTRACT

End-to-end automatic speech recognition directly maps input speech to characters. However, the mapping can be problematic when several different pronunciations should be mapped into one character or when one pronunciation is shared among many different characters. Japanese ASR suffers the most from such many-to-one and one-to-many mapping problems due to Japanese kanji characters. To alleviate the problems, we introduce explicit interaction between characters and syllables using Self-conditioned connectionist temporal classification (CTC), in which the upper layers are "self-conditioned" on the intermediate predictions from the lower layers. The proposed method utilizes character-level and syllable-level intermediate predictions as conditioning features to deal with mutual dependency between characters and syllables. Experimental results on Corpus of Spontaneous Japanese show that the proposed method outperformed the conventional multi-task and Self-conditioned CTC methods.

*Index Terms*— connectionist temporal classification, multi-task learning, self-condition

## 1. INTRODUCTION

End-to-end automatic speech recognition (ASR) simplifies the conventional ASR pipelines by mapping input speech into a text sequence directly, without using a hand-crafted pronunciation dictionary. Attention-based encoder-decoders [1] and recurrent neural network transducers [2] are popular end-to-end ASR methods in which an autoregressive decoder predicts the next token given previously estimated tokens with an encoded sequence of audio features. These autoregressive ASR methods have shown state-of-the-art performance with strong encoders such as Transformer [3] and Conformer [4].

Non-autoregressive ASR methods also have a lot of attention due to their efficient inference, which can predict all tokens simultaneously. Connectionist temporal classification (CTC) [5] is a fundamental approach to non-autoregressive ASR. Despite the fact that CTC assumes conditional independence between output labels, which was considered to be a drawback, various approaches on top of CTC [6, 7, 8, 9, 10, 11] have shown comparable performance with autoregressive models. Recently, self-supervised learning models such as

wav2vec 2.0 [12] and HuBERT [13] have shown further improvement by pretraining the encoder using a large amount of unlabeled data and finetuning with labeled data with CTC.

Although end-to-end ASR methods achieved sufficient performance in English, the challenge remains in languages such as Japanese, which face a large character vocabulary size with many homophones and multiple pronunciations [14]. The character vocabulary size is larger than that of phonogram languages such as English. Japanese has over three thousand characters, while English has at most about one hundred characters. Japanese ASR also suffers from homophones: many characters share the same pronunciation, e.g., 高, 公, 行 and other hundreds of characters have the same pronunciation "kou". Therefore, an acoustic feature should be mapped to different character labels considering language contexts. In addition, because most Japanese kanji characters have multiple pronunciations, e.g., the character 空 can be pronounced as "sora", "kara", "kuu", "a", "su" or others, multiple different acoustic features should be mapped to one character label. Because of such one-to-many/many-to-one mappings, end-to-end modeling of all these characters becomes problematic due to data scarcity.

For end-to-end modeling under the data scarcity situation, hierarchical multi-task learning with low-level auxiliary targets has been studied [15, 16, 17]. Loss functions for low-level auxiliary targets, i.e., phonemes or syllables, placed at intermediate layers can help regularize model training for the high-level targets, i.e., characters. However, there is no interaction between the low-level and high-level targets, although the low-level predictions could narrow down the candidates of the high-level predictions. In addition, if intermediate high-level predictions were given, the low-level predictions could be enhanced.

In this paper, we propose a multi-task learning-based ASR method with explicit interactions between the low-level and high-level targets. Our study is inspired by Self-conditioned CTC [10], which estimates intermediate CTC predictions at the lower encoder layer and feeds them back to the upper encoder layers: the upper-layer encoder is "self-conditioned" on the lower-layer predictions. The proposed method utilizes both character-level and syllable-level intermediate predictions as conditioning features fed into the next layers. The two-level conditioning layers can be placed alternately to deal

with *mutual* dependency; not only the syllable-to-character predictions but also the character-to-syllable predictions are considered in the proposed method.

Experimental results on Corpus of Spontaneous Japanese (CSJ) [18] show that the proposed alternate intermediate conditioning outperformed conventional multi-task learning-based and Self-conditioned CTC-based methods. The results also show that the proposed method further reduces syllable error rates than the conventional multi-task learning-based method. The accurate syllable information is particularly useful for downstream tasks of Japanese ASR, such as repeating a spoken query by the system with the correct pronunciation.

We also experiment with Mandarin and English, which also face the homophone and multiple pronunciation problems, not as much as Japanese. Experimental results with Mandarin AISHELL-1 [19] and LibriSpeech-100 [20] show that the proposed alternate intermediate conditioning also outperformed the conventional methods.

## 2. RELATED WORK

### 2.1. CTC-based non-autoregressive ASR

Non-autoregressive ASR methods developed on top of CTC can be categorized into iterative refinement decoding [6, 7, 8] and intermediate prediction objectives [9, 10].

Iterative refinement decoding relaxes the conditional independence assumption between output labels by using a non-autoregressive refinement decoder on top of the CTC-based encoder. The decoder is trained with the masked language model objective function of token sequences [7] or alignment paths [8]. The decoder iteratively refines an initial prediction from the encoder by editing low-confident characters with the masked language model. A drawback of the approach is that the iterative refinement processes increase the computational complexity.

In contrast, our proposed method does not use a decoder while adopting intermediate prediction objectives to the CTC-based encoder. Intermediate CTC [9] uses auxiliary CTC losses at the intermediate encoder layers, which helps faster optimization to the CTC task in lower layers. Self-conditioned CTC [10] further utilizes the intermediate CTC predictions to enhance the encoder by feeding the intermediate predictions to the next encoder layer. We extend Self-conditioned CTC approach by utilizing syllable-level auxiliary targets as well as character-level targets.

### 2.2. Multi-task ASR with low-level auxiliary targets

Multi-task learning with low-level tasks has been actively studied [15, 16, 17]. In typical hierarchical multi-task models, low-level auxiliary tasks are placed at intermediate layers, while the main task is always placed only at the highest layer. Our proposed method differs from these methods in that the predictions for the main task are also applied in intermediate layers, and both low-level and high-level intermediate predictions are utilized for conditioning the subsequent layers.

The grapheme and phoneme sequences can be simultaneously produced by a single-sequence one-to-one model [21] that combines the multiple sequences into a single sequence. This approach is attractive because no architectural modification from typical end-to-end ASR models is required. However, it was difficult to outperform a separate grapheme-only model, possibly due to increased data scarcity.

Joint phoneme-grapheme model [22] utilizes both phone-level and grapheme-level intermediate predictions with iterative refinement decoding. This model can handle mutual dependency between phoneme and grapheme sequences, while it requires two autoregressive decoders. Aiming to capture the mutual dependency, we propose a simpler non-autoregressive model without using autoregressive decoders or iterative refinement decoding.

Our study is inspired by HC-CTC [23], which extends Self-conditioned CTC by using multiple training targets. The multi-granular sub-word units are placed hierarchically according to the vocabulary sizes. Since their targets are character-based sub-words, it is not directly applicable to ideogram languages with large character vocabulary sizes. In this paper, we propose to use syllable-level targets instead of sub-word-level targets and examine non-hierarchical placements of the intermediate targets.

## 3. SELF-CONDITIONED CTC

In this section, we introduce Self-conditioned CTC [10] to describe necessary notations for the proposed method.

### 3.1. Conformer-CTC

End-to-end ASR models the posterior distribution of a $L$-length character sequence $Y = (y_l \in \mathcal{V} \mid l = 1, \ldots, L)$ given a $T$-length input sequence $X = (\mathbf{x}_t \in \mathbb{R}^D \mid t = 1, \ldots, T)$, where $\mathcal{V}$ is a character vocabulary and $D$ is the dimension of acoustic features. CTC [5] handles the posterior $p(Y|X)$ via frame-level alignment paths between $X$ and $Y$ with a blank symbol. The path is denoted by $A = (a_t \in \mathcal{V}' \mid t = 1, \ldots, T)$, where $\mathcal{V}' = \mathcal{V} \cup \{\text{blank}\}$. The path can be mapped to the character sequence by using the collapsing function $\mathcal{B}$ that removes all repeated characters and blank symbols. CTC trains a neural network that predicts the path. An output sequence of the neural network is denoted by $Z = (\mathbf{z}_t \in (0, 1)^{|\mathcal{V}'|} \mid t = 1, \ldots, T)$, where the element $z_{t,k}$ is interpreted as the posterior probability of the character $k$ at time $t$: $p(a_t = k|X)$. The neural network is trained by minimizing the following CTC loss between an output sequence $Z$ and the character sequence $Y$:

$$\mathcal{L}_{\text{ctc}}(Z, Y) = -\log \sum_{A \in \mathcal{B}^{-1}(Y)} \prod_t z_{t,a_t}, \qquad (1)$$

which is the negative log-likelihood over possible paths with the conditional independence assumption: $(a_t \perp\!\!\!\perp a_{\neq t} \mid X)$.

The neural network used in this paper has $N$ Conformer encoders [4] that accept an input sequence $X^{(n-1)}$ and outputs the encoded sequence:

$$X^{(n)} = \mathsf{Encoder}^{(n)}(X^{(n-1)}) \qquad (1 \leq n \leq N), \quad (2)$$

where $X^{(0)} = X$. The output sequence $Z$ is obtained by applying a linear and the softmax functions to the encoded sequence:

$$Z = \mathsf{Softmax}(\mathsf{Linear}_{D \to |\mathcal{V}'|}(X^{(N)})), \quad (3)$$

where $\mathsf{Linear}_{D \to |\mathcal{V}'|}(\cdot)$ maps a $D$-dimensional vector into a $|\mathcal{V}'|$-dimensional vector for each element in the input sequence.

### 3.2. Intermediate CTC

Intermediate CTC [9] introduces additional CTC predictions from intermediate encoder blocks. An intermediate prediction sequence for the $n$-th encoder block $Z^{(n)} = (\mathbf{z}_t^{(n)} \in (0,1)^{|\mathcal{V}'|}|t = 1, \ldots, T)$ is computed as:

$$Z^{(n)} = \mathsf{Softmax}(\mathsf{Linear}_{D \to |\mathcal{V}'|}(X^{(n)})). \quad (4)$$

Note that the linear layer is shared with the output layer (Eq. 3) so that no additional parameters are added. The losses for the intermediate predictions are computed as well as the original CTC loss (Eq. 1), and the total loss is a weighted sum of the original and the intermediate CTC losses:

$$\mathcal{L}_{\mathsf{ic}} = (1 - \lambda)\mathcal{L}_{\mathsf{ctc}}(Z, Y) + \frac{\lambda}{|\mathcal{N}|}\sum_{n \in \mathcal{N}} \mathcal{L}_{\mathsf{ctc}}(Z^{(n)}, Y), \quad (5)$$

where $\lambda \in (0, 1)$ is a mixing weight and $\mathcal{N}$ is a set of layer indices for intermediate loss computation.

### 3.3. Self-conditioned CTC

Self-conditioned CTC [10] further utilizes the intermediate CTC prediction as a condition for the next encoder input. Eq. 2 is modified as follows:

$$X^{(n)} = \mathsf{Encoder}^{(n)}(X'^{(n-1)}) \qquad (1 \leq n \leq N), \quad (6)$$

$$X'^{(n)} = \begin{cases} X^{(n)} + \mathsf{Linear}_{|\mathcal{V}'| \to D}(Z^{(n)}) & (n \in \mathcal{N}), \\ X^{(n)} & (n \notin \mathcal{N}), \end{cases} \quad (7)$$

where $\mathsf{Linear}_{|\mathcal{V}'| \to D}(\cdot)$ maps a $|\mathcal{V}'|$-dimensional vector into a $D$-dimensional vector for each element in the input sequence. This linear layer is shared among the $|\mathcal{N}|$ intermediate layers. In this way, the upper-layer encoder is "self-conditioned" on the lower-layer predictions. This intermediate conditioning relaxes the conditional independence assumption of CTC-based ASR models since the upper-layer encoder can see intermediate output tokens embedded into the encoder input.
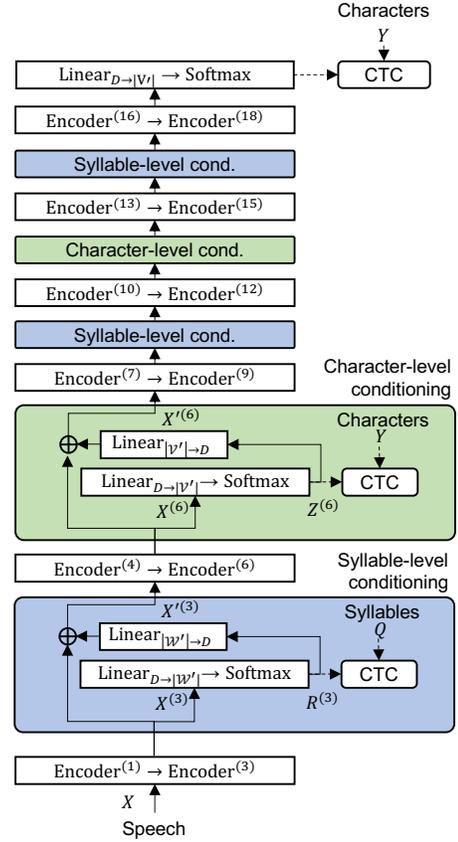


**Fig. 1**. The overview of the proposed alternate intermediate conditioning method. Intermediate conditioning is performed every three encoder blocks and syllable-level conditioning and character-level conditioning are placed alternately.

The proposed method is an extension of Self-conditioned CTC. We utilize the intermediate conditioning technique not only with the same character-based CTC objective but together with the low-level auxiliary targets.

## 4. PROPOSED METHOD: ALTERNATE INTERMEDIATE CONDITIONING

Figure 1 shows an overview of the proposed alternate intermediate conditioning method. A syllable-level auxiliary target sequence $Q = (q_m \in \mathcal{W} \mid m = 1, \ldots, M)$ is prepared for each training sample, where $\mathcal{W}$ is a set of syllables and $M$ is the length of the syllable sequence. Additional CTC predictions $R^{(n)} = (\mathbf{r}_t \in (0,1)^{|\mathcal{W}'|}|t = 0, \ldots, T)$ for the syllable sequence are computed from intermediate layers, similar to Eq. 4:

$$R^{(n)} = \mathsf{Softmax}(\mathsf{Linear}_{D \to |\mathcal{W}'|}(X^{(n)})), \quad (8)$$

where $\mathcal{W}' = \mathcal{W} \cup \{\mathsf{blank}\}$. The total loss is a weighted sum of the original CTC loss, the intermediate CTC loss for char-

acters, and the intermediate CTC loss for syllables:

$$\mathcal{L}_{\text{mic}} = (1 - \lambda)\mathcal{L}_{\text{ctc}}(Z, Y) \tag{9}$$
$$+ \frac{\lambda}{|\mathcal{N}| + |\mathcal{N}_q|} \sum_{n \in \mathcal{N}} \mathcal{L}_{\text{ctc}}(Z^{(n)}, Y)$$
$$+ \frac{\lambda}{|\mathcal{N}| + |\mathcal{N}_q|} \sum_{n \in \mathcal{N}_q} \mathcal{L}_{\text{ctc}}(R^{(n)}, Q),$$

where $\mathcal{N}_q$ is a set of layer indices for the intermediate loss for syllables. As with Self-conditioned CTC, the intermediate predictions are fed back to the next encoder, as follows:

$$X_t'^{(n)} = \begin{cases} X^{(n)} + \mathsf{Linear}_{|\mathcal{V}'| \to D}(Z^{(n)}) & (n \in \mathcal{N} \setminus \mathcal{N}_q), \\ X^{(n)} + \mathsf{Linear}_{|\mathcal{W}'| \to D}(R^{(n)}) & (n \in \mathcal{N}_q \setminus \mathcal{N}), \\ X^{(n)} + \mathsf{Linear}_{|\mathcal{V}'| \to D}(Z^{(n)}) \\ \qquad + \mathsf{Linear}_{|\mathcal{W}'| \to D}(R^{(n)}) & (n \in \mathcal{N} \cap \mathcal{N}_q), \\ X^{(n)} & (n \notin \mathcal{N} \cup \mathcal{N}_q). \end{cases} \tag{10}$$

The sets of intermediate layers $\mathcal{N}$ and $\mathcal{N}_q$ are important hyperparameters in the proposed method. In Figure 1, an example of the sets of layers for the proposed "alternate" intermediate conditioning is shown: $\mathcal{N} = \{6, 12\}$ and $\mathcal{N}_q = \{3, 9, 15\}$, where character-level and syllable-level conditioning layers are placed alternately to handle the mutual dependency between the two tasks.

We experiment with "hierarchical" intermediate conditioning for comparison: syllable-level predictions are placed at lower layers, while character-level predictions are at higher layers. We also investigate "parallel" intermediate conditioning. In this case, syllable-level and character-level conditioning are done at the same layer by adding two embeddings from both $Z^{(n)}$ and $R^{(n)}$ to the next encoder input.

## 5. EXPERIMENT

To verify the effectiveness of the proposed method, we conducted experiments conforming to the comparative study for non-autoregressive ASR models [11], by using ESPnet [24, 25] with almost the same hyperparameters.

### 5.1. Data

Our main results were obtained using CSJ [18]. The 271-hour subset of academic presentation speech (CSJ-APS) was used for training. Since the corpus includes not only character sequences but also corresponding pronunciation labels, syllable-level sequences were extracted from the pronunciation labels as the auxiliary training targets. The character vocabulary size $|\mathcal{V}|$ was 2753, and the syllable vocabulary size $|\mathcal{W}|$ was 256. We used the official evaluation sets: "eval1", "eval2", and "eval3" for testing. Although CSJ contains another 330-hour subset of simulated public speech on everyday

topics (CSJ-SPS) for training, we did not use CSJ-SPS to conform to the prior study [11]. Note that "eval1" and "eval2" were drawn from CSJ-APS, while "eval3" was drawn from CSJ-SPS, which was considered an out-of-domain test set.

Mandarin AISHELL-1 corpus [19] was also used for evaluating the proposed method. According to the official split of AISHELL-1, the 150-hour subset was used for training. For syllable-level targets, pinyin labels without tones were automatically generated from character sequences by using *pypinyin* [1]. The character vocabulary size $|\mathcal{V}|$ was 4231, and the syllable vocabulary size $|\mathcal{W}|$ was 404.

An additional experiment was conducted on 100-hour subset of LibriSpeech [20]. For the English experiment, 300 subwords tokenized with SentencePiece [26] were used instead of characters as the main output target: $|\mathcal{V}| = 300$. The low-level auxiliary targets are prepared using the CMU pronunciation dictionary [2]. The phoneme sequence within each word was concatenated and then tokenized into 300 subwords: $|\mathcal{W}| = 300$. [3]

For the input samples, 80-dimensional Mel-scale filterbank coefficients with three-dimensional pitch features were extracted using Kaldi toolkit [27]. Speed perturbation [28] and SpecAugment [29] were also applied to the training data.

### 5.2. Model configurations

We prepared four models with conventional methods and three variant models for our proposed method. Note that the number of parameters in these models was 32 million for CSJ, 31 million for LibriSpeech-100, and 52 million for AISHELL-1, respectively. The number of additional parameters introduced by the proposed method was small, at most 0.2 million, since the additional linear layers were shared among intermediate layers.

#### 5.2.1. Conventional models

**Baseline:** The Conformer-CTC model as described in Section 3.1 was used. The number of layers $N$ was 18, and the encoder dimension $D$ was 256. The convolution kernel size and the number of attention heads were 15 and 4, respectively. According to the character vocabulary size, the feed-forward layer dimension in the Conformer blocks was set differently: 1024 for CSJ and LibriSpeech-100, and 2048 for AISHELL-1. The CSJ and LibriSpeech-100 models were trained for 50 epochs, while the AISHELL-1 models were trained for 100 epochs. The final model was obtained by averaging model parameters over 10-best checkpoints in terms of validation loss values. The Adam optimizer [30] with $\beta_1 = 0.9$, $\beta_2 = 0.98$,

---

[1] https://github.com/mozillazg/python-pinyin

[2] http://www.speech.cs.cmu.edu/cgi-bin/cmudict

[3] The English experiment used the phoneme-based subwords instead of syllables because the number of English syllables was too large to form the output layer. Although the phoneme-based subwords were not actually syllables, they were derived from the phoneme sequences as well as syllables.

**Table 1**. Character error rates on CSJ. The results were obtained without language models.

| | Method | Layer indices for CTC loss | | Intermediate | CERs (%) | | |
|---|---|---|---|---|---|---|---|
| | | Character | Syllable | conditioning | eval1 | eval2 | eval3 |
| Conventional | Baseline | 18 | | N | 5.4 | 3.9 | 9.9 |
| | Multitask | 18 | 15 | N | 5.4 | 3.7 | 9.5 |
| | InterCTC | 3,6,9,12,15,18 | | N | 5.4 | 3.8 | 9.4 |
| | SelfCond | 3,6,9,12,15,18 | | Y | 5.3 | 3.7 | 9.2 |
| Proposed | Parallel | 6,12,18 | 6,12,18 | Y | 5.2 | 3.6 | **8.8** |
| | Hierarchical | 12,15,18 | 3,6,9 | Y | **5.1** | 3.6 | **8.8** |
| | Alternate | 6,12,18 | 3,9,15 | Y | **5.1** | **3.5** | **8.8** |
| Prior study | CTC/Att+BeamSearch [11] | | | | **5.1** | 3.8 | 9.0 |
| | Self-conditioned CTC [11] | | | | 5.3 | 3.7 | 9.1 |

**Table 2**. Examples of character-level and syllable-level intermediate predictions with the proposed model with alternate conditioning. The underlined characters indicate substitution errors and ** indicates deletion errors.

| Layer | Syl/Chr | Prediction |
|---|---|---|
| 3 | Syl | ダ イ イ チ ** レ ー ノ シ ケ ン シュ ー テ ン ト ダ イ ニ シ レ ー ノ シ テ ン ガ カ ナ ラ ズ イ ** キ ス ル |
| 9 | Syl | ダ イ イ チ シ レ ー ノ シ テ ン シュ ー テ ン ト ダ イ ニ シ レ ー ノ シ テ ン ガ カ ナ ラ ズ イッ チ ス ル |
| 15 | Syl | ダ イ イ チ シ レ ー ノ シ テ ン シュ ー テ ン ト ダ イ ニ シ レ ー ノ シ テ ン ガ カ ナ ラ ズ イッ チ ス ル |
| 6 | Chr | 第 一 ** 例 の 試 け 終 点 と 第 二 指 令 の 視 点 が 必 ず 一 致 す る |
| 12 | Chr | 第 一 ** 令 の 始 点 終 点 と 第 二 指 令 の 始 点 が 必 ず 一 致 す る |
| 18 | Chr | 第 一 指 令 の 始 点 終 点 と 第 二 指 令 の 始 点 が 必 ず 一 致 す る |

the Noam learning rate scheduling [31] with 25k warmup steps, a learning rate factor of 5.0 were used for training. For CSJ and AISHELL-1, four GPUs are used in parallel with a batch size of 128 per GPU. For LibriSpeech-100, one GPU is used with a batch size of 128. CTC greedy decoding [5] was used at inference, without using any language models.

**Multitask:** Based on the Baseline model, an additional syllable-level CTC prediction was placed at the 15-th layer, a few layers lower than the last layer, which is similar to conventional multi-task learning methods [15].

**InterCTC:** Five character-level intermediate CTC predictions were placed at $\mathcal{N} = \{3, 6, 9, 12, 15\}$ with $\lambda = 0.5$. Other configurations were identical to the Baseline model.

**SelfCond:** In addition to the InterCTC model, conditioning with the intermediate CTC predictions was applied.

*5.2.2. Proposed models*

We experimented with three strategies regarding the placement of intermediate predictions: $\mathcal{N}$ and $\mathcal{N}_q$. As with the SelfCond model, which had five intermediate CTC predictions, we placed the same number of CTC predictions for the proposed model.

**Parallel:** Character-level and syllable-level predictions were placed at the same layers: $\mathcal{N} = \{6, 12\}$ and $\mathcal{N}_q = \{6, 12, 18\}$ with $\lambda = 0.5$.

**Hierarchical:** Character-level predictions were placed at higher layers $\mathcal{N} = \{12, 15\}$, while syllable-level predictions were placed at lower layers $\mathcal{N}_q = \{3, 6, 9\}$ with $\lambda = 0.5$.

**Alternate:** Character-level and syllable-level CTC predictions were placed alternately as shown in Fig. 1: $\mathcal{N} = \{6, 12\}$ and $\mathcal{N}_q = \{3, 9, 15\}$ with $\lambda = 0.5$.

**5.3. Results on CSJ**

Table 1 shows the character error rates for CSJ evaluation sets. The results obtained from the conventional models show that the multi-task learning model with auxiliary syllable-level predictions performed better than the baseline CTC-based model. However, a similar performance was obtained by Intermediate CTC without using the syllable-level predictions. Self-conditioned CTC was the best model among the conventional methods.

The proposed models outperformed the conventional models. The results clearly show that the syllable-level intermediate predictions help improve the accuracy when they are used together with character-level intermediate predictions. We could not see a clear difference between the three conditioning strategies, while the alternate conditioning was slightly better than the others. The results suggest that the alternate conditioning can capture mutual dependency between syllables and characters. Table 2 shows examples of the intermediate predictions that gradually reduce errors layer by layer.

The proposed models showed significant improvement over conventional models on "eval3", which was considered as an out-of-domain test set. The results suggest that the auxiliary syllable-level prediction can improve robustness

**Table 3**. Syllable error rates on CSJ obtained from intermediate predictions.

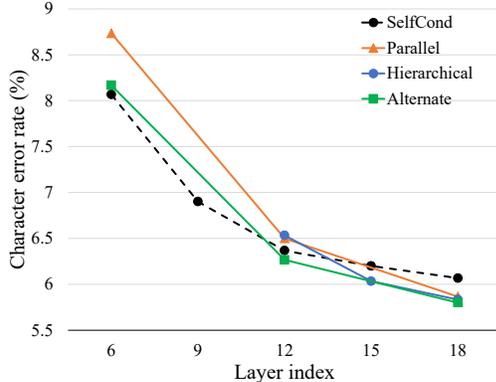| Method | Layer | eval1 | eval2 | eval3 |
|--------|-------|-------|-------|-------|
| Multitask | 15 | 4.0 | 2.6 | 4.8 |
| Alternate | 3 | 8.4 | 5.4 | 8.8 |
| | 9 | 4.2 | 2.5 | 4.8 |
| | 15 | **3.8** | **2.3** | **4.3** |



**Fig. 2**. Character error rates on CSJ obtained from intermediate and final predictions. The error rates on three datasets are averaged.

against domain mismatch, which is particularly important for large vocabulary ideogram languages with a relatively small number of training samples per character.

Table 3 shows syllable error rates on CSJ, which were obtained using the intermediate layer outputs. The proposed method achieved better accuracy on the low-level auxiliary task than the conventional multi-task learning method. Syllable error rates were reduced layer by layer. The results suggest that accurate syllable recognition information from the intermediate layer are fed back to the encoder on the upper layers and this conditioning helps improve the character-level predictions.

Fig. 2 shows character error rates on CSJ obtained from intermediate layers. Similar to syllable error rates shown in Table 3, character error rates were reduced layer by layer. Conventional self-conditioned CTC achieved slightly better results on lower than 12th layer, while the error reduction at the higher layer was smaller than the proposed method. At the 12th intermediate layer, we see that alternate intermediate conditioning was better than other strategies. The results suggest that conditioning with two different-level targets alternately gives better intermediate predictions and leads to better final performance.

**Table 4**. Character error rates on AISHELL-1. The results were obtained without language models.

| Method | dev | test |
|--------|-----|------|
| Baseline | 5.6 | 6.0 |
| Multitask | 4.9 | 5.4 |
| InterCTC | 4.4 | 4.8 |
| SelfCond | 4.2 | 4.6 |
| Parallel | 4.5 | 4.9 |
| Hierarchical | 4.1 | 4.5 |
| Alternate | **4.0** | **4.3** |
| Conformer-Transformer [32] | 4.4 | 4.7 |
| Gated Interlayer Collaboration [33] | **4.0** | 4.4 |

**Table 5**. Word error rates on Librispeech-100. The results were obtained without language models.

| Method | dev-clean | dev-other | test-clean | test-other |
|--------|-----------|-----------|------------|------------|
| SelfCond | 7.1 | 20.8 | 7.6 | 21.1 |
| Alternate | **6.7** | **19.9** | **7.1** | **20.0** |

### 5.4. Results on AISHELL-1 and LibriSpeech-100

Table 4 shows the results with Mandarin AISHELL-1. The proposed method with alternate conditioning strategy was the best among the proposed conditioning strategies. Contrary to the CSJ results, parallel conditioning showed worse performance than the others. The results suggest that the two-level conditioning placed at the same layer can hurt performance, while the alternate intermediate conditioning can give consistent performance improvement for both Japanese and Mandarin. Table 5 shows the results with LibriSpeech-100. Consistent performance improvement was observed even with the phonogram language.

### 5.5. Comparison with prior studies

Comparison results with the prior comparative study [11] for CSJ and recently published papers [32, 33] for AISHELL-1 are shown in Table 1 and 4, respectively. The proposed model with alternate conditioning outperformed the best non-autoregressive ASR methods and improved variants of Self-conditioned CTC. From these results, it can be concluded that the proposed method achieves state-of-the-art performance by utilizing syllable-level targets.

## 6. CONCLUSIONS

We proposed an end-to-end ASR method to model the interaction between the syllable-level and character-level predictions using intermediate conditioning technique with Self-conditioned CTC. Experimental results on CSJ show that the proposed method outperformed the conventional multi-task and Self-conditioned CTC methods.

# 7. REFERENCES

[1] Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio, "Attention-based models for speech recognition," in *Proc. NIPS*, 2015.

[2] Alex Graves, "Sequence transduction with recurrent neural networks," in *International Conference on Machine Learning: Representation Learning Workshop*, 2012.

[3] Qian Zhang, Han Lu, Hasim Sak, Anshuman Tripathi, Erik McDermott, Stephen Koo, and Shankar Kumar, "Transformer transducer: A streamable speech recognition model with transformer encoders and rnn-t loss," in *Proc. ICASSP*, 2020, pp. 7829–7833.

[4] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang, "Conformer: Convolution-augmented Transformer for Speech Recognition," in *Proc. Interspeech*, 2020, pp. 5036–5040.

[5] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, "Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks," in *Proc. ICML*, 2006, p. 369–376.

[6] Nanxin Chen, Shinji Watanabe, Jesús Villalba, Piotr Żelasko, and Najim Dehak, "Non-Autoregressive Transformer for Speech Recognition," *IEEE Signal Processing Letters*, vol. 28, pp. 121–125, 2021.

[7] Yosuke Higuchi, Shinji Watanabe, Nanxin Chen, Tetsuji Ogawa, and Tetsunori Kobayashi, "Mask CTC: Non-Autoregressive End-to-End ASR with CTC and Mask Predict," in *Proc. Interspeech*, 2020, pp. 3655–3659.

[8] Ethan A. Chi, Julian Salazar, and Katrin Kirchhoff, "Align-Refine: Non-Autoregressive Speech Recognition via Iterative Realignment," in *NAACL*, 2021.

[9] Jaesong Lee and Shinji Watanabe, "Intermediate Loss Regularization for CTC-Based Speech Recognition," in *Proc. ICASSP*, 2021, pp. 6224–6228.

[10] Jumon Nozaki and Tatsuya Komatsu, "Relaxing the Conditional Independence Assumption of CTC-Based ASR by Conditioning on Intermediate Predictions," in *Proc. Interspeech*, 2021, pp. 3735–3739.

[11] Yosuke Higuchi, Nanxin Chen, Yuya Fujita, Hirofumi Inaguma, Tatsuya Komatsu, Jaesong Lee, Jumon Nozaki, Tianzi Wang, and Shinji Watanabe, "A comparative study on non-autoregressive modelings for speech-to-text generation," in *Proc. ASRU*, 2021, pp. 47–54.

[12] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. NeurIPS*, 2020, vol. 33, pp. 12449–12460.

[13] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed, "HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.

[14] Hitoshi Ito, Aiko Hagiwara, Manon Ichiki, Takeshi Mishima, Shoei Sato, and Akio Kobayashi, "End-to-end speech recognition for languages with ideographic characters," in *2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2017, pp. 1228–1232.

[15] Shubham Toshniwal, Hao Tang, Liang Lu, and Karen Livescu, "Multitask Learning with Low-Level Auxiliary Tasks for Encoder-Decoder Based Speech Recognition," in *Proc. Interspeech*, 2017, pp. 3532–3536.

[16] Kanishka Rao and Haşim Sak, "Multi-accent speech recognition with hierarchical grapheme based models," in *Proc. ICASSP*, 2017, pp. 4815–4819.

[17] Ramon Sanabria and Florian Metze, "Hierarchical Multitask Learning With CTC," in *Proc. SLT*, 2018, pp. 485–490.

[18] Kikuo Maekawa, "Corpus of spontaneous japanese: Its design and evaluation," in *ISCA & IEEE Workshop on Spontaneous Speech Processing and Recognition*, 2003.

[19] Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng, "AISHELL-1: An open-source Mandarin speech corpus and a speech recognition baseline," in *20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment (O-COCOSDA)*, 2017, pp. 1–5.

[20] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "Librispeech: An asr corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.

[21] Motoi Omachi, Yuya Fujita, Shinji Watanabe, and Matthew Wiesner, "End-to-end ASR to jointly predict transcriptions and linguistic annotations," in *Proc. NAACL-HLT*. June 2021, pp. 1861–1871, Association for Computational Linguistics.

[22] Yotaro Kubo and Michiel Bacchiani, "Joint phoneme-grapheme model for end-to-end speech recognition," in *Proc. ICASSP*, 2020, pp. 6119–6123.

[23] Yosuke Higuchi, Keita Karube, Tetsuji Ogawa, and Tetsunori Kobayashi, "Hierarchical Conditional End-to-End ASR with CTC and Multi-Granular Subword Units," in *Proc. ICASSP*, 2021.

[24] Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson Enrique Yalta Soplin, Jahn Heymann, Matthew Wiesner, Nanxin Chen, Adithya Renduchintala, and Tsubasa Ochiai, "ESPnet: End-to-End Speech Processing Toolkit," in *Proc. Interspeech*, 2018, pp. 2207–2211.

[25] Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, Hirofumi Inaguma, Naoyuki Kamo, Chenda Li, Daniel Garcia-Romero, Jiatong Shi, Jing Shi, Shinji Watanabe, Kun Wei, Wangyou Zhang, and Yuekai Zhang, "Recent developments on espnet toolkit boosted by conformer," in *Proc. ICASSP*, 2021, pp. 5874–5878.

[26] Taku Kudo and John Richardson, "SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Brussels, Belgium, Nov. 2018, pp. 66–71, Association for Computational Linguistics.

[27] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely, "The Kaldi speech recognition toolkit," in *Proc. ASRU*, 2011.

[28] Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur, "Audio augmentation for speech recognition," in *Proc. Interspeech*, 2015, pp. 3586–3589.

[29] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," in *Proc. Interspeech*, 2019, pp. 2613–2617.

[30] Diederik P. Kingma and Jimmy Ba, "Adam: A Method for Stochastic Optimization," in *Proc. ICLR*, 2015.

[31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, p. 6000–6010.

[32] Xuankai Chang, Takashi Maekaku, Pengcheng Guo, Jing Shi, Yen-Ju Lu, Aswin Shanmugam Subramanian, Tianzi Wang, Shu-wen Yang, Yu Tsao, Hung-yi Lee, and Shinji Watanabe, "An Exploration of Self-Supervised Pretrained Representations for End-to-End Speech Recognition," in *Proc. ASRU*, 2021, pp. 228–235.

[33] Yuting Yang, Yuke Li, and Binbin Du, "Improving ctc-based asr models with gated interlayer collaboration," *ArXiv*, vol. abs/2205.12462, 2022.