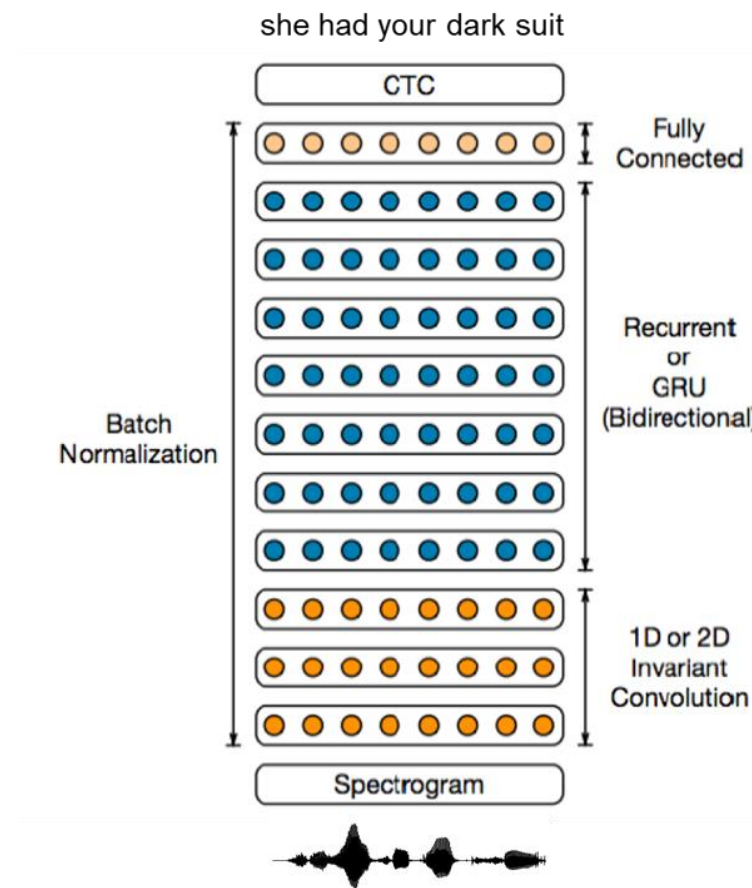


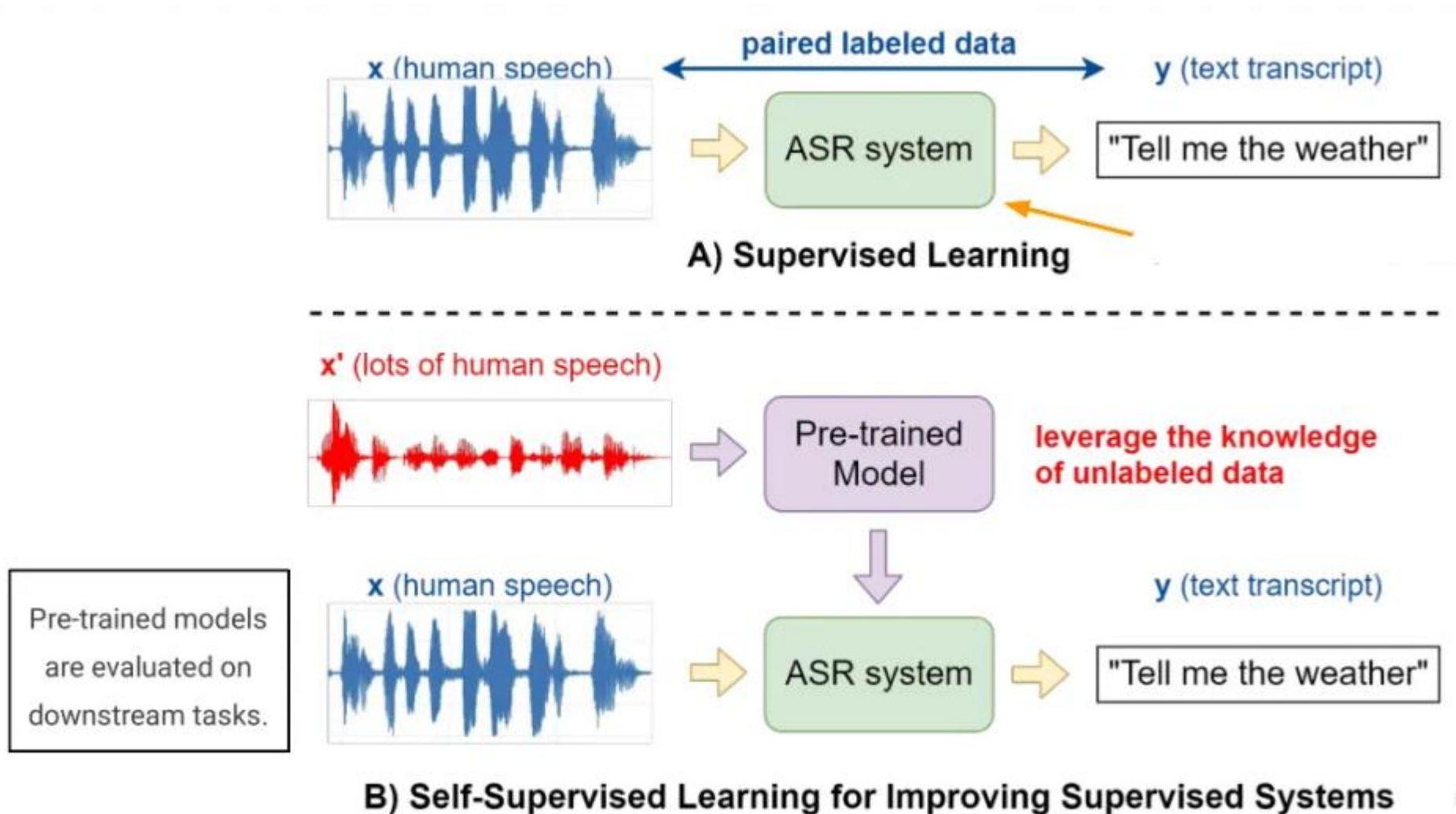
Wav2vec

Automatic Speech Recognition (ASR)



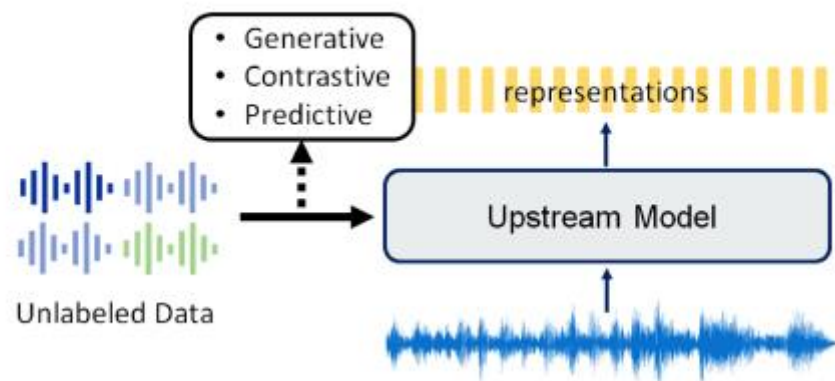
End-to-end ASR

Automatic Speech Recognition (ASR)

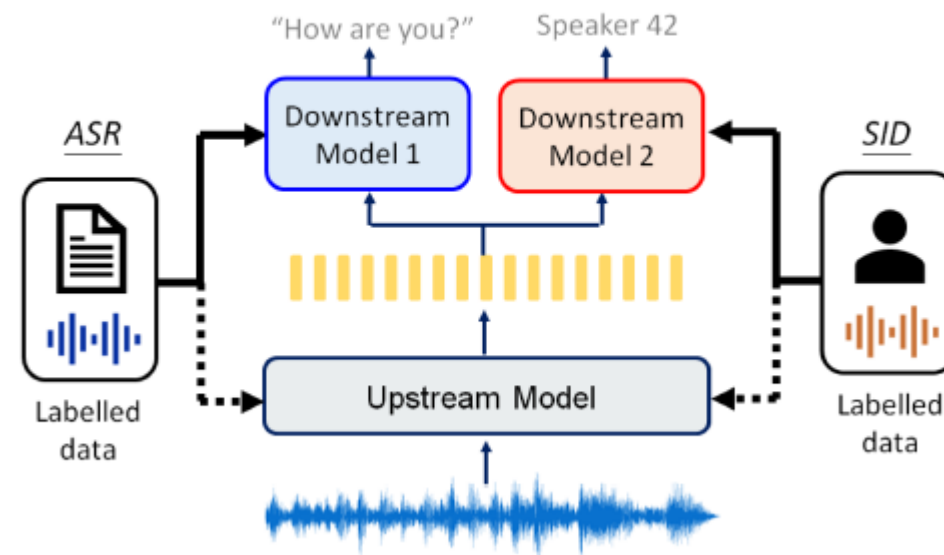


Foundation models in ASR

Phase 1: Pre-train

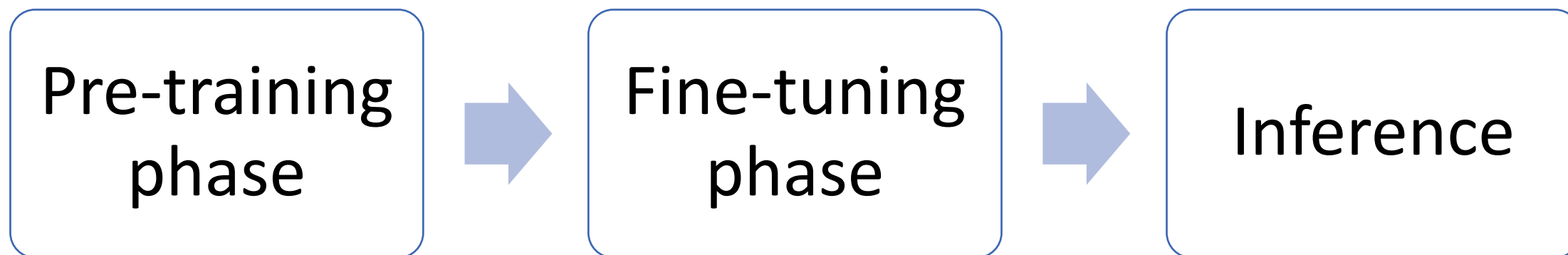


Phase 2: Downstream



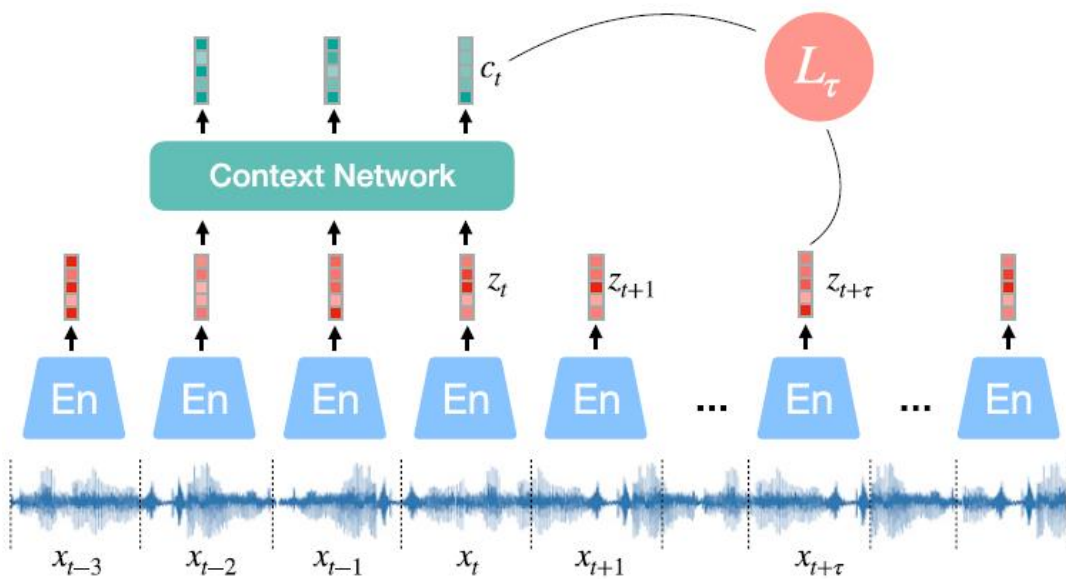
- In the first stage, we use SSL to pre-train a *representation model*, also called an *upstream model* or a *foundation model*.
- In the second stage, downstream tasks use either the learned representation from the frozen model, or fine-tune the entire pre-trained model in a supervised phase. Automatic speech recognition (ASR) and speaker identification (SID) are examples of downstream applications.

Foundation models in ASR



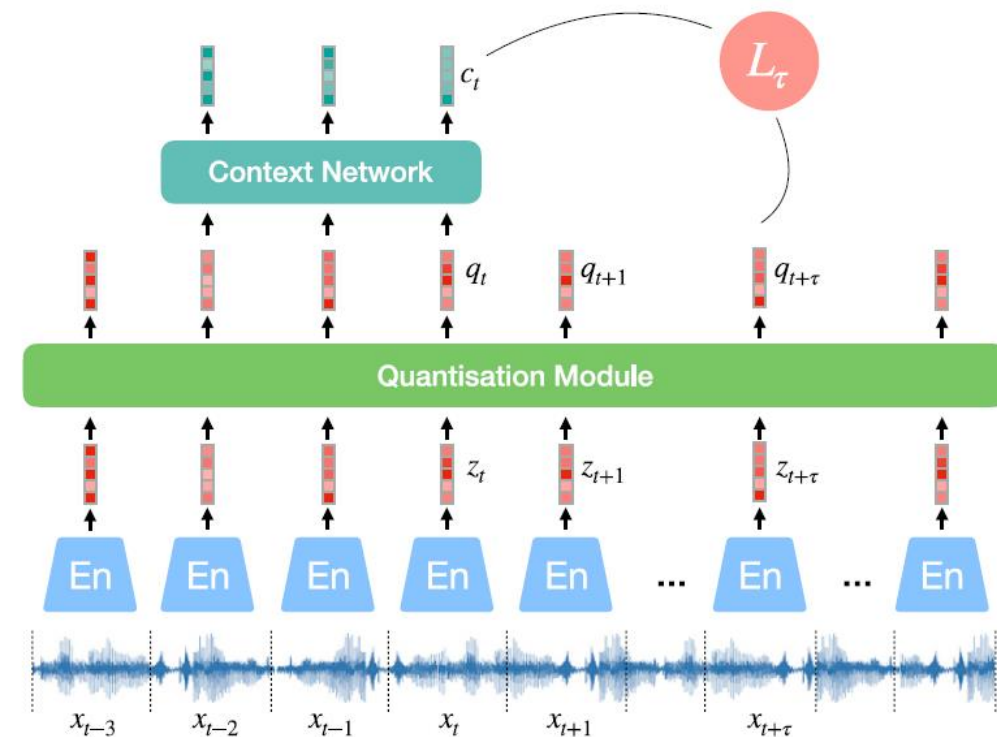
Foundation models in ASR

A



Wav2vec

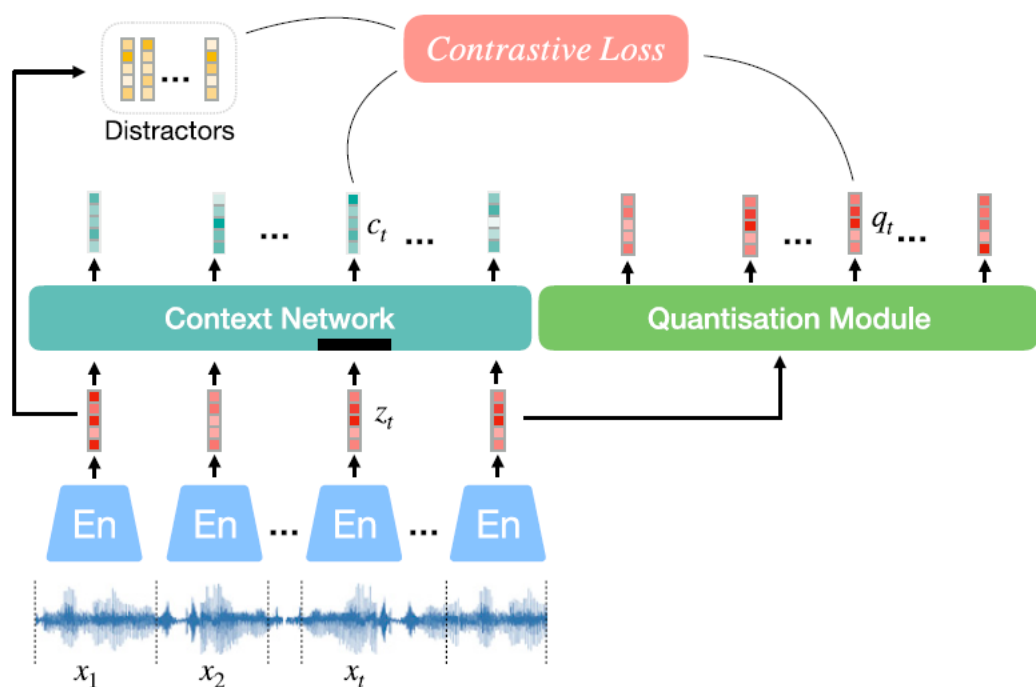
B



VQ-wav2vec

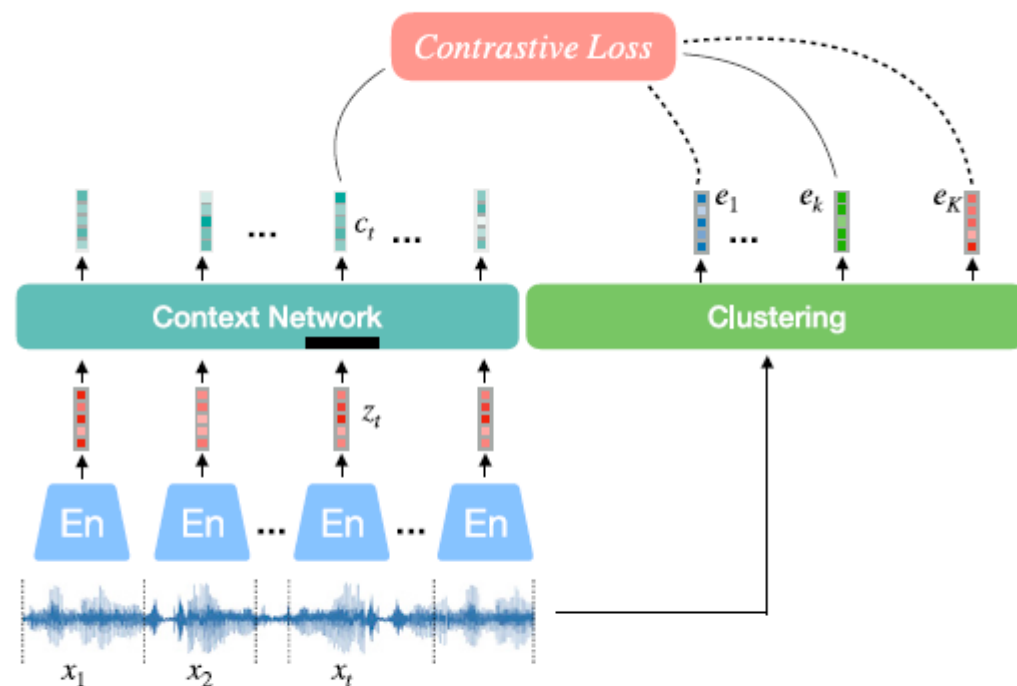
Foundation models in ASR

C

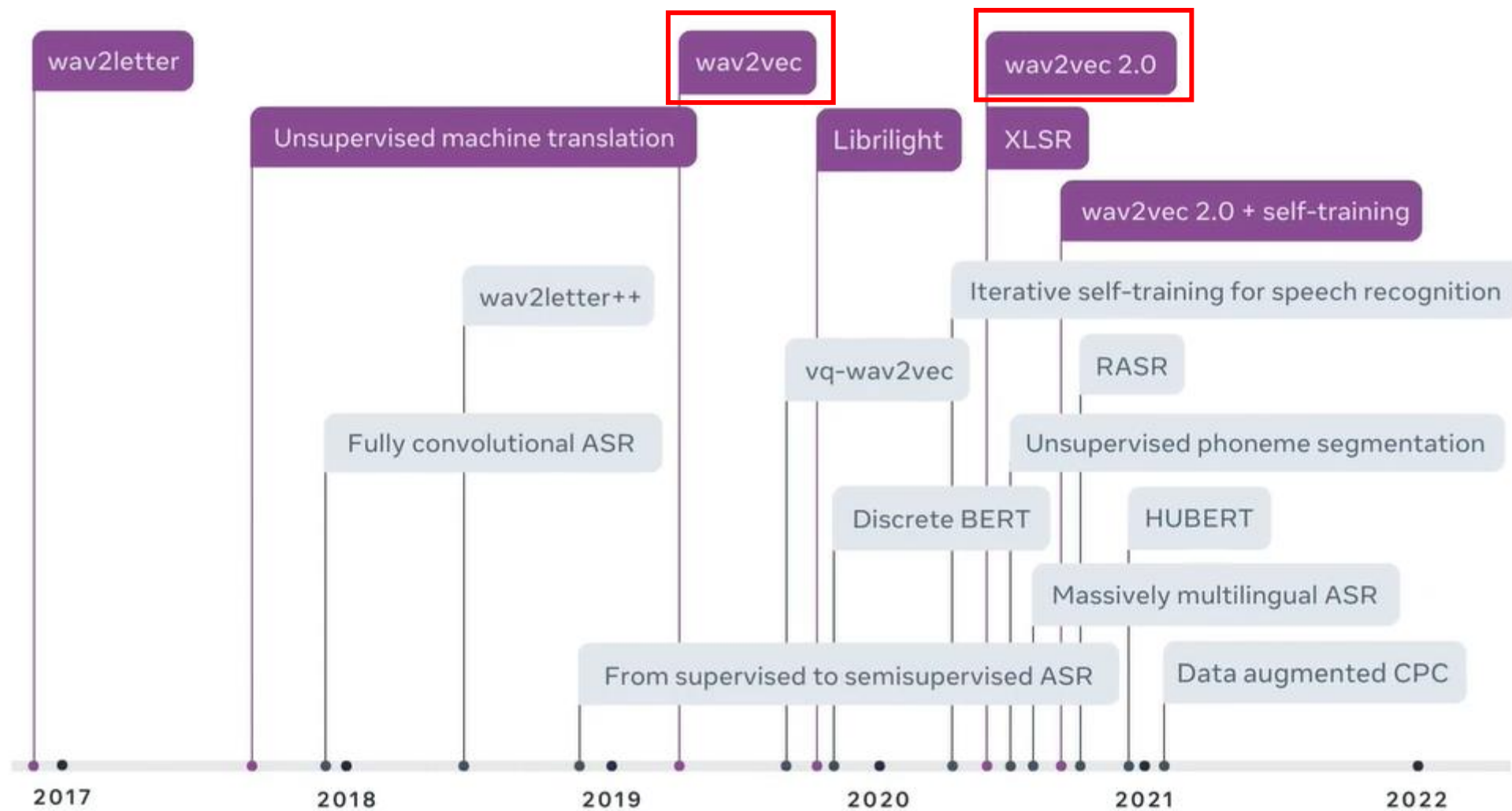


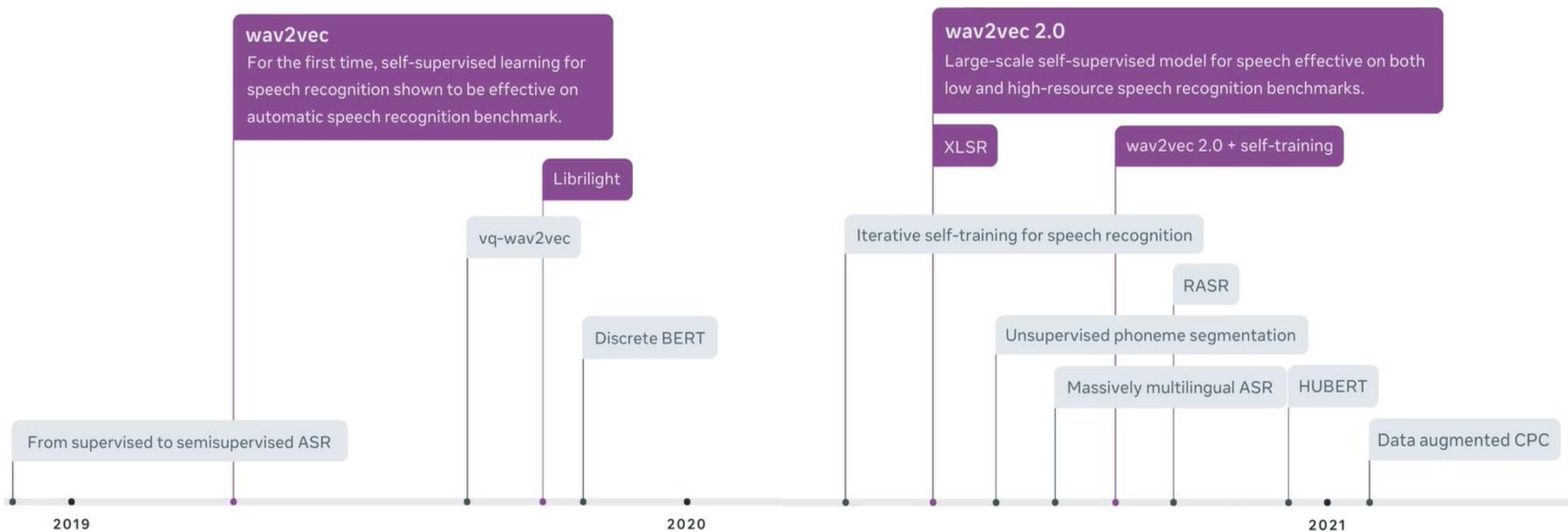
Wav2Vec 2.0

D



HuBERT





WAV2VEC: UNSUPERVISED PRE-TRAINING FOR SPEECH RECOGNITION

Steffen Schneider, Alexei Baevski, Ronan Collobert, Michael Auli
Facebook AI Research

ABSTRACT

We explore unsupervised pre-training for speech recognition by learning representations of raw audio. wav2vec is trained on large amounts of unlabeled audio data and the resulting representations are then used to improve acoustic model training. We pre-train a simple multi-layer convolutional neural network optimized via a noise contrastive binary classification task. Our experiments on WSJ reduce WER of a strong character-based log-mel filterbank baseline by up to 36 % when only a few hours of transcribed data is available. Our approach achieves 2.43 % WER on the nov92 test set. This outperforms Deep Speech 2, the best reported character-based system in the literature while using two orders of magnitude less labeled training data.¹

Wav2vec

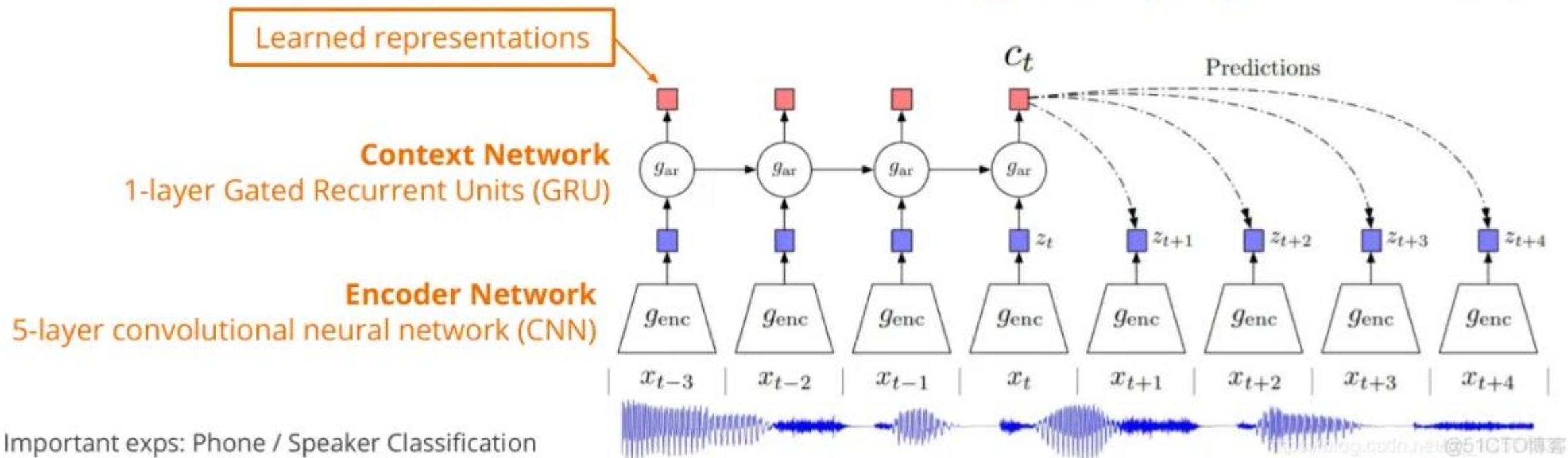
- Apply unsupervised pre-training to improve supervised speech recognition.
- Exploiting unlabeled audio data which is much easier to collect than labeled data.
- Wav2vec is trying to predict the future of an audio sequence.
 - First, pre-train a large network on unlabeled data to learn useful contextual representations of the text/audio sequence.
 - Second, use these pre-trained representations for a variety of tasks for which not enough data is available.
- This enables exploiting unlabeled audio data which is much easier to collect than labeled data.
- Wav2vec, is a convolutional neural network that takes raw audio as input and computes a general representation that can be input to a speech recognition system.
- The objective is a contrastive loss that requires distinguishing a true future audio sample from negatives.
- Different to previous work (van den Oord et al., 2018), we move beyond frame-wise phoneme classification and apply the learned representations to improve strong supervised ASR systems.
- Wav2vec relies on a fully convolutional architecture which can be easily parallelized over time on modern hardware compared to recurrent models used in previous work
- Wav2vec adjusts the CPC structure to a fully convolutional architecture, enabling easy parallelization over time on hardware. One CNN encodes the raw waveform into audio representations for each time step, and the other captures global context information into a context vector

Wav2vec

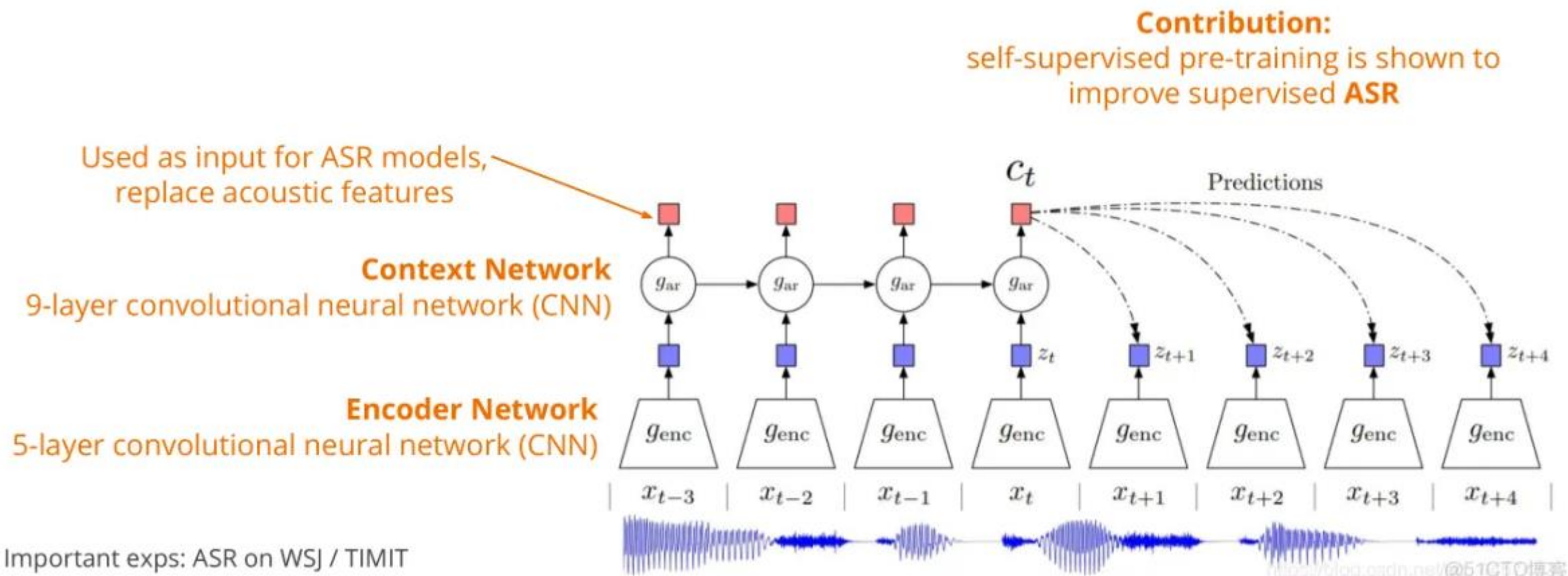
- At the core of wav2vec are two distinct networks:
 - the encoder network
 - the context network.
- Both are convolutional neural networks, albeit with different settings.
- Our model takes raw audio signal as input and then applies two networks.
- The encoder network embeds the audio signal in a latent space and
- the context network combines multiple time-steps of the encoder to obtain contextualized representations.
- **The encoder network** reduces the dimensionality of the speech data, by encoding 30 milliseconds of audio into a 512-dimensional feature vector \mathbf{z}_t at each timestep t , every 10 ms.
- **The context network** takes as input the encoder output features, encoding 210 ms of raw audio into another 512-dimensional feature vector \mathbf{c}_t .
- The objective is to aggregate information over a longer timeframe to model higher-order information. This network outputs *contextual representations* c_t that are used to predict future audio samples.

CPC

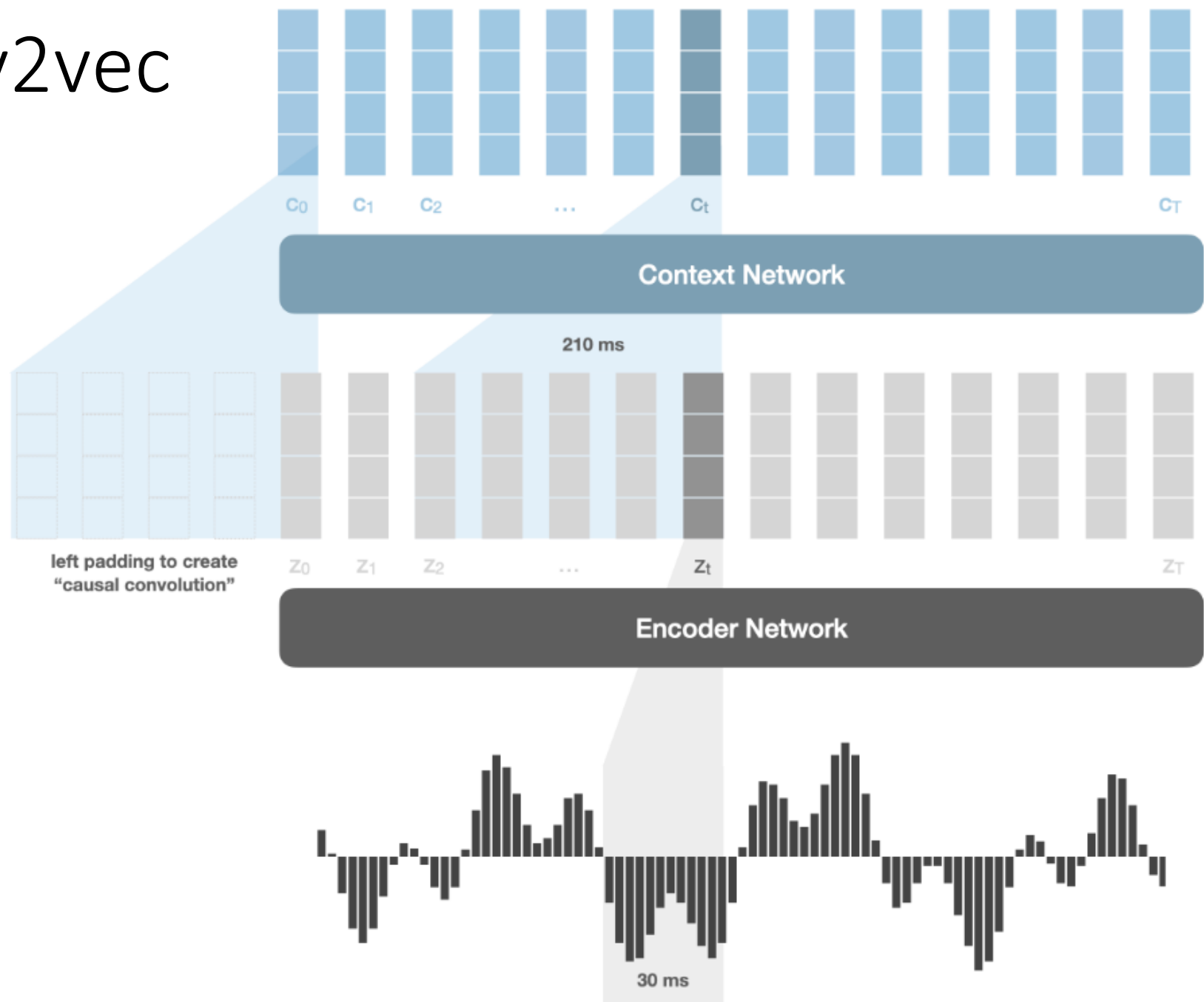
Intuition:
Pulls temporally nearby representations closer
and pushes temporally distant ones further.



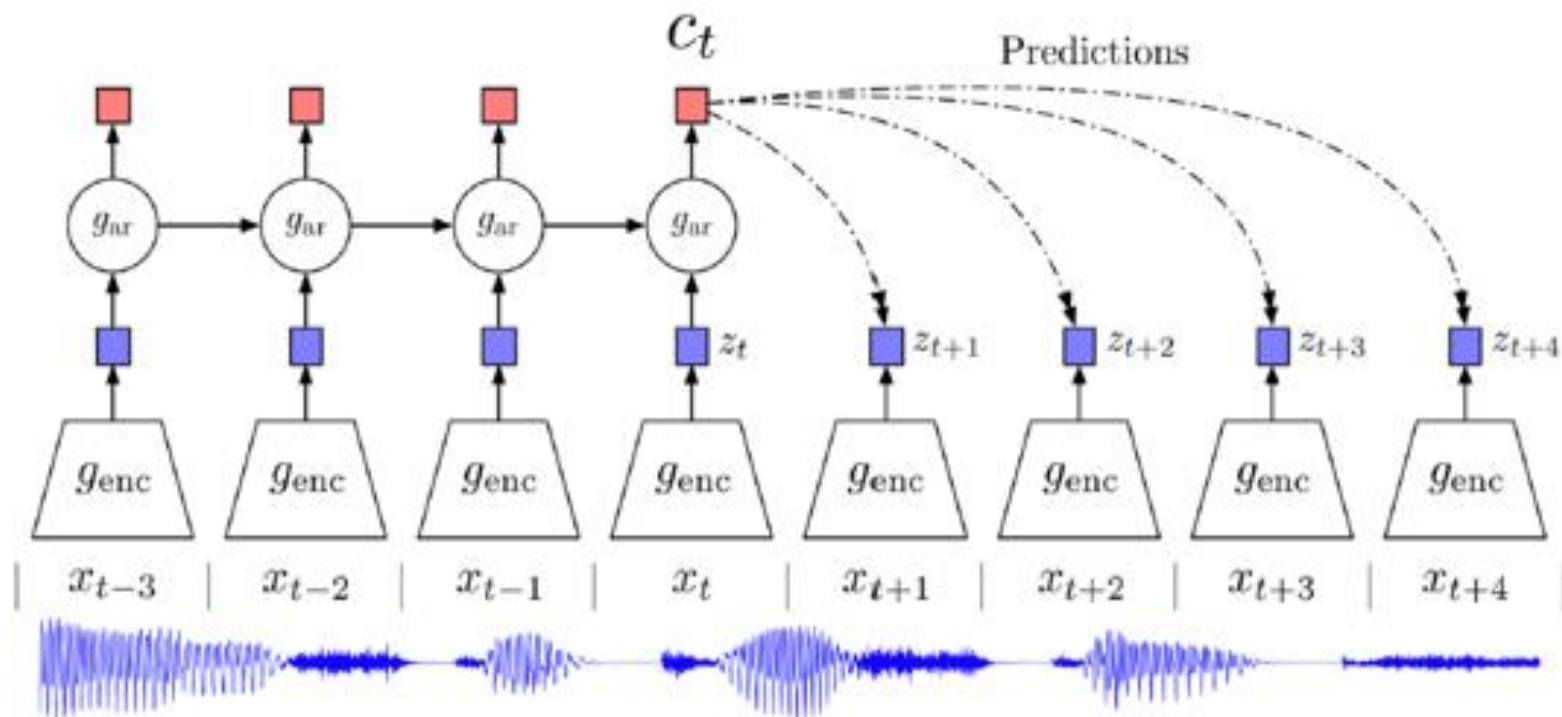
Wav2vec



Wav2vec

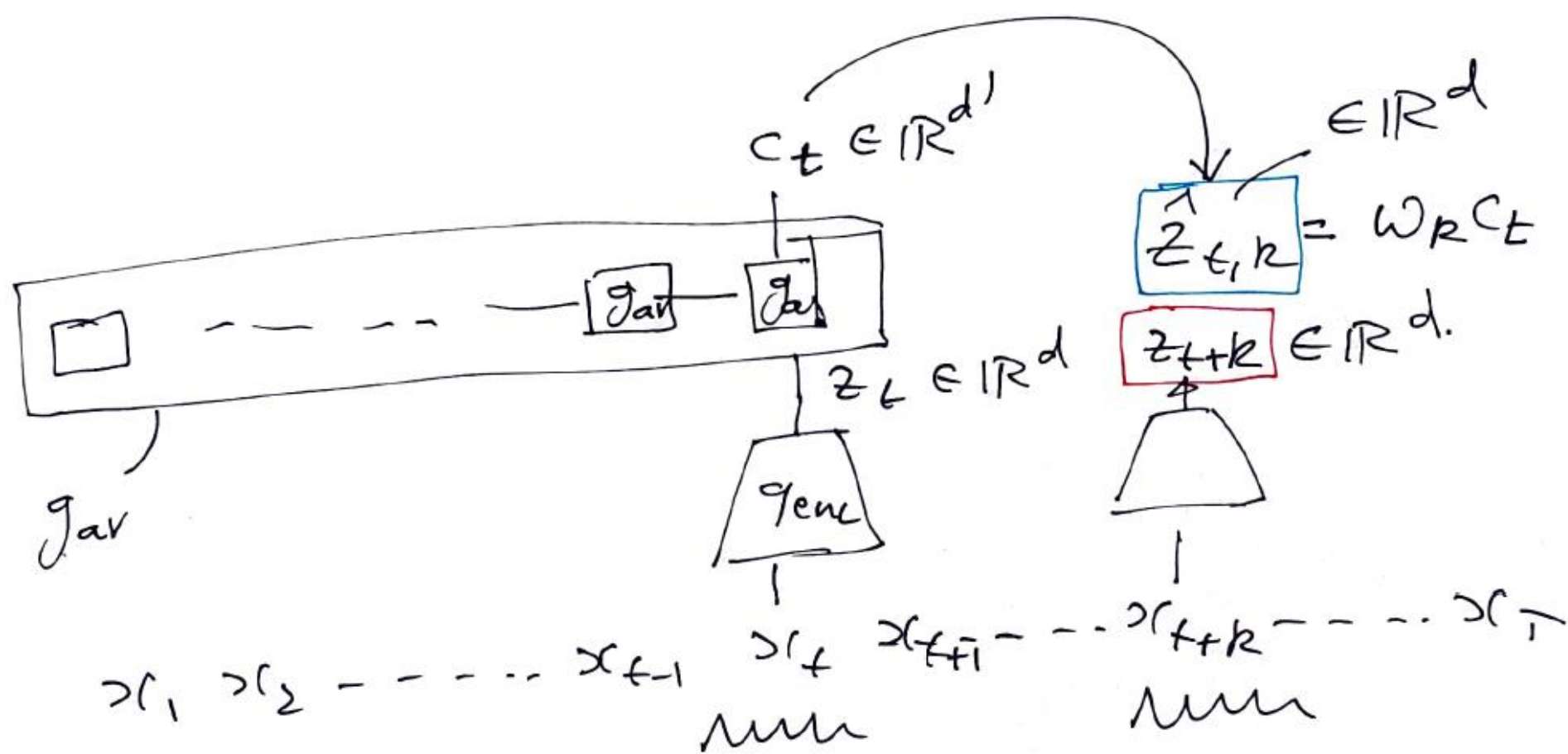


CPC



$$f_k(x_{t+k}, c_t) = \exp \left(z_{t+k}^T W_k c_t \right)$$

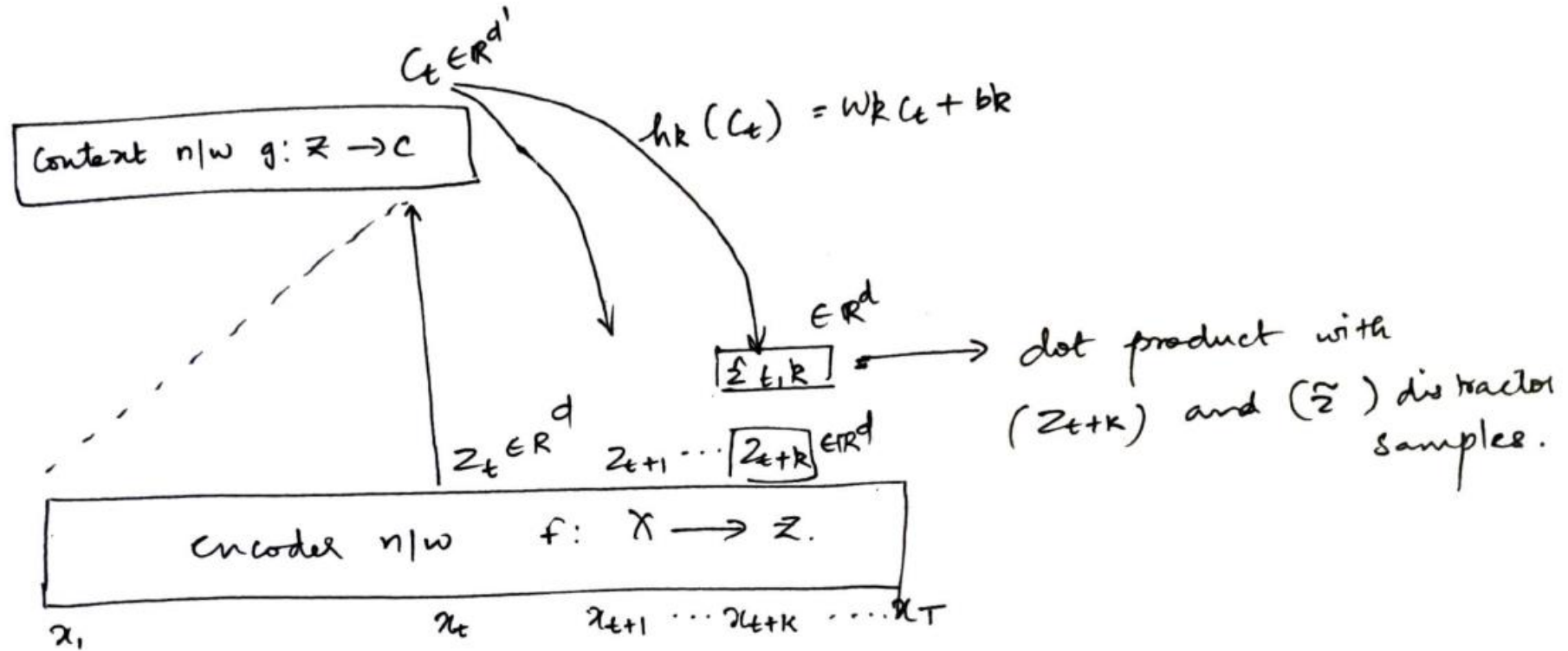
$$\mathcal{L}_N = -\mathbb{E}_X \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$



$$z_t = g_{enc}(x_t) \rightarrow 10ms \text{ latent variable}$$

$$c_t = g_{var}(z_{\leq t})$$

Notations \rightarrow wav2vec



$$z_t = f_{\text{enc}}(x_t).$$

$$c_t = g_{\text{cont}}(z_t \dots z_{t-v}) \text{ for receptive field size } v'.$$

$\boxed{z_{t+k}}$ \rightarrow distinguish the sample z_{t+k} ; that is k steps
in future from distractor samples \tilde{z} ;

\tilde{z} drawn from a proposal distribution p_n ;

where $p_n(z) = \frac{1}{T}$; where T is the
sequence length of the audio.

In practice, sample ten negative examples
by uniformly choosing distractors from each audio sequence.

distinguish z_{t+k} from \tilde{z} ; by minimizing
the contrastive loss for each step

$k = 1, \dots, K$.

Prediction process is
done ' K ' times at
each timestep ' t '.

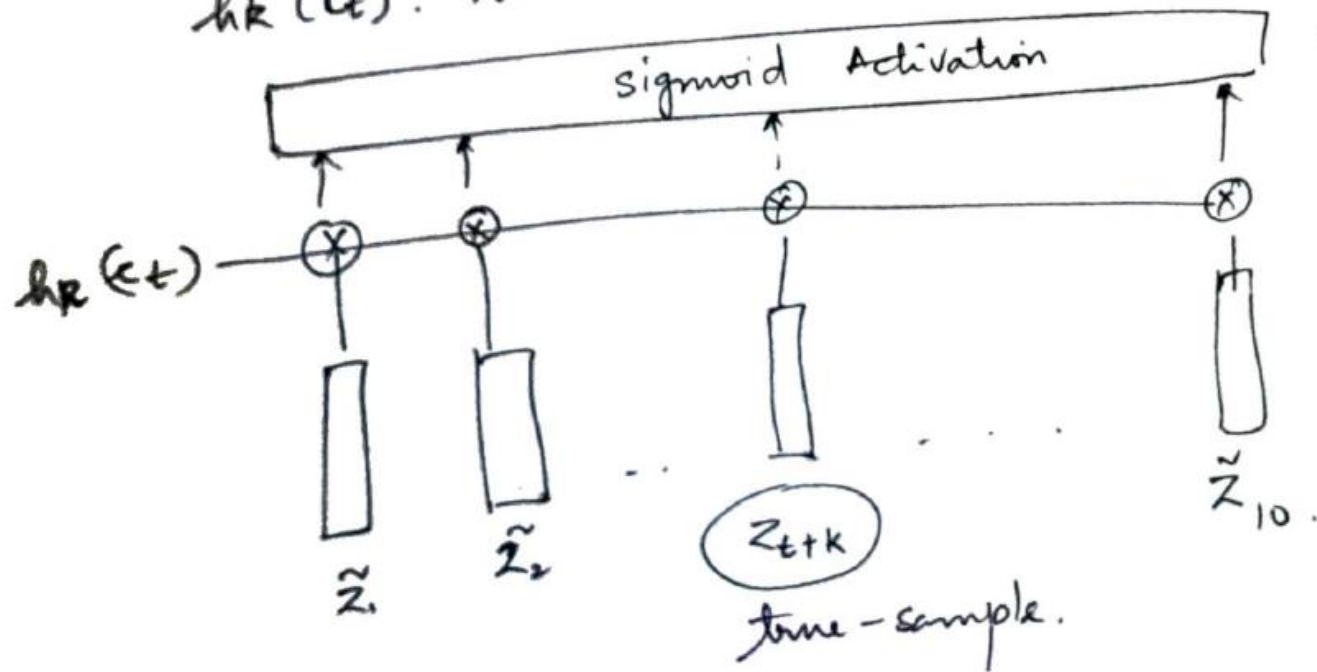
True sample : z_{t+k} .

distractor samples : $\tilde{z}_1, \tilde{z}_2 \dots \tilde{z}_{10}$.

for time-step t ; context vector c_t .

undergoes step-specific affine transformation $h_k c(t)$

$h_k(c_t)$. similarity with z_{t+k} and \tilde{z}



$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

dot - prod.

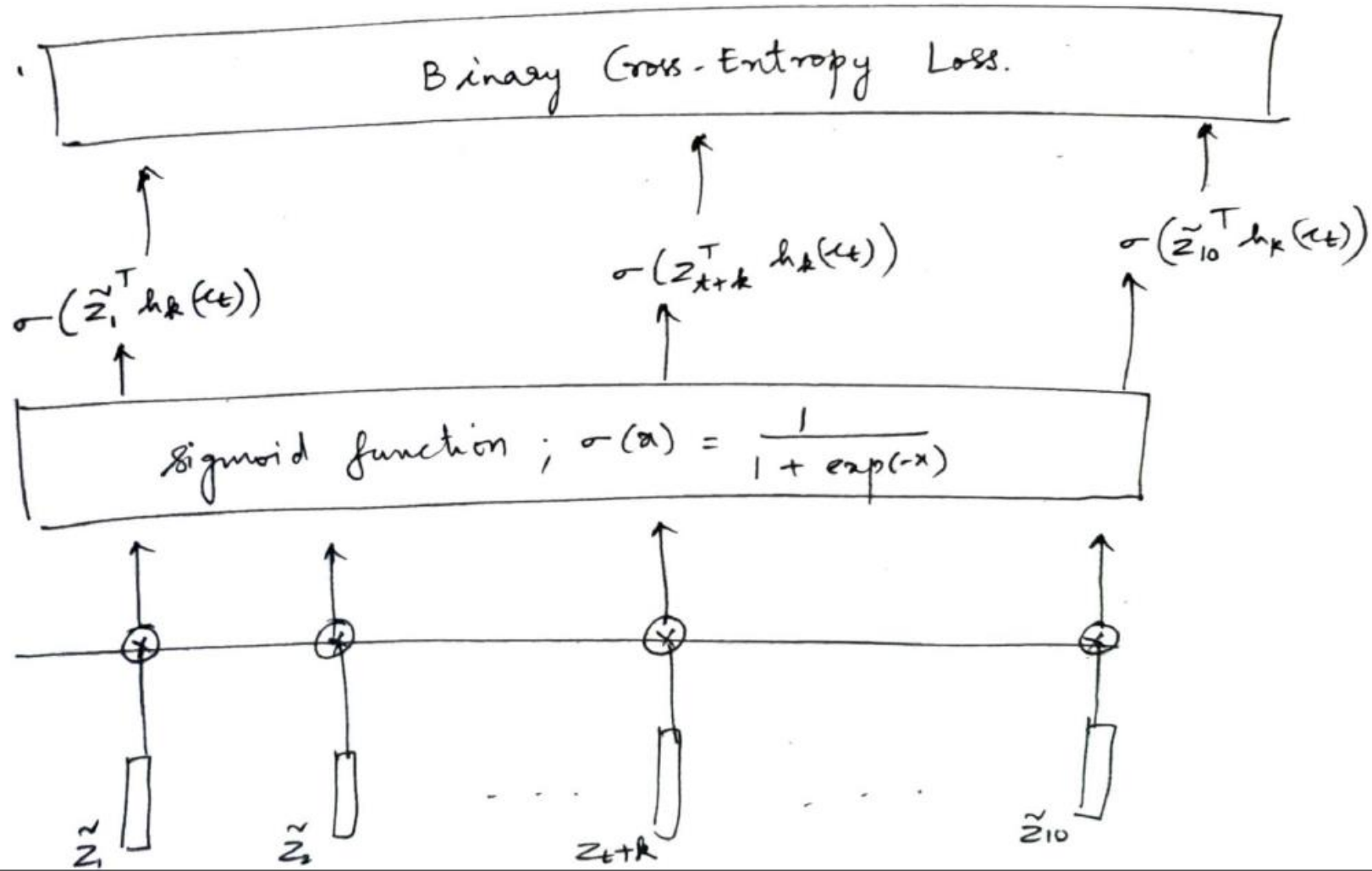
$$\left[\frac{z_{t+k}^T h_k(c_t)}{\left[\frac{\tilde{z}^T h_k(c_t)}{\right]} \right]$$

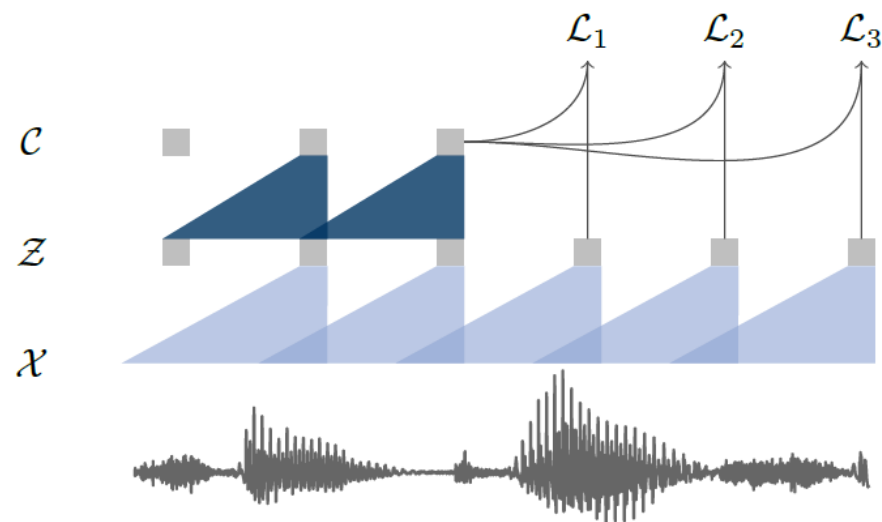
$$\left[\frac{\tilde{z}^T h_k(c_t)}{\right]$$

at time step 't'

$$L_k = -\left(\log \sigma(z_{t+k}^T h_k(t)) + \lambda \sum_{\tilde{z} \sim p_n} \left[\log \sigma(-\tilde{z}^T h_k(t)) \right] \right)$$

$L = \sum \text{loss of positive pairs} + \sum \text{loss of negative pairs}$





Total Loss

$$L_k = - \sum_{t=1}^{T-k} \left(\log \sigma(z_{t+k}^T h_k(t)) + \lambda E_{\tilde{z} \sim p_n} \left[\log \sigma(-\tilde{z}^T h_k(t)) \right] \right)$$

and for $k=1, \dots, K$

$$L = \sum_{k=1}^K L_k.$$

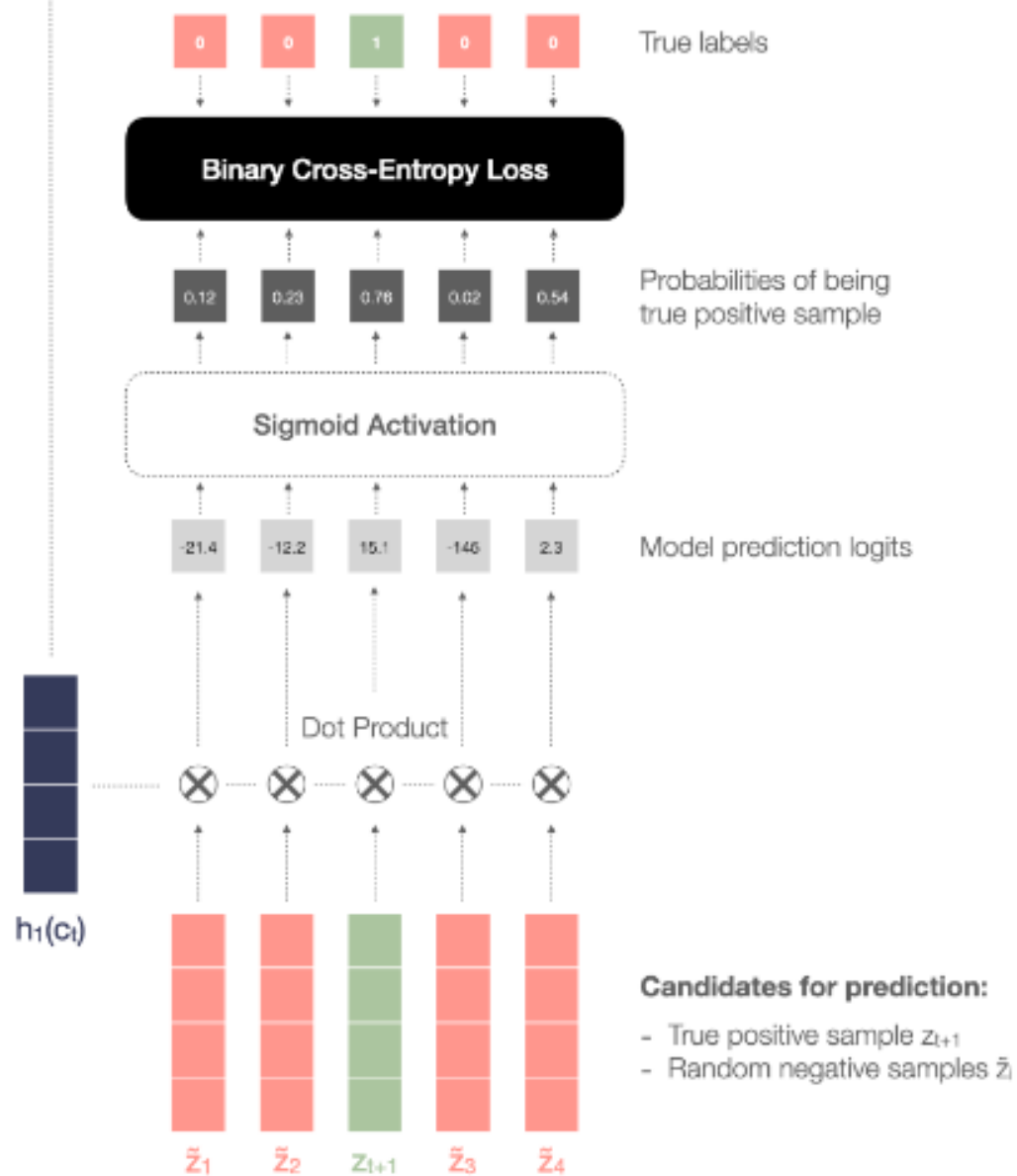
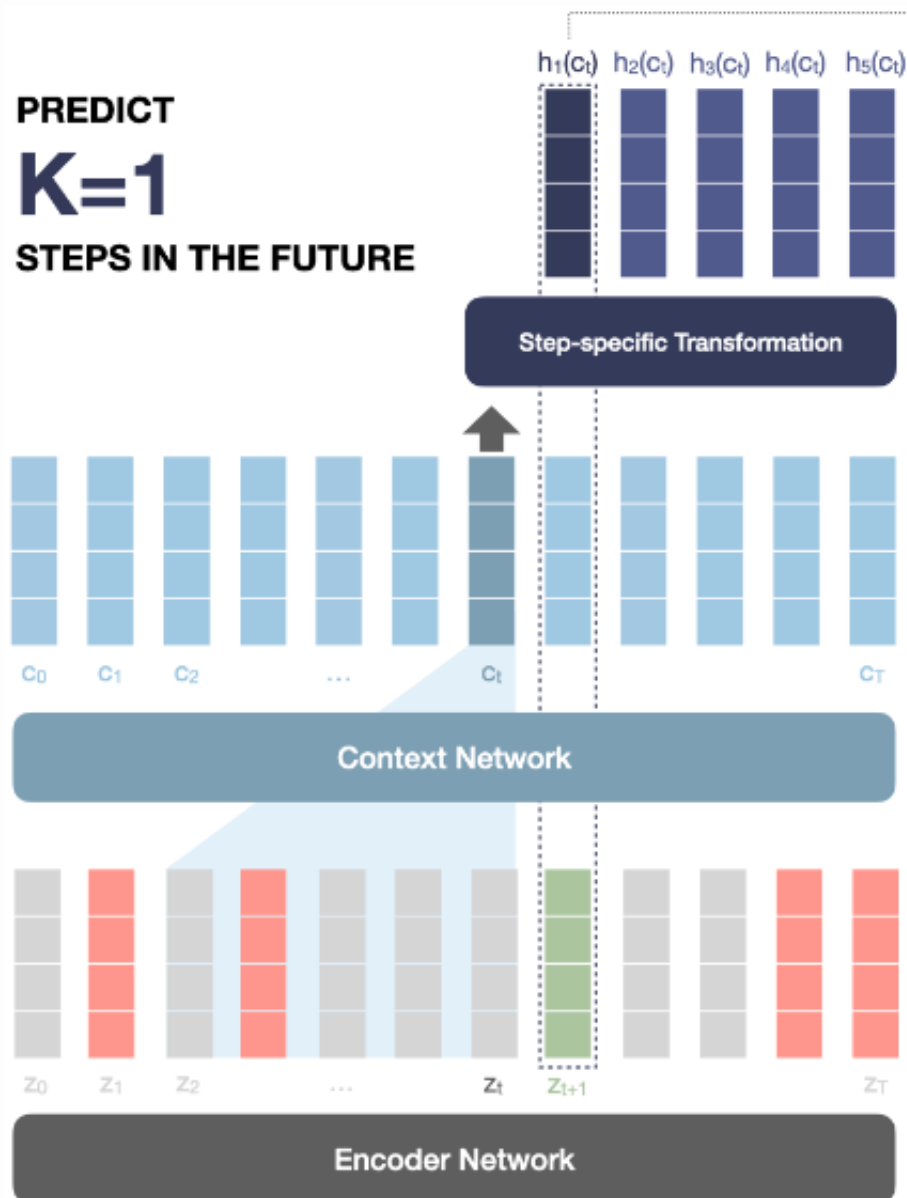
At a given timestep t , for each step $k = 1 \dots 12$, we do the following:

1. Extract the true audio sample \mathbf{z}_{t+k} at future step k
2. Pick 10 random *negative* samples $\tilde{\mathbf{z}}_{1..10}$ from the *same audio sequence*
3. Compute a step k -specific transformation $\mathbf{h}_k(\mathbf{c}_t)$ of the context vector at time t
4. Compute the similarity (dot product) of the transformed context vector with all \mathbf{z} candidates
5. Compute the final probabilities of positive/negative through a sigmoid activation
6. Compare with the ground truth and penalize the model for wrong predictions (binary cross-entropy loss)

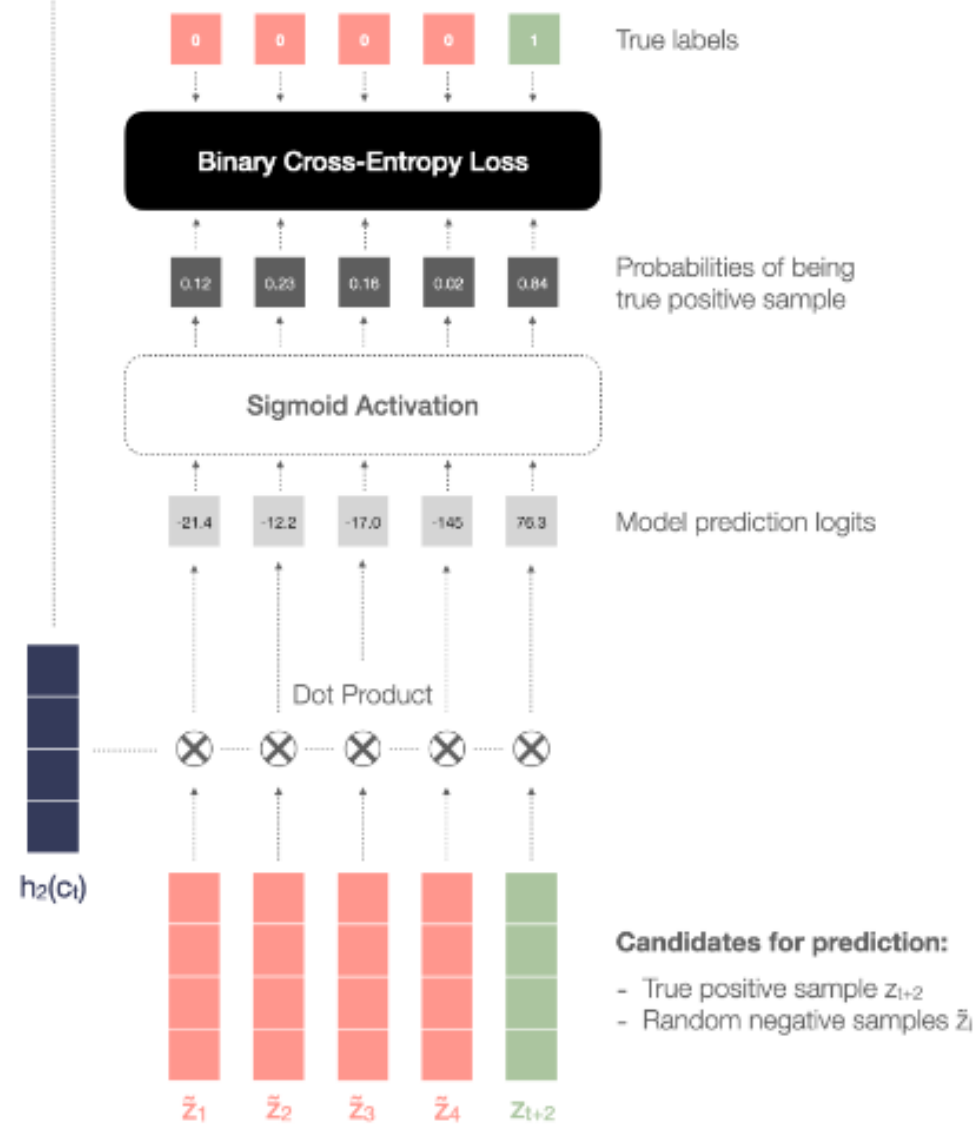
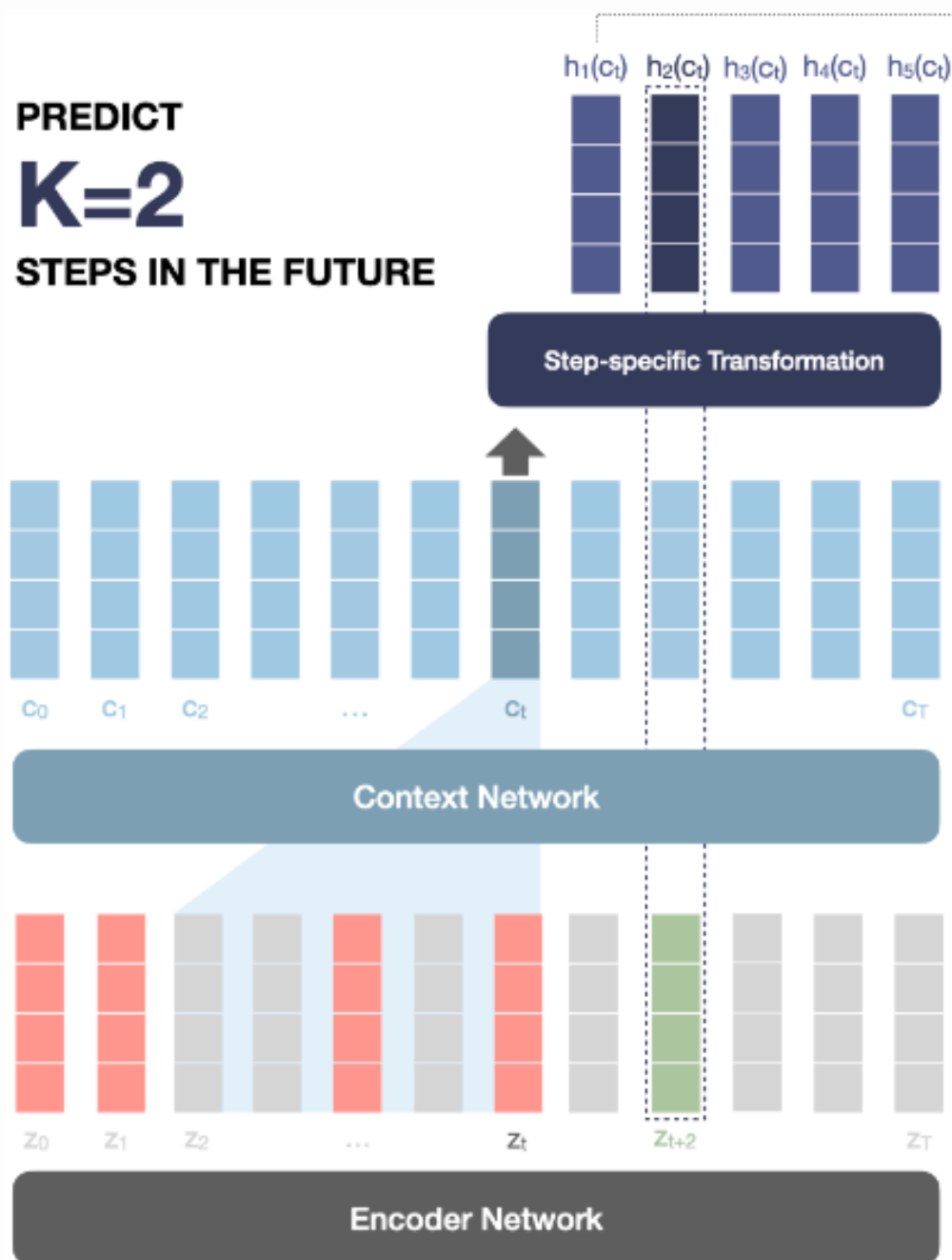
PREDICT

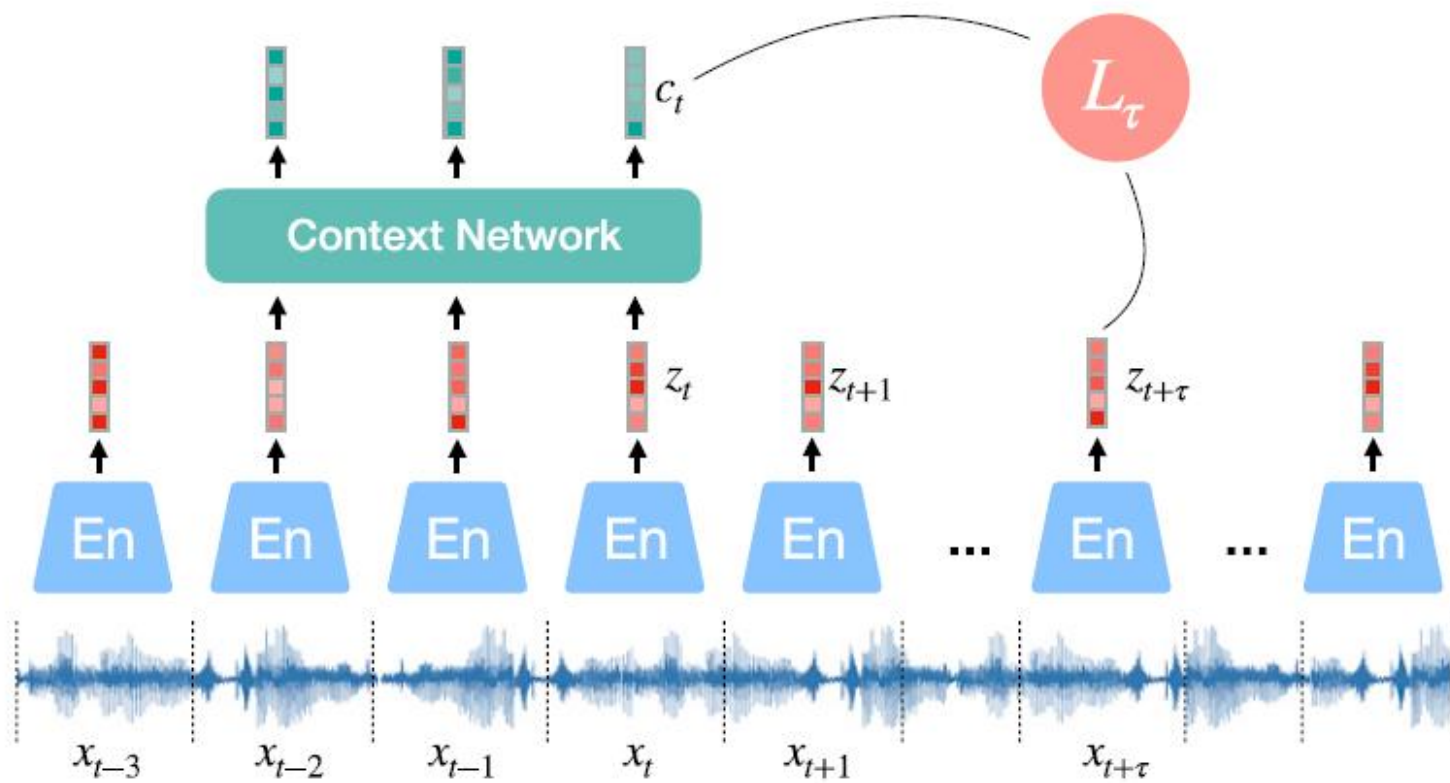
$K=1$

STEPS IN THE FUTURE



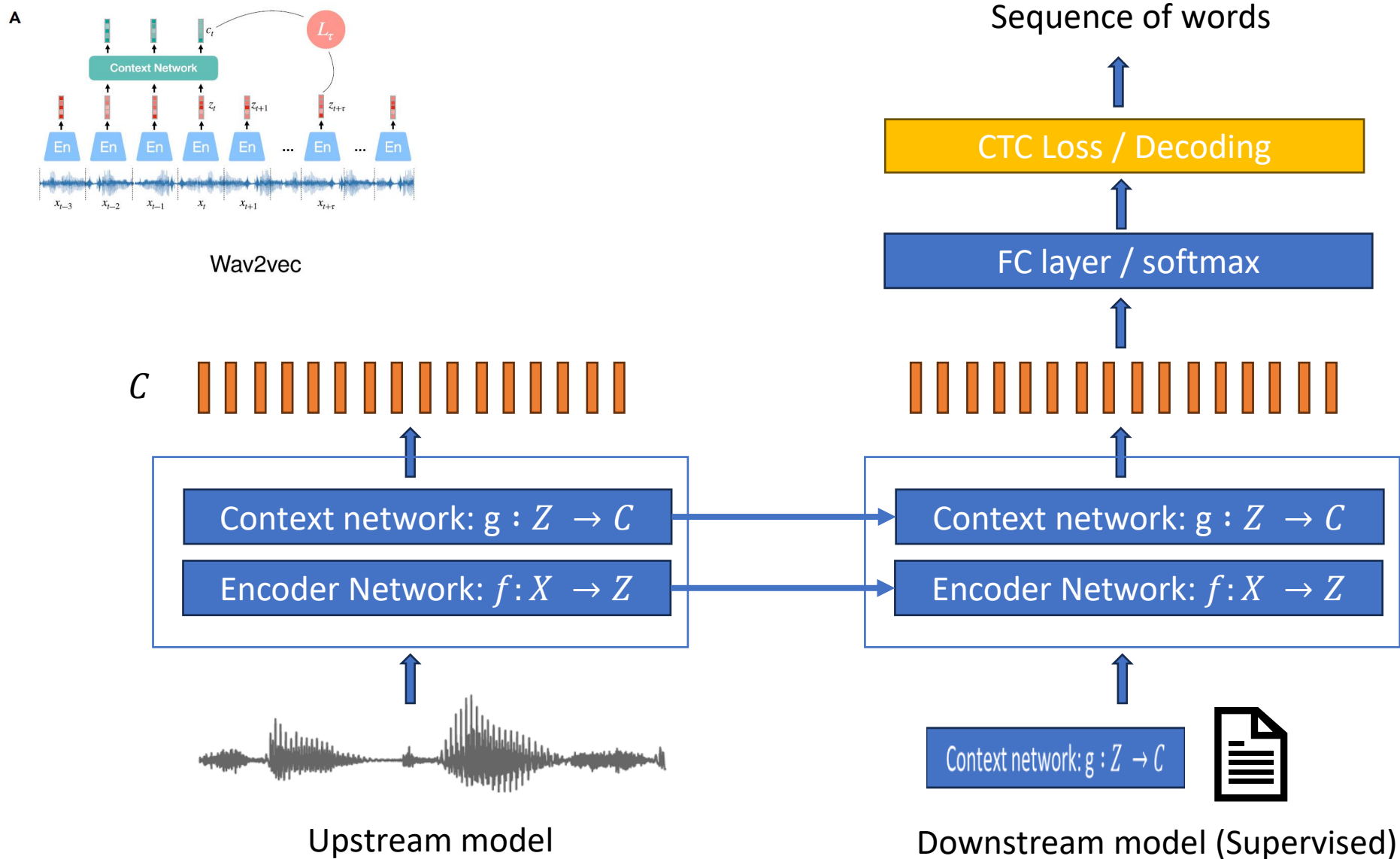
PREDICT
 $K=2$
STEPS IN THE FUTURE



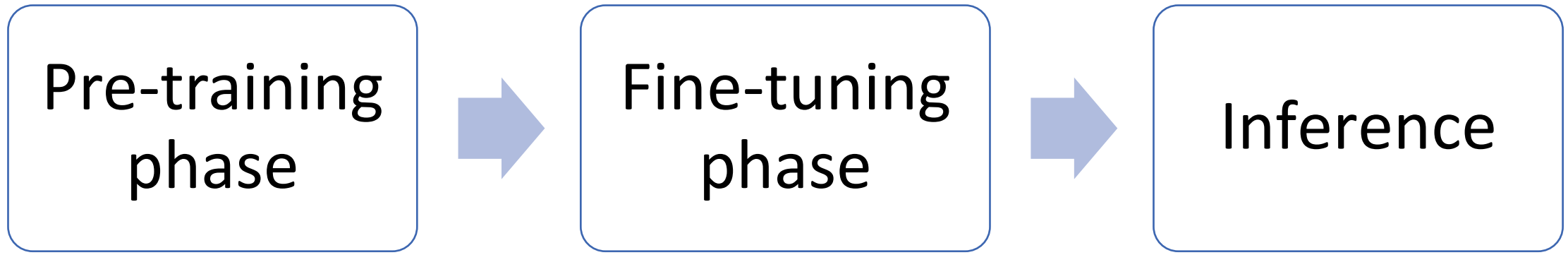
A

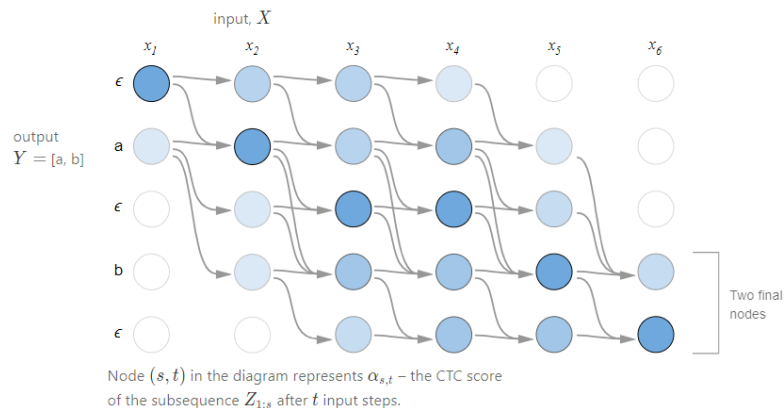
Wav2vec

After training, we input the representations c_i produced by the context network to the acoustic model instead of log-mel filterbank features.



Foundation models in ASR





$$Z = z_1 z_2 z_3 \dots z_M (M \leq T)$$

CTC Loss

$$\hat{y} = \hat{y}_1 \hat{y}_2 \hat{y}_3 \dots \hat{y}_t \dots \hat{y}_T$$

FC layer / softmax

\mathcal{C}

Context network: $g : Z \rightarrow \mathcal{C}$

Encoder Network: $f : X \rightarrow Z$

Initialized with
pre-trained
Weights

Context network: $g : Z \rightarrow \mathcal{C}$

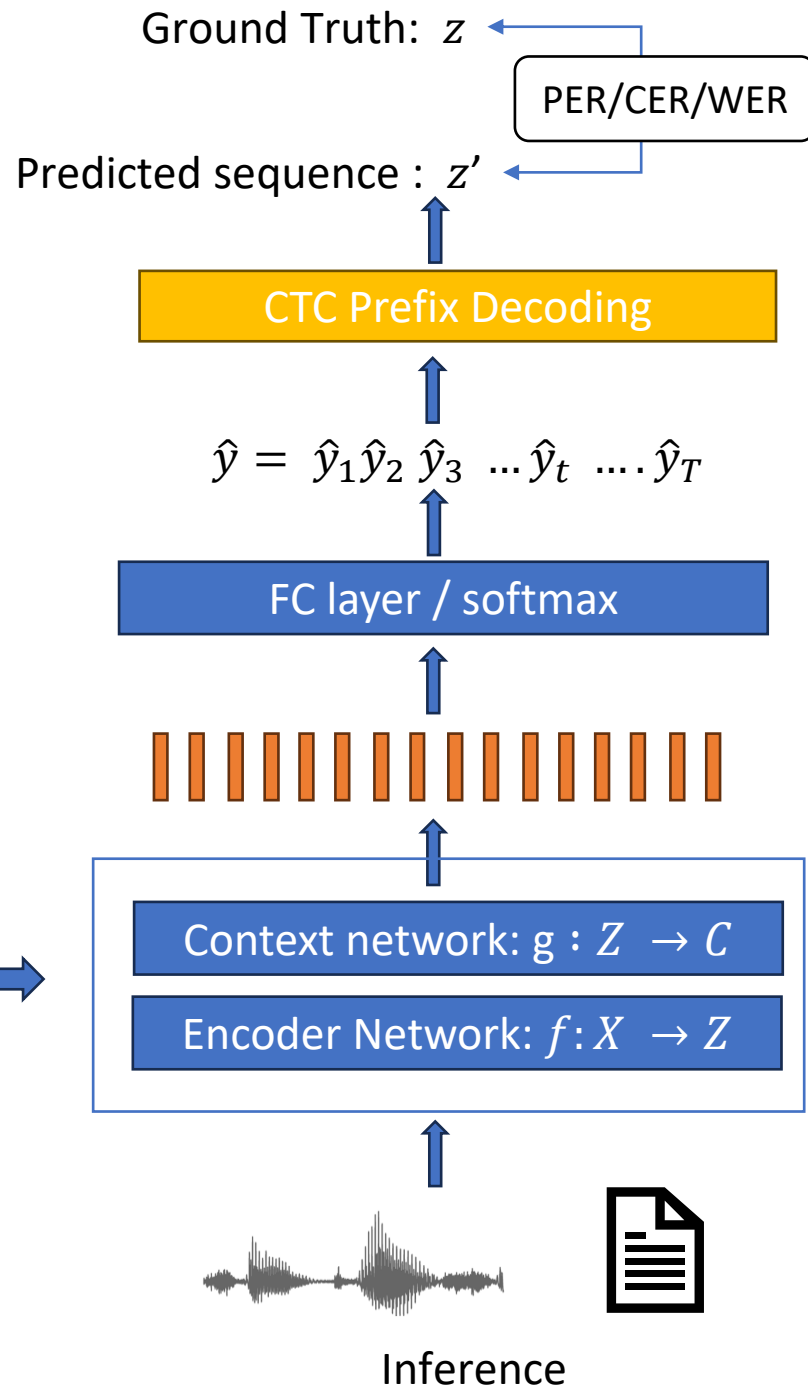
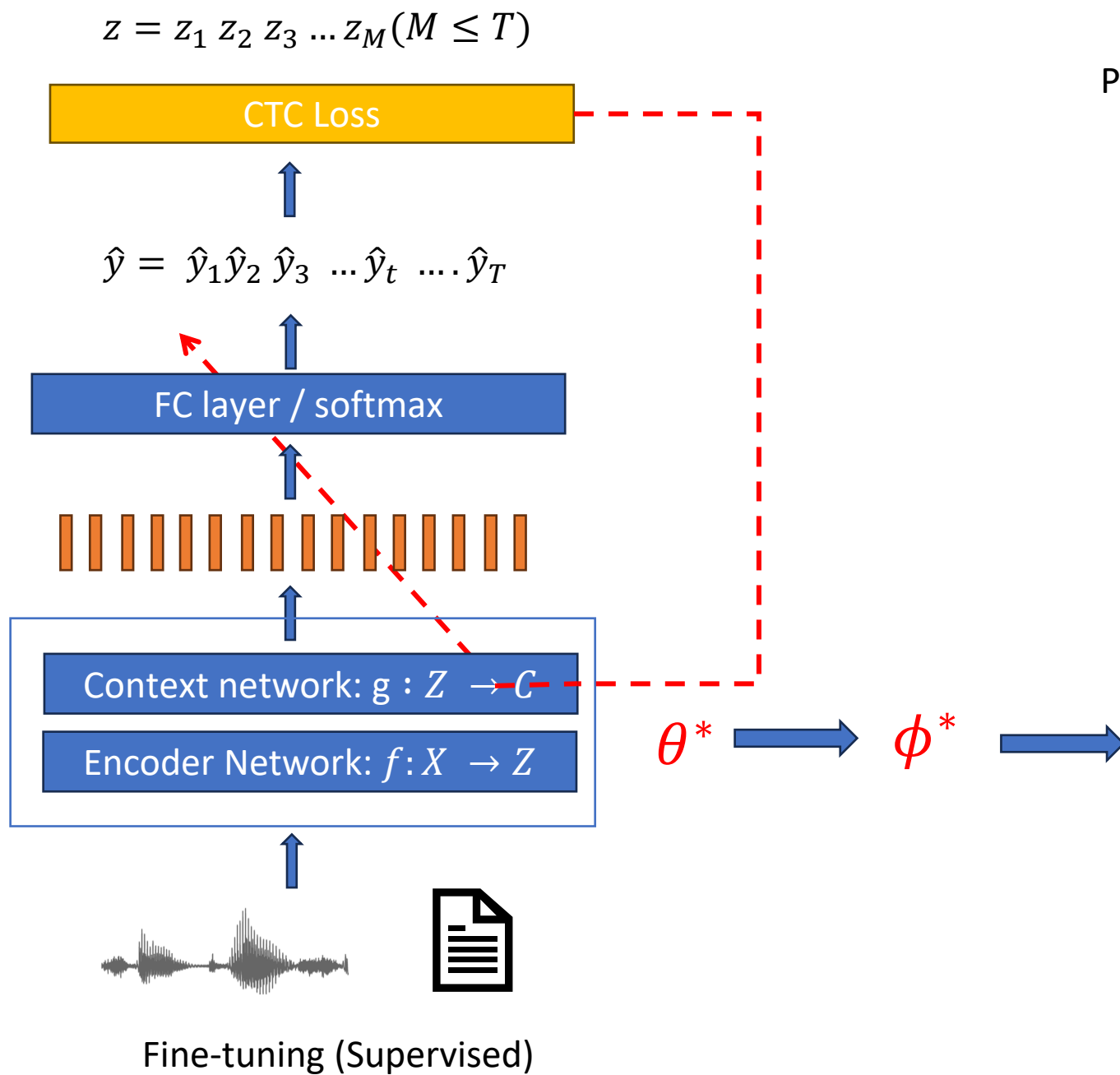
Encoder Network: $f : X \rightarrow Z$

θ^*

ϕ^*

Pre-trained model

Fine-tuning (Supervised)



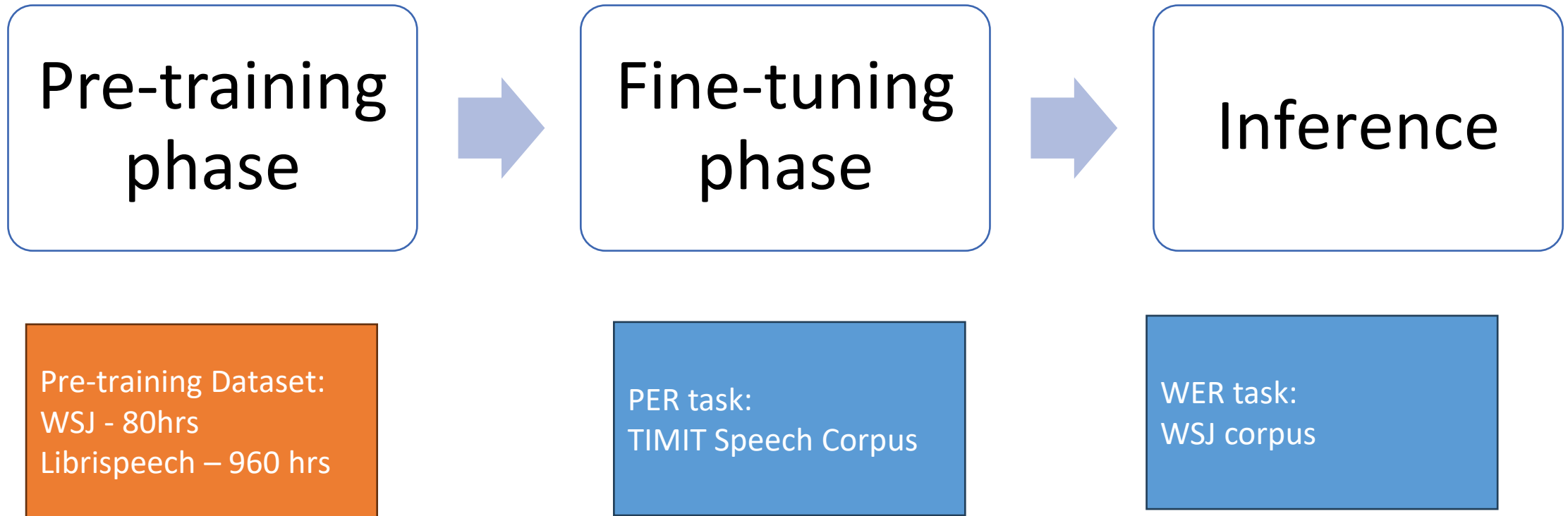
3 EXPERIMENTAL SETUP

3.1 DATA

We consider the following corpora: For phoneme recognition on TIMIT (Garofolo et al., 1993b) we use the standard train, dev and test split where the training data contains just over three hours of audio data. Wall Street Journal (WSJ; Garofolo et al. (1993a); Woodland et al. (1994)) comprises about 81 hours of transcribed audio data. We train on si284, validate on nov93dev and test on nov92. Librispeech (Panayotov et al., 2015) contains a total of 960 hours of clean and noisy speech for training. For pre-training, we use either the full 81 hours of the WSJ corpus, an 80 hour subset of clean Librispeech, the full 960 hour Librispeech training set or a combination of all of them.

To train the baseline acoustic model we compute 80 log-mel filterbank coefficients for a 25 ms sliding window with stride 10 ms. Final models are evaluated in terms of both word error rate (WER) and letter error rate (LER).

Results



PER Results

	dev	test
CNN + TD-filterbanks (Zeghidour et al., 2018a)	15.6	18.0
Li-GRU + MFCC (Ravanelli et al., 2018)	–	16.7 ± 0.26
Li-GRU + FBANK (Ravanelli et al., 2018)	–	15.8 ± 0.10
Li-GRU + fMLLR (Ravanelli et al., 2018)	–	14.9 ± 0.27
Baseline	16.9 ± 0.15	17.6 ± 0.11
wav2vec (Librispeech 80h)	15.5 ± 0.03	17.6 ± 0.12
wav2vec (Librispeech 960h)	13.6 ± 0.20	15.6 ± 0.23
wav2vec (Librispeech + WSJ)	12.9 ± 0.18	14.7 ± 0.42

Table 2: Results for phoneme recognition on TIMIT in terms of PER. All our models use the CNN-8L-PReLU-do0.7 architecture (Zeghidour et al., 2018a).

negatives	dev PER	train time (h)
1	16.3	6.1
2	15.8	6.3
5	15.9	8.2
10	15.5	10.5
20	15.7	15.3

Table 3: Effect of different number of negative samples during pre-training for TIMIT on the development set.

WER Results

			nov93dev		nov92	
			LER	WER	LER	WER
Deep Speech 2 (12K h labeled speech; Amodei et al., 2016)			-	4.42	-	3.1
Trainable frontend (Zeghidour et al., 2018a)			-	6.8	-	3.5
Lattice-free MMI (Hadian et al., 2018)			-	5.66 [†]	-	2.8 [†]
Supervised transfer-learning (Ghahremani et al., 2017)			-	4.99 [†]	-	2.53 [†]
4-GRAM LM (Heafield et al., 2013)						
Baseline	–	–	3.32	8.57	2.19	5.64
wav2vec	Librispeech	80 h	3.71	9.11	2.17	5.55
wav2vec	Librispeech	960 h	2.85	7.40	1.76	4.57
wav2vec	Libri + WSJ	1,041 h	2.91	7.59	1.67	4.61
wav2vec large	Librispeech	960 h	2.73	6.96	1.57	4.32
WORD CONVLM (Zeghidour et al., 2018b)						
Baseline	–	–	2.57	6.27	1.51	3.60
wav2vec	Librispeech	960 h	2.22	5.39	1.25	2.87
wav2vec large	Librispeech	960 h	2.13	5.16	1.02	2.53
CHAR CONVLM (Likhomanenko et al., 2019)						
Baseline	–	–	2.77	6.67	1.53	3.46
wav2vec	Librispeech	960 h	2.14	5.31	1.15	2.78
wav2vec large	Librispeech	960 h	2.11	5.10	0.99	2.43

Thank you

<https://ai.meta.com/blog/wav2vec-state-of-the-art-speech-recognition-through-self-supervision/>