# End-to-End Text-to-Speech using Latent Duration based on VQ-VAE

Yusuke Yasuda, Xin Wang, Junichi Yamagishi
National institute of informatics, Japan
The graduate university for advanced studies, SOKENDAI, Japan

ICASSP 2021

# Background

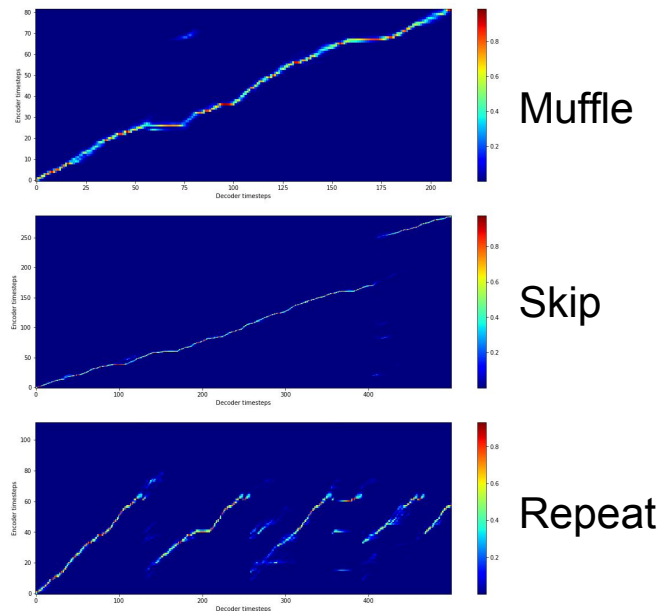End-to-end TTS is a method of converting text to speech directly in single model.

☆Issues

1) Alignments predicted from soft-attention tend to be unstable.

☆Solutions
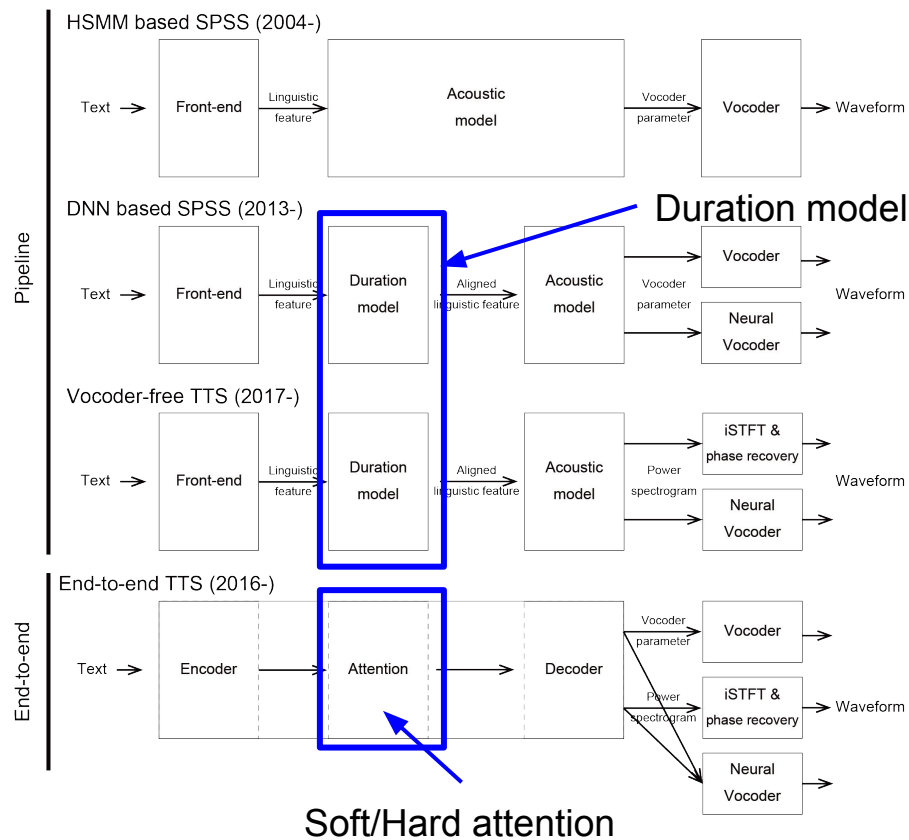
1) Construct monotonic alignments from phoneme duration.



Muffle

Skip

Repeat

Errors from soft-attention.

# Our approach: simplification of duration based TTS

| Method | Teacher-Student | Training phases | Aligner | Aligner is external | Duration form | Latent duration |
|---|---|---|---|---|---|---|
| FastSpeech | ✓ | 3 | Soft-attention | ✓ | continuous | |
| DurIAn | | 3 | HMM? | ✓ | continuous | |
| FastSpeech2 | | 3 | HMM-GMM | ✓ | continuous | |
| AlignTTS | | 4 | SSNT-like MDN | ✓ | discrete | |
| JDI-T | ✓ | 1 | Soft-attention + CTC | | continuous | |
| Glow-TTS | | 1 | MAS | | continuous | |
| Non-Attentive Tacotron | | 2 | HMM | ✓ | continuous | ? |
| VQ-VAE | | 1 | CTC | | discrete | ✓ |

- **All modules are jointly trainable**
  - Single training phase
- **Discrete duration**
  - Conform with forced aligner and upsampling
  - No length regularizer or ceiling of duration
- **Duration is modeled as latent variables**
  - Simple training criterion based on VAE

# Alignment methods in TTS frameworks
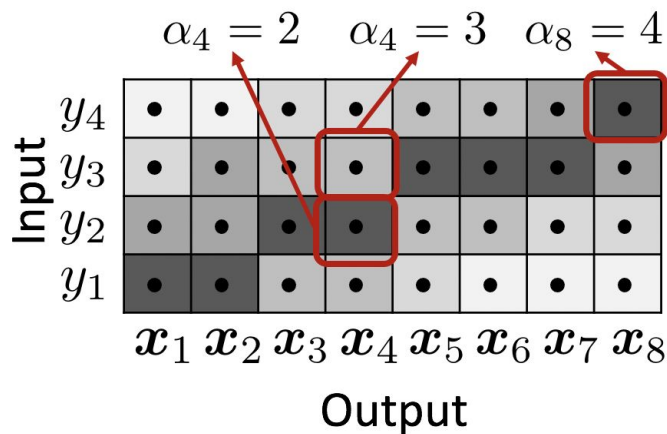


- Pipeline TTS
  - Consists of multiple models to convert texts into speech.
  - Each model is dedicated to a single function.
  - Duration model forms alignments between source texts and target speech.
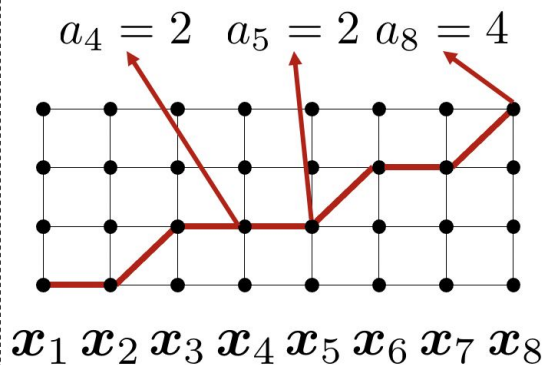
- End-to-end TTS
  - Consists of a single model to convert texts into speech
  - Attention mechanism forms alignments between source texts and target speech.
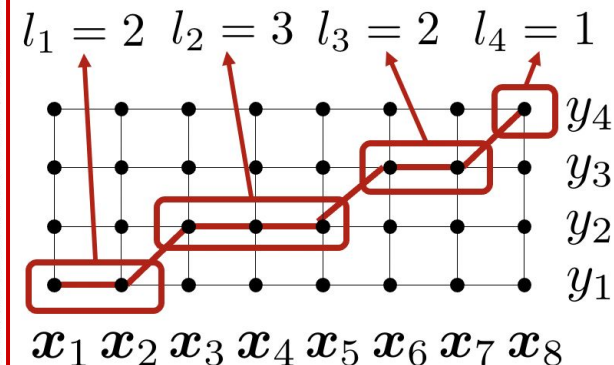
# Design of a latent alignments based on duration

## Soft-attention

$\alpha_4 = 2 \quad \alpha_4 = 3 \quad \alpha_8 = 4$

Input

$y_4$
$y_3$
$y_2$
$y_1$

$\boldsymbol{x}_1 \, \boldsymbol{x}_2 \, \boldsymbol{x}_3 \, \boldsymbol{x}_4 \, \boldsymbol{x}_5 \, \boldsymbol{x}_6 \, \boldsymbol{x}_7 \, \boldsymbol{x}_8$

Output

## Hard-attention

$a_4 = 2 \quad a_5 = 2 \quad a_8 = 4$

$\boldsymbol{x}_1 \, \boldsymbol{x}_2 \, \boldsymbol{x}_3 \, \boldsymbol{x}_4 \, \boldsymbol{x}_5 \, \boldsymbol{x}_6 \, \boldsymbol{x}_7 \, \boldsymbol{x}_8$

## Proposed model

$l_1 = 2 \quad l_2 = 3 \quad l_3 = 2 \quad l_4 = 1$

$y_4$
$y_3$
$y_2$
$y_1$

$\boldsymbol{x}_1 \, \boldsymbol{x}_2 \, \boldsymbol{x}_3 \, \boldsymbol{x}_4 \, \boldsymbol{x}_5 \, \boldsymbol{x}_6 \, \boldsymbol{x}_7 \, \boldsymbol{x}_8$
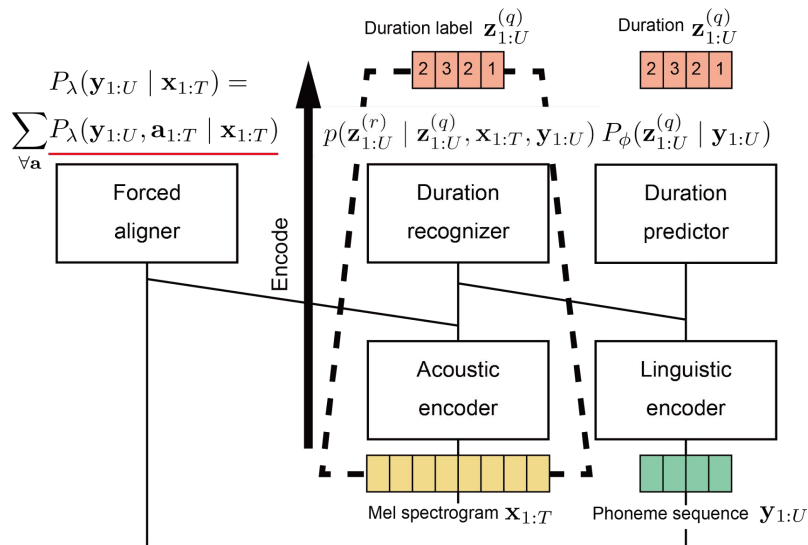
- Using phoneme duration as a latent variable
- Let the phoneme duration length be a discrete variable (the number of frames of acoustic features)

# Modeling latent duration with variational autoencoder (VAE)



$$\mathbf{z}_{1:U}^{(q)} = \arg\max_{\mathbf{a}_{1:T}} P_\lambda(\mathbf{y}_{1:U}, \mathbf{a}_{1:T} \mid \mathbf{x}_{1:T})$$

Approximate posterior $\sim Q_\psi(\mathbf{z}_{1:U}^{(r)} \mid \mathbf{z}_{1:U}^{(q)}, \mathbf{x}_{1:T}, \mathbf{y}_{1:U})$

$$P_\lambda(\mathbf{y}_{1:U} \mid \mathbf{x}_{1:T}) = \sum_{\forall \mathbf{a}} P_\lambda(\mathbf{y}_{1:U}, \mathbf{a}_{1:T} \mid \mathbf{x}_{1:T})$$

Duration label $\mathbf{z}_{1:U}^{(q)}$

Duration $\mathbf{z}_{1:U}^{(q)}$

Mel spectrogram $\mathbf{x}_{1:T}$

$p(\mathbf{z}_{1:U}^{(r)} \mid \mathbf{z}_{1:U}^{(q)}, \mathbf{x}_{1:T}, \mathbf{y}_{1:U})$   $P_\phi(\mathbf{z}_{1:U}^{(q)} \mid \mathbf{y}_{1:U})$

$p_\theta(\mathbf{x}_{1:T} \mid \mathbf{z}_{1:U}^{(q)}, \mathbf{y}_{1:U})$

Forced aligner

Encode

Duration recognizer

Duration predictor

Acoustic decoder

Decode

Acoustic encoder

Linguistic encoder

Mel spectrogram $\mathbf{x}_{1:T}$

Phoneme sequence $\mathbf{y}_{1:U}$

- Introduce phoneme duration as discrete latent variables ( $\mathbf{z}_{1:U}^{(q)}, \mathbf{z}_{1:U}^{(r)}$ )
- Sample phoneme duration from forced aligner
- Define approximate posterior with samples of phoneme duration
- Use duration predictor as prior for VAE
- Use duration recognizer as VAE's encoder
- Use TTS decoder as VAE's decoder

6

# Objective function of conditional VQ-VAE

$$\log p(\mathbf{x}_{1:T} \mid \mathbf{y}_{1:U})$$

TTS as a probabilistic model

$$= \log \sum_{\forall \mathbf{z}_{1:U}^{(q)}} \int_{\mathbf{z}_{1:U}^{(r)}} p(\mathbf{x}_{1:T}, \mathbf{z}_{1:U}^{(q)}, \mathbf{z}_{1:U}^{(r)} \mid \mathbf{y}_{1:U}) d\mathbf{z}_{1:U}^{(r)} \quad (7.1)$$

$$\geq \mathbb{E}_{Q_\lambda(\mathbf{z}_{1:U}^{(q)})} [\underbrace{\log p_\theta(\mathbf{x}_{1:T} \mid \mathbf{z}_{1:U}^{(q)}, \mathbf{y}_{1:U})}_{\text{Decoder}}] \quad (7.2)$$

TTS decoder

$$- \mathrm{KL}[Q_\lambda(\mathbf{z}_{1:U}^{(q)} \mid \mathbf{x}_{1:T}, \mathbf{y}_{1:U}) \| \underbrace{P_\phi(\mathbf{z}_{1:U}^{(q)} \mid \mathbf{y}_{1:U})}_{\text{Prior}}] \quad (7.3)$$

Duration predictor (prior)

Duration recognizer (vector quantization)

$$- \mathbb{E}_{Q_\lambda(\mathbf{z}_{1:U}^{(q)})} \left\{ \mathrm{KL}[Q_\psi(\mathbf{z}_{1:U}^{(r)} \mid \mathbf{z}_{1:U}^{(q)}, \mathbf{x}_{1:T}, \mathbf{y}_{1:U}) \| \underbrace{p(\mathbf{z}_{1:U}^{(r)} \mid \mathbf{z}_{1:U}^{(q)})}_{} ] \right\}. \quad (7.4)$$
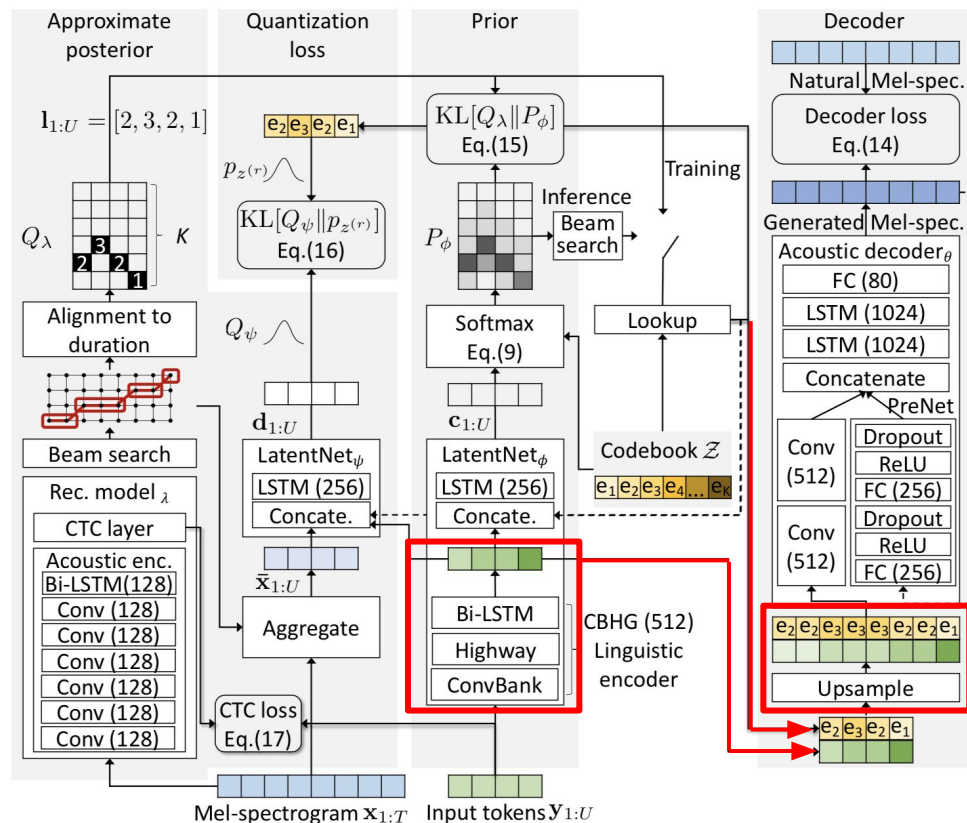
$$+ \gamma \log \sum_{\forall \mathbf{a}_{1:T}} P_\lambda(\mathbf{y}_{1:U}, \mathbf{a}_{1:T} \mid \mathbf{x}_{1:T}), \quad (7.17)$$

Forced aligner

Duration samples from forced aligner (approximate posterior)

- TTS can be model with the conditional probability (x: speech, y: texts).
- The marginal probability can be approximated with ELBO.
- Modules in TTS can be incorporated to VQ-VAE:
  - Prior: duration predictor
  - VQ: duration recognizer
  - Approximate posterior: duration samples
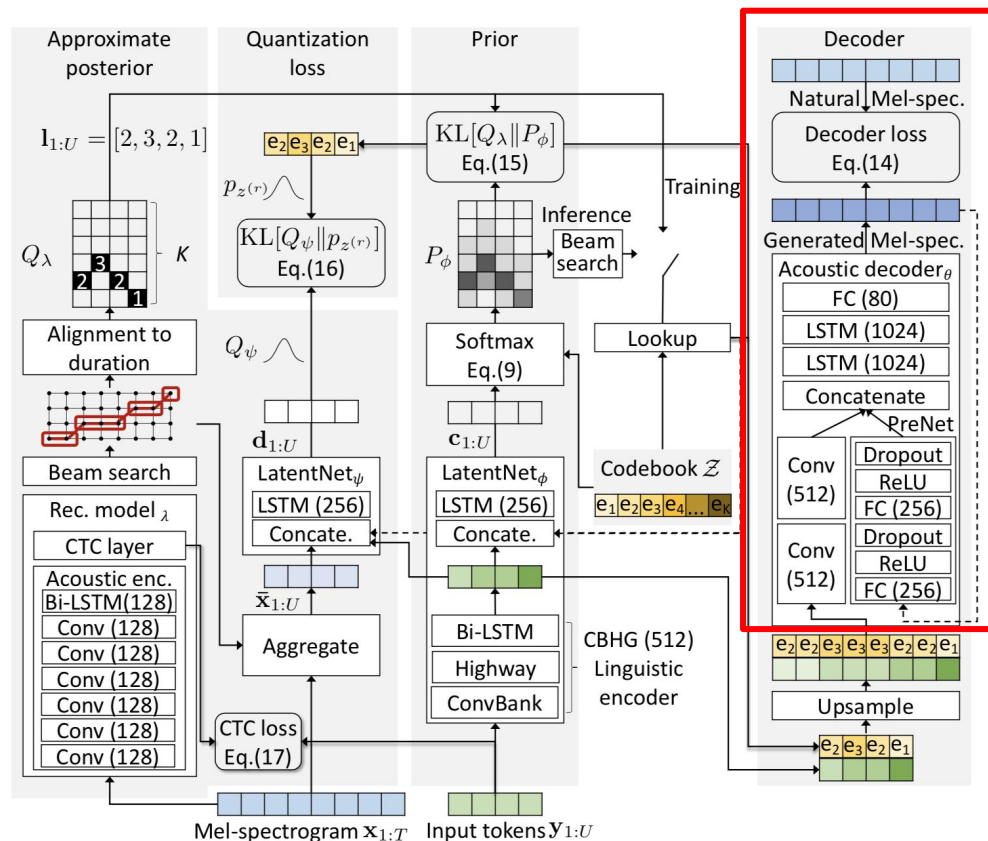  - Decoder: TTS decoder
  - Sampler: forced aligner

7

# An architecture of the proposed method (1)
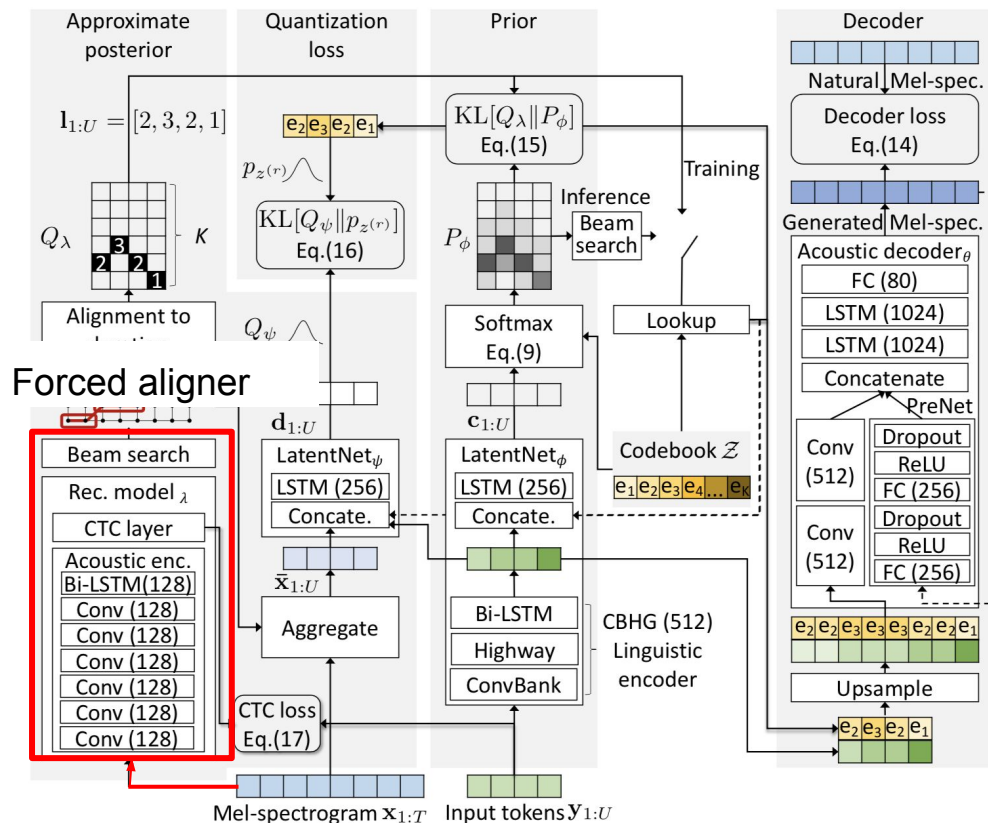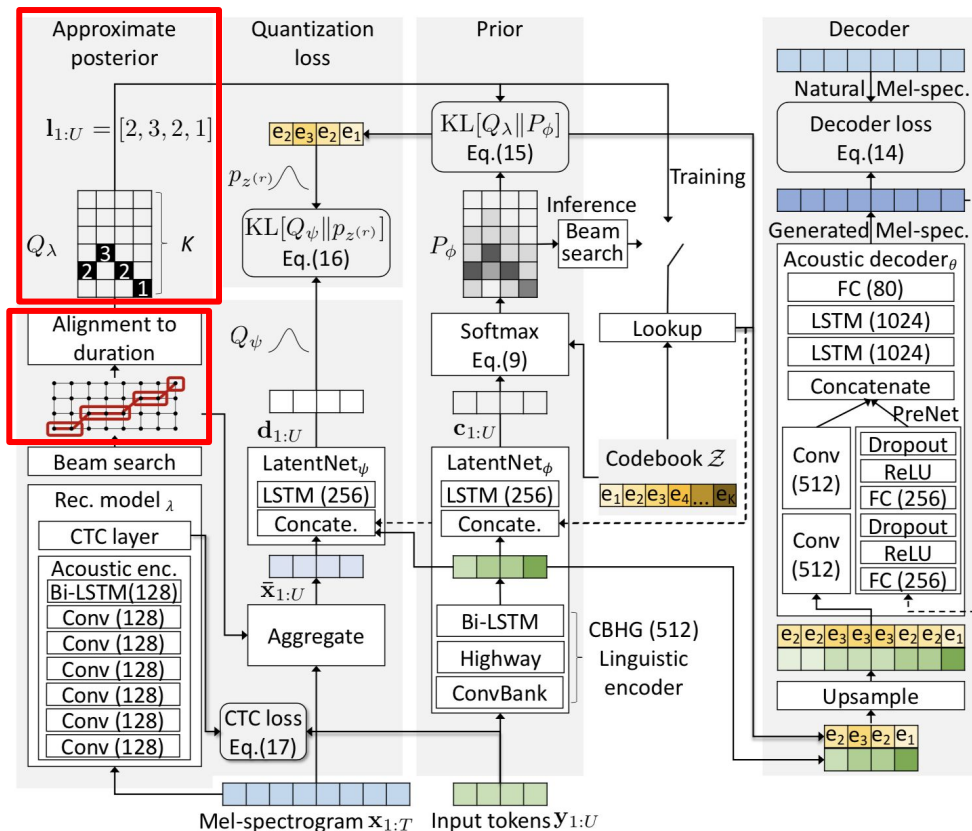


- Common behaviors
  a. Upsample linguistic features based on phoneme duration
  b. Decoder decodes the upsampled linguistic features into acoustic features.

# An architecture of the proposed method (2)



- ● Common behaviors
  - a. Upsample linguistic features based on phoneme duration
  - b. Decoder decodes the upsampled linguistic features into acoustic features.

# An architecture of the proposed method (3)



- Training phase
  a. Sample alignments from forced aligener.
  b. Convert alignments into phoneme duration. The duration samples defines approximate posterior.
  c. Duration recognizer predicts distributions by encoding linguistic and acoustic features. Duration codebook is optimized by minimizing KLD between the distribution and codebook.
  d. The samples are used to train duration recognizer and predictor.
  e. Duration predictor predicts distribution about duration from linguistic features. Minimizing KLD between the distribution and the approximate posterior optimizes the duration predictor.
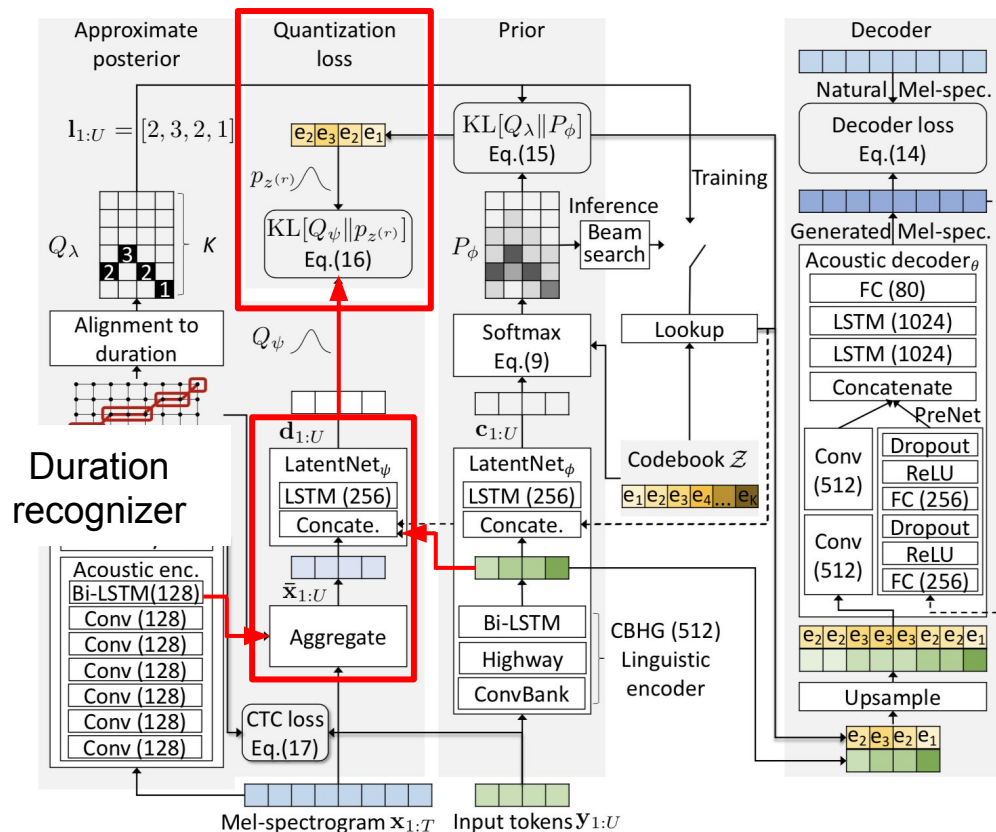
# An architecture of the proposed method (4)
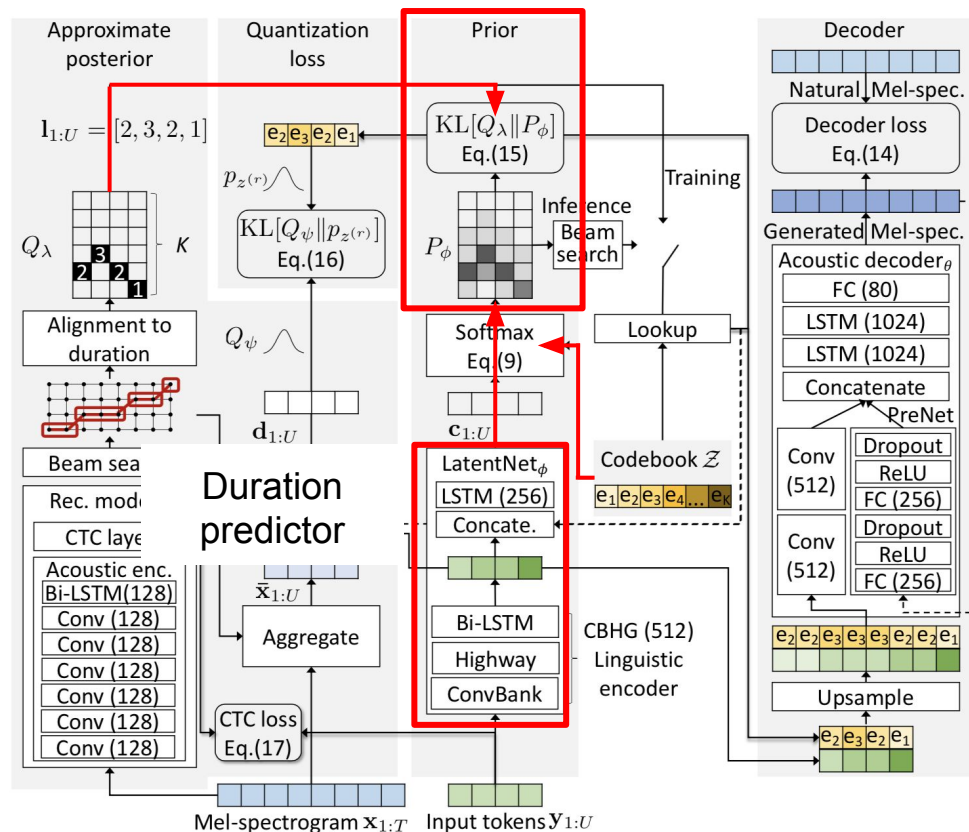


- ● Training phase
  a. Sample alignments from forced aligner.
  b. Convert alignments into phoneme duration. The duration samples defines approximate posterior.
  c. Duration recognizer predicts distributions by encoding linguistic and acoustic features. Duration codebook is optimized by minimizing KLD between the distribution and codebook.
  d. The samples are used to train duration recognizer and predictor.
  e. Duration predictor predicts distribution about duration from linguistic features. Minimizing KLD between the distribution and the approximate posterior optimizes the duration predictor.

# An architecture of the proposed method (5)



- ● Training phase
  a. Sample alignments from forced aligner.
  b. Convert alignments into phoneme duration. The duration samples defines approximate posterior.
  c. Duration recognizer predicts distributions by encoding linguistic and acoustic features. Duration codebook is optimized by minimizing KLD between the distribution and codebook.
  d. The samples are used to train duration recognizer and predictor.
  e. Duration predictor predicts distribution about duration from linguistic features. Minimizing
  f. KLD between the distribution and the approximate posterior optimizes the duration predictor.
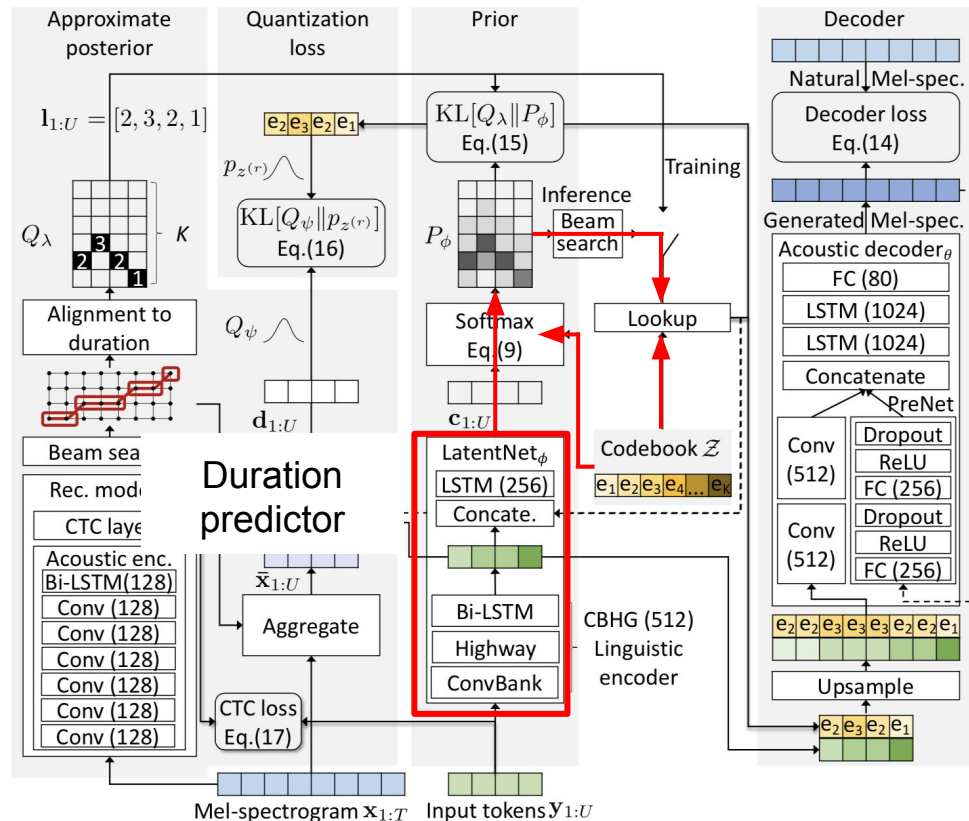
# An architecture of the proposed method (6)



- ● Training phase
  - a. Sample alignments from forced aligener.
  - b. Convert alignments into phoneme duration. The duration samples defines approximate posterior.
  - c. Duration recognizer predicts distributions by encoding linguistic and acoustic features. Duration codebook is optimized by minimizing KLD between the distribution and codebook.
  - d. Duration predictor predicts distribution about duration from linguistic features. Minimizing KLD between the distribution and the approximate posterior optimizes the duration predictor.
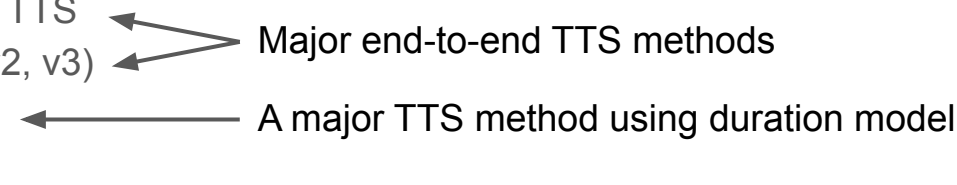
# An architecture of the proposed method (7)



- Prediction phase
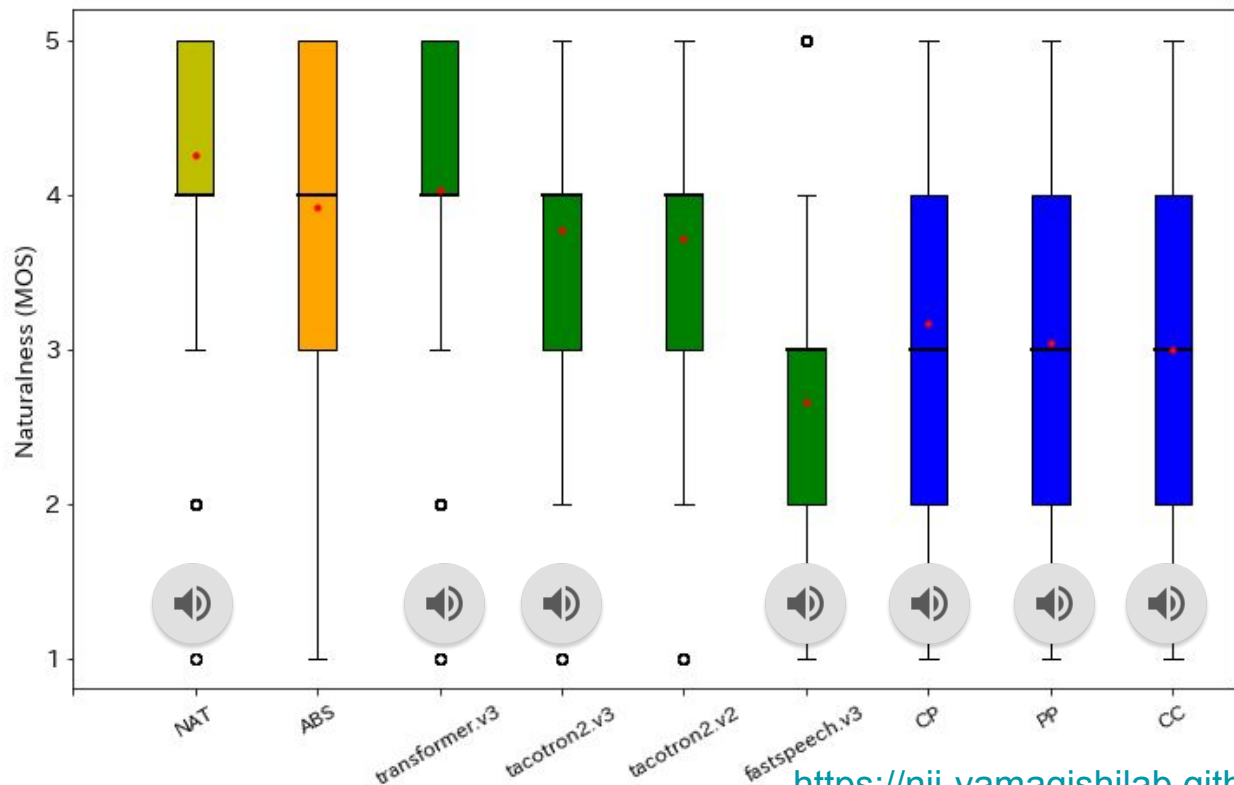  a. Sample phoneme duration from duration predictor with beam search.

# Experiment

- Corpus: LJSpeech (English, 1 female speaker, 24h)
- Proposed systems
  - CP (Character inputs for TTS, phoneme inputs for forced aligner)
  - PP (Phoneme inputs)
  - CC (Character inputs)
- Baselines:
  - Transformer TTS
  - Tacotron2 (v2, v3)    Major end-to-end TTS methods
  - FastSpeech    A major TTS method using duration model
  - ABS
  - Natural
- Evaluation
  - Listening test about naturalness (5-grade MOS)
  - 200 listeners

# Experimental results



Naturalness (MOS) box plot for NAT, ABS, transformer.v3, tacotron2.v3, tacotron2.v2, fastspeech.v3, CP, PP, CC

- Naturalness:
  - Natural speech
  - \> End-to-end TTS
    - Transformer
    - Tacotron v3, v2
  - \> Proposed systems
    - CP, PP, CC
  - \> TTS using duration model
    - FastSpeech
- Proposed systems
  - CP (character for TTS, phoneme for CTC) was the best

# Pros & Cons of the proposed method

- Pros
  - Single training phase.
  - Simple training criterion.

| Method | Teacher-Student | Training phases | Aligner | Aligner is external | Duration form | Latent duration |
|---|---|---|---|---|---|---|
| FastSpeech | ✓ | 3 | Soft-attention | ✓ | continuous | |
| DurIAn | | 3 | HMM? | ✓ | continuous | |
| FastSpeech2 | | 3 | HMM-GMM | ✓ | continuous | |
| AlignTTS | | 4 | SSNT-like MDN | ✓ | discrete | |
| JDI-T | ✓ | 1 | Soft-attention + CTC | | continuous | |
| Glow-TTS | | 1 | MAS | | continuous | |
| Non-Attentive Tacotron | | 2 | HMM | ✓ | continuous | ? |
| VQ-VAE | | 1 | CTC | | discrete | ✓ |

- Cons
  - Sensitive to design of linguistic feature labels
    - Absence of pauses
    - Symbols which is not straightforwardly related to duration
  - Duration sampled from forced aligner (CTC) is not good enough.
    - CTC assumes conditional independence across time steps.
    - It may not be suitable for segmentation.