

VAE-based SSE

Phenomena Label

Mel

Phenomena
Encoder

Reference
Encoder

$L2\ Loss$

Phone
Average

Phenomena
Predictor

Projection

μ_ϕ σ_ϕ

z

Mask

Flow-based SSP

$f_\theta(z_s)$

Flow

$KL\ Loss$

μ_θ , σ_θ

Prosody
Predictor

FastSpeech

BN

BN
Decoder

LR

Duration
Predictor

Text
Encoder

Text

----- Training Only
..... Inference Only

//
Stop
Gradient