

# HYPERTTS: Parameter Efficient Adaptation in Text to Speech using Hypernetworks

Yignting Li\*, Rishabh Bhardwaj\*, Ambuj Mehrish\*, Bo Cheng, Soujanya Poria

Beijing University of Posts and Telecommunications, China

Singapore University of Technology and Design, Singapore



# Introduction

## Motivations:

- Fine-tuning adaptation is resource-intensive, so use parameter-efficient domain adapters can make the adaptation scalable.
- Static adapter parameters can't perform well across multiple speakers of the adaptation domain due to underparameterization.

# Introduction

## Contributions:

- **Dynamic Adapters**  
Use HYPERTTS learns speaker-adaptative adapters conditioned on speaker embeddings.
- **Parameter Sampling**  
Employ parameter sampling from a continuous distribution defined by a learnable hypernetwork.
- **Parameter Efficiency**  
Achieves competitive results with less than 1% of fine-tuning parameters, making it highly practical and resource-friendly for scalable applications.

# Methodology

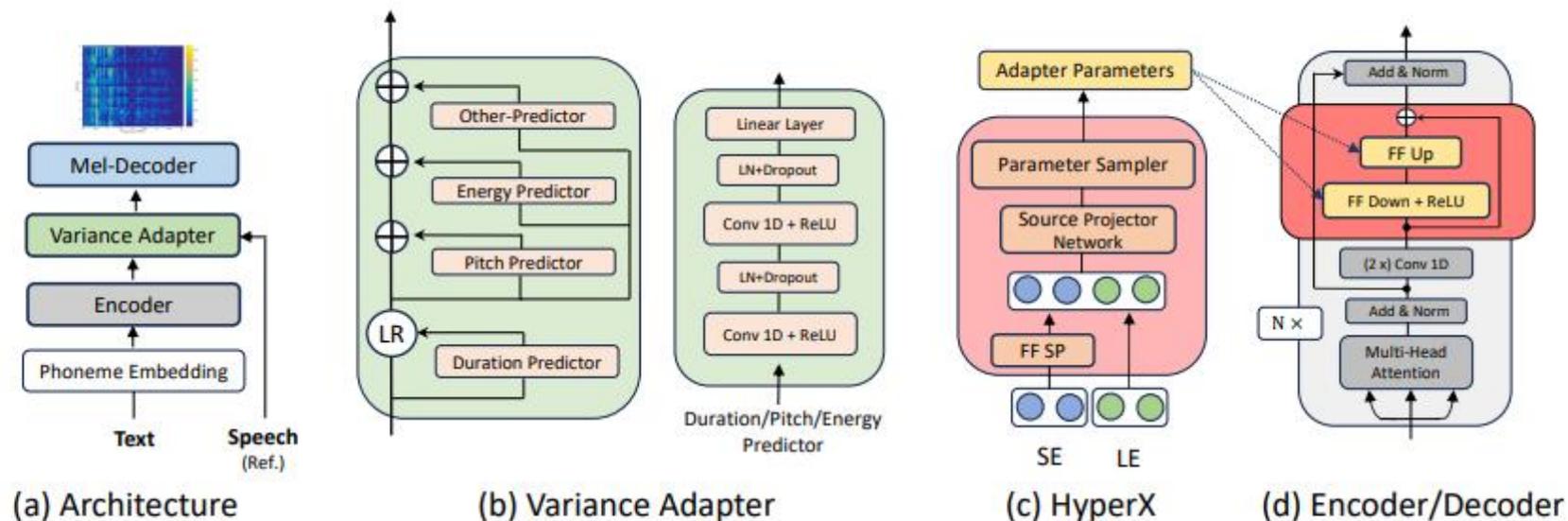


Figure 2: An overview of the HYPERTTS. SE and LE denote speaker embedding and layer embedding.

The TTS backbone architecture comprises a text encoder, Variance Adapter and Mel-Decoder.

# Methodology

## Encoder

- Input: A phoneme sequence  $(p_1, \dots, p_n)$  obtained from text.
- Architecture: Four Feed-Forward Transformer(FFT) blocks, with each block comprised of two multi-head attention modules and two 1D-convolutions.
- output: phoneme embeddings.

## Mel-Decoder and Postnet

- Input: Variance adaptor's hidden sequence.
- Architecture: Six Feed-Forward Transformer(FFT) blocks, same as Encoder FFT block.
- Mel-Decoder Output: Mel-spectrogram.
- Postnet Output: Mel-spectrogram reduced artifacts and distortions.

# Methodology

## Variance Adapter

Transforms the phoneme embeddings (length  $n$ ) into mel-spectrogram embeddings (length  $m$ ) where  $m \gg n$ .

- Duration Predictor

Alignment between phoneme sequences and mel-frames using Viterbi algorithm.

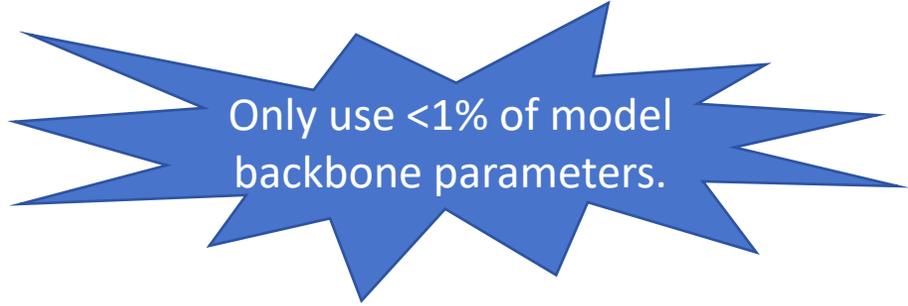
- Pitch Predictor

Predict pitch contour for each element of the length-regulated phoneme sequence using continuous wavelet transform(CWT) .

- Energy Predictor

Predict energy values for each Short-Time Fourier Transform(STFT) frame.

# Methodology



Only use <1% of model backbone parameters.

## Hypernetwork

$$h = h + \text{ReLU}(hW_d)W_u \quad (1)$$

- $W_d$  is down-projection and  $W_u$  is up-projection.

$$W_d = f_d(v_s, v_l) \quad (2)$$

$$W_u = f_u(v_s, v_l) \quad (3)$$

- $f_d$  and  $f_u$  denote parameter generators,  $v_s$  denote speaker embedding and  $v_l$  denote layer embedding.

# Experiments

## Baselines

- TTS-0 --- lower-bound
  - Represents the zero-shot performance.
- Reference and Rerference (Voc.)
  - Reconstruct the speech singnals using HiFi-GAN from ground-truth mel-spectrograms.
- TTS-FT (Full fine-tuning) ---- upper-bound
  - The performance after fine-tune all parameters of the backbone model on target dataset.
- AdapterTTS
  - The performance only fine-tune bottleneck adapter parameters on target dataset.
- HYPERTTS
  - The performance only fine-tune hypernetwork parameters on target dataset.

# Experiments

## Baselines explanation

- AdapterTTS<sub>e</sub>, AdapterTTS<sub>v</sub>, AdapterTTS<sub>d</sub> mean bottleneck adapter block inserted in the encoder, VA, and decoder, respectively.
- AdapterTTS<sub>e/v/d</sub> is a combination.
- HYPERTTS<sub>e</sub>, HYPERTTS<sub>v</sub>, HYPERTTS<sub>d</sub> mean bottleneck adapter block inserted in the encoder, VA, and decoder, respectively.
- HYPERTTS<sub>e/v/d</sub> is a combination.

## Datasets

- Pretrained Dataset: LibriTTS100
- Target Datasets: VCTK and LTS2( LibriTTS dev-clean and test-clean)

# Evaluation Metrics

## Object Metrics

- Cosine Similarity(COS)
- F0 Frame Error(FFE)
- Mel cepstral distortion(MCD)
- Word Error Rate(WER)

## Subjective Metrics

- Naturalness(MOS)
- XAB test

# Main Results ---- Object

Model	LTS → VCTK				
	COS ↑	FFE ↓	WER ↓	MCD ↓	Params
Reference	100.000 <sub>(±0.000)</sub>	00.00 <sub>(±0.00)</sub>	0.2055	—	—
Reference (Voc.)	95.027 <sub>(±0.001)</sub>	22.10 <sub>(±0.03)</sub>	0.2074	—	—
TTS-0	73.794 <sub>(±0.004)</sub>	39.19 <sub>(±0.02)</sub>	0.2035	5.9232	—
TTS-FT	80.443 <sub>(±0.003)</sub>	34.63 <sub>(±0.02)</sub>	0.2027	5.2387	35.7M (100%)
AdapterTTS <sub>e</sub>	73.769 <sub>(±0.004)</sub>	38.73 <sub>(±0.02)</sub>	<b>0.2075</b>	5.9002	66.6K (0.186%)
AdapterTTS <sub>v</sub>	73.131 <sub>(±0.004)</sub>	42.87 <sub>(±0.02)</sub>	0.2258	6.2733	33.3K (0.095%)
AdapterTTS <sub>d</sub>	76.180 <sub>(±0.004)</sub>	39.14 <sub>(±0.03)</sub>	0.2101	6.0092	100K (0.280%)
AdapterTTS <sub>e/d</sub>	72.703 <sub>(±0.004)</sub>	37.99 <sub>(±0.02)</sub>	0.2141	5.8804	166.7K (0.466%)
AdapterTTS <sub>e/v/d</sub>	77.298 <sub>(±0.006)</sub>	35.53 <sub>(±0.03)</sub>	0.2234	<b>5.2971</b>	200K (0.559%)
HyperTTS <sub>e</sub>	75.432 <sub>(±0.004)</sub>	36.07 <sub>(±0.02)</sub>	0.2367	5.3930	151.1K (0.423%)
HyperTTS <sub>v</sub>	73.731 <sub>(±0.004)</sub>	38.47 <sub>(±0.02)</sub>	0.2367	5.9137	150.9K (0.422%)
HyperTTS <sub>d</sub>	77.590 <sub>(±0.004)</sub>	38.55 <sub>(±0.02)</sub>	0.2090	5.9641	151.2K (0.423%)
HyperTTS <sub>e/d</sub>	79.232 <sub>(±0.003)</sub>	35.02 <sub>(±0.02)</sub>	0.2168	5.3650	302.3K (0.846%)
HyperTTS <sub>e/v/d</sub>	<b>79.464</b> <sub>(±0.003)</sub>	<b>34.47</b> <sub>(±0.02)</sub>	0.2340	5.3293	453.3K (1.269%)

Table 1: Domain adaptation performance on VCTK. TTS-0 denotes the zero-shot performance of the backbone TTS model evaluated on the VCTK test set. TTS-FT is a fine-tuned backbone model on the VCTK train set and evaluated on its test set. Where subscript  $m \in \{e, v, d\}$  in  $\text{HYPER}\text{TTS}_m$  or  $\text{Adapter}\text{TTS}_m$  denotes hypernetwork-adapter or adapter inserted to the encoder, variance adapter, and decoder of the backbone model, respectively.

# Main Results ---- Object

LTS → LTS2					
Model	COS ↑	FFE ↓	WER ↓	MCD ↓	Params
Reference	100.000 <sub>(±0.000)</sub>	00.00 <sub>(±0.00)</sub>	0.2046	—	—
Reference (Voc.)	96.919 <sub>(±0.000)</sub>	19.72 <sub>(±0.02)</sub>	0.2089	—	—
TTS-0	78.784 <sub>(±0.005)</sub>	43.31 <sub>(±0.02)</sub>	0.2129	7.7039	—
TTS-FT	82.351 <sub>(±0.004)</sub>	41.26 <sub>(±0.02)</sub>	0.2135	7.5843	35.7M (100%)
AdapterTTS	77.989 <sub>(±0.005)</sub>	42.28 <sub>(±0.02)</sub>	<b>0.2143</b>	7.6581	66.6K (0.186%)
HYPERTTS <sub>e</sub>	78.302 <sub>(±0.006)</sub>	41.38 <sub>(±0.02)</sub>	0.2232	<b>7.4746</b>	151.1K (0.423%)
HYPERTTS <sub>v</sub>	78.853 <sub>(±0.005)</sub>	43.19 <sub>(±0.02)</sub>	0.2180	7.7055	150.9K (0.422%)
HYPERTTS <sub>d</sub>	81.021 <sub>(±0.005)</sub>	41.77 <sub>(±0.02)</sub>	0.2159	7.5942	151.2K (0.423%)
HYPERTTS <sub>e/d</sub>	81.360 <sub>(±0.005)</sub>	41.95 <sub>(±0.02)</sub>	0.2180	7.5865	302.3K (0.846%)
HYPERTTS <sub>e/v/d</sub>	<b>81.742</b> <sub>(±0.005)</sub>	<b>41.03</b> <sub>(±0.02)</sub>	0.2337	7.5876	453.3K (1.269%)

Table 2: Speaker adaptation on LTS2 is assessed. "TTS-0" represents the base TTS model's zero-shot performance on the LTS2 test set. "TTS-FT" is the fine-tuned base model on the LTS2 training set, evaluated on its test set. The subscript  $m \in e, v, d$  in HYPERTTS<sub>m</sub> or AdapterTTS<sub>m</sub> denotes the integration of hypernetwork-adapter or adapter in the encoder, variance adapter, or decoder of the base model, respectively.

# Main Results ---- Subjective

Model	MOS $\uparrow$	Params
Reference	4.62( $\pm 0.02$ )	-
TTS-FT	3.70( $\pm 0.16$ )	35.7M (100%)
AdapterTTS <sub>e</sub>	3.47( $\pm 0.07$ )	66.6K (0.186%)
HyperTTS <sub>e</sub>	3.43( $\pm 0.10$ )	151.1K (0.423%)
HyperTTS <sub>d</sub>	3.64( $\pm 0.06$ )	150.9K (0.422%)

Table 4: MOS comparison among Reference, TTS-FT, AdapterTTS<sub>e</sub>, HYPERTTS<sub>e</sub> and HYPERTTS<sub>d</sub> for samples randomly selected from VCTK validation set.

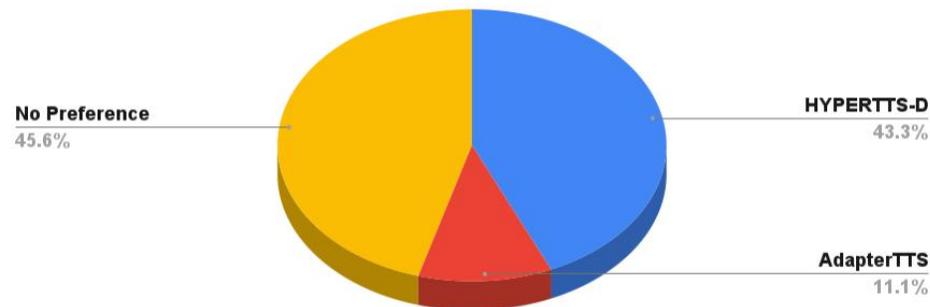


Figure 3: XAB speaker similarity test results between AdapterTTS, and *HyperTTS<sub>d</sub>*.

# Impact of Parameter Efficiency

Model	COS $\uparrow$	FFE $\downarrow$	WER $\downarrow$	MCD $\downarrow$	Params
HYPERTTS <sub>d</sub> (2)	75.89 <sub>(<math>\pm 0.0040</math>)</sub>	38.99 <sub>(<math>\pm 0.0211</math>)</sub>	0.2096	6.0016	50K (0.14%)
HYPERTTS <sub>d</sub> (8)	77.59 <sub>(<math>\pm 0.0036</math>)</sub>	38.55 <sub>(<math>\pm 0.0208</math>)</sub>	<b>0.2090</b>	<b>5.9641</b>	151K (0.42%)
HYPERTTS <sub>d</sub> (32)	79.38 <sub>(<math>\pm 0.0033</math>)</sub>	<b>38.05</b> <sub>(<math>\pm 0.0210</math>)</sub>	0.2152	6.0024	554K (1.55%)
HYPERTTS <sub>d</sub> (128)	<b>80.26</b> <sub>(<math>\pm 0.0033</math>)</sub>	38.15 <sub>(<math>\pm 0.0210</math>)</sub>	0.2186	5.9820	2.17M (6.06%)

Table 3: Varying number of parameters of hypernetwork in decoder with VCTK as the target domain. HYPERTTS<sub>d</sub>( $n$ ) denotes hypernetwork with  $n$ -dimensional source projection.

Model	COS $\uparrow$	FFE $\downarrow$	WER $\downarrow$	MCD $\downarrow$	Params
HYPERTTS <sub>e</sub> (2)	74.64 <sub>(<math>\pm 0.0038</math>)</sub>	37.27 <sub>(<math>\pm 0.0178</math>)</sub>	<b>0.2170</b>	5.4041	50K (0.14%)
HYPERTTS <sub>e</sub> (8)	75.43 <sub>(<math>\pm 0.0035</math>)</sub>	36.07 <sub>(<math>\pm 0.0193</math>)</sub>	0.2367	5.3930	151K (0.42%)
HYPERTTS <sub>e</sub> (32)	<b>76.02</b> <sub>(<math>\pm 0.0035</math>)</sub>	<b>35.17</b> <sub>(<math>\pm 0.0190</math>)</sub>	0.2449	<b>5.3779</b>	554K (1.5%)
HYPERTTS <sub>e</sub> (128)	75.97 <sub>(<math>\pm 0.0036</math>)</sub>	35.93 <sub>(<math>\pm 0.0193</math>)</sub>	0.2612	5.4210	2.17M (6.06%)

Table 5: Varying number of parameters of hypernetwork in encoder with VCTK as the target domain.

Model	COS $\uparrow$	FFE $\downarrow$	WER $\downarrow$	MCD $\downarrow$	Params
HYPERTTS <sub>e/v/d</sub> (2)	77.26 <sub>(<math>\pm 0.0034</math>)</sub>	35.10 <sub>(<math>\pm 0.0180</math>)</sub>	<b>0.2314</b>	5.3436	150K (0.42%)
HYPERTTS <sub>e/v/d</sub> (8)	79.46 <sub>(<math>\pm 0.0030</math>)</sub>	34.47 <sub>(<math>\pm 0.0198</math>)</sub>	0.2340	5.3293	453K (1.27%)
HYPERTTS <sub>e/v/d</sub> (32)	80.67 <sub>(<math>\pm 0.0027</math>)</sub>	33.85 <sub>(<math>\pm 0.0209</math>)</sub>	0.2513	<b>5.2761</b>	1.66M (4.65%)
HYPERTTS <sub>e/v/d</sub> (128)	<b>81.49</b> <sub>(<math>\pm 0.0031</math>)</sub>	<b>33.58</b> <sub>(<math>\pm 0.0223</math>)</sub>	0.2622	5.2879	6.50M (18.19%)

Table 6: Varying number of parameters of hypernetwork in the encoder, decoder, and VA with VCTK as the target domain.

# Conclusion

Dynamic adapter parameter, which enable input-conditioned parameter sampling. And Hypernetwork allows the generation of adapter parameters for numerous speakers!

HYPERTTS is an excellent choice for multi-speaker TTS adaptation, surpassing traditional adapter limitations.

Thanks for Listening!