



# Neural Speech Synthesis



Xu Tan and Tao Qin  
Microsoft Research Asia

Tutorial slides: <https://github.com/tts-tutorial/ijcai2021>

Survey paper: <https://arxiv.org/pdf/2106.15561>

# Outline

1. Evolution and taxonomy of TTS, Tao Qin, 20'
2. Key models in TTS, Xu Tan, 30'
3. Advanced topics in TTS, Xu Tan, 30'
4. Summary and future directions, Tao Qin, 5'
5. QA

# Part 1: Evolution and Taxonomy

-- Evolution, basic concepts, taxonomies



alexa



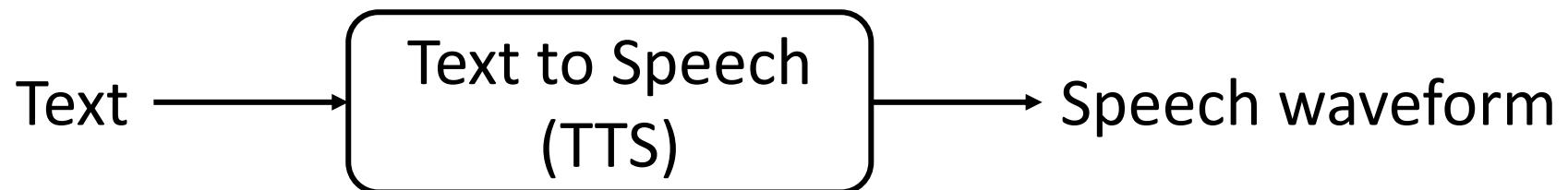
Hi, I'm Cortana.



Siri

# Text to speech synthesis

- The artificial production of human speech from text



- Disciplines: acoustics, linguistics, digital signal processing, statistics and deep learning
- The quality of the synthesized speech is measured by
  - Intelligibility and naturalness

# Formant TTS

## How does it work?

- produce speech segments by generating artificial signals based on a set of specified rules mimicking the formant structure and other spectral properties of natural speech
- using additive synthesis and an acoustic model (with parameters like voicing, fundamental frequency, noise levels)

## Advantages:

- highly intelligible, even at high speeds
- well-suited for embedded systems, with limited memory and computation power

## Limitations:

- not natural, produces artificial, robotic-sounding speech, far from human speech
- difficult to design rules that specify model parameters

# Concatenative TTS

How does it work?

- a very large database of short and high-quality speech fragments are recorded from a single speaker
- speech fragments are recombined to form complete utterances

Advantages: intelligible

Limitations:

- require huge databases and hard-coding the combination
- emotionless, not natural
- difficult to modify the voice (e.g., switching to a different speaker, or altering the emphasis or emotion) without recording a whole new database

# Parametric TTS

How does it work?

- using learning based parametric models, e.g., HMM
- all the information required to generate speech is stored in the parameters of the model

Advantages: lower data cost and more flexible

Limitations: less intelligible than concatenative TTS

# Neural TTS

How does it work?

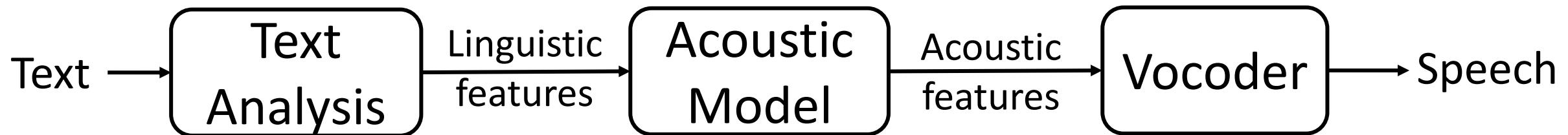
- a special kind of parametric models
- text to waveform mapping is modeled by (deep) neural networks
- Advantages:
  - huge quality improvement, in terms of both intelligibility and naturalness
  - less human preprocessing and feature engineering

# Examples

Concatenative	Parametric	Neural
		

# Basic components of parametric/neural TTS systems

- Text analysis, acoustic model, and vocoder



- Text analysis: text → linguistic features
- Acoustic model: linguistic features → acoustic features
- Vocoder: acoustic features → speech

# Text analysis

- Transforms input text into linguistic features, including
  - Text normalization
    - 1999 → nineteen ninety-nine, *Jan. 24<sup>th</sup>* → *January twenty-fourth*
  - Homograph disambiguation
    - Do you **live** (/l ih v/) near a zoo with **live** (/l ay v/) animals?
  - Phrase/word/syllable segmentation
    - synthesis → syn-the-sis
  - Part of speech (POS) tagging
    - Mary went to the store → noun, verb, prep, noun,
  - ToBI (Tones and Break Indices)
    - Mary went to the store ? → Mary' store' H%
  - Grapheme-to-phoneme conversion
    - *Speech* → s p i y ch

# Text analysis: linguistic features

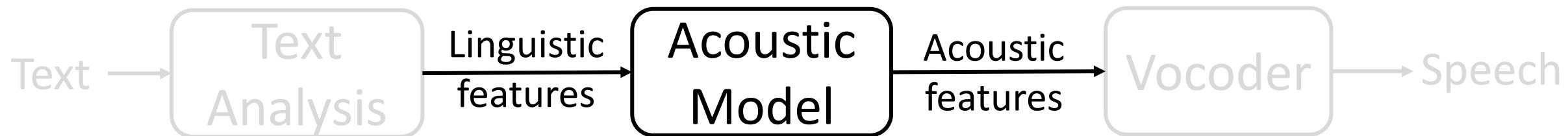
- Phoneme, syllable, word, phrase and sentence-level features, e.g.,
  - The phonetic symbols of the previous before the previous, the previous, the current, the next or the next after the next;
  - Whether the previous, the current or the next syllable is stressed;
  - The part of speech (POS) of the previous, the current or the next word;
  - The prosodic annotation of the current phrase;
  - The number of syllables, words or phrases in the current sentence.

# Text analysis: linguistic features

- phoneme:
  - current phoneme
  - preceding and succeeding two phonemes
  - position of current phoneme within current syllable
- syllable:
  - numbers of phonemes within preceding, current, and succeeding syllables
  - stress<sup>3</sup> and accent<sup>4</sup> of preceding, current, and succeeding syllables
  - positions of current syllable within current word and phrase
  - numbers of preceding and succeeding stressed syllables within current phrase
  - numbers of preceding and succeeding accented syllables within current phrase
  - number of syllables from previous stressed syllable
  - number of syllables to next stressed syllable
  - number of syllables from previous accented syllable
  - number of syllables to next accented syllable
  - vowel identity within current syllable
- word:
  - guess at part of speech of preceding, current, and succeeding words
  - numbers of syllables within preceding, current, and succeeding words
  - position of current word within current phrase
  - numbers of preceding and succeeding content words within current phrase
  - number of words from previous content word
  - number of words to next content word
- phrase:
  - numbers of syllables within preceding, current, and succeeding phrases
  - position of current phrase in major phrases
  - ToBI endtone of current phrase
- utterance:
  - numbers of syllables, words, and phrases in utterance<sup>14</sup>

# Acoustic model

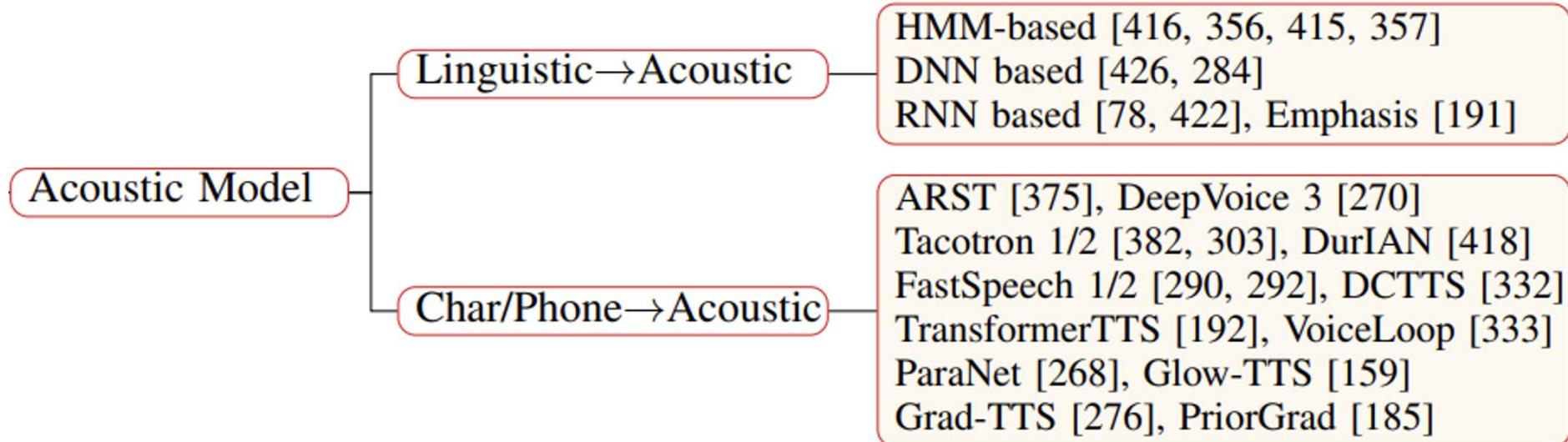
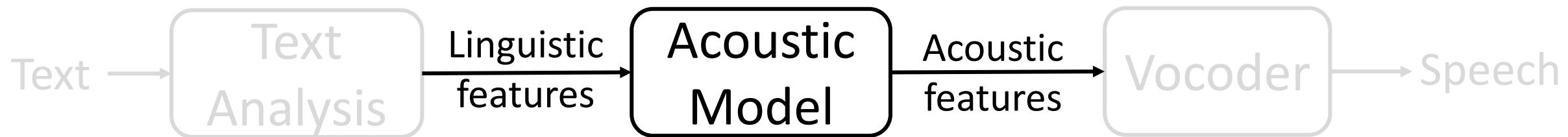
- Generate acoustic features from linguistic features



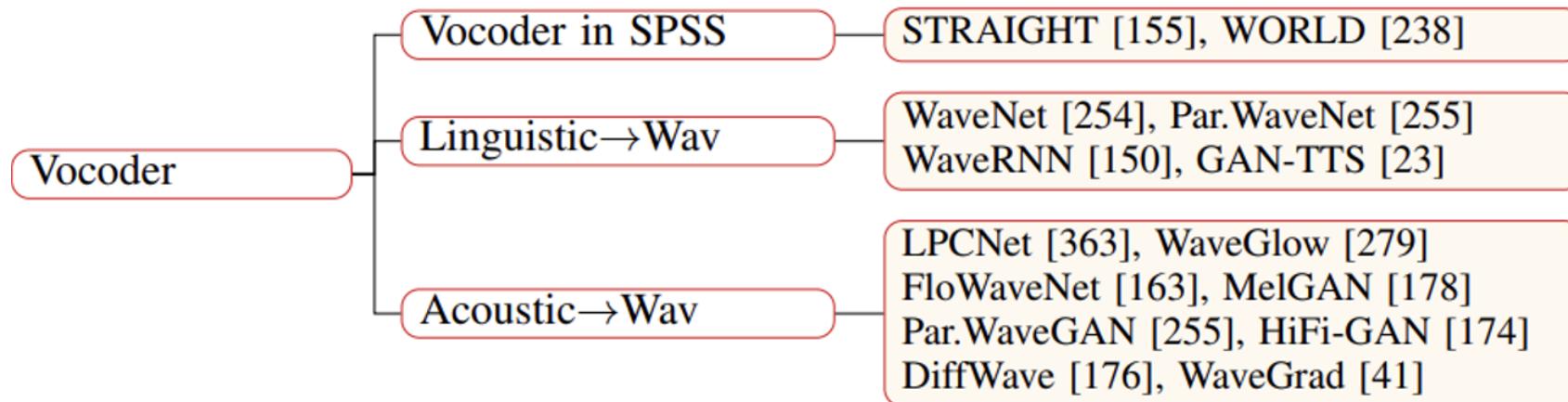
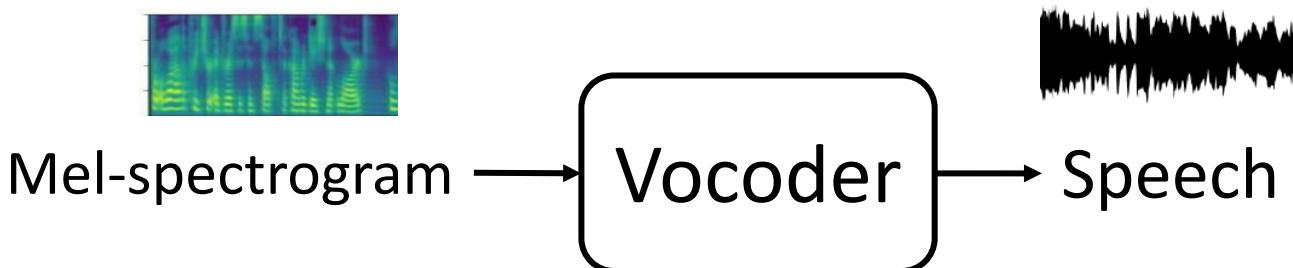
- F0, V/UV, energy
- Mel-scale Frequency Cepstral Coefficients (MFCC), Bark-Frequency Cepstral Coefficients (BFCC)
- Mel-generalized coefficients (MGC), band aperiodicity (BAP),
- Linear prediction coefficients (LPC),
- Mel-spectrograms
  - Pre-emphasis, Framing, Windowing, Short-Time Fourier Transform (STFT), Mel filter

# Acoustic model

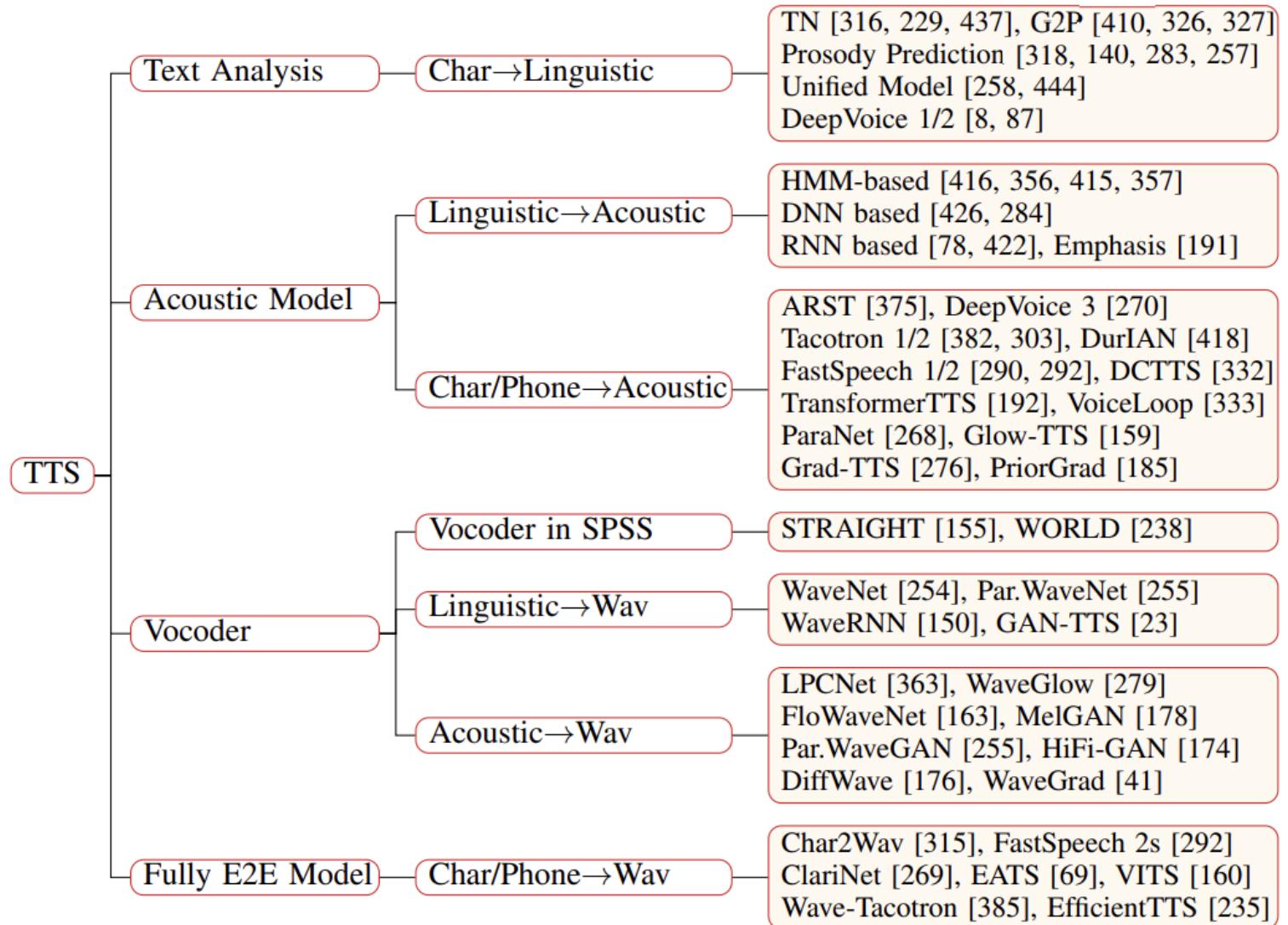
- Predict acoustic features from linguistic features



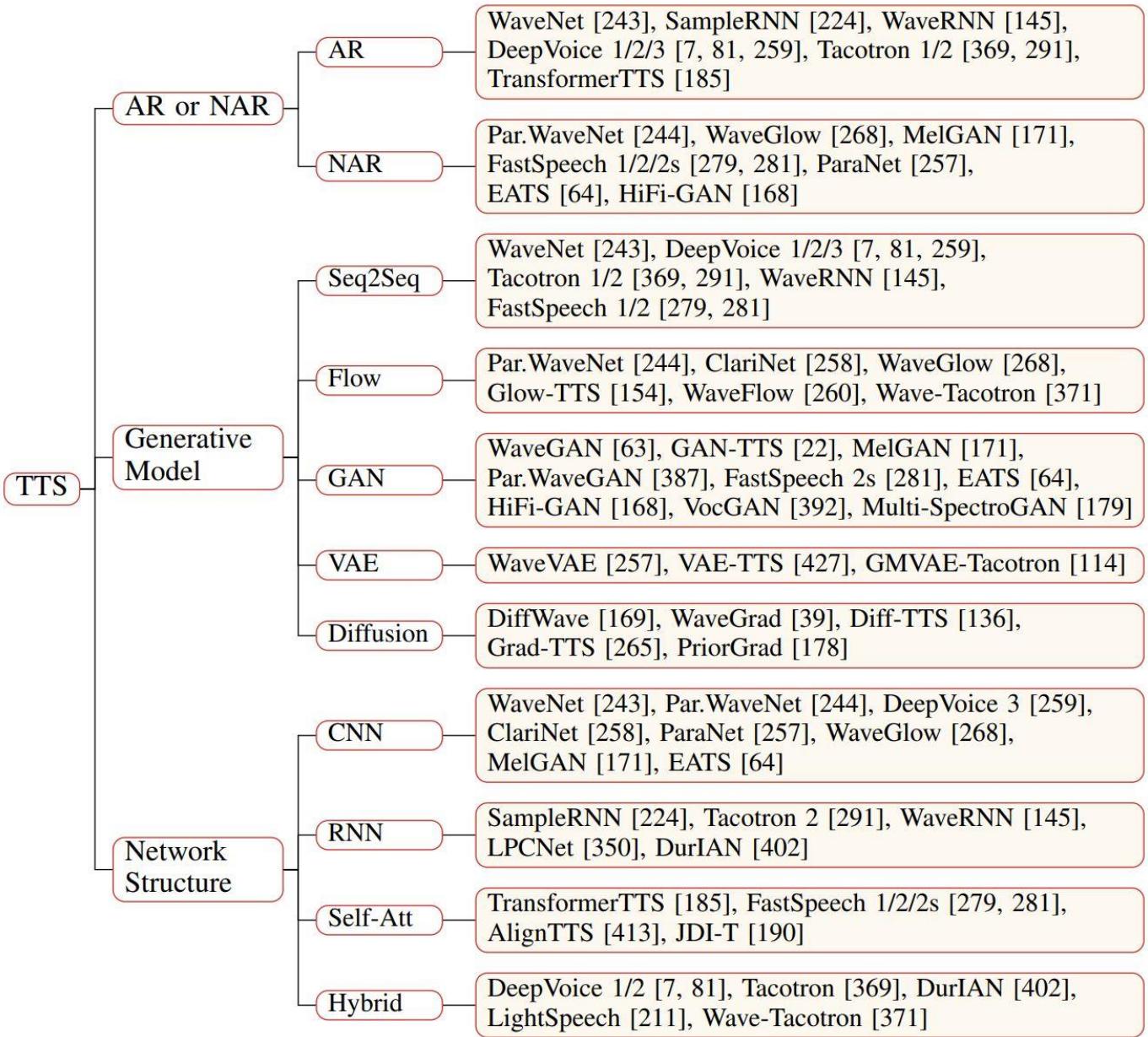
# Vocoder



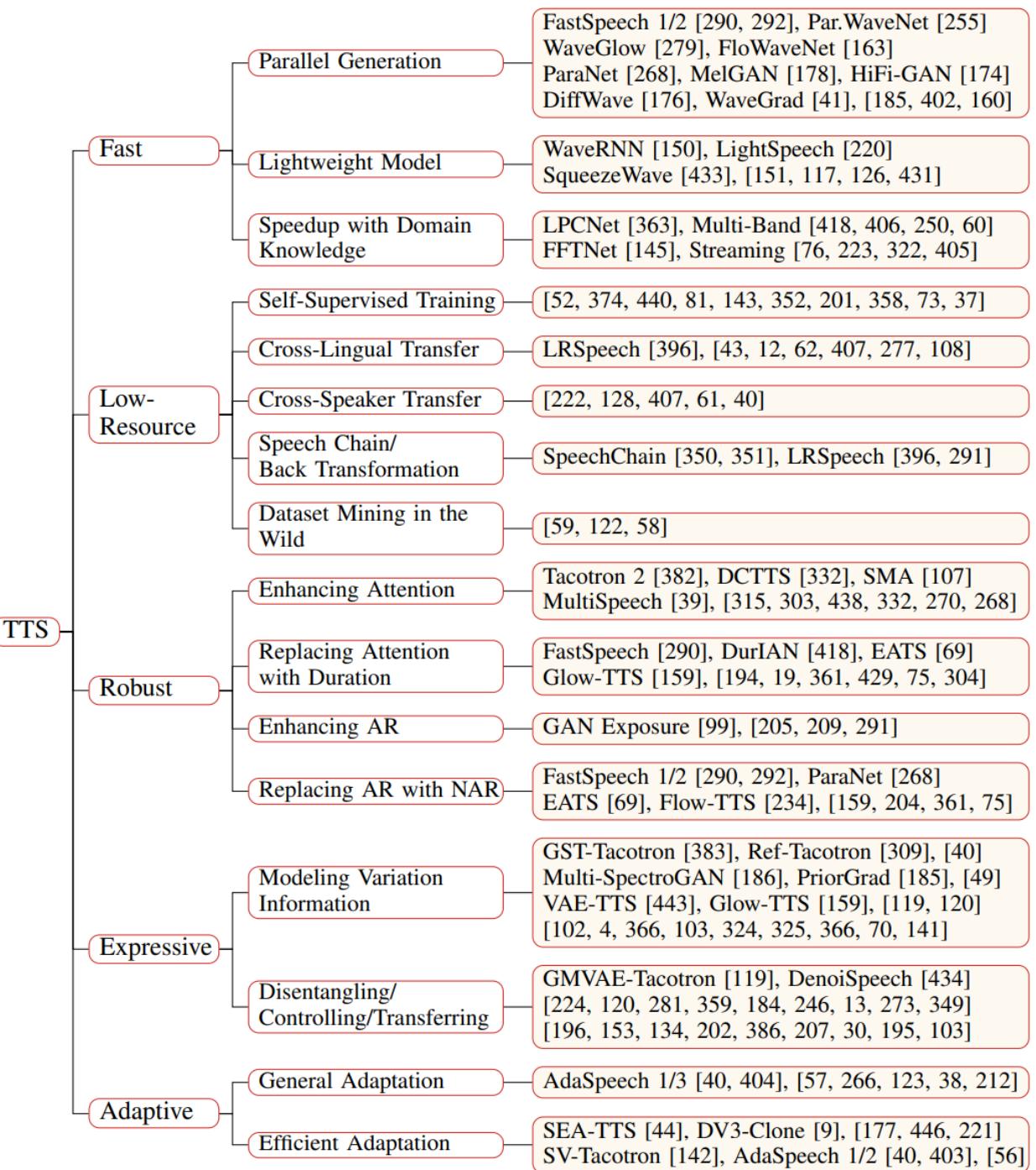
# Taxonomy from the perspective of components



# Taxonomy from the perspective of models



# Advanced topics



## Part 2: Key Components in TTS

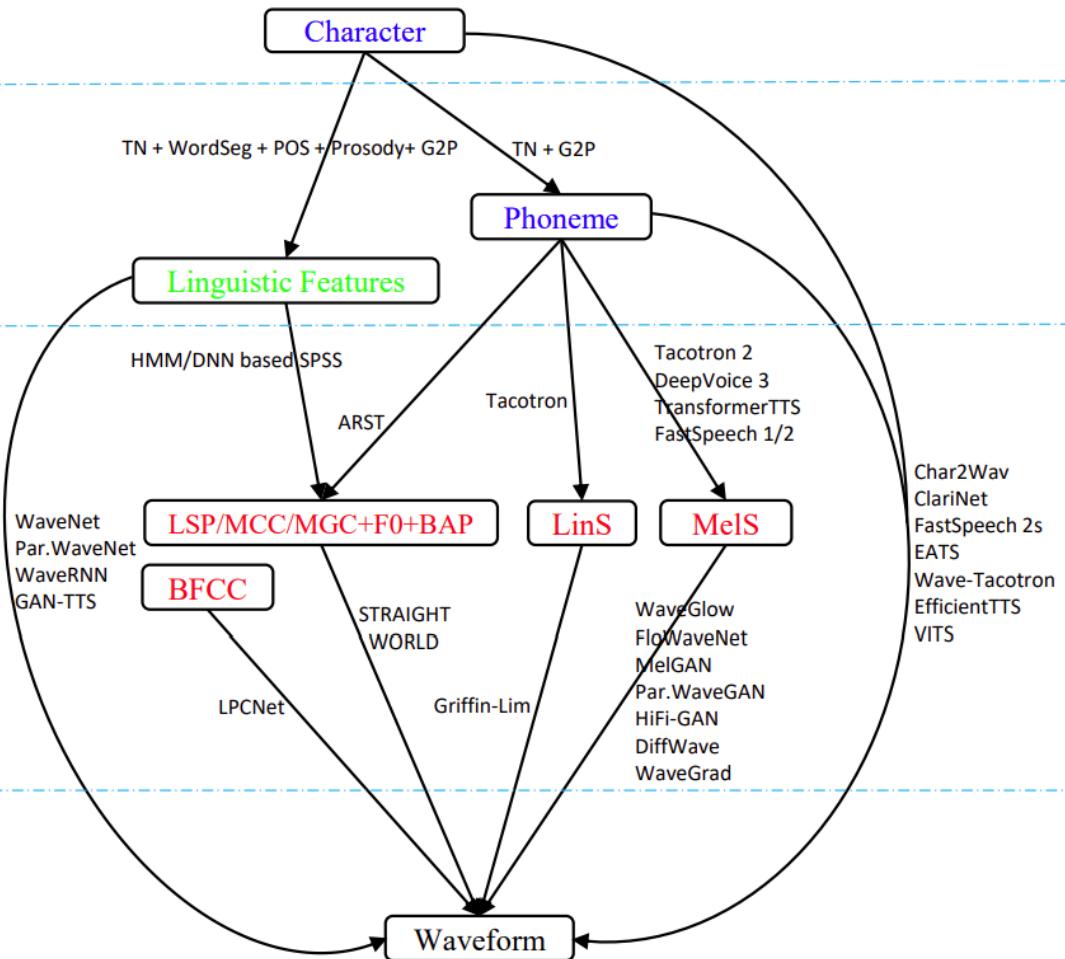
# Data conversion pipeline

Text

Linguistic Features

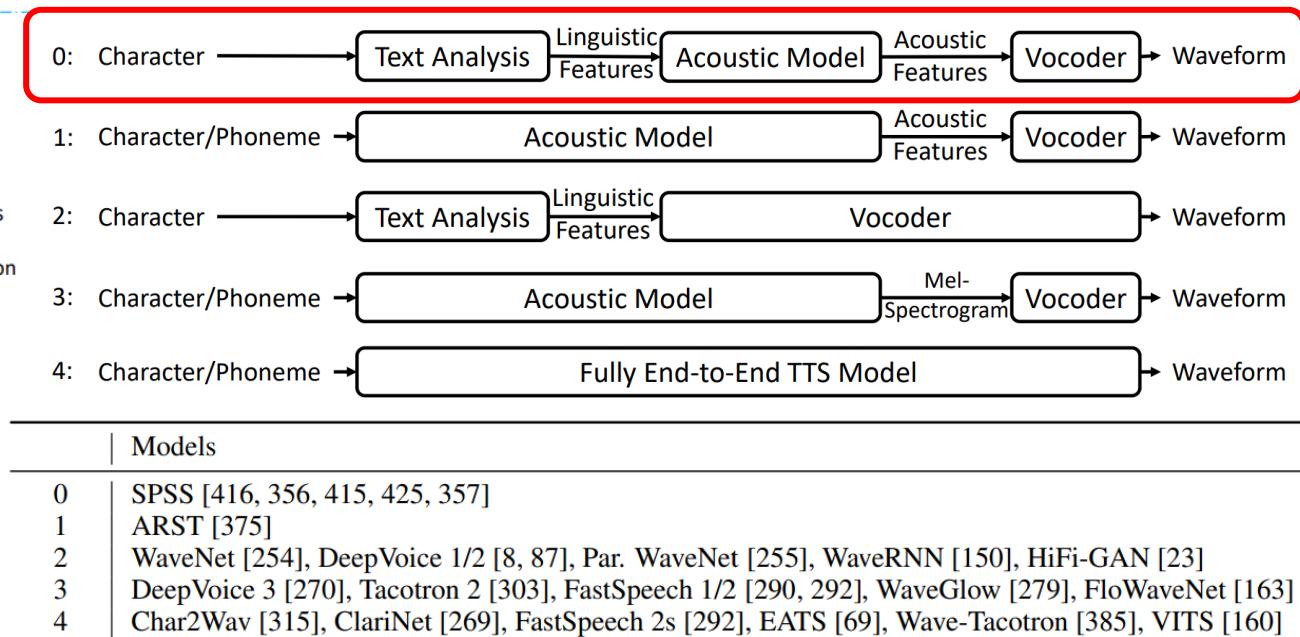
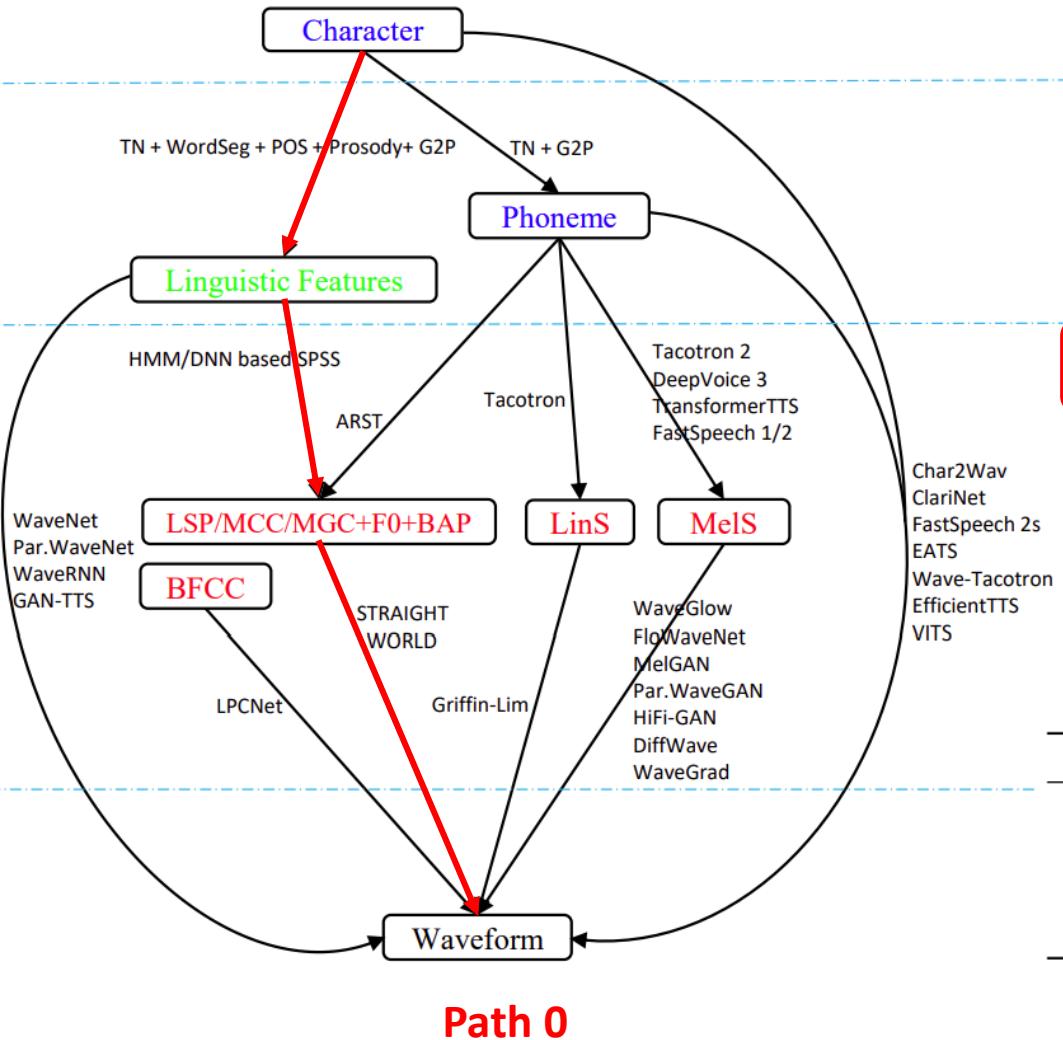
Acoustic Features

Waveform

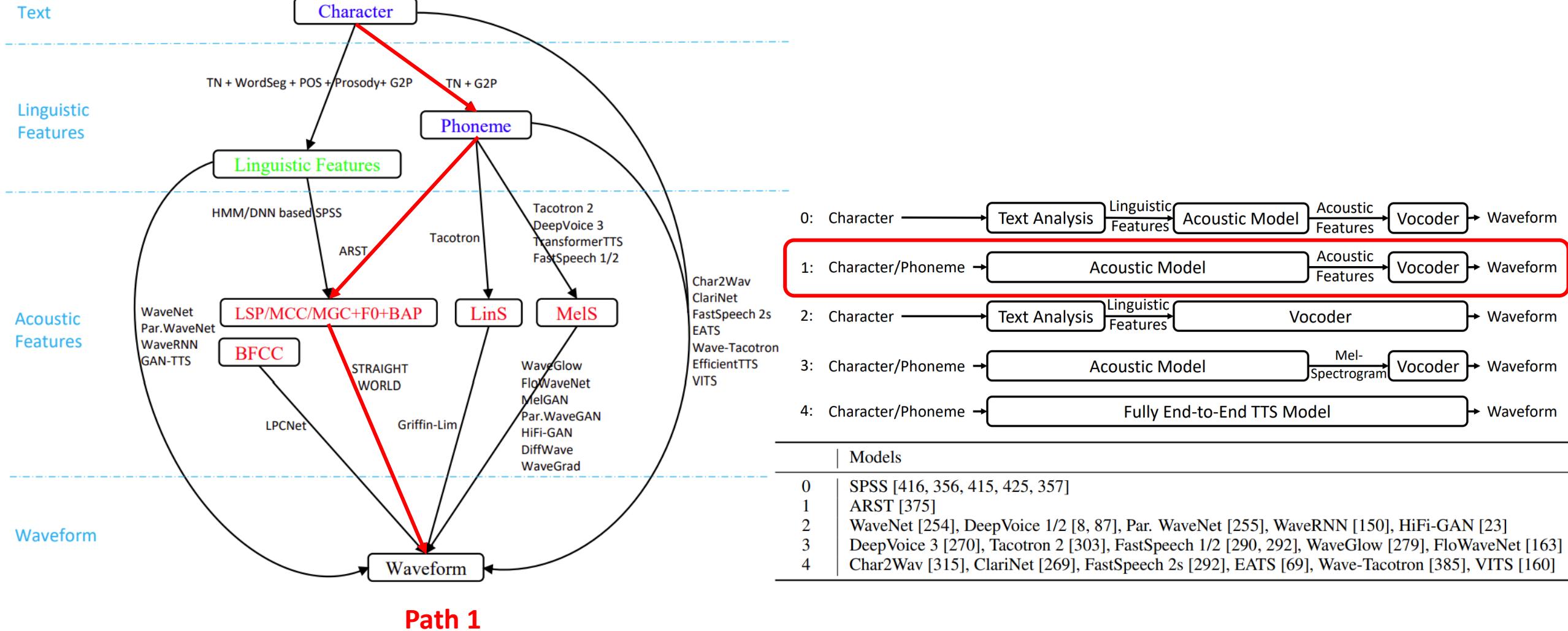


# Data conversion pipeline

Text

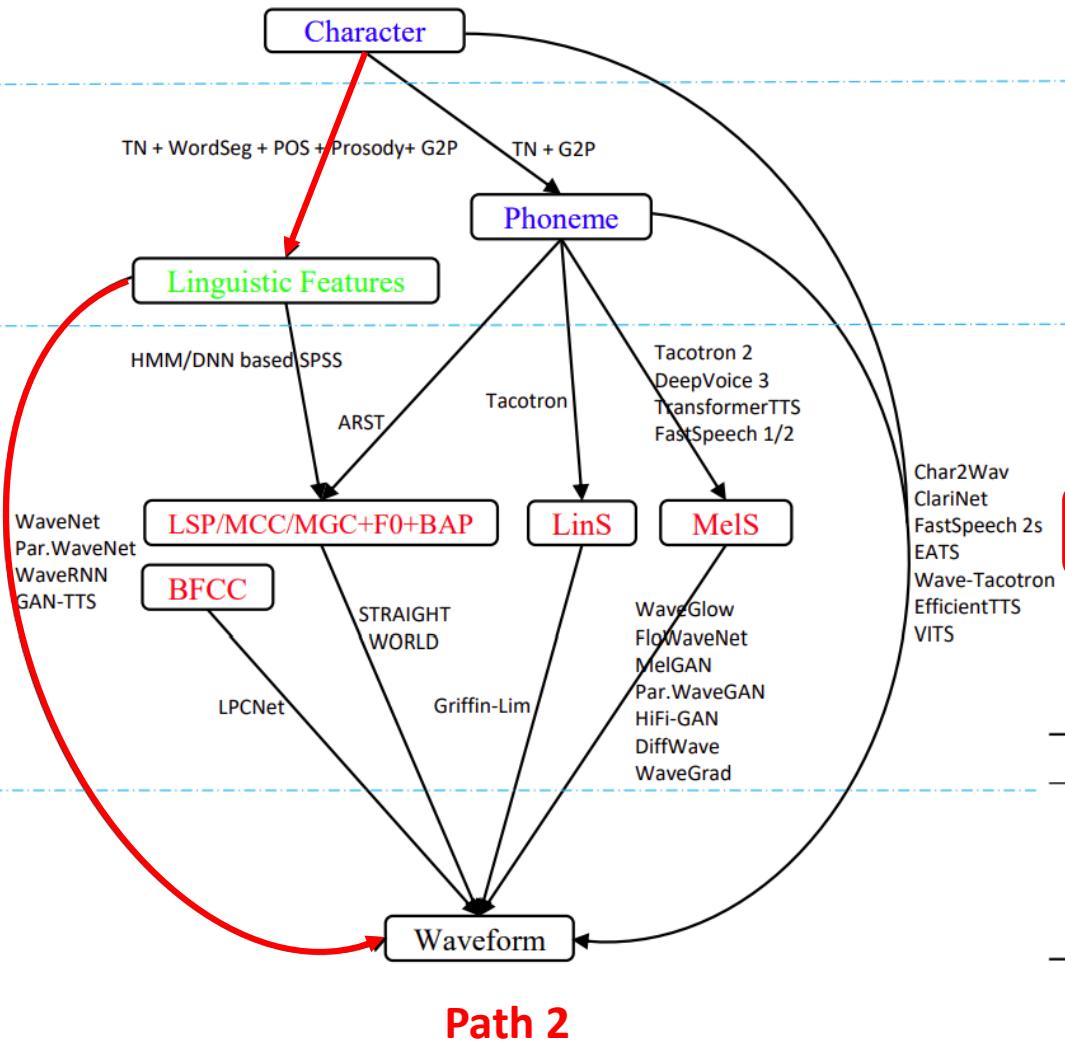


# Data conversion pipeline



# Data conversion pipeline

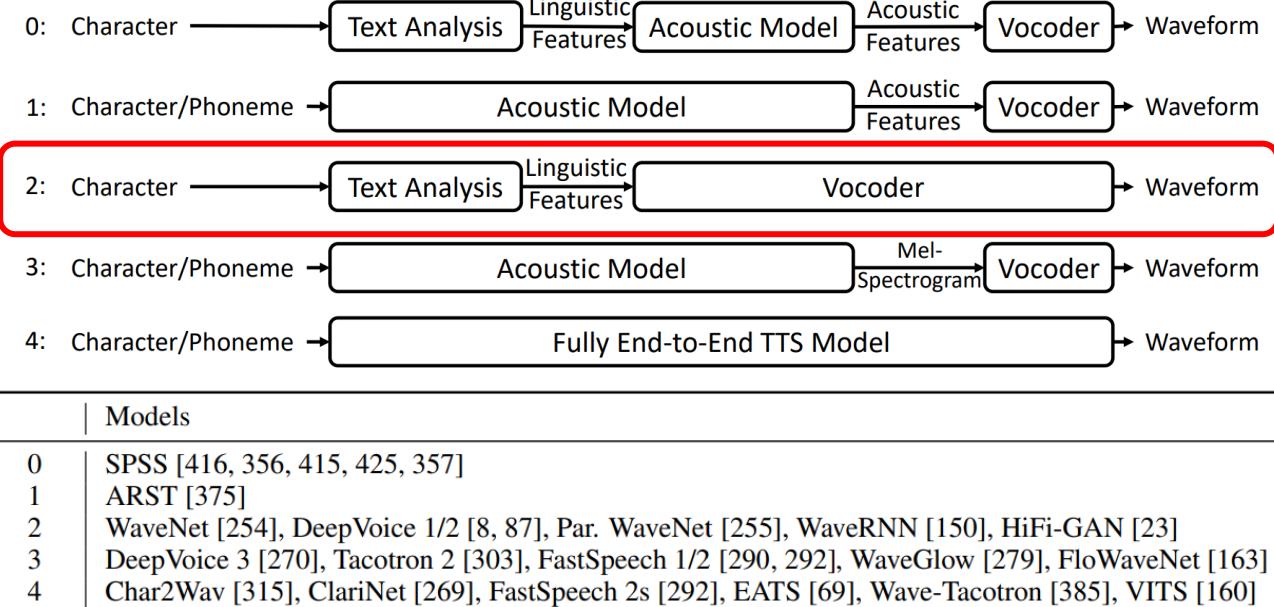
Text



Linguistic Features

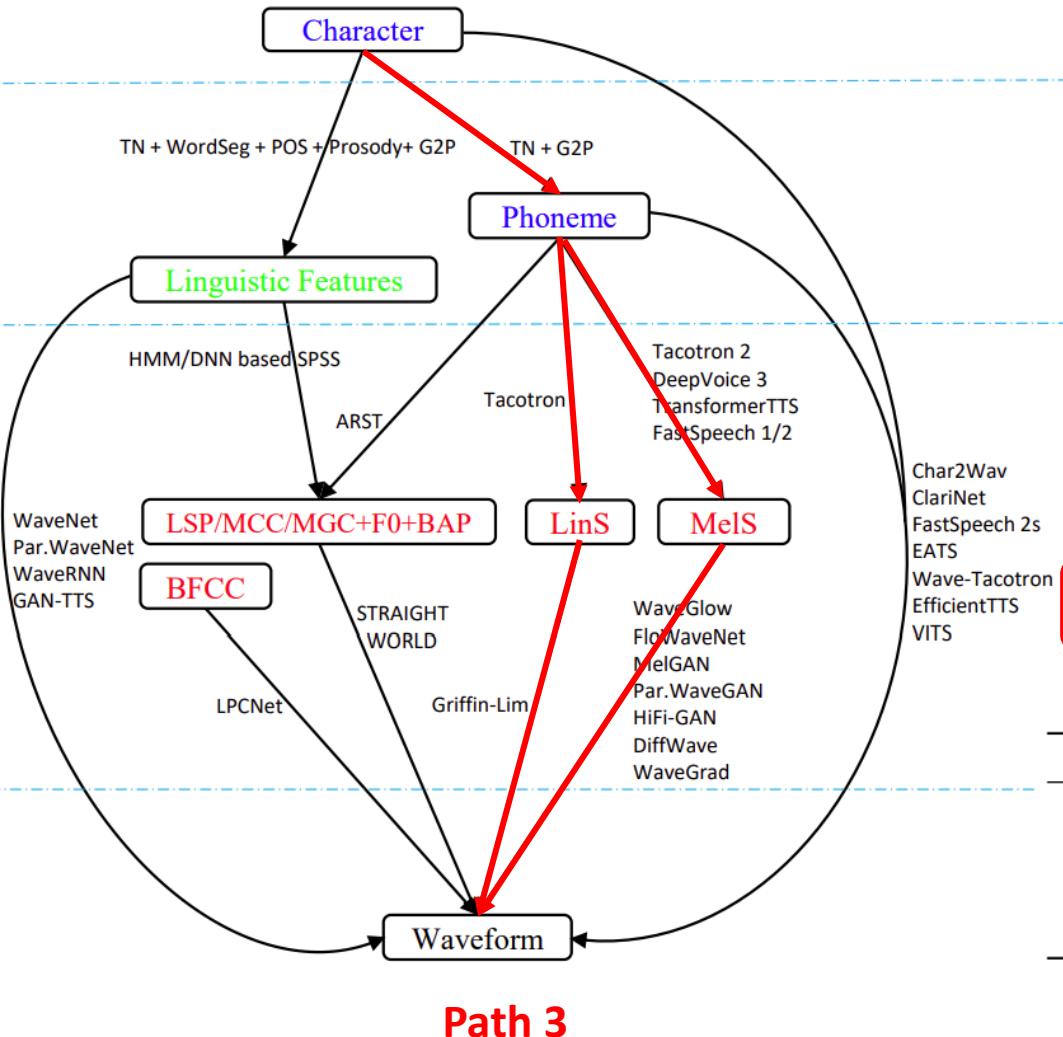
Acoustic Features

Waveform



# Data conversion pipeline

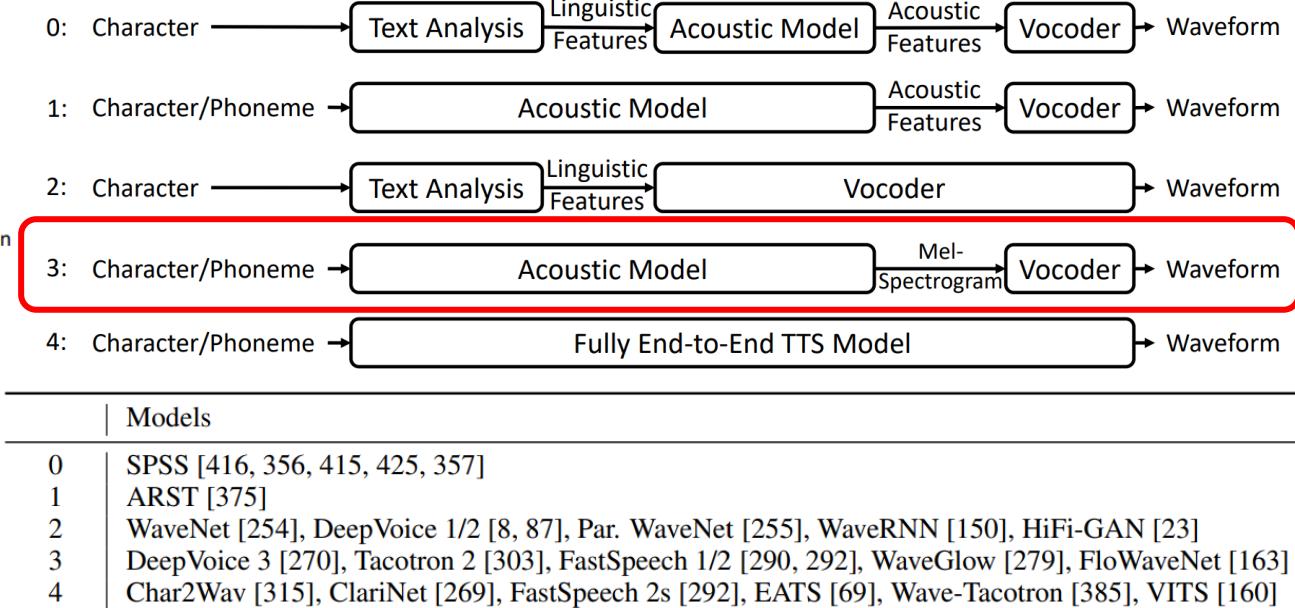
Text



Linguistic Features

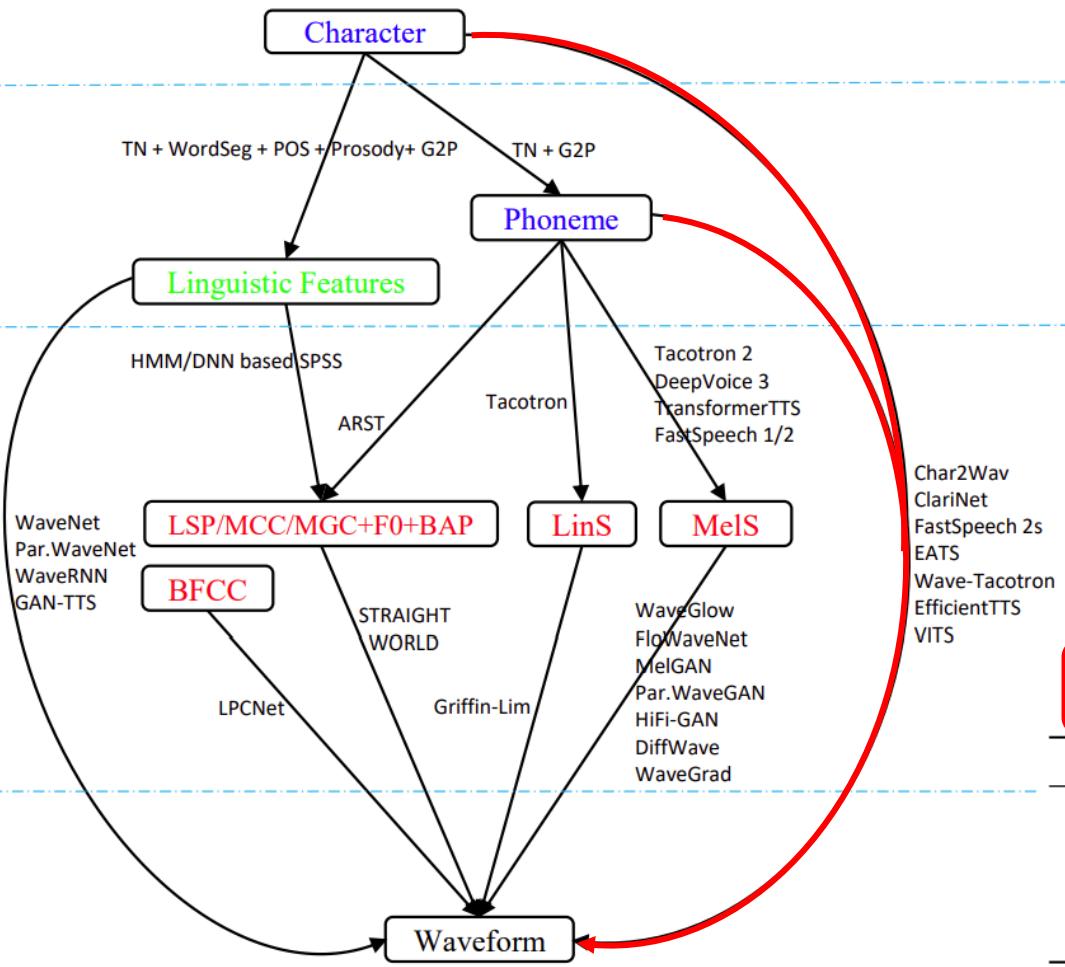
Acoustic Features

Waveform



# Data conversion pipeline

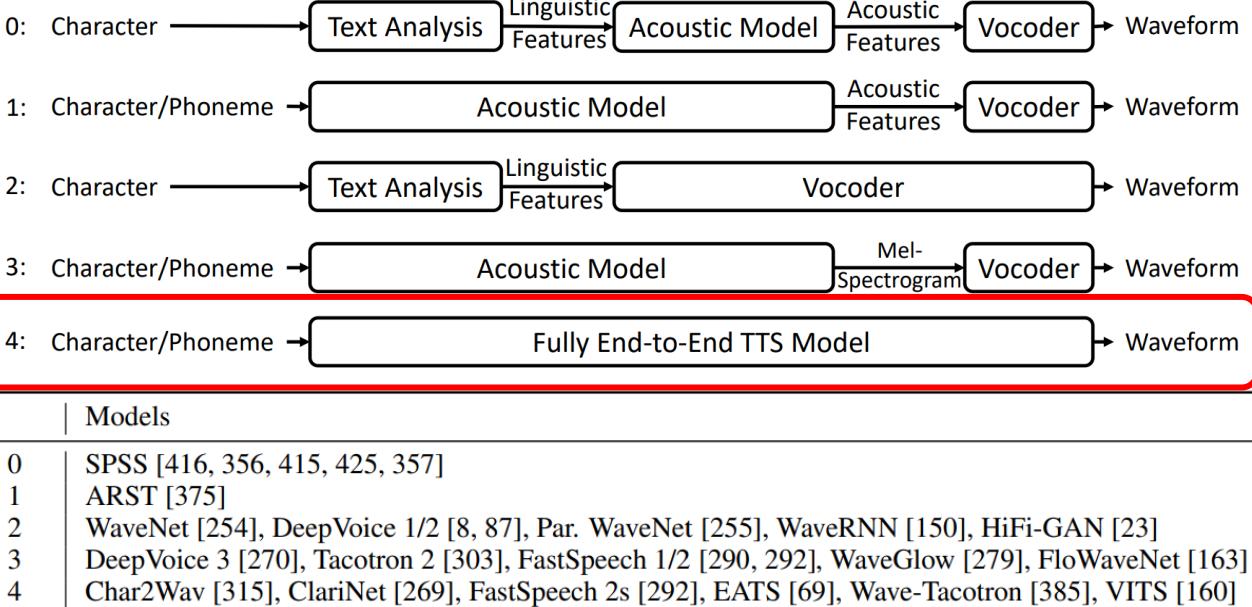
Text



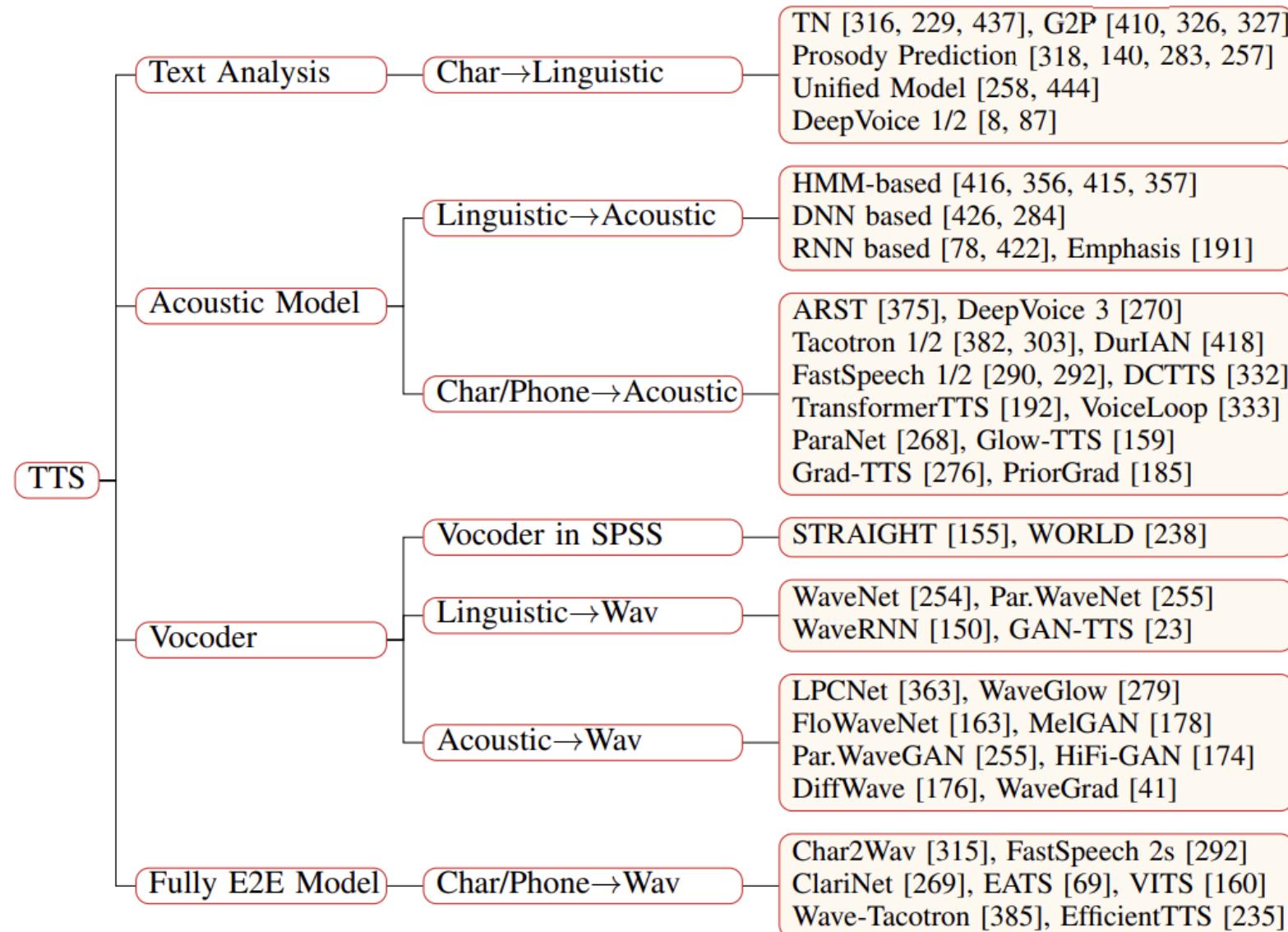
Linguistic Features

Acoustic Features

Waveform



# Key components in TTS



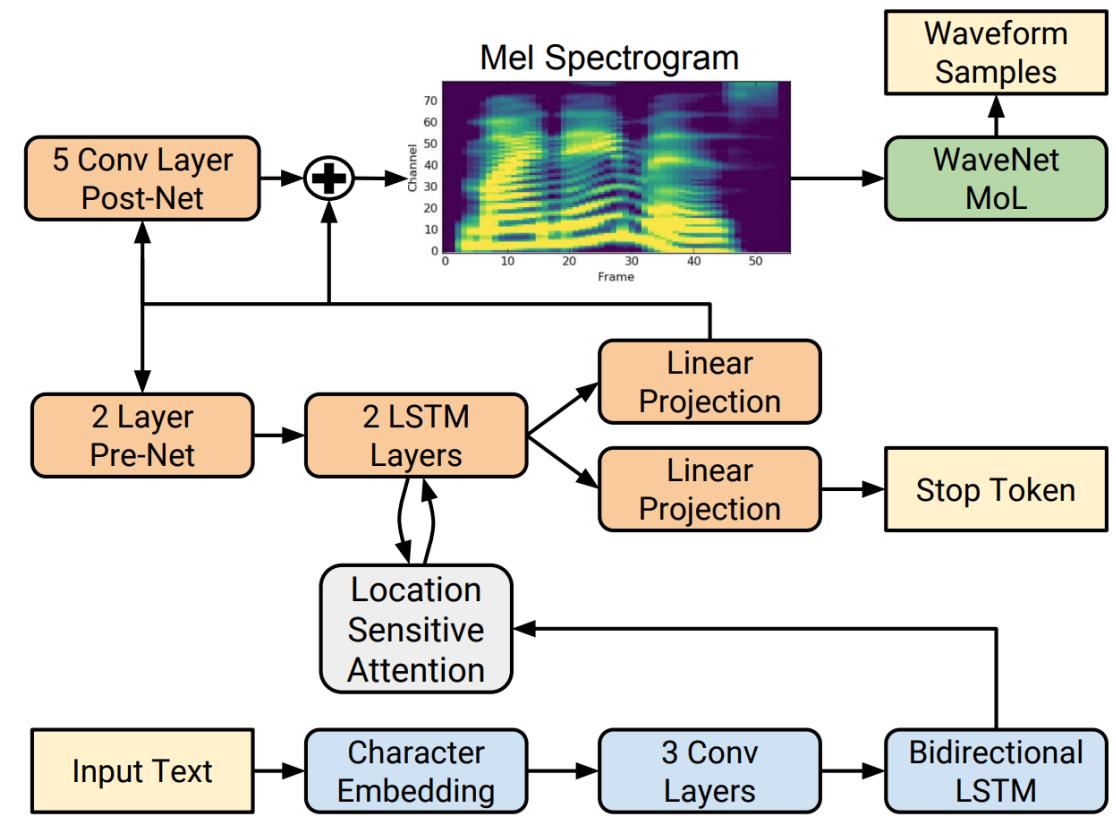
# Acoustic model

- Acoustic model in SPSS
- Acoustic models in end-to-end TTS
  - RNN-based (e.g., Tacotron series)
  - CNN-based (e.g., DeepVoice series)
  - Transformer-based (e.g., FastSpeech series)
  - Other (e.g., Flow, GAN, VAE, Diffusion)

	Acoustic Model	Input→Output	AR/NAR	Modeling	Structure
SPSS	HMM-based [416, 356]	Ling→MCC+F0	/	/	HMM
	DNN-based [426]	Ling→MCC+BAP+F0	NAR	/	DNN
	LSTM-based [78]	Ling→LSP+F0	AR	/	RNN
	EMPHASIS [191]	Ling→LinS+CAP+F0	AR	/	Hybrid
	ARST [375]	Ph→LSP+BAP+F0	AR	Seq2Seq	RNN
	VoiceLoop [333]	Ph→MGC+BAP+F0	AR	/	hybrid
RNN	Tacotron [382]	Ch→LinS	AR	Seq2Seq	Hybrid/RNN
	Tacotron 2 [303]	Ch→MelS	AR	Seq2Seq	RNN
	DurIAN [418]	Ph→MelS	AR	Seq2Seq	RNN
	Non-Att Tacotron [304]	Ph→MelS	AR	/	Hybrid/CNN/RNN
	Para. Tacotron 1/2 [74, 75]	Ph→MelS	NAR	/	Hybrid/Self-Att/CNN
	MelNet [367]	Ch→MelS	AR	/	RNN
CNN	DeepVoice [8]	Ch/Ph→MelS	AR	/	CNN
	DeepVoice 2 [87]	Ch/Ph→MelS	AR	/	CNN
	DeepVoice 3 [270]	Ch/Ph→MelS	AR	Seq2Seq	CNN
	ParaNet [268]	Ph→MelS	NAR	Seq2Seq	CNN
	DCTTS [332]	Ch→MelS	AR	Seq2Seq	CNN
	SpeedySpeech [361]	Ph→MelS	NAR	/	CNN
Transformer	TalkNet 1/2 [19, 18]	Ch→MelS	NAR	/	CNN
	TransformerTTS [192]	Ph→MelS	AR	Seq2Seq	Self-Att
	MultiSpeech [39]	Ph→MelS	AR	Seq2Seq	Self-Att
	FastSpeech 1/2 [290, 292]	Ph→MelS	NAR	Seq2Seq	Self-Att
	AlignTTS [429]	Ch/Ph→MelS	NAR	Seq2Seq	Self-Att
	JDIT-T [197]	Ph→MelS	NAR	Seq2Seq	Self-Att
Flow	FastPitch [181]	Ph→MelS	NAR	Seq2Seq	Self-Att
	AdaSpeech 1/2/3 [40, 403, 404]	Ph→MelS	NAR	Seq2Seq	Self-Att
	DenoiSpeech [434]	Ph→MelS	NAR	Seq2Seq	Self-Att
	DeviceTTS [126]	Ph→MelS	NAR	/	Hybrid/DNN/RNN
	LightSpeech [220]	Ph→MelS	NAR	/	Hybrid/Self-Att/CNN
	Flow-TTS [234]	Ch/Ph→MelS	NAR*	Flow	Hybrid/CNN/RNN
VAE	Glow-TTS [159]	Ph→MelS	NAR	Flow	Hybrid/Self-Att/CNN
	Flowtron [366]	Ph→MelS	AR	Flow	Hybrid/RNN
	EfficientTTS [235]	Ch→MelS	NAR	Flow	Hybrid/CNN
	GMVAE-Tacotron [119]	Ph→MelS	AR	VAE	Hybrid/RNN
GAN	VAE-TTS [443]	Ph→MelS	AR	VAE	Hybrid/RNN
	BVAE-TTS [187]	Ph→MelS	NAR	VAE	CNN
	GAN exposure [99]	Ph→MelS	AR	GAN	Hybrid/RNN
Diffusion	TTS-Stylization [224]	Ch→MelS	AR	GAN	Hybrid/RNN
	Multi-SpectroGAN [186]	Ph→MelS	NAR	GAN	Hybrid/Self-Att/CNN
	Diff-TTS [141]	Ph→MelS	NAR*	Diffusion	Hybrid/CNN
GradTTS	GradTTS [276]	Ph→MelS	NAR	Diffusion	Hybrid/Self-Att/CNN
	PriorGrad [185]	Ph→MelS	NAR	Diffusion	Hybrid/Self-Att/CNN

# Acoustic model——RNN based

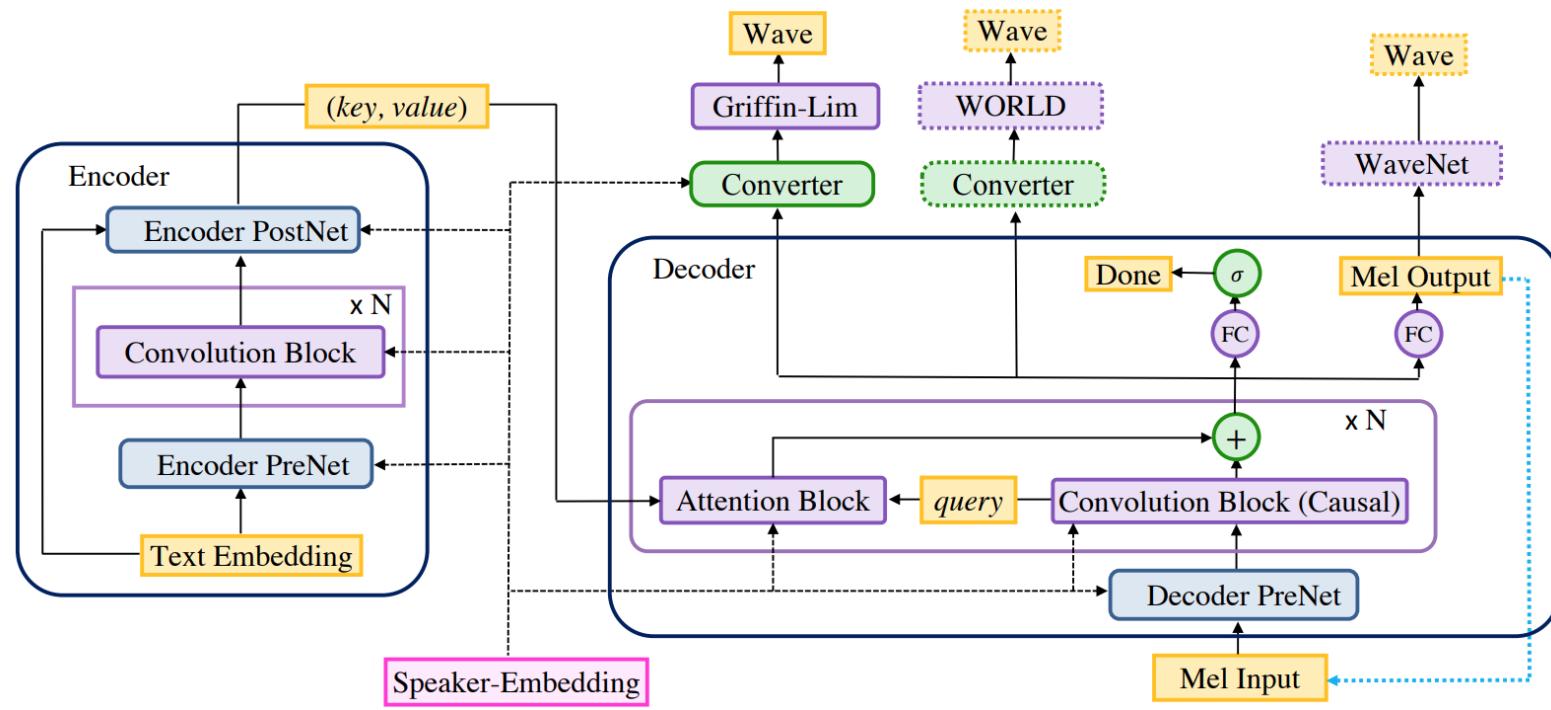
- Tacotron 2 [303]
  - Evolved from Tacotron [382]
  - Text to mel-spectrogram generation
  - LSTM based encoder and decoder
  - Location sensitive attention
  - WaveNet as the vocoder
- Other works
  - GST-Tacotron [383], Ref-Tacotron [309]
  - DurlAN [418]
  - Non-Attentative Tacotron [304]
  - Parallel Tacotron 1/2 [74, 75]
  - WaveTacotron [385]



# Acoustic model——CNN based

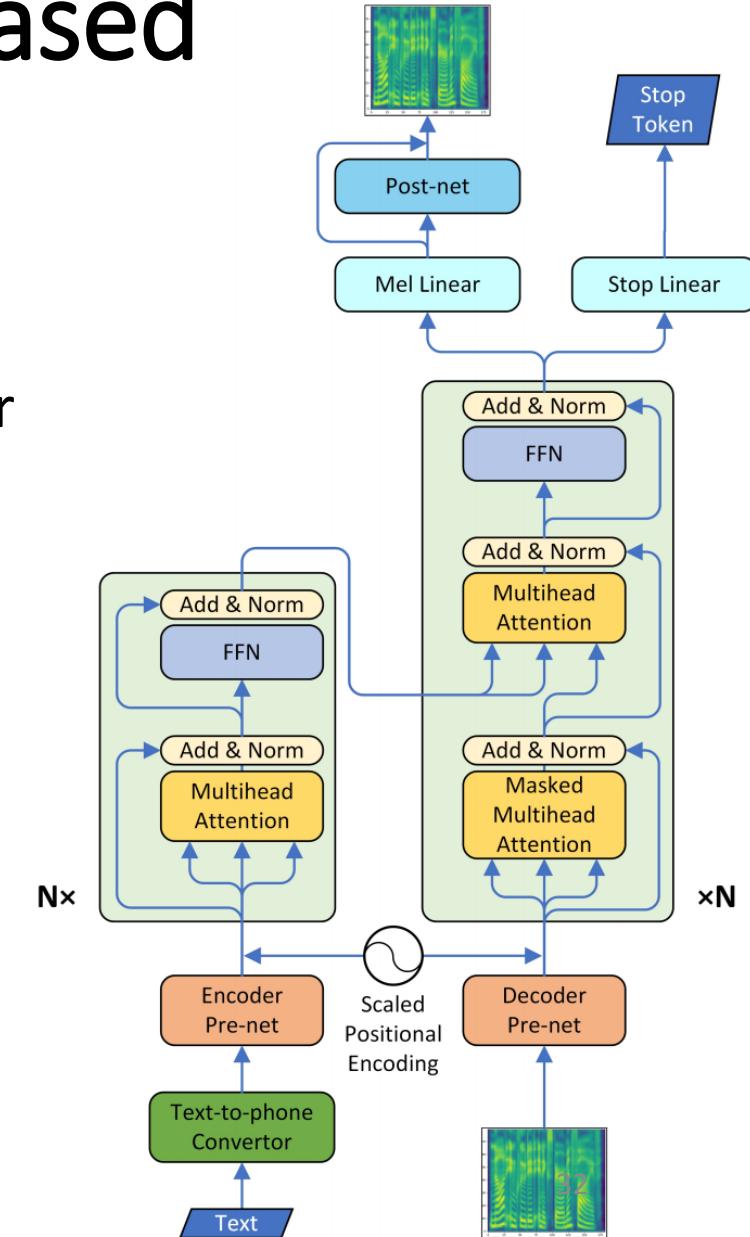
- DeepVoice 3 [270]
  - Evolved from DeepVoice 1/2 [8, 87]
  - Enhanced with purely CNN based structure
  - Support different acoustic features as output
  - Support multi-speakers

- Other works
  - DCTTS [332] (Contemporary)
  - ClariNet [269]
  - ParaNet [268]



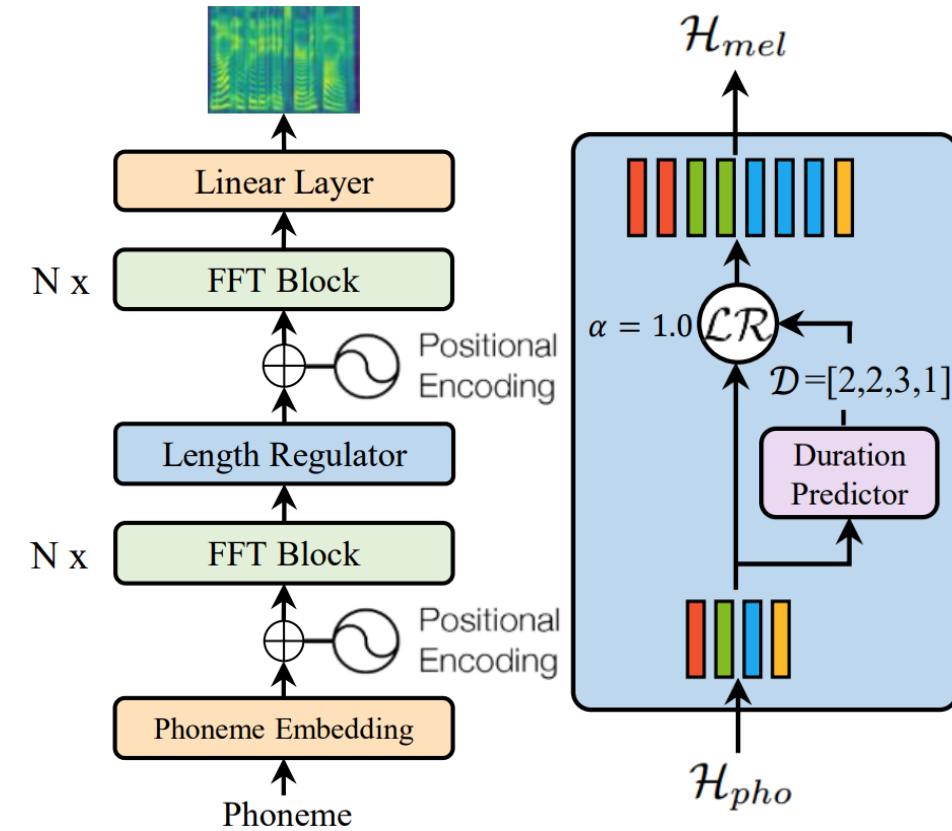
# Acoustic model——Transformer based

- TransformerTTS [192]
  - Framework is like Tacotron 2
  - Replace LSTM with Transformer in encoder and decoder
  - Parallel training, quality on par with Tacotron 2
  - Attention with more challenges than Tacotron 2, due to parallel computing
- Other works
  - MultiSpeech [39]
  - Robutrans [194]



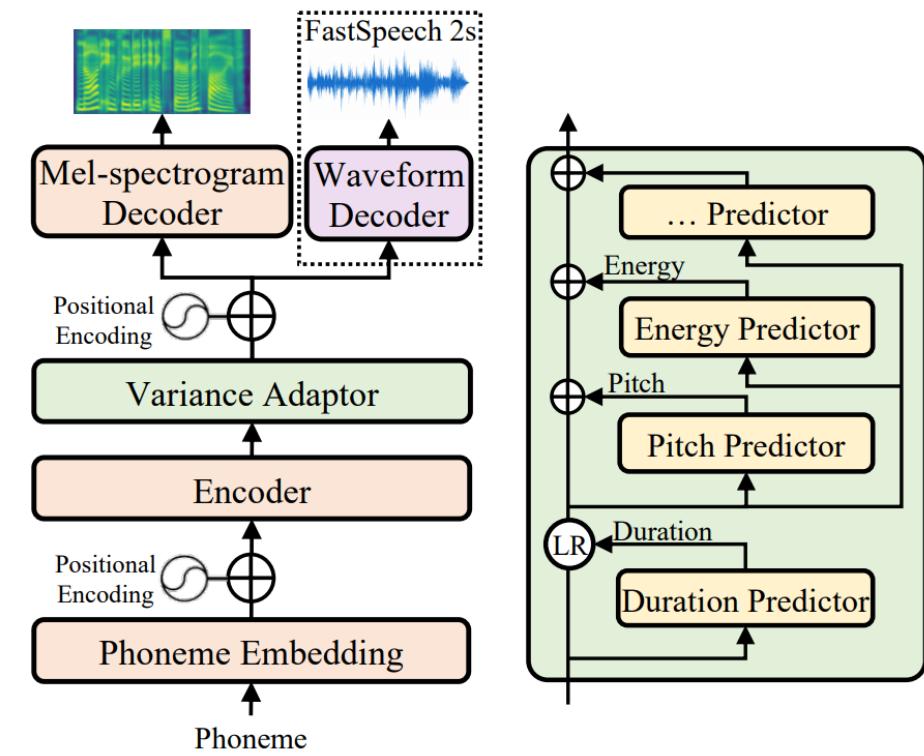
# Acoustic model——Transformer based

- FastSpeech [290]
  - Generate mel-spectrogram in parallel (for speedup)
  - Remove the text-speech attention mechanism (for robustness)
  - Feed-forward transformer with length regulator (for controllability)



# Acoustic model——Transformer based

- FastSpeech 2 [292]
  - Improve FastSpeech
  - Use variance adaptor to predict duration, pitch, energy, etc
  - Simplify training pipeline of FastSpeech (KD)
  - FastSpeech 2s: a fully end-to-end parallel text to wave model
- Other works
  - FastPitch [181]
  - JDI-T [197], AlignTTS [429]



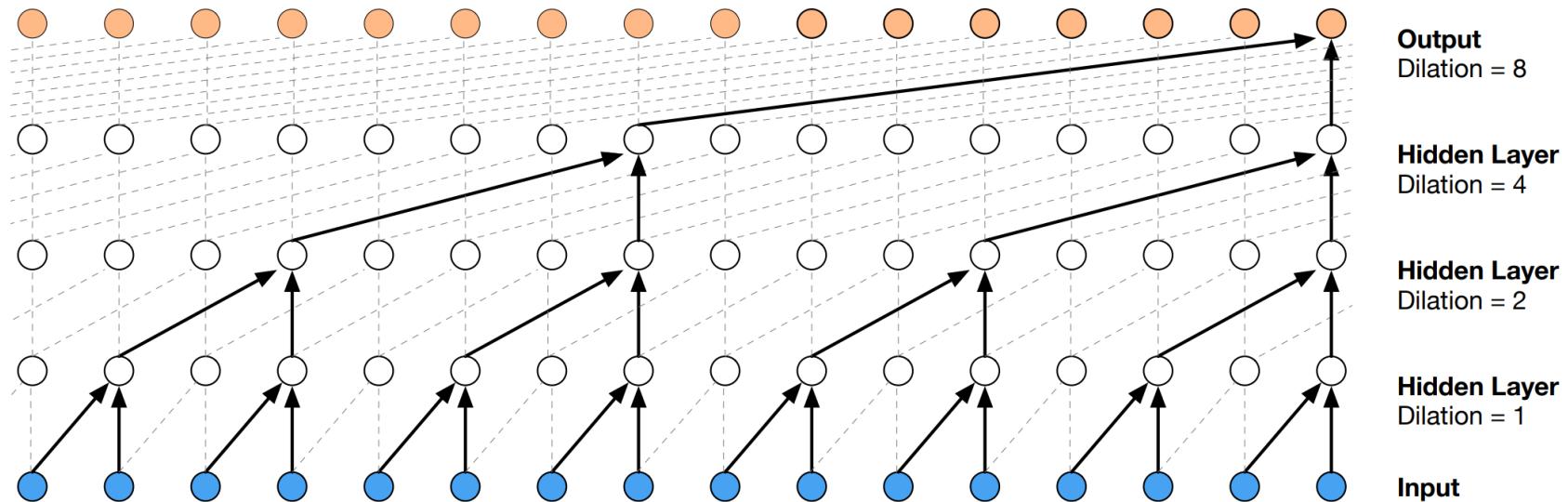
# Vocoder

- Autoregressive vocoder
- Flow-based vocoder
- GAN-based vocoder
- VAE-based vocoder
- Diffusion-based vocoder

	Vocoder	Input	AR/NAR	Modeling	Architecture
AR	WaveNet [254]	Linguistic Feature	AR	/	CNN
	SampleRNN [233]	/	AR	/	RNN
	WaveRNN [150]	Linguistic Feature	AR	/	RNN
	LPCNet [363]	BFCC	AR	/	RNN
	Univ. WaveRNN [215]	Mel-Spectrogram	AR	/	RNN
	SC-WaveRNN [265]	Mel-Spectrogram	AR	/	RNN
	MB WaveRNN [418]	Mel-Spectrogram	AR	/	RNN
	FFTNet [145]	Cepstrum	AR	/	CNN
Flow	Par. WaveNet [255]	Linguistic Feature	NAR	Flow	CNN
	WaveGlow [279]	Mel-Spectrogram	NAR	Flow	Hybrid/CNN
	FloWaveNet [163]	Mel-Spectrogram	NAR	Flow	Hybrid/CNN
	WaveFlow [271]	Mel-Spectrogram	AR	Flow	Hybrid/CNN
	SqueezeWave [433]	Mel-Spectrogram	NAR	Flow	CNN
GAN	WaveGAN [68]	/	NAR	GAN	CNN
	GELP [149]	Mel-Spectrogram	NAR	GAN	CNN
	GAN-TTS [23]	Linguistic Feature	NAR	GAN	CNN
	MelGAN [178]	Mel-Spectrogram	NAR	GAN	CNN
	Par. WaveGAN [402]	Mel-Spectrogram	NAR	GAN	CNN
	HiFi-GAN [174]	Mel-Spectrogram	NAR	GAN	Hybrid/CNN
	VocGAN [408]	Mel-Spectrogram	NAR	GAN	CNN
	GED [96]	Linguistic Feature	NAR	GAN	CNN
	Fre-GAN [161]	Mel-Spectrogram	NAR	GAN	CNN
VAE	Wave-VAE [268]	Mel-Spectrogram	NAR	VAE	CNN
Diffusion	WaveGrad [41]	Mel-Spectrogram	NAR	Diffusion	Hybrid/CNN
	DiffWave [176]	Mel-Spectrogram	NAR	Diffusion	Hybrid/CNN
	PriorGrad [185]	Mel-Spectrogram	NAR	Diffusion	Hybrid/CNN

# Vocoder——AR

- WaveNet: autoregressive model with dilated causal convolution [254]



- Other works
  - WaveRNN [150]
  - LPCNet [363]

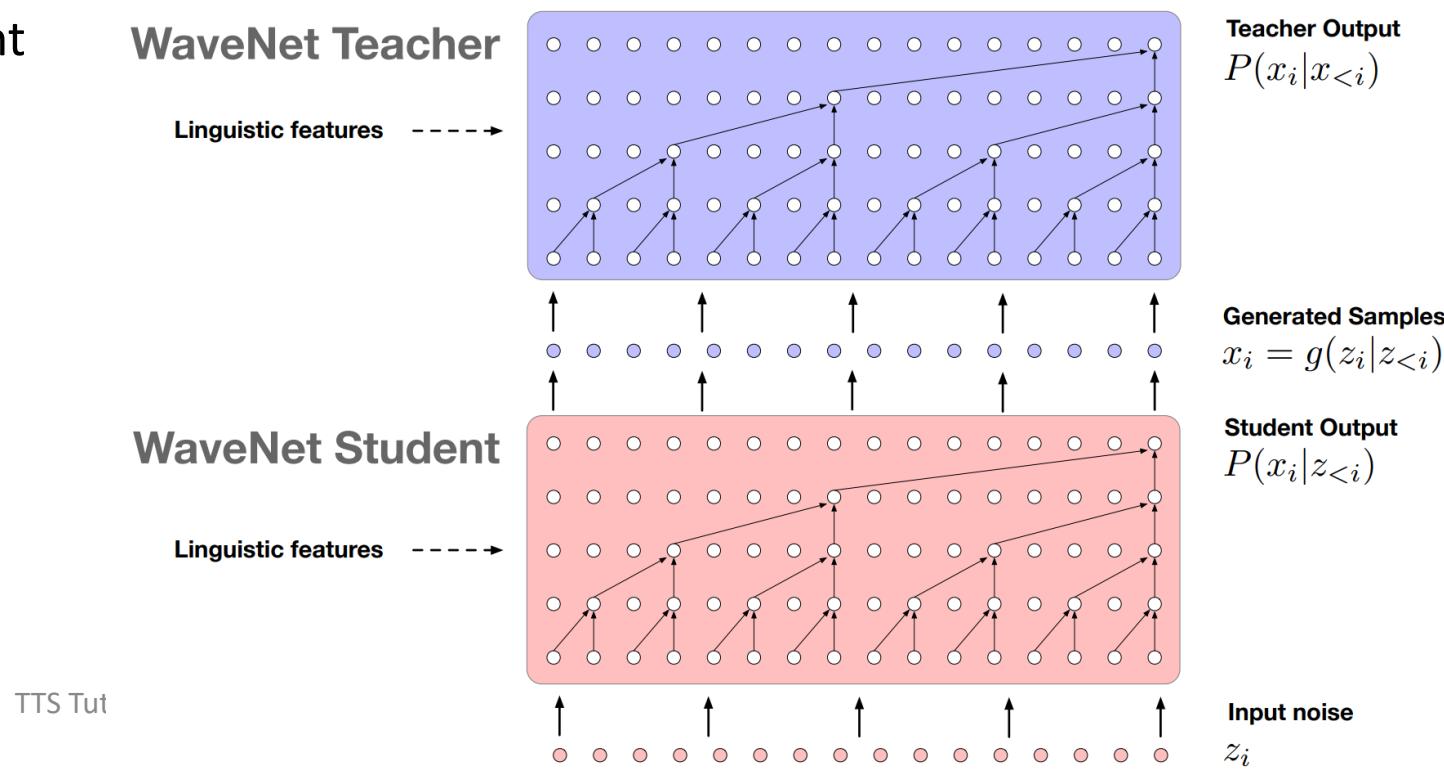
# Vocoder—Flow

- Map between data distribution  $x$  and standard (normalizing) prior distribution  $z$       Evaluation  $z = f^{-1}(x)$       Synthesis  $x = f(z)$
- Category of normalizing flow
  - AR (autoregressive): AF (autoregressive flow) and IAF (inverse autoregressive flow)
  - Bipartite: RealNVP and Glow

Flow	Evaluation $z = f^{-1}(x)$	Synthesis $x = f(z)$
AR	AF [261] $z_t = x_t \cdot \sigma_t(x_{<t}; \theta) + \mu_t(x_{<t}; \theta)$	$x_t = \frac{z_t - \mu_t(x_{<t}; \theta)}{\sigma_t(x_{<t}; \theta)}$
	IAF [169] $z_t = \frac{x_t - \mu_t(z_{<t}; \theta)}{\sigma_t(z_{<t}; \theta)}$	$x_t = z_t \cdot \sigma_t(z_{<t}; \theta) + \mu_t(z_{<t}; \theta)$
Bipartite	RealNVP [66] $z_a = x_a,$	$x_a = z_a,$
	Glow [167] $z_b = x_b \cdot \sigma_b(x_a; \theta) + \mu_b(x_a; \theta)$	$x_b = \frac{z_b - \mu_b(x_a; \theta)}{\sigma_b(x_a; \theta)}$

# Vocoder—Flow

- Parallel WaveNet [255] (AR)
  - Knowledge distillation: Student (IAF), Teacher (AF)
  - Combine the best of both worlds
    - Parallel inference of IAF student
    - Parallel training of AF teacher
- Other works
  - ClariNet [269]



TTS Tut

# Vocoder—Flow

- WaveGlow [279] (Bipartite)

- Flow based transformation

$$z = f_k^{-1} \circ f_{k-1}^{-1} \circ \dots f_0^{-1}(x) \quad x = f_0 \circ f_1 \circ \dots f_k(z)$$

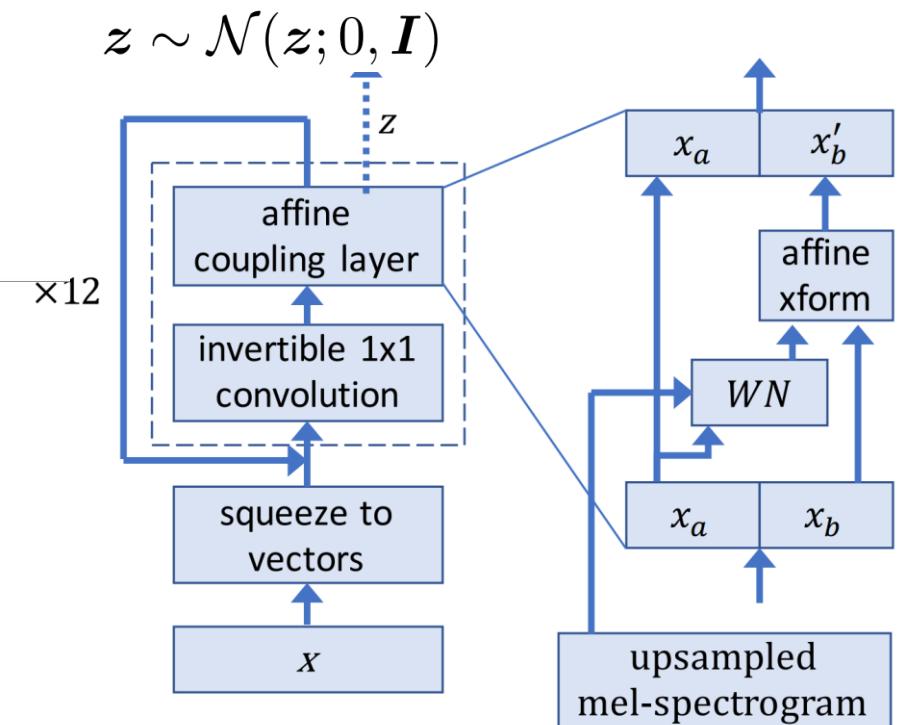
- Affine Coupling Layer

$$\begin{aligned} x_a, x_b &= \text{split}(x) \\ (\log s, t) &= WN(x_a, \text{mel-spectrogram}) \end{aligned}$$

$$x_{b'} = s \odot x_b + t$$

$$f_{coupling}^{-1}(x) = \text{concat}(x_a, x_{b'})$$

- Other works
  - FloWaveNet [163]
  - WaveFlow [271]



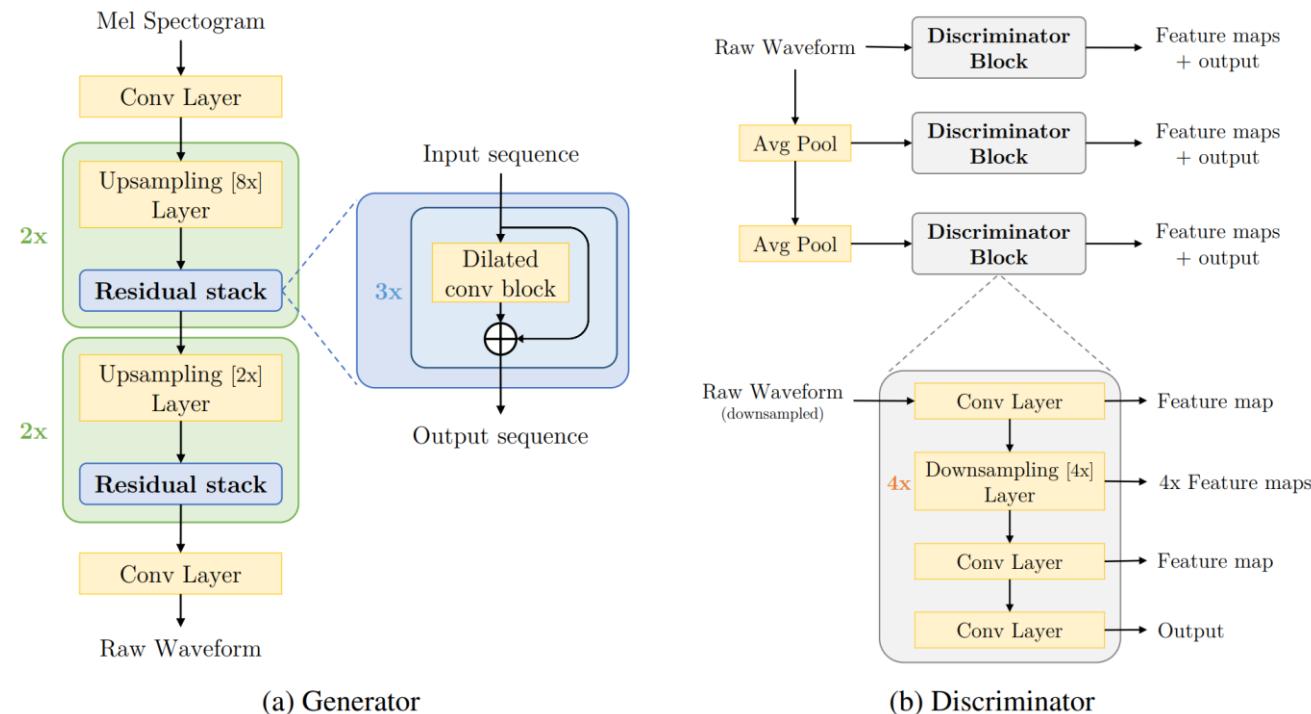
# Vocoder—GAN

- Category of GAN based vocoders

GAN	Generator	Discriminator	Loss
WaveGAN [68]	DCGAN [287]	/	WGAN-GP [97]
GAN-TTS [23]	/	Random Window D	Hinge-Loss GAN [198]
MelGAN [178]	/	Multi-Scale D	LS-GAN [231] Feature Matching Loss [182]
Par.WaveGAN [402]	WaveNet [254]	/	LS-GAN, Multi-STFT Loss
HiFi-GAN [174]	Multi-Receptive Field Fusion	Multi-Period D, Multi-Scale D	LS-GAN, STFT Loss, Feature Matching Loss
VocGAN [408]	Multi-Scale G	Hierarchical D	LS-GAN, Multi-STFT Loss, Feature Matching Loss
GED [96]	/	Random Window D	Hinge-Loss GAN, Repulsive loss

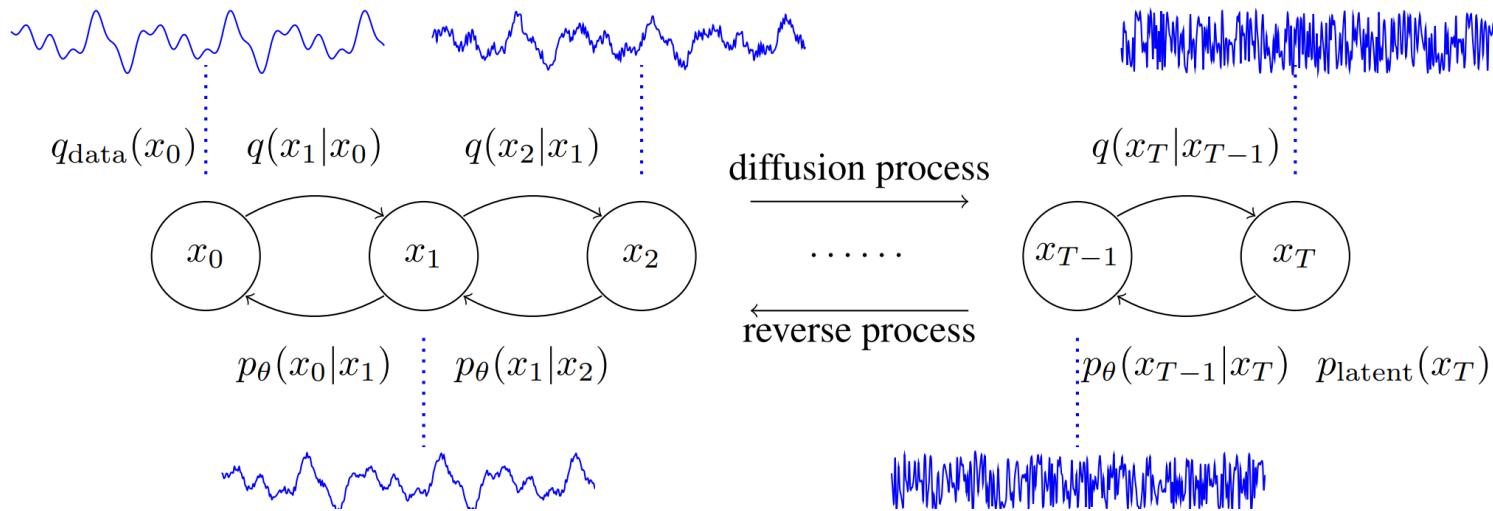
# Vocoder—GAN

- MelGAN [68]
  - Generator: Transposed conv for upsampling, dilated conv to increase receptive field
  - Discriminator: Multi-scale discrimination



# Vocoder—Diffusion

- Diffusion probabilistic model: DiffWave [176], WaveGrad [41]




---

**Algorithm 1** Training

---

```

for  $i = 1, 2, \dots, N_{\text{iter}}$  do
    Sample  $x_0 \sim q_{\text{data}}$ ,  $\epsilon \sim \mathcal{N}(0, I)$ , and
     $t \sim \text{Uniform}(\{1, \dots, T\})$ 
    Take gradient step on
     $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|_2^2$ 
    according to Eq. (7)
end for

```

---

**Algorithm 2** Sampling

---

```

Sample  $x_T \sim p_{\text{latent}} = \mathcal{N}(0, I)$ 
for  $t = T, T - 1, \dots, 1$  do
    Compute  $\mu_{\theta}(x_t, t)$  and  $\sigma_{\theta}(x_t, t)$  using Eq. (5)
    Sample  $x_{t-1} \sim p_{\theta}(x_{t-1}|x_t) =$ 
         $\mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \sigma_{\theta}(x_t, t)^2 I)$ 
end for
return  $x_0$ 

```

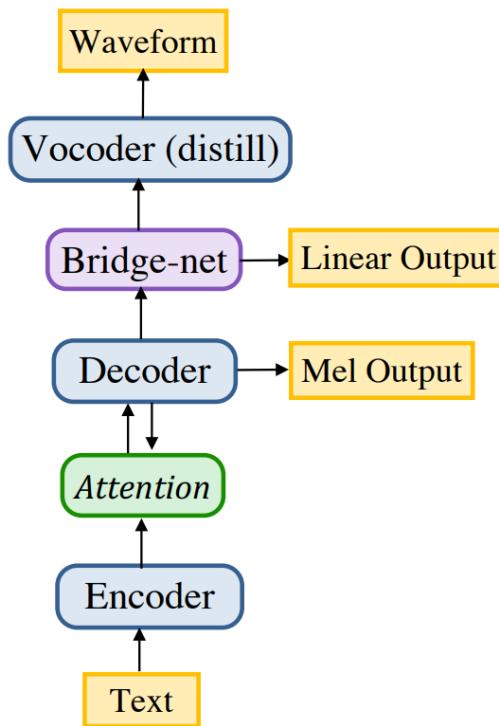
# Vocoder

- A comparison among different vocoders
  - Simplicity in math formulation and optimization
  - Support parallel generation
  - Support latent manipulation
  - Support likelihood estimation

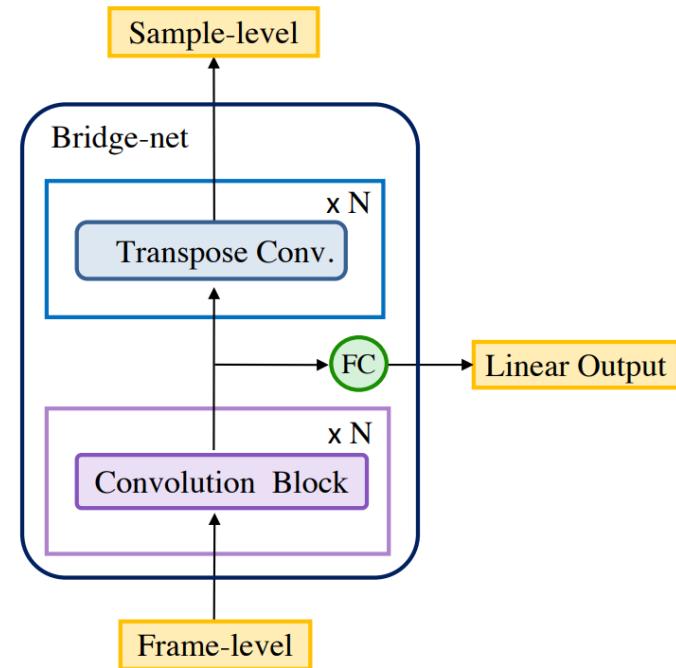
Generative Model	AR	VAE	Flow/AR	Flow/Bipartite	Diffusion	GAN
Vocoder (e.g.)	WaveNet	WaveVAE	Par.WaveNet	WaveGlow	DiffWave	MelGAN
Simple	Y	N	N	N	N	N
Parallel	N	Y	Y	Y	Y	Y
Latent Manipulate	N	Y	Y	Y	Y	Y*
Likelihood Estimate	Y	Y	Y	Y	Y	N

# Fully End-to-End TTS

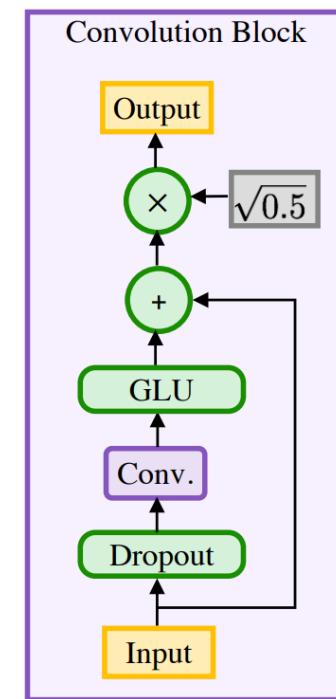
- ClariNet: AR acoustic model and NAR vocoder [269]



(a) Text-to-wave architecture



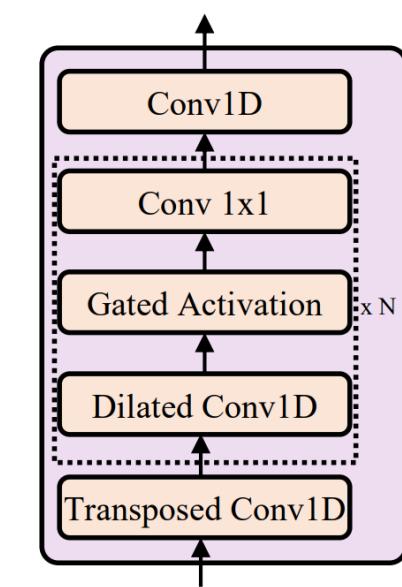
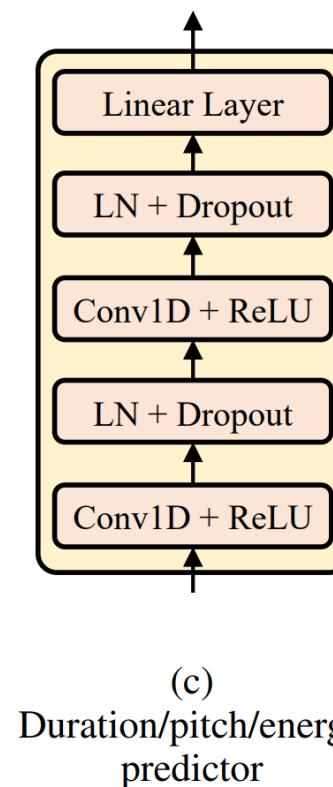
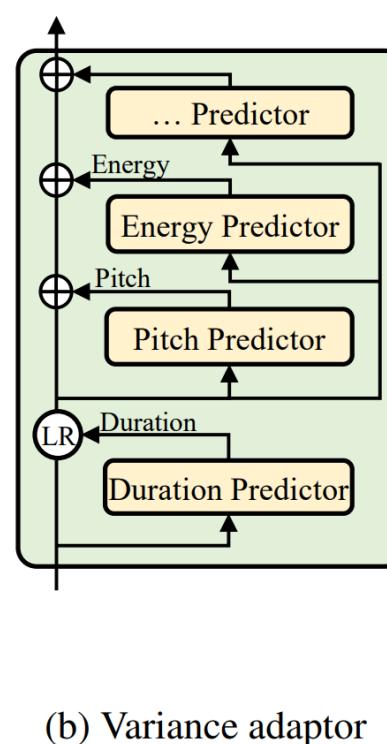
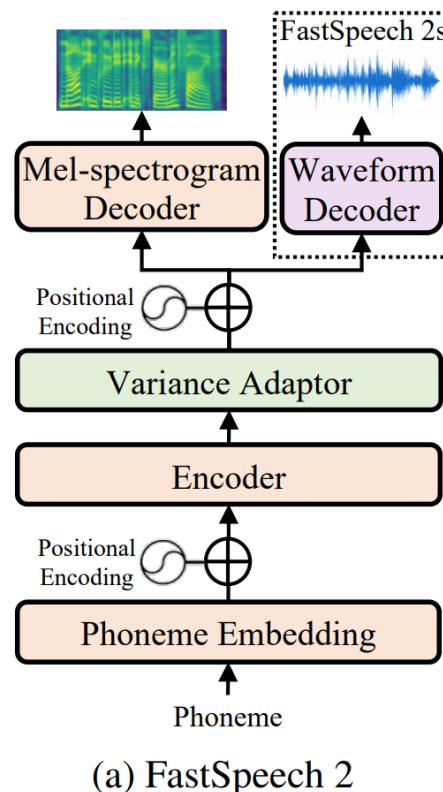
(b) Bridge-net



(c) Convolution block

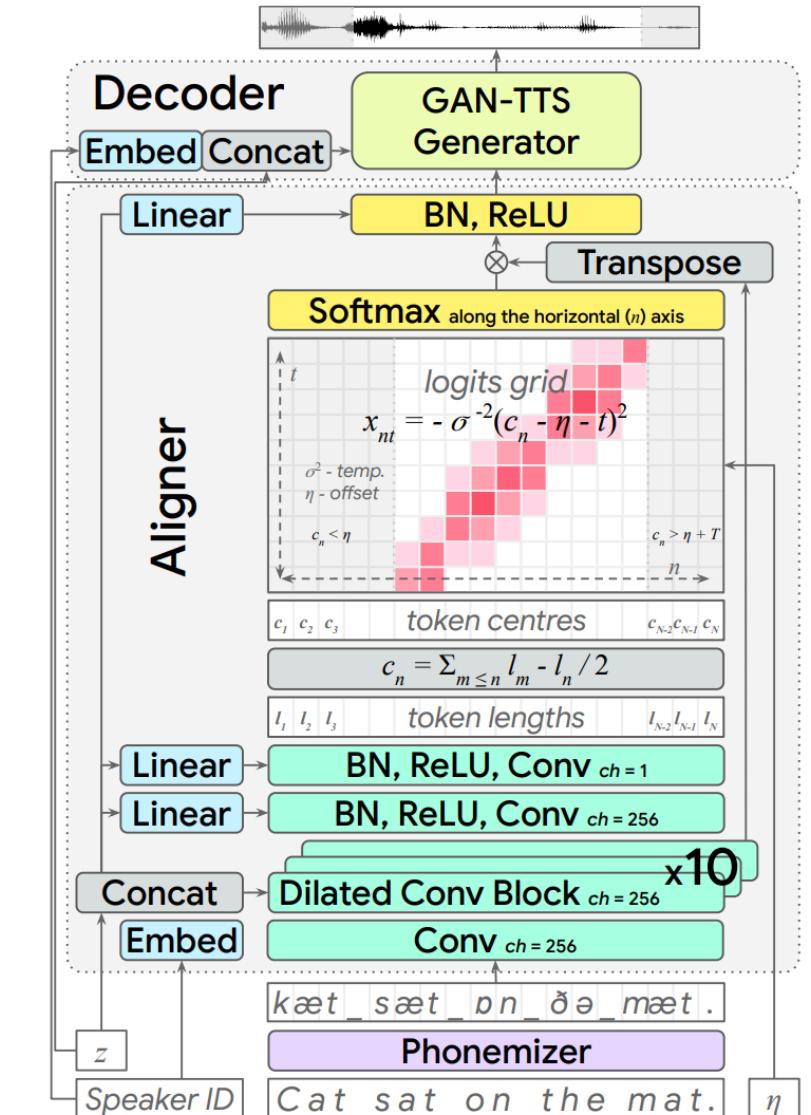
# Fully End-to-End TTS

- FastSpeech 2s: fully parallel text to wave model [292]



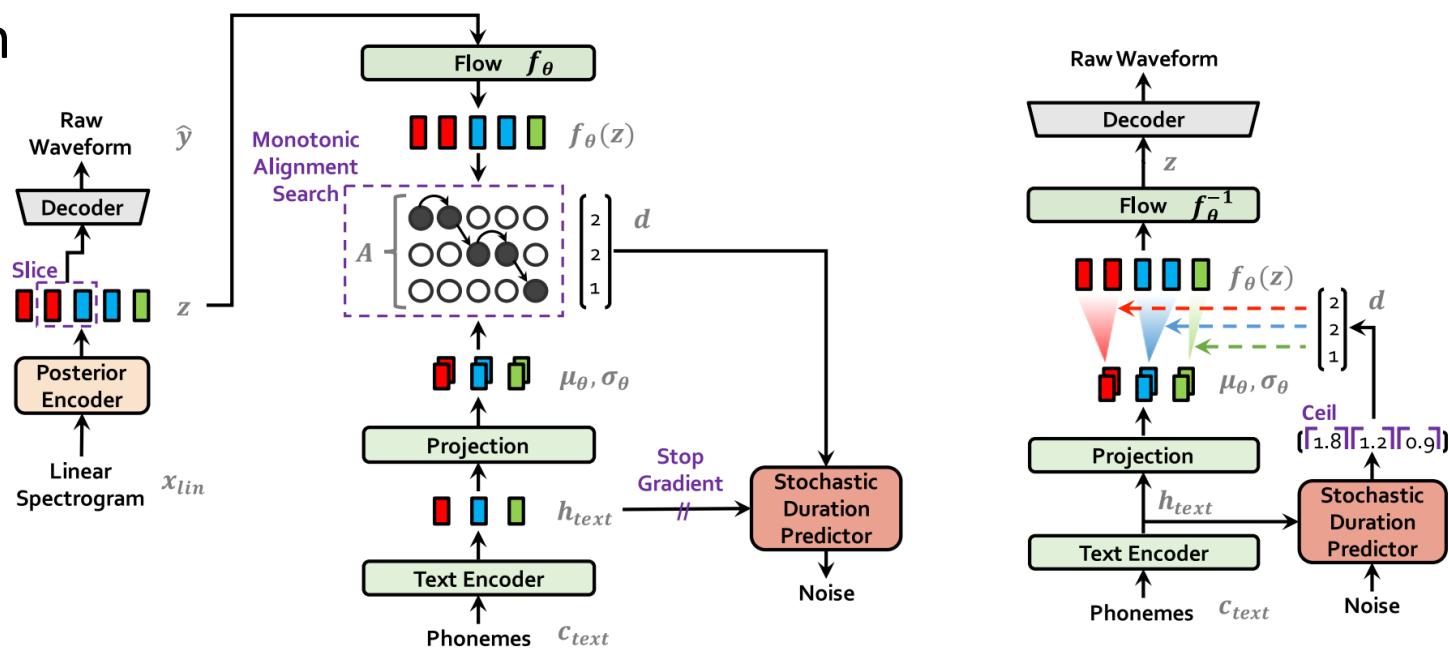
# Fully End-to-End TTS

- EATS: fully parallel text to wave model [69]
  - Duration prediction
  - Monotonic interpolation for upsampling
  - Soft dynamic time warping loss
  - Adversarial training



# Fully End-to-End TTS

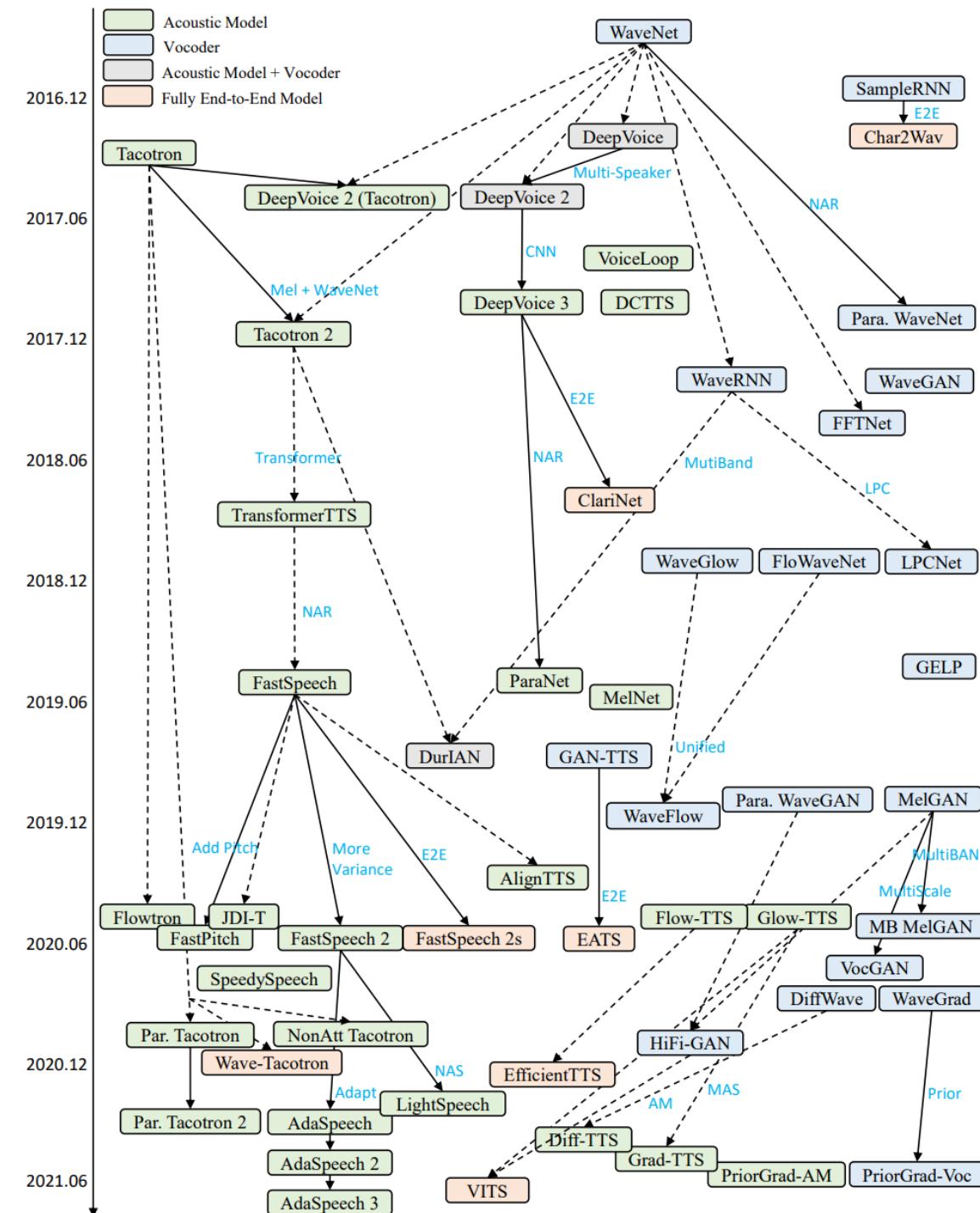
- VITS [160]
  - VAE, Flow, GAN
  - VAE: input mel, output waveform
  - Flow for VAE prior
  - GAN for waveform generation
  - Monotonic alignment search



(a) Training procedure

(b) Inference procedure

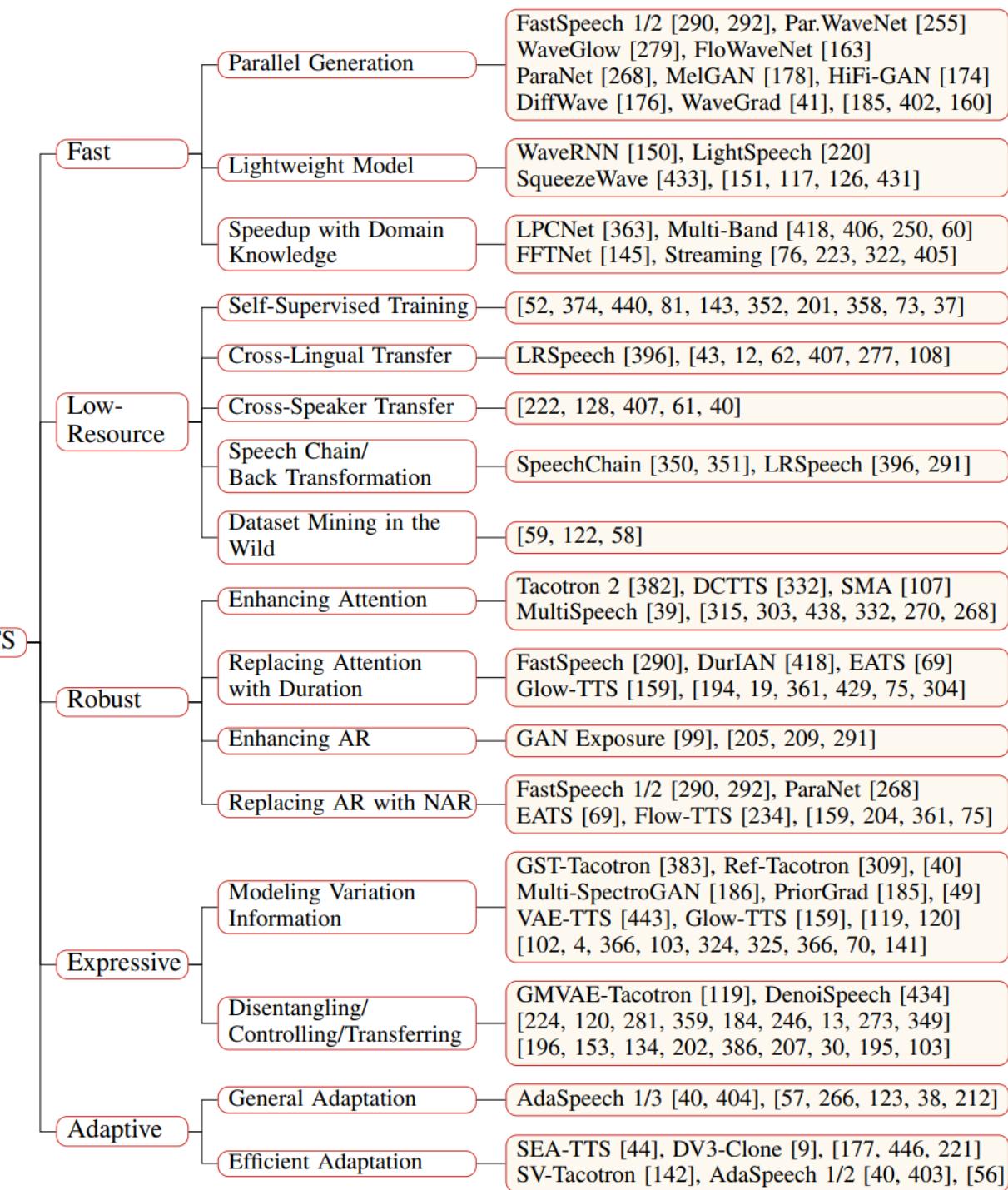
# Evolution of TTS



# Part 3: Advanced Topics in TTS

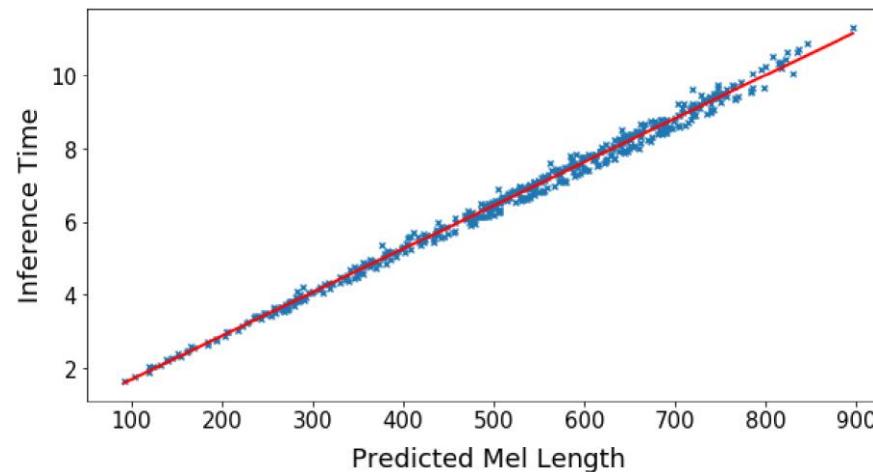
# Advanced topics in TTS

- Fast TTS
- Low-resource TTS
- Robust TTS
- Expressive TTS
- Adaptive TTS



# Fast TTS

- The model usually adopts autoregressive mel and waveform generation
  - Sequence is very long, e.g., 1s speech, 100 mel, 24000 waveform points
  - Slow inference speed



- The model size is usually large
  - Slow in low-end GPU and edge device

# Fast TTS

- Parallel generation

Modeling Paradigm	TTS Model	Training	Inference
AR (RNN)	Tacotron 1/2, SampleRNN, LPCNet	$\mathcal{O}(N)$	$\mathcal{O}(N)$
AR (CNN/Self-Att)	DeepVoice 3, TransformerTTS, WaveNet	$\mathcal{O}(1)$	$\mathcal{O}(N)$
NAR (CNN/Self-Att)	FastSpeech 1/2, ParaNet	$\mathcal{O}(1)$	$\mathcal{O}(1)$
NAR (GAN/VAE)	MelGAN, HiFi-GAN, FastSpeech 2s, EATS	$\mathcal{O}(1)$	$\mathcal{O}(1)$
Flow (AR)	Par. WaveNet, ClariNet, Flowtron	$\mathcal{O}(1)$	$\mathcal{O}(1)$
Flow (Bipartite)	WaveGlow, FloWaveNet, Glow-TTS	$\mathcal{O}(T)$	$\mathcal{O}(T)$
Diffusion	DiffWave, WaveGrad, Grad-TTS, PriorGrad	$\mathcal{O}(T)$	$\mathcal{O}(T)$

- Lightweight model
  - pruning, quantization, knowledge distillation, and neural architecture search
- Speedup with domain knowledge
  - linear prediction, multiband modeling, subscale prediction, multi-frame prediction, streaming synthesis

# Low-resource TTS

- There are **7,000+** languages in the world, but popular commercialized speech services only support **dozens of** languages
  - There is strong business demand to support more languages in TTS.



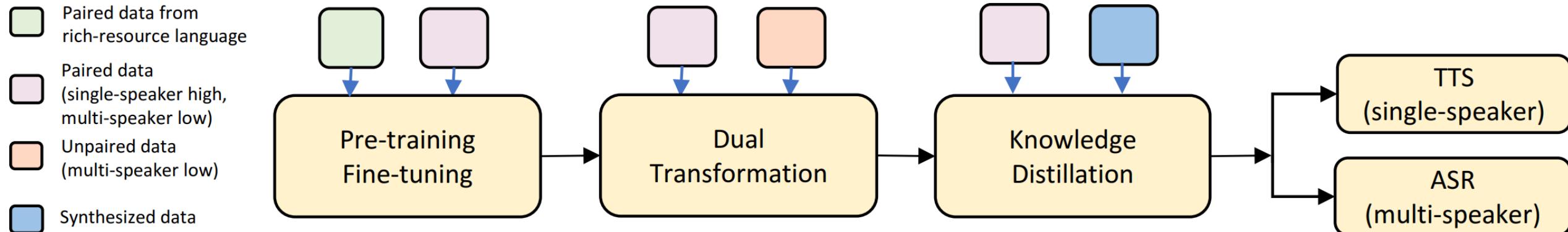
- However, lack of data in low-resource languages and the data collection cost is high.

# Low-resource TTS

Techniques	Data	Work
Self-supervised Training	Unpaired text or speech	[52, 374, 440, 81, 143, 352, 201, 358, 73]
Cross-lingual Transfer	Paired text and speech	[43, 396, 12, 407, 62, 277, 108]
Cross-speaker Transfer	Paired text and speech	[222, 128, 61, 407, 40]
Speech chain/Back transformation	Unpaired text or speech	[291, 396, 350, 351]
Dataset mining in the wild	Paired text and speech	[59, 122, 58]

- Self-supervised training
  - Text pre-training, speech pre-training, discrete token quantization
- Cross-lingual transfer
  - Languages share similarity, phoneme mapping/re-initialization/IPA/byte
- Cross-speaker transfer
  - Voice conversion, voice adaptation
- Speech chain/back transformation
  - TTS  $\leftrightarrow$  ASR
- Dataset mining in the wild
  - Speech enhancement, denoising, disentangling

# Low-resource TTS---LRSpeech [396]



- **Step 1:** Language transfer
  - Human languages share similar pronunciations; Rich-resource language data is “free”
- **Step 2:** TTS and ASR help with each other
  - Leverage the task duality with unpaired speech and text data
- **Step 3:** Customization for product deployment with knowledge distillation
  - Better accuracy by data knowledge distillation
  - Customize multi-speaker TTS to a target-speaker TTS, and to small model

# Robust TTS

- Robustness issues
  - Word skipping, repeating, attention collapse

*You can call me directly at xx(deleted due to privacy issue)xx or my cell xx(deleted due to privacy issue)xx or send me a meeting request with all the appropriate information.*
- The cause of robustness issues
  - The difficulty of alignment learning between text and mel-spectrograms
  - Exposure bias and error propagation in AR generation
- The solutions
  - Enhance attention
  - Replace attention with duration prediction
  - Enhance AR
  - Replace AR with NAR

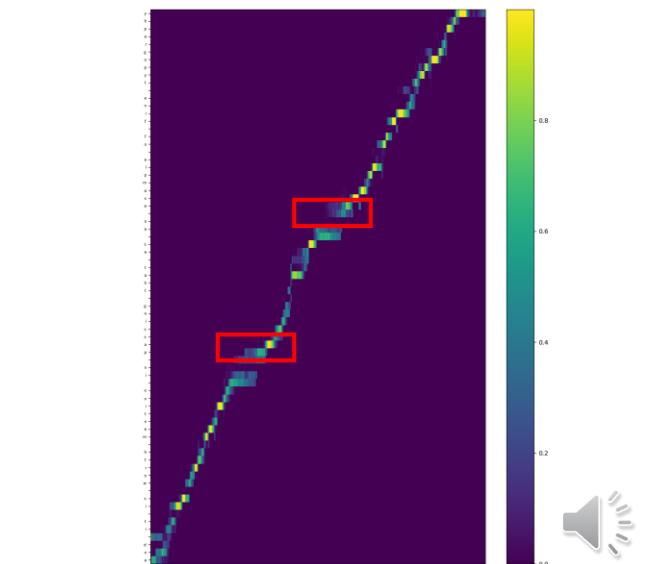
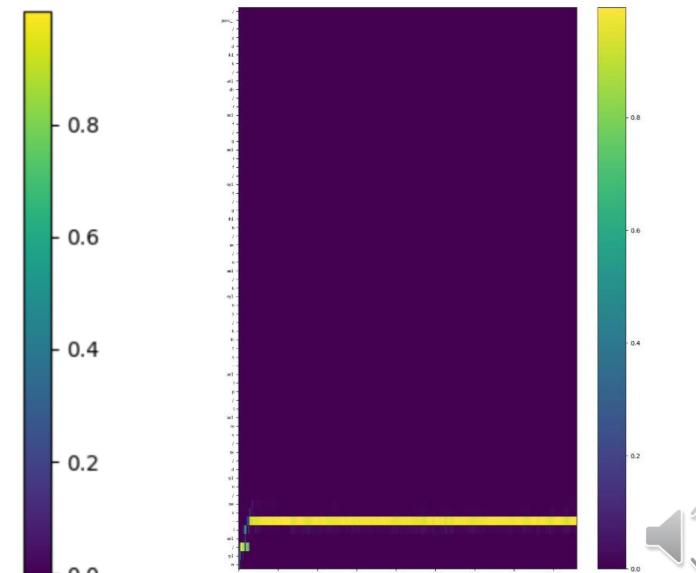
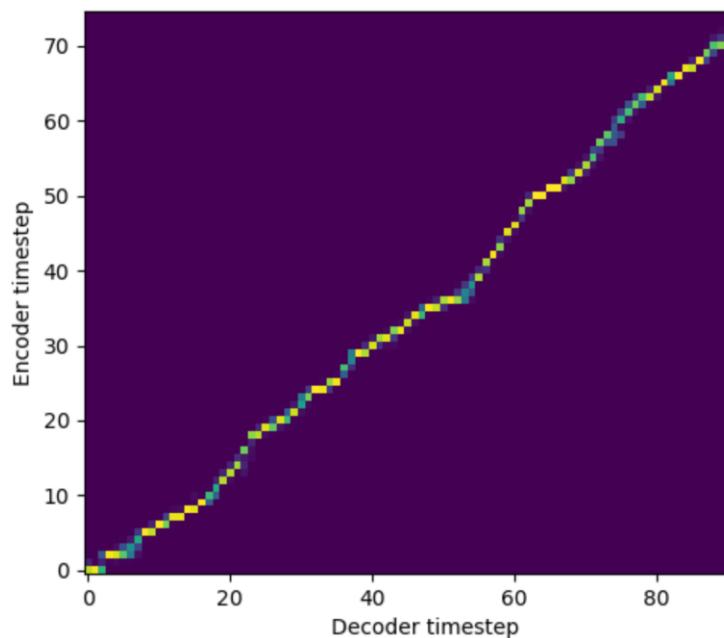


# Robust TTS

Category	Technique	Work
Enhancing Attention	Content-based attention Location-based attention Content/Location hybrid attention Monotonic attention Windowing or off-diagonal penalty Enhancing enc-dec connection Positional attention	[382, 192] [315, 333, 367, 17] [303] [438, 107, 411] [332, 438, 270, 39] [382, 303, 270, 203, 39] [268, 234, 204]
Replacing Attention with Duration Prediction	Label from encoder-decoder attention Label from CTC alignment Label from HMM alignment Dynamic programming Monotonic alignment search Monotonic interpolation with soft DTW	[290, 361, 197, 181] [19] [292, 418, 194, 252, 74, 304] [429, 193, 235] [159] [69, 75]
Enhancing AR	Professor forcing Reducing training/inference gap Knowledge distillation Bidirectional regularization	[99, 205] [361] [209] [291, 452]
Replacing AR with NAR	Parallel generation	[290, 292, 268, 69]

# Robust TTS——Attention improvement

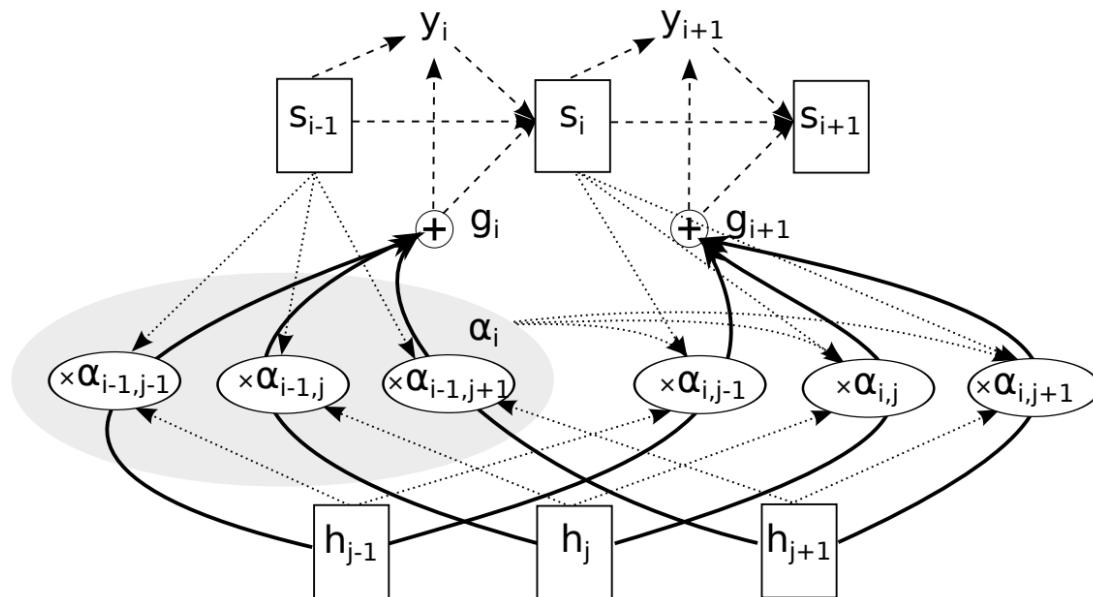
- Encoder-decoder attention: alignment between text and mel
  - Local, monotonic, and complete



And it is worth mention **in** passing that,  
**as an** example of fine typography

# Robust TTS——Attention improvement

- Location sensitive attention [50, 303]
  - Use previous alignment to compute the next attention alignment



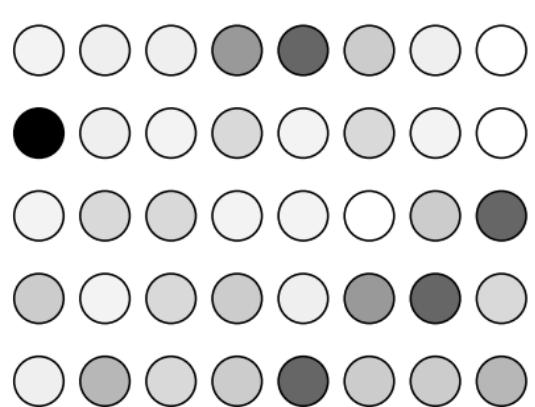
$$\alpha_i = \text{Attend}(s_{i-1}, \alpha_{i-1}, h)$$

$$g_i = \sum_{j=1}^L \alpha_{i,j} h_j$$

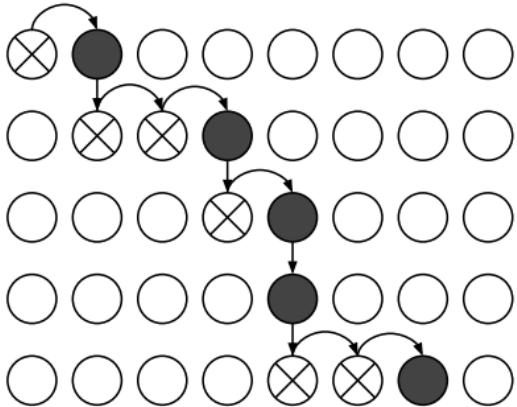
$$y_i \sim \text{Generate}(s_{i-1}, g_i),$$

# Robust TTS——Attention improvement

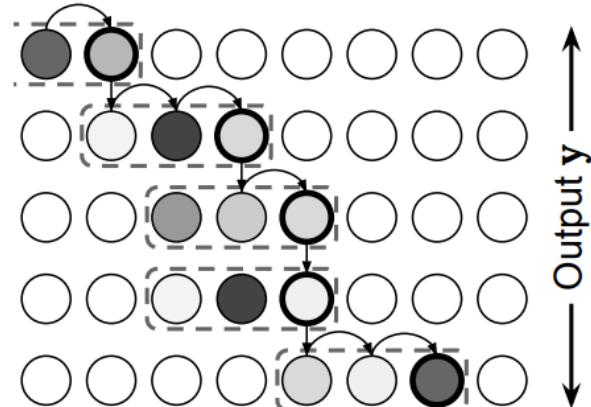
- Monotonic attention [288, 47]
  - The attention position is monotonically increasing



(a) Soft attention.



(b) Hard monotonic attention.



(c) Monotonic chunkwise attention.

$$e_{i,j} = \text{MonotonicEnergy}(s_{i-1}, h_j)$$

$$p_{i,j} = \sigma(e_{i,j})$$

$$z_{i,j} \sim \text{Bernoulli}(p_{i,j})$$

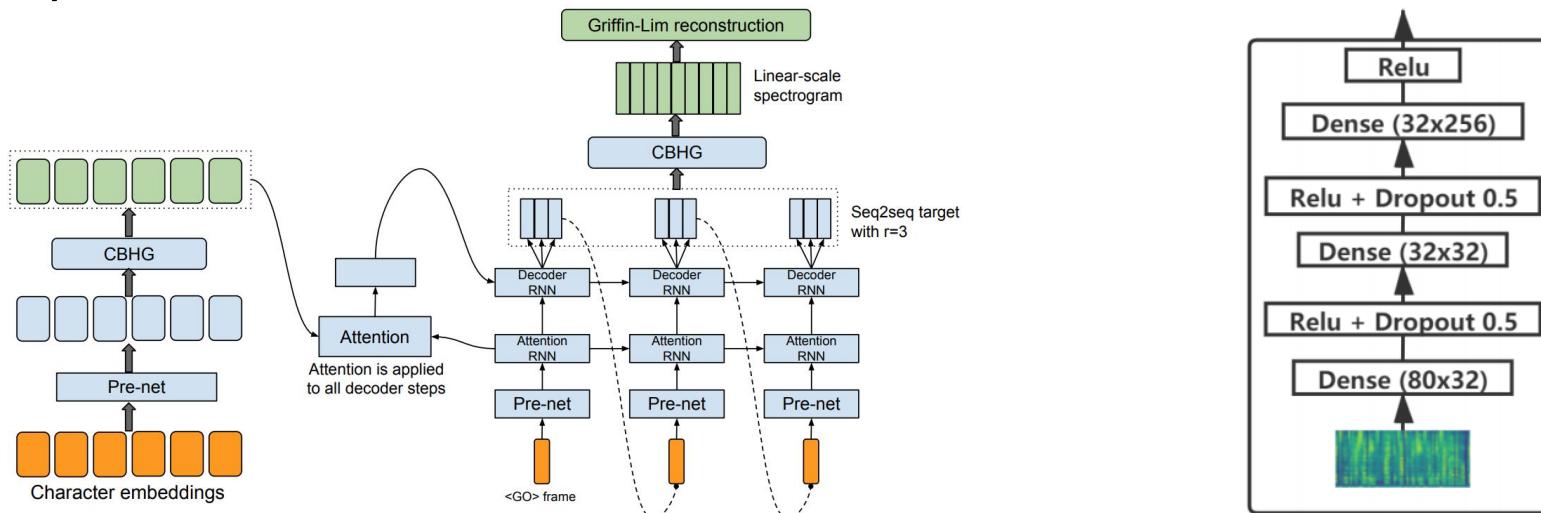
# Robust TTS—Attention improvement

- Windowing [332, 438]
  - Only a subset of the encoding results  $\hat{\mathbf{x}} = [\mathbf{x}_{p-w}, \dots, \mathbf{x}_{p+w}]$  are considered at each decoder timestep when using the windowing technique
- Penalty loss for off-diagonal attention distribution [39]
  - Guided attention loss with diagonal band mask



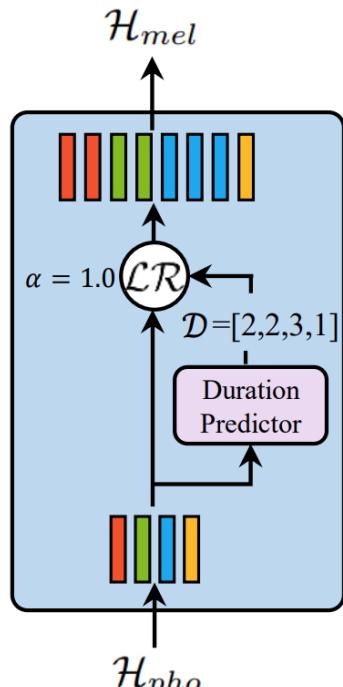
# Robust TTS—Attention improvement

- Multi-frame prediction [382]
  - Predicting multiple, non-overlapping output frames at each decoder step
  - Increase convergence speed, with a much faster (and more stable) alignment learned from attention
- Decoder prenet dropout/bottleneck [382,39]
  - 0.5 dropout, small hidden size as bottleneck

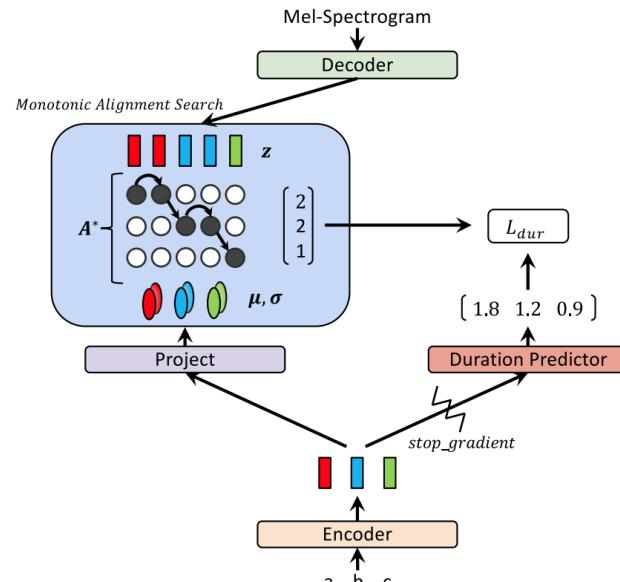


# Robust TTS——Duration Prediction

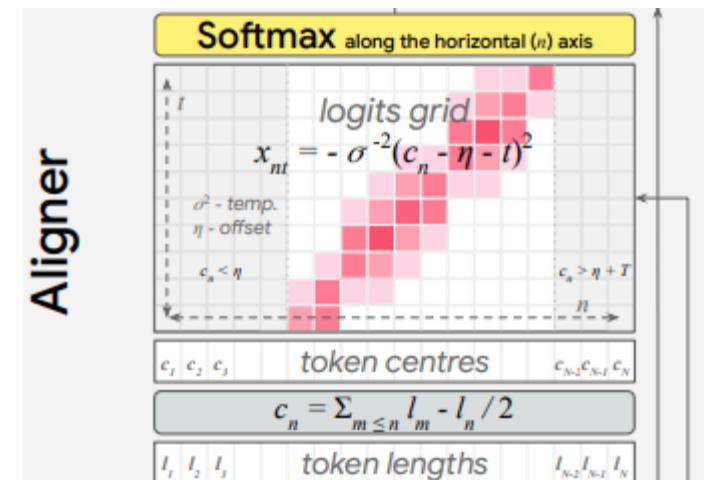
- Duration prediction and expansion
  - SPSS → Seq2Seq model with attention → Non-autoregressive model
  - Duration → attention, no duration → duration prediction (technique renaissance)



FastSpeech 1/2



Glow-TTS



EATS

# Robust TTS

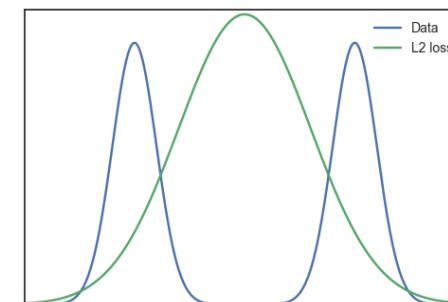
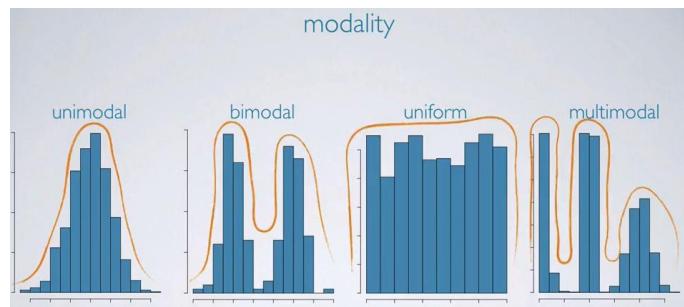
- A new taxonomy of TTS

AR?		AR	Non-AR
Attention?	Attention	Tacotron 2 [303], DeepVoice 3 [270]	ParaNet [268], Flow-TTS [234]
	Non-Attention	DurIAN [418], Non-Att Tacotron [304]	FastSpeech [290, 292], EATS [69]

# Expressive TTS

- Expressiveness
  - Characterized by content (what to say), speaker/timbre (who to say), prosody/emotion/style (how to say), noisy environment (where to say), etc
- Over-smoothing prediction
  - One to many mapping in text to speech:  $p(y|x)$  multimodal distribution

Text  
↓  
multiple speech variations  
(duration, pitch, sound volume, speaker, style, emotion, etc)



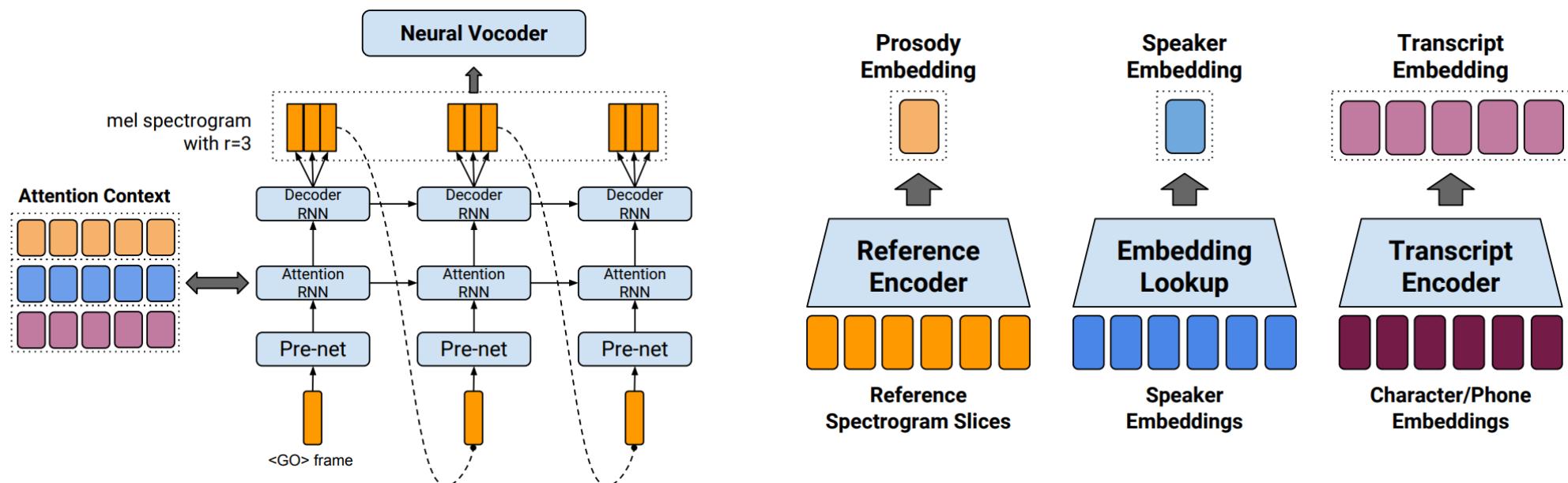
# Expressive TTS

- Modeling variation information

Perspective	Category	Description	Work
Information Type	Explicit	Language/Style/Speaker ID	[445, 247, 195, 162, 39]
		Pitch/Duration/Energy	[290, 292, 181, 158, 239, 365]
	Implicit	Reference encoder	[309, 383, 224, 142, 9, 49, 37, 40]
		VAE	[119, 4, 443, 120, 324, 325, 74]
		GAN/Flow/Diffusion	[224, 186, 366, 234, 159, 141]
		Text pre-training	[81, 104, 393, 143]
Information Granularity	Language/Speaker Level	Multi-lingual/speaker TTS	[445, 247, 39]
	Paragraph Level	Long-form reading	[11, 395, 376]
	Utterance Level	Timbre/Prosody/Noise	[309, 383, 142, 321, 207, 40]
	Word/Syllable Level		[325, 116, 45, 335]
	Character/Phoneme Level	Fine-grained information	[188, 324, 430, 325, 45, 40, 189]
	Frame Level		[188, 158, 49, 434]

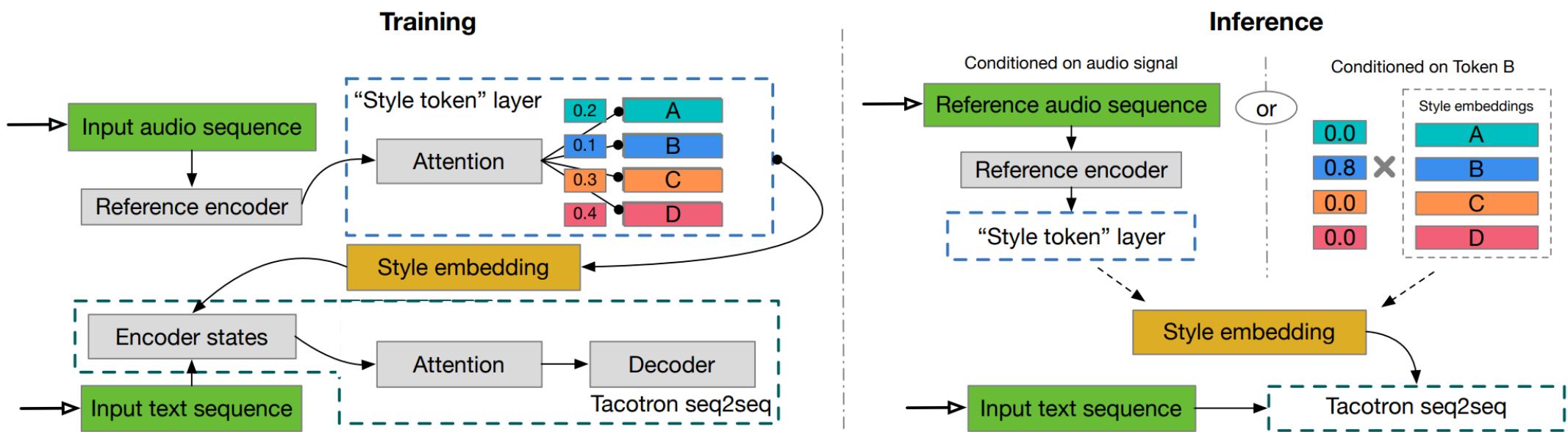
# Expressive TTS——Reference encoder

- Prosody embedding from reference audio [309]



# Expressive TTS——Reference encoder

- Style tokens [383]
  - Training: attend to style tokens
  - Inference: attend to style tokens or simply pick style tokens



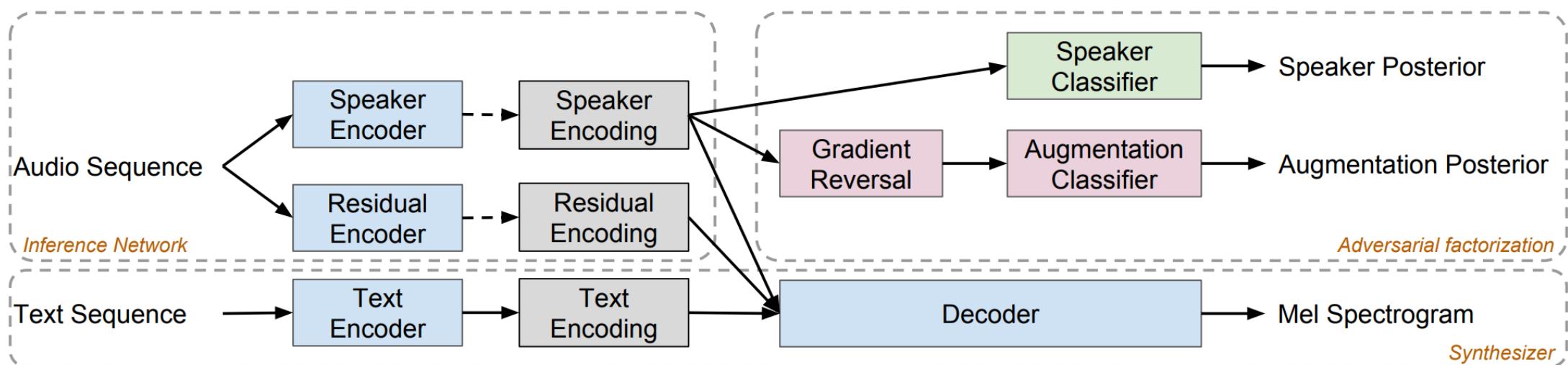
# Expressive TTS—Disentangling, Controlling and Transferring

- Disentangling
  - Content/speaker/style/noise, e.g., adversarial training
- Controlling
  - Cycle consistency/feedback loss, semi-supervised learning for control
- Transferring
  - Changing variance information for transfer

Technique	Description	Work
Disentangling with Adversarial Training	Disentanglement for control	[224, 120, 281, 434]
Cycle Consistency/Feedback for Control	Enhance style/timbre generation	[202, 386, 207, 30, 195]
Semi-Supervised Learning for Control	Use VAE and adversarial training	[103, 119, 120, 434, 302]
Changing Variance Information for Transfer	Different information in inference	[309, 383, 142, 443, 40]

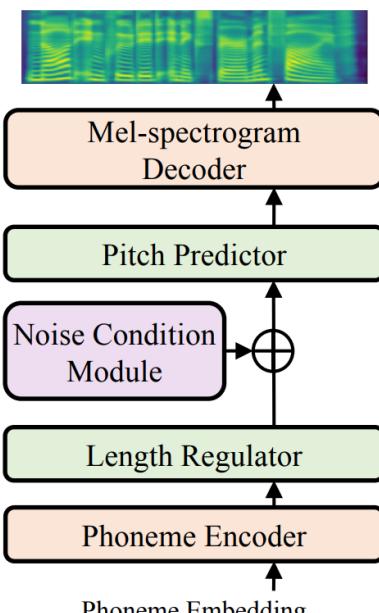
# Expressive TTS—Disentangling, Controlling and Transferring

- Disentangling correlated speaker and noise [120]
  - Synthesize clean speech for noisy speakers

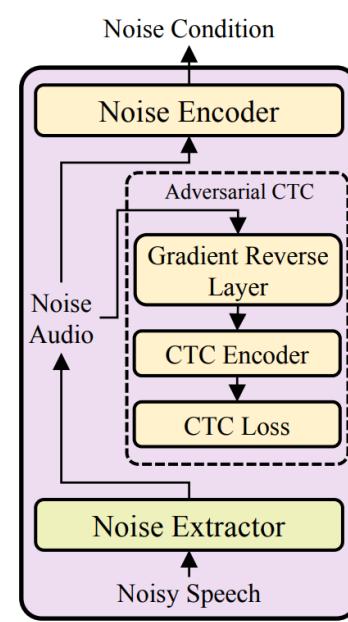


# Expressive TTS—Disentangling, Controlling and Transferring

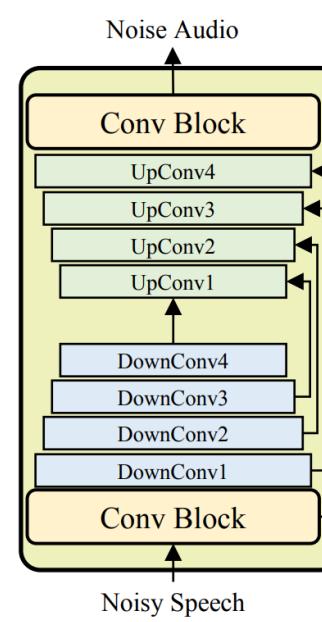
- Disentangling correlated speaker and noise with frame-level modeling [434]
  - Synthesize clean speech for noisy speakers



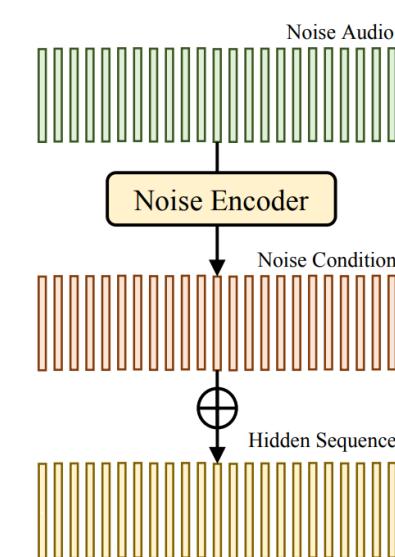
(a) DenoSpeech



(b) Noise Condition Module



(c) Noise Extractor



(d) Noise Encoder

# Adaptive TTS

- Voice adaptation, voice cloning, custom voice
- Empower TTS for everyone
  - Pre-training on multi-speaker TTS model
  - Fine-tuning on speech data from target speaker
  - Inference speech for target speaker
- Challenges
  - To support diverse customers, the source model needs to be generalizable enough, the target speech may be diverse (different acoustics/styles/languages)
  - To support many customers, the adaptation needs to be data and parameter efficient

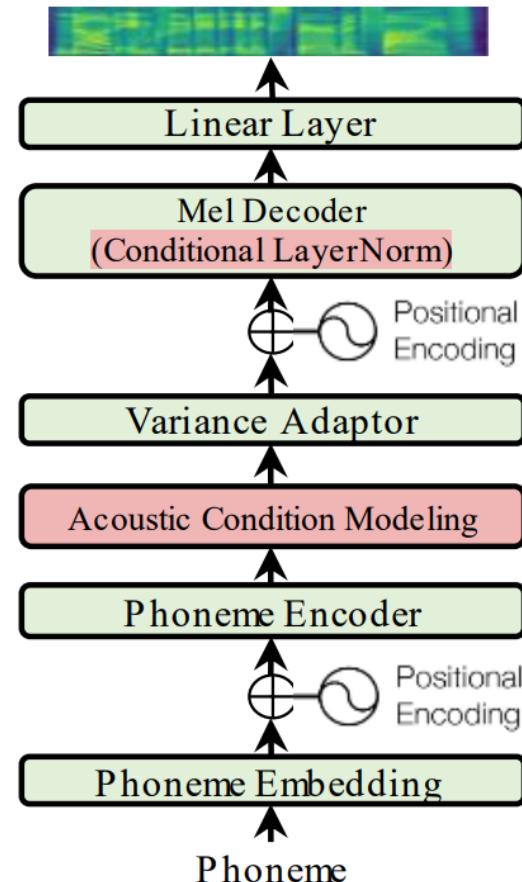
# Adaptive TTS

- A taxonomy on adaptive TTS

Category	Topic	Work
General Adaptation	Modeling Variation Information	[40]
	Increasing Data Coverage	[57, 407]
	Cross-Acoustic Adaptation	[40, 54]
Efficient Adaptation	Cross-Style Adaptation	[404, 266, 123]
	Cross-Lingual Adaptation	[445, 38, 212]
	Few-Data Adaptation	[44, 9, 177, 240, 446, 49, 40, 236]
	Untranscribed Data Adaptation	[403, 133, 221]
Few-Parameter Adaptation	Few-Parameter Adaptation	[9, 44, 40]
	Zero-Shot Adaptation	[9, 44, 142, 56]

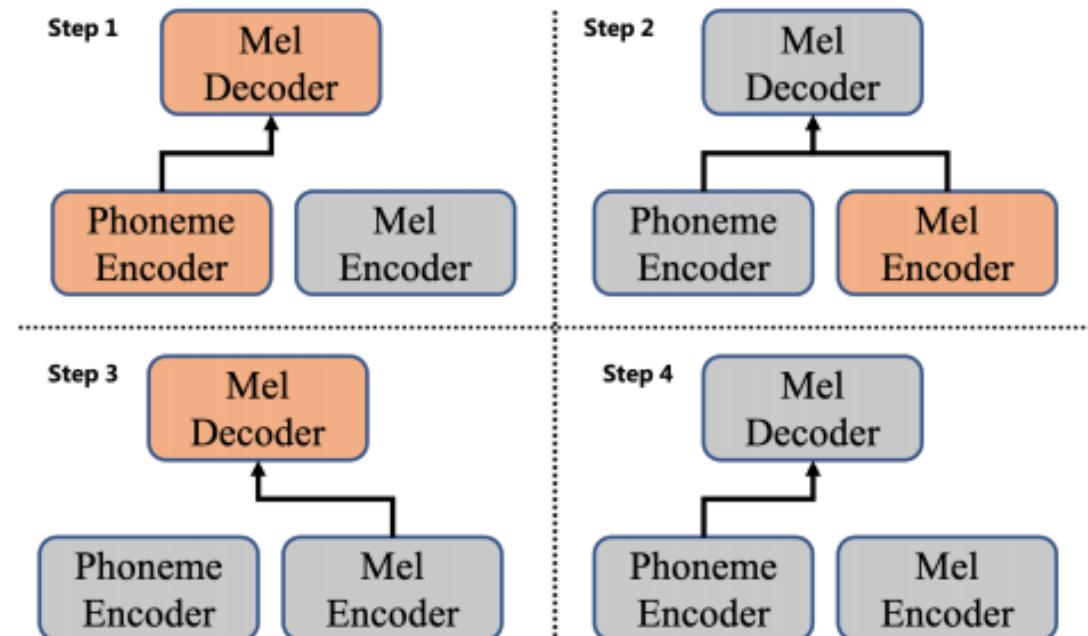
# Adaptive TTS——AdaSpeech [40]

- AdaSpeech
  - Acoustic condition modeling
    - Model diverse acoustic conditions at speaker/utterance/phoneme level
    - Support diverse conditions in target speaker
  - Conditional layer normalization
    - To fine-tune as small parameters as possible while ensuring the adaptation quality



# Adaptive TTS——AdaSpeech 2 [403]

- Only untranscribed data, how to adapt?
  - In online meeting, only speech can be collected, without corresponding transcripts
- AdaSpeech 2, speech reconstruction with latent alignment
  - Step 1: source TTS model training
  - Step 2: speech reconstruction
  - Step 3: speaker adaptation
  - Step 4: inference

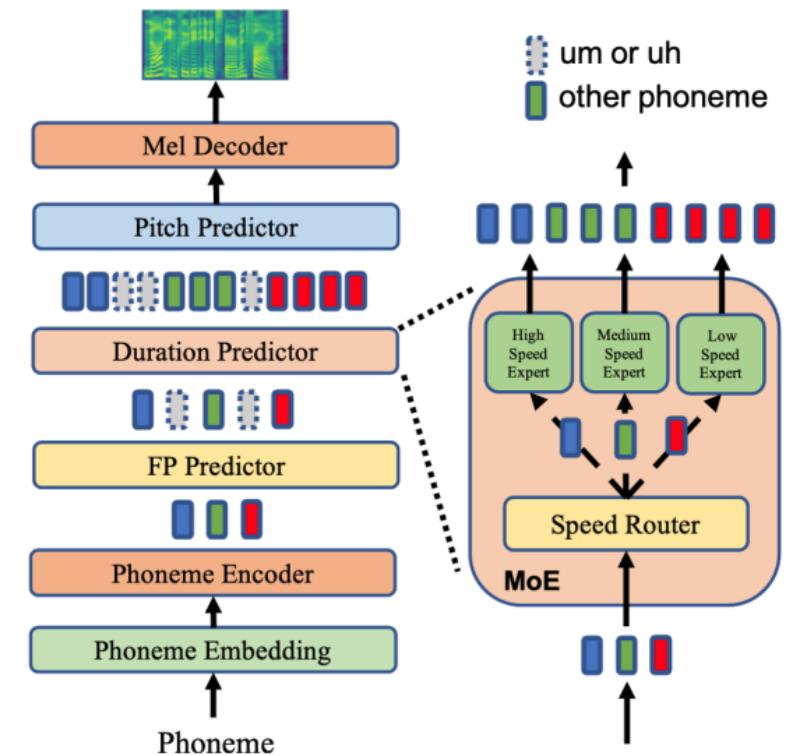


# Adaptive TTS——AdaSpeech 3 [404]

- Spontaneous style
  - Current TTS voices mostly focus on reading style.
  - Spontaneous-style voice is useful for scenarios like podcast, conversation, etc.
- AdaSpeech 3
  - Construct spontaneous dataset
  - Modeling filled pauses (FP, um and uh) and diverse rhythms

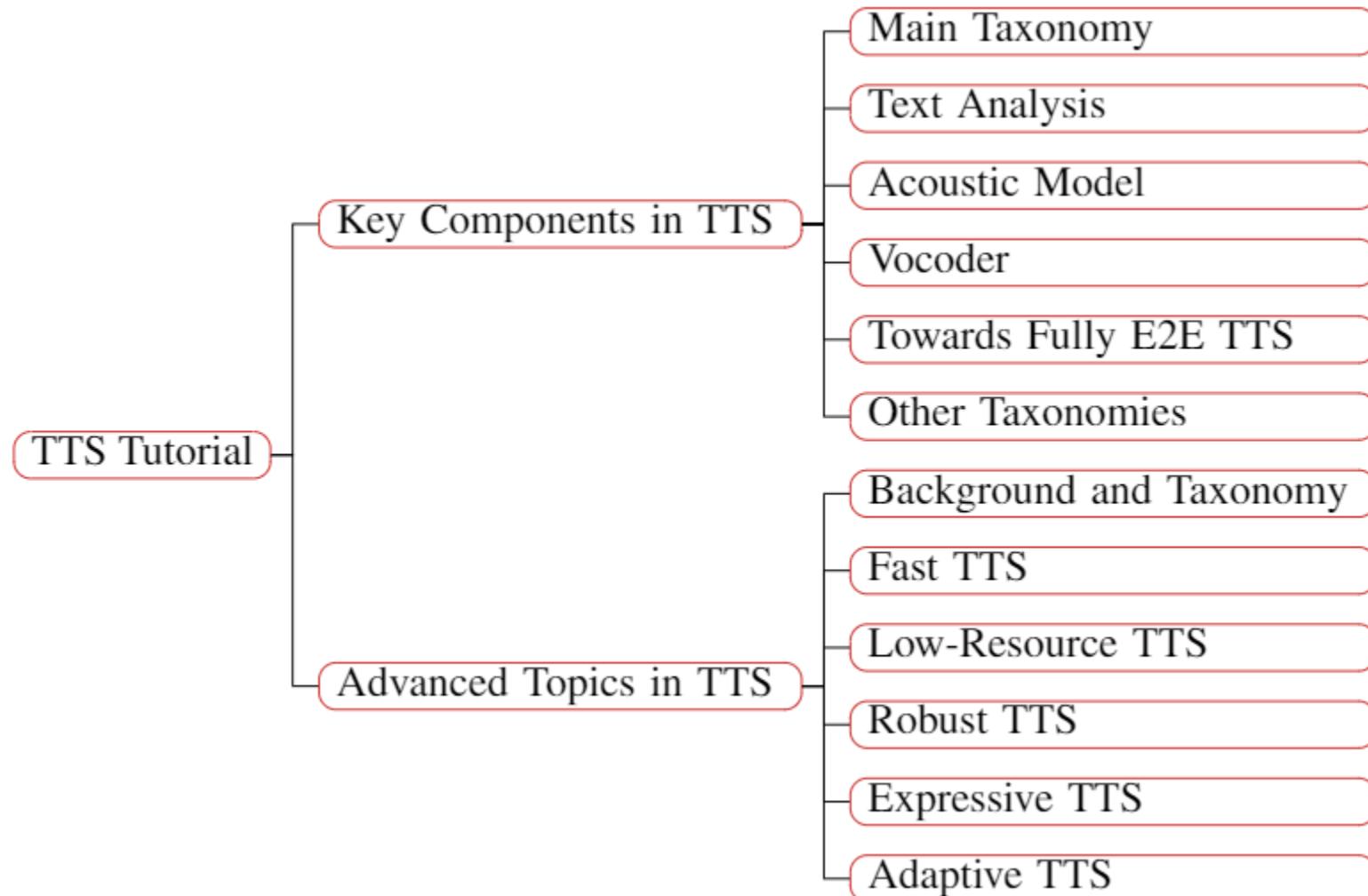


*Cecily package in all of that **um yeah** so ...*



## Part 4: Summary and Future Directions

# Summary



# Outlook: higher-quality synthesis

- Powerful generative models
- Better representation learning
- Robust speech synthesis
- Expressive/controllable/transferrable speech synthesis
- More human-like speech synthesis

# Outlook: more efficient synthesis

- Data-efficient TTS
- Parameter-efficient TTS
- Energy-efficient TTS

# We're hiring!

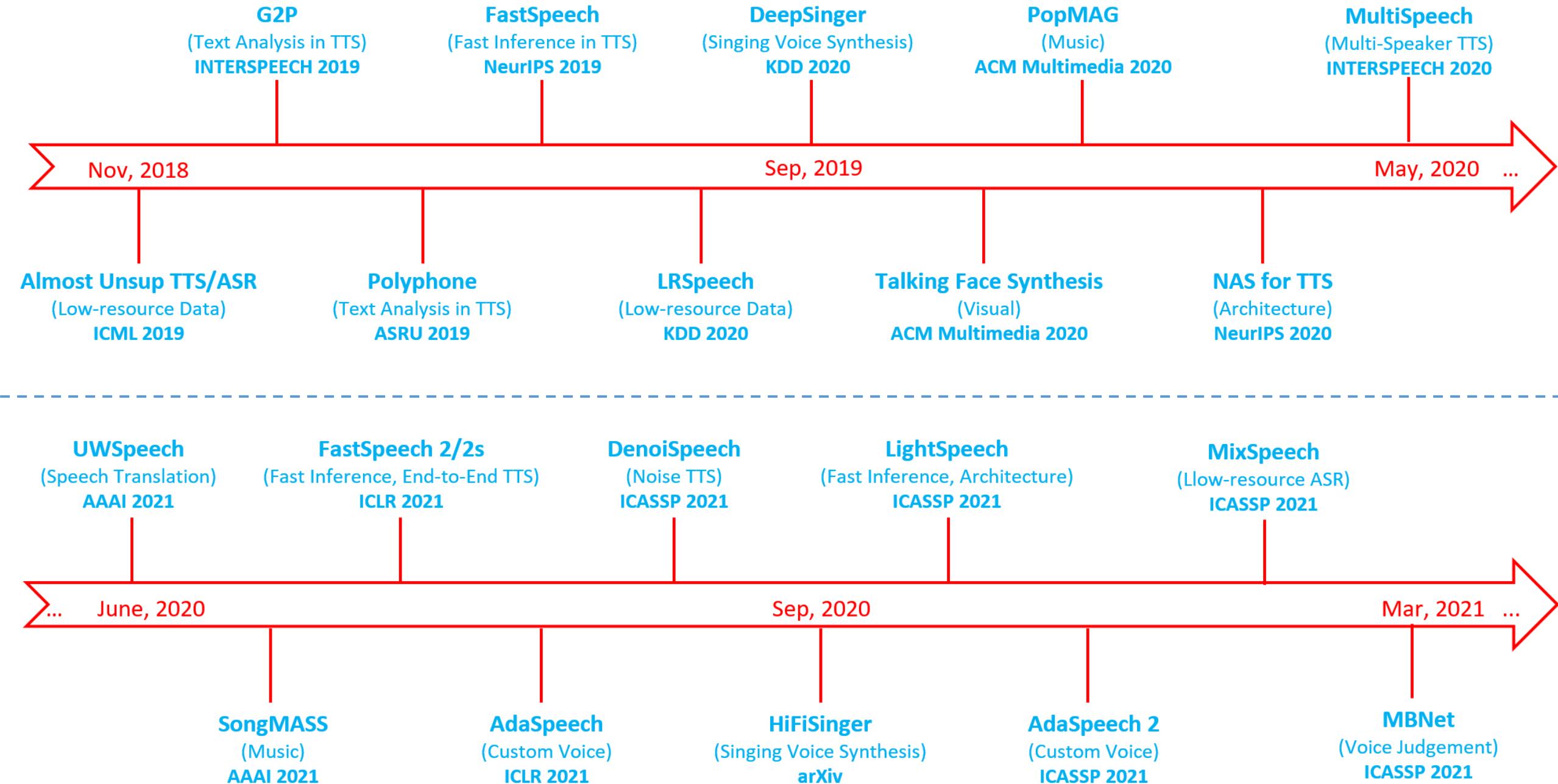
If you are passionate about machine learning research, especially **deep learning** and **reinforcement learning**, welcome to join us!!

Contact: [taoqin@Microsoft.com](mailto:taoqin@Microsoft.com)  
<http://research.Microsoft.com/~taoqin>

# Thank You!

<https://speechresearch.github.io/>

# Our research on speech



# Reference

See the reference in:

A Survey on Neural Speech Synthesis

<https://arxiv.org/pdf/2106.15561v3.pdf>