

TABLE – SotA: TTS (Text – To- Speech).

ID	Features	Databases	Classifier	Methodology Details	Best Results	Link
1	<ul style="list-style-type: none"><li>• Mel-spectrograms</li><li>• Phoneme-embeddings</li><li>• Positional Encoding</li></ul>	<ul style="list-style-type: none"><li>• LJSpeech</li></ul>	<ul style="list-style-type: none"><li>• Feed-forward Transformer consists of Phoneme Embeddings + Positional Enmcodings, 6xFFT Block, Length Regulator with Duration Predictor (2xConv1D + 1xLinear Layer), 6xFFT Block, Linear Layer</li></ul>	<ul style="list-style-type: none"><li>• first train the autoregressive Transformer TTS model on 4 NVIDIA V100 GPUs</li><li>• batch_size=16, Adam with <math>\beta_1=0.9</math>, <math>\beta_1=0.98</math>, <math>\epsilon=10^{-9}</math></li><li>• train the duration predictor -&gt; feed the text and speech pairs in the training set to the model again to obtain the encoder-decoder attention alignments</li><li>• source text sequence + generated mel-spectrograms with the autoregressive Transformer TTS model =&gt; the paired data for FastSpeech model training</li><li>• train the FastSpeech model together with the duration prediction</li><li>• output mel-spectrograms are transformed into audio samples using the pretrained WaveGlow</li></ul>	<ul style="list-style-type: none"><li>• MOS = <math>3.84 \pm 0.08</math></li><li>• Latency = <math>0.18 \pm 0.078s</math></li></ul>	<a href="#">ART 2019</a>
2	<ul style="list-style-type: none"><li>• Raw audio waveform</li></ul>	<ul style="list-style-type: none"><li>• VCTK</li></ul>	<ul style="list-style-type: none"><li>• Causal Conv + Residual Blocks (Dilated Conv + Tanh x Dilated Conv + Sigmoid ) + 2x Linear + Softmax</li></ul>	<ul style="list-style-type: none"><li>• Causal conv keeps the order of the samples</li><li>• Softmax layer has 65,536 probabilities corresponding to 16-bit integer representation for a sample</li><li>• Condition the model not in text, but speaker</li><li>• Receptive field: 240ms</li></ul>	<ul style="list-style-type: none"><li>• MOS = <math>4.21 \pm 0.081</math> (American English)</li><li>• MOS = <math>4.08 \pm 0.085</math> (Mandarin Chinse)</li></ul>	<a href="#">ART 2016</a>
3	<ul style="list-style-type: none"><li>• 80-dimensional log-mel filter bank coefficients</li><li>• Phoneme-embeddings</li></ul>	<ul style="list-style-type: none"><li>• LJ Speech</li><li>• VCTK</li></ul>	<ul style="list-style-type: none"><li>• VAE (Posterior Encoder – residual blocks from WaveGlow + Prior Encore – transformer encoder + Decoder – HiFi GAN v1 Generator + Stochastic Duration Predictor - residual blocks with dilated and depth-separable convolutional layers)</li></ul>	<ul style="list-style-type: none"><li>• AdamW optimizer: <math>\beta_1 = 0.8</math>, <math>\beta_2 = 0.99</math>, weight decay <math>\lambda = 0.01</math></li><li>• Lr = <math>2 \times 10^{-4}</math></li><li>• Batch_size=64</li><li>• 800k steps</li></ul>	<ul style="list-style-type: none"><li>• MOS = <math>4.43 (\pm 0.06)</math> LJSpeech</li><li>• MOS = <math>4.38 (\pm 0.06)</math> VCTK</li></ul>	<a href="#">ART 2021</a>
4	<ul style="list-style-type: none"><li>• Phoneme-embeddings</li><li>• High-dimensional speech compressed representation.</li></ul>	<ul style="list-style-type: none"><li>• LJSpeech</li><li>• A large-scale text corpus with 200 million sentences for phoneme pre-training</li></ul>	<ul style="list-style-type: none"><li>• Phoneme Encoder: 6-layer Feed-Forward Transformer (FFT) blocks+ Duration predictor with upsampling layer + Posterior Encoder (16-layer WaveNet) + Bidirectional Prior and Posterior Module + Waveform Decoder(Residual convolution blocks with upsampling).</li></ul>	<ul style="list-style-type: none"><li>• Pre-train phoneme embedding model.</li><li>• Memory-based VAE to simplify the posterior.</li><li>• Bidirectional prior/posterior modeling for enhancing prior and reducing the complexity of posterior.</li><li>• Trained on 8 8 NVIDIA V100 GPUs</li><li>• Optimizer: AdamW</li><li>• lr=20e-42, <math>\gamma=0.999875</math></li><li>• <math>\beta_1=0.8</math>, <math>\beta_2=0.99</math></li></ul>	<ul style="list-style-type: none"><li>• MOS = <math>4.56 \pm 0.13</math></li></ul>	<a href="#">ART 2022</a>
5	<ul style="list-style-type: none"><li>• Phoneme-embeddings</li><li>• Mel-spectrograms</li><li>• Scalar Quantization Codec</li></ul>	<ul style="list-style-type: none"><li>• Multilingual LibriSpeech (MLS)</li><li>• WenetSpeech</li></ul>	<ul style="list-style-type: none"><li>• Text Encoder(T5 arihitecture) + Speaker Encoder + SQ-Codec(consists of an encoder, decoder, and scalar quantization for compact latent representations) + Transformer Diffusion Model + Dense Layer</li></ul>	<ul style="list-style-type: none"><li>• lr=10e-41 for transformer diffusion, with a cosine scheduler and warmup of 1k steps.</li><li>• optimizer: Adam optimizer used for both SQ-Codec and transformer diffusion model.</li><li>• Diffusion steps=25</li><li>• Sentence duration predictors are used to control the length of generated speech.</li></ul>	<ul style="list-style-type: none"><li>• MOS = <math>4.06 \pm .052</math></li></ul>	<a href="#">ART 2024</a>

6	<ul style="list-style-type: none"> <li>speaker latent vectors from the reference mel spectrogram.</li> <li>local frame-level features.</li> </ul>	<ul style="list-style-type: none"> <li>LibriTTS</li> <li>AiShell3</li> </ul>	<ul style="list-style-type: none"> <li>FastSpeech 2 Integration(Encoder - phoneme encoder integrated with TSCM to incorporate speaker control + Variance Adaptor - adds duration, pitch, and energy information + Mel-spectrogram Decoder - utilizes TSCM for speaker-specific adaptations) + VITS Integration(Text Encoder, Duration Predictor, Generator: Integrated with TSCM to improve speaker control).</li> </ul>	<ul style="list-style-type: none"> <li>batch_size=50</li> <li>Training Steps=250,000</li> <li>Optimizer: Adam</li> <li>Pretrained HiFi-GAN vocoder used to convert mel spectrograms to audio.</li> </ul>	<ul style="list-style-type: none"> <li>TSCM-FastSpeech 2: MOS = 3.81</li> <li>TSCM-VITS: MOS = 4.44</li> </ul>	<a href="#">ART</a> 2024
7	<ul style="list-style-type: none"> <li>Mel-spectrograms</li> <li>Spectrogram Tokens (Vector Quantized)</li> <li>Phoneme Embeddings</li> </ul>	<ul style="list-style-type: none"> <li>LJSpeech</li> </ul>	<ul style="list-style-type: none"> <li>Spectrogram VQ Model(consists of an encoder, decoder, and discrete codebook) + Text Encoder(EfficientSpeech encoder with 2 transformer blocks) + Discrete Diffusion Model (12-layer transformer with 8 heads) with Contrastive Learning.</li> </ul>	<ul style="list-style-type: none"> <li>Compresses mel-spectrograms into discrete tokens to reduce computational costs.</li> <li>Diffusion Model Training: Uses a forward and reverse diffusion process with a Mask-and-Replace strategy.</li> <li>Contrastive Learning: Text-wise Contrastive Learning Loss added to improve text-spectrogram alignment.</li> <li>Inference: Generated mel-spectrograms are converted into audio using the Griffin-Lim algorithm.</li> </ul>	<ul style="list-style-type: none"> <li>MOS: <math>3.64 \pm 0.05</math></li> <li>mRTF (Real-Time Factor): 73.9 (GPU), 17.6 (CPU)</li> </ul>	<a href="#">ART</a> 2023
8	<ul style="list-style-type: none"> <li>Mel-spectrograms</li> <li>Phoneme Representations</li> <li>Dynamic Quantized Representation</li> </ul>	<ul style="list-style-type: none"> <li>LJSpeech</li> </ul>	<ul style="list-style-type: none"> <li>Sequential Autoencoder (Encoder - convolution blocks and LSTMs + Dynamic Codebook Module + Decoder: Tacotron 2-based)</li> </ul>	<ul style="list-style-type: none"> <li>Dynamic Quantized Representation Learning: Quantization with a dynamic codebook that expands based on unpaired data using pseudo-labels generated by a pre-trained ASR.</li> <li>Train with a mix of 120 minutes of paired data and 600 minutes of unpaired data.</li> <li>batch_Size=64</li> <li>Optimizer: Adam (<math>\beta_1=0.9</math>, <math>\beta_2=0.999</math>, <math>\text{lr}=10\text{e-}3</math>).</li> <li>Use Connectionist Temporal Classification (CTC) for recognition loss.</li> </ul>	<ul style="list-style-type: none"> <li>MOS: <math>3.12 \pm 0.32</math> with mixed data (120 min paired + 600 min unpaired).</li> </ul>	<a href="#">ART</a> 2024
9	<ul style="list-style-type: none"> <li>Mel-spectrograms</li> <li>Phoneme embeddings</li> <li>Pitch, duration, and energy features</li> </ul>	<ul style="list-style-type: none"> <li>LJSpeech</li> <li>VCTK</li> <li>LibriTTS</li> </ul>	<ul style="list-style-type: none"> <li>Phoneme Encoder(Transformer blocks) + Variance Adaptor + CM-Decoder (non-causal WaveNet-like structure) + Vocoder (HiFi-GAN).</li> </ul>	<ul style="list-style-type: none"> <li>CM-TTS employs a consistency model-based approach for real-time mel-spectrogram generation.</li> <li>Utilizes weighted samplers to improve model training by incorporating dynamic probabilities.</li> <li>Model trained for 300K steps with exponential learning rate decay and a batch size of 32.</li> </ul>	<ul style="list-style-type: none"> <li>MOS: <math>3.9618 \pm 0.0186</math></li> <li>Latency: Real-time capability with fewer synthesis steps (1, 2, 4 steps tested).</li> </ul>	<a href="#">ART</a> 2024
10	<ul style="list-style-type: none"> <li>Prosody features (pitch, duration)</li> <li>Phoneme embeddings</li> </ul>	<ul style="list-style-type: none"> <li>LibriSpeech:</li> <li>MLS corpus</li> </ul>	<ul style="list-style-type: none"> <li>Phoneme Encoder + Codec Decoder (SoundStream) + VALL-E (12-layer transformers)</li> </ul>	<ul style="list-style-type: none"> <li>Prosody Tokens are predicted using a Chain-of-Thought (CoT) prompting technique, stabilizing pitch and duration before speech token prediction</li> <li>Trained on 8 NVIDIA V100 and 16 AMD MI200 GPUs with Adam optimizer</li> <li>Utilized nucleus sampling for phoneme, pitch, and duration prediction</li> <li>Tested multiple window sizes for duration-guided masking, with optimal size being 1 for WER improvements</li> <li>Evaluation on hard sentences to assess robustness, with error types classified as mispronunciation, omission, repetition, and hallucination</li> </ul>	<ul style="list-style-type: none"> <li>Word Error Rate (WER): 2.5% (RALL-E), 5.6% (VALL-E) on standard LibriSpeech test-clean set</li> <li>MOS:nUTMOS = 4.00 for RALL-E,</li> </ul>	<a href="#">ART</a> 2024
11	<ul style="list-style-type: none"> <li>Mel-spectrograms</li> <li>Vector Quantized Variational codes</li> </ul>	<ul style="list-style-type: none"> <li>XTTS</li> <li>LibriTTS-R</li> <li>Common Voice</li> </ul>	<ul style="list-style-type: none"> <li>VQ-VAE (13M parameters) encodes mel-spectrograms to 1024 codebook vectors + GPT-2 encoder (443M parameters) predicts VQ-VAE audio codes from text input + Conditioning Encoder (6 layers) for generating speaker embeddings, producing 32 embeddings per audio + HiFi-GAN based decoder (26M parameters) reconstructs audio from latent vectors</li> </ul>	<ul style="list-style-type: none"> <li>The model was trained on multilingual datasets using a language batch balancer</li> <li>XTTS was trained for approximately 2.5M steps on 4 NVIDIA A100 GPUs (80GB)</li> <li>AdamW optimizer, betas 0.9 and 0.96, with MultiStepLR learning rate scheduler</li> <li>XTTS improves speaker cloning capability by conditioning the encoder on multiple embeddings, rather than a single embedding</li> </ul>	<ul style="list-style-type: none"> <li>English Evaluation: CER: 0.5425 UTMOS: <math>4.007 \pm 0.25</math> SECS: 0.6423</li> </ul>	<a href="#">ART</a> 2024

12	<ul style="list-style-type: none"><li>• Phoneme monotonic alignment</li><li>• Merged codec with reduced sampling rate</li><li>• Phoneme prediction during training</li></ul>	<ul style="list-style-type: none"><li>• LibriSpeech</li></ul>	<ul style="list-style-type: none"><li>• Encoder-decoder (convolution-based encoder) + Residual Vector Quantizer module (8-layer) +</li><li>• Transformer-based architecture (12-layer) + Vocoder</li></ul>	<ul style="list-style-type: none"><li>• Two-stage training: autoregressive model predicts acoustic tokens from phonemes and aligned phoneme sequences, and NAR model iteratively generates tokens for higher layers.</li><li>• The merged codec reduces the number of autoregressive steps by downsampling in the first layer without retraining the codec.</li><li>• Monotonic alignment ensures that phoneme and acoustic tokens align, improving robustness by preventing repetition or skipping.</li><li>• Experiments used 3-second acoustic prompts and phoneme sequences for zero-shot TTS tasks like speech continuation and cross-sentence synthesis.</li></ul>	<ul style="list-style-type: none"><li>• QMOS = 4.02, SMOS = 3.89</li><li>• Latency: achieved a generation time of 3.67s for 10s of speech</li></ul>	<a href="#">ART 2024</a>
13	<ul style="list-style-type: none"><li>• Mel-spectrograms</li><li>• Time-Invariant and Time-Variant Style Representations</li></ul>	<ul style="list-style-type: none"><li>• VCTK</li><li>• Emotional Speech Dataset.</li></ul>	<ul style="list-style-type: none"><li>• Text Encoder(8 layers of Transformer encoders with relative positional embedding and adaptive layer normalization) + Aligner(Convolution-based Duration Predictor) + Diffusion Decoder (convolution blocks and DiT blocks) + Time-Invariant and Time-Variant Encoders (uses cross-attention).</li></ul>	<ul style="list-style-type: none"><li>• Batch_size=32</li><li>• Optimizer: Adam (lr=10e-4)</li><li>• 1000 epochs for VCTK, 1500 epochs for ESD</li><li>• Diffusion Process: Incorporates Gaussian noise into input data and iteratively refines to generate mel-spectrograms.</li></ul>	<ul style="list-style-type: none"><li>• VCTK: Seen Speakers: 3.75 (MOS-N), 3.88 (MOS-S) Unseen Speakers: 3.76 (MOS-N), 3.81 (MOS-S)</li><li>• ESD: Seen Speakers: 3.73 (MOS-N), 3.84 (MOS-S) Unseen Speakers: 3.57 (MOS-N), 3.52 (MOS-S)</li></ul>	<a href="#">ART 2024</a>
14	<ul style="list-style-type: none"><li>• Mel-spectrograms</li><li>• Phoneme-embeddings</li><li>• Scalar Quantization Codec</li></ul>	<ul style="list-style-type: none"><li>• Multilingual LibriSpeech</li><li>• WenetSpeech</li></ul>	<ul style="list-style-type: none"><li>• Text Encoder (ByT5 model with 2 transformer blocks) + Speaker Encoder (Pre-trained FAcoder) + SQ-Codec (Encoder, decoder, and scalar quantization) + Transformer Diffusion Model (flow-based scalar latent transformer diffusion with 12 layers and 8 heads).</li></ul>	<ul style="list-style-type: none"><li>• Batch_size =32</li><li>• Optimizer: Adam (lr=10e-4)</li><li>• 400,000 training steps</li><li>• Diffusion steps: 25</li><li>• Sentence Duration Prediction:4 strategies explored: ByT5-based, ChatGPT-based, FS2-based, and AR-based duration predictors.</li><li>• Generated mel-spectrograms converted to audio using SQ-Codec decoder.</li></ul>	<ul style="list-style-type: none"><li>• MOS: 4.28 ± 0.12</li></ul>	<a href="#">ART 2024</a>
15	<ul style="list-style-type: none"><li>• Mel-spectrograms</li><li>• Style Embeddings (text and audio)</li><li>• Speaker Embeddings</li></ul>	<ul style="list-style-type: none"><li>• Emotional Speech Dataset</li></ul>	<ul style="list-style-type: none"><li>• General Style Fusion Encoder (CLIP-based text encoder and audio encoder) + Hierarchical Conformer Two-Branch Style Control Module (fuses style and speaker control embeddings into the VITS-based TTS architecture for optimal control of both speaker and style) + Backbone (VITS).</li></ul>	<ul style="list-style-type: none"><li>• 1,000,000 training steps</li><li>• Multimodal Input Processing: combines style prompts and audio references to control style and speaker embeddings.</li><li>• Gradient Reversal Layer: used to disentangle speaker and style information.</li><li>• HiFi-GAN is used as a vocoder, with Speech Super Resolution upsampling for enhanced quality.</li></ul>	<ul style="list-style-type: none"><li>• Speaker-MOS: 4.19</li><li>• Emotion-MOS: 4.28</li></ul>	<a href="#">ART 2024</a>
16	<ul style="list-style-type: none"><li>• Phoneme embeddings</li><li>• Pitch</li><li>• Duration</li><li>• Energy</li></ul>	<ul style="list-style-type: none"><li>• BEAT2</li></ul>	<ul style="list-style-type: none"><li>• Phoneme Encoder (Causal Transformer Encoder) + Rhythmic Predictors (Separate CNN-based predictors for pitch, duration, and energy) + Shared Rhythm Predictors + Speech Decoder (1D dilated convolutions) + Gesture Decoder (Pretrained VQ-VAE for gesture reconstruction using semantic and rhythmic latent features) + Neural Architecture Search</li></ul>	<ul style="list-style-type: none"><li>• Joint Generation: the model jointly generates speech and gestures by sharing intermediate rhythmic features (pitch, duration, and energy).</li><li>• Causal Network: redesigned to avoid dependencies on future inputs for real-time applications.</li><li>• Neural Architecture Search: used to fine-tune model components (e.g., layers and convolutional kernels) for performance improvement.</li><li>• Training on NVIDIA A100 GPUs, leveraging reinforcement learning-based NAS to optimize network hyperparameters.</li></ul>	<ul style="list-style-type: none"><li>• MOS: 3.93 (speech quality)</li><li>• Latency: 0.17 seconds per second (speech and gesture generation on NVIDIA 3090)</li></ul>	<a href="#">ART 2024</a>

17	<ul style="list-style-type: none"> <li>• Mel-spectrograms</li> <li>• Prosody (pitch, duration)</li> <li>• Phoneme embeddings</li> <li>• Style diffusion modeling (latent random variable for speech style)</li> </ul>	<ul style="list-style-type: none"> <li>• LJSpeech</li> <li>• VCTK</li> <li>• LibriTTS</li> </ul>	<ul style="list-style-type: none"> <li>• Text Encoder (Causal Transformer for phoneme representation) + Prosodic Text Encoder (BERT-based) + Style Encoder + Duration and Prosody Predictors + Waveform Decoder (iSTFTNet or HiFiGAN) + SLM Discriminators (WavLM-based adversarial training, with convolutional head)</li> </ul>	<ul style="list-style-type: none"> <li>• Training Strategy: two-stage process; first pre-train acoustic modules (100 epochs on LJSpeech), then jointly optimize all components with differentiable duration modeling.</li> <li>• Adversarial Training: utilizes large pre-trained Speech Language Models (SLMs) like WavLM as discriminators for human-like quality synthesis.</li> <li>• End-to-End Training: joint optimization of text encoder, style encoder, prosody predictor, and waveform decoder for direct waveform generation.</li> <li>• Diffusion-Based Sampling: style diffusion model samples a latent style vector conditioned on text, enabling diverse and expressive speech generation without reference audio.</li> <li>• Zero-Shot Speaker Adaptation: fine-tuned on LibriTTS for zero-shot speaker adaptation using only 3-second reference clips.</li> </ul>	<ul style="list-style-type: none"> <li>• MOS: <ul style="list-style-type: none"> <li>3.83 (LJSpeech, surpasses ground truth with CMOS +0.28)</li> <li>4.15 (LibriTTS, zero-shot)</li> <li>4.03 (Similarity score for zero-shot speaker adaptation on LibriTTS)</li> </ul> </li> <li>• Latency: 0.0185 Real-Time Factor</li> </ul>	<a href="#">ART 2024</a>
18	<ul style="list-style-type: none"> <li>• Mel-spectrograms</li> <li>• Latent space control for speech variation (pitch, tone, speech rate, cadence, accent)</li> </ul>	<ul style="list-style-type: none"> <li>• LJSpeech</li> <li>• LibriTTS</li> </ul>	<ul style="list-style-type: none"> <li>• Text Encoder (modified Tacotron encoder -instance-norm replaces batch-norm) + Affine Coupling Layer + Latent Space Control (invertible mapping between mel-spectrograms and latent space, modeled using a Gaussian Mixture Model) + Waveform Decoder (WaveGlow)</li> </ul>	<ul style="list-style-type: none"> <li>• Maximize the likelihood of training data to train the autoregressive flow model.</li> <li>• Train the model on the LSH dataset for 1,000 epochs and fine-tune for 500 epochs on LibriTTS.</li> <li>• Models trained on NVIDIA DGX-1 with 8 GPUs.</li> <li>• Variability Control: Adjust the amount of variation in speech output by sampling from a Gaussian prior with different variances (<math>\sigma^2 = 0.0, 0.5, 1.0</math>).</li> <li>• Posterior Sampling: Style transfer between seen and unseen speakers by sampling from a posterior distribution conditioned on prior evidence (e.g., expressive vs. monotonic speech styles).</li> <li>• Interpolation: Smooth interpolation between samples and speakers by manipulating latent space (z-space).</li> </ul>	<ul style="list-style-type: none"> <li>• MOS: <ul style="list-style-type: none"> <li>3.665 <math>\pm</math> 0.1634 (Flowtron)</li> <li>3.521 <math>\pm</math> 0.1721 (Tacotron 2)</li> <li>4.274 <math>\pm</math> 0.1340 (real data)</li> </ul> </li> </ul>	<a href="#">ART 2020</a>
19	<ul style="list-style-type: none"> <li>• Mel-spectrograms</li> <li>• Hierarchical latent variables</li> <li>• Phoneme embeddings</li> <li>• Positional encoding (sinusoidal)</li> </ul>	<ul style="list-style-type: none"> <li>• LJSpeech</li> <li>• Multi-speaker internal Mandarin Chinese corpus (55 hours, 7 female speakers)</li> </ul>	<ul style="list-style-type: none"> <li>• Very Deep Variational Autoencoder(residual blocks containing 4 convolution layers with GELU activation) + Residual Attention Mechanism + Text Encoder (Feature-wise Linear Modulation + 4 convolution layers and positional encodings) + Speaking Speed Predictor</li> </ul>	<ul style="list-style-type: none"> <li>• Batch size of 32, trained on 2 NVIDIA V100 GPUs for 90k iterations using the Adam optimizer (<math>\beta_1=0.9, \beta_2=0.999</math>), learning rate maxing at 1.5e-4 with 10k warm-up steps.</li> <li>• Combined loss including reconstruction loss, Kullback-Leibler divergence, and speaking speed prediction loss.</li> <li>• Inference is 16x faster than Tacotron 2 due to the non-autoregressive nature of the model.</li> </ul>	<ul style="list-style-type: none"> <li>• MOS on LJSpeech: 3.88 <math>\pm</math> 0.20</li> <li>• MOS on the multi-speaker Mandarin dataset: 4.49 <math>\pm</math> 0.11</li> <li>• VARA-TTS inference speed: 32.01ms</li> </ul>	



20	<ul style="list-style-type: none"><li>• Mel-spectrograms (80-dimensional)</li><li>• Speaker embeddings</li><li>• Phoneme-level input for text processing</li></ul>	<ul style="list-style-type: none"><li>• VCTK</li><li>• LibriTTS</li></ul>	<ul style="list-style-type: none"><li>• Speaker Encoder (ECAPA-TDNN) + Squeeze-and-Excitation blocks for channel interdependencies + Res2Net blocks with skip connections for feature aggregation + Improved pooling with channel/context-dependent frame attention + Acoustic Model (FastSpeech 2, non-autoregressive model) + Transformer-based encoder and decoder + Variant adaptor with duration, pitch, and energy predictors + Postnet (Conv1D blocks) added after the decoder for fine-tuning.+ Vocoder (HiFi-GAN)</li></ul>	<ul style="list-style-type: none"><li>• Data downsampled to 22,050 Hz for experiments</li><li>• Pretrained speaker encoder on VoxCeleb1 and VoxCeleb2 datasets</li><li>• Acoustic model trained for 400k steps with a batch size of 16 using an RTX 3090 GPU</li><li>• Ground-truth phoneme durations were obtained via Montreal Force Aligner</li><li>• Testing includes seen speakers (from VCTK) and unseen speakers (from LibriTTS and VCTK).</li><li>• Average speaker embedding calculated for each speaker to improve stability and speaker similarity.</li></ul>	<ul style="list-style-type: none"><li>• MOS-N:  Ground truth: ~4.19</li><li>• ECAPA-TDNN: 3.62 (seen), 3.47 (unseen LibriTTS)</li><li>• x-vector: 3.51 (seen), 3.38 (unseen LibriTTS)</li></ul>	<a href="#">ART 2022</a>
21	<ul style="list-style-type: none"><li>• Phoneme embeddings (raw phonetic input)</li><li>• Explicit modeling of prosody (pitch, duration)</li><li>• BERT-based contextualized word embeddings</li><li>• Hybrid grapheme-to-phoneme conversion with punctuation modeling</li></ul>	<ul style="list-style-type: none"><li>• French language dataset from Blizzard Challenge 2023</li><li>• Speaker 1 (NEB): 50 hours of aligned, high-quality audiobook data</li><li>• Speaker 2 (AD): 2 hours of aligned data</li></ul>	<ul style="list-style-type: none"><li>• Custom network for prosody prediction (handles duration and pitch) + CamemBERT for contextualized word embeddings + HiFi-GAN for vocoding + 3 parallel stacks of BiLSTMs for duration, pitch, and vocoder conditioning + Shared sub-network with 3 convolutional layers + BiLSTM for phonetic embeddings + HiFi-GAN conditioned on the predicted duration and pitch embeddings + Grapheme-to-Phoneme Module (Stack of 3 convolutional layers followed by BiLSTMs)</li></ul>	<ul style="list-style-type: none"><li>• Training the phonemizer with a split dataset (90% train, 10% validation)</li><li>• Early stopping with sentence accuracy rate</li><li>• Fine-tuned CamemBERT along with the custom prosody network and HiFiGAN</li><li>• BERT optimized with a fixed learning rate of 10e−6, 1M steps</li><li>• Phonetic alignment and pitch annotations used for forced alignment</li><li>• Two layers of non-uniform upsampling to match phonemes with BERT embeddings and duration predictions</li><li>• Training performed on an NVIDIA RTX 3090 with a batch size of 16 over 3 weeks</li></ul>	<ul style="list-style-type: none"><li>• MOS for Speaker NEB: Speech Experts: 4.0 (±1.48) Non-speech experts: 4.3 (±0.74)</li><li>• MOS for Speaker AD: Speech Experts: 3.3 (±1.00) Non-speech experts: 4.1 (±0.83)</li></ul>	<a href="#">ART 2023</a>
22	<ul style="list-style-type: none"><li>• Self-supervised speech representations</li><li>• Mel-spectrograms</li><li>• Phoneme embeddings</li></ul>	<ul style="list-style-type: none"><li>• VCTK</li><li>• LibriTTS</li></ul>	<ul style="list-style-type: none"><li>• Hierarchical Conditional Variational Autoencoder: Self-supervised speech representations from XLS-R (12th layer for linguistic information) + Text Encoder (Transformer with relative positional encoding) + Linguistic Encoder (Bi-directional WaveNet) + Acoustic Encoder (Non-causal WaveNet residual blocks for acoustic latent variables) + HiFi-GAN vocoder + Flow-based Stochastic Duration Predictor + Monotonic Alignment Search</li></ul>	<ul style="list-style-type: none"><li>• Training with AdamW optimizer (<math>\beta_1 = 0.8</math>, <math>\beta_2 = 0.99</math>, weight decay = 0.01)</li><li>• HierSpeech trained on 4 NVIDIA A100 GPUs for 600k steps, batch size 256</li><li>• Untranscribed speech training (HierSpeech-U): Speaker adaptation without text transcripts using a style encoder</li><li>• Evaluation using fine-tuned wav2vec 2.0 for phoneme and word error rates (PER, WER)</li></ul>	<ul style="list-style-type: none"><li>• VCTK: MOS = 4.04 (N), 3.22 (S)</li><li>• LibriTTS: MOS = 3.98 (N), 3.26 (S)</li><li>• Untranscribed Speech: MOS (naturalness): 4.08</li></ul>	<a href="#">ART 2022</a>