



Towards End-to-End Prosody Transfer for Expressive Speech Synthesis with Tacotron



RJ Skerry-Ryan, Eric Battenberg, Ying Xiao, Yuxuan Wang, Daisy Stanton, Joel Shor, Ron J. Weiss, Rob Clark, Rif A. Saurous
Google AI

Prosody in Speech



What's Prosody?

- Intonation, rhythm, pitch, stress, loudness.
- Conveys emotion, emphasis, and additional meaning.
- Aids understanding.

How should we say this text? It depends.

- The **cat sat** on the **mat**. 🎧
- End-to-end TTS sounds **pretty good**. 🎧

Our working definition:

Definition. Prosody is the variation in speech signals that remains after accounting for variation due to phonetics, speaker identity, and channel effects (i.e. the recording environment).

Prosody Transfer

How to Control Prosody:

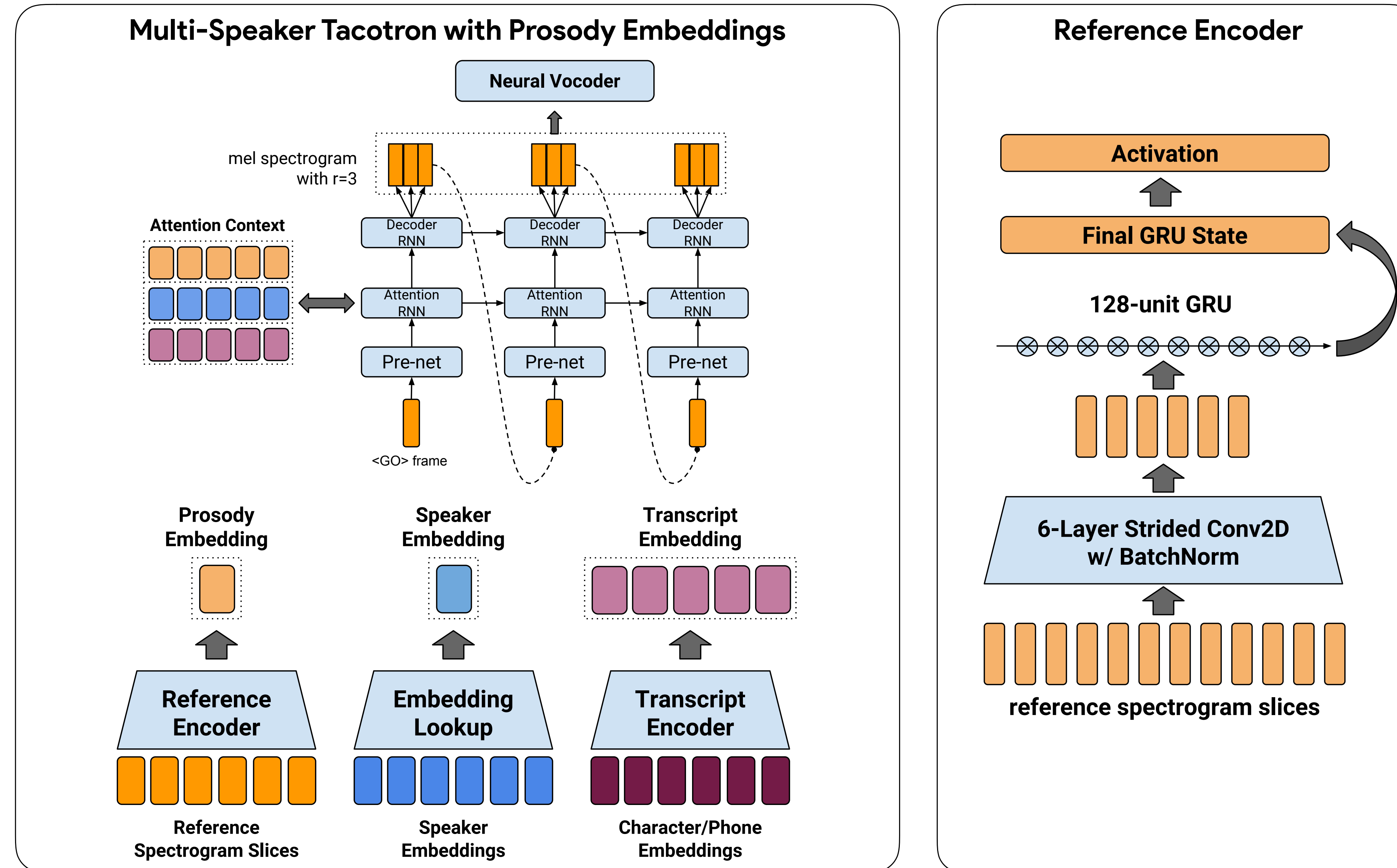
- Prosody annotations (e.g., ToBI)
- Phoneme-wise pitch, energy, duration.
- “Say it like this” (prosody transfer)

Prosody Transfer Desiderata:

- Pitch relative transfer (output within a speaker's natural range).
- Robust to text modifications (makes it scalable).

Tacotron Architecture

Specs: Phoneme inputs, GMM attention, CBHG transcript encoder, Griffin-Lim or WaveNet sample generation.



Quantitative / Subjective Results

Quantitative Similarity Measures

- **Mel Cepstral Distortion (MCD_K):** Mean squared error over first K MFCCs.

$$\frac{1}{T} \sum_{t=0}^{T-1} \sqrt{\sum_{k=1}^K (c_{t,k} - c'_{t,k})^2}$$

- **F0 Frame Error (FFE):** Percentage of frames with either a >20% pitch error or a voicing decision error.

$$\frac{\sum_{t=0}^{T-1} 1[|p_t - p'_t| > 0.2p_t] 1[v_t] 1[v'_t] + 1[v_t \neq v'_t]}{T}$$

Subjective Similarity Measure

- **Subjective:** Anchored side-by-side prosody similarity comparisons on a scale of [-3 to 3]

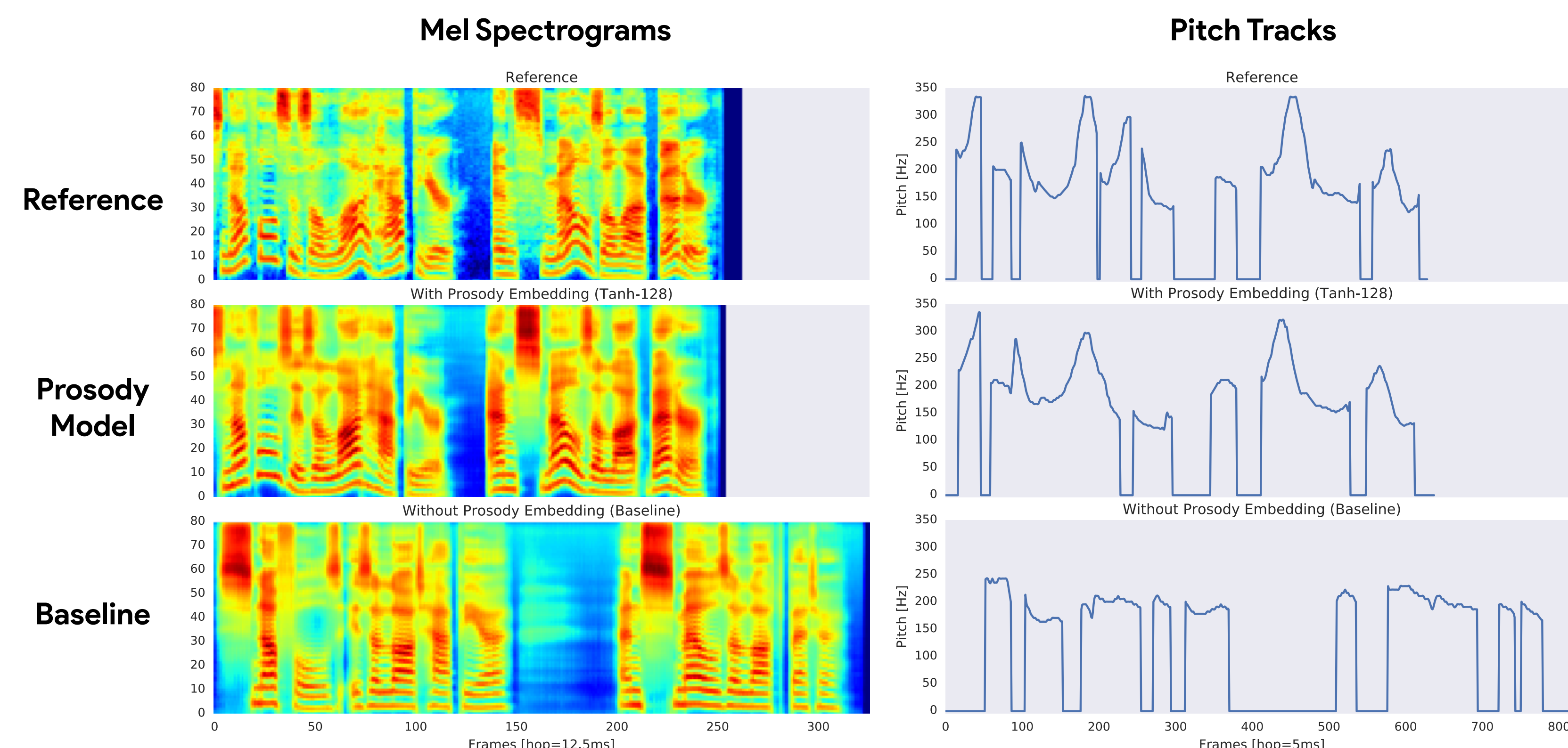
Results Table

VOICE	MODEL	REFERENCE	MCD ₁₃	FFE	SUBJECTIVE
SINGLE-SPEAKER	BASELINE	SAME SPEAKER	10.63	53.2%	
SINGLE-SPEAKER	TANH-128	SAME SPEAKER	7.92	28.1%	1.611 ± 0.164
SINGLE-SPEAKER	BASELINE	UNSEEN SPEAKER	11.22	59.6%	
SINGLE-SPEAKER	TANH-128	UNSEEN SPEAKER	8.89	38.0%	1.465 ± 0.132
MULTI-SPEAKER	BASELINE	SAME SPEAKER	9.93	48.5%	
MULTI-SPEAKER	TANH-128	SAME SPEAKER	6.99	27.5%	1.307 ± 0.127
MULTI-SPEAKER	BASELINE	SEEN SPEAKER	12.37	64.2%	
MULTI-SPEAKER	TANH-128	SEEN SPEAKER	9.51	37.1%	0.871 ± 0.138
MULTI-SPEAKER	BASELINE	UNSEEN SPEAKER	11.84	60.0%	
MULTI-SPEAKER	TANH-128	UNSEEN SPEAKER	10.87	41.3%	1.146 ± 0.246

Visualizing Prosody Transfer

Single-speaker model with unseen reference speaker

Text: *Snuffles is a lot happier. And smells a lot better.*



Experiment Setup

Datasets

- **Single-speaker dataset:** Audio book dataset, 147 hours, 49 books, read in an animated and emotive storytelling style (2013 Blizzard Challenge speaker).
- **Multi-speaker dataset:** Proprietary assistant-style dataset, 296 hours, 44 speakers (5 Australian, 6 British, 1 Indian, 2 Singaporean, and 30 American).

Training

- Train for >200k steps with Adam, batch size 256.
- Learning rate annealed from 1E-3 to 5E-5.
- Converges in 3-4 days.



Robustness to Text Modifications

- Prosody embeddings work even with modified text input.
- Example modifications:
 - Reference: *For the **first** time in **her** life she had been **danced** **tired**.*
 - Modified: *For the **last** time in his life he had been **handily** **embarrassed**.*
- Audio samples show that prosody is transferred to modified text.

Preservation of Target Speaker Identity

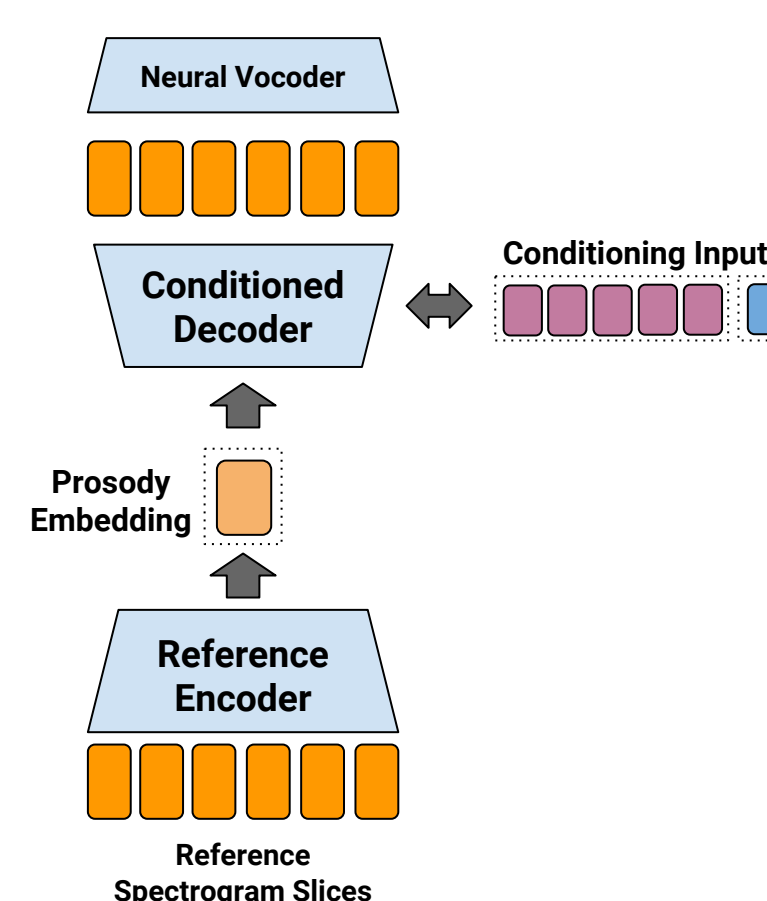
Speaker Classifier

- A simple speaker classifier is 99% accurate on ground truth and baseline synthesized audio.
- However, for a prosody model, it **picks the target speaker only 20% of the time**.
- And **picks the reference speaker 61% of the time**.

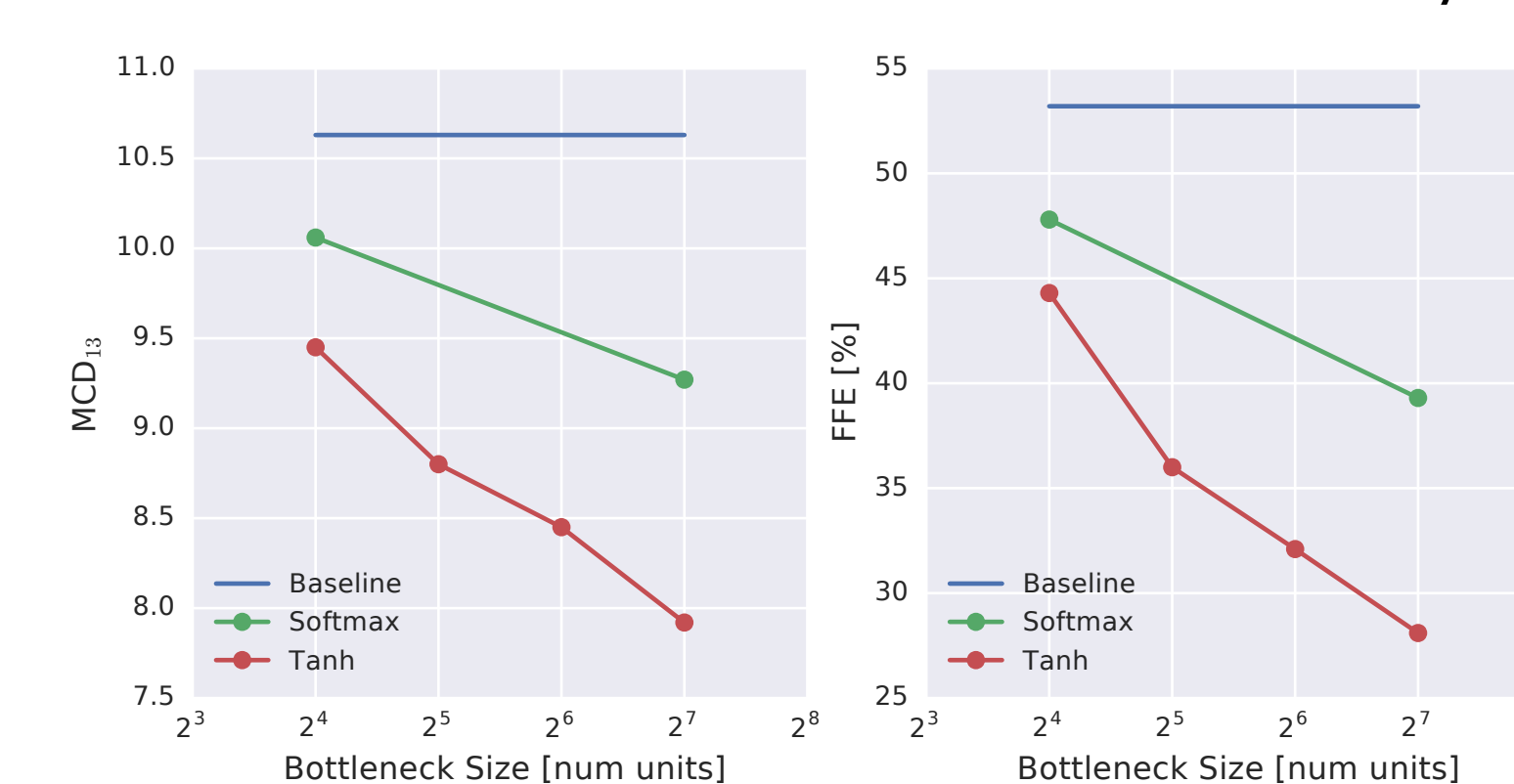
Disentangling Speaker and Prosody

- Prosody is a part of speaker identity.
- Which aspects of speaker identity should be preserved during prosody transfer?
- **Pitch-relative transfer is a feasible goal.**

Conditional Autoencoder Interpretation



Bottleneck determines reconstruction accuracy



Audio Samples

Please visit our demo page for audio samples:

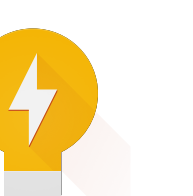
https://google.github.io/tacotron/publications/end_to_end_prosody_transfer/

Or ask to hear a demo!

0:00 / 0:06



Takeaways



- Prosody is a very important aspect of speech.
- Prosody transfer is a natural interface for prosody control.
- End-to-end prosody transfer works well and is robust to text transformations.
- Pitch-relative prosody transfer is a goal for future work.