

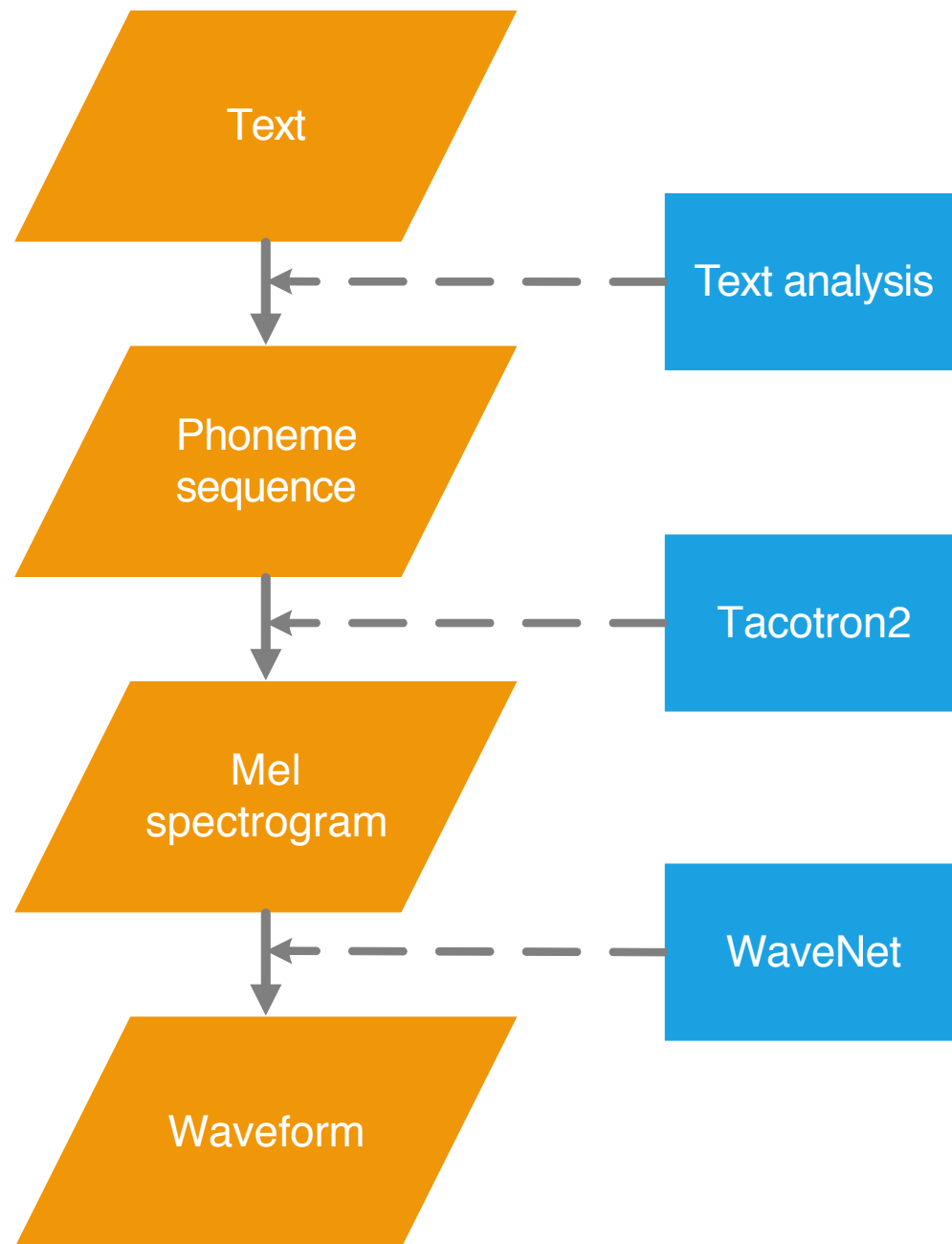
# • The Mobvoi TTS System for BC2019

Bing Liu, Yunlin Chen, Hao Yin,  
Yongqiang Li, Xin Lei, Lei Xie

Presenter: Shan Yang

- System
- Data processing
- Front-end
- Back-end: Tacotron2
- Vocoder: WaveNet
- Results

# System



- Source data: 480 audio files, 48 kHz, 8 hours, MP3 format
- Data processing:
  - Convert to WAV format
  - downsample: 48 kHz -> 16 kHz
  - Segmentation:  $\leq 10$ s per segment
  - Control silence at the beginning and end
  - Energy-based normalization
  - Clean-up text
- Final data: 4187 audio files, 16k Hz, WAV format

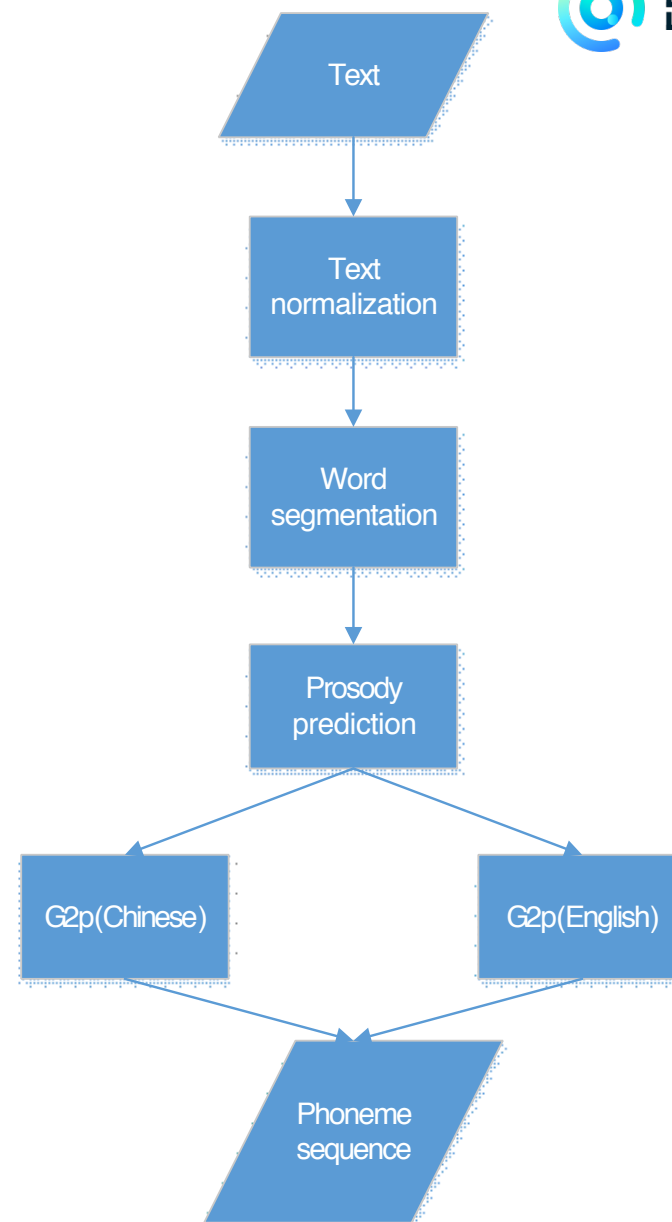
## Front-end

Text:

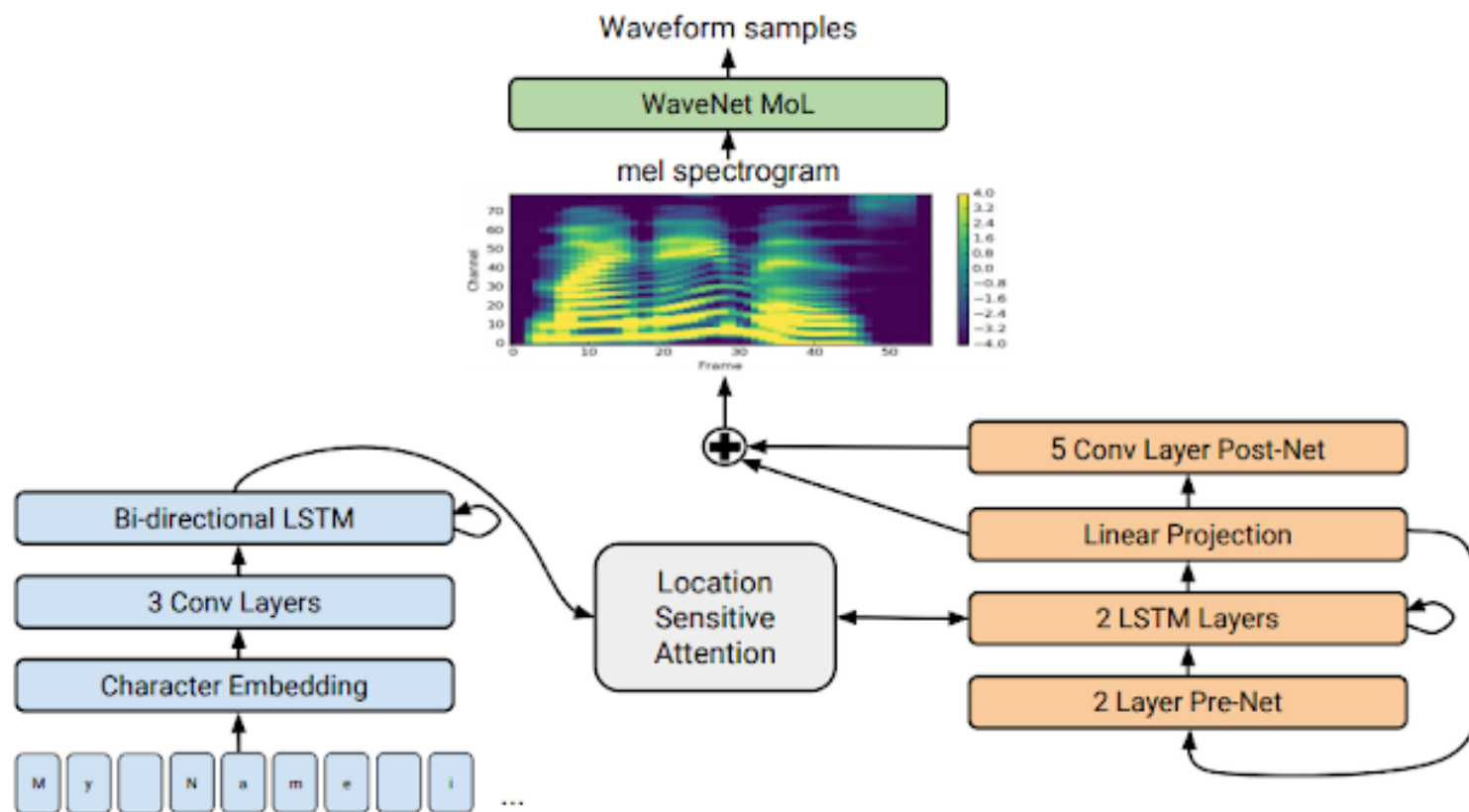
今晚八点,在得到App的直播间,吴军老师会做一场直播。

Phoneme sequence:

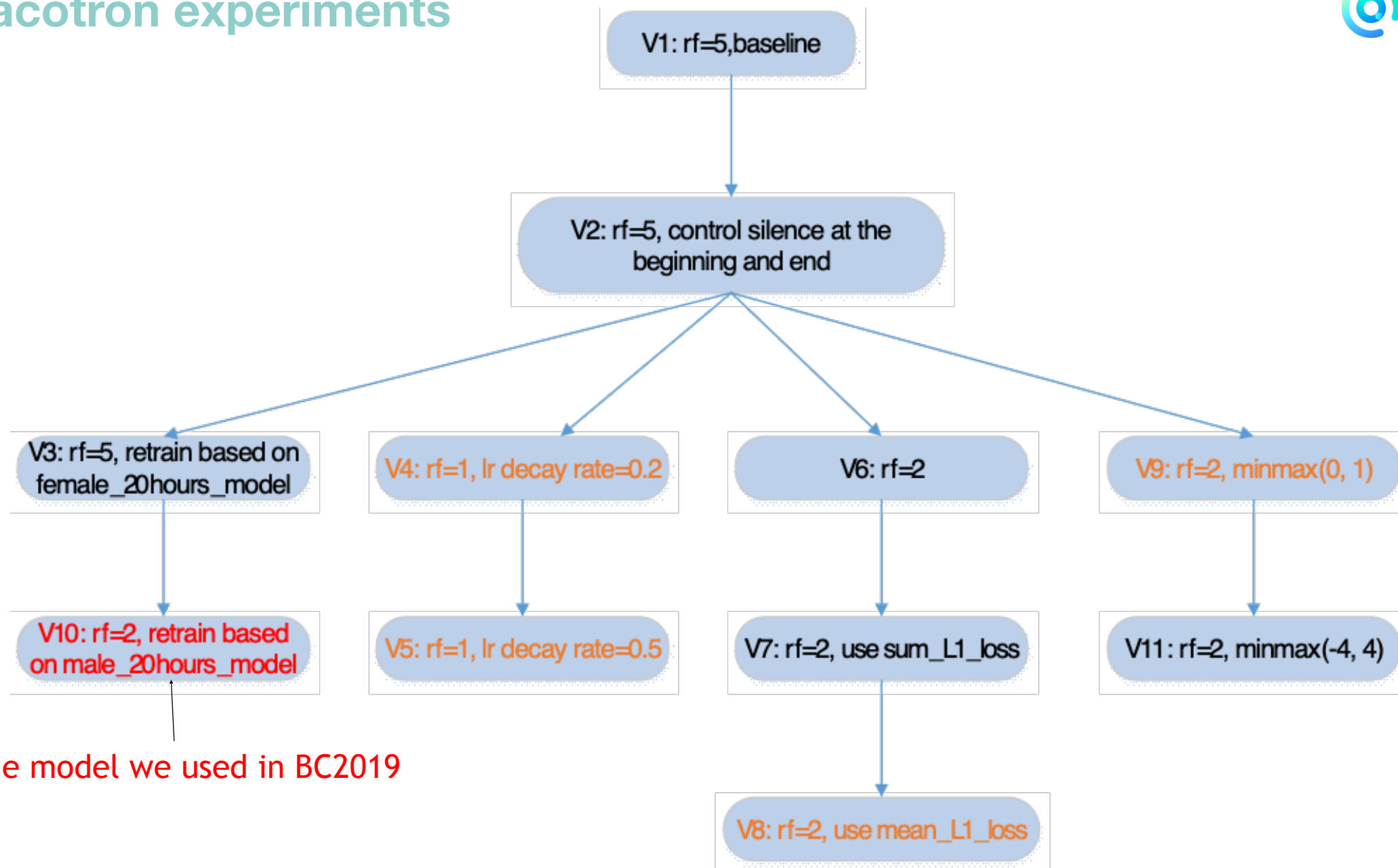
j in 1 w an 3 #1 b a 1 d ian 3 #3 z ai 4 #1 d e 2  
d ao 4 #1 / EY 1 . P IY 1 . P IY 1 / d e 5 #1 zh i  
2 b o 1 j ian 1 #3 w u 2 j vn 1 #1 l ao 3 sh i 1  
#2 h ui 4 z uo 4 #1 y i 1 ch ang 3 #1 zh l 2 b o  
1 SIL



# Back-end: Tacotron2



# Tacotron experiments



The model we used in BC2019

# Tacotron experiments



control silence at the beginning and end

V1:**V2**

To control silence at the beginning and end is better

Adaptation

V2:**V3**, V6:**V10**

Adaptation is better, **great improvement**

Reduce factor (rf)

V2(rf=5):**V6**(rf=2)

rf=2 is better

V4, V5 (rf=1)

No alignment

L1 loss

V6:**V7**(sum L1 loss)

So close

V8(mean L1 loss)

No alignment

Minmax normalization

V6:**V11**(minmax(-4,4))

Minmax is better

V9(minmax(0,1))

No alignment

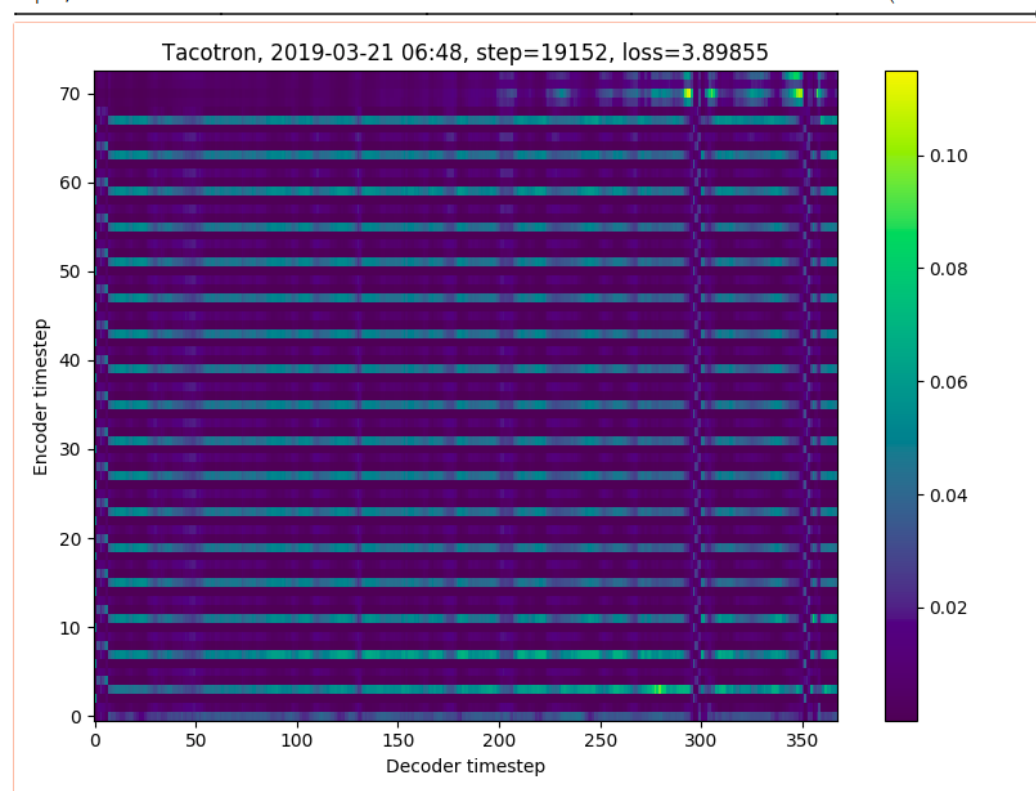


# Tacotron experiments

- No alignment

eval/alignment/image/0  
step 19,152

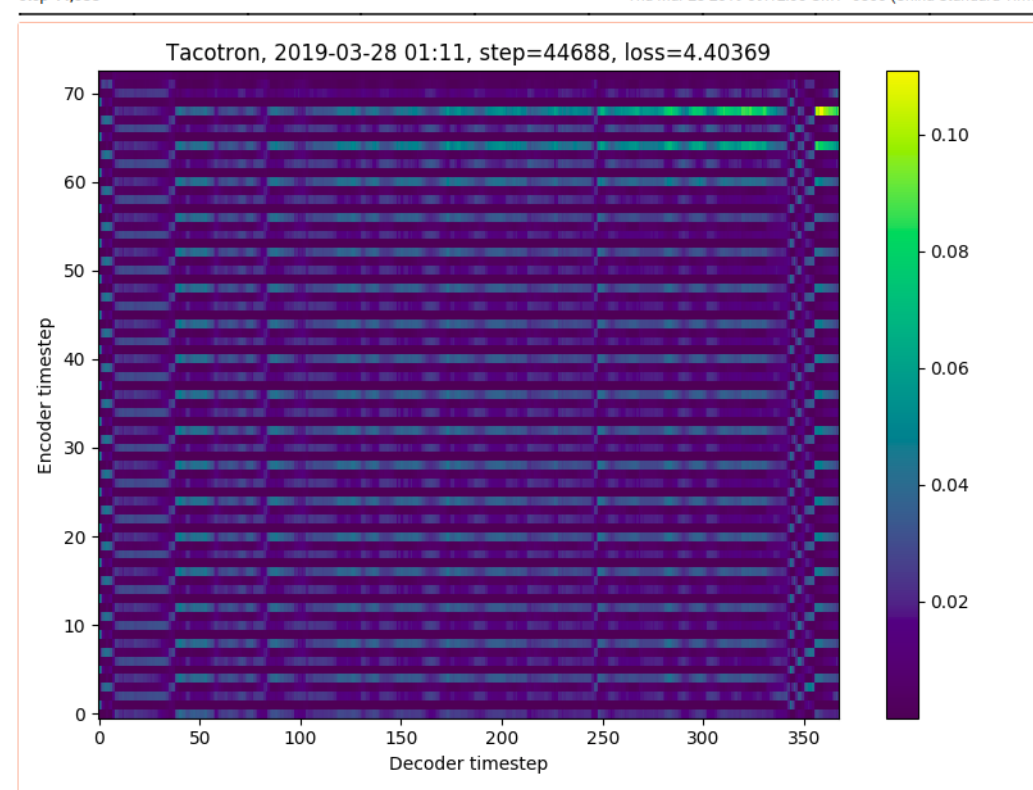
logs-tacotron/files  
Thu Mar 21 2019 14:48:51 GMT+0800 (China Standard Time)



V4: rf=1,  
lr\_decay\_rate=0.2

eval/alignment/image/0  
step 44,688

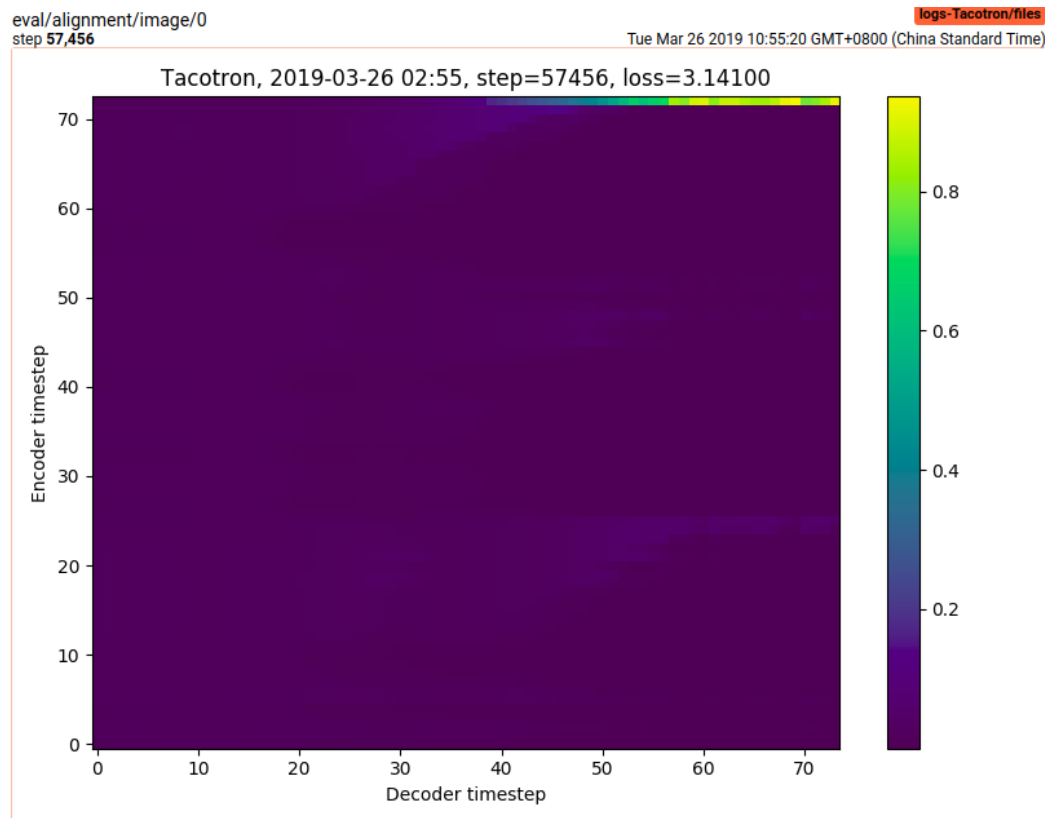
logs-Tacotron/files  
Thu Mar 28 2019 09:12:05 GMT+0800 (China Standard Time)



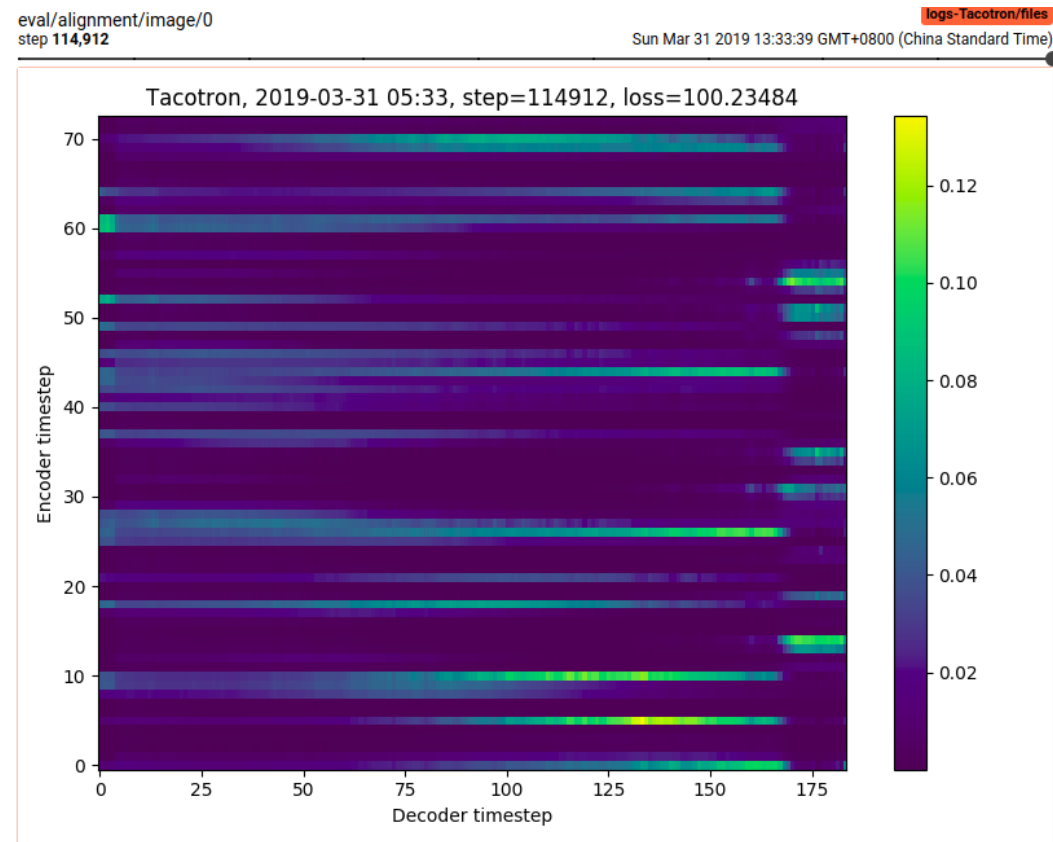
V5: rf=1,  
lr\_decay\_rate=0.5

# Tacotron experiments

- No alignment



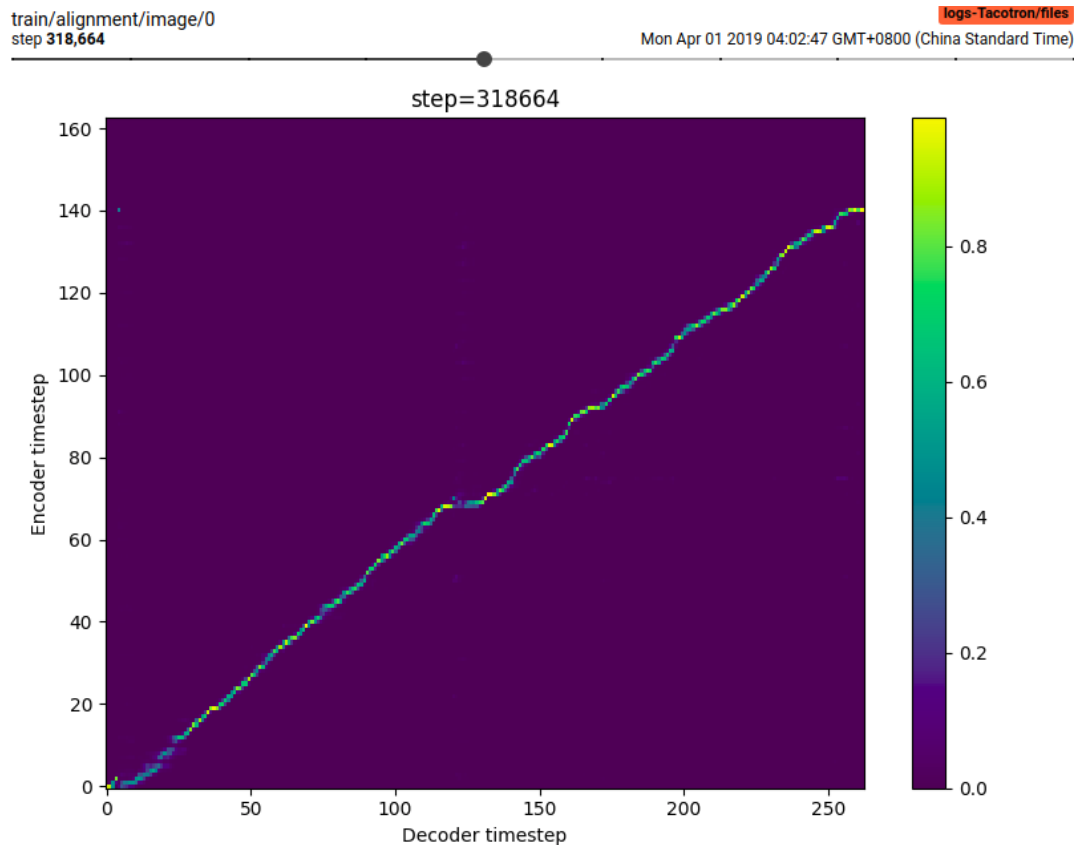
**V8:** rf=2, use mean L1  
loss



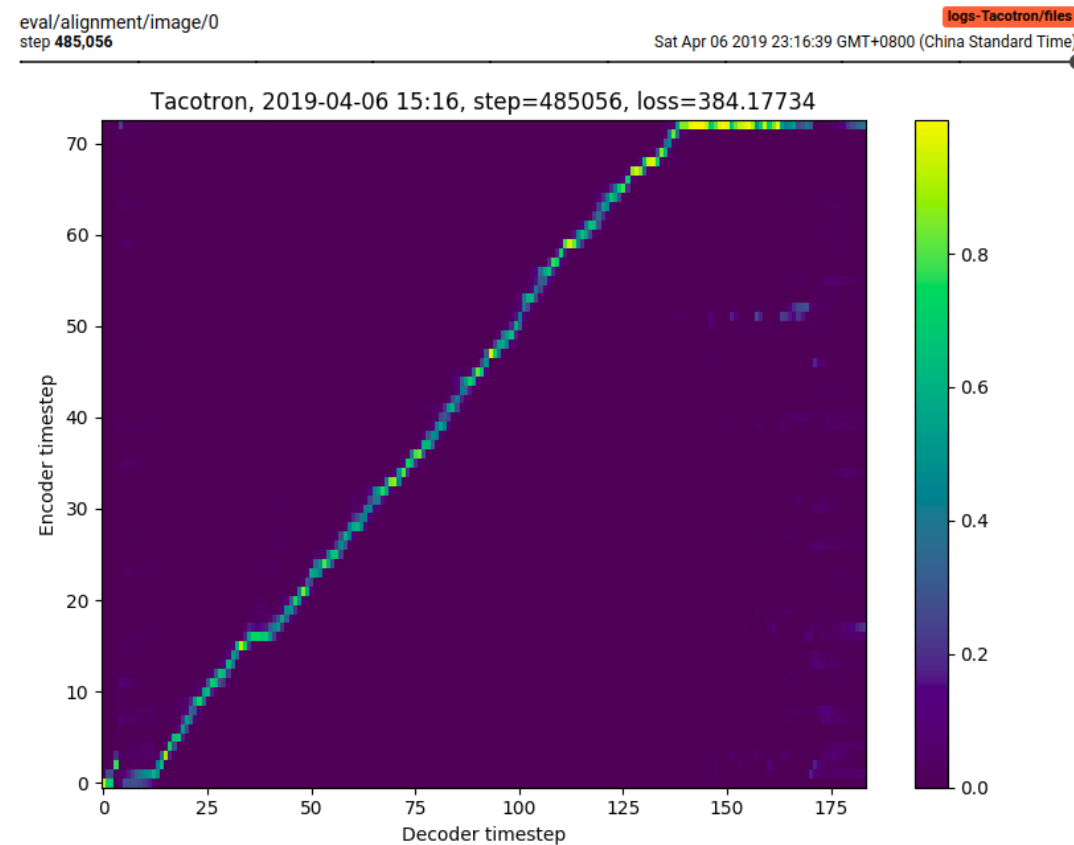
**V9:** rf=2, minmax(0, 1)

# Tacotron experiments

- Normal alignment

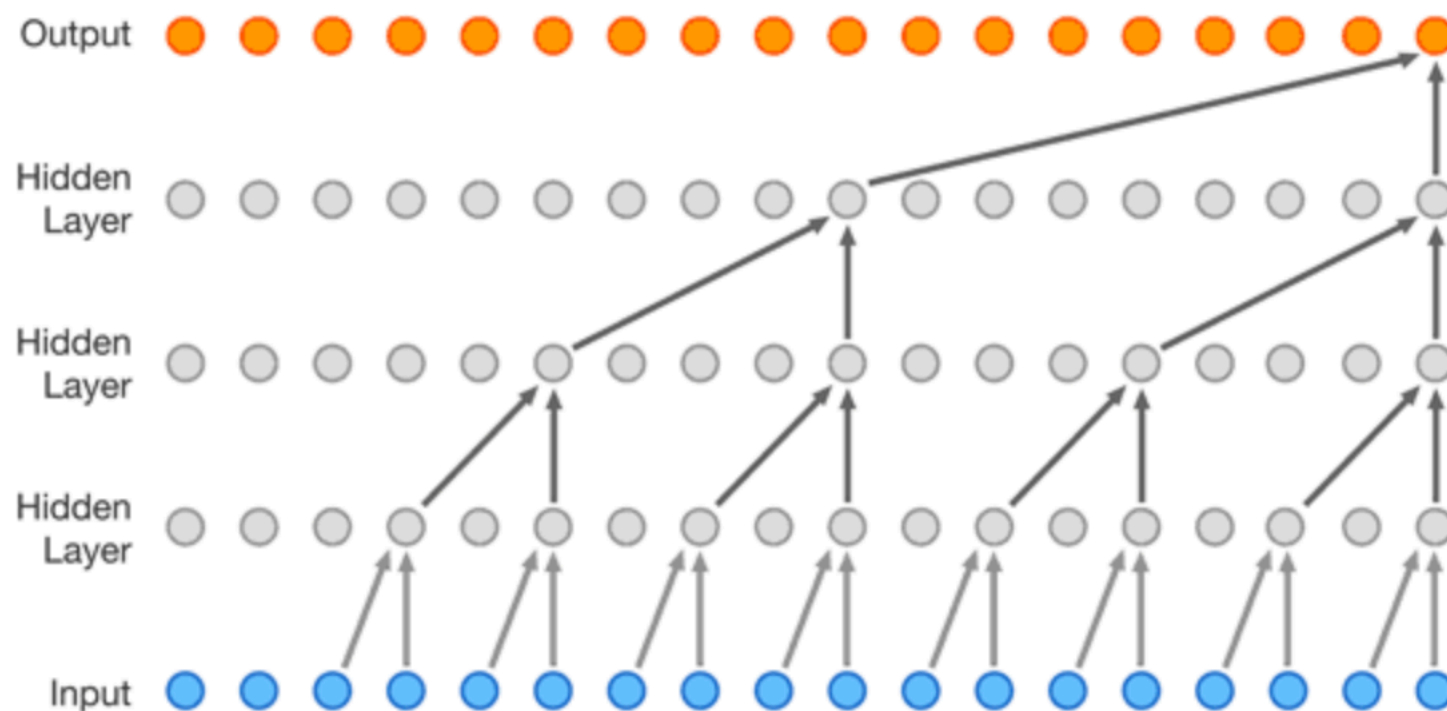


V10: rf=2, retrain, training set



V10: rf=2, retrain, valid set

# Vocoder: Wavenet



1. 20 dilated convolution layers, grouped into 2 dilation cycles
2. Predict 10-component mixture of logistic distributions (MoL)

- WaveNet experiments:
  - **V1**: Baseline
    - Conditioning in ground truth mel spectrogram
    - Stop\_step: 1,000,000
  - **V2**: GTA(taco v10) model:
    - Conditioning in mel spectrogram generated from tacotron v10 GTA inference
    - Stop\_step: 2,000,000, retrain based on v1

- Ground Truth Alignment(GTA) training

Training	Synthesis	
	Predicted	Ground truth
Predicted	$4.526 \pm 0.066$	$4.449 \pm 0.060$
Ground truth	$4.362 \pm 0.066$	$4.522 \pm 0.055$

**Table 2.** Comparison of evaluated MOS for our system when WaveNet trained on predicted/ground truth mel spectrograms are made to synthesize from predicted/ground truth mel spectrograms.

# Results

Original Wave		
WaveNet_1M(Ground Truth feature)		
WaveNet_1M(taco_v10)		
WaveNet_retrain_2M (taco_v10)	Random test	
	BC2019 test	

- To improve
  - **Rhotic accent:** some are somewhat blurred.
  - **Interrogative sentences:** Insufficient tone.
  - **Rhythm:** some sentences are read faster, and some sentences are unreasonably paused (too long or too short).
  - **English:** english adaptive training with less data.
  - **Digital:** missing or blurring of consecutive identical numbers.
  - **Missing:** miss word at the end of the sentence.
  - **Ancient poetry:** the rhythm of ancient poetry needs to be improved.





 **Thanks**