

Kuan Yu                      kuanyu@uni-potsdam.de

Simon Untergasser      untergasser@uni-potsdam.de

Jörg Schwartz            jschwartz@uni-potsdam.de

July 18, 2018

# outline

data

model

status

plan

# LJSpeech-1.0<sup>1</sup>

- ▶ 13,100 audio clips
- ▶ 1.11 sec to 10.10 sec
- ▶ single speaker
- ▶ with normalized transcription
- ▶ first 12 clips for validation

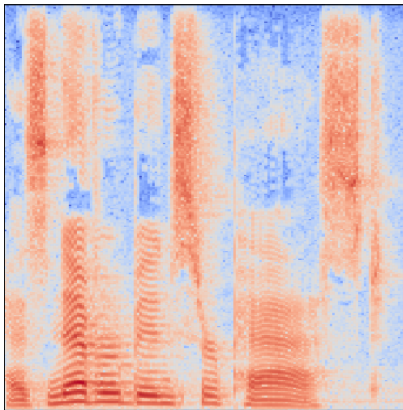
---

<sup>1</sup><https://keithito.com/LJ-Speech-Dataset/>

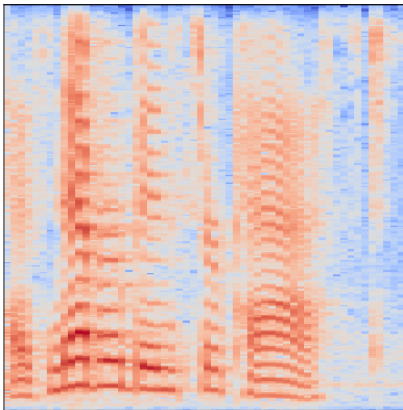
# clean up transcription

- ▶ inconsistent use of quotes and brackets
- ▶ non-ascii characters
- ▶ case

## downsampling



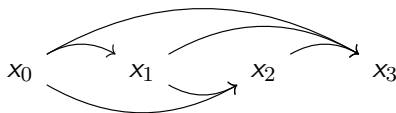
22050



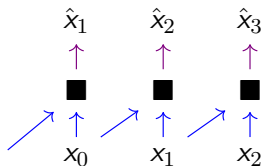
8000

# autoregressive model

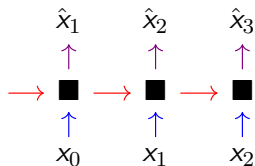
- ▶ number of frequency bins  $f = 256$
- ▶ frame  $x_i : \mathbb{C}^f \simeq \mathbb{R}^{2f}$
- ▶ model  $x_{0...i} \rightarrow x_{i+1}$



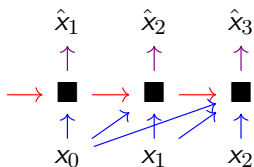
# choices



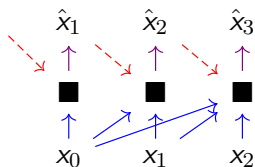
convolution



recurrent

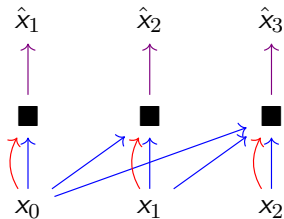


recurrent with attention

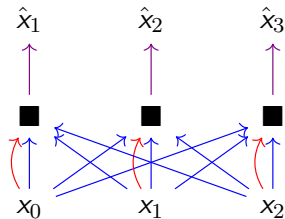


redirect recurrent connection

## self-attention<sup>2</sup>



decoder



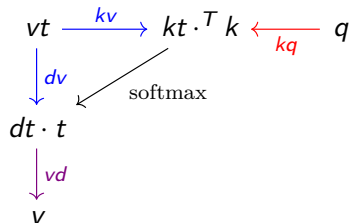
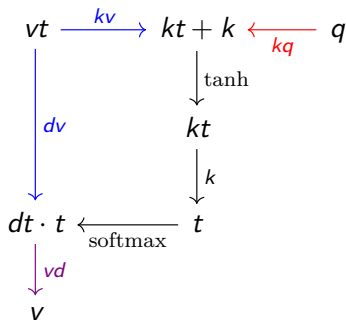
encoder

<sup>2</sup><https://arxiv.org/abs/1703.03130>



# attention: additive<sup>3</sup>, dot-product<sup>4</sup>, key-value<sup>5</sup>

time steps  $t$ , value dim  $v$ , key dim  $k$ , query dim  $q$ , model dim  $d$



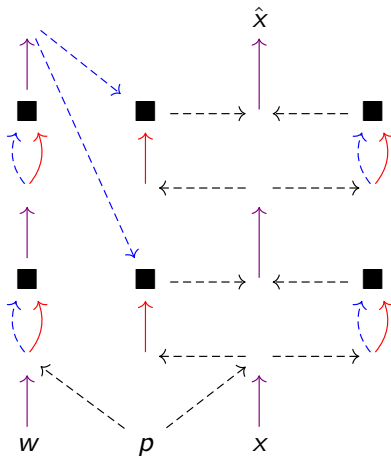
<sup>3</sup><https://arxiv.org/abs/1409.0473>

<sup>4</sup><https://arxiv.org/abs/1508.04025>

<sup>5</sup><https://arxiv.org/abs/1702.04521>

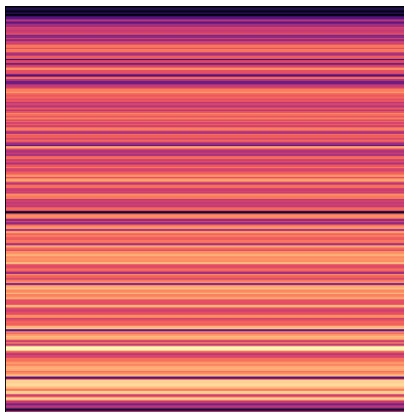
## architecture based on transformer<sup>6</sup>

- ▶  $w$ : transcript (characters)
- ▶  $p$ : position encoding (sinusoids)

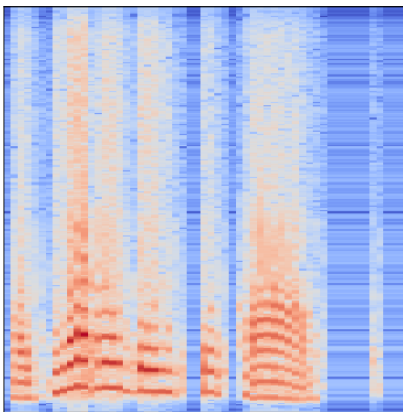


<sup>6</sup><https://arxiv.org/abs/1706.03762>

## results

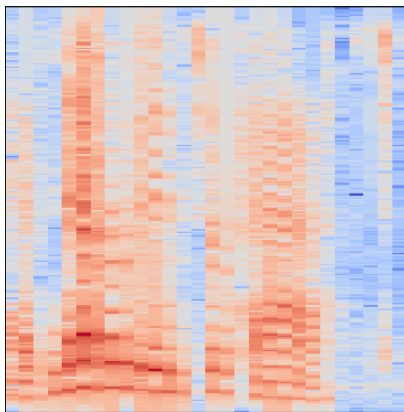


autoregressive

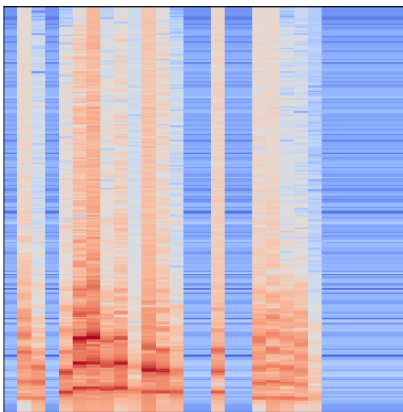


teacher forcing

# boxcar

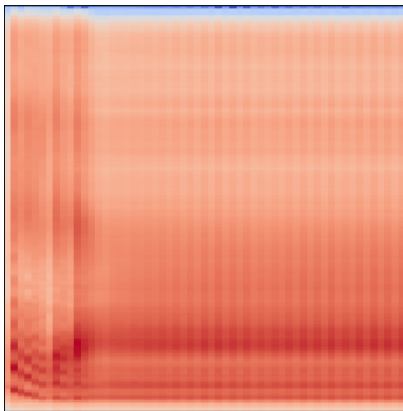


boxcar

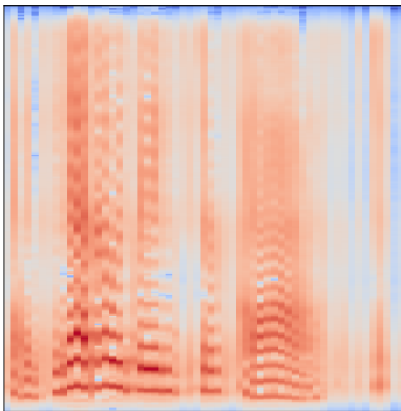


boxcar teacher forcing

## magnitude only



autoregressive



teacher forcing

# plan

- ▶ steal from tacotron<sup>7</sup>
- ▶ complex arithmetics

---

<sup>7</sup><https://github.com/keithito/tacotron>