



## Deep Speech 2: End-to-End Speech Recognition in English and Mandarin

Baidu Silicon Valley AI Lab, Baidu Speech Technology Group [Presented by Eric Battenberg]

Baidu Research Silicon Valley Al Lab



# Speech Interfaces: A New/Old Input Paradigm

- Speaking is human
- Speed: speak 150 wpm vs type 40 wpm
- Hands-free
- Skip the confusing menus
- Small footprint
- The huge difference between 95% vs 99% accuracy

Start













# The Deep Learning Bet

- SVAIL is an AI lab, not a speech lab.
- Instead of more domain expertise...
- Scale the model, scale the data
  - -> Improved performance.
- -> Superhuman Chinese speech recognition





![](_page_3_Picture_0.jpeg)

#### The Evolution of Speech Recognition Systems

![](_page_3_Figure_2.jpeg)

![](_page_4_Picture_0.jpeg)

#### The Evolution of Speech Recognition Systems

![](_page_4_Figure_2.jpeg)

![](_page_5_Picture_0.jpeg)

# Deep Speech 1 to Deep Speech 2

- Multiple layers of 2D convolution
- Up to 7 recurrent layers
- Trained with Batch Norm
- + Lots of systems expertise

![](_page_5_Picture_6.jpeg)

#### Deep Speech 2

![](_page_5_Picture_8.jpeg)

6

![](_page_6_Picture_0.jpeg)

### Batch Normalization For RNNs

- Sequence-wise Batch Norm on upward connections
- More effective in deeper networks
- Speeds up training
- Improves test accuracy

![](_page_6_Figure_6.jpeg)

![](_page_7_Picture_0.jpeg)

# Scaling Up Deep Speech 2

- Tens of <u>exaflops</u> required to train model
- Cluster with 8 TitanX GPUs per node
  - Partition minibatch across GPUs
  - Synchronous SGD

![](_page_7_Picture_6.jpeg)

![](_page_7_Picture_7.jpeg)

**Baidu Research** 

![](_page_8_Picture_0.jpeg)

# Scaling Up Deep Speech 2

- Custom All-Reduce code
  - 4x-21x speedup over OpenMPI's
- Fast GPU CTC (Warp-CTC)
  - Reduced training time by 10%-20%
- Fastest available kernels
- Overall, sustained 45% peak performance on each node.

![](_page_8_Picture_8.jpeg)

# English Results

- English training data:
  - 11,940 hours
  - 8 million utterances
  - ...and growing everyday

![](_page_9_Figure_5.jpeg)

![](_page_10_Picture_0.jpeg)

# English Results

- English training data:
  - 11,940 hours
  - 8 million utterances
  - ...and growing everyday

#### **Accented Speech - VoxForge [English]**

![](_page_10_Figure_7.jpeg)

![](_page_11_Picture_0.jpeg)

# Porting to Mandarin

- Increase softmax size:
  - from 29 to ~6000
- Feed in Mandarin data
- Tweak some hyperparameters

Deep Speech 2 English	 Deep Speech 2 Mandarin
Softmax	Softmax
Fully Connected	Fully Connected
Bidir ReLU RNN	Bidir ReLU RNN
Bidir ReLU RNN	Bidir ReLU RNN
Bidir ReLU RNN	Bidir ReLU RNN
Bidir ReLU RNN	Bidir ReLU RNN
Bidir ReLU RNN	Bidir ReLU RNN
Bidir ReLU RNN	 Bidir ReLU RNN
Bidir ReLU RNN	Bidir ReLU RNN
2D Convolution	2D Convolution
2D Convolution	2D Convolution
2D Convolution	2D Convolution
Spectrogram	Spectrogram

![](_page_12_Picture_0.jpeg)

# Mandarin Results

- Mandarin training data:
  - 9,400 hours
  - 11 million utterances
  - More diverse than English data
- (...growing everyday)

![](_page_12_Figure_7.jpeg)

![](_page_13_Picture_0.jpeg)

## To Production-Ready Models

- Production models require causality
- Unidirectional recurrent layers
- GRUs instead of ReLU RNNs
- 1D lookahead convolution

![](_page_13_Figure_6.jpeg)

14

![](_page_14_Picture_0.jpeg)

# Lookahead Convolutions

- Some future context is useful.
- Difficult to force network to delay predictions.
- We use 1D convolution after recurrent layers
- Within 5% relative performance of bidirectional models

![](_page_14_Figure_6.jpeg)

![](_page_14_Figure_7.jpeg)

![](_page_15_Picture_0.jpeg)

### Research to Deployment

- Deep Speech 2 is now in production!
- Baidu Silicon Valley AI Lab + Baidu Speech Technology Group
- From research idea to deployed product in 2 years.
- Combined efforts of 69 researchers and engineers

![](_page_15_Figure_6.jpeg)

![](_page_16_Picture_0.jpeg)

- The story of Deep Speech 2
  - End-to-end speech recognition
  - Scaled up
  - Ported to Mandarin
  - Deployed to production

![](_page_16_Picture_6.jpeg)

• Deep Speech 3, coming to a conference near you in 2017!