[Video on Vimeo Here](#)

# Interactive Long-Form Conversational Speech & Language Application

# Roadmap

1. Fundamental Challenges
2. Case Study: Closed Loop Feedback in Speech Recognition Training
3. Case Study: Latency and Randomness in Speech Synthesis Inference
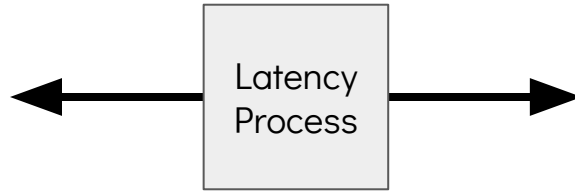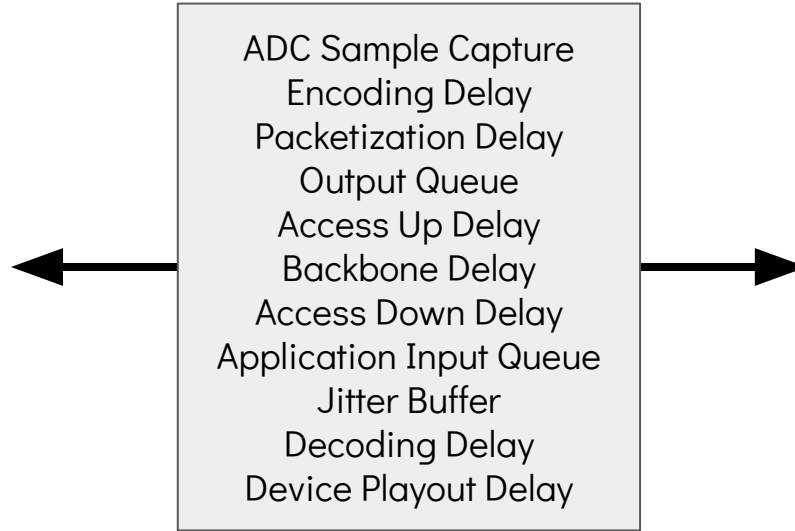
# FUNDAMENTAL CHALLENGES

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⏳ | Stochasticity 🎲 | Scale ⚖️ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity 🦋 | Interactivity 🤖 |

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⏳ | Stochasticity 🎲 | Scale ⚖️ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity 🦋 | Interactivity 🤖 |

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⧗ | Stochasticity 🎲 | Scale ⚖️ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity 🦋 | Interactivity |

# Simple function

# The Tent Map

$$F(x) = \mu \min( x , 1 - x )$$

# The Tent Map

$$F(x) = \mu \min(x, 1 - x)$$

# The Tent Map

# The Tent Map (μ < 1)



μ/2

1/2

# The Tent Map (μ > 1)



μ/2

1/2

# The Tent Map (μ >> 1)



μ/2

1/2

# The Tent Map
# Bifurcation Diagram

# More Complex Function

# Full System Feedback Loop

```
┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐
│   ASR    │──▶ │ Dialogue │──▶ │  Speech  │──▶ │   The    │──▶ │   PSTN   │──▶ │LTE Tower │
│          │    │  System  │    │Synthesis │    │ Internet │    │          │    │ Networks │
└──────────┘    └──────────┘    └──────────┘    └──────────┘    └──────────┘    └──────────┘
     ▲                                                                                │
     │                                                                                ▼
┌──────────┐                                                                    ┌──────────┐
│ Phones & │                                                                    │  iPhone  │
│ Internet │                                                                    │          │
│  Again   │                                                                    │          │
└──────────┘                                                                    └──────────┘
     ▲                                                                                │
     │                                                                                ▼
┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐
│Microphone│◀── │  Sound   │◀── │  Vocal   │◀── │  Brain   │◀── │ Cochlea  │◀── │  Sound   │
│          │    │  Waves   │    │  Tract   │    │          │    │          │    │  Waves   │
└──────────┘    └──────────┘    └──────────┘    └──────────┘    └──────────┘    └──────────┘
```
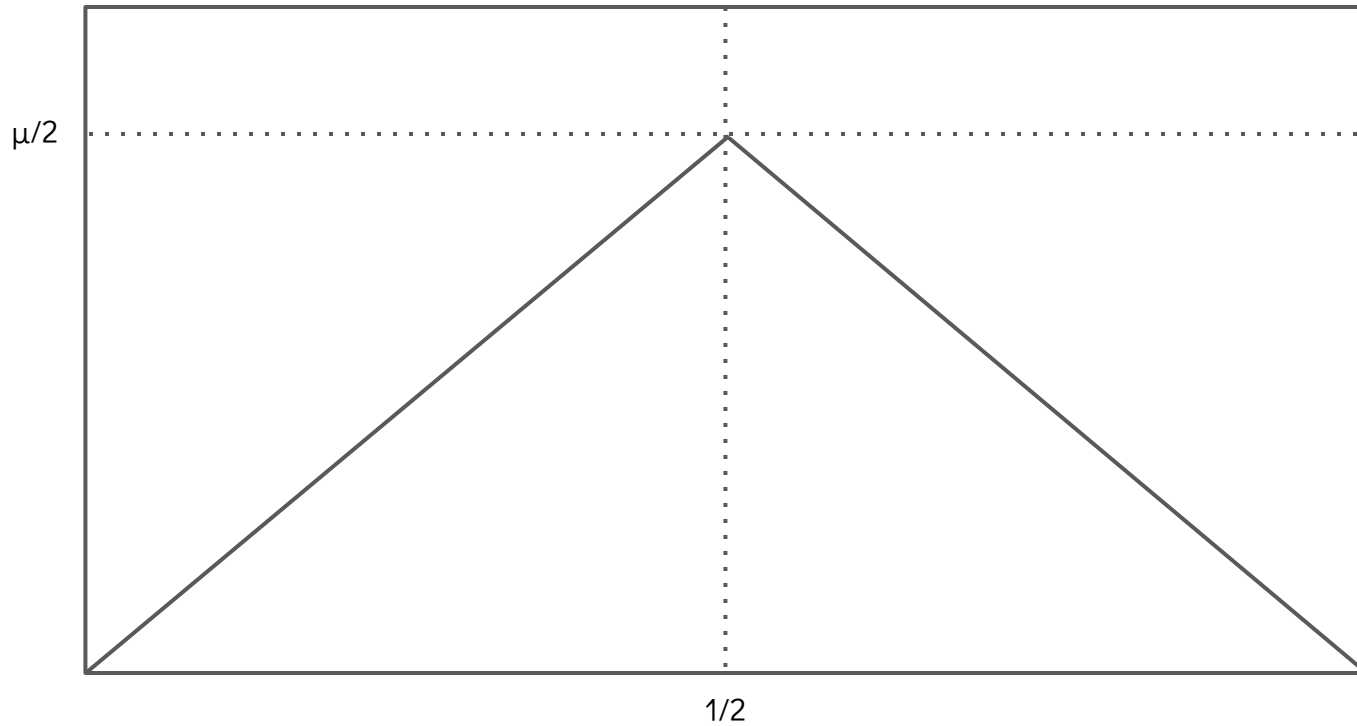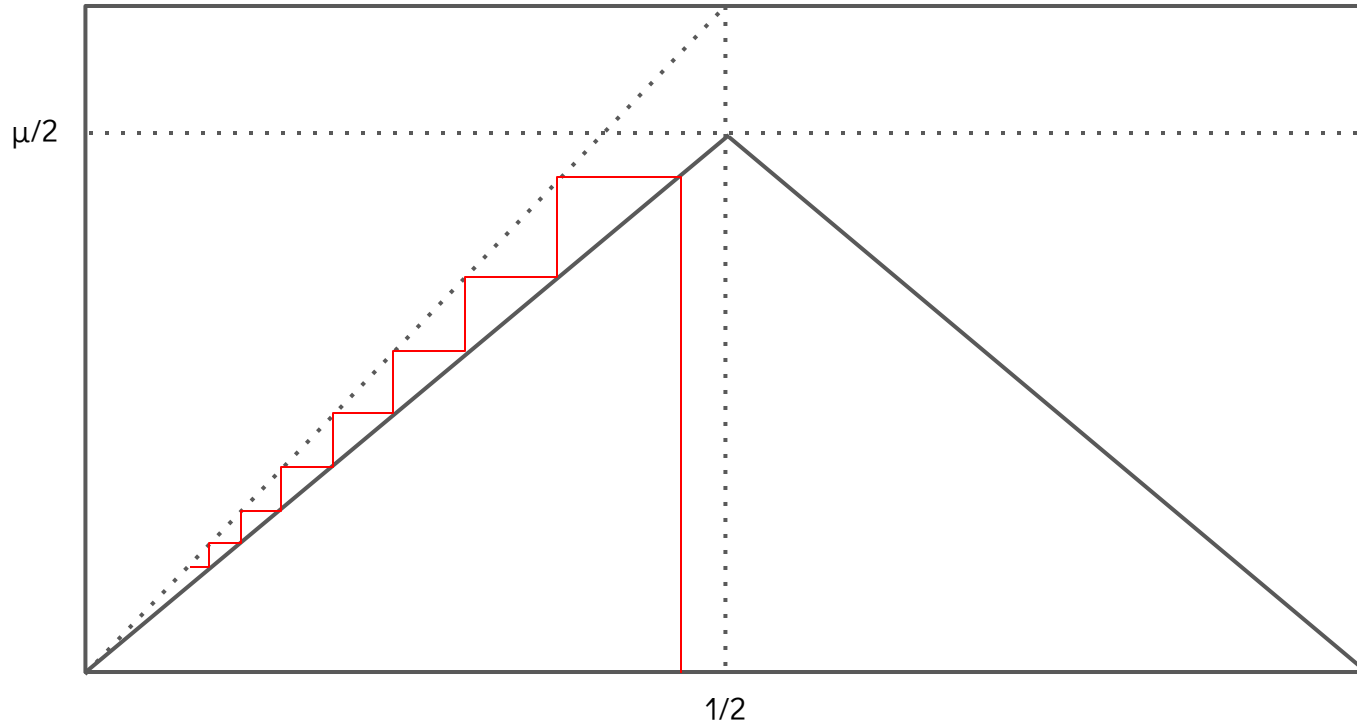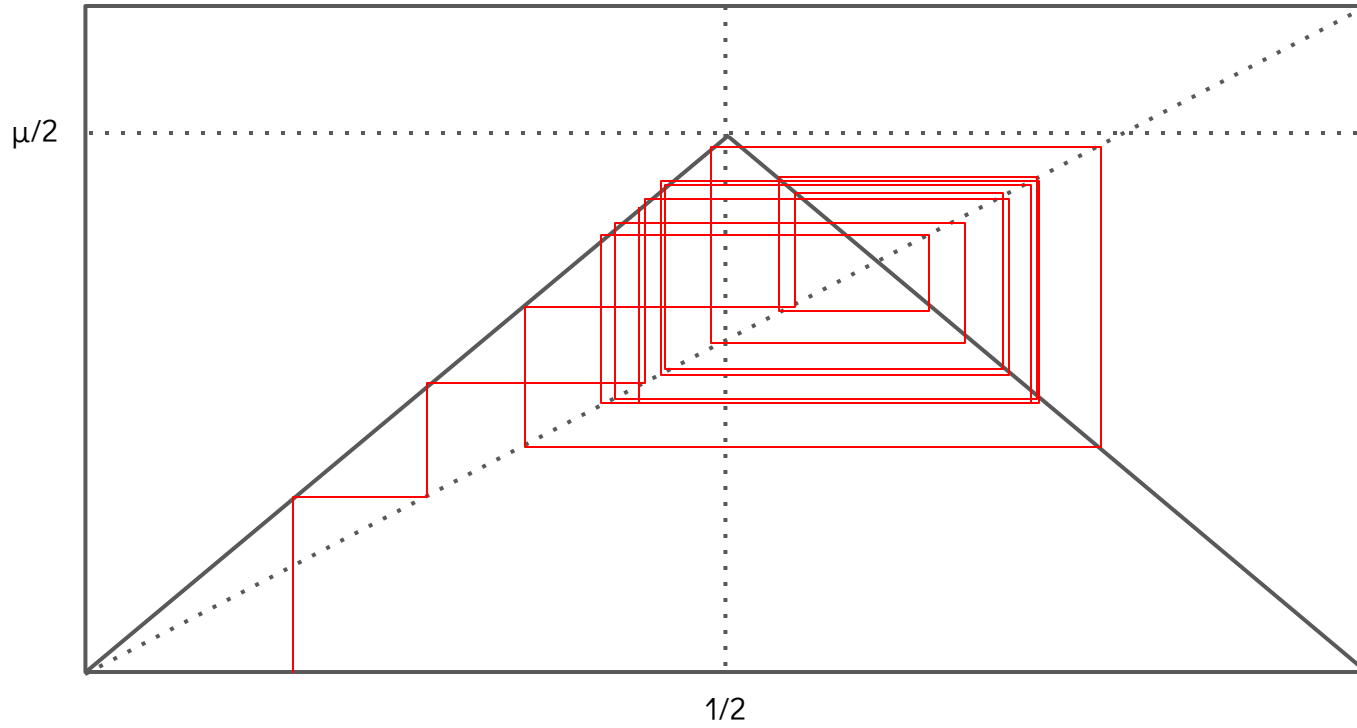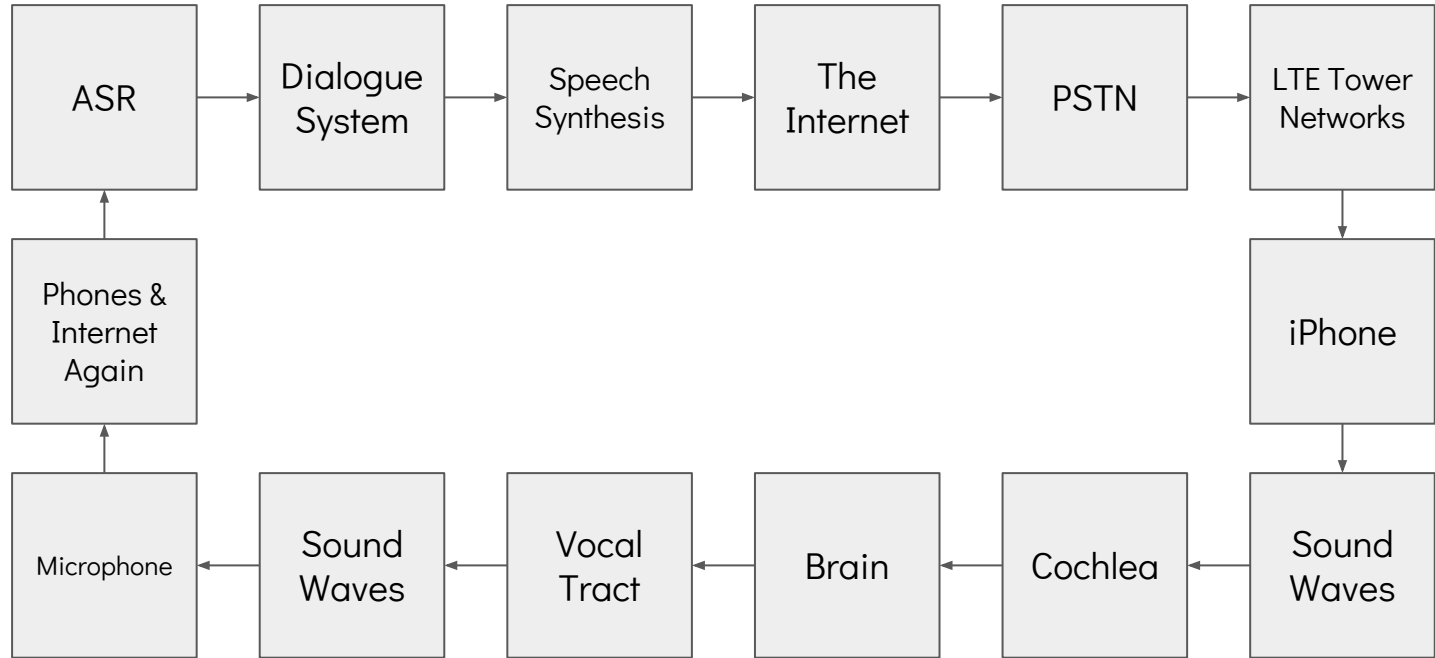
# Feature a Neural Network



ASR → Dialogue System → Speech Synthesis → The Internet → PSTN → LTE Tower Networks → iPhone → Sound Waves → Cochlea → Brain → Vocal Tract → Sound Waves → Microphone → Phones & Internet Again → ASR

# Physical Effects

# Noisy & Stochastic

# Nonlinearity and Distortion



ASR → Dialogue System → Speech Synthesis → The Internet → PSTN → LTE Tower Networks → iPhone → Sound Waves → Cochlea → Brain → Vocal Tract → Sound Waves → Microphone → Phones & Internet Again → ASR

# Information Loss

# Consumes Latency Budget



ASR → Dialogue System → Speech Synthesis → The Internet → PSTN → LTE Tower Networks → iPhone → Sound Waves → Cochlea → Brain → Vocal Tract → Sound Waves → Microphone → Phones & Internet Again → ASR

# NONLINEARITY

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⌛ | Stochasticity 🎲 | Scale ⚖️ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity | Interactivity 🤖 |

$$I. \ F(x + y) = F(x) + F(y)$$

# F(x)



1x   ⟲10   ▶    🗨   00:00 /00:04      🔊 ▬▬▬▬

🔍 Search transcript    ✕     ☑ Show Annotations

0:00:00   San Jose is a city in California.

# F(y)



1x  ⟲10  ▶  💬  00:00 /00:04                                          🔊 ────

🔍 Search transcript                                    ✕     ☑ Show Annotations

0:00:00   Machine learning as a statistical computing
          technique.

# F(x+y)

1x ⟲ ▶ 💬 00:00 /00:04 🔊 ━━━

🔍 Search transcript ✕

☑ Show Annotations

0:00:01 That needs to be easy to California.

Linearity

# II. F(ax) = aF(x)

# F(x)



1x  ⟲10  ▶  🗨  00:00 /00:04                                                🔊 ———

🔍 Search transcript                                          ✕        ☑ Show Annotations

0:00:00   San Jose is a city in California.

# F(ax)

1x  00:00 /00:04

Search transcript ✕

☑ Show Annotations

0:00:03  Yeah.

STOCHASTICITY

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⏳ | Stochasticity 🎲 | Scale ⚖️ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity 🦋 | Interactivity 🤖 |

JITTER & APERIODICITY

New England Load - Super Bowls 49 and 50

# Typical US Call Center Diurnal (Weekday)



Lunch break

Start of the work day
(eastern time zone)

# Typical US Call Center Diurnal (First of month)



*Lunch break*

*Start of the work day
(eastern time zone)*

# SCALE

$$f(\mathbf{V})$$

$$f(\mathbf{V})$$

API Web Server ↔ $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$

$f(\mathbf{V})$  $f(\mathbf{V})$

# CASE STUDY:
# CLOSED LOOP FEEDBACK IN SPEECH RECOGNITION TRAINING
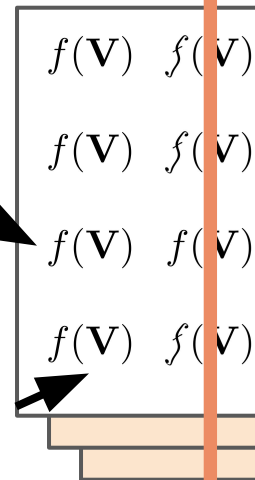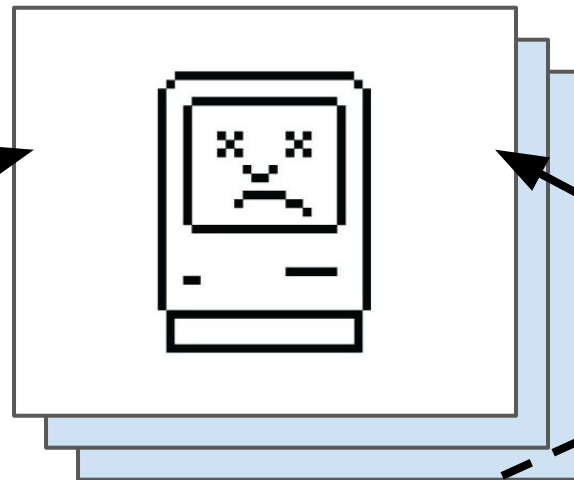
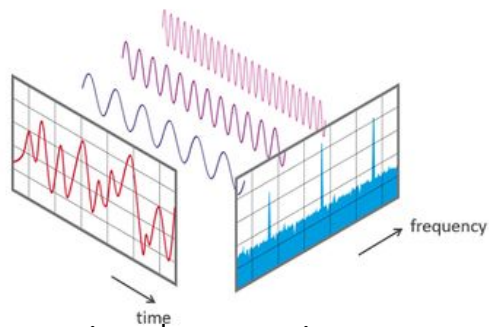|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⌛ | Stochasticity 🎲 | Scale ⚖️ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity 🦋 | Interactivity 🤖 |

# Speech Recognition



"Welcome to Gridspace"
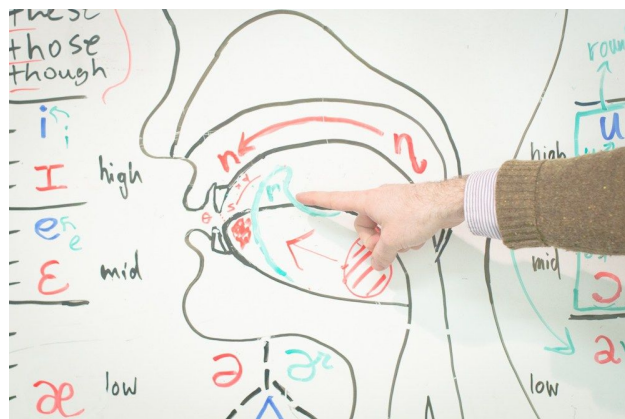
# Speech Recognition


< signal processing >
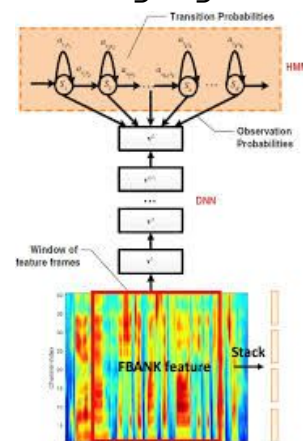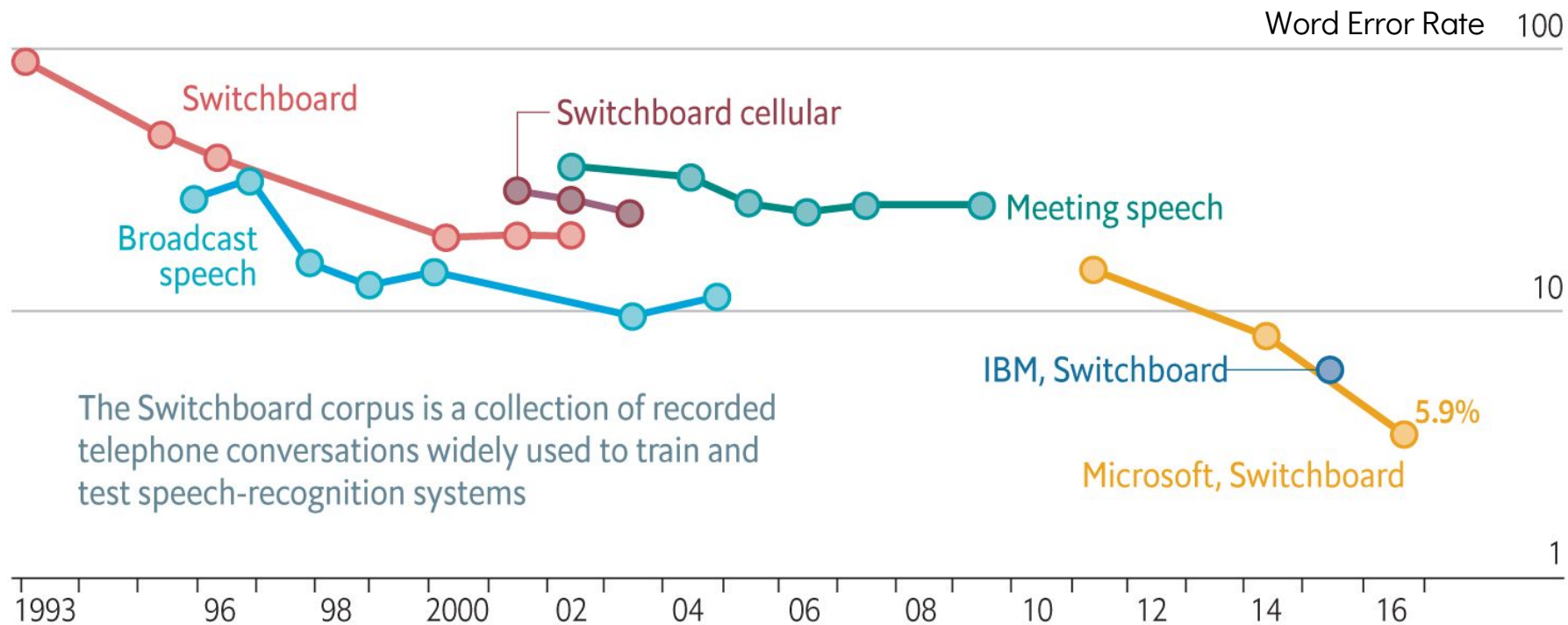

< Language>


< pronunciation model>


< acoustic model>

# Speech Recognition - History

# Data Boost



Data(Hrs) vs. Year

# Computation Boost



GPU Performance (FP32, single precision floating point)

(from: Grigory Sapunov, Intento blogs)

# Model Boost



(from: Ming Sun et el., IS 2017 paper)

input layer   hidden layer 1   hidden layer 2   hidden layer 3

output layer

(from: Kyu Han et el., IS 2018 paper)

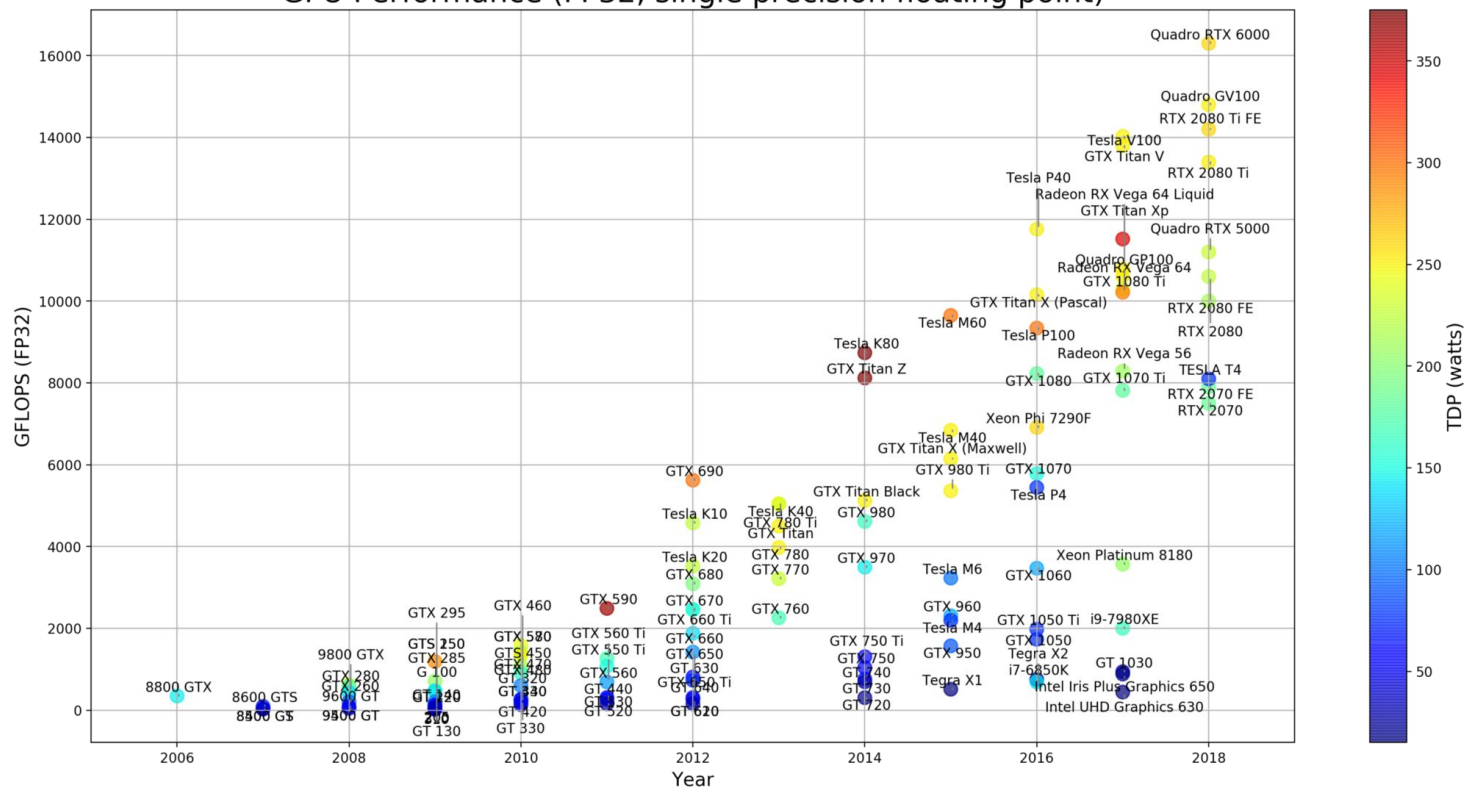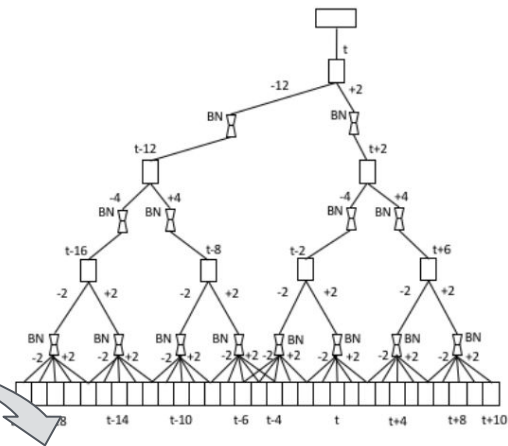# Speech recognition in the Wild - Still Challenging

"When my oldest child was a baby..."

"It didn't send me any kind of verification code. It did pop up a message saying that the account was locked"





| Switchboard | DATA | Gridspace Call Center Data |
|---|---|---|
| normal vocab, normal emotion | Speech Type | domain specific vocab, emotional modulated speech |
| 7.5% | Word Error Rate | 18% (without optimization) |

# Speech recognition in the Wild - Still Challenging



# Each domain has each specific language.

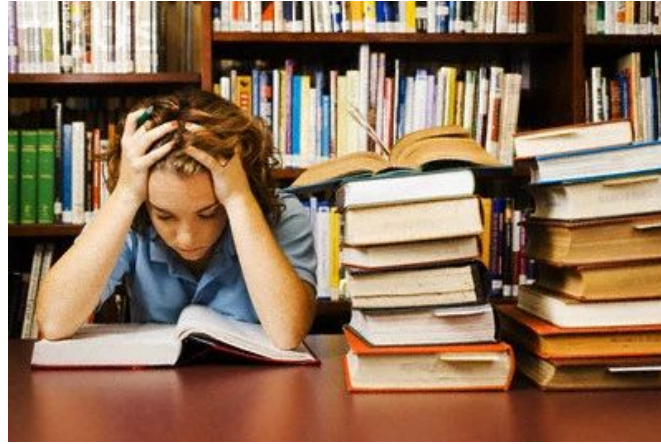- command query
- financial
- media
- customer support

# Speech recognition in the Wild - Still Challenging



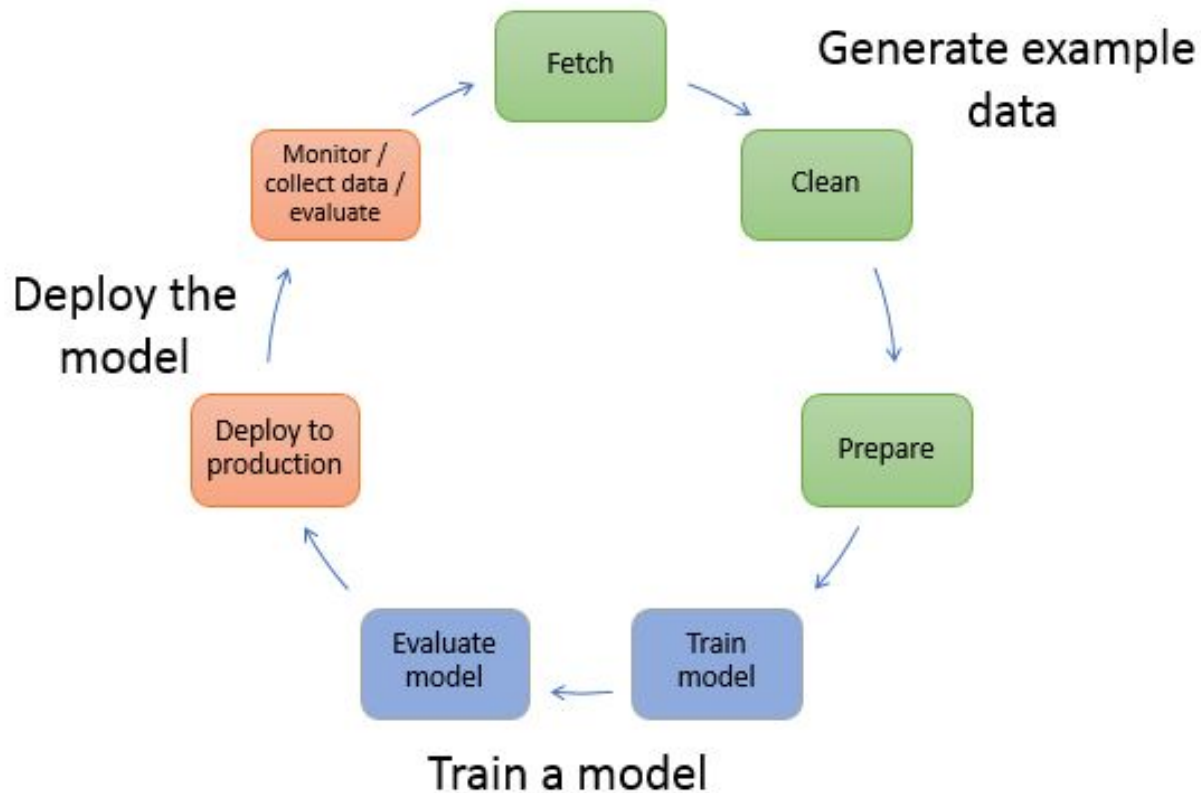# There is lots of variabilities in speech

- accent
- noise
- emotion modulated speech
- mis-pronunciation
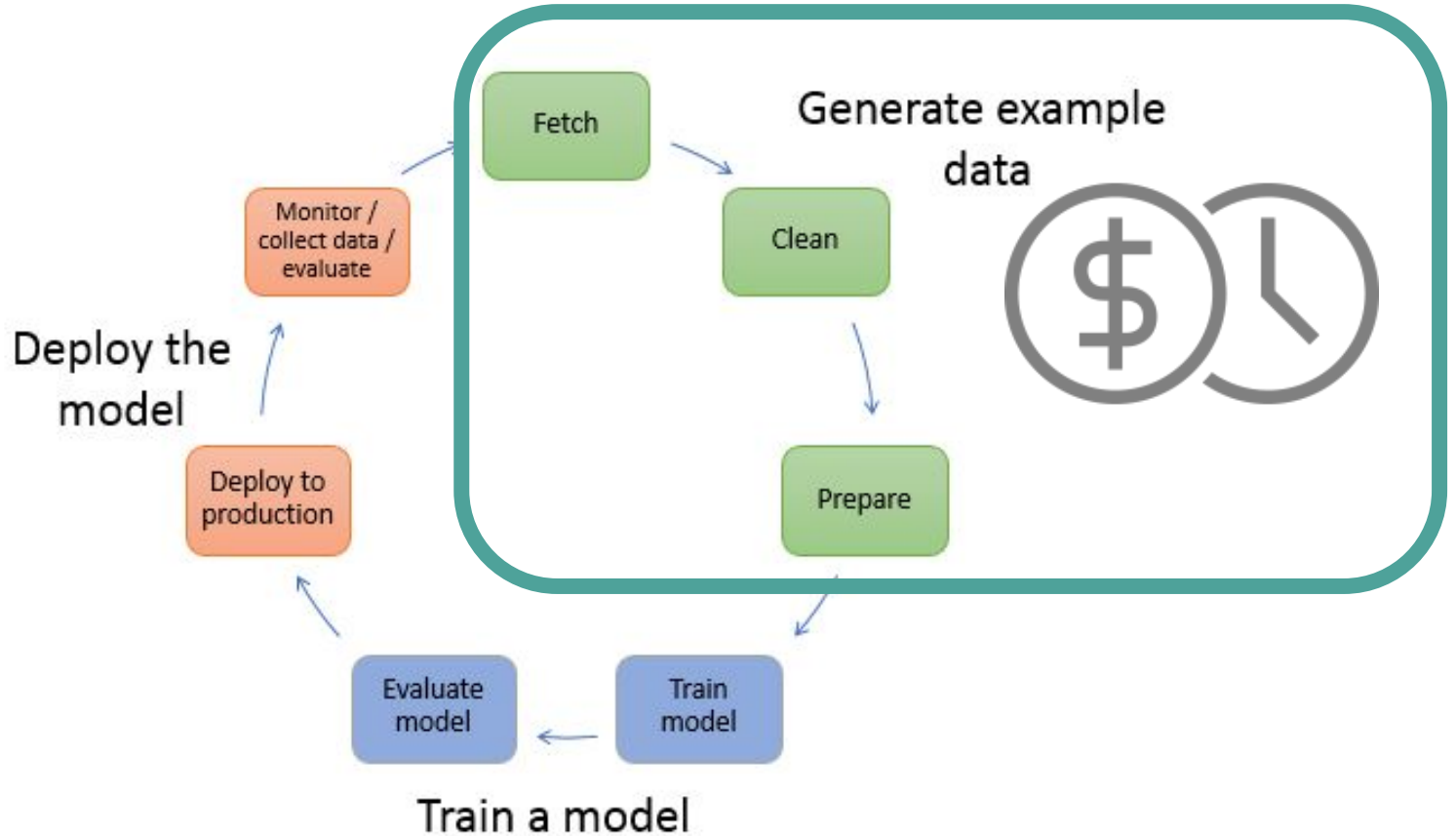
# Speech recognition in the Wild - Still Challenging



# It's hard to learn all beforehand!
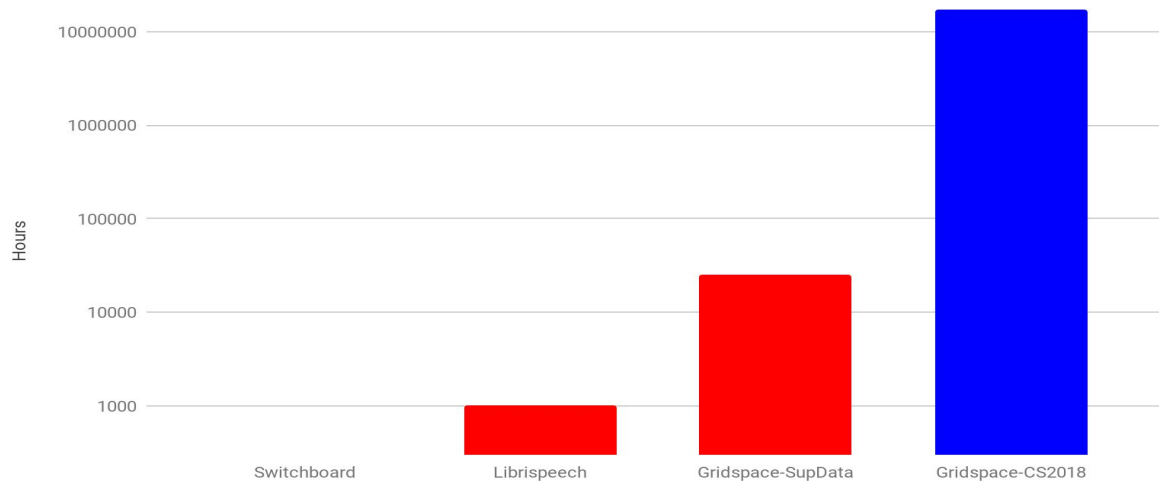
# Speech recognition in the Wild - in Practice

# Speech recognition in the Wild - in Practice

# Use of Unsupervised Data

Hours vs. Data



| | SWBD | Librispeech | GS-SupData | GS-CS2018 |
|---|---|---|---|---|
| Hours | 0.3k | 1k | 30k | 17M (2000 years) |
| Data | Supervised | Supervised | Supervised | Unsupervised |

# Use of Unsupervised Data
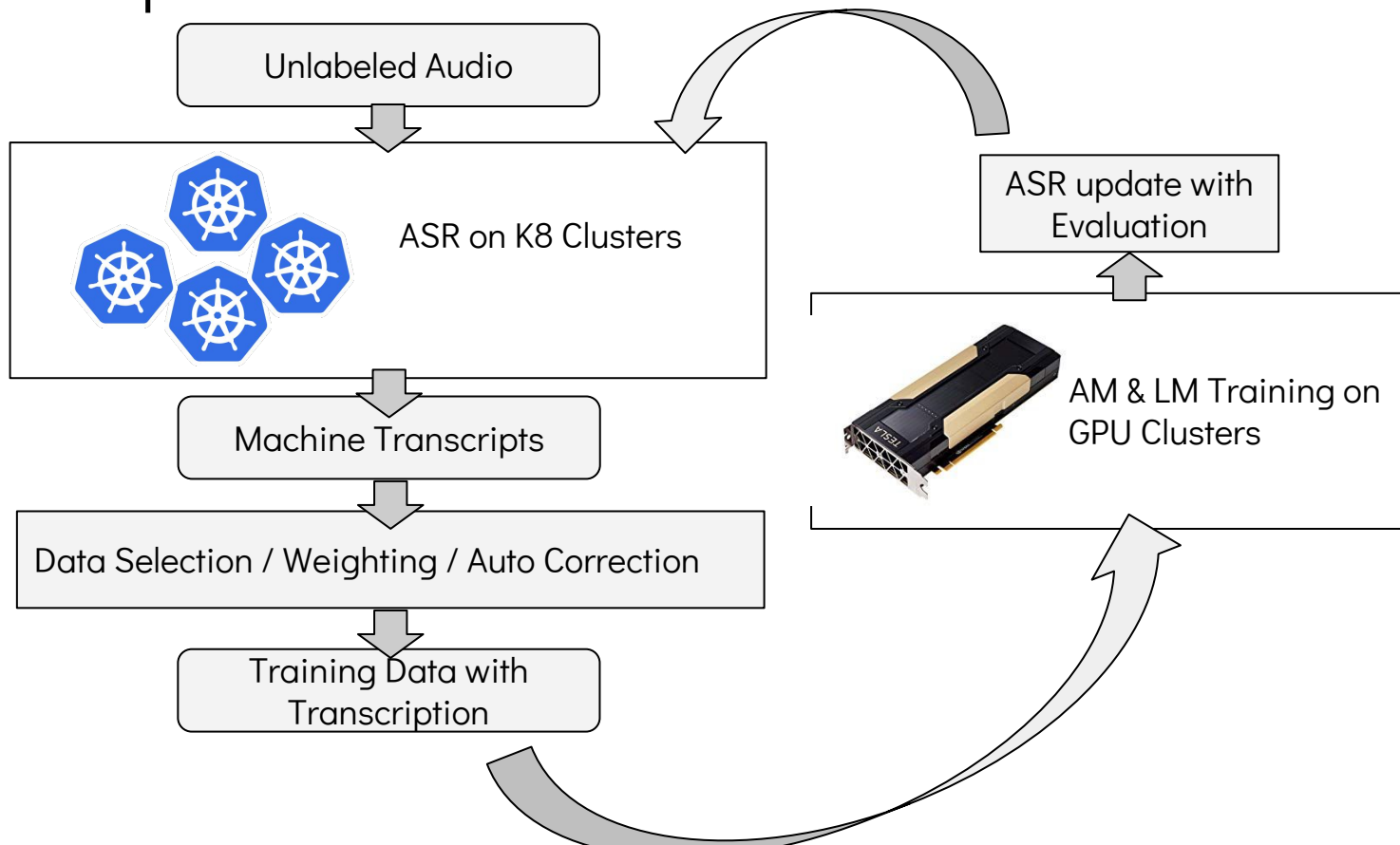
## Semi-supervised training(SST)

- a good way to use unsupervised data for supervised tasks

- It has to deal with uncertainties

- We can update AM and LM iteratively

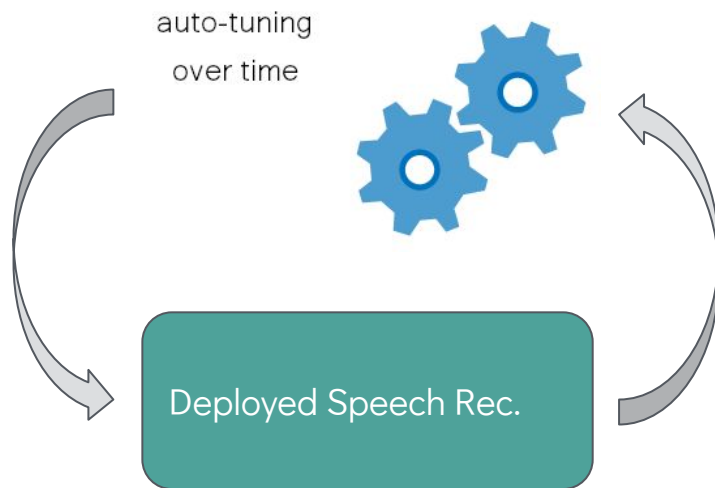# Use of Unsupervised Data

## Data selection for SST

- 100% use(no filtering) of Unsupervised data would cause model's degradation on accuracy

- knowledge based selection helps

  - confidence score, length, topic, speaker info

# Use of Unsupervised Data



Unlabeled Audio

ASR on K8 Clusters

ASR update with Evaluation

Machine Transcripts

AM & LM Training on GPU Clusters

Data Selection / Weighting / Auto Correction

Training Data with Transcription

# Continuous Learner

- accent learning
- noise learning
- language/grammar learning

auto-tuning
over time

Deployed Speech Rec.

# ... and become better learner over time

# Continuous Learner

## Throughput

| | P-100 | V-100 |
|---|---|---|
| SST training / 1 sec / 1 gpu | 450 sec | 580 sec |
| SST training / 24 hours / 1 gpu | 10800 hours | 13920 hours |

* Training for Acoustic Model(Resent TDNN), 150 frames per example, 64 example per minibatch
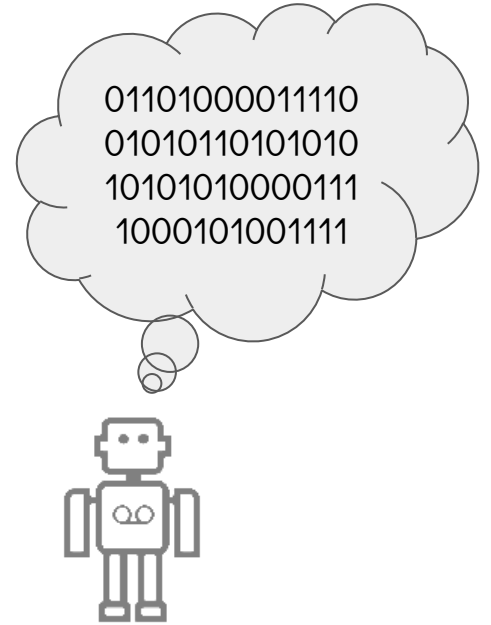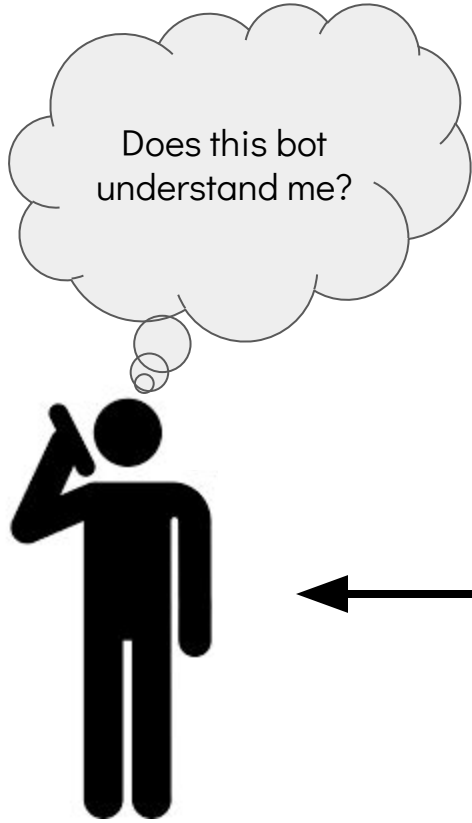
# Continuous Learner



( Word Error Rate )

Call Volume (days of audio)

Time

Ideal Error Graph?

# Continuous Learner

Experimental Results

|  | Unmatched Domain | Supervised Training | Semi-Sup Training |
|---|---|---|---|
| Word Error Rate | 14.29% | 9.52% | 7.83% |
| System Building Hours | | 3 months | < 1 DAY |

# CASE STUDY: STOCHASTICITY AND LATENCY IN SPEECH SYNTHESIS INFERENCE

# Dialog System Evaluators
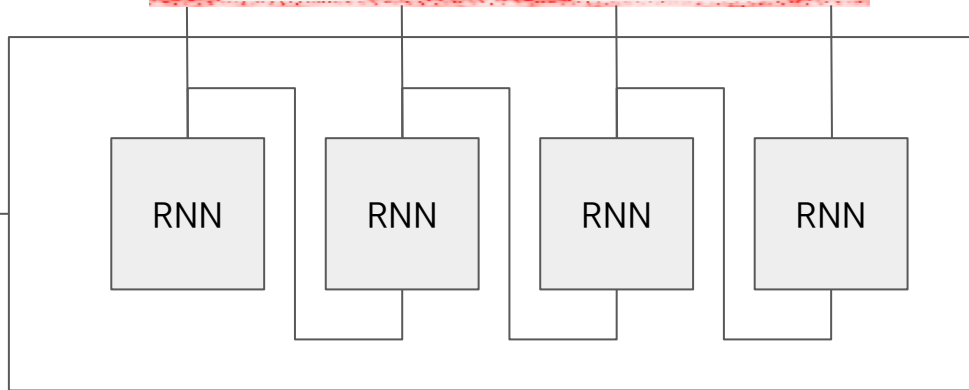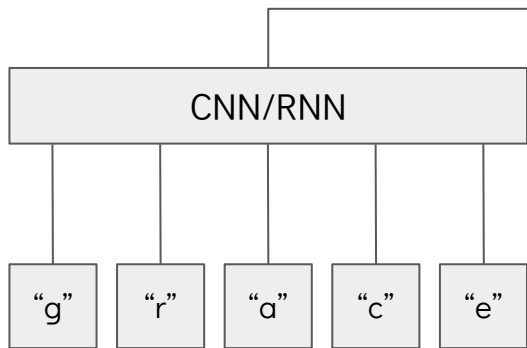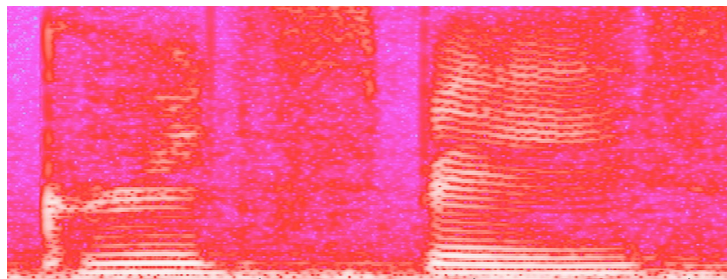
Content

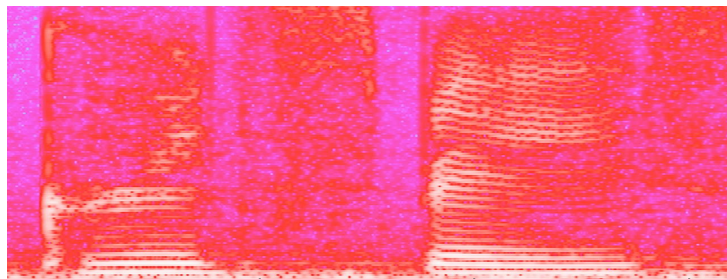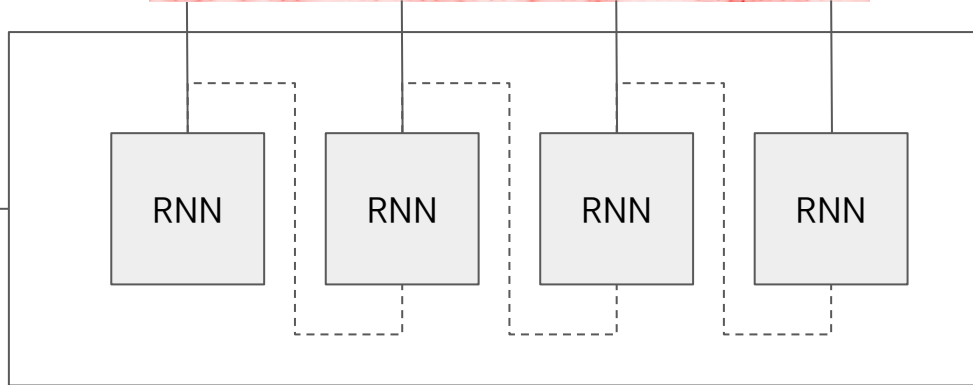Naturalness

Cadence

# Influenced by TTS
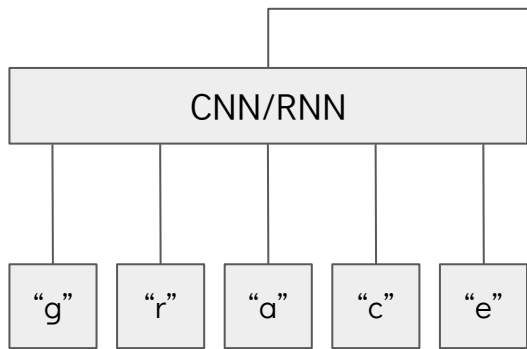
Content

Naturalness

Cadence
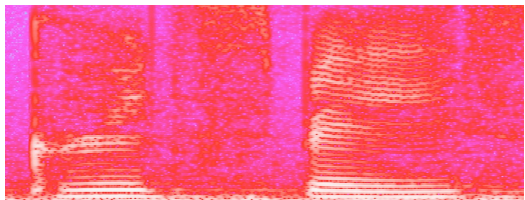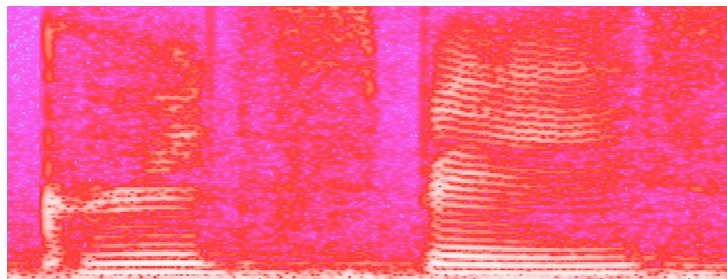
# Modern Neural TTS

"I'm speech that came from a big neural network"

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⌛ | Stochasticity 🎲 | Scale ⚖ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity | Interactivity |

"g"  "r"  "a"  "c"  "e"

CNN/RNN

RNN  RNN  RNN  RNN

Hsu, Wei-Ning, et al. "Hierarchical generative modeling for controllable speech synthesis." *arXiv preprint arXiv:1810.07217* (2018).

# Influenced by TTS

Content

Naturalness

Cadence

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | Latency ⧗ | Stochasticity 🎲 | Scale ⚖ |
| **Manage** | Jitter & Aperiodicity 🌊 | Nonlinearity ▢ | Interactivity ▢ |

# Latency



Imperceptible     Acceptable     Unusable

0ms     150ms     500ms     1000ms

# Where does the latency come from?

# WaveNet Vocoding



Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

Oord, Aaron van den, et al. "Wavenet: A generative model for raw audio."
*arXiv preprint arXiv:1609.03499* (2016).

# WaveNet

## Training

## Inference

# Vocoder Evolution

WaveNet
2016

→

Parallel
WaveNet
2017

→

WaveGlow
2018

# Vocoder Evolution

```
┌─────────────┐      ┌─────────────┐      ┌─────────────┐
│             │      │   Parallel  │      │             │
│   WaveNet   │ ───> │   WaveNet   │ ───> │   WaveGlow  │
│    2016     │      │    2017     │      │    2018     │
│             │      │             │      │             │
└─────────────┘      └─────────────┘      └─────────────┘
```
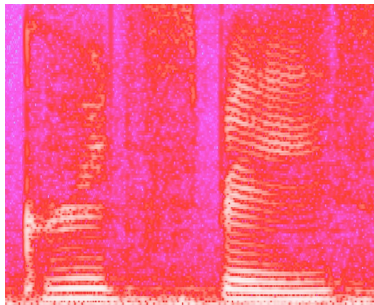
# WaveGlow

## Training

## Inference

```python
from tensorflow.python.compiler.tensorrt.trt_convert import TrtGraphConverter

converter = TrtGraphConverter(
    input_saved_model_dir='my_saved_model',
    precision_mode=FP16
)
converter.convert()
converter.save('trt_saved_model')
```

# Tesla V100 Latency (batch size=1)

| Precision | Latency /ms | Samples per second /Hz | Speed Up |
|-----------|-------------|------------------------|----------|
| FP32 | 277 | 520 | 1x |
| FP16 (TRT) | 196 | 735 | 1.4x |

# Can we go faster?

# INT8 Calibration



2^31 - 1                                    -2^31 + 1

128                          -128

```python
from tensorflow.python.compiler.tensorrt.trt_convert import TrtGraphConverter

converter = TrtGraphConverter(
    input_saved_model_dir='my_saved_model',
    precision_mode='INT8'
)
converter.convert()
converter.calibrate(
    fetch_names=['output:0'],
    num_runs=1000,
    input_map_fn=get_examples
)
converter.save('trt_saved_model')
```
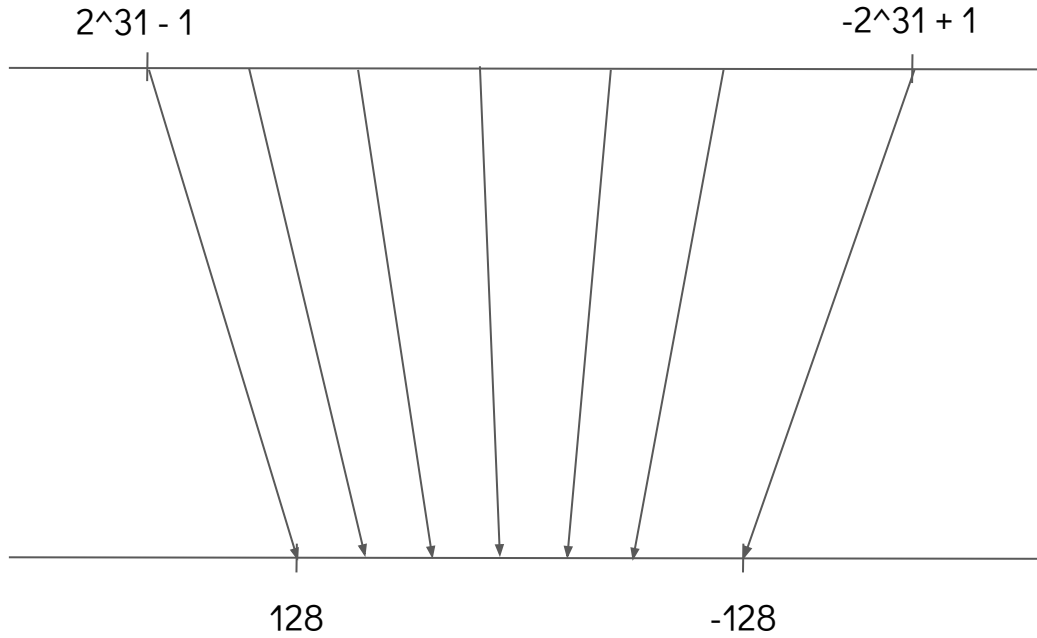
# Tesla V100 Latency (batch size=1)

| Precision | Latency /ms | Samples per second /Hz | Speed Up |
|---|---|---|---|
| FP32 | 277 | 520 | 1x |
| FP16 (TRT) | 196 | 735 | 1.4x |
| INT8 (TRT) | 164 | 878 | 1.7x |

# Deployment



Google Container Engine

Tensorflow Serving

|  | Time | Predictability | Complexity |
|---|---|---|---|
| **Optimize** | ✓Latency ⌛ | ✓Stochasticity 🎲 | ✓Scale ⚖️ |
| **Manage** | ✓Jitter & Aperiodicity 🏄 | ✓Nonlinearity🦋 | ✓Interactivity 🤖 |

# Enables...

THE FOLLOWING IS A REAL INTERACTION

Access the demo video:
[Video on Vimeo Here](#)

Want to demo Grace at GTC?
Do you operate a call center?
Do you need speech processing or automation?
**gtc@gridspace.com**

Do you want to work at Gridspace?
**hiring@gridspace.com**

Thank you everyone!