Tencent AI Lab

# Encrypted Speech Recognition
# using deep polynomial networks

Austin Zhang, **Dong Yu**

Tencent AI Lab

Yifan Gong

Microsoft

## Table of contents
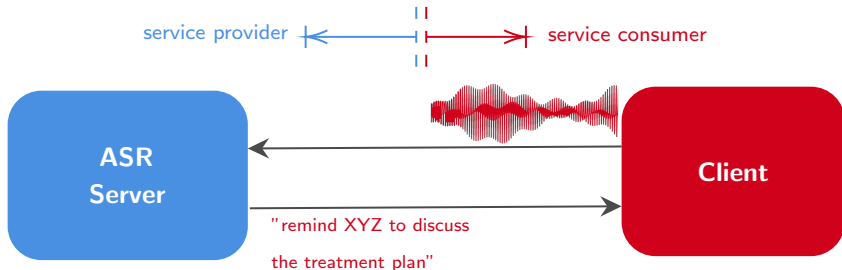
# Motivation
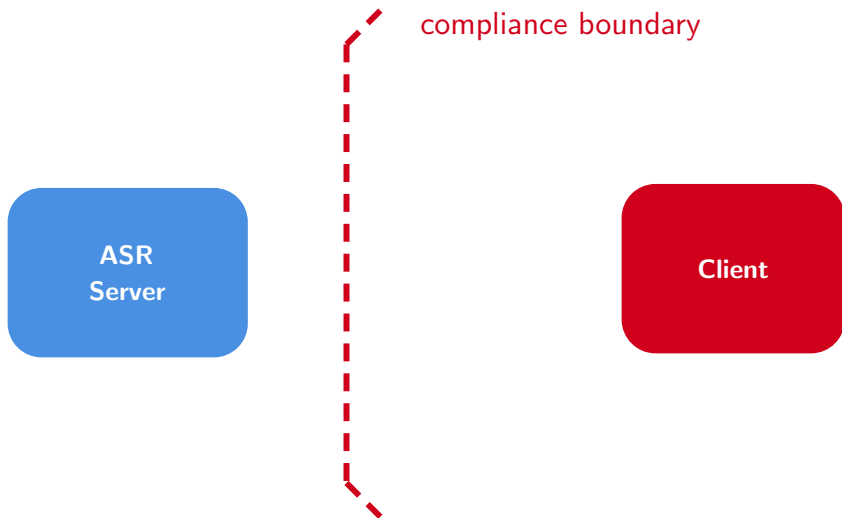
- State-of-the-art speech recognition services are running on cloud
- However, this will leak the client's private information to the server. e.g., medical/financial/enterprise/sensitive data
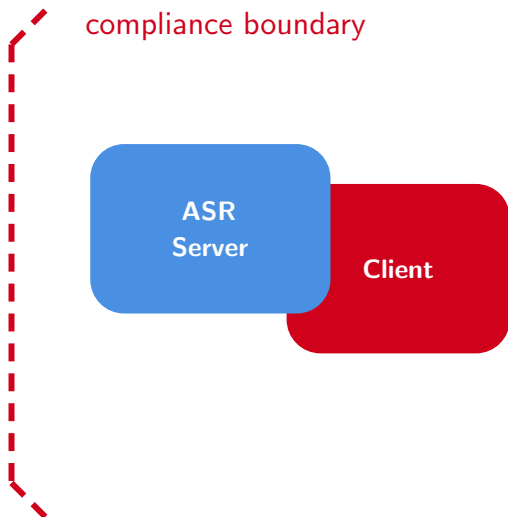
- Step 1: build a compliance boundary → prevent data to leak out

compliance boundary

ASR
Server

Client

- Step 2: move ASR service inside of compliance boundary



compliance boundary

ASR
Server

Client

# How to protect privacy? — A private-cloud solution

- issues 1: hard to deploy an update to the private cloud
- issues 2: costly for some small business/individual users
- issues 3: service provider may divulge the model and decoder to the service consumer who may resale to others
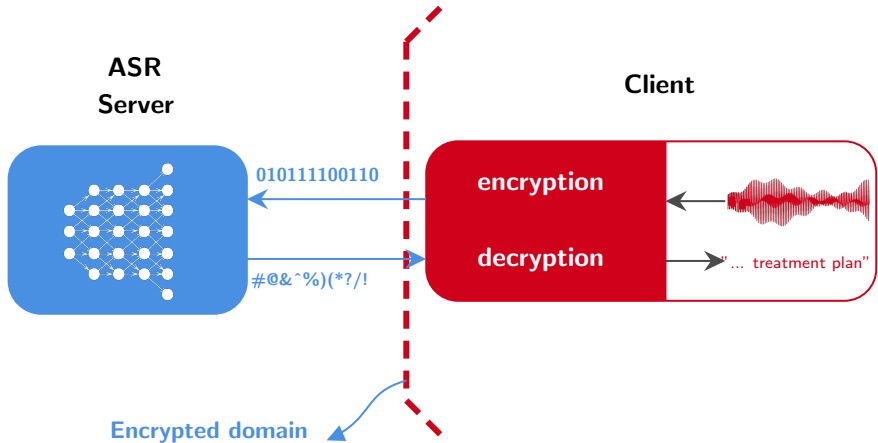


compliance boundary

ASR Server
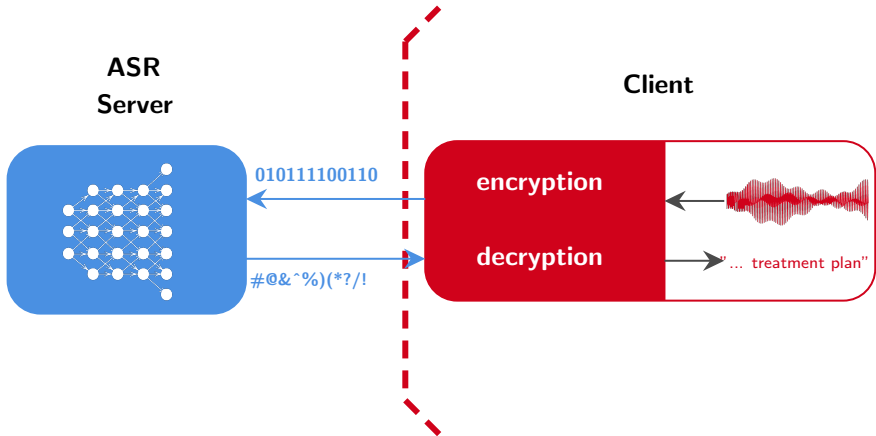
Client

# Encrypted Speech Recognition

- Ideally, we want ...



ASR Server

Client

010111100110

encryption

decryption

#@&^%)(*?/!

"... treatment plan"

Encrypted domain

- Does this encryption exist?

## Homomorphic Encryption

$\mathbf{E}^{-1}\Big[f\big(\mathbf{E}[x]\big)\Big] \equiv f(x) \rightarrow$ An elegant solution for all above questions

Homomorphic Encryption

$\mathbf{E}[x] \longleftarrow\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim \quad x$

$f(\ ) \Big\Downarrow \qquad\qquad\qquad\qquad\qquad f(\ ) \Big\Downarrow$

Decryption    $\mathbf{E}^{-1}[\ ]$

$f(\mathbf{E}[x]) \sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\sim\hspace{-0.5em}\longrightarrow f(x)$

$$x = \binom{7}{3}$$

$$f(\ ) = x_1 \times x_2$$

$$f(x) = 21$$

# Homomorphic Encryption — example



Encryption withpub.key(23,143)

$$7^{23} \bmod 143 = 2$$

$$\binom{2}{126} \longleftarrow \cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots \longrightarrow x = \binom{7}{3}$$

$$3^{23} \bmod 143 = 126$$

# Homomorphic Encryption — example



Encryption with pub. key (23,143)

$7^{23} \mod 143 = 2$

$\binom{2}{126} \longleftarrow x = \binom{7}{3}$

$3^{23} \mod 143 = 126$

$f(\ ) = x_1 \times x_2$

252

# Homomorphic Encryption — example



Encryption withpub.key(23,143)

$7^{23} \bmod 143 = 2$

$\left(\begin{smallmatrix} 2 \\ 126 \end{smallmatrix}\right) \quad \longleftarrow \quad x = \left(\begin{smallmatrix} 7 \\ 3 \end{smallmatrix}\right)$

$3^{23} \bmod 143 = 126$

$f(\ )=x_1 \times x_2$

$252 \quad \longrightarrow \quad f(x) = 21$

$252^{47} \bmod 143 = 21$

Decryption withprivate.key(47,143)

# Homomorphic Encryption — example



model privacy

data privacy

$\binom{2}{126}$ $\xleftarrow{\hspace{1cm}}$ Encryption $\hspace{1cm}$ $x = \binom{7}{3}$

$f(\ ) = x_1 \times x_2$

252 $\xrightarrow{\hspace{1cm}}$ Decryption $\hspace{1cm}$ $f(x) = 21$
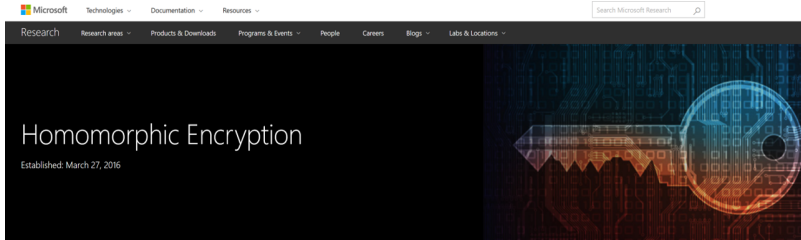
- Only AM scores are computed on server side.
- Original result guaranteed after decryption.
- No need to retrain the DNN on encrypted data.

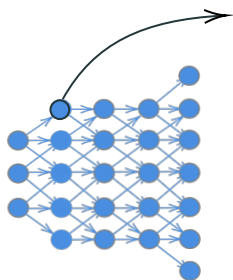It is extremely slow and not feasible before, until …



Microsoft researchers smash homomorphic encryption speed barrier!

- But $f(\cdot)$ must be polynomial
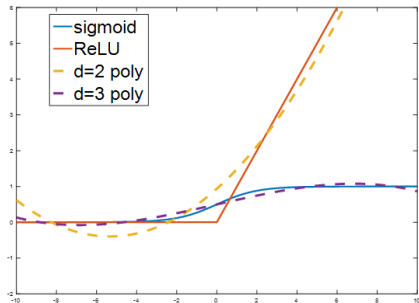- must be fixed point operation
- open source $\rightarrow$ http://sealcrypto.org/

# Deep Polynomial Network
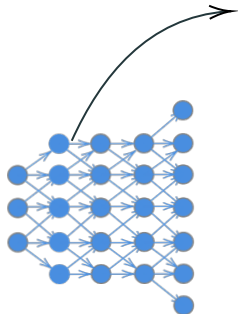
# Deep polynomial network



Replace ReLU and Sigmoid as polynomials

- unbounded polynomial approximation $\rightarrow$ batch norm is a must.

## Deep polynomial network

**Dense layer** (polynomial)
$$\mathbf{E}[\mathbf{W}]^\top \mathbf{E}[\mathbf{x}] \xrightarrow{\mathbf{E}^{-1}} \mathbf{W}^\top \mathbf{x}$$

**Convolution layer** (polynomial)
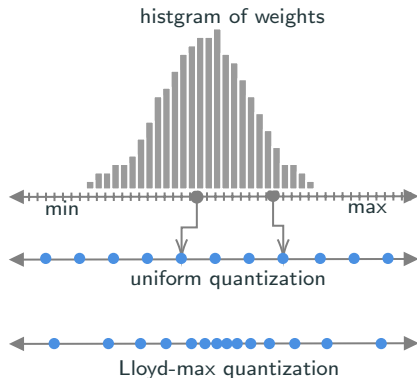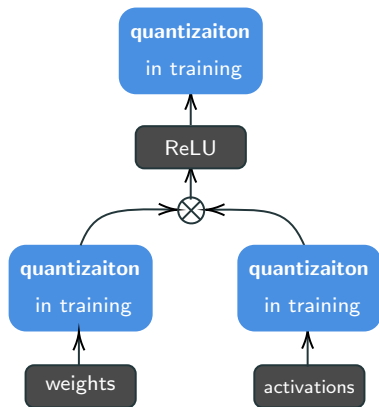
**Batch norm** (merged to dense layer)
$$\mathbf{W}^\top \left( \gamma \frac{\mathbf{z}-\mu}{\sigma} + \beta \right) + \mathbf{b} = \mathbf{W'}^\top \mathbf{z} + \mathbf{b'}$$

**Max pooling layer** (approximate)
$$\max(x_1, \ldots x_n) = \lim_{d \to \infty} \left( \sum_{i=1}^{n} x_i^d \right)^{\frac{1}{d}}$$

low-bit model is critical for encryption speed
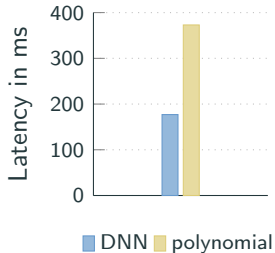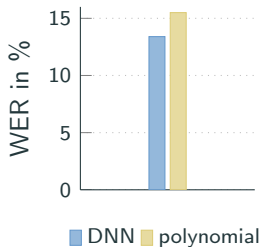
# Experimental Results

## WER on Switchboard

| WER in % | | 16-bit | 8-bit | 4-bit | 2-bit |
|---|---|---|---|---|---|
| DNN | quantized train | 14.7% | 14.7% | 14.9% | 30.3% |
| | $\rightarrow$ polynomial | 15.8% | 15.8% | 16.1% | 30.8% |
| CNN | quantized train | 12.2% | 12.3% | 12.7% | – |
| | $\rightarrow$ polynomial | 13.5% | 13.6% | 14.0% | – |

- with proper quantized training, 4-bit is sufficient.
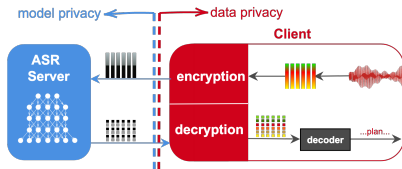- the polynomial networks increase WERs by a little as a cost.

# WER and Latency on Cortana Task

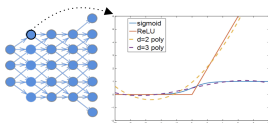| | 16-bit | 4-bit | avg. latency per utterance | | |
| --- | --- | --- | --- | --- | --- |
| | | | encryption | decryption | overall |
| DNN | 12.9% | 13.4% | – | – | 177ms |
| polynomial | 14.8% | 15.5% | 202ms | 16ms | 373ms |

- a framework that enables privacy-preserving speech recognition



- a polynomial network that can make predictions over the encrypted speech in real time.



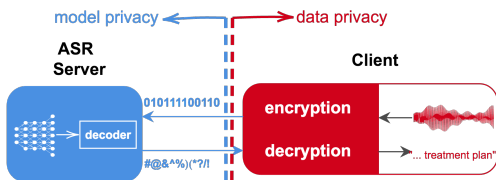- with quantized training, 4-bit is sufficient for DNN/CNN.

# Thanks.

# Questions?

Tencent AI Lab

# Future work

- make the decoder also work on encrypted domain, so that we could run everything on the cloud.



- investigate training on encrypted data so that multiple parties (e.g.Microsoft, Google and Amazon) can encrypt and combine their data together to train models without sacrificing users privacy.