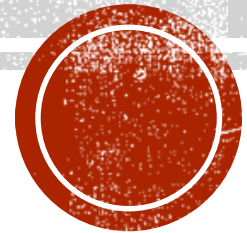# AN OVERVIEW OF SPOKEN LANGUAGE UNDERSTANDING

Xuedong Huang
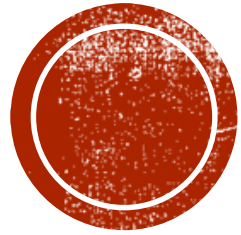
Chief Speech Scientist & Distinguished Engineer

Microsoft Corporation

xdh@microsoft.com

An **invisible** revolution is coming

# ARE WE READY?

Cloud-enabled multimodal NUI with speech, gesture, gaze…

# TODAY'S STATE OF THE ART: CORTANA

# TODAY'S STATE OF THE ART: SKYPE TRANSLATOR

# SPEECH RECOGNITION – APPROACHING HUMAN PARITY

BY XUEDONG HUANG, JAMES BAKER, AND RAJ REDDY

# A Historical Perspective of Speech Recognition

WITH THE INTRODUCTION of Apple's Siri and similar voice search services from Google and Microsoft, it is natural to wonder why it has taken so long for voice recognition technology to advance to this level. Also, we wonder, when can we expect to hear a more human-level performance? In 1976, one of the authors (Reddy) wrote a comprehensive review of the state of

http://cacm.acm.org/magazines/2014/1/170863-a-historical-perspective-of-speech-recognition

» key insights

■ The insights gained from the speech recognition advances over the past 40 years are explored, originating

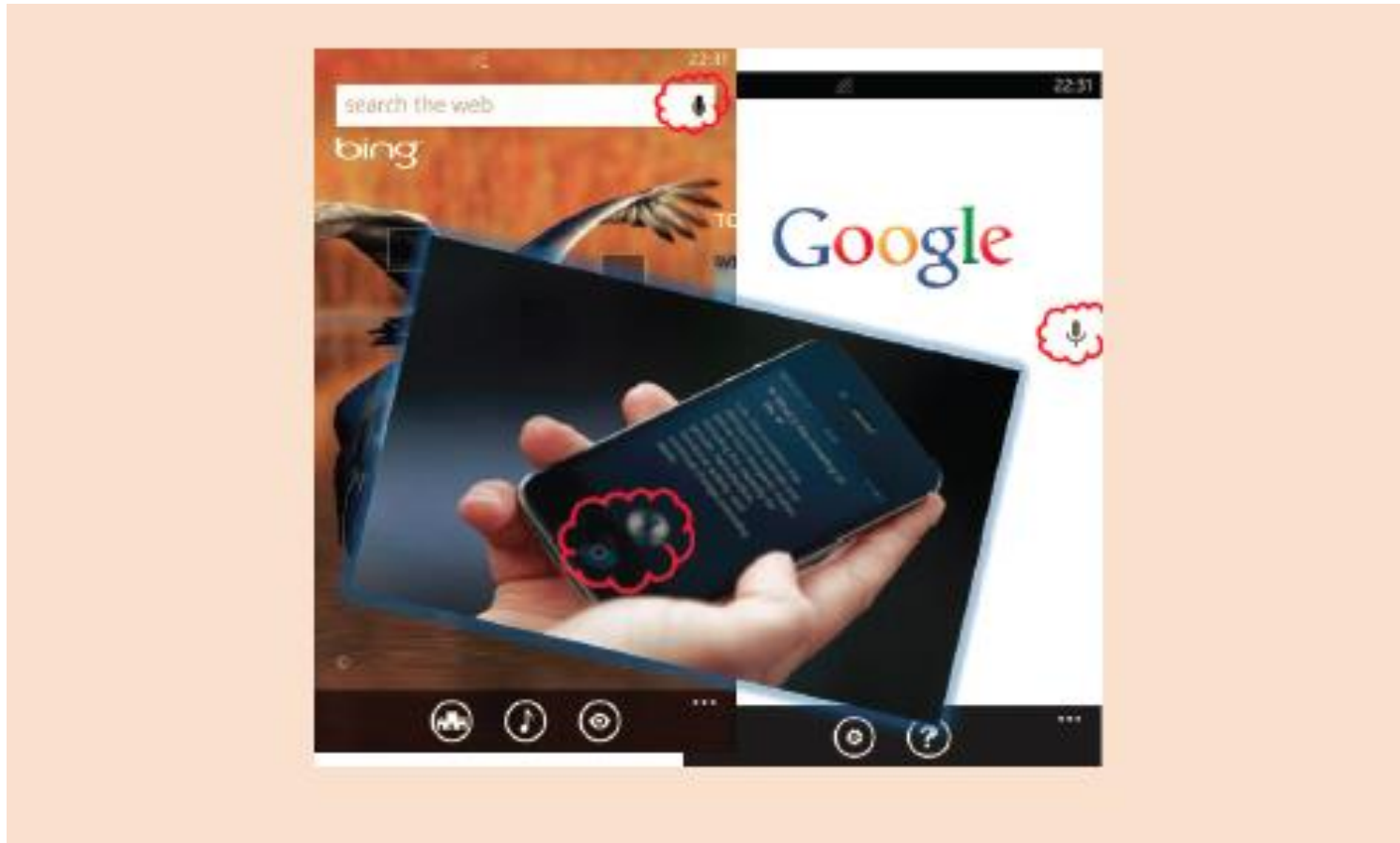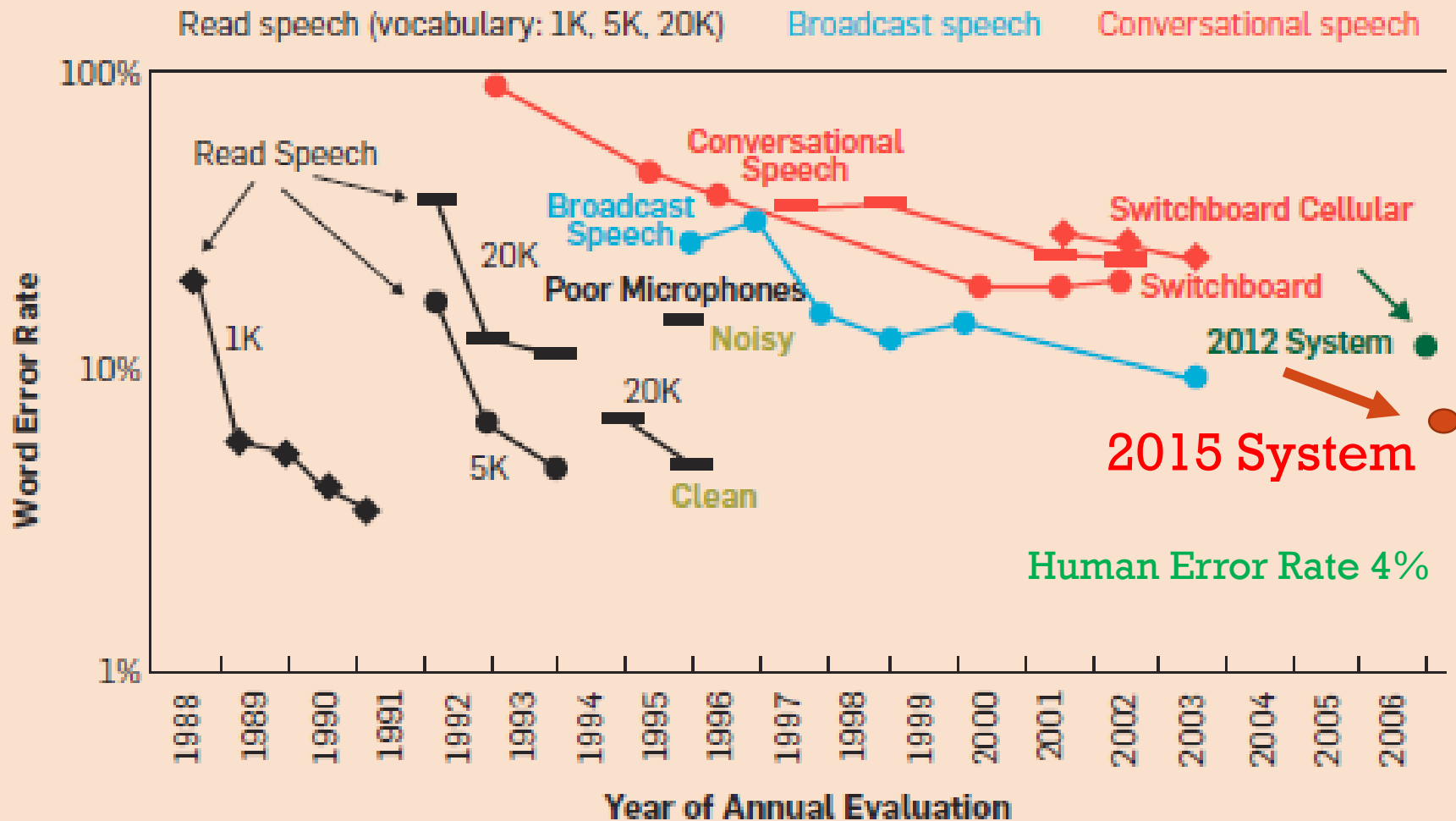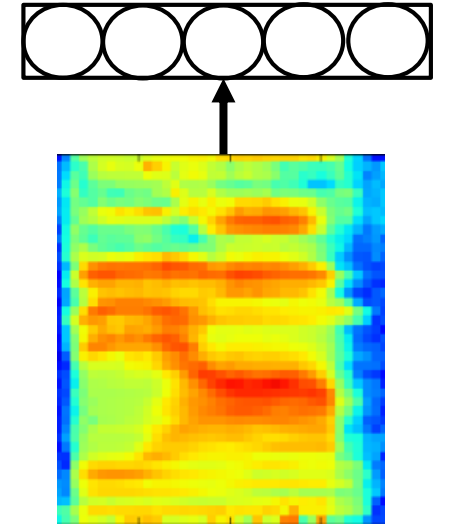# SECRET SAUCE: BIG DATA + MACHINE LEARNING + INFRASTRUCTURE

Figure 1. Historical progress of speech recognition word error rate on more and more difficult tasks.[10] The latest system for the switchboard task is marked with the green dot.
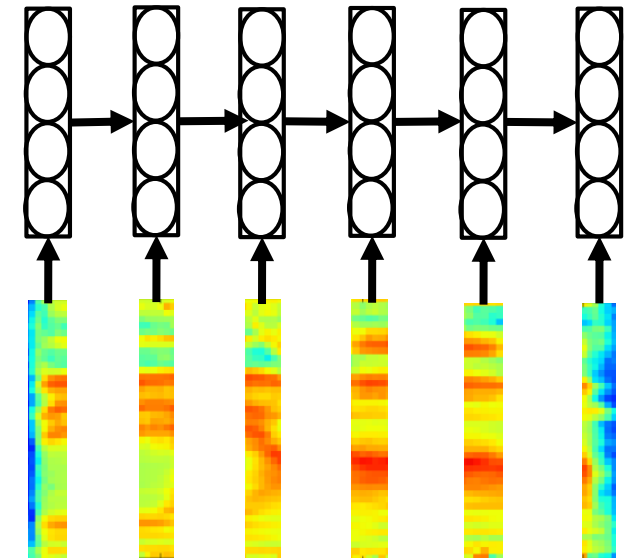
# BEYOND FF-DNNS

- Speech is a sequential process while FF-DNNs are not sequential in nature.

- FF-DNNs are not efficient in modeling temporal/spectral compression/stretch.

- Theoretically, recurrent models should better model speech data.
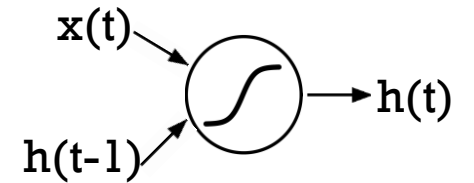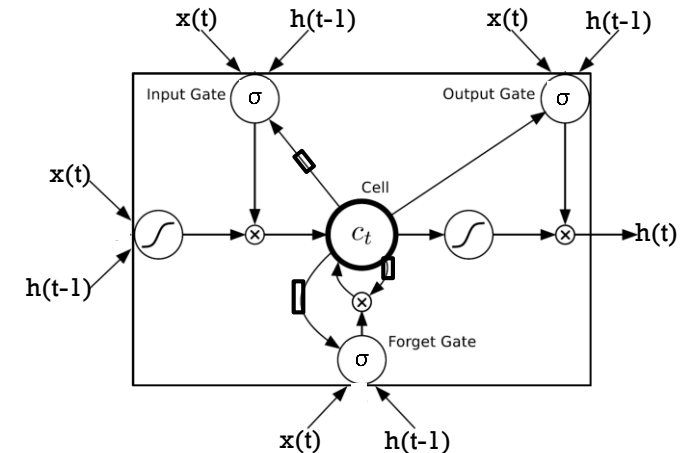  - In RNNs, hidden representations are conditioned on all previously seen frames.

(a) FF-DNN

(b) RNN

# FROM SIMPLE RNNS TO LSTM RNNS

▪ At each time step, each hidden node in an RNN consumes the current input vector and the previous hidden vector.



(a) RNN node

▪ RNNs have a limited capacity to remember important distant past events. (Vanishing gradient problem)

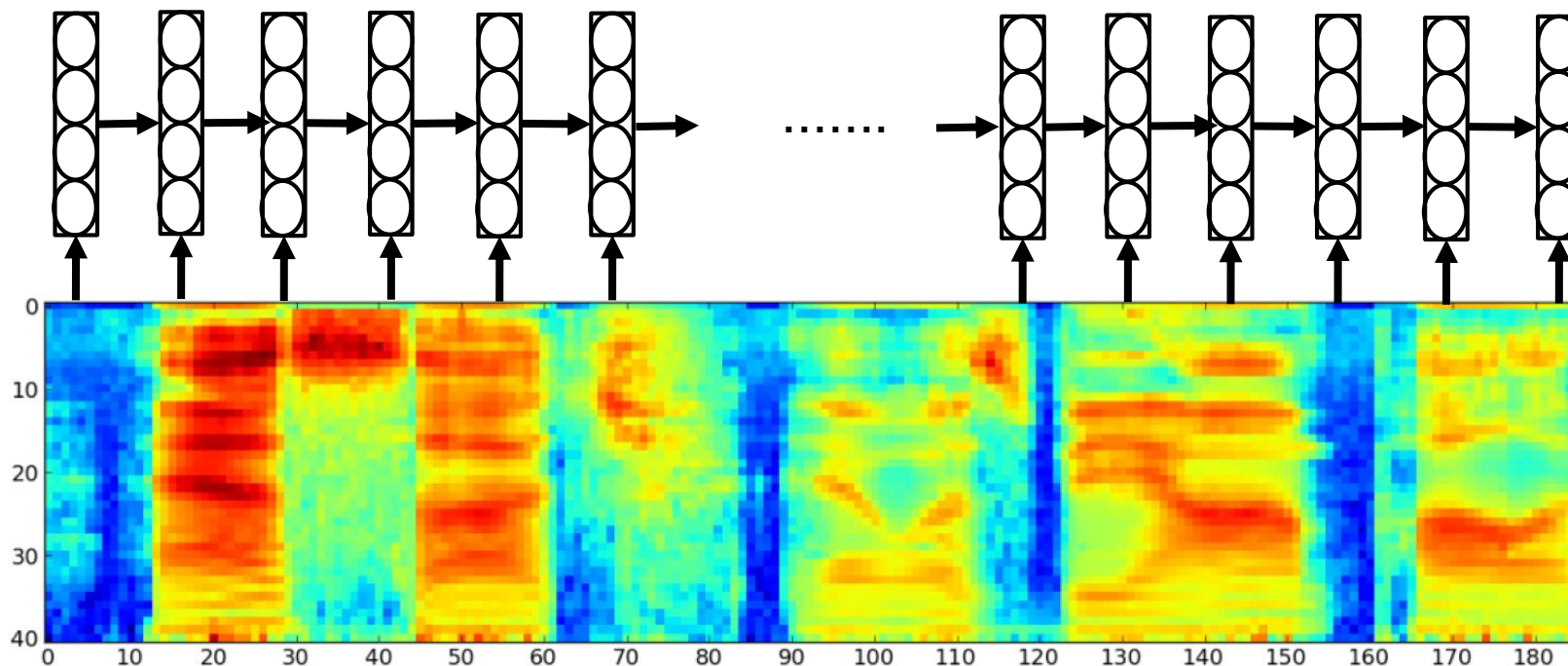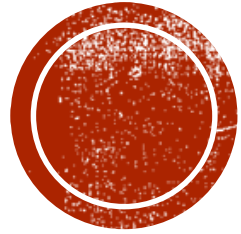▪ LSTM input, output, and forget gates combat vanishing gradient problem.



(b) LSTM node

11

# ACOUSTIC MODELING USING LSTMS

- Models long-span phenomena well

- Good performance improvements observed

- But less convenient than fixed window methods

# UNDERSTANDING LANGUAGE — MUCH HARDER

Machine learning enables nearly every value proposition of web search.

# CORTANA

*A personal assistant built around you*

Cortana gets to know you over time, building a relationship with you that's based on trust.

She tracks the stuff you care about, looks out for you throughout the day, and helps filter out the noise so you can stay on top of what matters.

Throughout, Cortana is delightful and easy to use.

# CORTANA – CORE PILLARS

## 1 PERSONAL

**Cortana…**

…is your truly personal assistant

…gets to know you

…is transparent

**Proof points:**

- Learning
- Notebook
- Personal suggestions
- Transparency & control

## 2 LOOKS OUT FOR YOU

**Cortana…**

…looks out for you

…filters out the noise

…reminds you of what's important

**Proof points:**

- Useful and relevant alerts
- Planners
- Event scheduling
- Quiet hours and inner circle
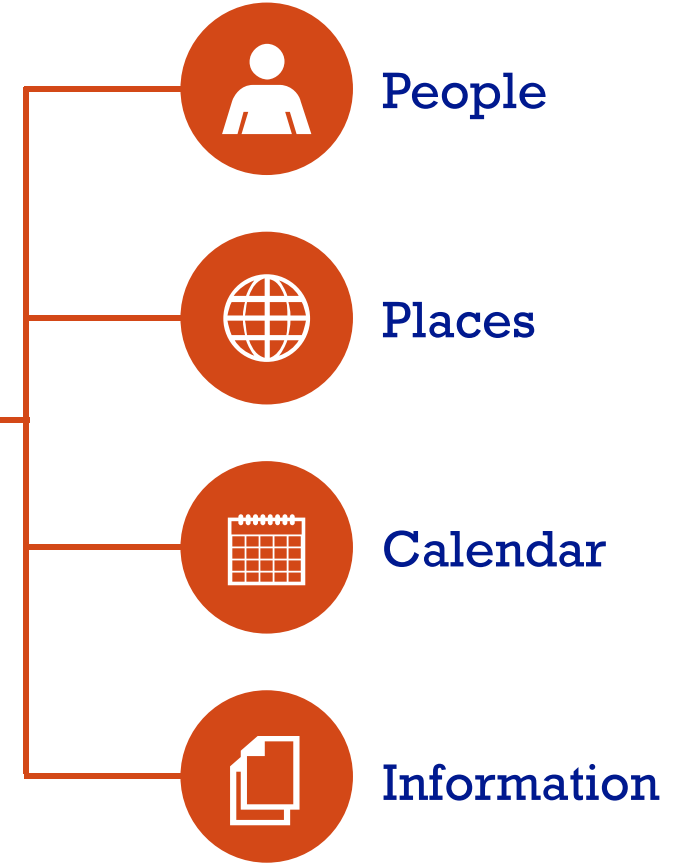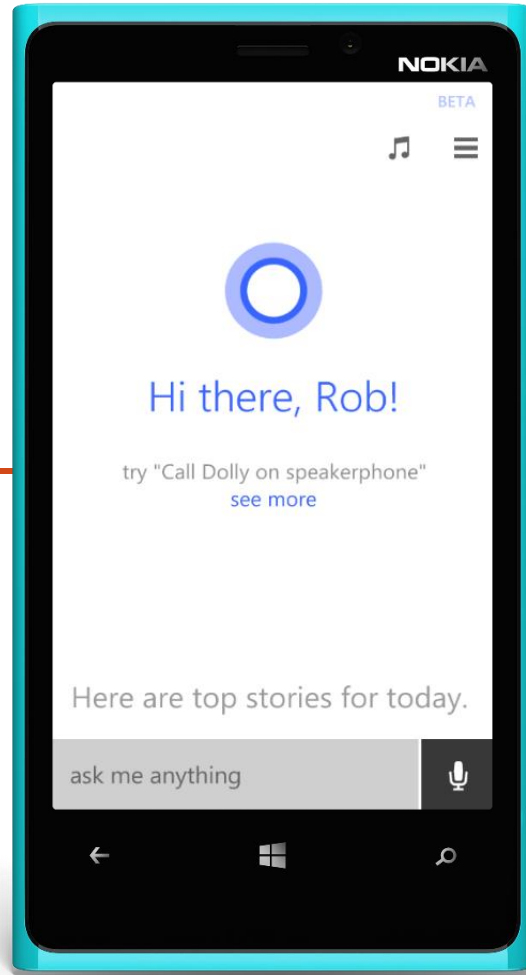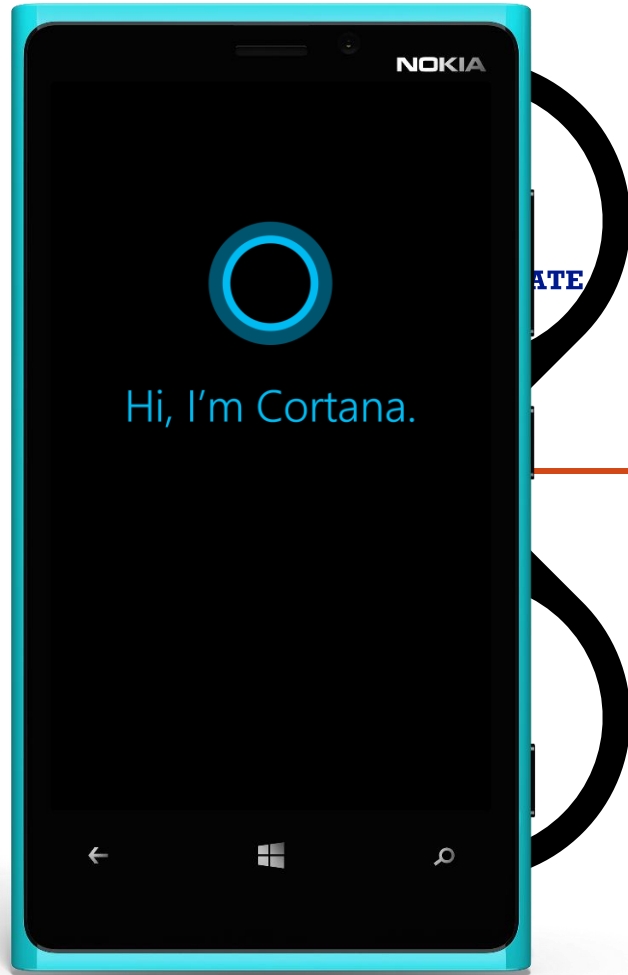- Reminders

## 3 DELIGHTFUL & EASY TO USE

**Cortana…**

…"just works"

…lets you interact on your terms

…has a fun & engaging personality

**Proof points:**

- Voice & natural language
- Text input
- Personality (visual, spoken voice, and behavior)
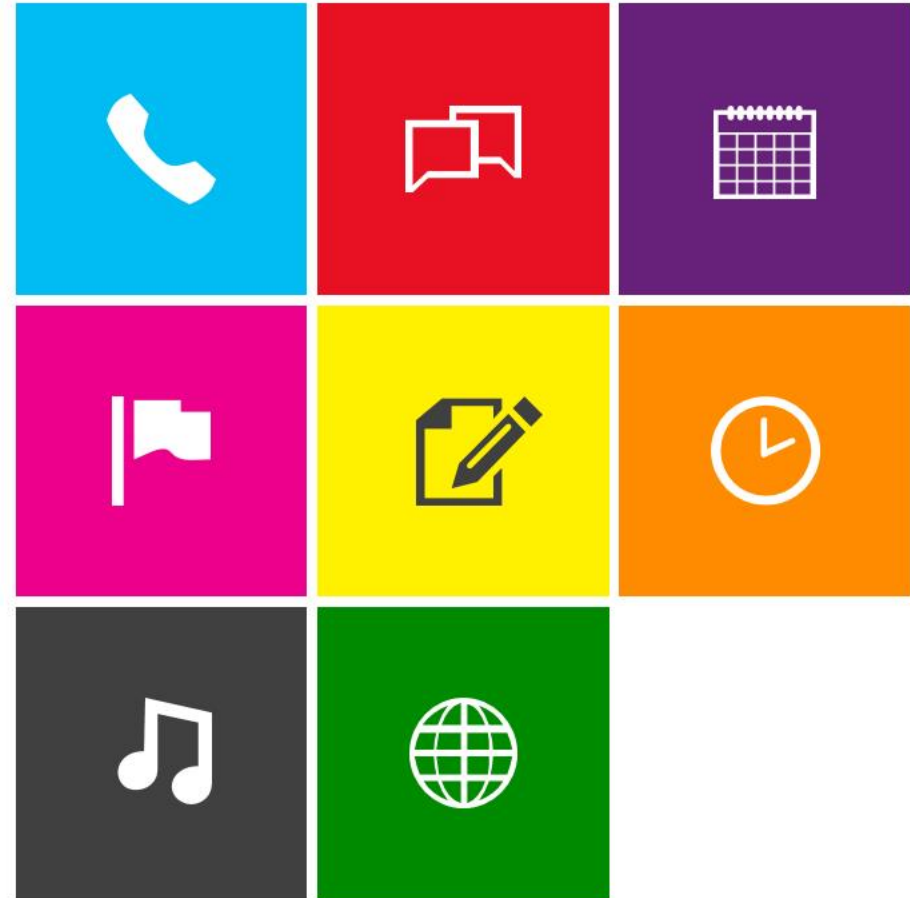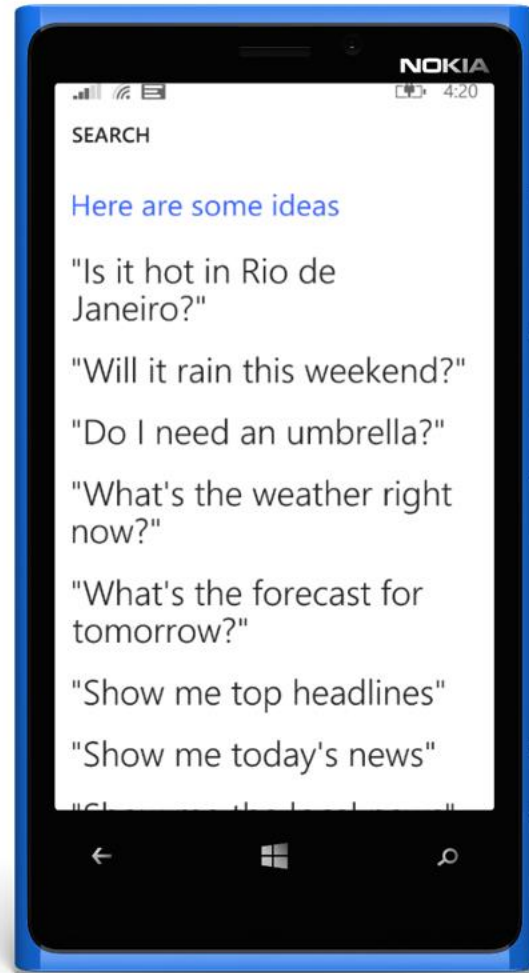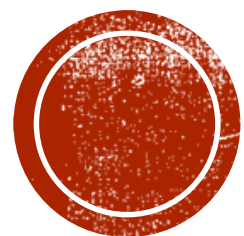
# Cortana: What Can You Do?



- People
- Places
- Calendar
- Information

# Cortana: What Can You Do?

- 📞 Phone
- 🗨 Messaging
- 📅 Calendar
- 🚩 Reminders
- 📝 Notes
- 🕐 Alarms
- 🎵 Music
- 🌐 Places
- 🔍 Search

PROJECT OXFORD

# MICROSOFT PROJECT OXFORD SERVICES

**Compute APIs**

**Face APIs**

**Speech APIs**

Emotion APIs

*BETA*

Understand your users with Emotion Recognition

Spell Check APIs

*BETA*

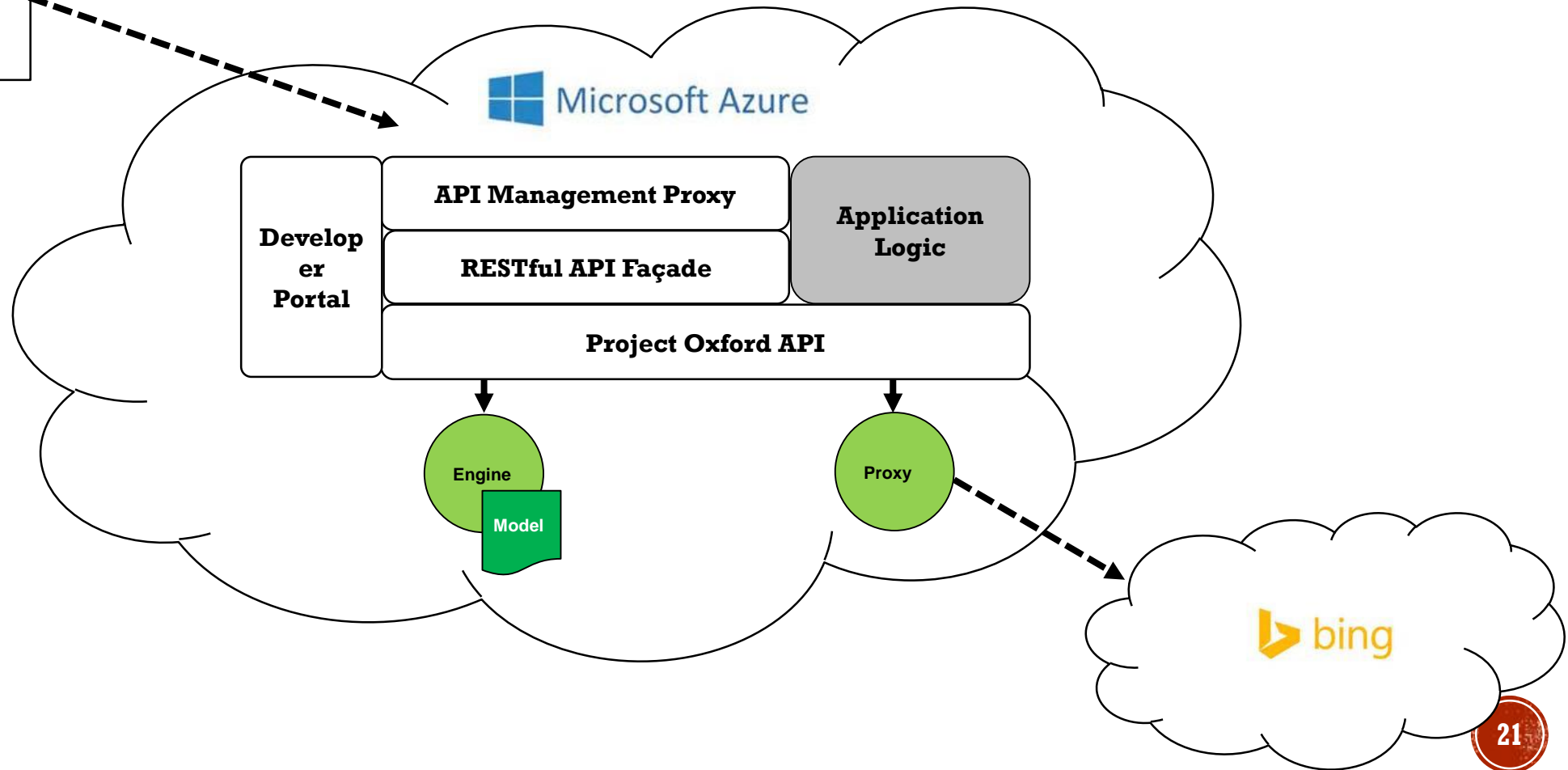Detect and correct common and uncommon spelling errors, via the Bing document index
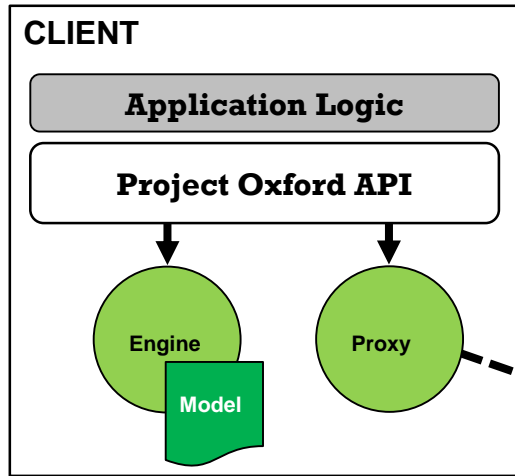
Language Understanding Intelligent Service (LUIS)

*BETA*

Understand natural language commands tailored to your application

# ARCHITECTURE OVERVIEW

**CLIENT**

Application Logic

Project Oxford API

Engine

Model

Proxy

Microsoft Azure

Developer Portal

API Management Proxy

RESTful API Façade

Application Logic

Project Oxford API

Engine

Model

Proxy

bing

21

# OXFORD'S LUIS - BEYOND CORTANA

News about flight delays
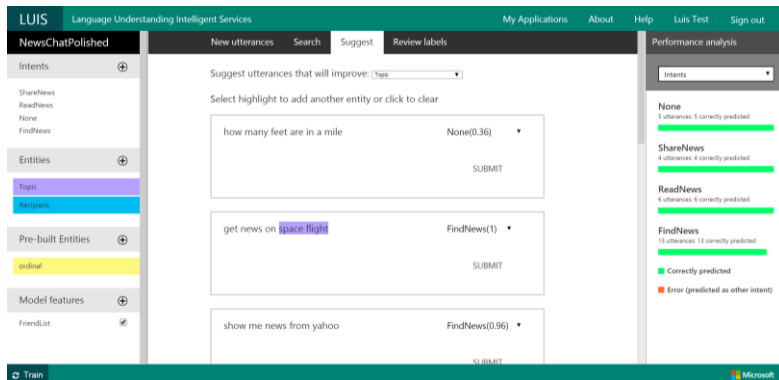
```
{
  "entities": [
    {
      "entity": "flight delays",
      "type": "Topic"
    }
  ],
  "intents": [
    {
      "intent": "FindNews",
      "score": 0.9985384
    },
    {
      "intent": "None",
      "score": 0.07289317
    },
    {
      "intent": "ReadNews",
      "score": 0.0167122427
    },
    {
      "intent": "ShareNews",
      "score": 1.0919299E-06
    }
  ]
}
```
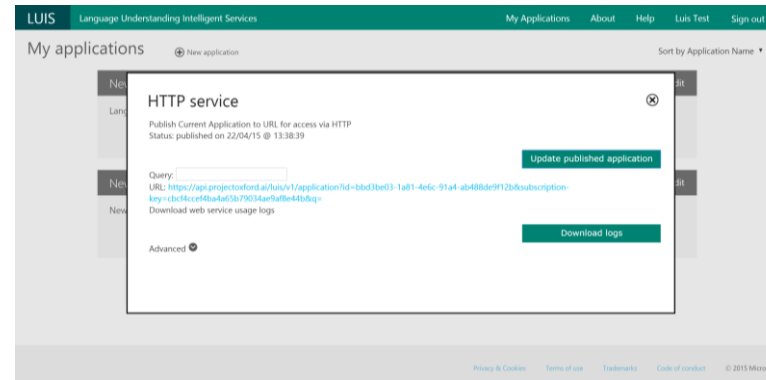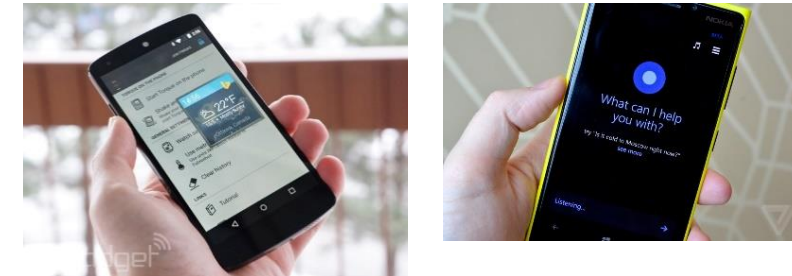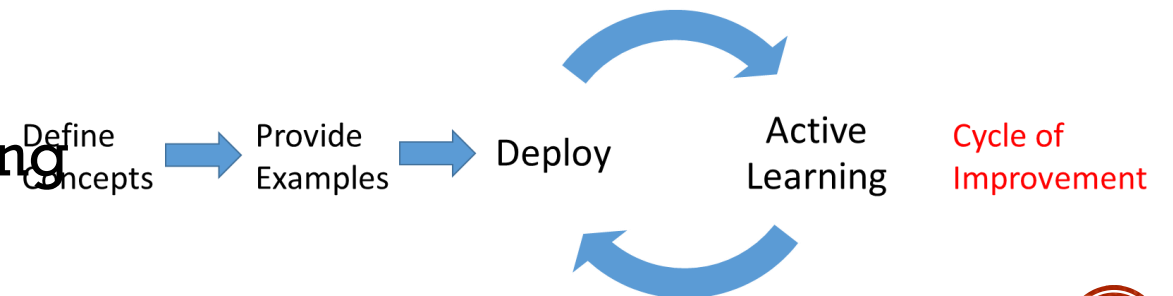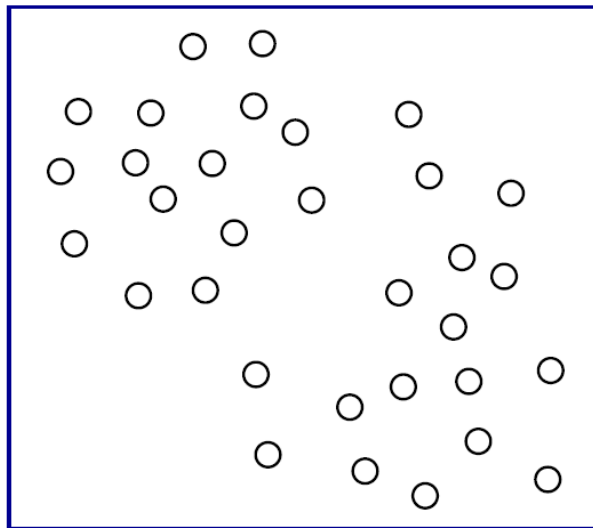
http://www.luis.ai

# LUIS OVERVIEW

Train

Deploy

Access



- Learns models from a few examples
- Continuous refinement with active learning
- Existing Bing models for common cases

Define Concepts → Provide Examples → Deploy → Active Learning
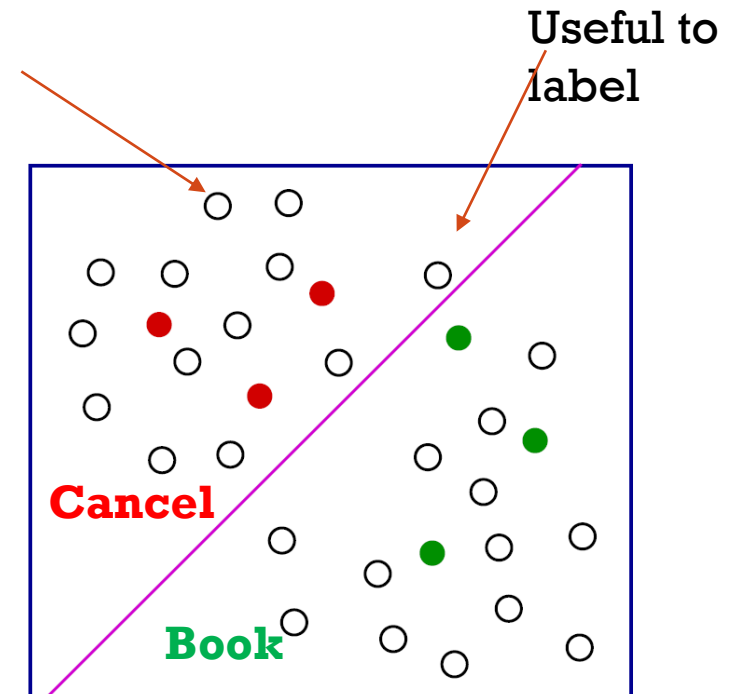
Cycle of Improvement

# ACTIVE LEARNING

- Model parameters are estimated with labeled data

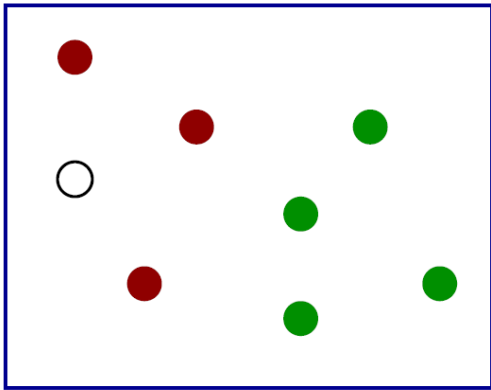- Labeling data is expensive

- Want to infer decision boundaries quickly

No sense labeling this!
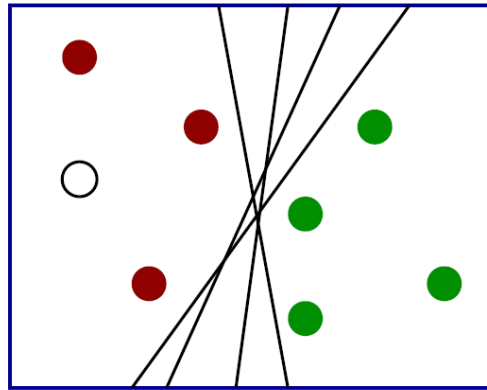
Useful to label
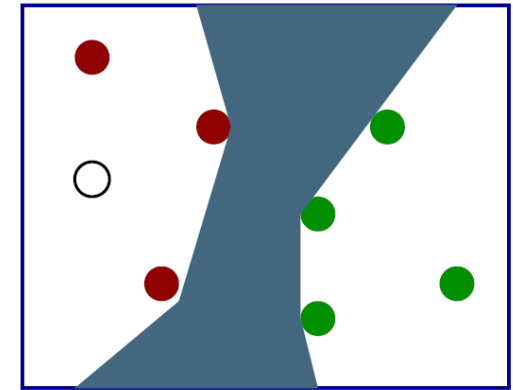
Cancel

Book

# ACTIVE LEARNING (2)

- Only label examples near the decision boundary



Is a label needed?

$H_t$ = current candidate hypotheses

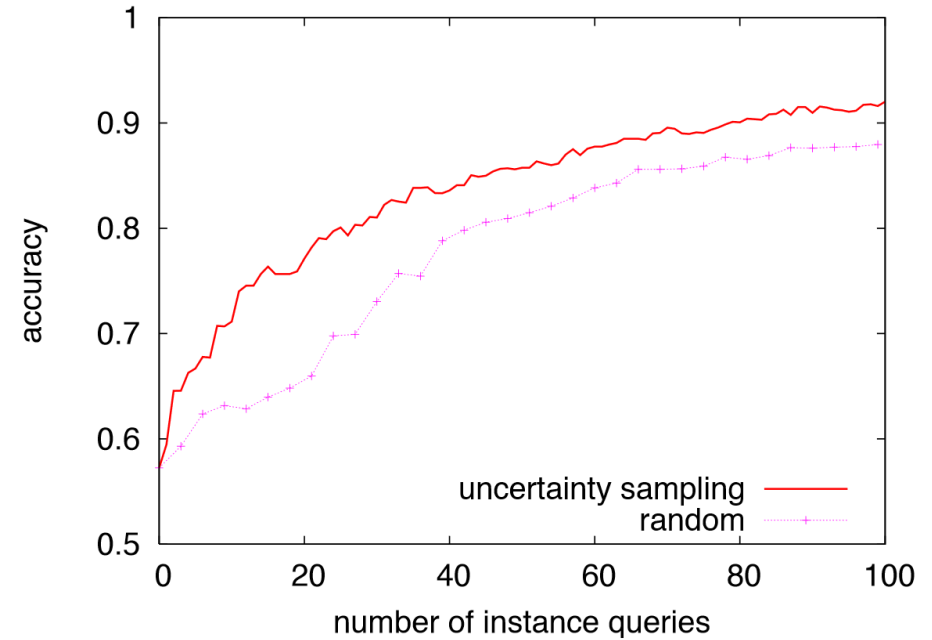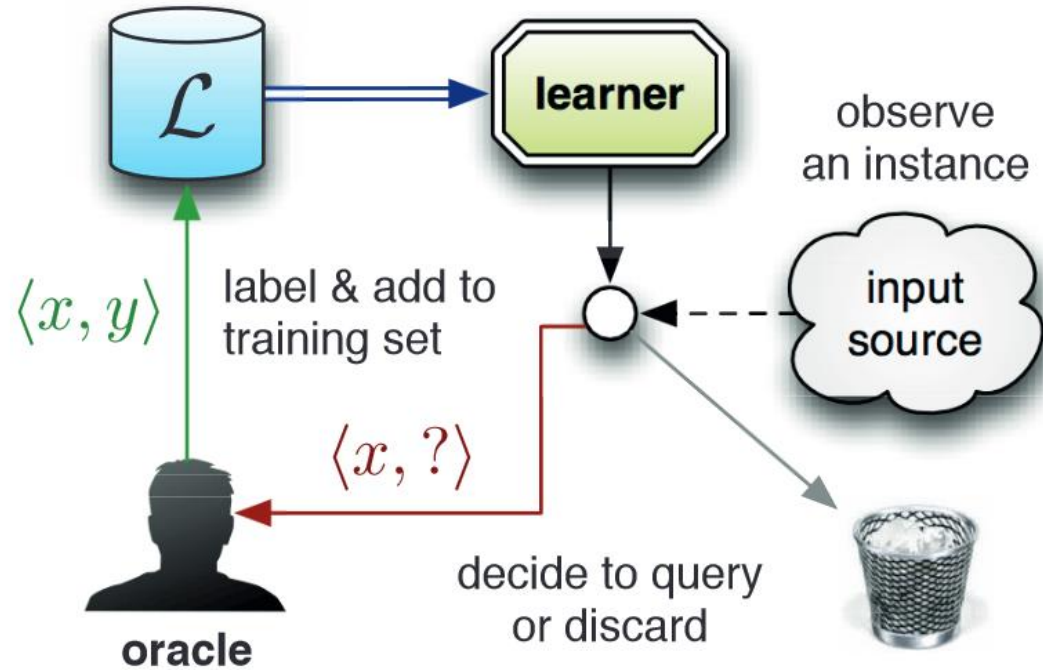Region of uncertainty

[Dasgupta & Langford 2009]

# ACTIVE LEARNING == FAST LEARNING

# SPOKEN LANGUAGE UNDERSTANDING DEMO



```
8
9    #import <UIKit/UIKit.h>
10   #import <SpeechSDK/SpeechRecognitionService.h>
11
12   @interface ViewController : UIViewController <SpeechRecognitionProtocol, UIWebViewDelegate>
13   @property (weak, nonatomic) IBOutlet UIWebView * mainWebUI;
14
15   @end
16
```

```
54
55   - (void)initSpeech // call when initialization
56   {
57       // set up speech session
58       MyLuisPreferences* prefs = [[MyLuisPreferences alloc] init];
59       prefs.LuisAppId = @"2e8ba91b-9491-                    ";
60       prefs.LuisSubscriptionId = @"7468c7cae122438699          ";
61       prefs.Locale = @"zh-cn"; // Chinese
62       // additional preferences settings for authentication or customized service endpoints
63
64       conversation = [ConversationBase alloc];
65       [conversation initWithPrefs:prefs withProtocol:(self)];
66       [conversation createConversation];
67   }
68
69   - (void)onPartialResponseReceived:(NSString*)value // be invoked while speech recognition is going on
70   {
71       NSMutableString* js = [[NSMutableString alloc] init];
72       [js appendFormat:(@"chat.updateVoiceInput('%@');"), value];
73       [self webRunJavascript:js];
74   }
75
76   -(void)onIntentReceived:(IntentResult*)intent
77   {
78       NSMutableString* js = [[NSMutableString alloc] init];
79       [js appendFormat:(@"chat.handleIntent('%@');"), intent.Body];
80       [self webRunJavascript:js];
81   }
82
```

# CREATING LUIS MODELS

# PROJECT PHILLY

Coming soon

# Computational Network Toolkit (CNTK)

## CNTK – a flexible and open source deep learning toolkit

- Networks: CNN, RNN, Bidirectional LSTM, DSSM, CRF…
- Problems: Speech, NLP, Ads, Search, large scale
- Learning algorithms: SGD, Adagrad, ADMM

## Network definition language

- Provides a simple yet powerful way to define a network

## CNTK simplifies deep learning experiments

- Design the model
- Derive the learning algorithm
- Implement the model
- Run the experiments
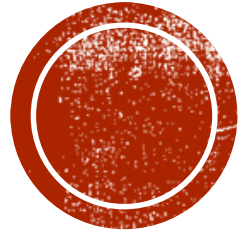
# AZURE GPU CLUSTER LAB

## Improved speech and image recognition workloads

- Deep Learning dramatically improved both speech and image recognition
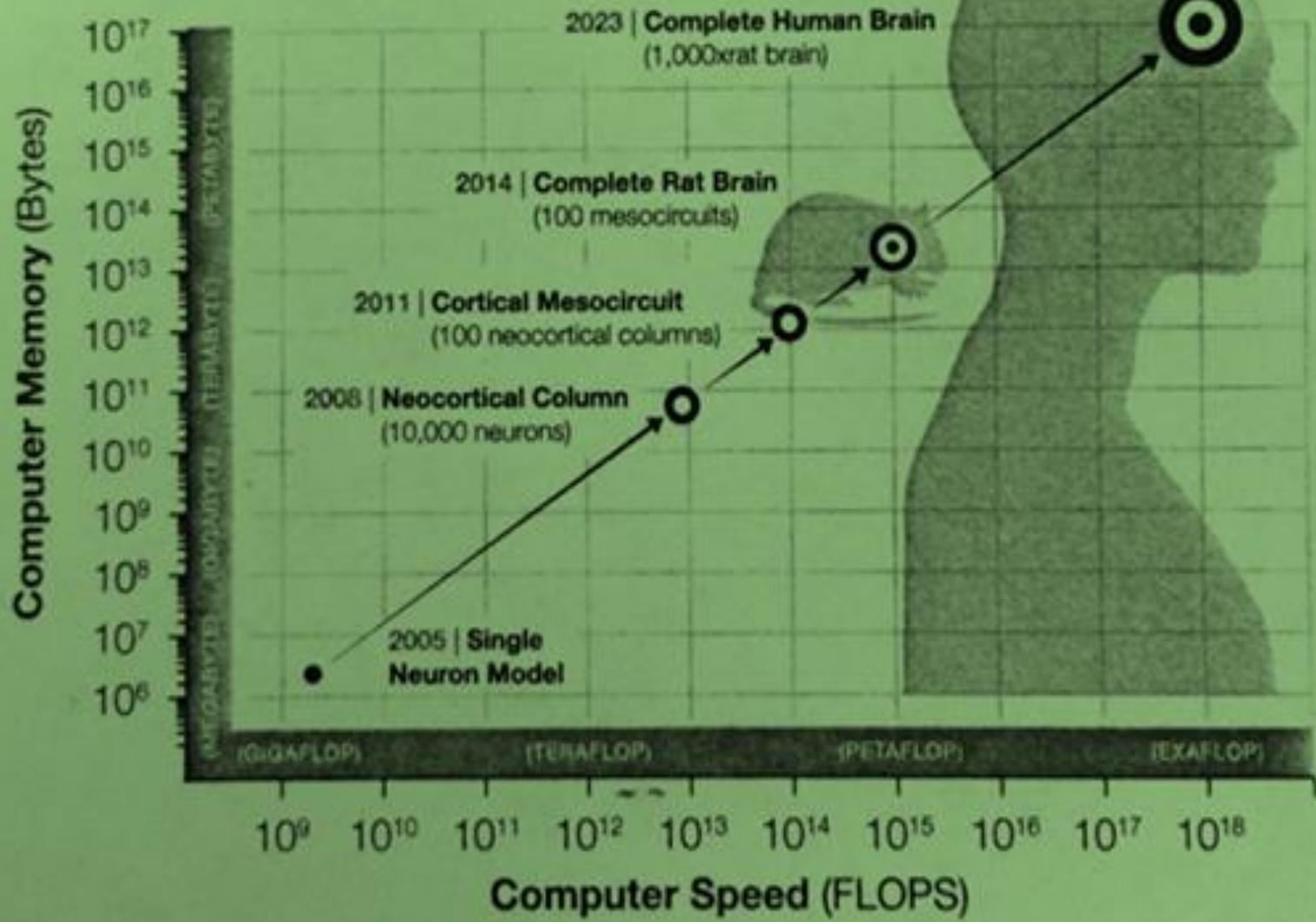- Progress hampered by distributed computing infrastructure

## AZURE GPU Labs

- Team up with researchers to scale out on deep learning
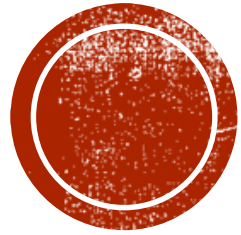- Optimized for CNTK

# WHAT IS NEXT?

# An **invisible** revolution is coming

# Q&A?