# The Deep Learning Revolution

Alex Acero

# Agenda

The Deep Learning Revolution

Fundamentals of Deep Learning

Why now? A brief history

Transforming our Digital Lives

# Acknowledgments

John Bridle and the Siri team

Josh Suskind, Sofien Bouaziz, and Apple's Video Team

# The Deep Learning Revolution

# Technology Disruptions

Content Creation:

    Text

    Photography

Content Consumption:

    Text

    Photography

    Music

    Video
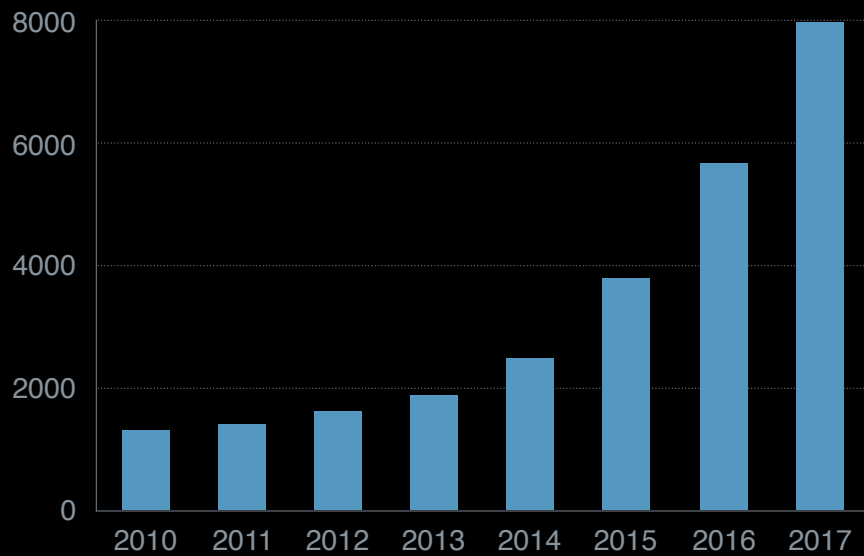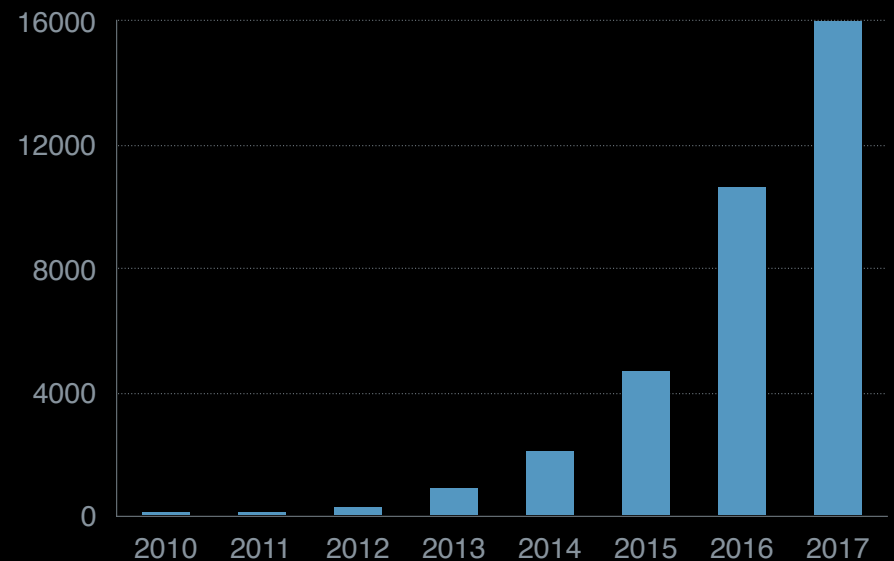
Our daily lives:

    Transportation

    Communication

    Shopping

    Travel

# The Deep Learning Revolution



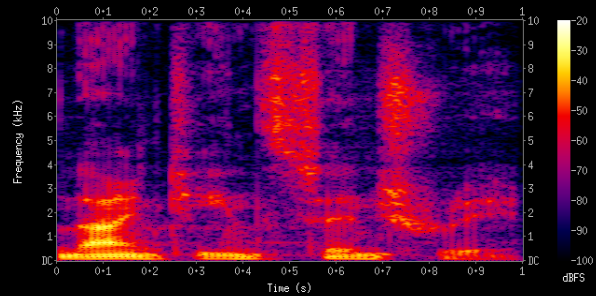Neural Information Processing Systems (NIPS) Attendees

Papers with "Deep Neural Networks"
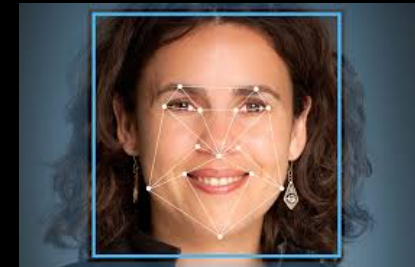
# Fundamentals of Deep Learning
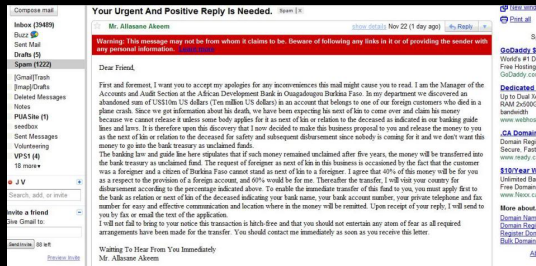
# Binary Classification



TouchID



Speaker Verification



Face ID



Email Spam



Motion Detection



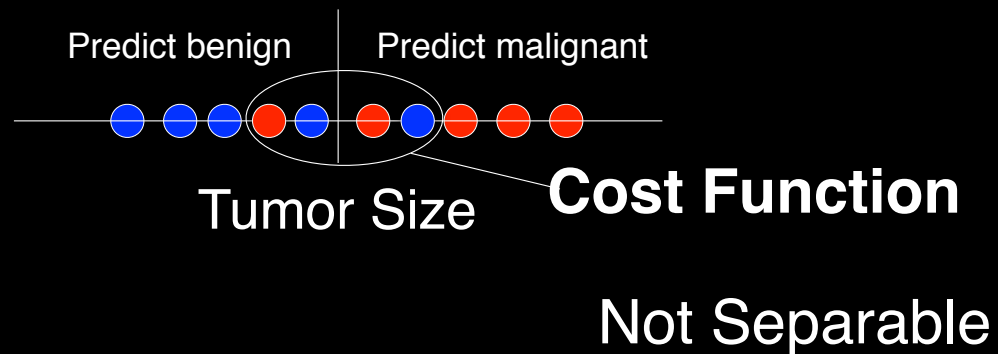Credit Card Fraud
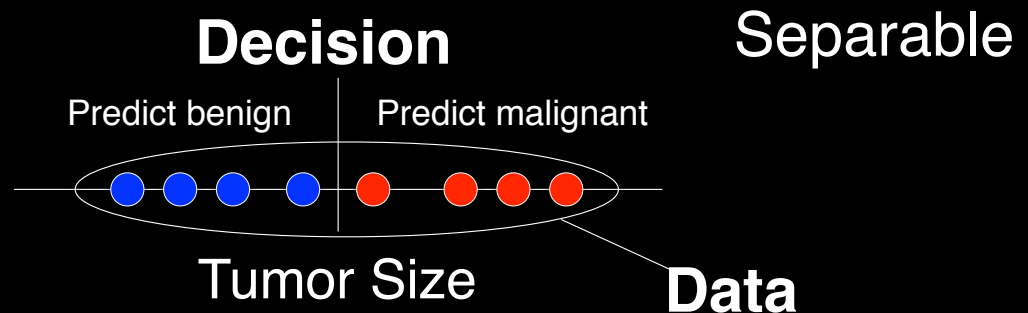
# Binary Classification

## Output Labels

Breast Cancer
- 🔵 Benign
- 🔴 Malignant

## Input Features
Tumor Size

**Decision**

Predict benign | Predict malignant

Tumor Size

**Data**

Separable

Predict benign | Predict malignant

Tumor Size

**Cost Function**

Not Separable

# Binary Classification

**Output Labels**

Breast Cancer
- 🔵 Benign
- 🔴 Malignant

**Input Features**

Tumor Size
Age

**More features:**

Clump thickness
Uniformity of cell size
…

Age

Predict malignant
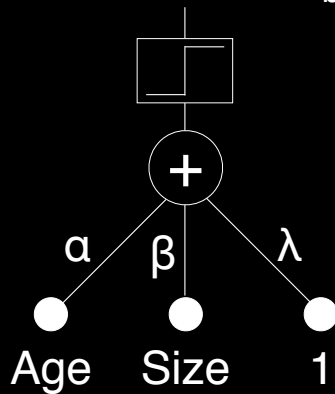
Predict benign

Tumor Size

Linearly Separable

$\alpha \cdot \text{Age} + \beta \cdot \text{Size} + \lambda > 0$

**Learning**

# Perceptron Learning

Rosenblatt, 1958

$$\alpha \cdot Age + \beta \cdot Size + \lambda \underset{benign}{\overset{malignant}{\gtrless}} 0$$



$\alpha(i) = \alpha(i - 1) + \eta \cdot \{Target(i) - Output(i)\} \cdot Age(i)$

$\beta(i) = \beta(i - 1) + \eta \cdot \{Target(i) - Output(i)\} \cdot Size(i)$

$\lambda(i) = \lambda(i - 1) + \eta \cdot \{Target(i) - Output(i)\}$

# Stochastic Gradient Descent (SGD)



$$y = \sigma(x)$$

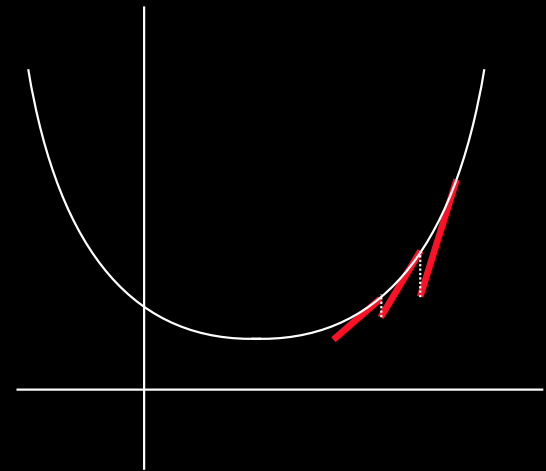$$x = c + \mathbf{v}^T\mathbf{w}$$

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$p(t|\mathbf{v}) = y^t(1-y)^{1-t}$$
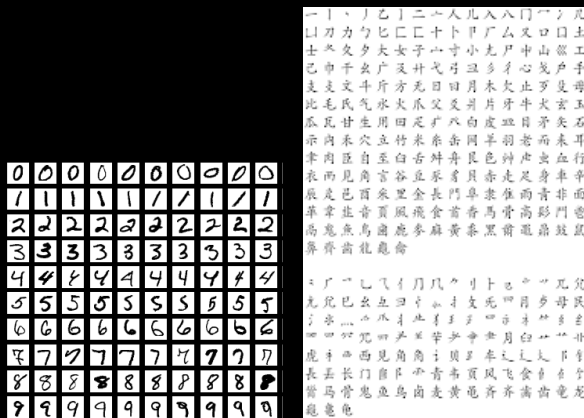
$$L = \ln p(t|\mathbf{v}) = t\ln y + (1-t)\ln(1-y)$$

$$\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial y}\right)\left(\frac{\partial y}{\partial x}\right)\left(\frac{\partial x}{\partial w_1}\right) = (y-t)v_1$$

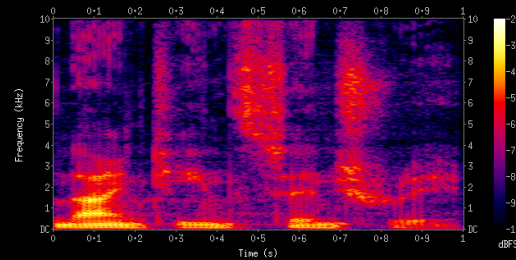$$w_j^{(i)} = w_j^{(i-1)} - \eta\frac{\partial L}{\partial w_j} = w_j^{(i-1)} + \eta v_j\big(t - y^{(i-1)}\big)$$
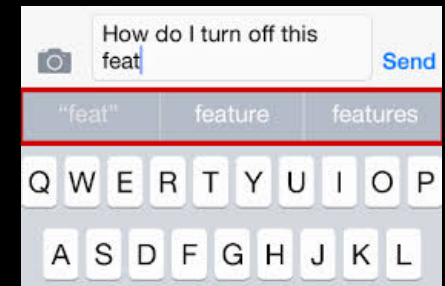
# N-ary Classification
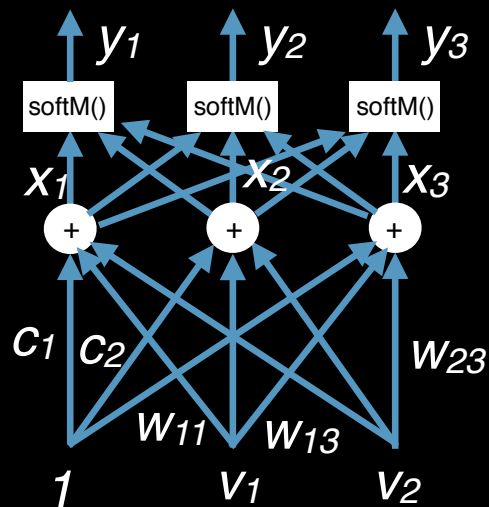


Handwriting Recognition



Speaker Identification



Word prediction

# N-ary Classification



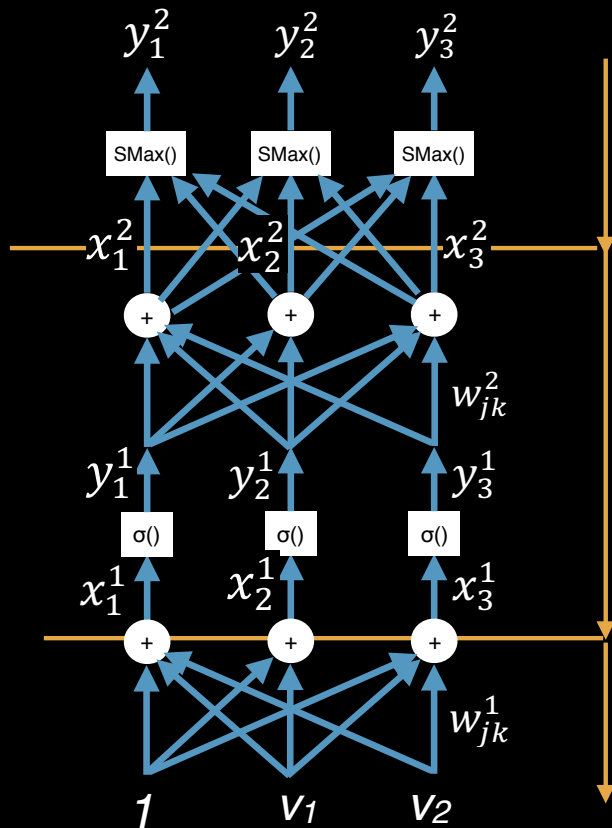$$y_i = p(i|\mathbf{v}) = \frac{e^{x_i}}{\sum_{l=1}^{N} e^{x_l}}$$

Softmax

$$L = \sum_{i=1}^{N} t_i \ln y_i$$

$$w_{nj}^{(i)} = w_{nj}^{(i-1)} + \eta v_n^{(i-1)}\left(t_j - y_j^{(i-1)}\right)$$

# Perceptron Learning

Werbos, 1974; Rumelhart, Hinton, Williams 1986



Two-layers

2 input features

3 output labels

$$\nabla_n^2(m) = y_n^2(m) - t_n(m)$$

$$\left[w_{jn}^2\right]^{(i)} = \left[w_{jn}^2\right]^{(i-1)} - \eta \frac{1}{M} \sum_{m=1}^{M} y_n^1(m)\nabla_n^2(m)$$

$$\nabla_n^1(m) = y_n^1(m)\left(1 - y_n^1(m)\right) \sum_{k=1}^{N} w_{nk}^2 \nabla_k^2(m)$$

$$\left[w_{jn}^1\right]^{(i)} = \left[w_{jn}^1\right]^{(i-1)} - \eta \frac{1}{M} \sum_{m=1}^{M} v_j(m)\nabla_n^1(m)$$

backpropagation

Mini-batch

# CNN on Face Images
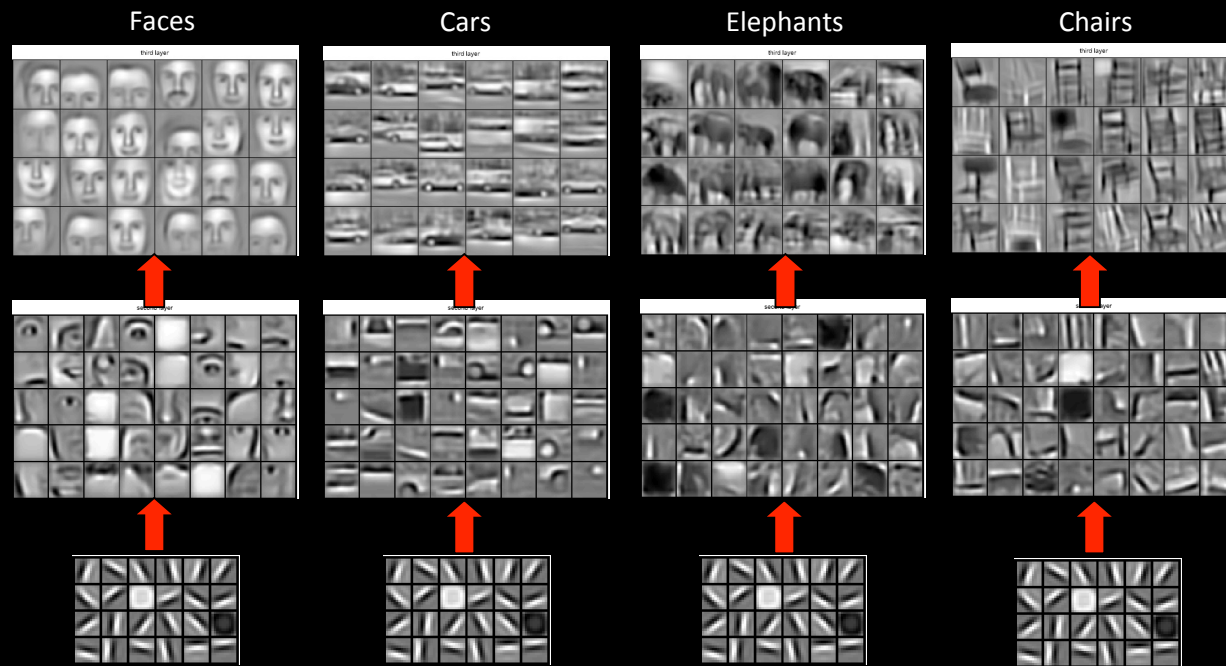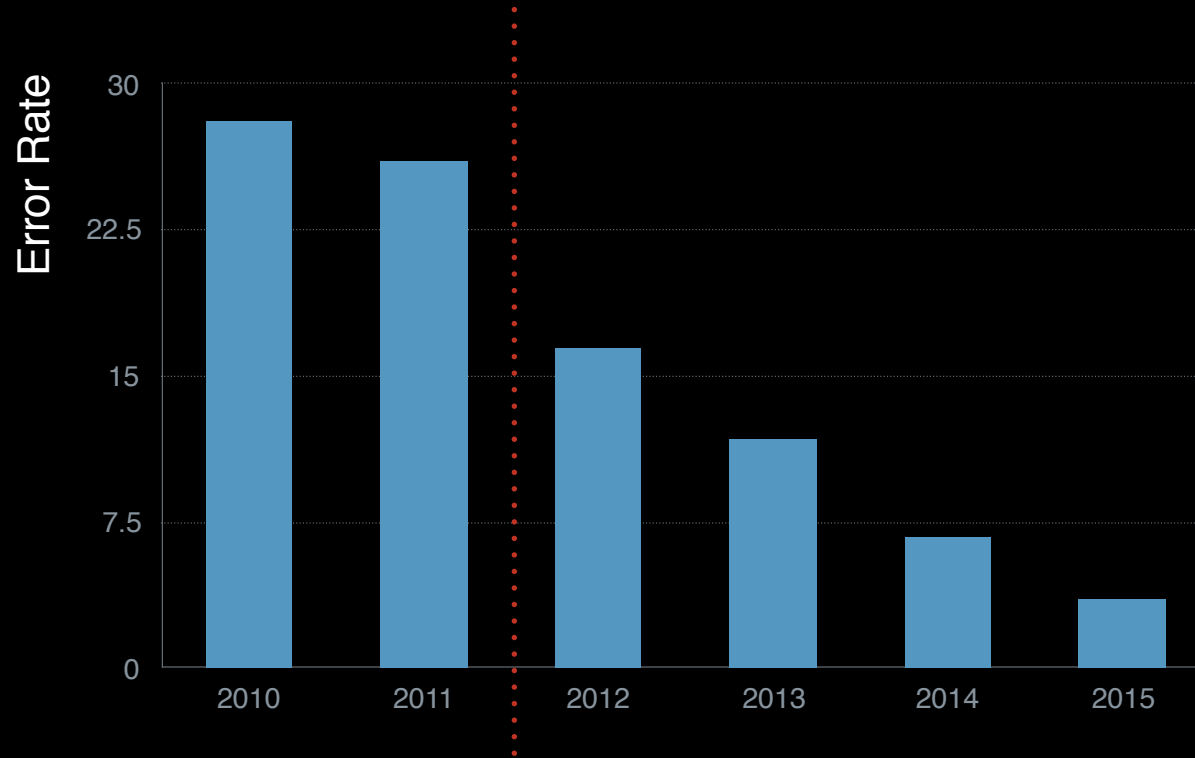## 2012



Deep Belief Net on Face Images

Based on materials by Andrew Ng

16

# ImageNet Large Scale
## Visual Recognition Challenge, 2012

Examples of learned object parts from object categories

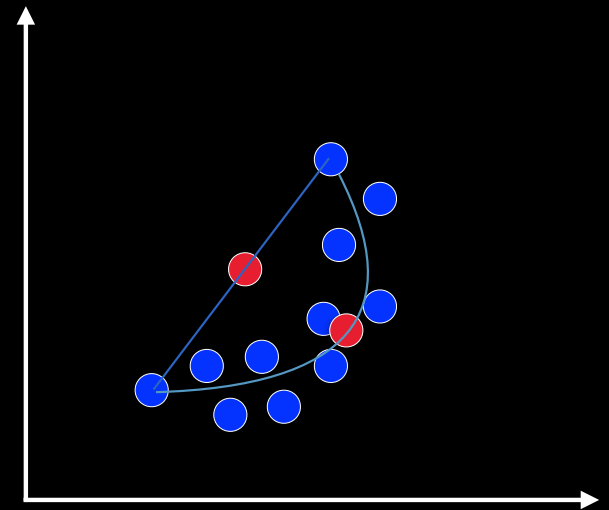| Faces | Cars | Elephants | Chairs |
|-------|------|-----------|--------|

# ImageNet Task Progress

# Non-Linear Manifolds

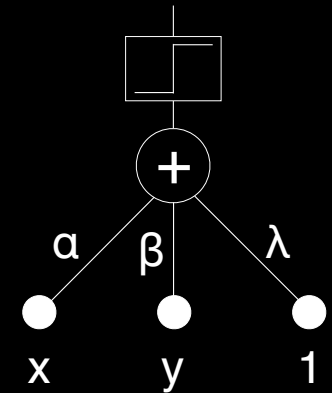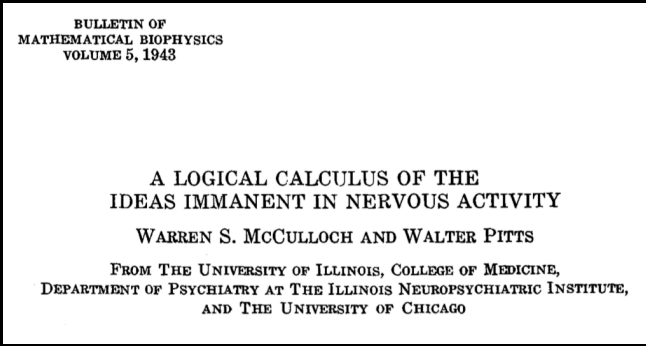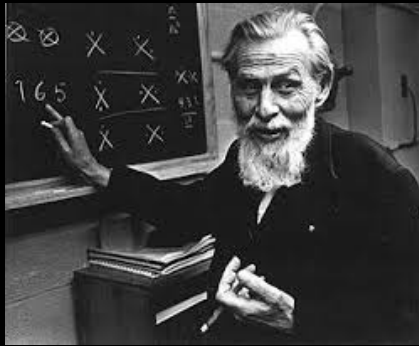Such non-linearity requires multiple layers

# Why Now?

A brief history

# McCulloch-Pitts Neurons

1943

# Norbert Wiener

Wiener–Khinchin Theorem (1930)
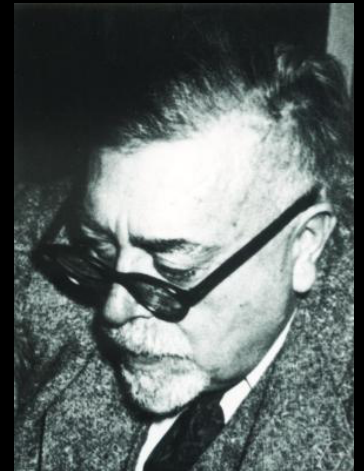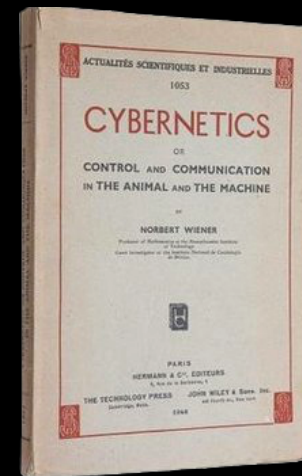Wiener Filter (1949)

McCulloch & Pitts joined Wiener at MIT (1943)

Cybernetics (1948)
   5. Computing Machines and the Nervous System
   10. Brain Waves and Self-Organising Systems

Suggested chess playing programs
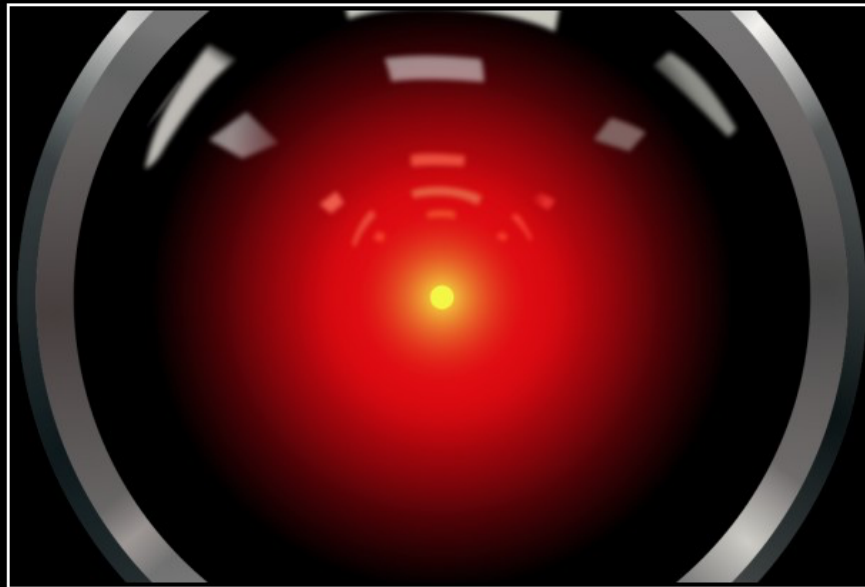
# Checkers
Arthur Samuel, 1956

# ELIZA
## Weizenbaum,1966

```
========================================================================
          EEEEEEEE  L          IIIIIII  ZZZZZZZ      AAA
          E         L             I          Z      A   A
          E         L             I         Z      A     A
          EEEEE     L             I        Z       A     A
          E         L             I       Z        AAAAAAA
          E         L             I      Z         A     A
          EEEEEEEE  LLLLLLLL   IIIIIII  ZZZZZZ     A     A
========================================================================
ELIZA > (PLEASE TYPE IN ALL CAPS) WHAT'S YOUR NAME DEAR ?
AMIT
ELIZA > HI, I'M ELIZA. WHAT DO YOU WANT TO TALK ABOUT ?
AMIT > CAN I TALK ABOUT MY PROBLEM ?
ELIZA > SURE... !
AMIT > DUE TO ON GOING EXAMS IN COLLEGE, I HAVE A LOT OF STRESS.
ELIZA > PLEASE RELAX A LITTLE. SLEEP WELL.
AMIT > THANKS FOR YOUR ADVICE.
ELIZA > NO MENTION.
AMIT > BYE
ELIZA > BYE AND KEEP IN TOUCH...

========================================================================
```

# 2001 Space Odyssey
Stanley Kubrick, 1968



I'M SORRY, DAVE.
I'm afraid I can't do that.

# AI Winter

# Fundamental Equation of Speech Recognition



$$\widehat{W} = \underset{W}{\mathrm{argmax}}\, p(W|A) = \underset{W}{\mathrm{argmax}}\, p(A|W)p(W) = \underset{W}{\mathrm{argmax}}\{\ln p(A|W) + \ln p(W)\}$$
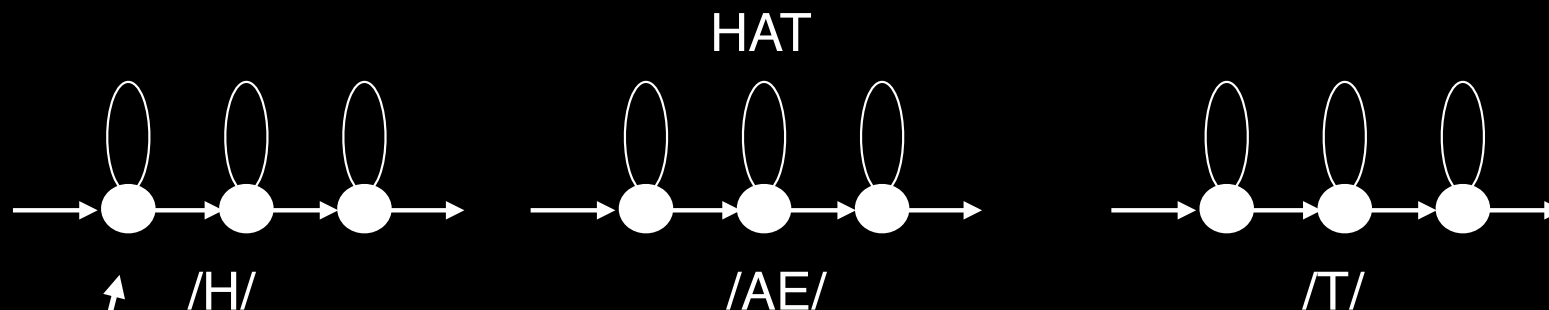
$$\widehat{W} = \underset{W}{\mathrm{argmax}}\{\lambda \ln p(A|W) + \ln p(W)\}$$

Acoustic Model          Language Model

# Acoustic Model
Hidden Markov Models

HAT



/H/          /AE/          /T/

$$p(a_t|s_j) = \sum_{i=1}^{I} \alpha_i \mathcal{N}\left(a_t, \mu_{ij}, \Sigma_{ij}\right)$$

24                IEEE TRANSACTIONS ON ACOUSTICS, SPEECH, AND SIGNAL PROCESSING, VOL. ASSP-23, NO. 1, FEBRUARY 1975

## The DRAGON System—An Overview

JAMES K. BAKER

PROCEEDINGS OF THE IEEE, VOL. 64, NO. 4, APRIL 1976

## Continuous Speech Recognition by Statistical Methods

FREDERICK JELINEK, FELLOW, IEEE

# Neural Networks for Speech Recognition in the 1990's

# Neural Network Winter for Speech Recognition

# Open Challenge Tasks
DARPA



Courtesy NIST 1999 DARPA HUB-4 Report, Pallett et al. & new updates from DARPA

# Deep Learning

# Deep Belief Networks ➡ Deep Neural Networks



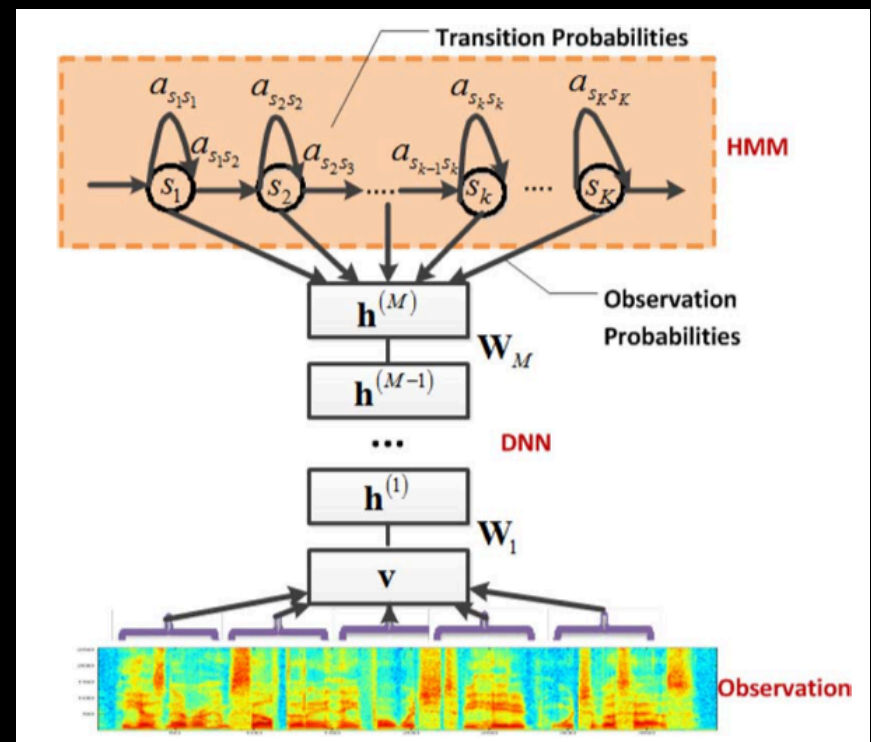Papers with "Deep Neural Networks"

# Deep Learning for Speech
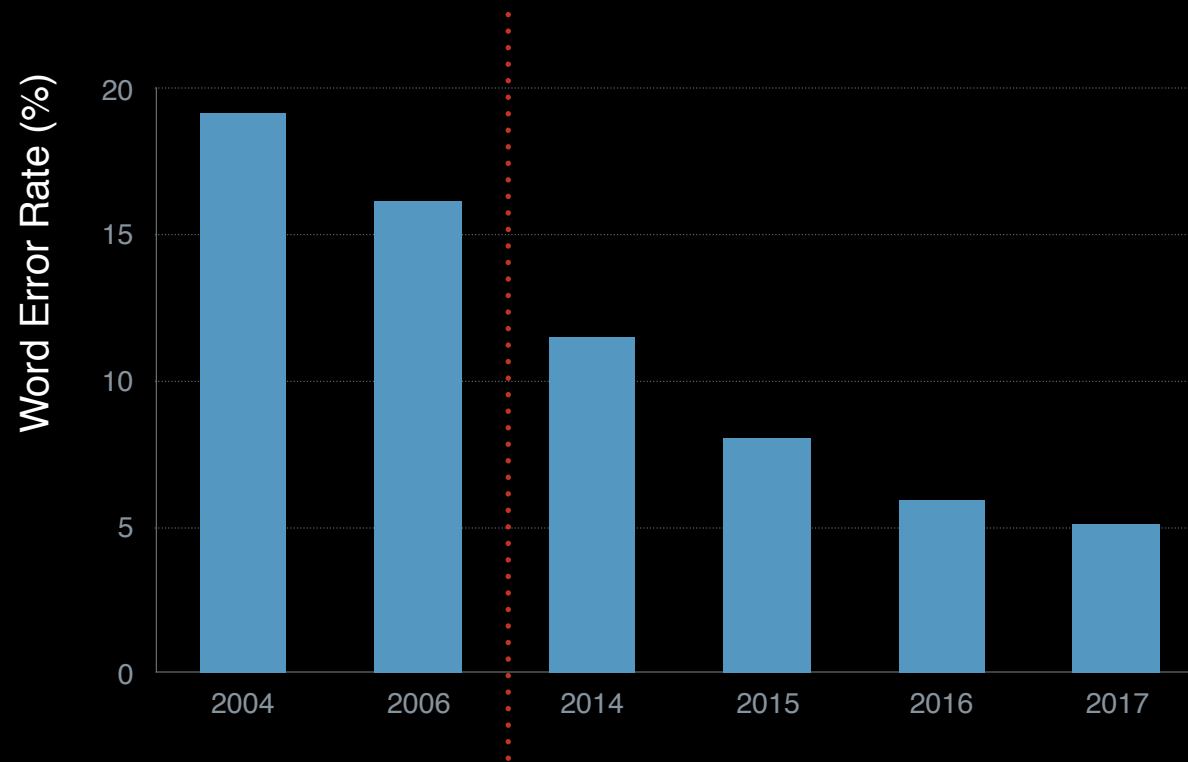Deng et al., 2010

DNNs for large vocabulary

- 800 input features
- 5 layer network
- 1000 neurons per layer
- 8000 output labels
- 12 Million weights

Training

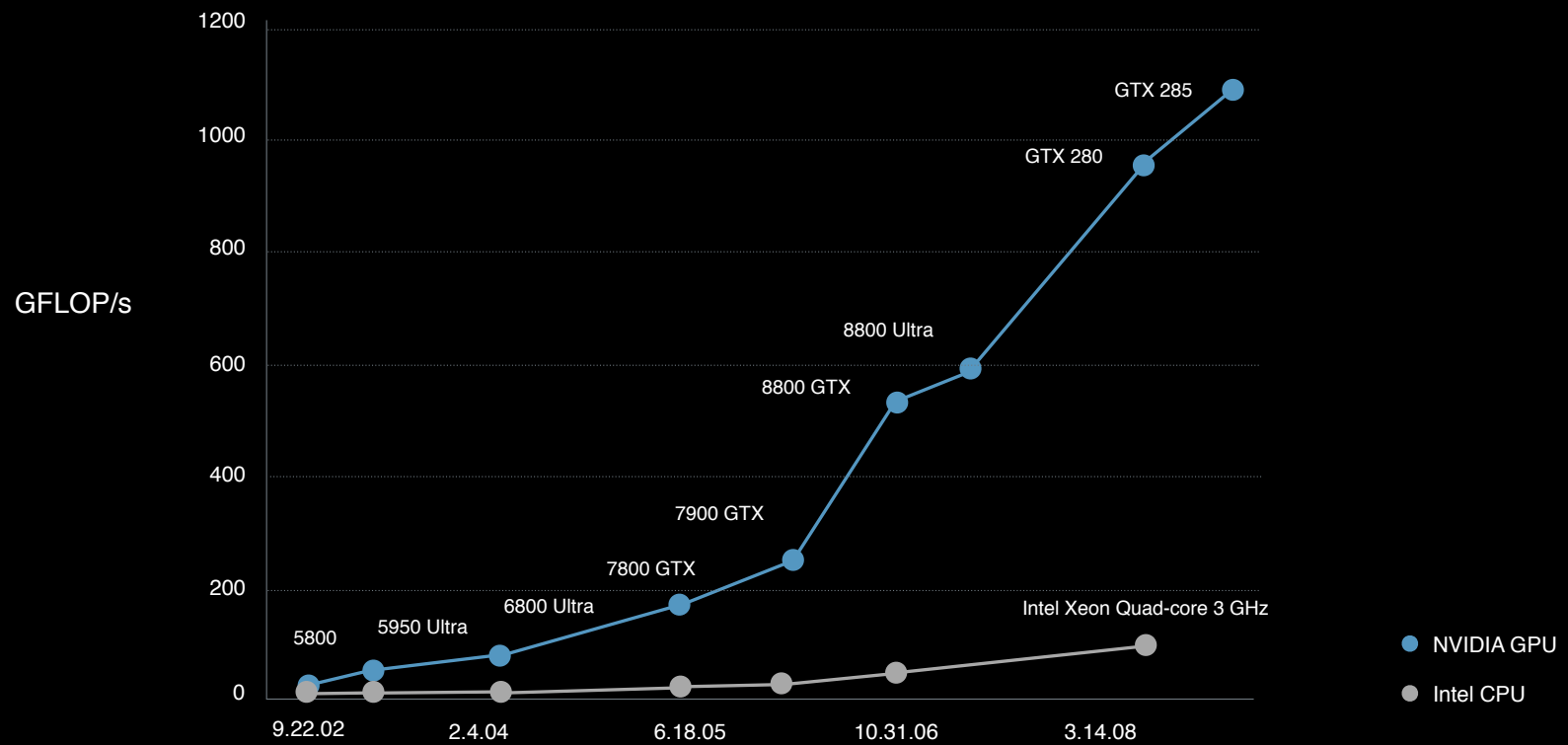- 300 hours of speech with transcripts
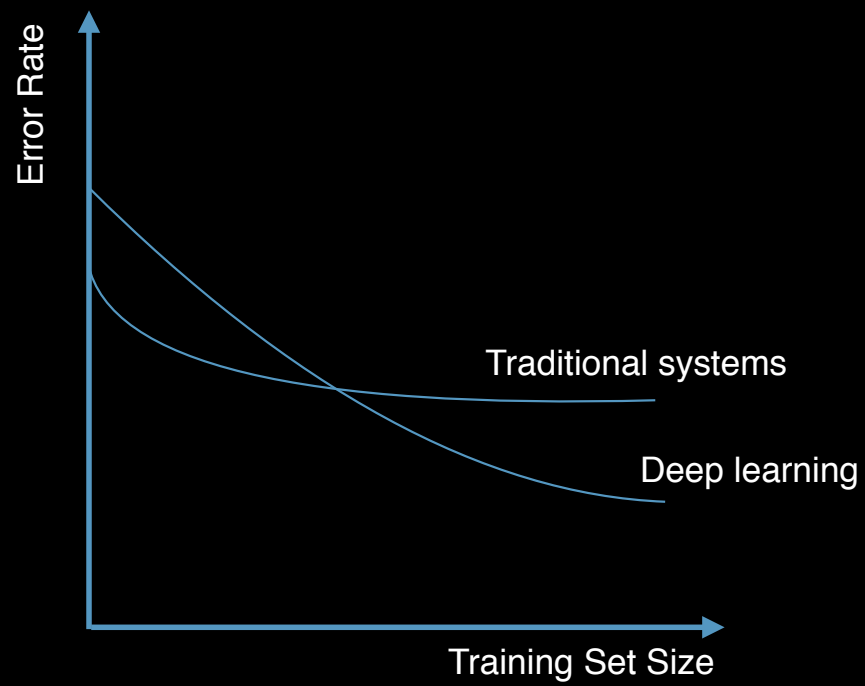- 1 week training time on a GPU

# Switchboard

# Why Now?
## GPUs

### Raw Performance Trends



GFLOP/s

- GTX 285
- GTX 280
- 8800 Ultra
- 8800 GTX
- 7900 GTX
- 7800 GTX
- 6800 Ultra
- Intel Xeon Quad-core 3 GHz
- 5950 Ultra
- 5800

- ● NVIDIA GPU
- ● Intel CPU

X-axis: 9.22.02, 2.4.04, 6.18.05, 10.31.06, 3.14.08

Y-axis: 0, 200, 400, 600, 800, 1000, 1200

# Why Now?
## Large Amounts of Data

# Why Now?
Algorithms

- Direct modeling of context-dependent (tied triphone states) through the DNN

- ~~Unsupervised Pre-training~~

- Deeper networks

# Why Now?
Open sharing

U. Toronto
Microsoft
Google
IBM



Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N. Sainath, and Brian Kingsbury

## Deep Neural Networks for Acoustic Modeling in Speech Recognition

[The shared views of four research groups]

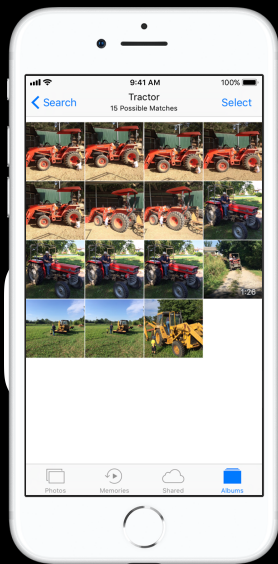# Why Now?
Tools

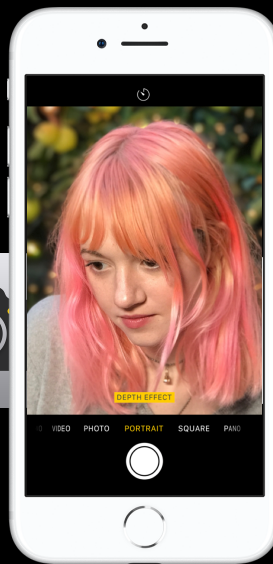# Deep Learning Has Roots in Signal Processing

ICASSP Attendees

NIPS Attendees

# Transforming Our Digital Lives
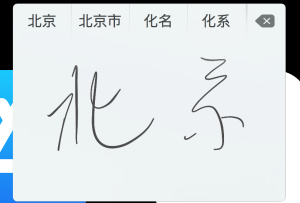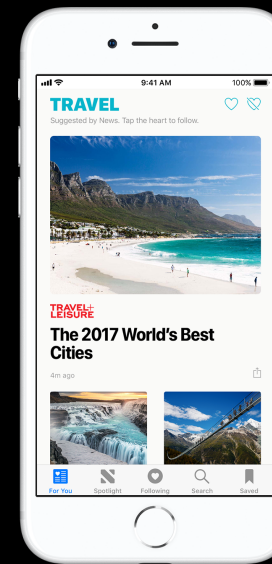
# ML Becomes Mainstream



On-device scene recognition

Portrait Mode

Language modeling

Handwriting recognition

News recommendation

Intelligent assistant

# Siri
Apple, 2011
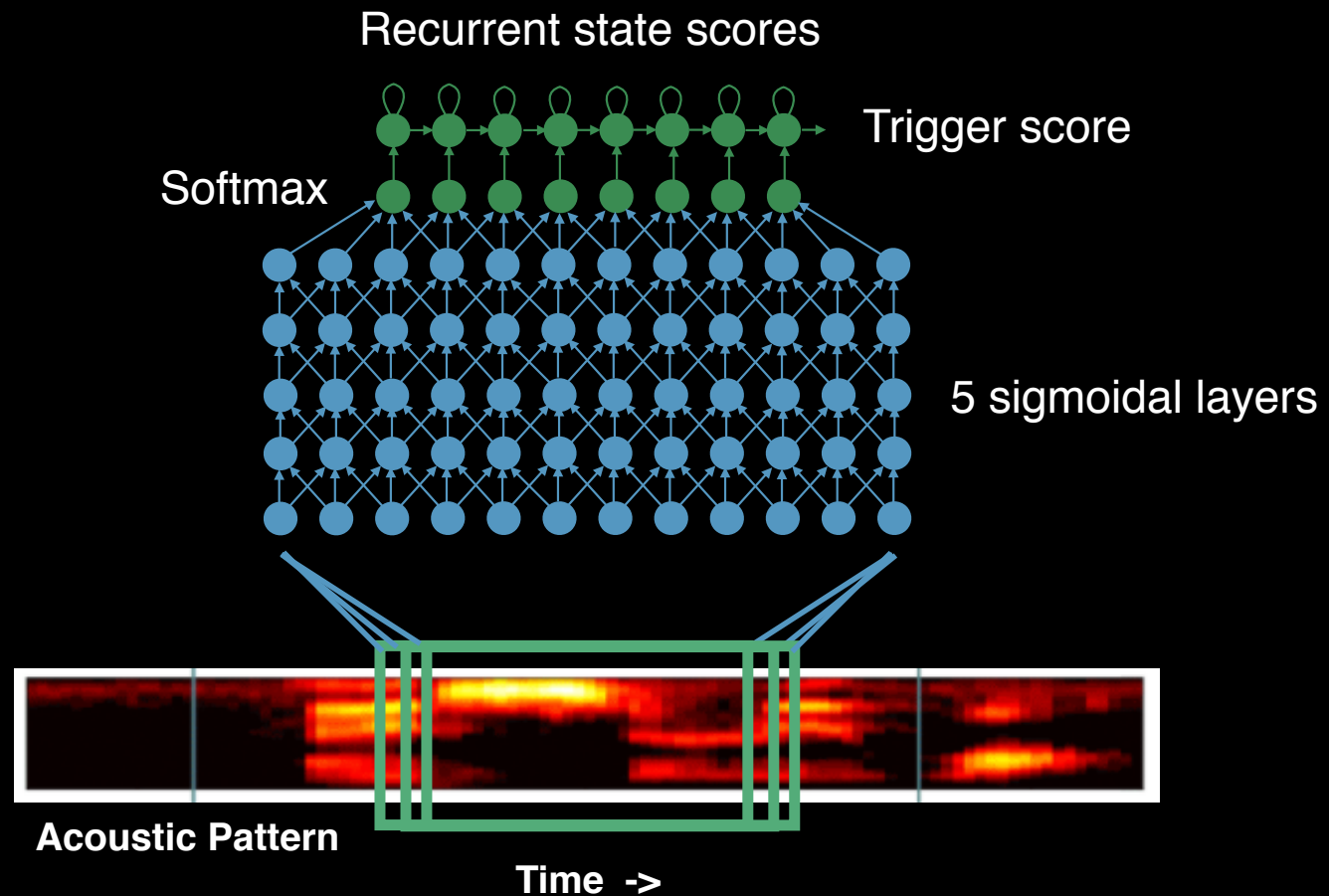
# Hands-Free Siri
## Design of the Voice Trigger

# Hey Siri DNN

Recurrent state scores

Trigger score

Softmax

5 sigmoidal layers

Acoustic Pattern

Time ->

# Multi-Pass Detection

**Always on processor (low compute)**

MFCC Computation

frame buffer

small DNN

HMM scorer

**Main processor (more accurate)**

large DNN

HMM scorer

Personalized Model

**Server processor (full speech recognition system)**

X large DNN

# Two-Pass Detection



FRR: Percentage of attempts rejected

25%
20%
15%
10%
5%
0%

Small: 5x32

Large: 5x192

Normal

Lower

Primary

1

10

100

FAR: False Alarms per 100 hours

# Computing for Deep Learning



DNN Runtime

**DNN for HeySiri runs in WATCH**

**Training done offline in a large GPU grid**

Speech transcripts

DNN Training

Models

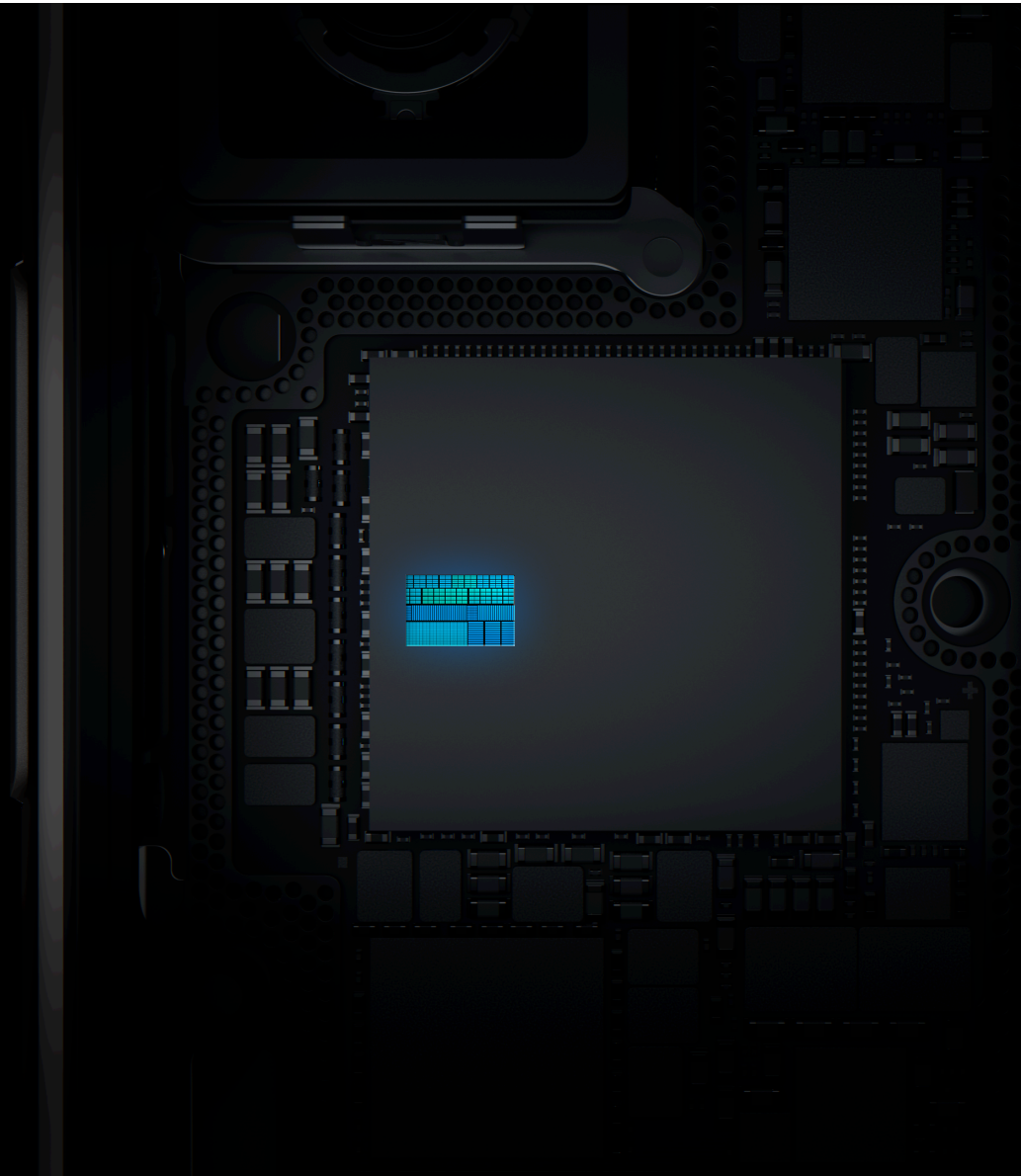**Training data**

Face ID

Apple, 2017

Neural engine

Dual-core design

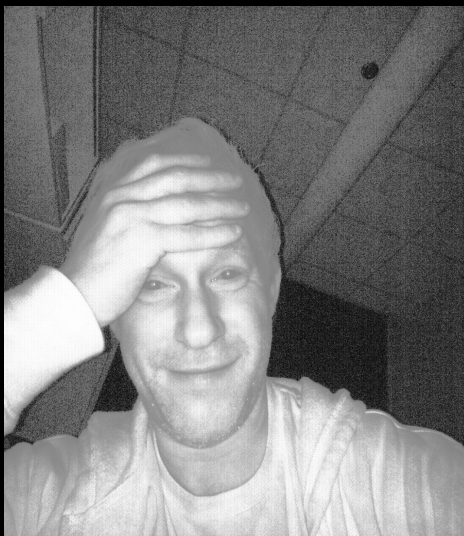600 billion operations per second

Real-time processing

# Unconstrained Face Matching

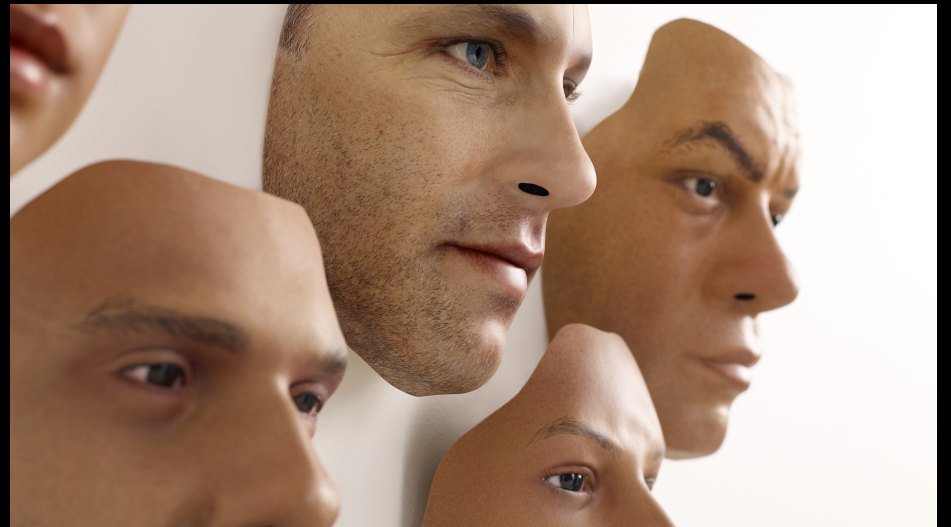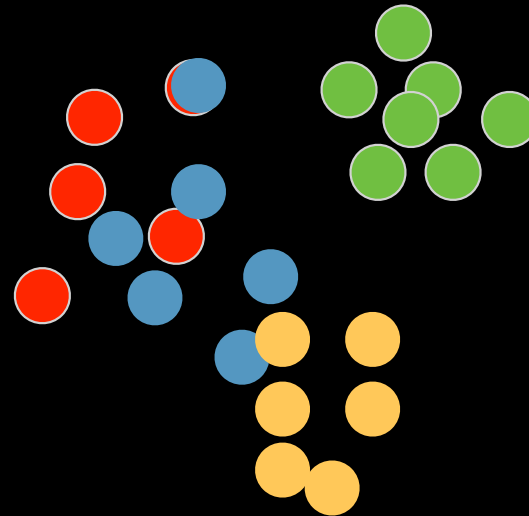# Works in Bright Sunlight and Shadows

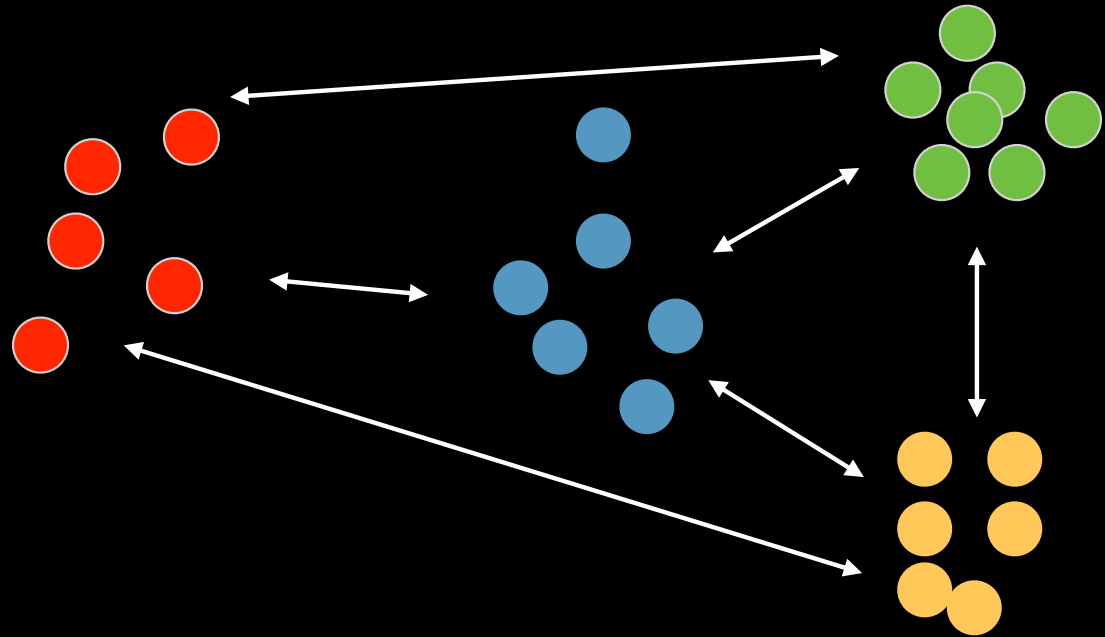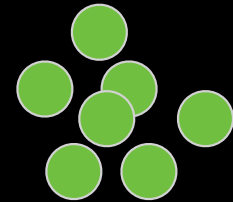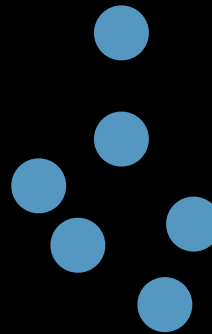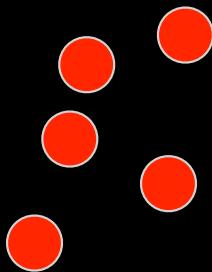# Robust to Occlusions

# Anti-Spoofing

# Face ID is a Machine Learning Problem

Goal is to pull same identity
pairs together and push
different identity pairs apart

# Face ID is a Machine Learning Problem

Goal is to pull same identity
pairs together and push
different identity pairs apart
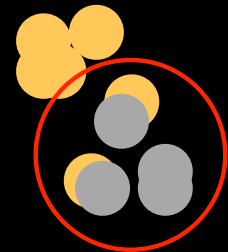
# Face ID is a Machine Learning Problem

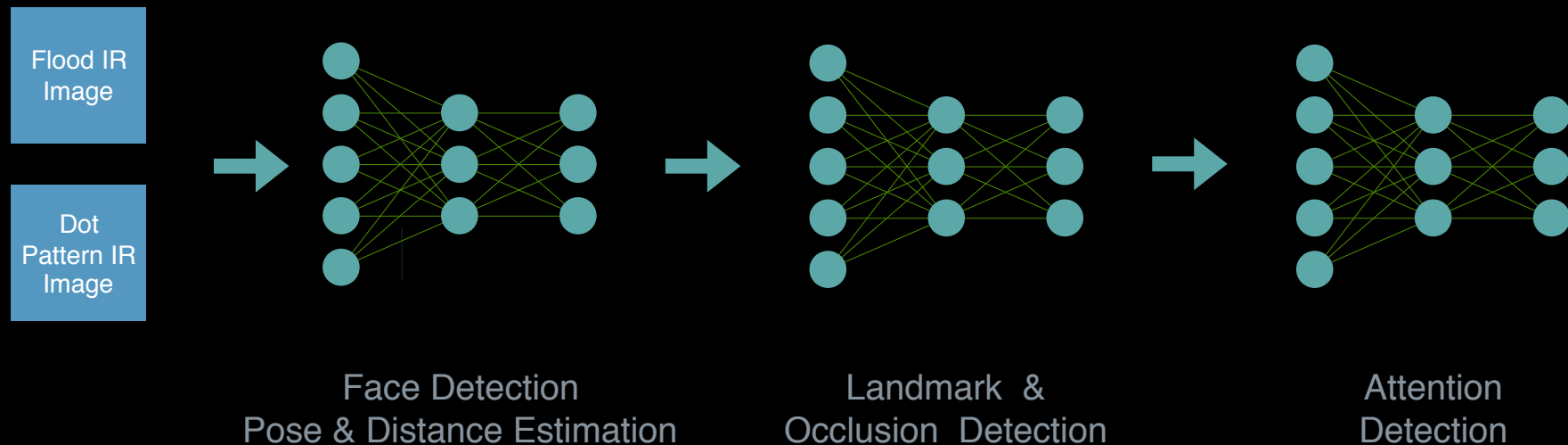Most faces are not similar—
needles in a haystack

# Face ID is a Machine Learning Problem

Sometimes it is very hard to find patterns that separate people that are not spurious

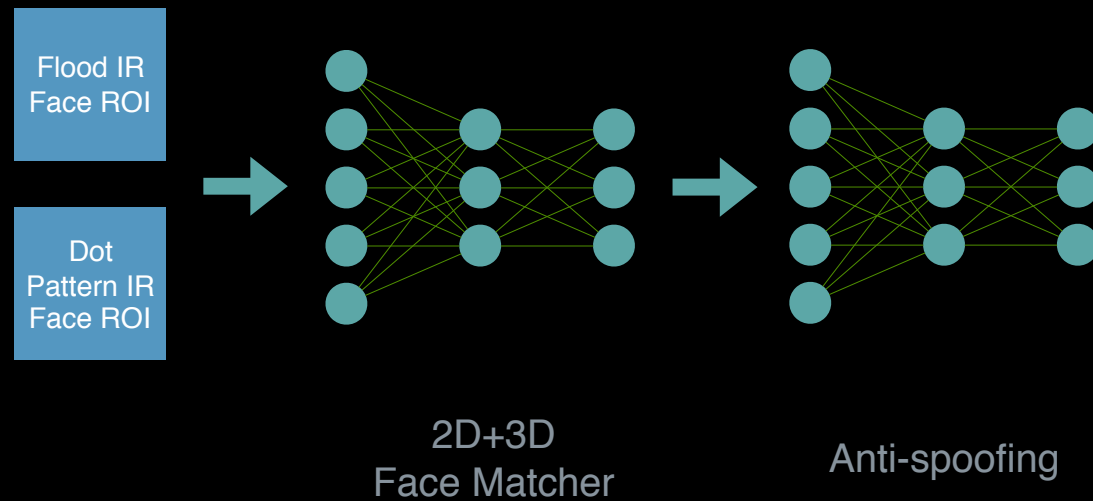A model that is better at all the easy cases is not necessarily better at solving the hard cases
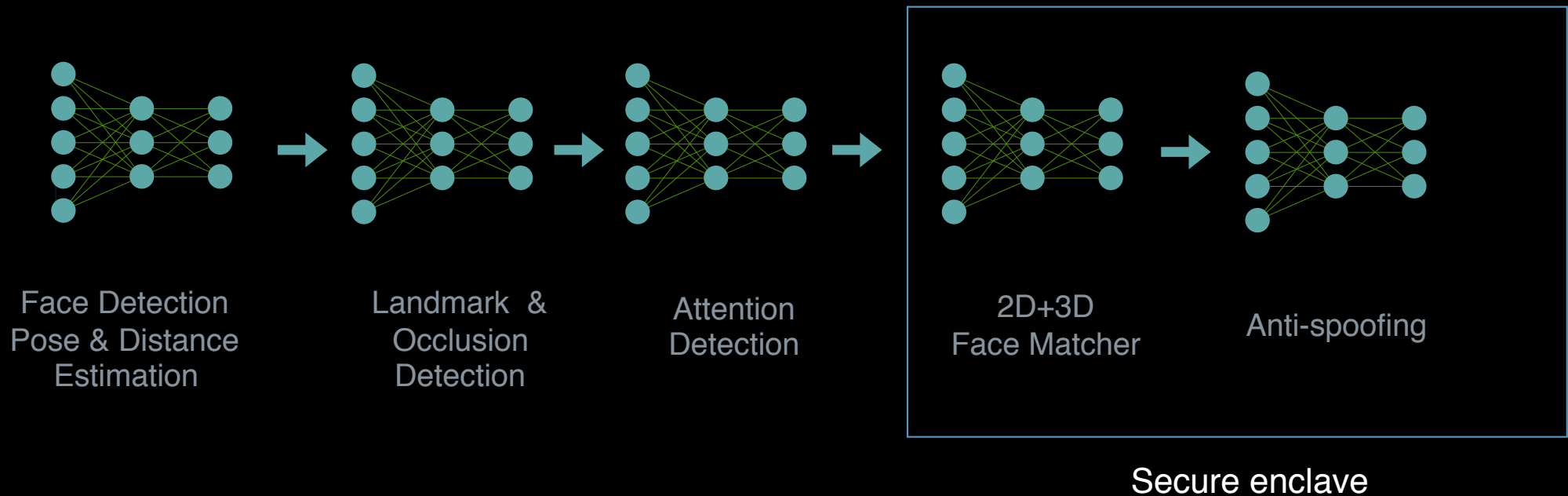
# Neural Network Face Matching Pipeline: Detection



Flood IR Image

Dot Pattern IR Image

Face Detection
Pose & Distance Estimation

Landmark &
Occlusion Detection

Attention
Detection

Makes decision at any point (no face, out of spec, inattention)

Localizes faces for matching

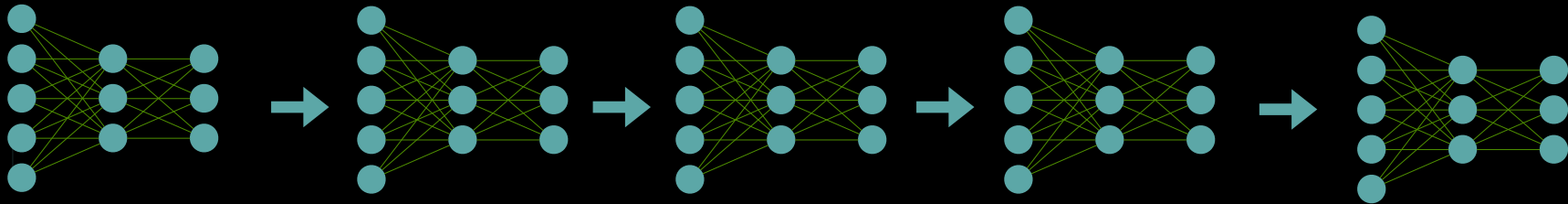# Neural Network Face Matching Pipeline: Verification



Flood IR Face ROI

Dot Pattern IR Face ROI

2D+3D Face Matcher

Anti-spoofing

Multimodal learning problem (how to fuse 2D and 3D representations)

# Neural Network Face Matching Pipeline: End-To-End



Face Detection
Pose & Distance
Estimation

Landmark &
Occlusion
Detection

Attention
Detection

2D+3D
Face Matcher

Anti-spoofing

Secure enclave

# Neural Network Face Matching Pipeline: End-To-End



Has to be really fast

Small memory footprint

Limited power impact

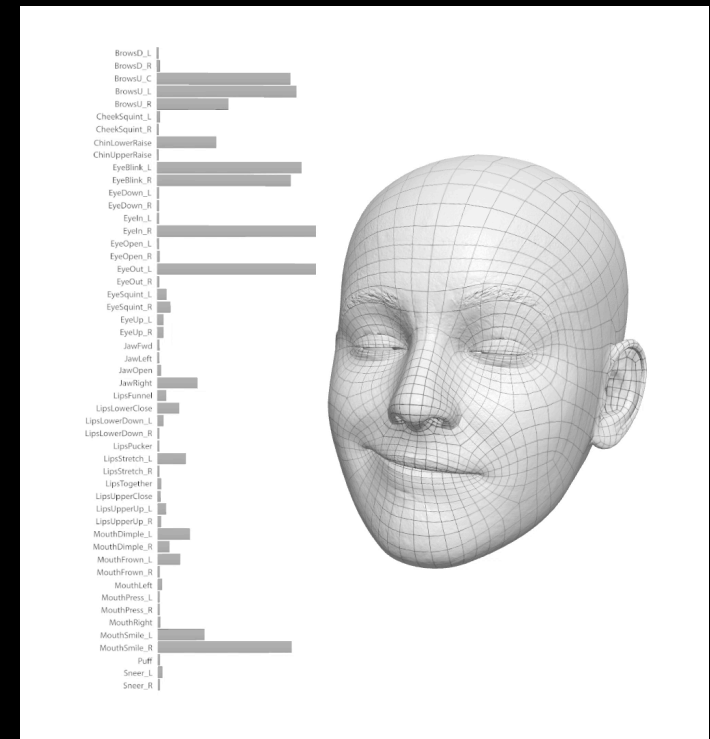Optimized for full system performance

# Animoji
Apple, 2017

# Realtime Facial Animation

# Blendshape Model

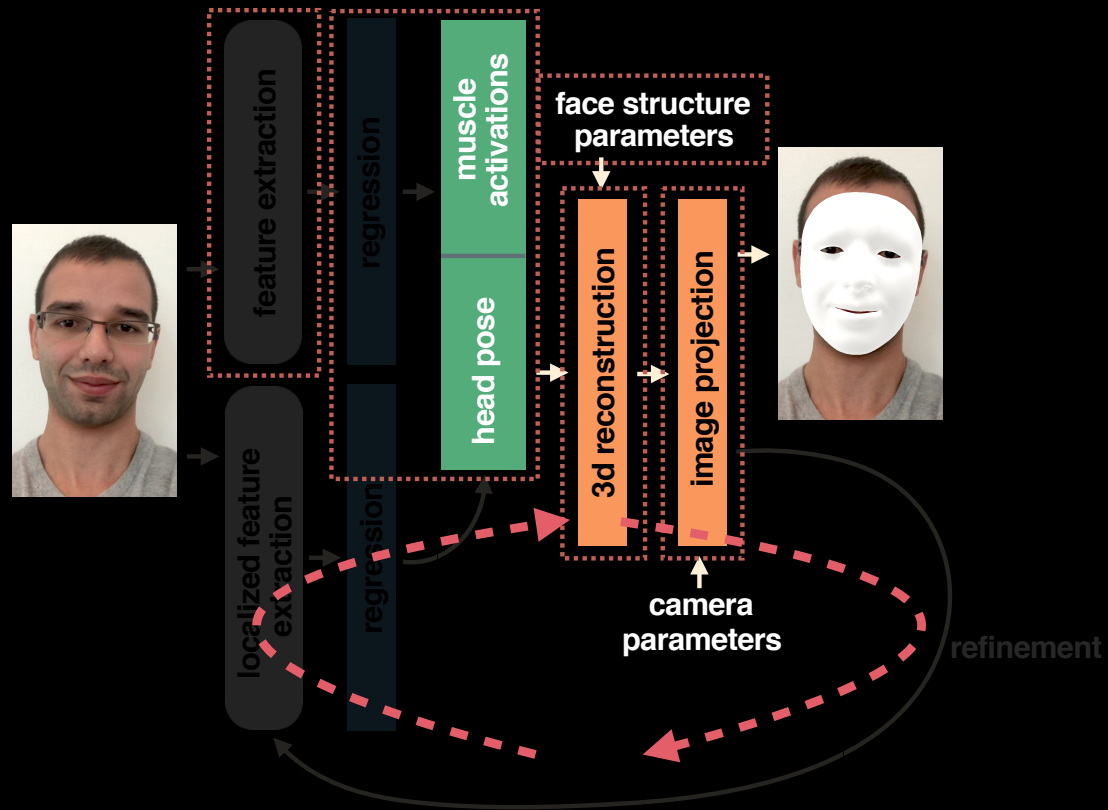51 blendshapes ("muscles") driving more
than 100 shapes

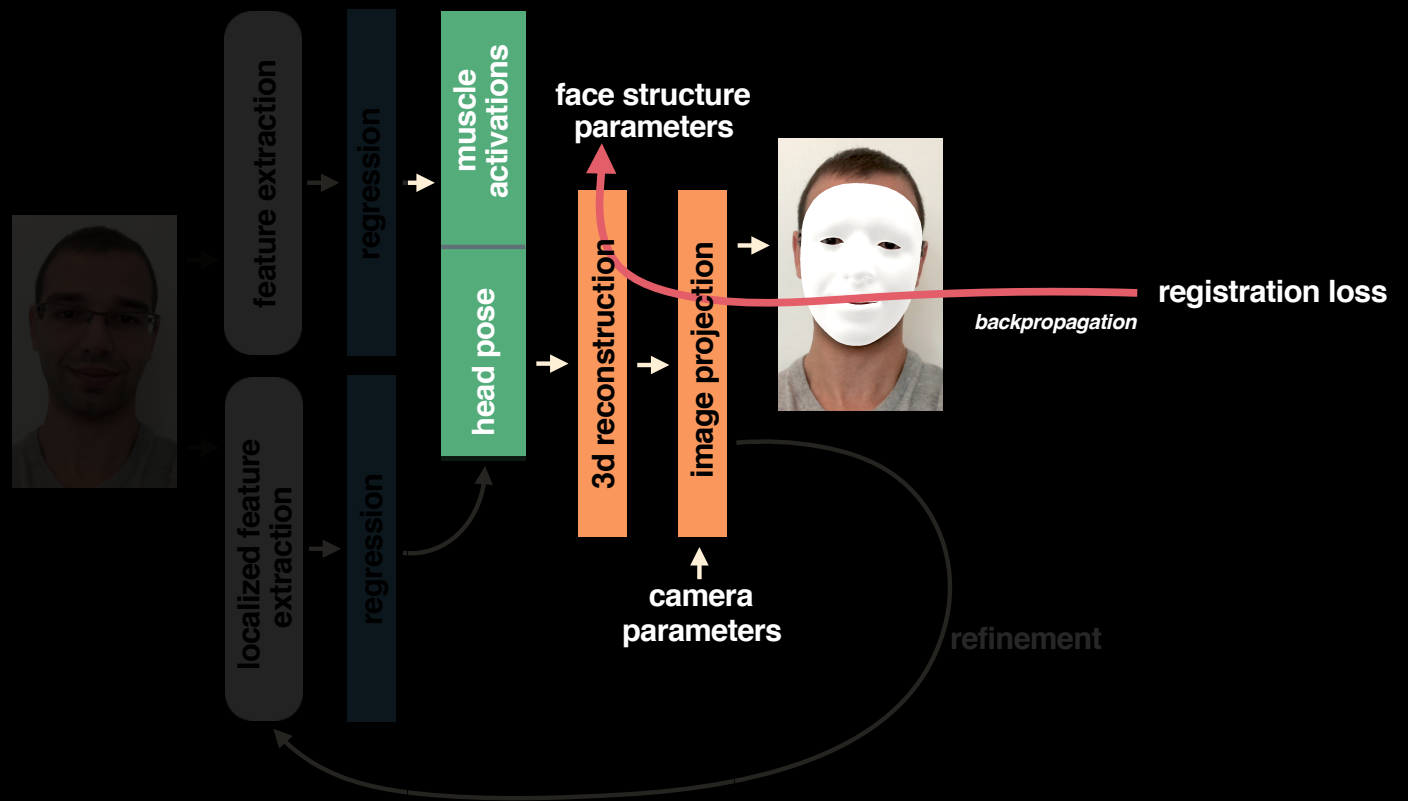# Animojis Driven by Blendshape Model
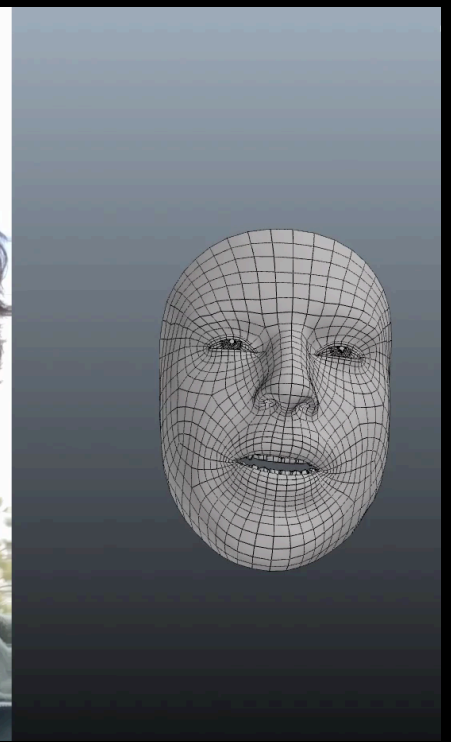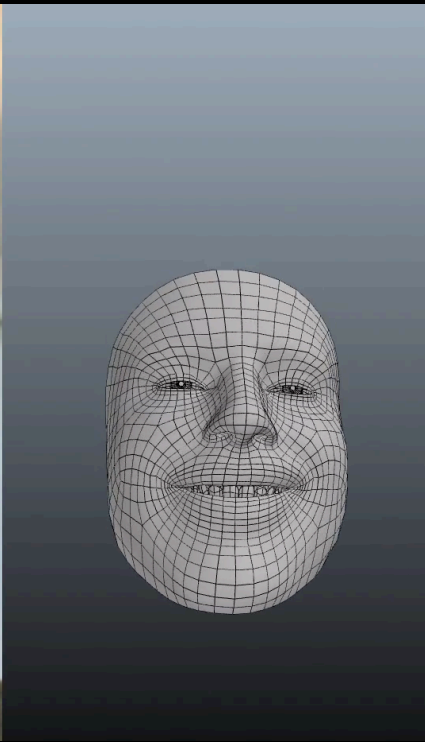
# Realtime Facial Animation
Model-based RNN

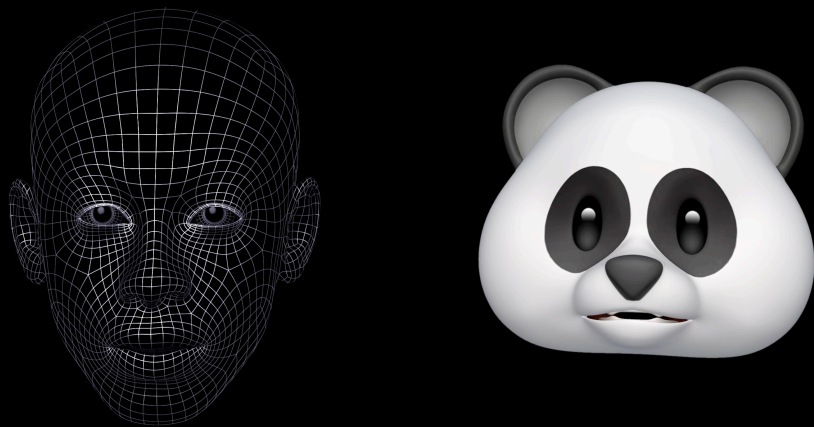# Online Identity Adaptation

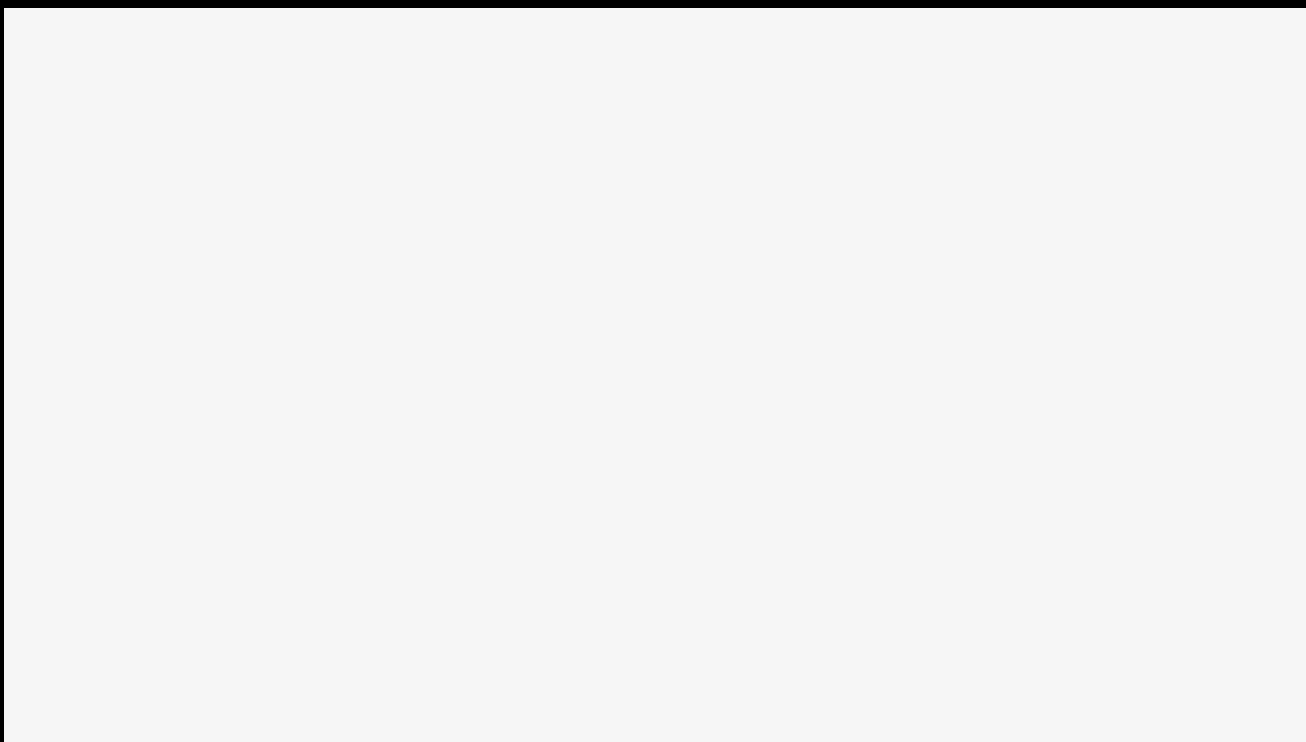Geometric backpropagation
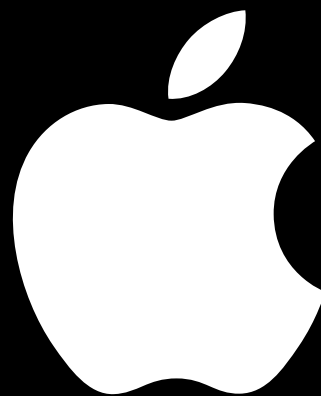
# Results

Indoor and outdoor

# Animoji
Performance



The animation runs sustainably at 60fps

And of course…

Animoji Karaoke