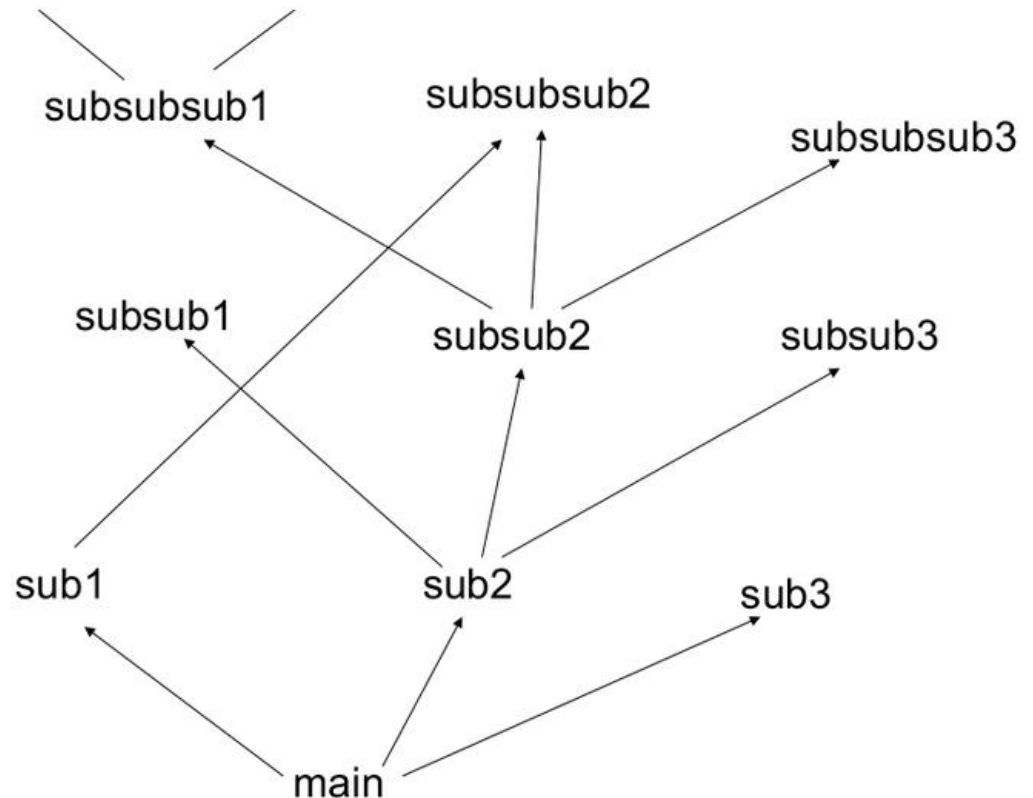# Why Deep Learning?

# Modularization

- Deep → Modularization

Don't put everything in your main function.



subsubsub1  subsubsub2  subsubsub3

subsub1  subsub2  subsub3

sub1  sub2  sub3

main

# Modularization



- Deep→Modularization

- Each basic classifier can have sufficient training examples.
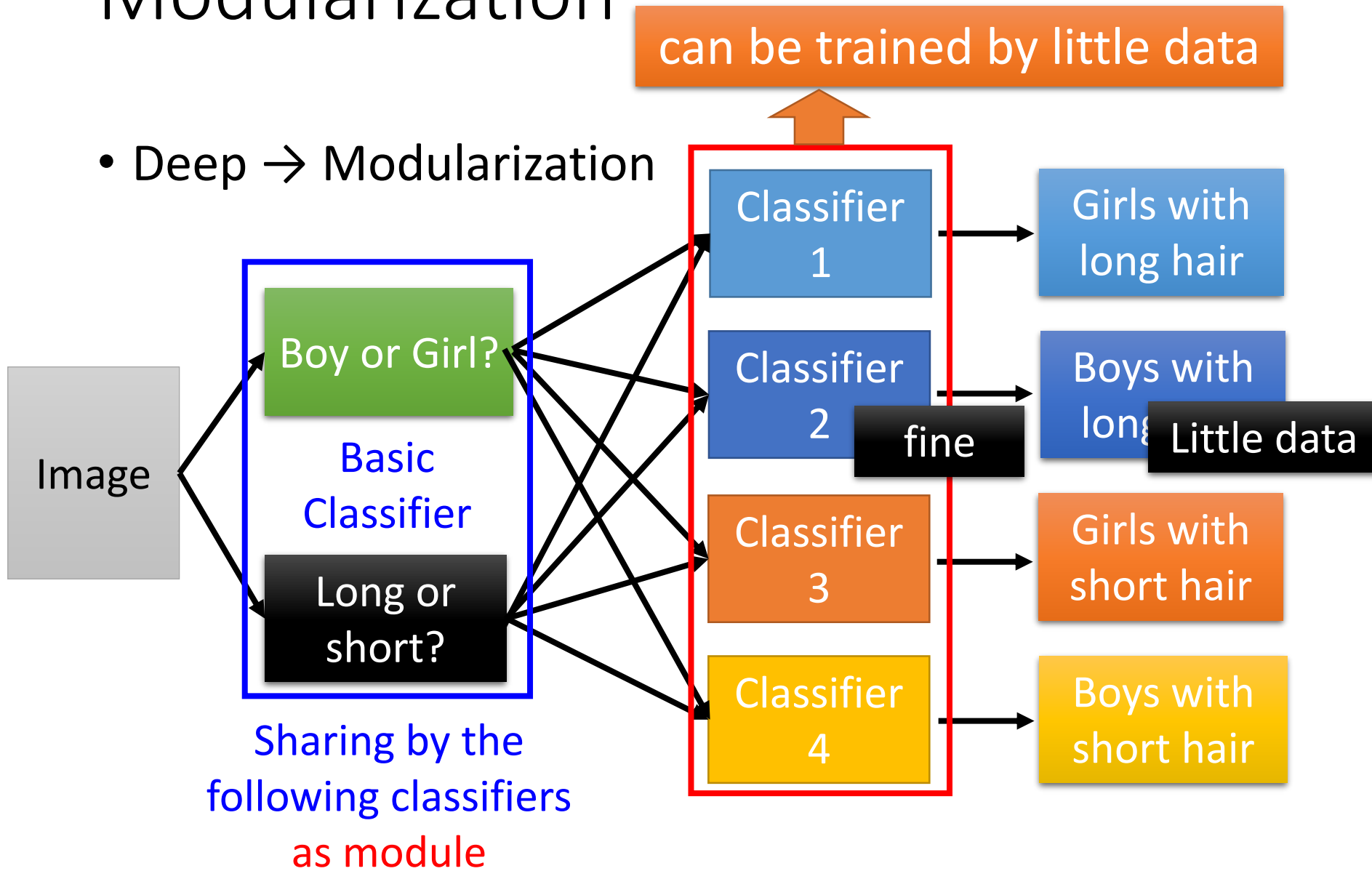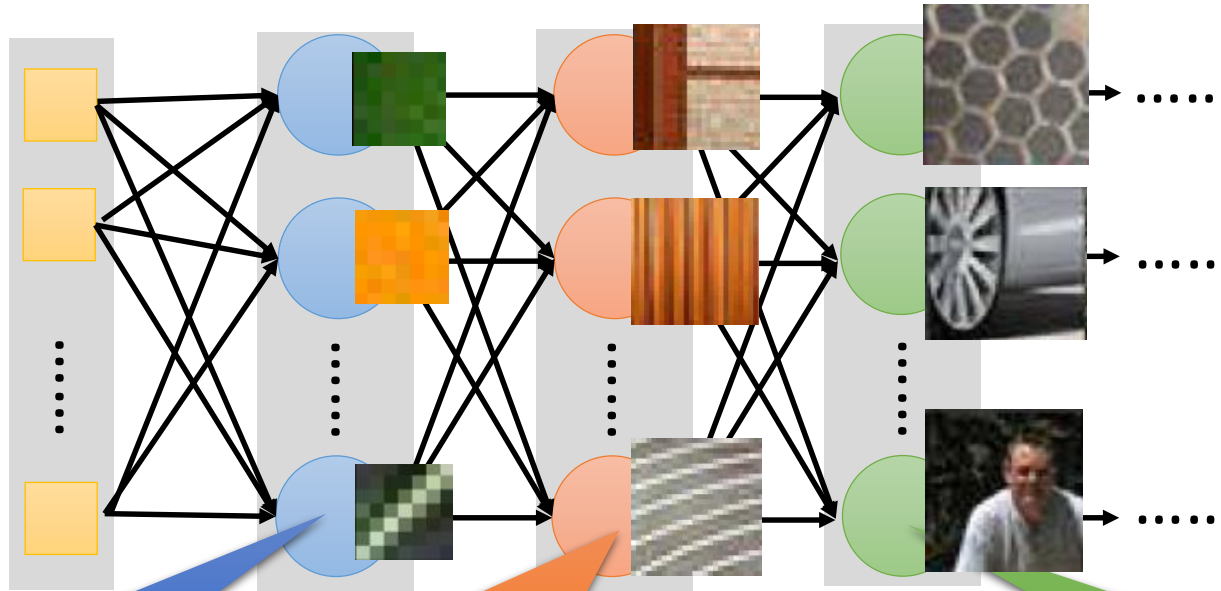
# Modularization

- Deep → Modularization

# Modularization - Image

- Deep → Modularization



The most basic classifiers

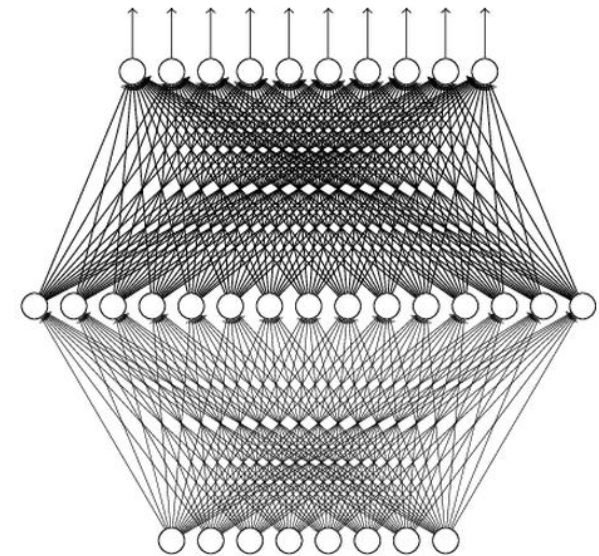Use 1st layer as module to build classifiers

Use 2nd layer as module ……

Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014* (pp. 818-833)

# Universality Theorem

Any continuous function f

$$f : R^N \rightarrow R^M$$

Can be realized by a network
with one hidden layer

(given **enough** hidden neurons)



Reference for the reason:
http://neuralnetworksandde
eplearning.com/chap4.html

Yes, shallow network can represent any function.

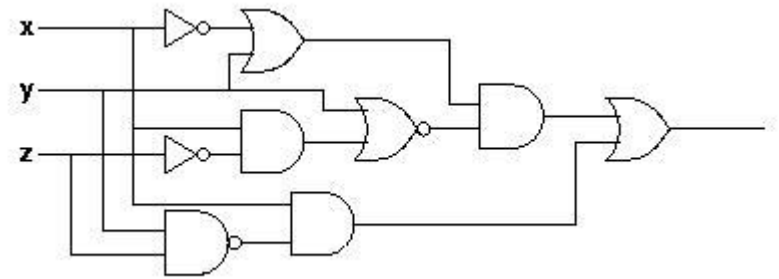However, using deep structure is more effective.

# Analogy



**Logic circuits**

**Neural network**

- Logic circuits consists of **gates**

- **A two layers of logic gates** can represent **any Boolean function.**

- Using multiple layers of logic gates to build some functions are much simpler

- Neural network consists of **neurons**

- **A hidden layer network** can represent **any continuous function.**

- Using multiple layers of neurons to represent some functions are much simpler
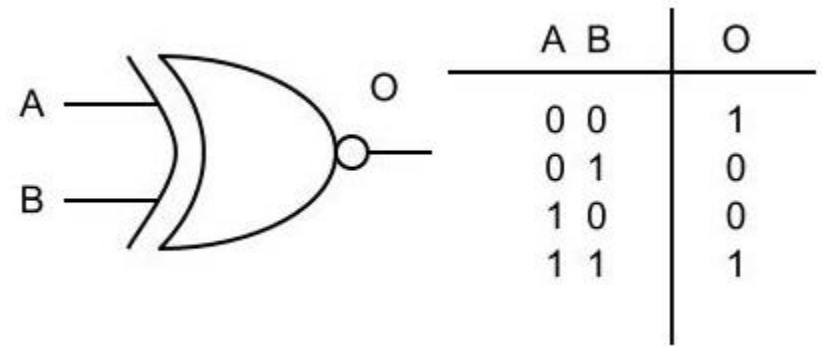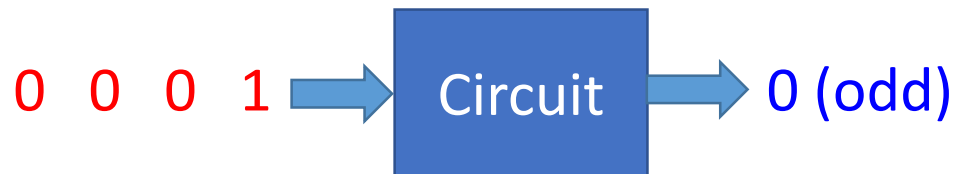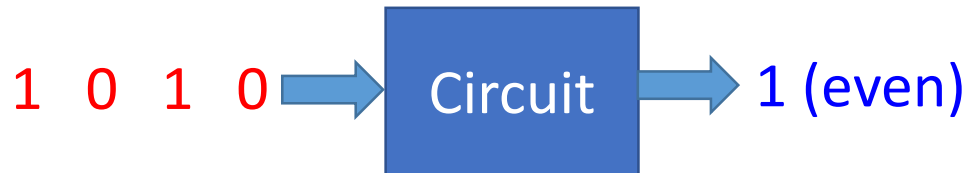
less gates needed

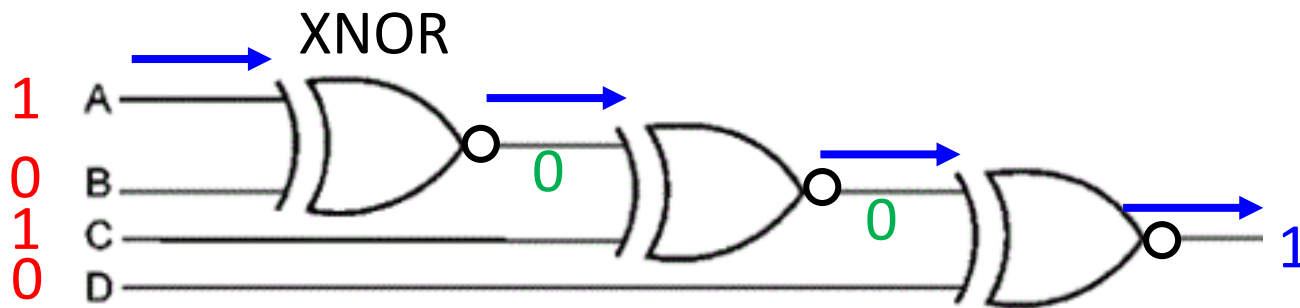less parameters

less data?

This page is for EE background.

# Analogy



| A B | O |
|-----|---|
| 0 0 | 1 |
| 0 1 | 0 |
| 1 0 | 0 |
| 1 1 | 1 |

- E.g. ***parity check***

1  0  1  0  → Circuit → 1 (even)

0  0  0  1  → Circuit → 0 (odd)

For input sequence with d bits,

Two-layer circuit need $O(2^d)$ gates.

XNOR

1 A
0 B
1 C
0 D

0    0    1
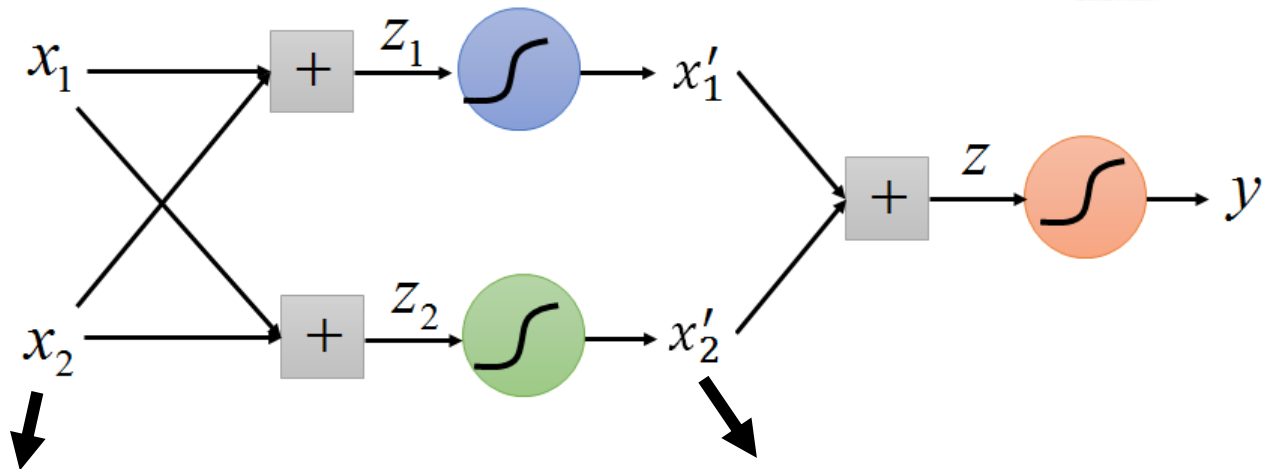
With multiple layers, we need only $O(d)$ gates.

# More Analogy

$$x_1 \rightarrow \boxed{+} \xrightarrow{z_1} \bigcirc \rightarrow x_1'$$

$$x_2 \rightarrow \boxed{+} \xrightarrow{z_2} \bigcirc \rightarrow x_2'$$

$$\boxed{+} \xrightarrow{z} \bigcirc \rightarrow y$$
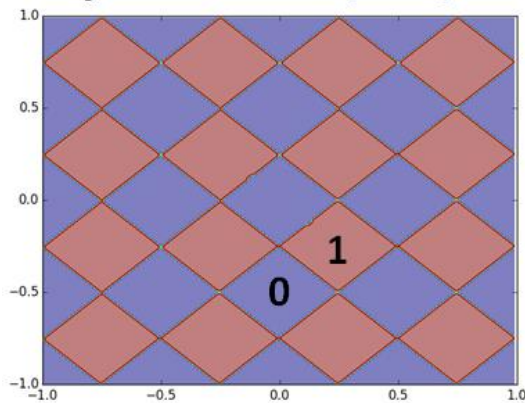
Folding the space

(0.73, 0.05)

(0.27, 0.27)

(0.05, 0.73)

# More Analogy - Experiment

**Different numbers of training examples**

10,0000                                    2,0000
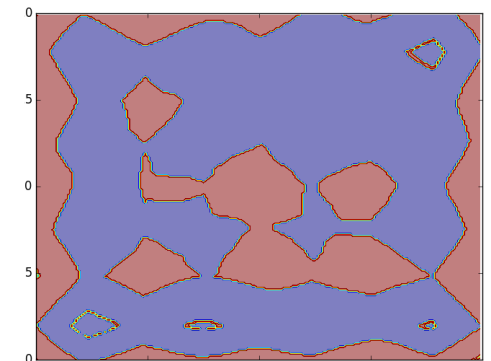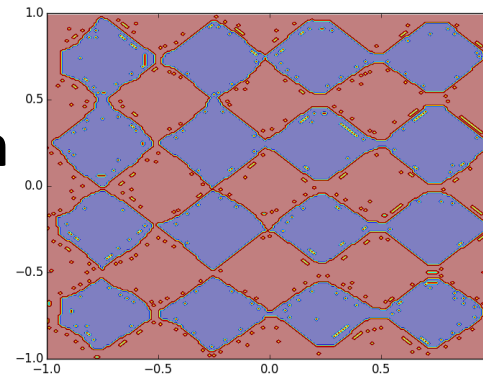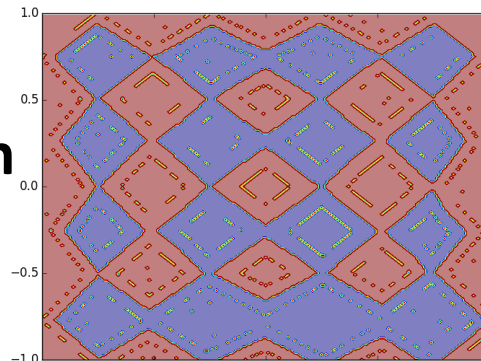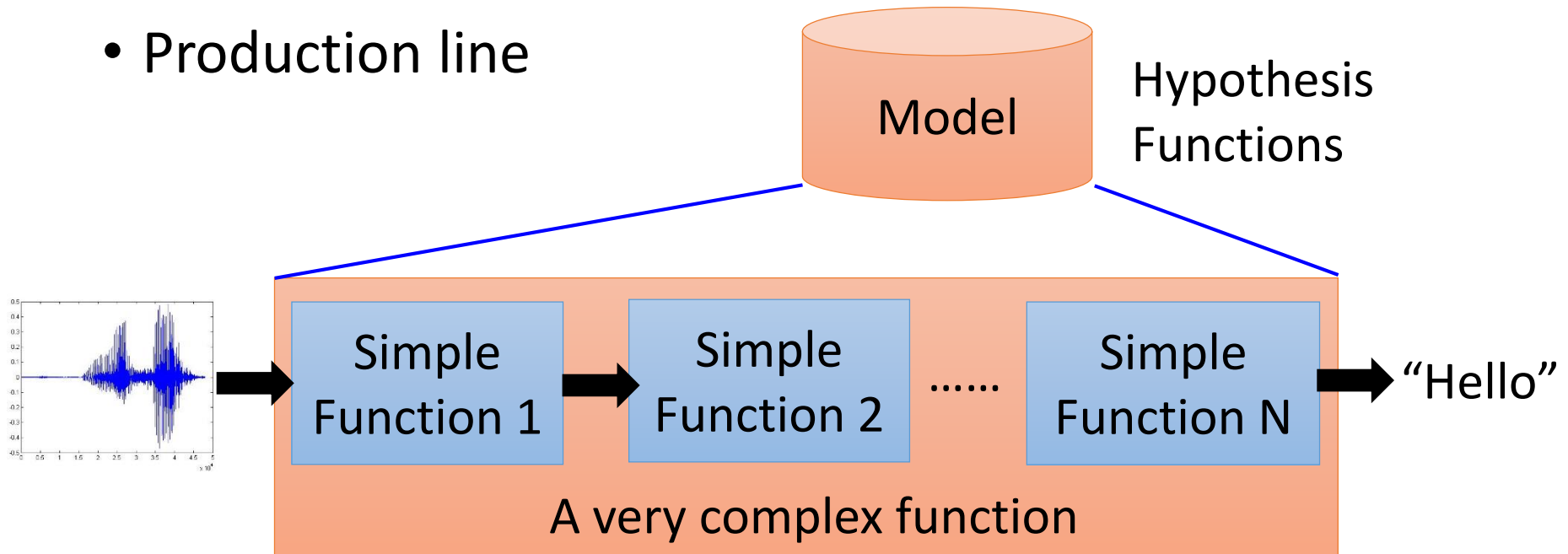
$f : R^2 \rightarrow \{0,1\}$

**1 hidden layer**

**3 hidden layers**

# End-to-end Learning

- Production line
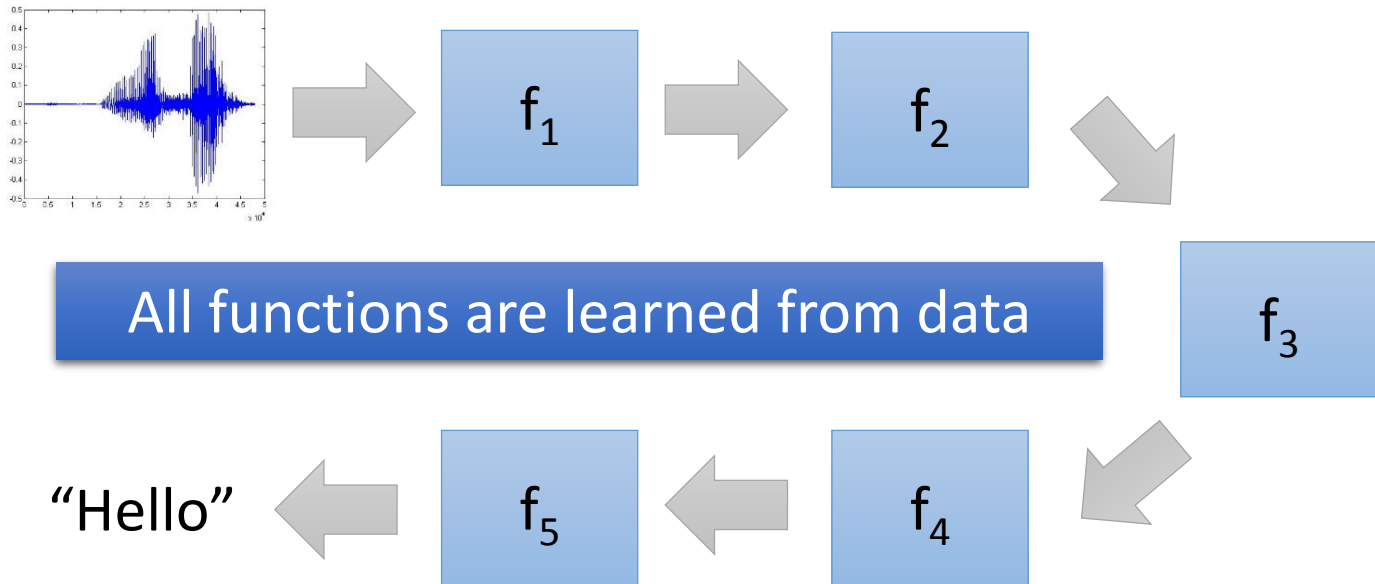


End-to-end training:
    What each function should do is learned automatically

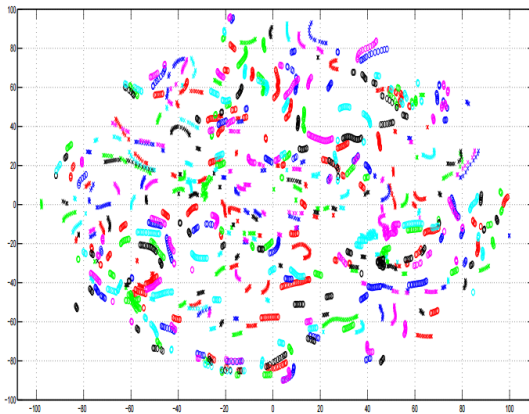# End-to-end Learning
# - Speech Recognition

• Deep Learning



$f_1$ → $f_2$ → $f_3$ → $f_4$ → $f_5$ → "Hello"

All functions are learned from data
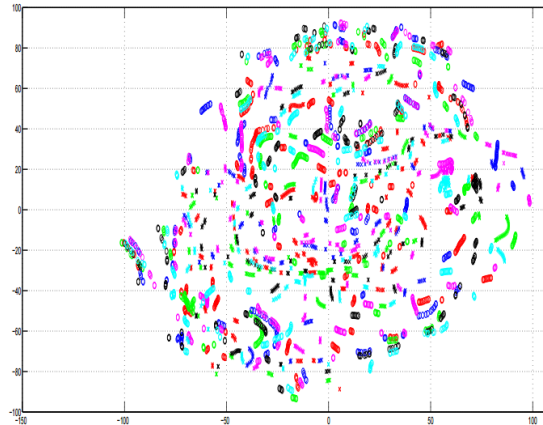
Less engineering labor, but machine learns more

# Complex Task …

A. Mohamed, G. Hinton, and G. Penn, "Understanding how Deep Belief Networks Perform Acoustic Modelling," in ICASSP, 2012.
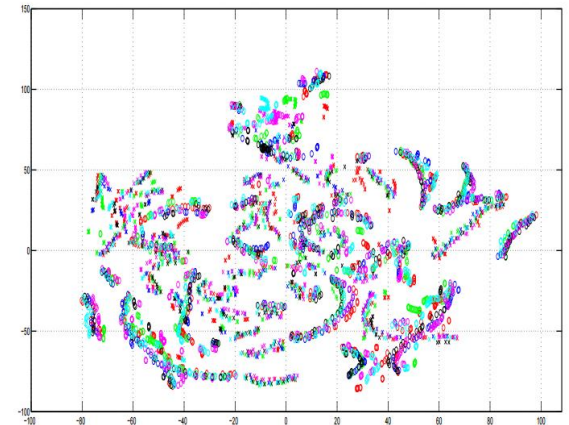
- Speech recognition: Speaker normalization is automatically done in DNN



Input Acoustic Feature (MFCC)

1-st Hidden Layer

8-th Hidden Layer