deeplearning.ai

# Error Analysis

Carrying out error analysis

# Look at dev examples to evaluate ideas



90% accuracy
→ 10% error

Should you try to make your cat classifier do better on dogs?

Error analysis: → 5-10 min
- Get ~100 mislabeled dev set examples.
- Count up how many are dogs.

"ceiling"

→ 5%
5/100
10%
9.5%

→ 50%
50/100

10%
↓
5%

Andrew Ng

# Evaluate multiple ideas in parallel

Ideas for cat detection:

- Fix pictures of dogs being recognized as cats ←

- Fix great cats (lions, panthers, etc..) being misrecognized ←
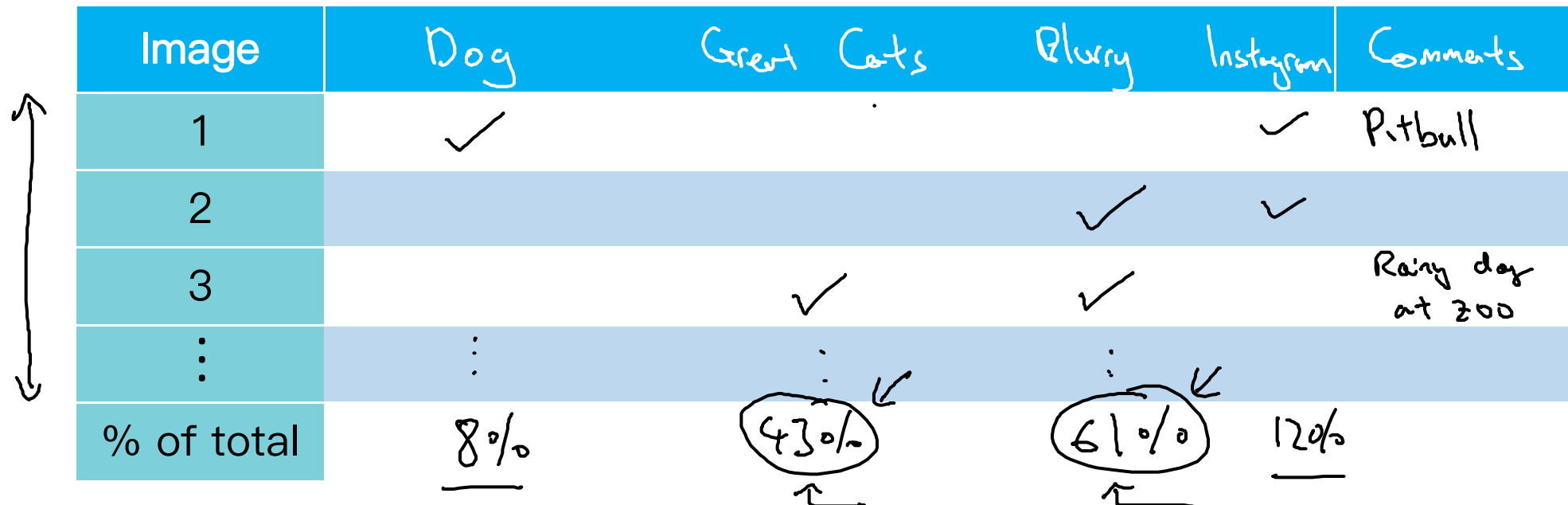
- Improve performance on blurry images ←

| Image | Dog | Great Cats | Blurry | Instagram | Comments |
|---|---|---|---|---|---|
| 1 | ✓ | | | ✓ | Pitbull |
| 2 | | | ✓ | ✓ | |
| 3 | | ✓ | ✓ | | Rainy day at zoo |
| ⋮ | ⋮ | ⋮ | ⋮ | | |
| % of total | 8% | 43% | 61% | 12% | |

Andrew Ng

deeplearning.ai

# Error Analysis

---

# Cleaning up Incorrectly labeled data

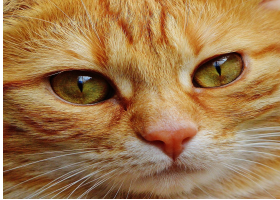# Incorrectly labeled examples

x 

y     1     0     1     1     0     ①     1

Training Set.

DL algorithms are quite robust to random errors in the training set.

Systematic errors

Andrew Ng

# Error analysis

| Image | Dog | Great Cat | Blurry | Incorrectly labeled | Comments |
|-------|-----|-----------|--------|---------------------|----------|
| ... | | | | | |
| 98 | | | | ✓ | Labeler missed cat in background |
| 99 | | ✓ | | | |
| 100 | | | | ✓ | Drawing of a cat; Not a real cat. |
| % of total | 8% | 43% | 61% | 6% | |

Overall dev set error ............... 10%          2%

Errors due incorrect labels ........ 0.6% ←        0.6%

Errors due to other causes ........ 9.4% ←        1.4%

                                                   2.1%      1.9%

Goal of dev set is to help you select between two classifiers A & B.

Andrew Ng

# Correcting incorrect dev/test set examples

- Apply same process to your dev and test sets to make sure they continue to come from the same distribution

- Consider examining examples your algorithm got *98.%* right as well as ones it got wrong. *2.%*

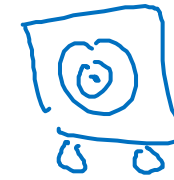- Train and dev/test data may now come from slightly different distributions.

deeplearning.ai

Error Analysis

Build your first system quickly, then iterate

# Speech recognition example

- Noisy background
  - Café noise
  - Car noise
- Accent
- Far from microphone
- Young children's speech
- Stuttering
- ...

Guideline:

**Build your first system quickly, then iterate**

- Set up dev/test set and metric
- Build initial system quickly
- Use Bias/Variance analysis & Error analysis to prioritize next steps.

Andrew Ng

deeplearning.ai
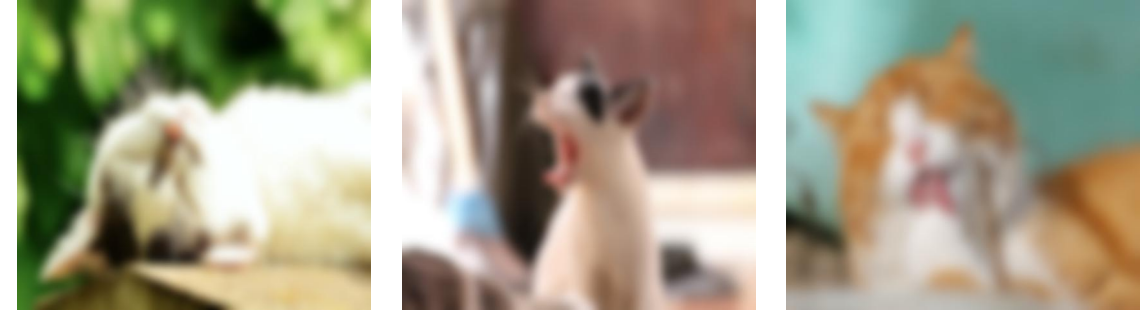
# Mismatched training and dev/test data

# Training and testing on different distributions

# Cat app example

## Data from webpages

## Data from mobile app

Care about this



$\approx 200,000$     210,000     $\approx 10,000$

(shuffle

**Option 1:**

| train | dev | test |
| --- | --- | --- |

205,000     2,500    2,500

$\frac{200K}{210K}$

2381 - web
119 - mobile app

**Option 2:**

| train | dev | test |
| --- | --- | --- |

web

train: 205,000

app    app 2500   app 2500

Andrew Ng

# Speech recognition example



Speech activated rearview mirror

## Training

Purchased data   $x, y$

Smart speaker control

Voice keyboard

...

500,000 utterances

## Dev/test

Speech activated rearview mirror

→ 20,000

10K   5K 5K  D T

train

500 K

510K   D T

10K mirror   5K 5K
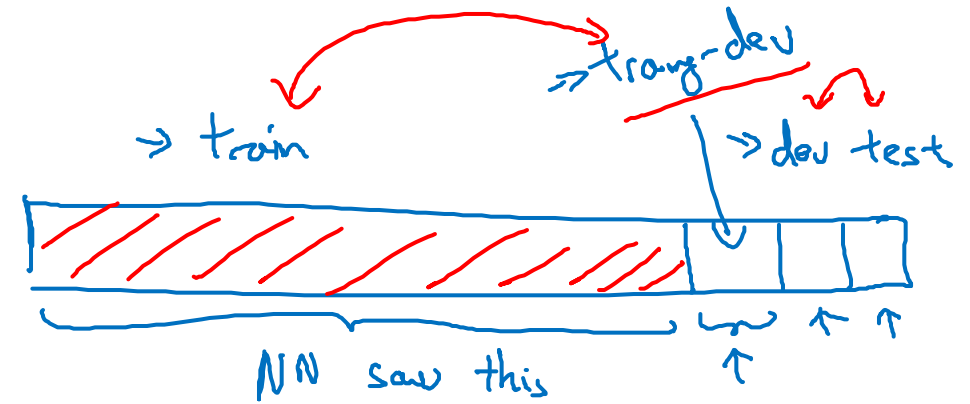
deeplearning.ai

# Mismatched training and dev/test data

Bias and Variance with mismatched data distributions

# Cat classifier example

Assume humans get ≈ 0% error.

Training error ..... 1%

Dev error ........... 10% ↓ 9%

Training–dev set: Same distribution as training set, but not used for training

→ train

→ traing-dev

→ dev test

NN saw this

| | | |
|---|---|---|
| Traing error | 1% ↑ | 1% ↑ |
| → Traing-dev error | 9% ↓ Variance | 1.5% ↓ |
| → Dev error | 10% | 10% ↓ data mismatch |

Variance

| | | | |
|---|---|---|---|
| Human error | ---- 0% ↓ Avoidable | | ↓ Avoidable bias |
| Traing error | 10% ↓ bias | 10% ↑ | |
| Traing-dev error | 11% | 11% ↓ Variance | |
| Dev error | 12% | 20% ↑ Data mismatch | |
| | Bias | Bias + Data mismatch | |

Andrew Ng

# Bias/variance on mismatched training and dev/test sets

Human level    4%

Traing set error   7%

Traing - dev set error 10%

→ Dev error    12%

→ Test error    12%

avoidable bias

variance

data mismatch

degree of overfitting to dev set.

4%

7% }

10% }

6% }

6% }

Andrew Ng

# More general formulation



Rearview Mirror

|  | General speech recognition | Rearview mirror speech data |
|---|---|---|
| Human level | "Human level" 4% | 6% |
| Error on examples trained on | "Training error" 7% | 6% |
| Error on examples not trained on | "Training-dev error" 10% | Dev/Test error 6% |

Avoidable bias

Variance

data mismatch

Andrew Ng

deeplearning.ai

Mismatched training and dev/test data

Addressing data mismatch

# Addressing data mismatch

→ • Carry out manual error analysis to try to understand difference between training and dev/test sets

E.g. noisy — car noise          street numbers

→ • Make training data more similar; or collect more data similar to dev/test sets

E.g. Simulate noisy in-car data

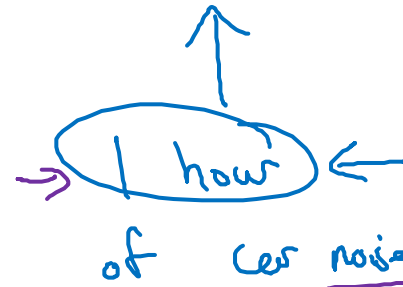Andrew Ng

# Artificial data synthesis


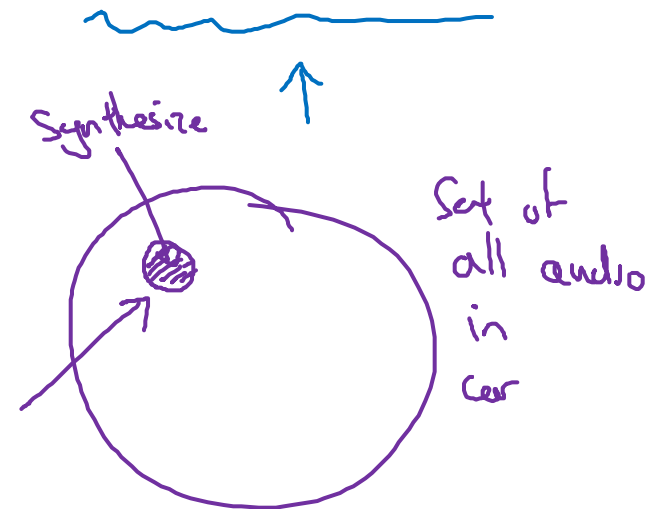
"The quick brown ← fox jumps over the lazy dog." + Car noise = Synthesized in-car audio

10,000 hours

1 hour of car noise

Overfit to 1 hour of car noise

10,000 hours

Synthesize

Set of all audio in car

Andrew Ng

# Artificial data synthesis

Car recognition:



$\approx 20$ cars

Synthesized

All cars

Learning from multiple tasks

Transfer learning

deeplearning.ai

# Transfer learning



image recognition $(x, y)$
↑
pre-training

→ $(x, y)$ — fine-tuning
↑          ↖
radiology   diagnoses
image

x → □□□ → □□□□□ → □□□□□ → □□□ → □□□ → □□□ → ⊗ → $\hat{y}$

$w^{[L]}, b^{[L]}$

image recognition
→ [1,000,000]   [100] ←
                ↓
radiology diagnosis
→ [100]   [1000]

x → □□□ → □□□□□ → □□□□□ → □□□ → □□□ → □□□ → ⊗ → $\hat{y}$

audio

Speech recognition
10h   10,000h

wakeword/triggerword detection   1h
50h

Andrew Ng

# When transfer learning makes sense

Transfer from A → B

- Task A and B have the same input x.

- You have a lot more data for Task A than Task B.

- Low level features from A could be helpful for learning B.

deeplearning.ai

# Learning from multiple tasks

---

# Multi-task learning

# Simplified autonomous driving example



$X^{(i)}$

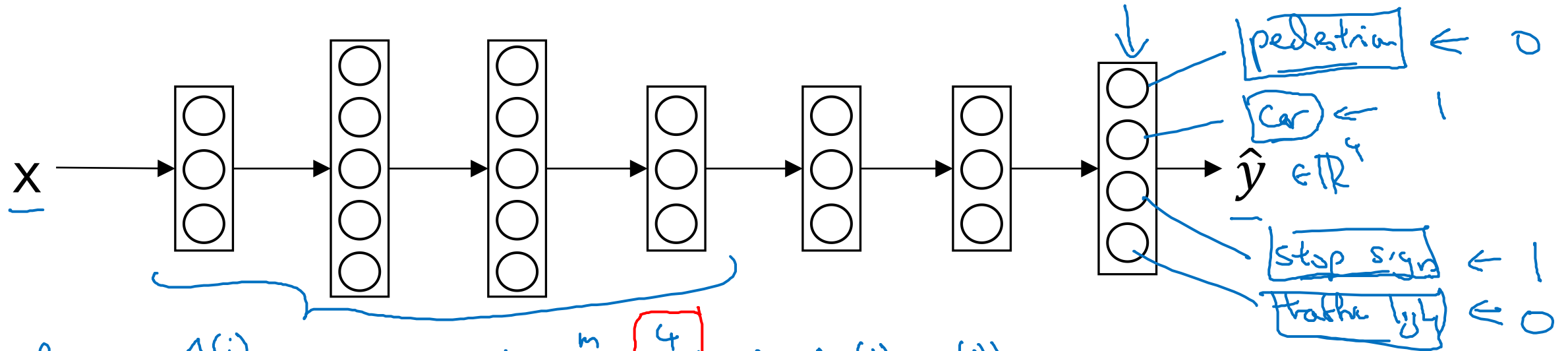Pedestrians

Cars

Stop signs

Traffic lights

$y^{(i)}$

$\begin{matrix} 0 \\ 1 \\ 1 \\ 0 \end{matrix}$  $(4,1)$

$$Y = \begin{bmatrix} | & | & | & & | \\ y^{(1)} & y^{(2)} & y^{(3)} & \cdots, & y^{(m)} \\ | & | & | & & | \end{bmatrix}$$

$(4,m)$

Andrew Ng

# Neural network architecture



$x \rightarrow \hat{y} \in \mathbb{R}^4$

pedestrian $\leftarrow 0$
Car $\leftarrow 1$
stop sign $\leftarrow 1$
traffic light $\leftarrow 0$

$$\text{Loss:} \quad \overset{\wedge(i)}{Y}_{(4,1)} \quad \rightarrow \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{4} \mathcal{L}\left(\hat{y}_j^{(i)}, y_j^{(i)}\right)$$
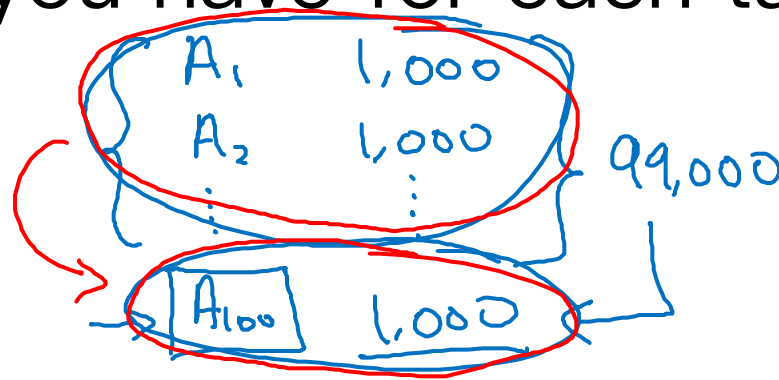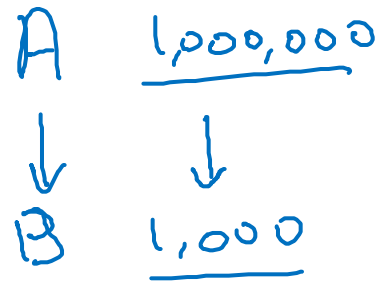
Sum only over value of $j$ with $0/1$ label.

$\rightarrow$ Usual logistic loss
$-y_j^{(i)} \log \hat{y}_j^{(i)} - (1-y_j^{(i)}) \log (1 - \hat{y}_j^{(i)})$

Multi-task learning $\leftarrow$

Unlike softmax regression:
One image can have multiple labels

$$Y = \begin{bmatrix} 1 & 1 & 0 & ? \\ 0 & 1 & 1 & ? \\ ? & ? & 1 & ? \\ ? & ? & 0 & ? \end{bmatrix} \leftarrow$$

Andrew Ng

# When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features.
- Usually: Amount of data you have for each task is quite similar.

$A \quad 1,000,000$

$B \quad 1,000$

$A_1 \quad 1,000$
$A_2 \quad 1,000$
$\vdots$

$99,000$

$A_{100} \quad 1,000$

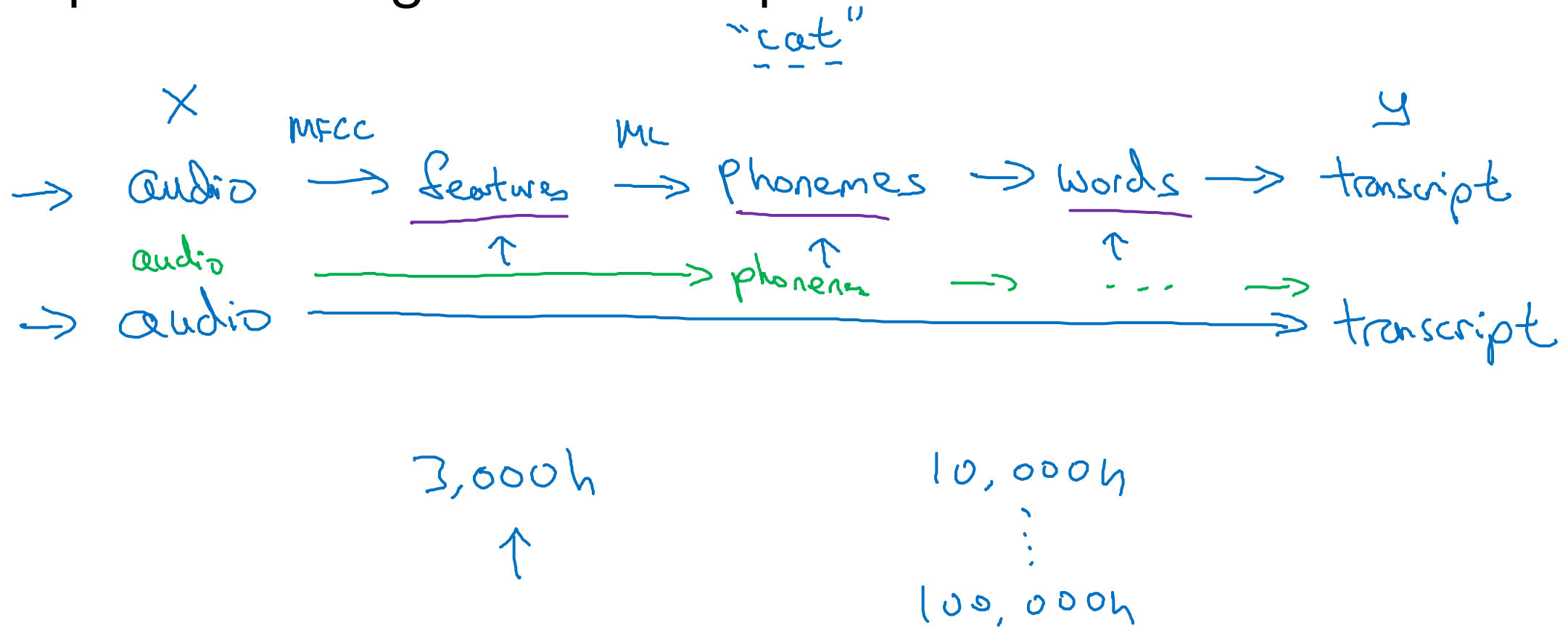- Can train a big enough neural network to do well on all the tasks.

deeplearning.ai

# End-to-end deep learning

---

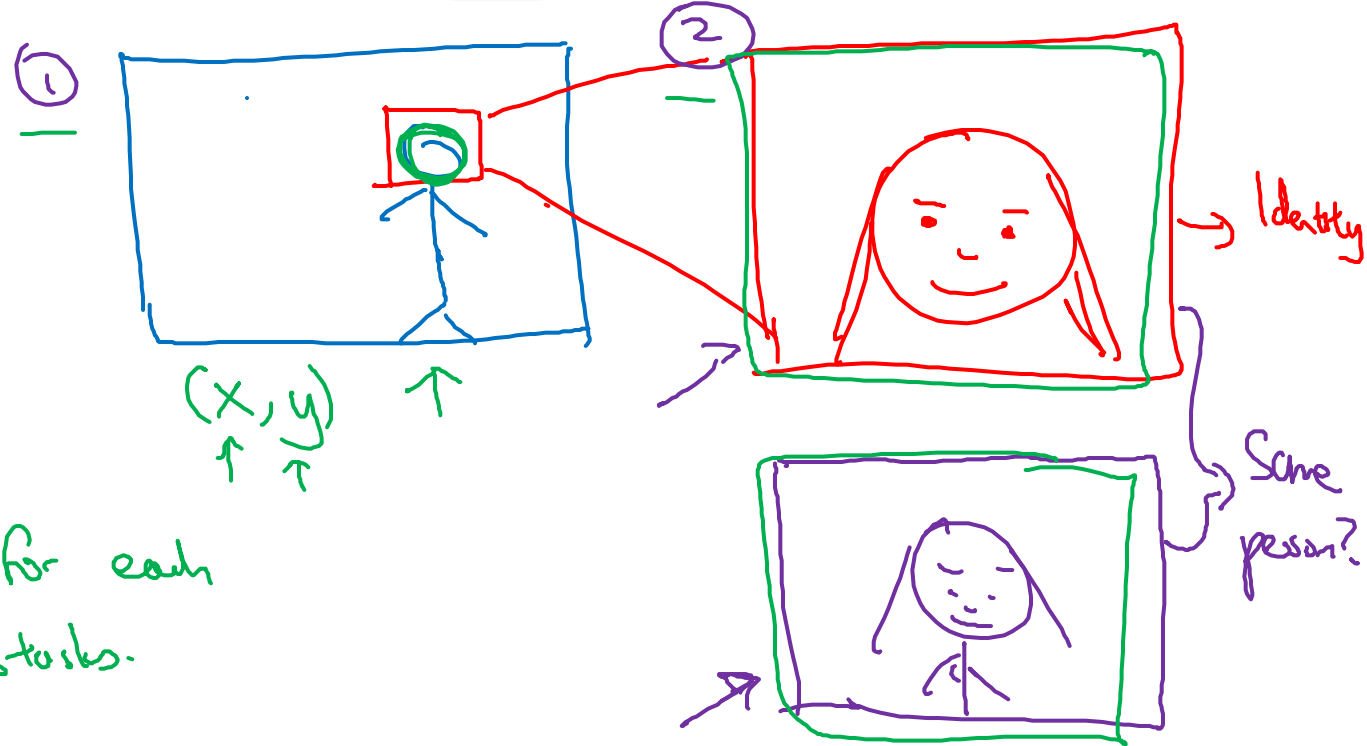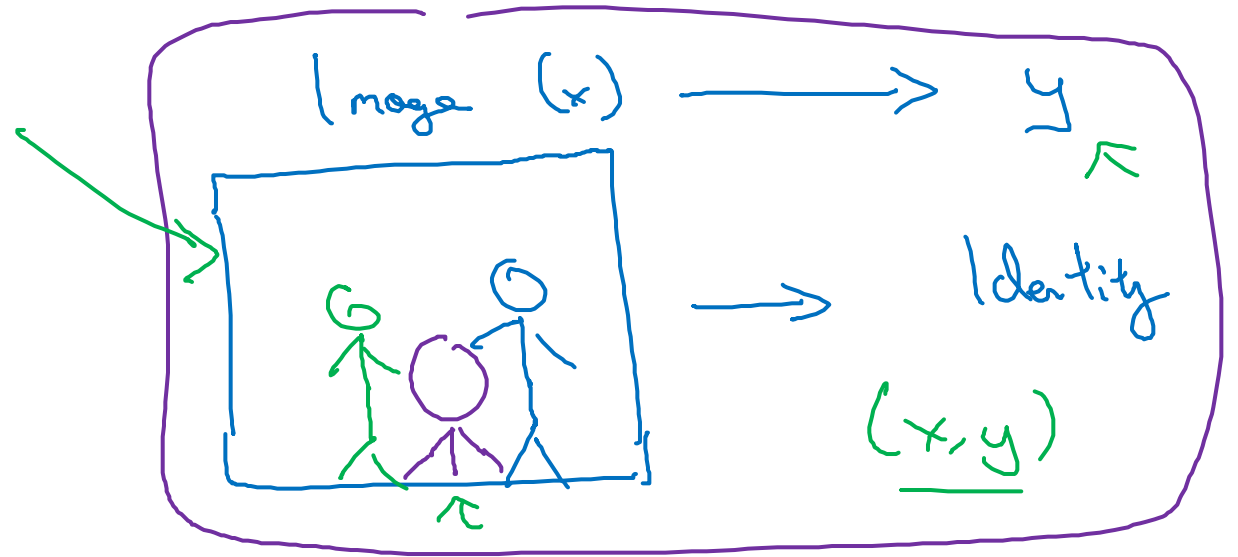## What is end-to-end deep learning

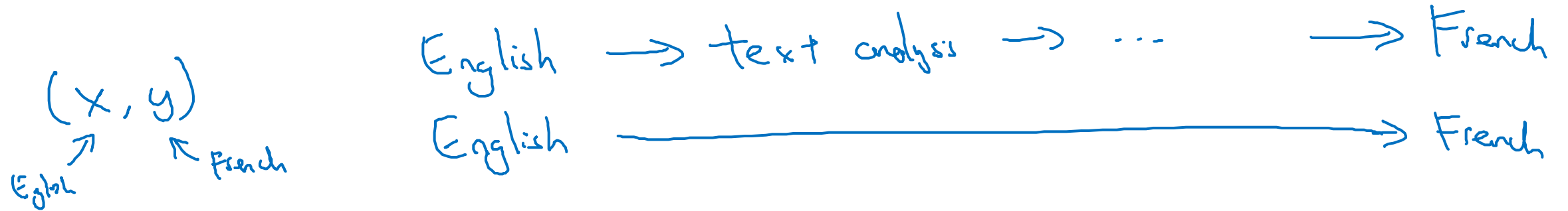# What is end–to–end learning?

Speech recognition example



"cat"

$X$     MFCC     ML     $y$

→ audio → features → Phonemes → Words → transcript

audio ————————→ phonemes —→ ···· —→ transcript

→ audio ——————————————————————→ transcript

3,000h

10,000h
⋮
100,000h

Andrew Ng

# Face recognition



[Image courtesy of Baidu]

Image (x) ⟶ y

Identity

(x, y)

① (x, y)

② Identity

Same person?

Have data for each of 2 subtasks.

Andrew Ng

# More examples

## Machine translation

$(x, y)$

English → French

English → text analysis → ... → French

English ⟶ French

## Estimating child's age:



Image ①→ bones ②→ age

Image ⟶ age ←

deeplearning.ai

# End-to-end deep learning

---

# Whether to use end-to-end learning

# Pros and cons of end–to–end deep learning

Pros:

- Let the data speak $\quad x \longrightarrow y$

- Less hand–designing of components needed

$\to$ "phonemes"
$\quad$ cat

$x - - \quad - \quad - \to y$

input end $\quad$ output end

Cons:

- May need large amount of data $\quad x \longrightarrow y \quad (x, y)$

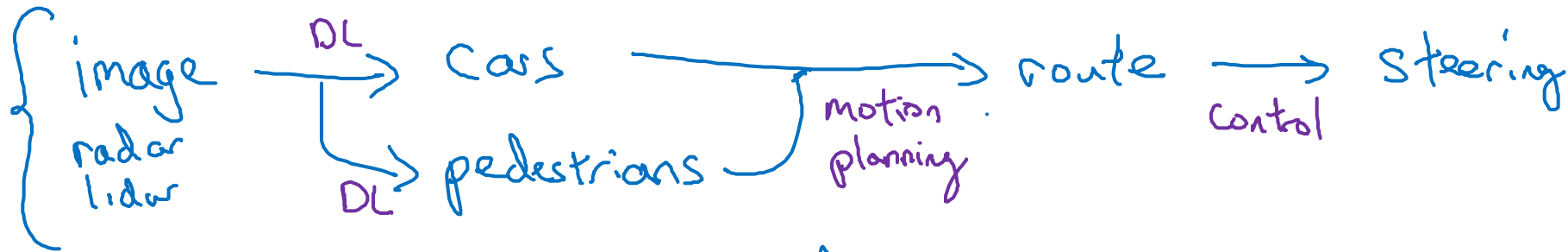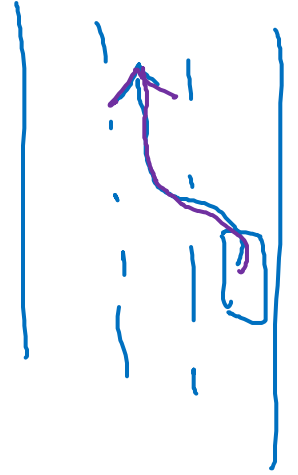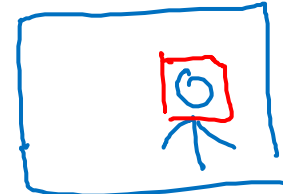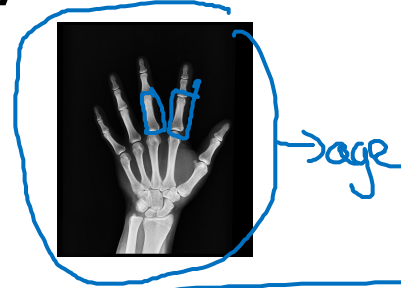- Excludes potentially useful hand–designed components

Data. $\qquad$ Hand–design

# Applying end–to–end deep learning

Key question: Do you have sufficient data to learn a function of the complexity needed to map x to y?

$X \rightarrow y$

$\rightarrow$ age

image
radar
lidar $\xrightarrow{DL}$ cars

$\xrightarrow{DL}$ pedestrians

motion planning $\rightarrow$ route $\xrightarrow{Control}$ steering

- Use DL to learn individual components
- Carefully choose $X \rightarrow Y$ depending what tasks you can get data for.

$\rightarrow$ image $\xrightarrow{\hspace{6cm}}$ steering

Andrew Ng