



deeplearning.ai

Error Analysis

Carrying out error analysis

analysis

Look at dev examples to evaluate ideas



90% occur
→ 10% error

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

Error analysis:

5-10 mis

"ceiling"

- Get ~100 mislabeled dev set examples.
- Count up how many are dogs.

$$\begin{array}{r} \text{5\%} \\ \text{10\%} \\ \text{50\%} \\ \text{50/100} \\ \text{5\%} \\ \hline \text{5/100} \end{array}$$

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		✓	✓	✓	Rainy day at 200
3	✓		✓		
:	:	⋮	⋮	⋮	
% of total	8%	43%	61%	12%	



Error Analysis

Cleaning up
Incorrectly labelled
data

Incorrectly labeled examples



y

— 1 — 0 —

— 1 — 1 —

— 0 —

1

Training set.

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

Systematic errors

Error analysis



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



Overall dev set error

10%

2%

Errors due to 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.

Correcting incorrect dev/test set examples

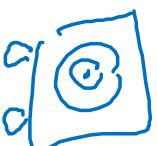
- Apply same process to your dev and test sets to make sure they continue to come from the same distribution
- Consider examining examples your algorithm got right as well as ones it got wrong. { 2%
- Train and dev/test data may now come from slightly different distributions.



Error Analysis

Build your first system
quickly, then iterate
deeplearning.ai

Speech recognition example



- Noisy background
 - Café noise
 - Car noise
- Accent Far from Guideline:
- Young
- Stutter
- ...

Build your first system quickly, then iterate

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



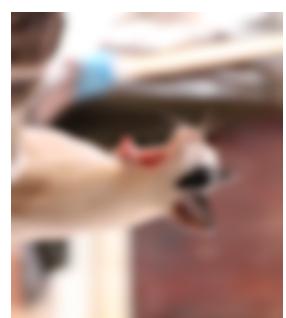
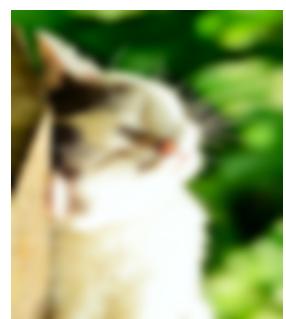
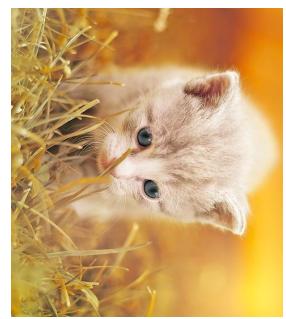
Mismatched training
and dev/test data

Training and testing
on different
distributions

deeplearning.ai

Cat app example

Data from webpages



Data from mobile app

care about this

~~X~~ Option 1:

train: $\approx 200,000$

test: $\approx 10,000$

(shuffle)

train: $\approx 205,000$

test: $\approx 2,500$

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

app: ≈ 2500

web: $\approx 2381 - \text{web}$

119 - mobile app

Option 2:

train: $\approx 205,000$

test: $\approx 2,500$

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

app: ≈ 2500

web: $\approx 2381 - \text{web}$

119 - mobile app

Speech recognition example

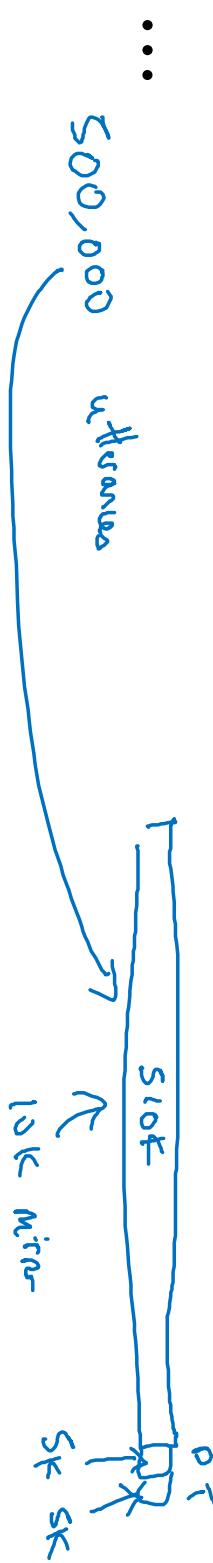
Speech activated review mirror

Training

Purchased data $\downarrow \downarrow$
 x, y

Smart speaker control

Voice keyboard



...

500,000 utterances

Dev/test

Speech activated
rearview mirror





Mismatched training
and dev/test data

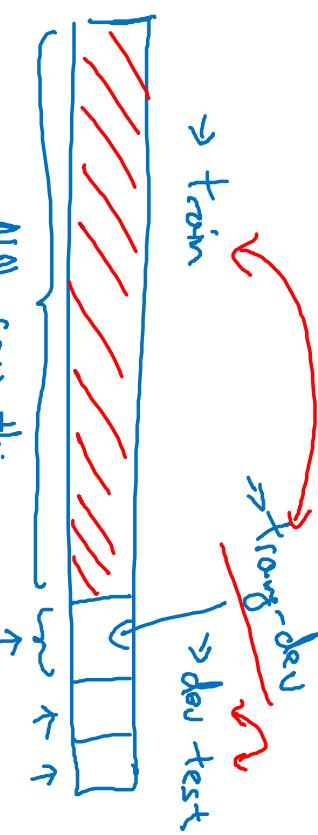
Bias and Variance with
mismatched data
distributions

deeplearning.ai

Cat classifier example

Assume humans get $\approx 0\%$ error.

Training error 10% $\downarrow 1\%$
Dev error 10%



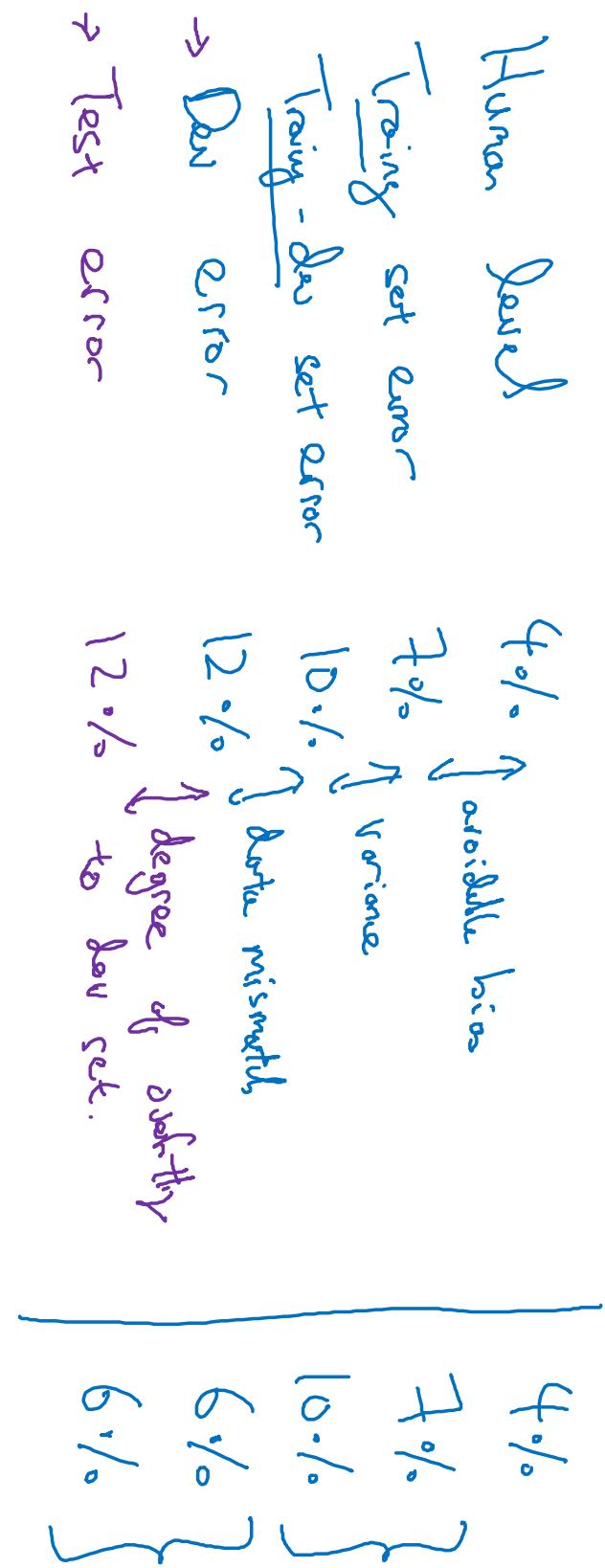
Training error 10% \uparrow variance
 \rightarrow Train-dev error 9% \uparrow bias
 \rightarrow Dev error 10% \uparrow data mismatch

Training-dev set: Same

distribution as training set, but not used for training

	Variance	Bias
Human error	0%	
Training error	$10\% \uparrow$ Available bias	$10\% \uparrow$ Available bias
Train-dev error	$1\% \uparrow$ Variance	$1\% \uparrow$ Variance
Dev error	$12\% \uparrow$ Data mismatch	$20\% \uparrow$ Data mismatch

Bias/variance on mismatched training and dev/test sets



More general formulation

Resumes mirror

General speech
recognition

Resumes mirror
speech data.

General speech
recognition

Resumes mirror
speech data.

Human level

"Human level" 4%

6%

Error on
examples tested

"Training error" 7%

6%

Error on
examples tested

"Training - dev
error"

10%

Dev / Test 6%

Actual mismatch

Avoidable bias

Variance



Mismatched training
and dev/test data

Addressing data
mismatch

deeplearning.ai

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

Artificial data synthesis



+



=



“The quick brown fox jumps over the lazy dog.”

←
10,000 hours

→
1 hour
of car noise

↓
Draft to 1 hour &
car noise
10,000 hours

Synthesized
in-car audio

Synthesize →

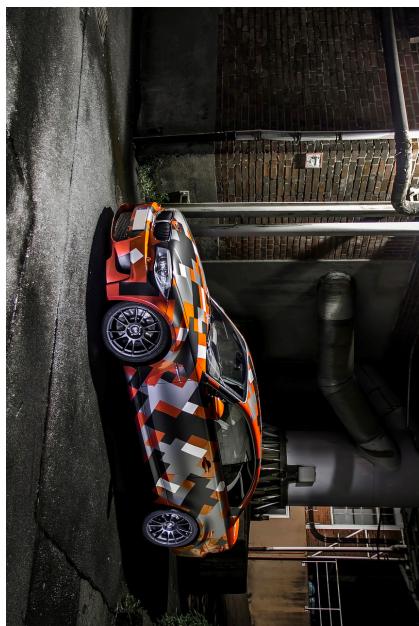
Set of
all
audio
in
car

Artificial data synthesis

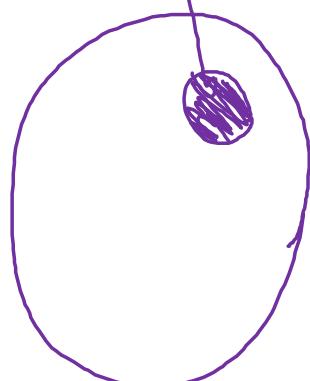
Car recognition:



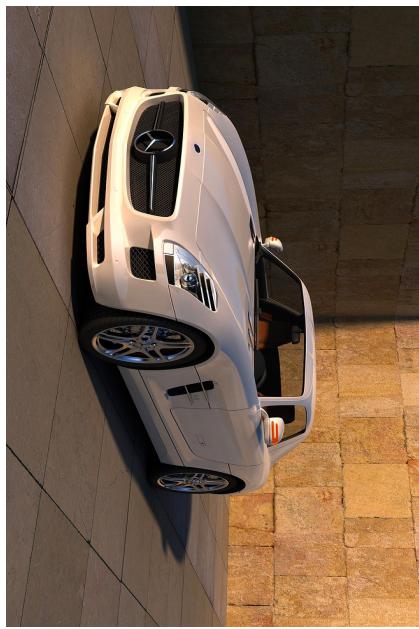
≈ 20 cars



Synthesized



All cars



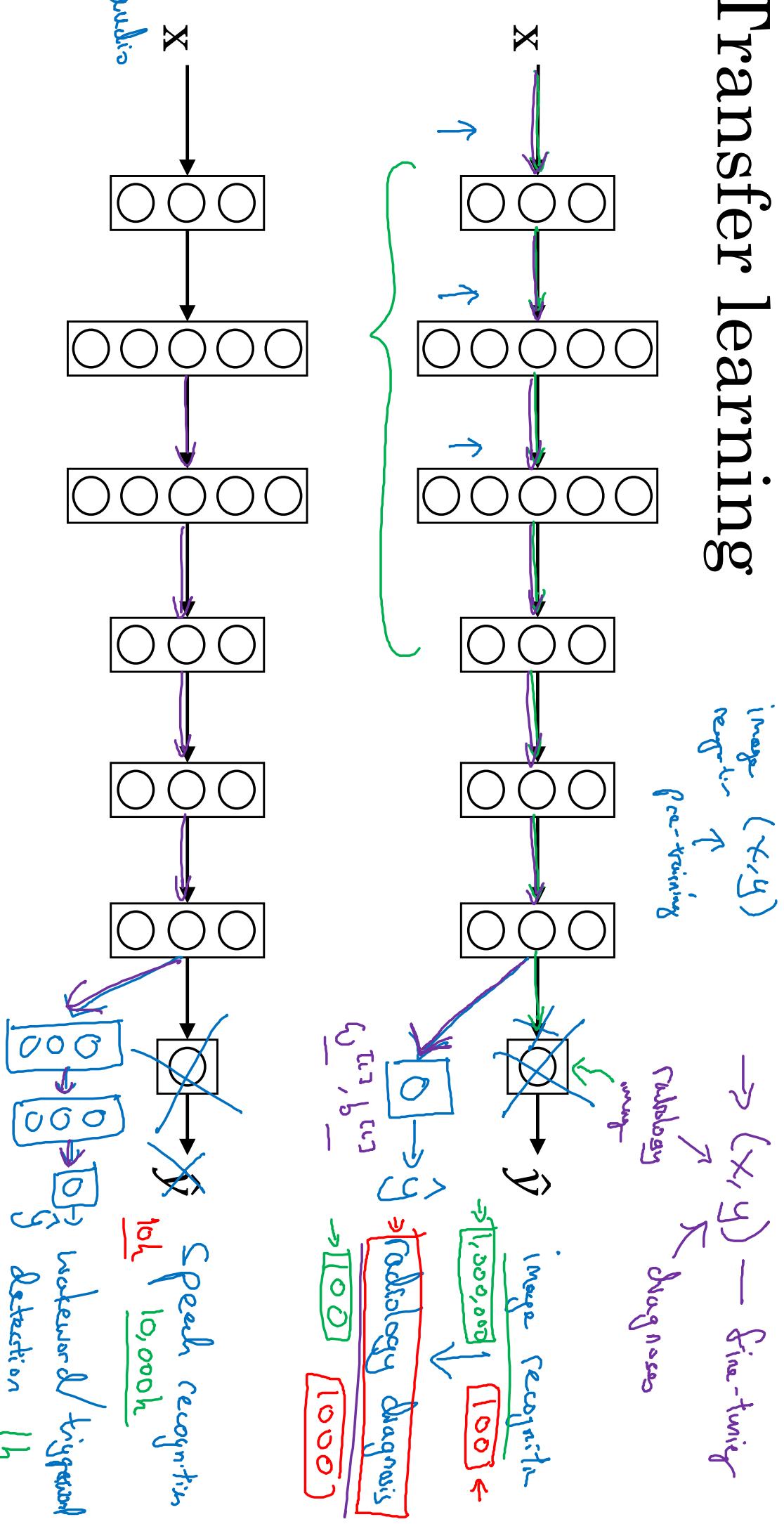


deeplearning.ai

Learning from
multiple tasks

Transfer learning

Transfer learning



50h

Andrew Ng

When transfer learning makes sense

Task from A \rightarrow 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.



Learning from
multiple tasks

Multi-task
learning

deeplearning.ai

Simplified autonomous driving example

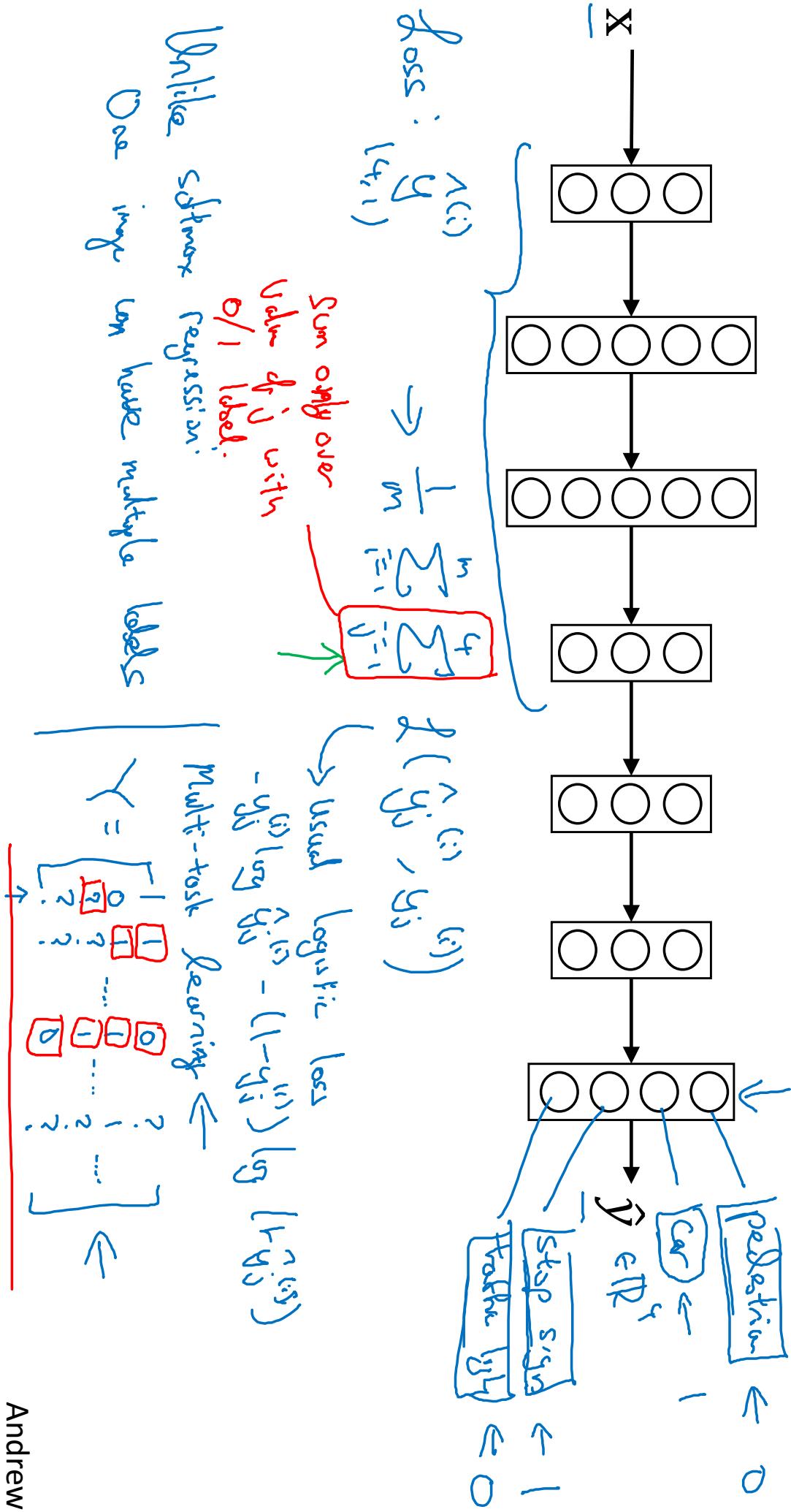


$$X^{(i)} = \begin{bmatrix} -y_{(1)}^{(i)} \\ -y_{(2)}^{(i)} \\ -y_{(3)}^{(i)}, \dots, -y_{(m)}^{(i)} \end{bmatrix}$$

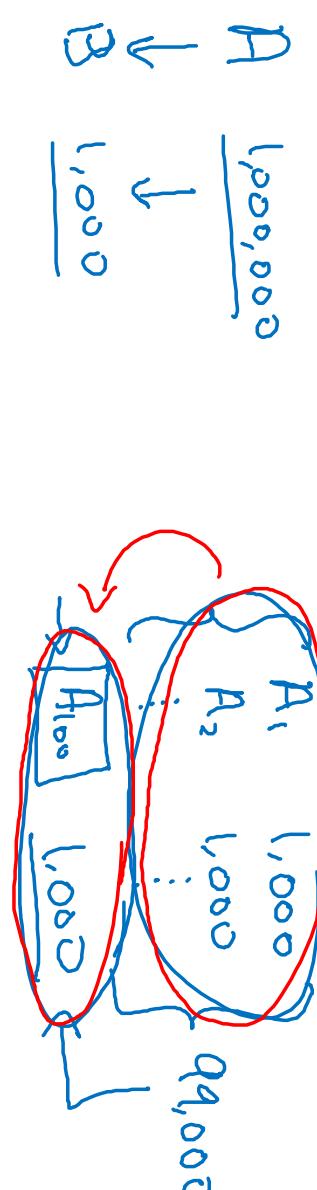
Pedestrians
Cars
Stop signs
Traffic lights
...

$$\cdots \quad 0 \quad - \quad - \quad 0 \quad y^{(i)}_{(4-i)}$$

Neural network architecture



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.

The diagram illustrates the distribution of data points across multiple tasks. At the top, there is a large blue arrow pointing down to a horizontal bar labeled 'A' with a value of 1,000,000. Below this, another blue arrow points down to a horizontal bar labeled 'B' with a value of 1,000. To the right of these bars, there are three ovals representing different tasks. The first oval contains a blue arrow pointing down to a horizontal bar labeled 'A₁' with a value of 1,000. The second oval contains a blue arrow pointing down to a horizontal bar labeled 'A₂' with a value of 1,000. The third oval contains a blue arrow pointing down to a horizontal bar labeled 'A₁₀₀' with a value of 1,000. A red bracket groups all three ovals together, indicating they represent the same type of data (lower-level features) across different tasks.
- Can train a big enough neural network to do well on all the tasks.



End-to-end deep
learning

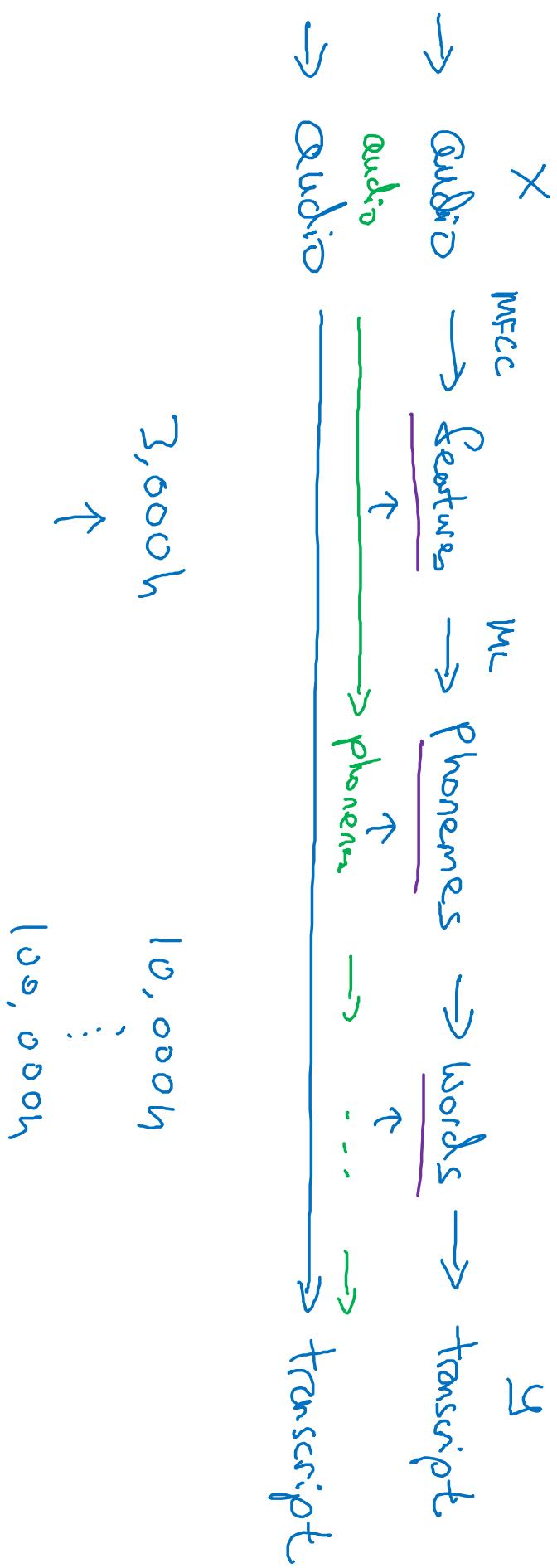
What is
end-to-end
deep learning

deeplearning.ai

What is end-to-end learning?

Speech recognition example

"cat"

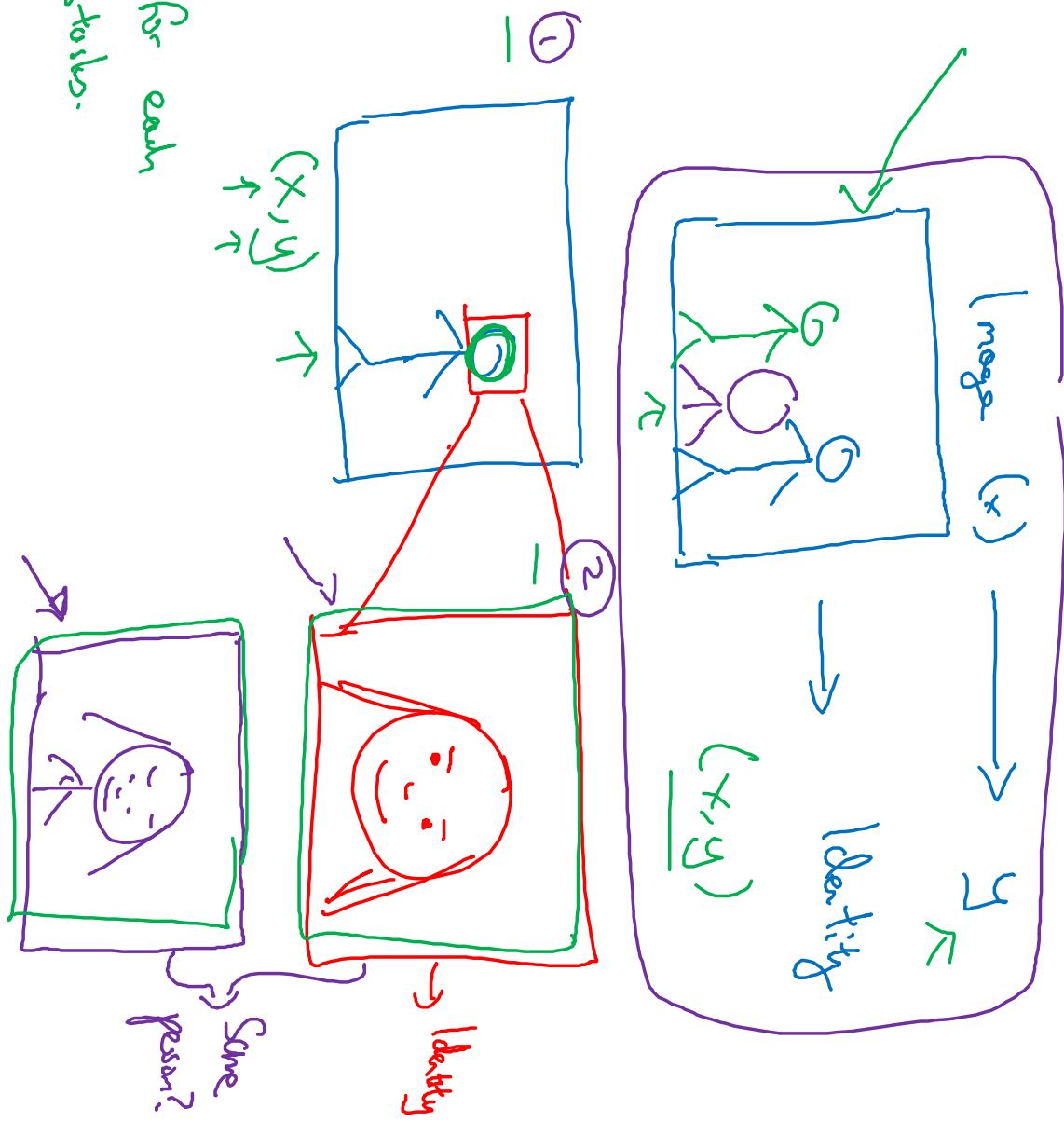


Face recognition



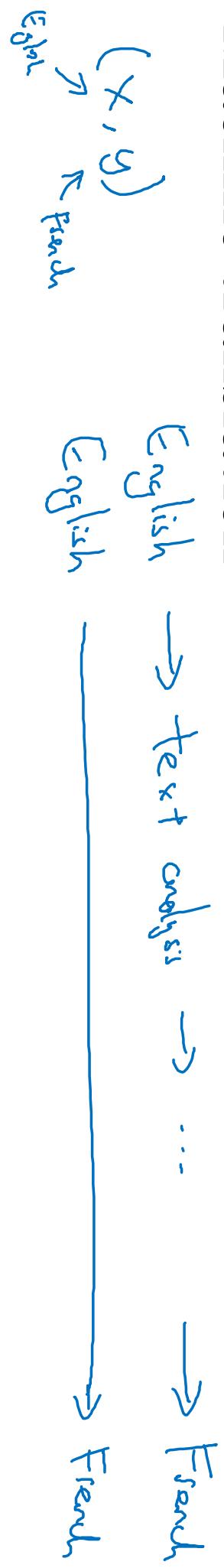
[Image courtesy of Baidu]

Have photo for each
of 2 subjects.



More examples

Machine translation



Estimating child's age:





End-to-end deep
learning

Whether to use
end-to-end learning

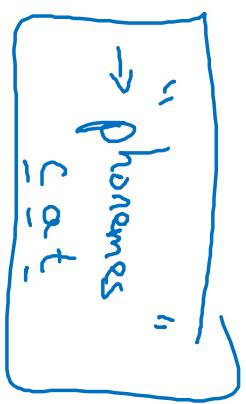
deeplearning.ai

Pros and cons of end-to-end deep learning

Pros:

- Let the data speak
- Less hand-designing of components needed

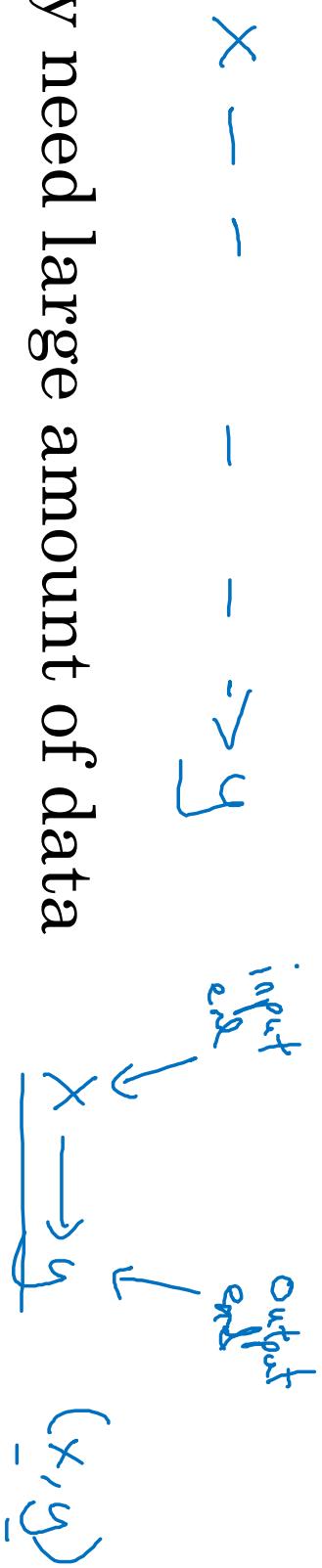
$$x \rightarrow y$$



Cons:

- May need large amount of data
- Excludes potentially useful hand-designed components

Data:
—
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 ?

