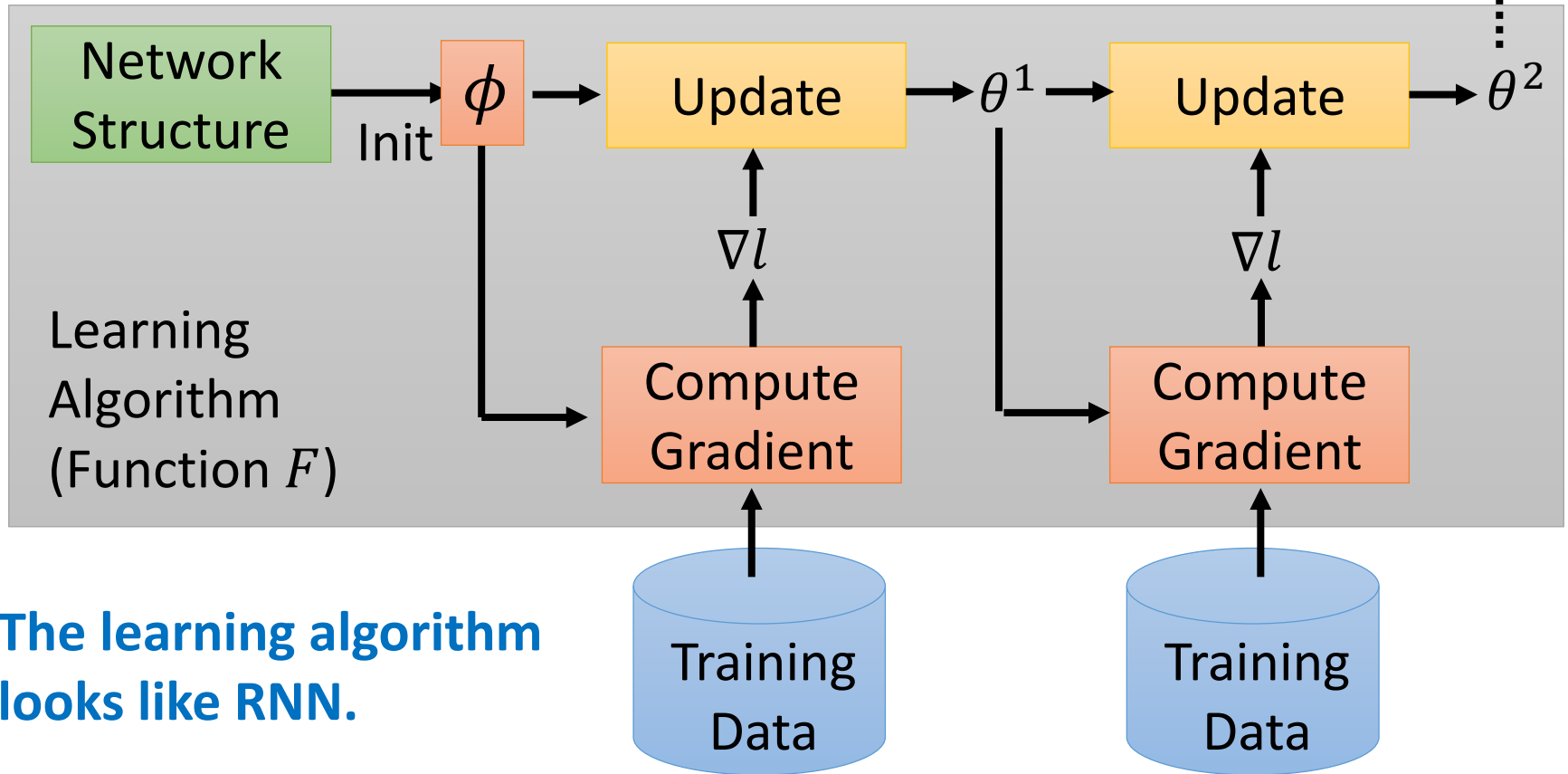




# Meta Learning (Part 2): **Gradient Descent as LSTM**

Hung-yi Lee

# Can we learn more than initialization parameters?



The learning algorithm looks like RNN.

OPTIMIZATION AS A MODEL FOR FEW-SHOT LEARNING

Sachin Ravi\* and Hugo Larochelle

Twitter, Cambridge, USA

{sachinr, hugo}@twitter.com

Learning to learn by gradient descent by gradient descent

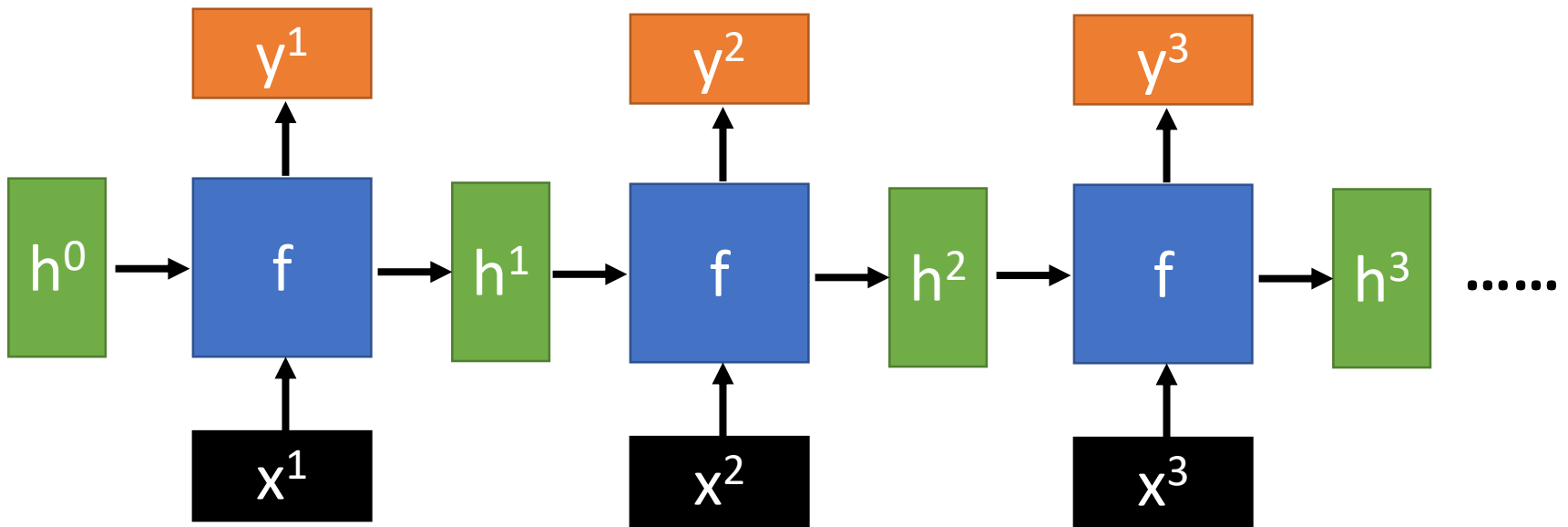
Marcin Andrychowicz<sup>1</sup>, Misha Denil<sup>1</sup>, Sergio Gómez Colmenarejo<sup>1</sup>, Matthew W. Hoffman<sup>1</sup>, David Pfau<sup>1</sup>, Tom Schaul<sup>1</sup>, Brendan Shillingford<sup>1,2</sup>, Nando de Freitas<sup>1,2,3</sup>

<sup>1</sup>Google DeepMind <sup>2</sup>University of Oxford <sup>3</sup>Canadian Institute for Advanced Research

# Recurrent Neural Network

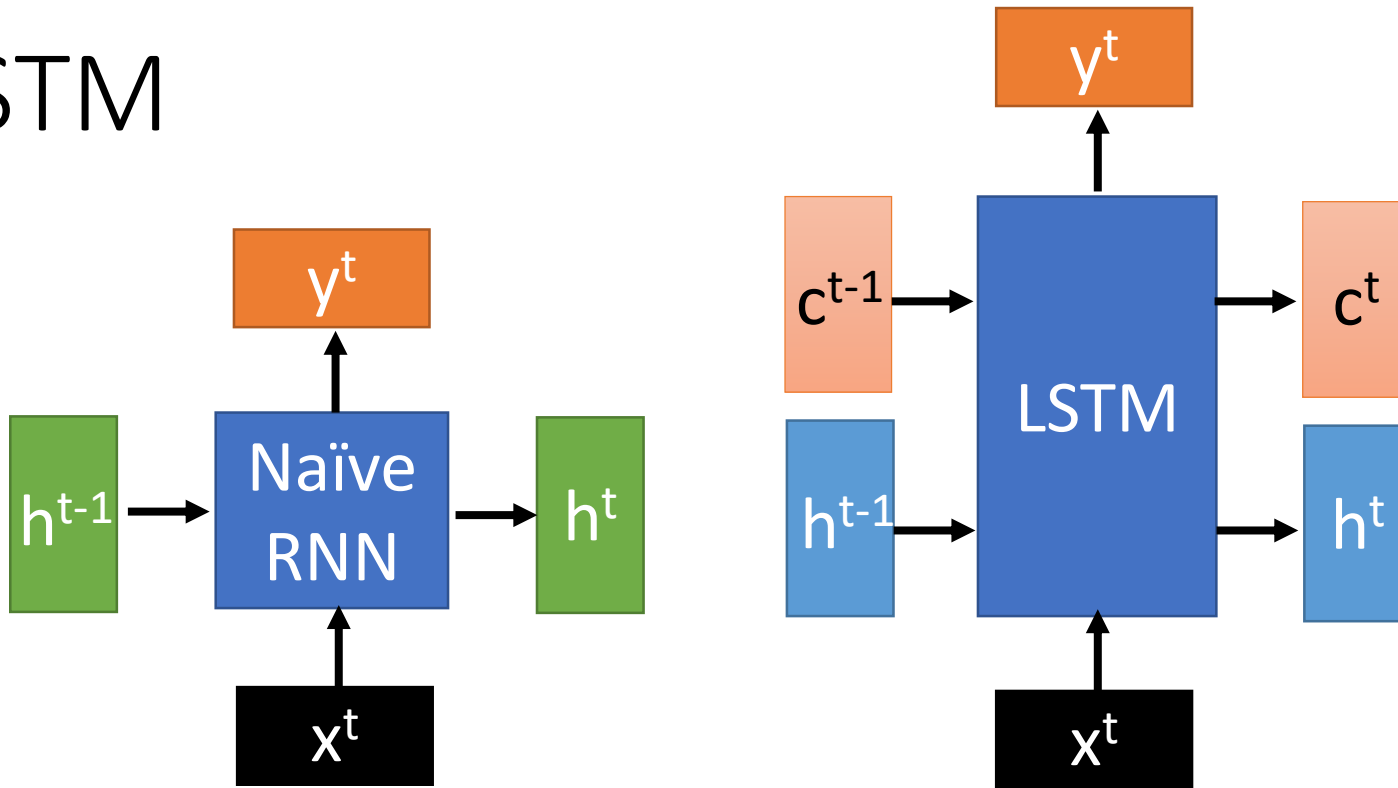
- Given function  $f: h', y = f(h, x)$

$h$  and  $h'$  are vectors with the same dimension



No matter how long the input/output sequence is, we only need one function  $f$

# LSTM

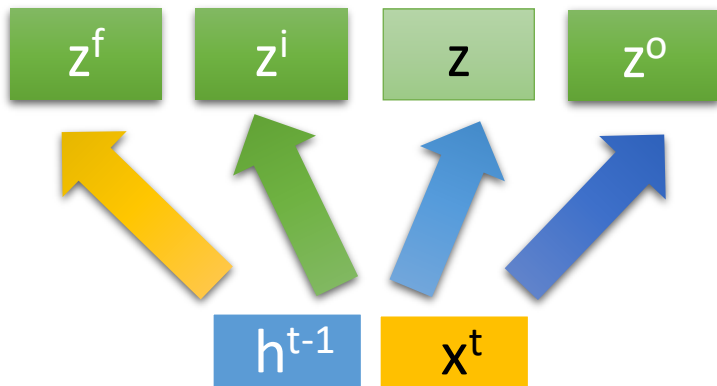


$c$  change slowly  $\Rightarrow c^t$  is  $c^{t-1}$  added by something

$h$  change faster  $\Rightarrow h^t$  and  $h^{t-1}$  can be very different

# Review: LSTM

$c^{t-1}$



$$z = \tanh\left( W \begin{array}{c} x^t \\ h^{t-1} \end{array} \right)$$

$$z^i = \sigma\left( W^i \begin{array}{c} x^t \\ h^{t-1} \end{array} \right)$$

input

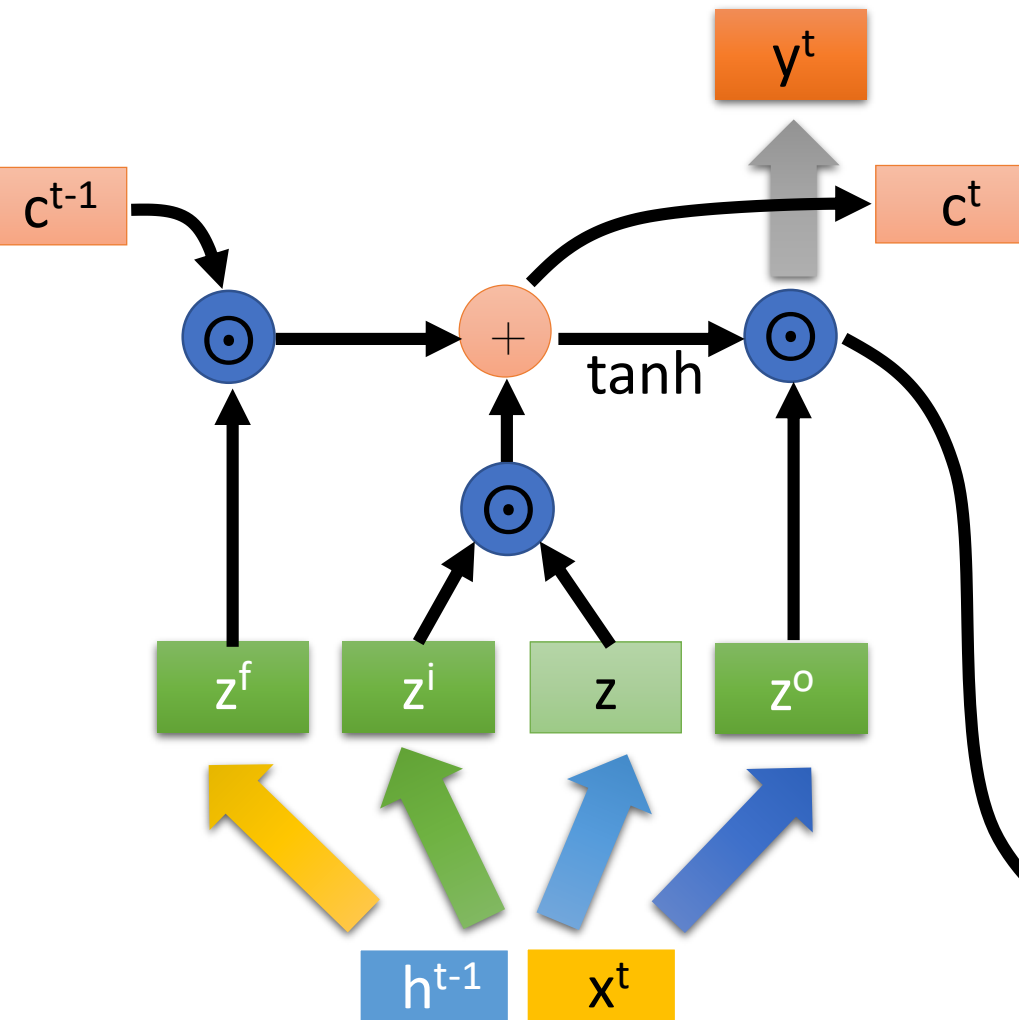
$$z^f = \sigma\left( W^f \begin{array}{c} x^t \\ h^{t-1} \end{array} \right)$$

forget

$$z^o = \sigma\left( W^o \begin{array}{c} x^t \\ h^{t-1} \end{array} \right)$$

output

# Review: LSTM

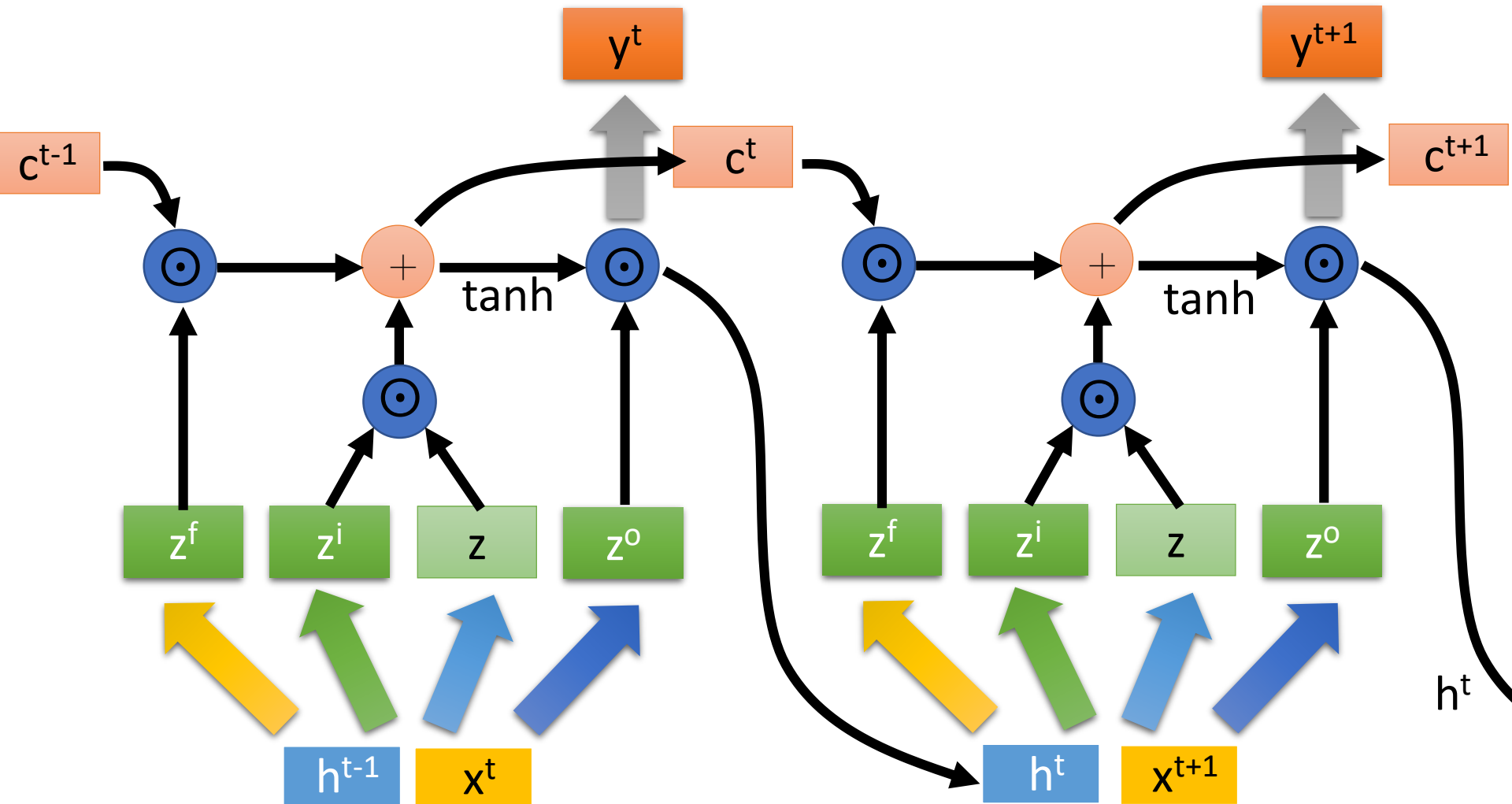


$$c^t = z^f \odot c^{t-1} + z^i \odot z$$

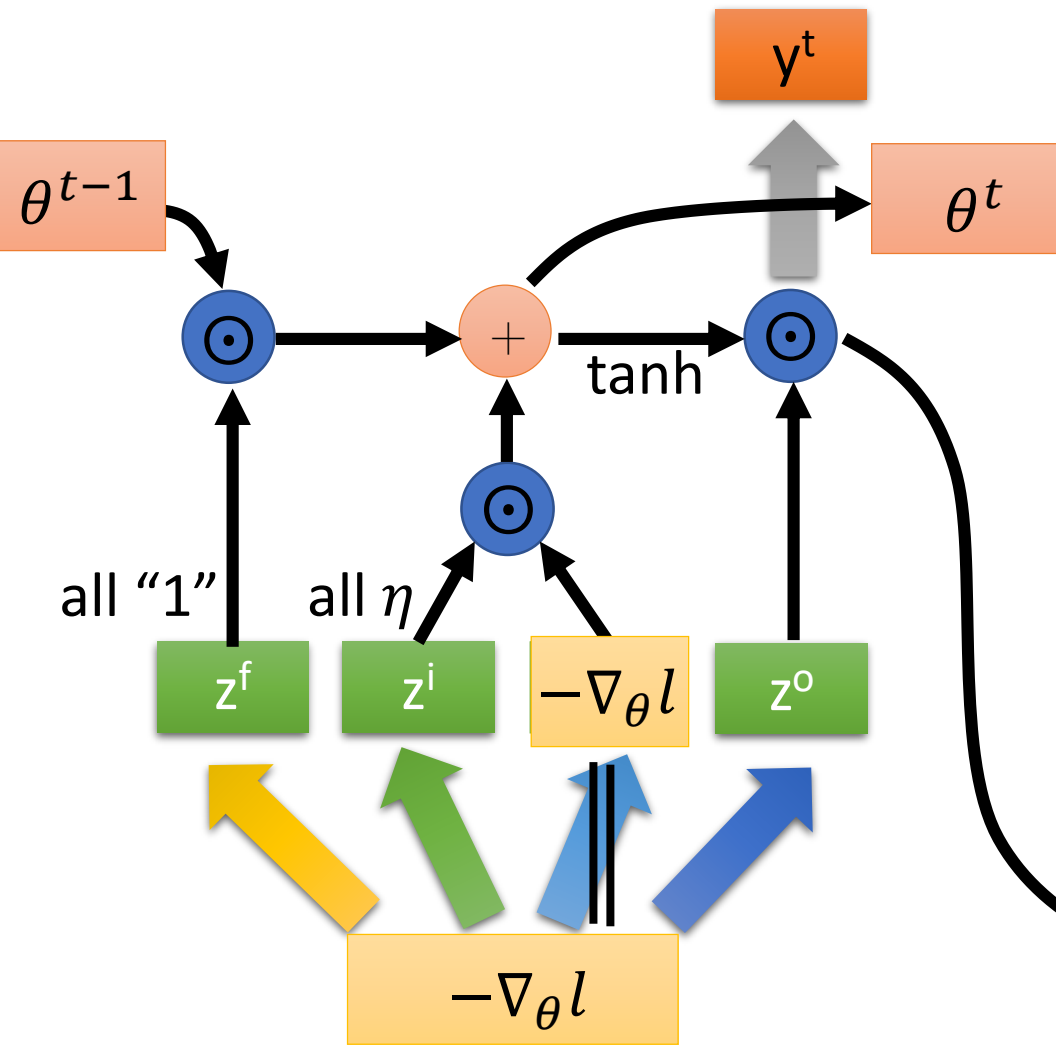
$$h^t = z^o \odot \tanh(c^t)$$

$$y^t = \sigma(W' h^t)$$

# Review: LSTM



# Similar to gradient descent based algorithm



$$\theta^t = \theta^{t-1} - \eta \nabla_{\theta} l$$

$$\theta^t = z^f \odot \begin{bmatrix} 1 \\ 1 \\ \vdots \end{bmatrix} \odot \theta^{t-1} + z^i \odot \begin{bmatrix} \eta \\ \eta \\ \vdots \end{bmatrix} \odot -\nabla_{\theta} l$$

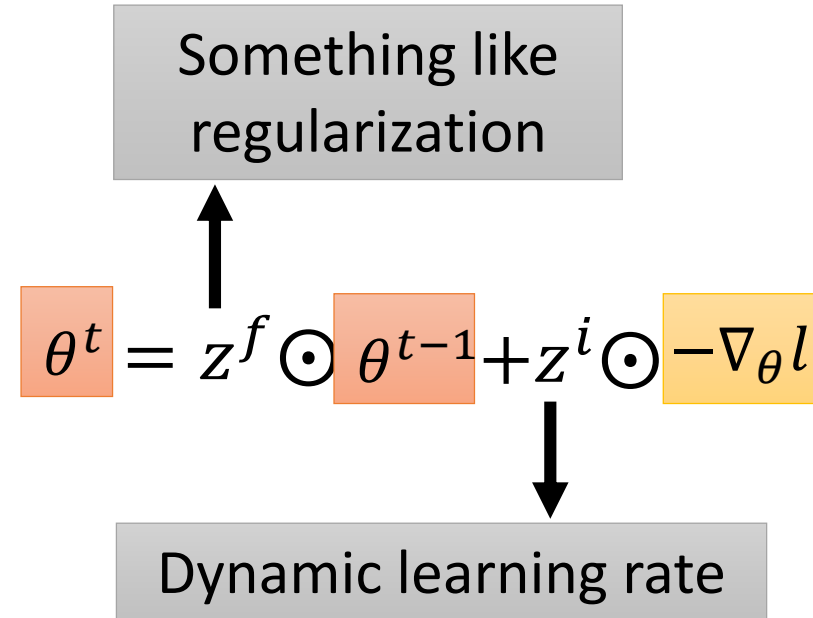
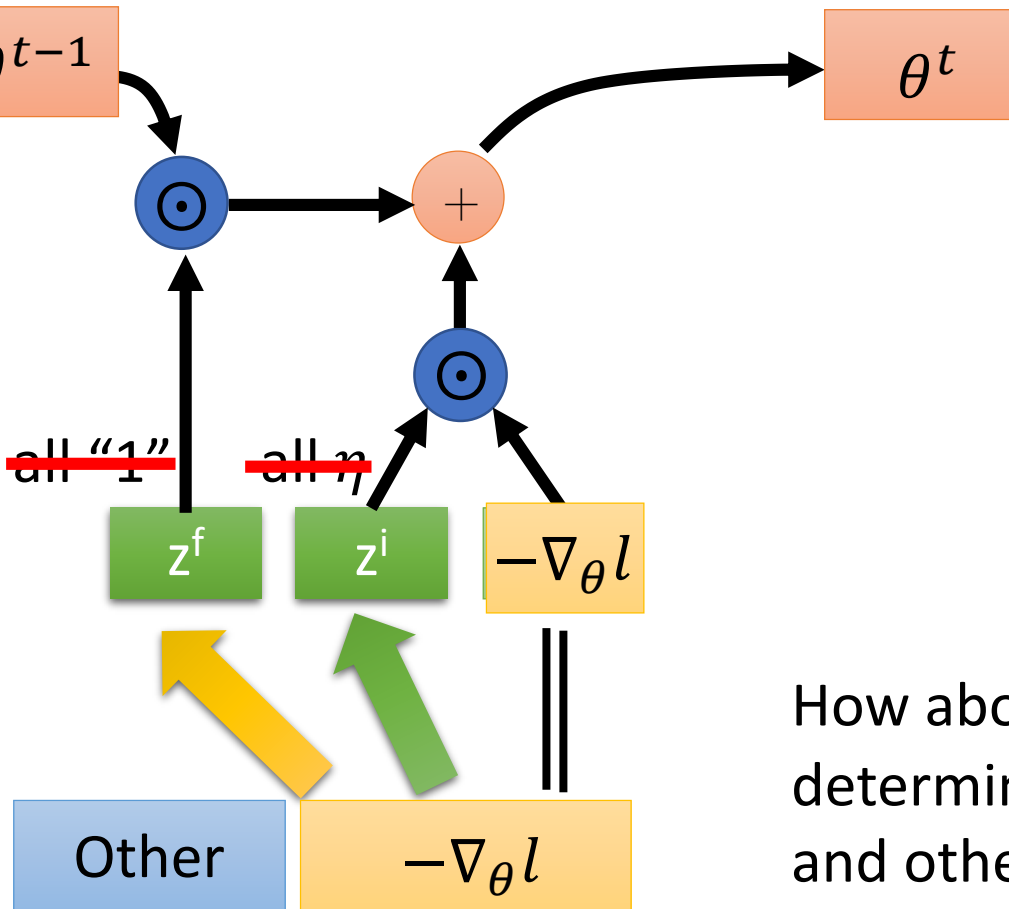
$$h^t = z^o \odot \tanh(c^t)$$

$$y^t = \sigma(W' h^t)$$



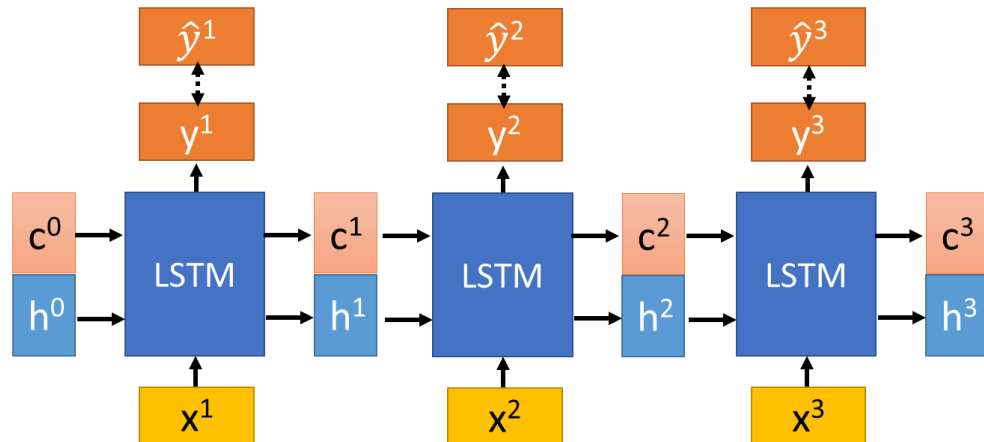
# Similar to gradient descent based algorithm

$$\theta^t = \theta^{t-1} - \eta \nabla_{\theta} l$$



How about machine learn to determine  $z^f$  and  $z^i$  from  $-\nabla_{\theta} l$  and other information?

# Typical LSTM



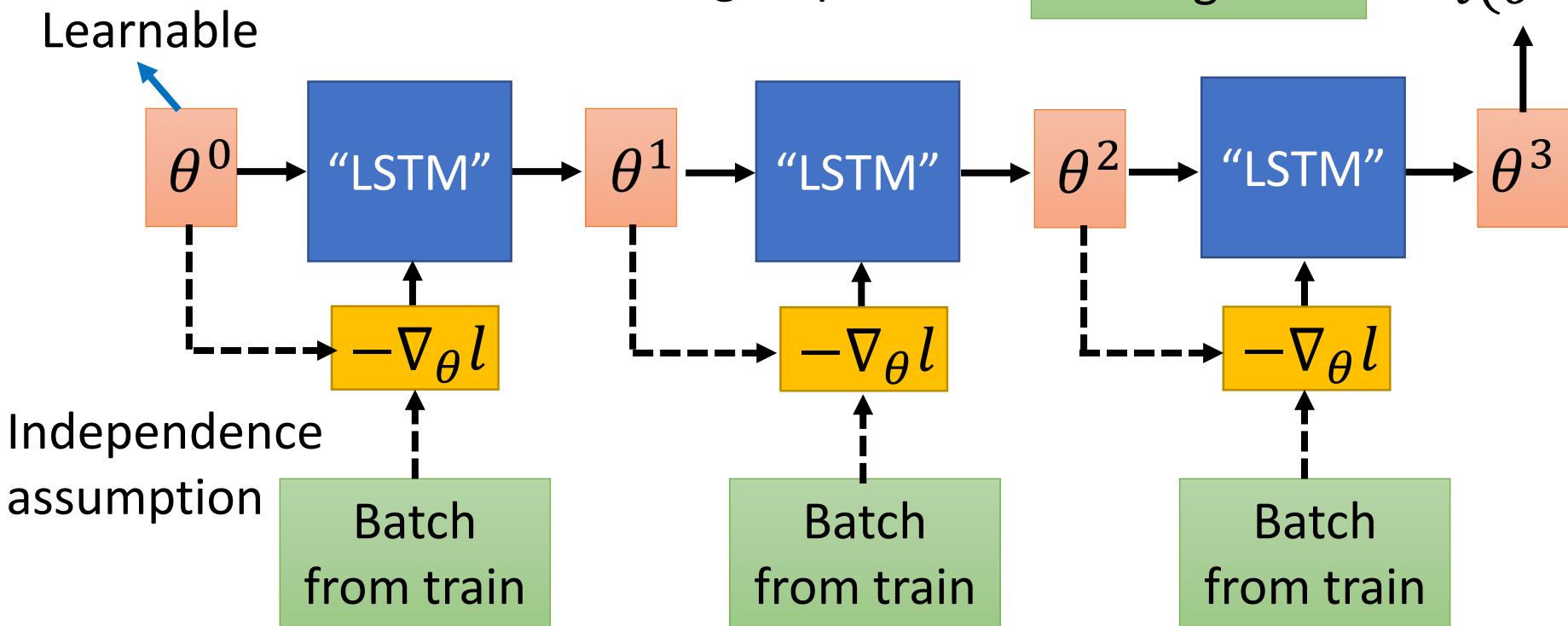
Learn to minimize

Testing Data  $\rightarrow l(\theta^3)$

3 training steps

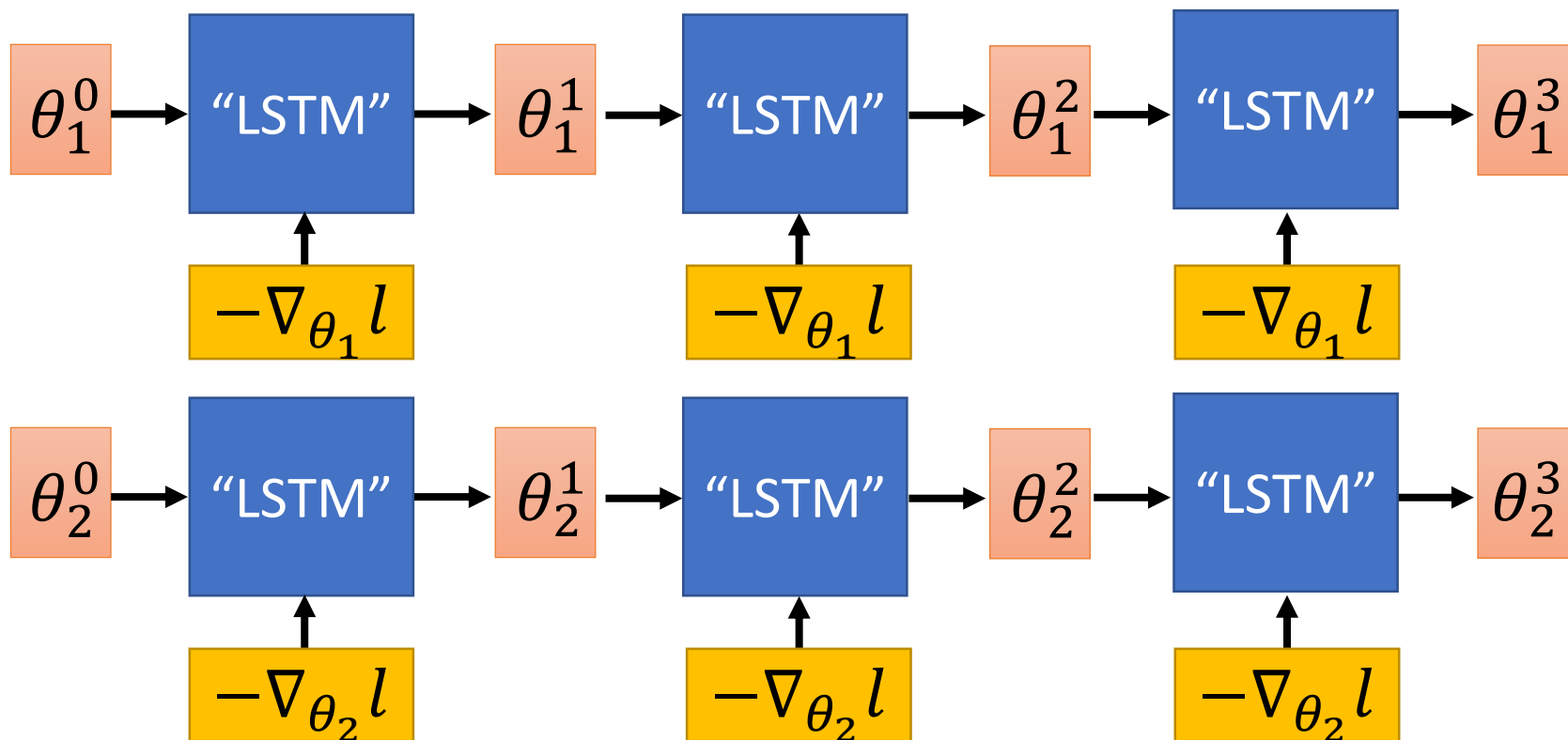
# LSTM for Gradient Descent

$$\theta^t = z^f \odot \theta^{t-1} + z^i \odot -\nabla_{\theta} l$$



# Real Implementation

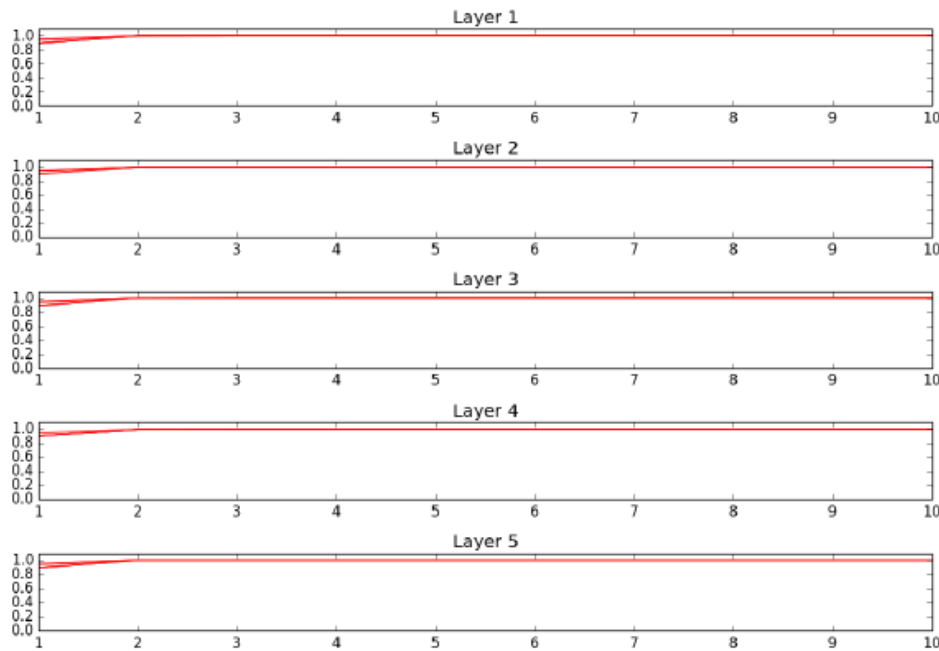
The LSTM used only has one cell. Share across all parameters



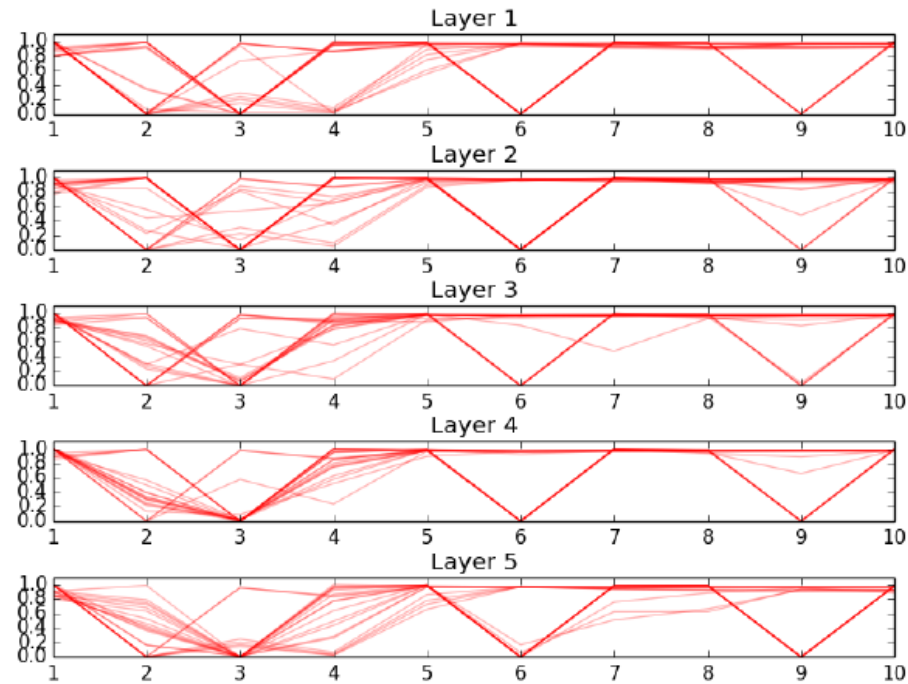
- Reasonable model size
- In typical gradient descent, all the parameters use the same update rule
- Training and testing model architectures can be different.

# Experimental Results

$$\theta^t = z^f \odot \theta^{t-1} + z^i \odot -\nabla_{\theta} l$$



(a) Forget gate values for 1-shot meta-learner



(b) Input gate values for 1-shot meta-learner

Parameter update depends on not only current gradient, but *previous gradients*.

RMSProp

$$w^1 \leftarrow w^0 - \frac{\eta}{\sigma^0} g^0 \quad \sigma^0$$

$$w^2 \leftarrow w^1 - \frac{\eta}{\sigma^1} g^1 \quad \sigma^1$$

$$w^3 \leftarrow w^2 - \frac{\eta}{\sigma^2} g^2 \quad \sigma^2$$

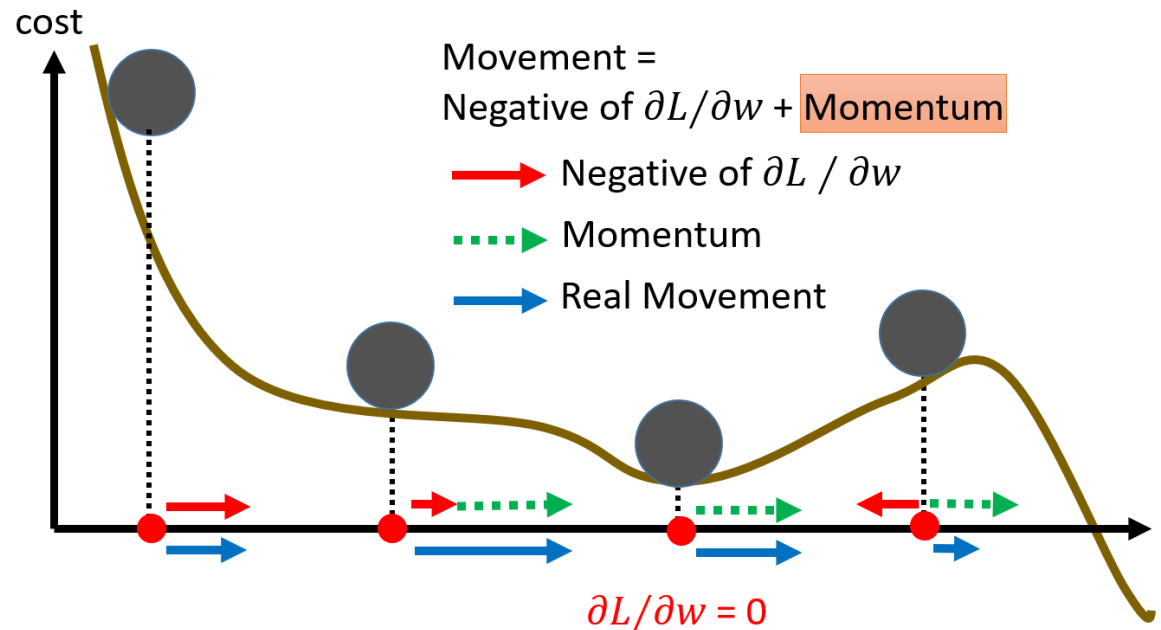
⋮

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sigma^t} g^t \quad \sigma^t$$

Re  
with

Momentum

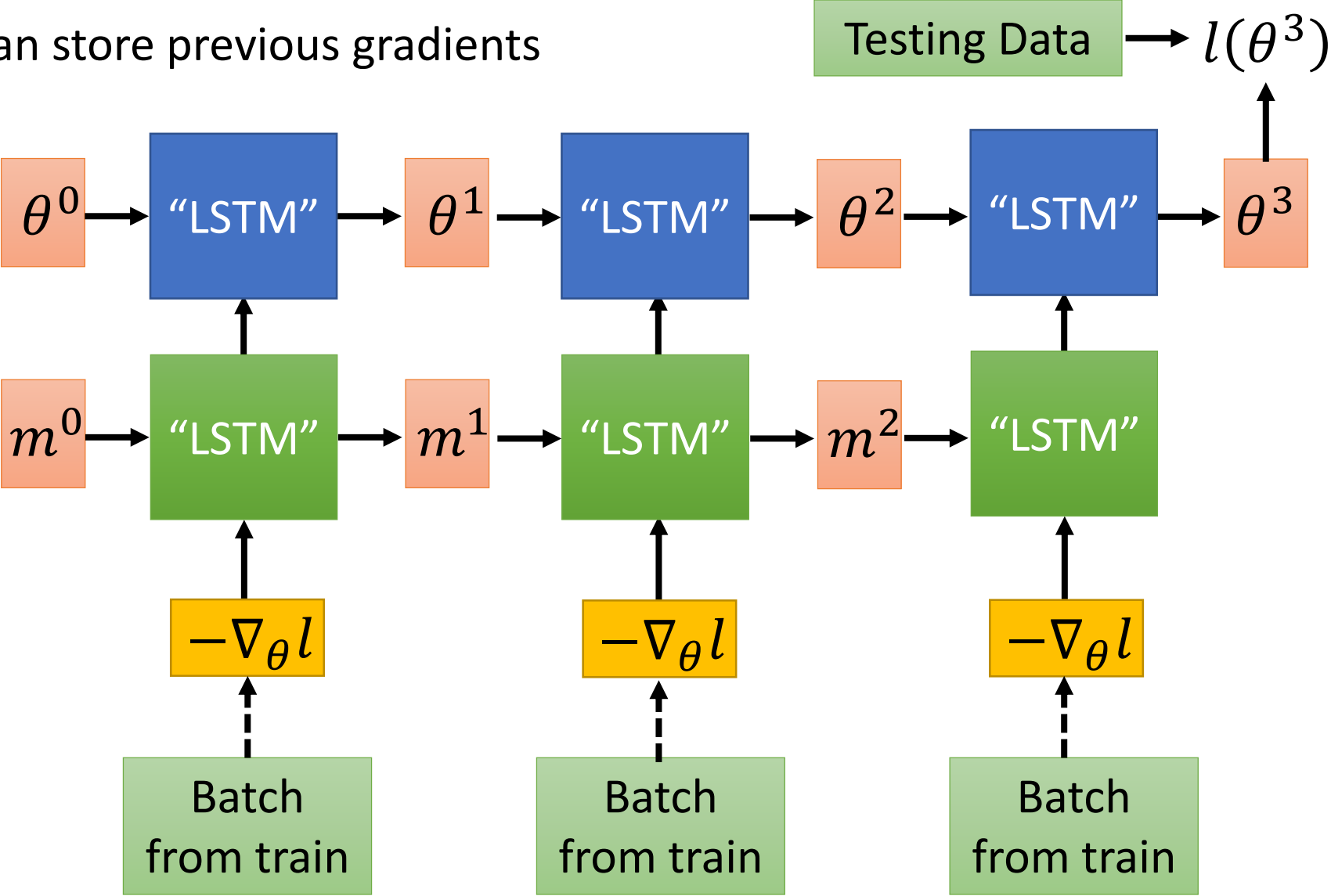
Still not guarantee reaching global minima, but give some hope .....



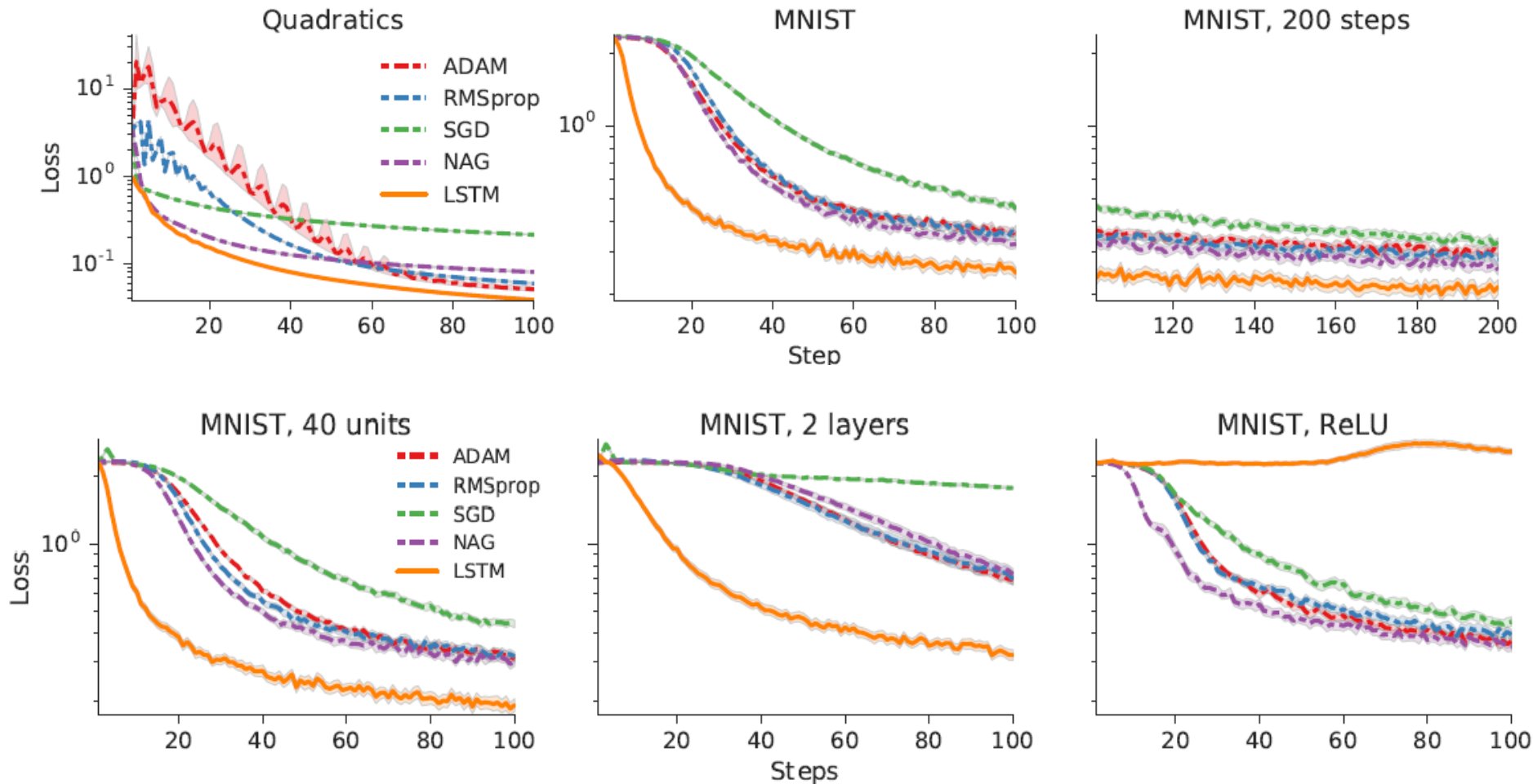
# LSTM for Gradient Descent (v2)

3 training steps

$m$  can store previous gradients



# Experimental Results



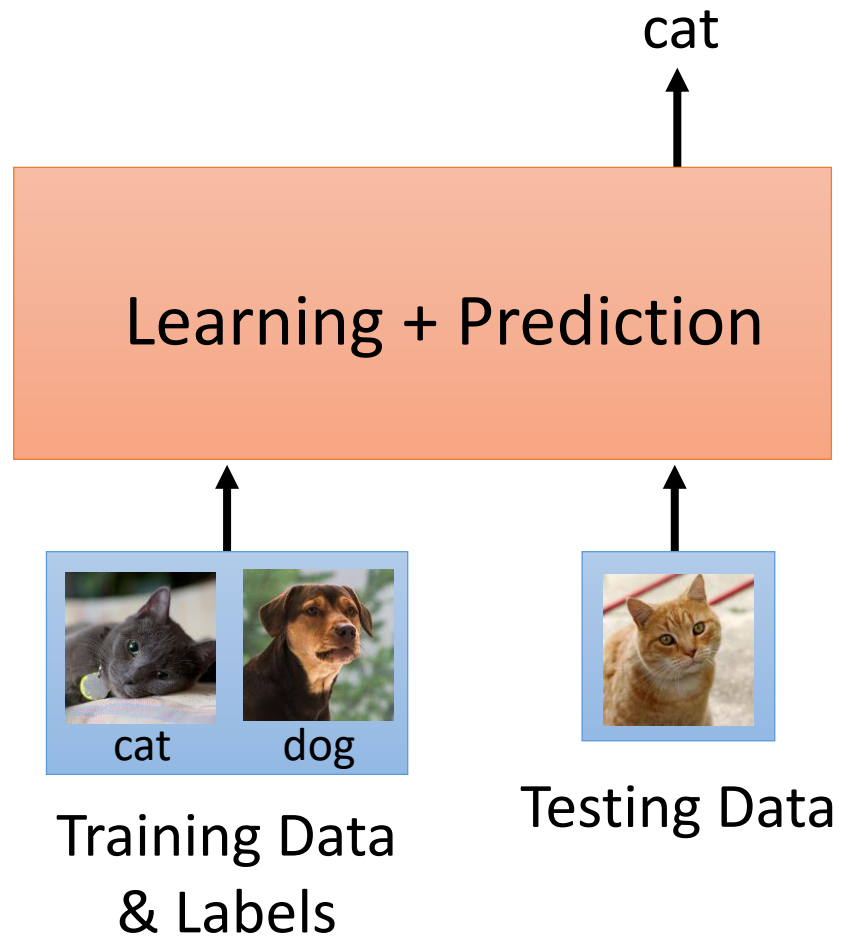
# Meta Learning (Part 3)

Hung-yi Lee



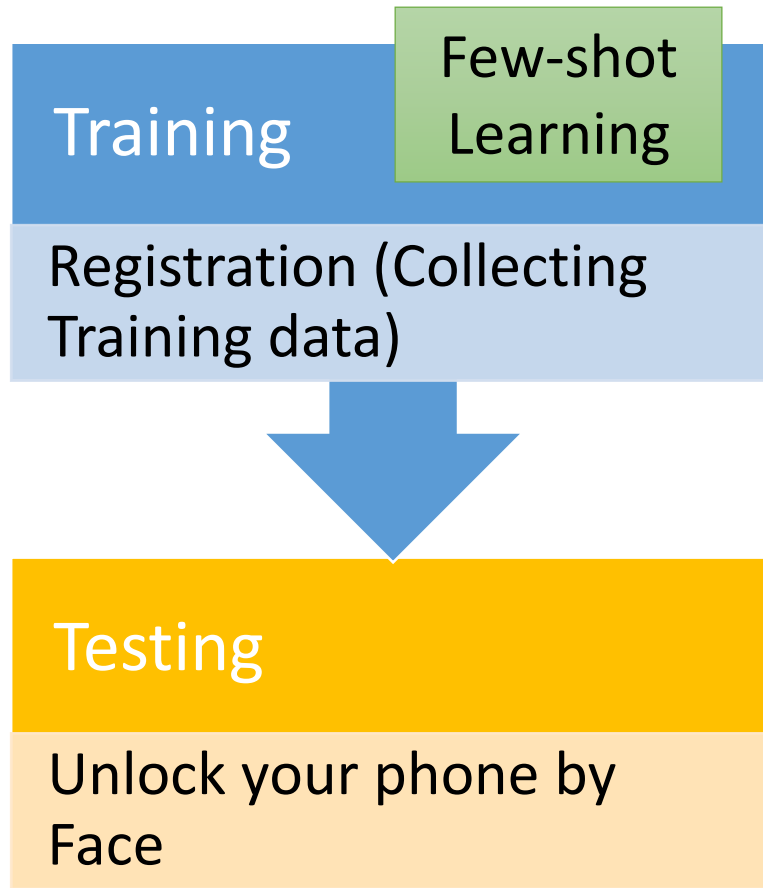
# Even more crazy idea ...

- Input:
  - Training data and their labels
  - Testing data
- Output:
  - Predicted label of testing data





# Face Verification


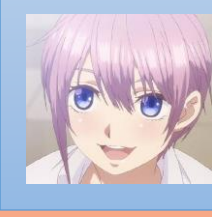
In each task:





# Meta Learning



Same approach for Speaker Verification

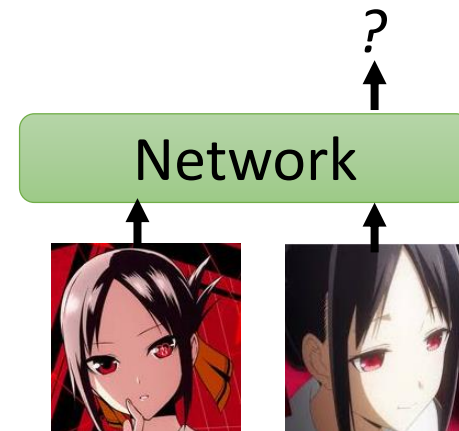
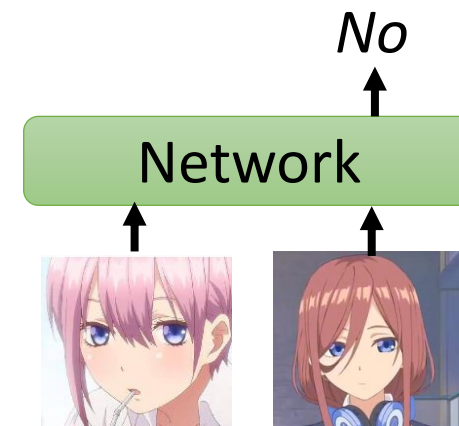
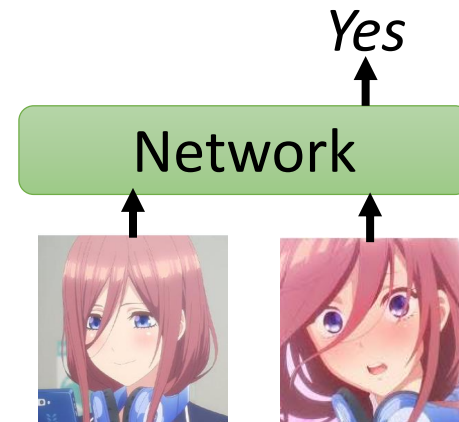
*Train*  *Test*  *Yes*

*Train*  *Test*  *No*

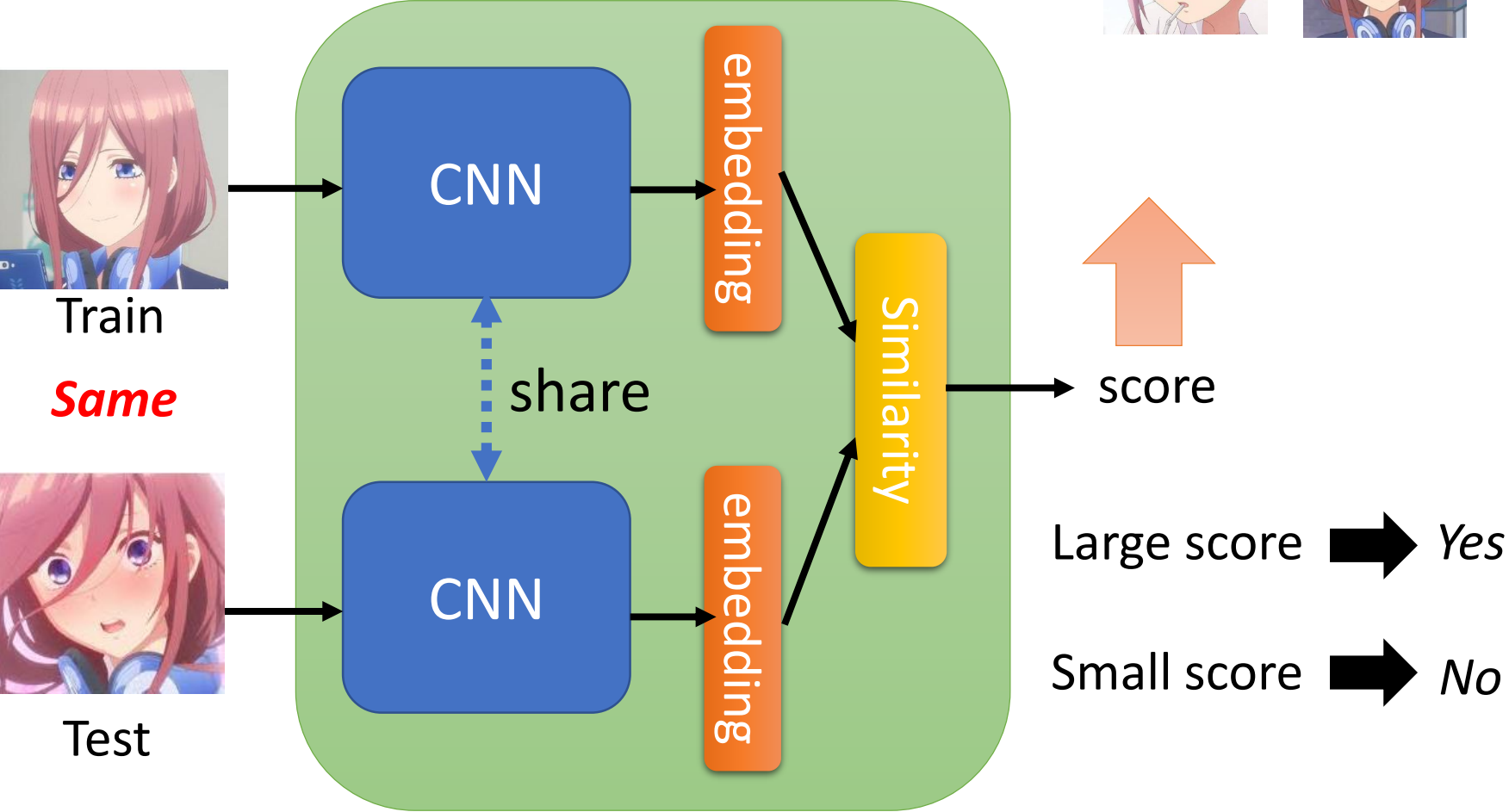
*Train*  *Test*  *No*

## Testing Tasks

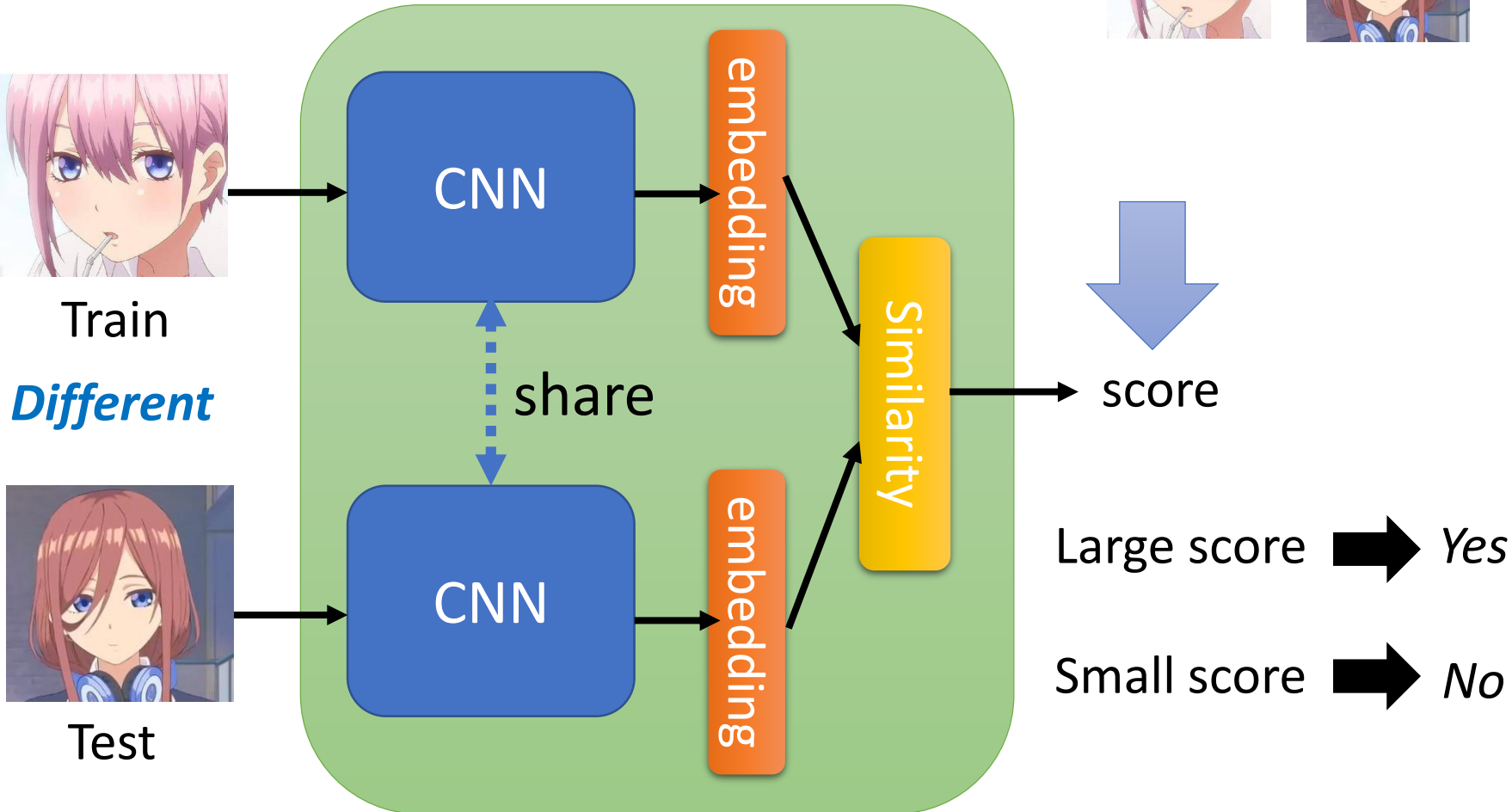
*Train*  *Test*  *Yes or No*



# Siamese Network



# Siamese Network



Train  
*Different*

Test

CNN

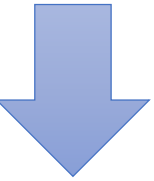
embedding

share

CNN

embedding

Similarity



score

Large score → Yes

Small score → No

Network



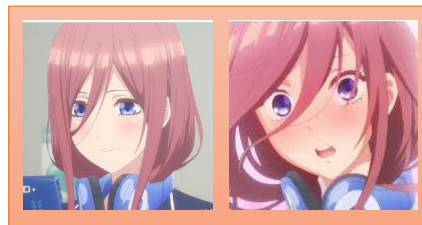
No

# Siamese Network

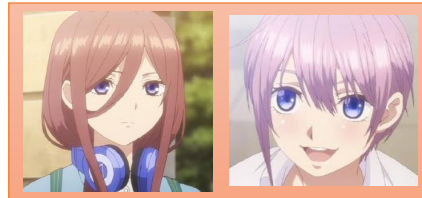
## - Intuitive Explanation

- Binary classification problem: “Are they the same?”

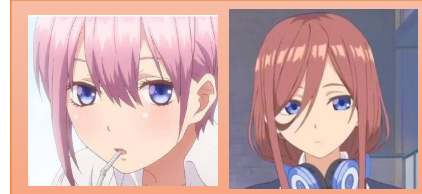
Training  
Set



same

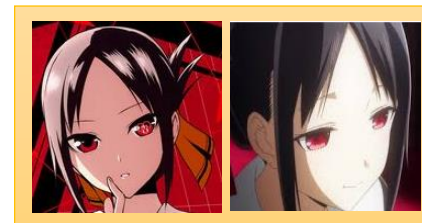


different



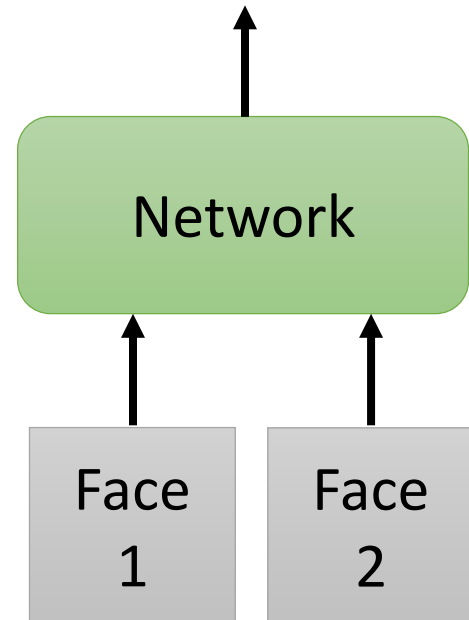
different

Testing  
Set



?

Same or Different

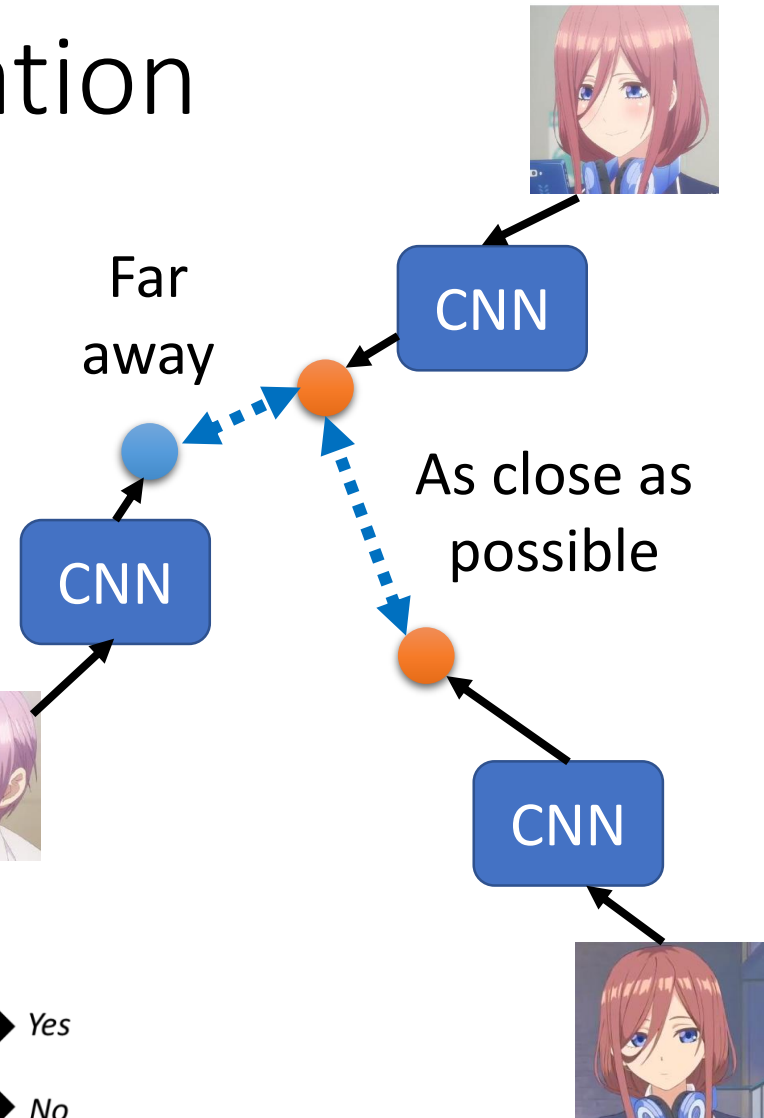
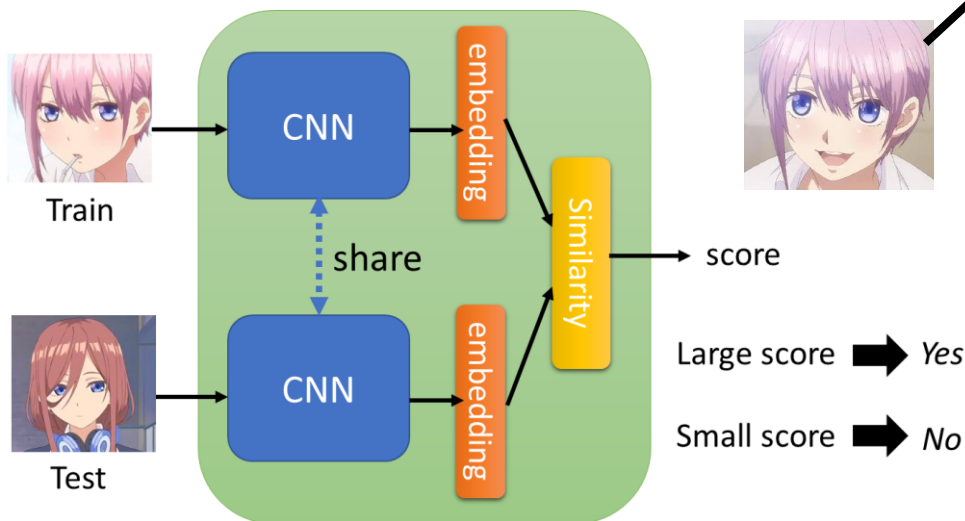


# Siamese Network

## - Intuitive Explanation

Learning embedding for faces

e.g. learn to ignore the background



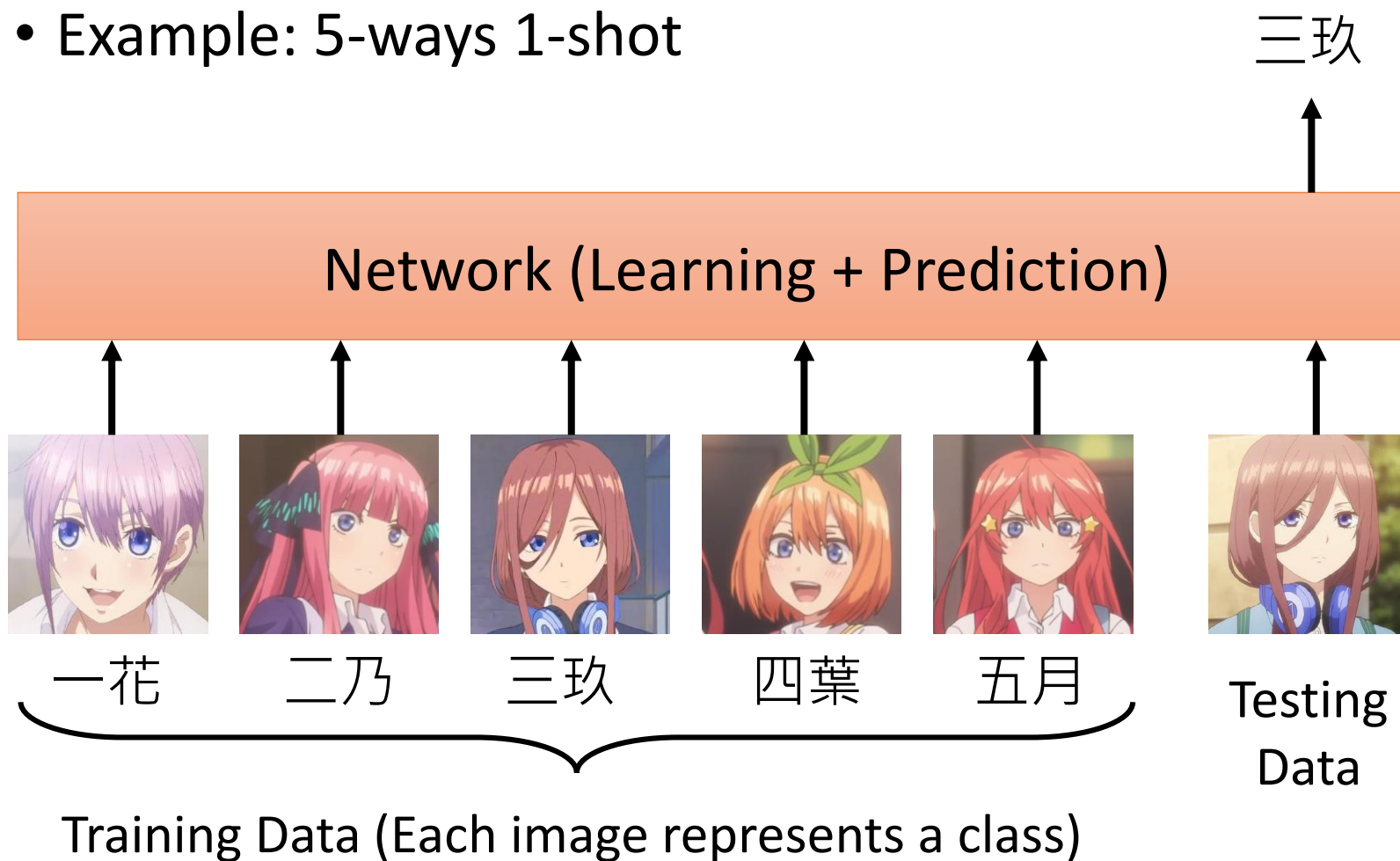
# To learn more ...

- What kind of distance should we use?
  - SphereFace: Deep Hypersphere Embedding for Face Recognition
  - Additive Margin Softmax for Face Verification
  - ArcFace: Additive Angular Margin Loss for Deep Face Recognition
- Triplet loss
  - Deep Metric Learning using Triplet Network
  - FaceNet: A Unified Embedding for Face Recognition and Clustering

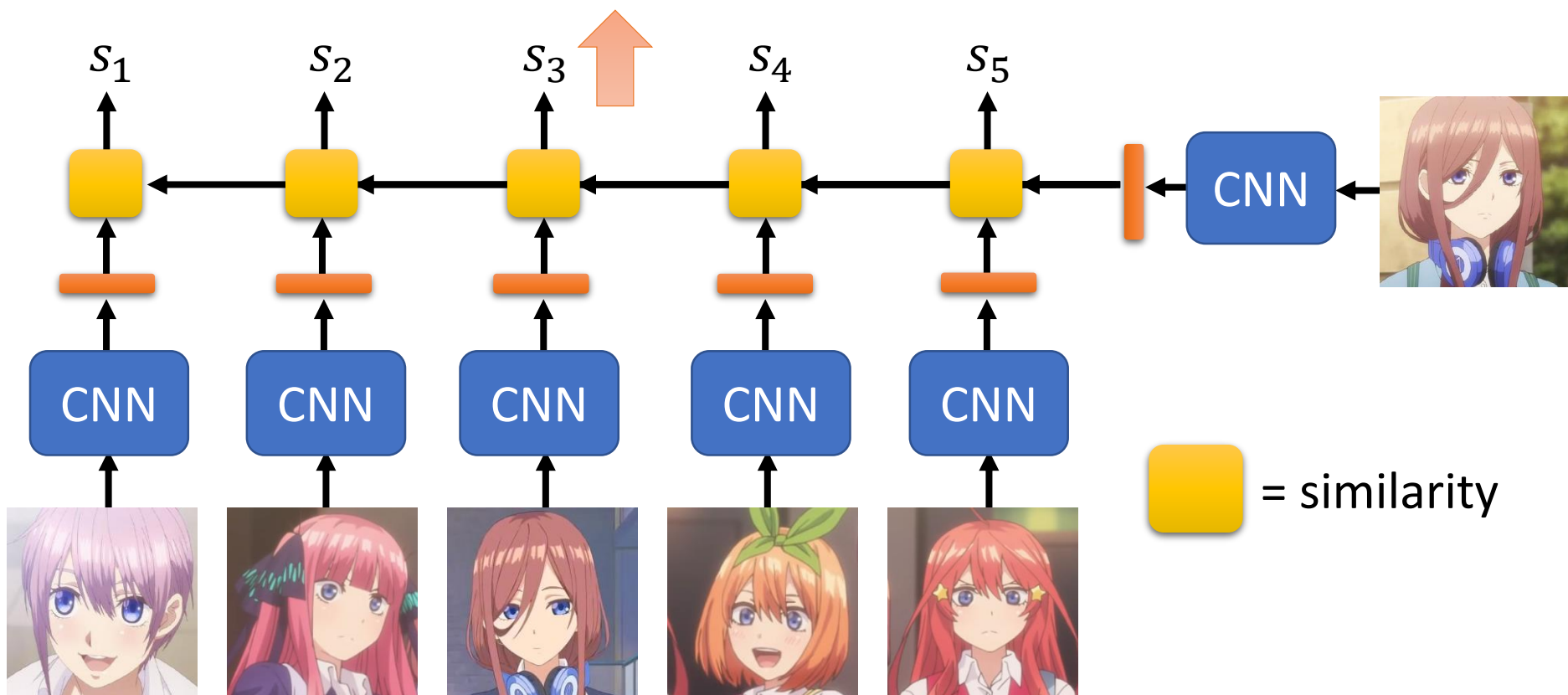
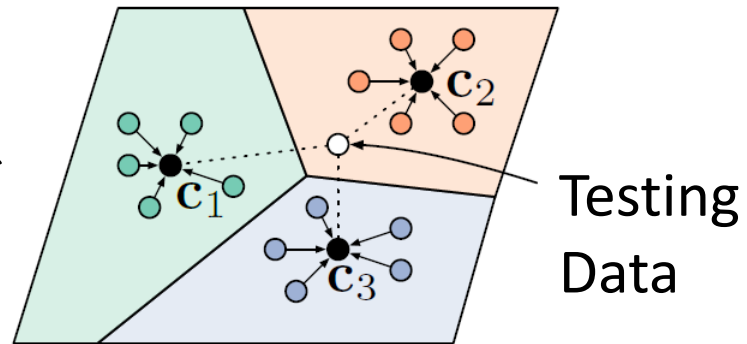


# N-way Few/One-shot Learning

- Example: 5-ways 1-shot

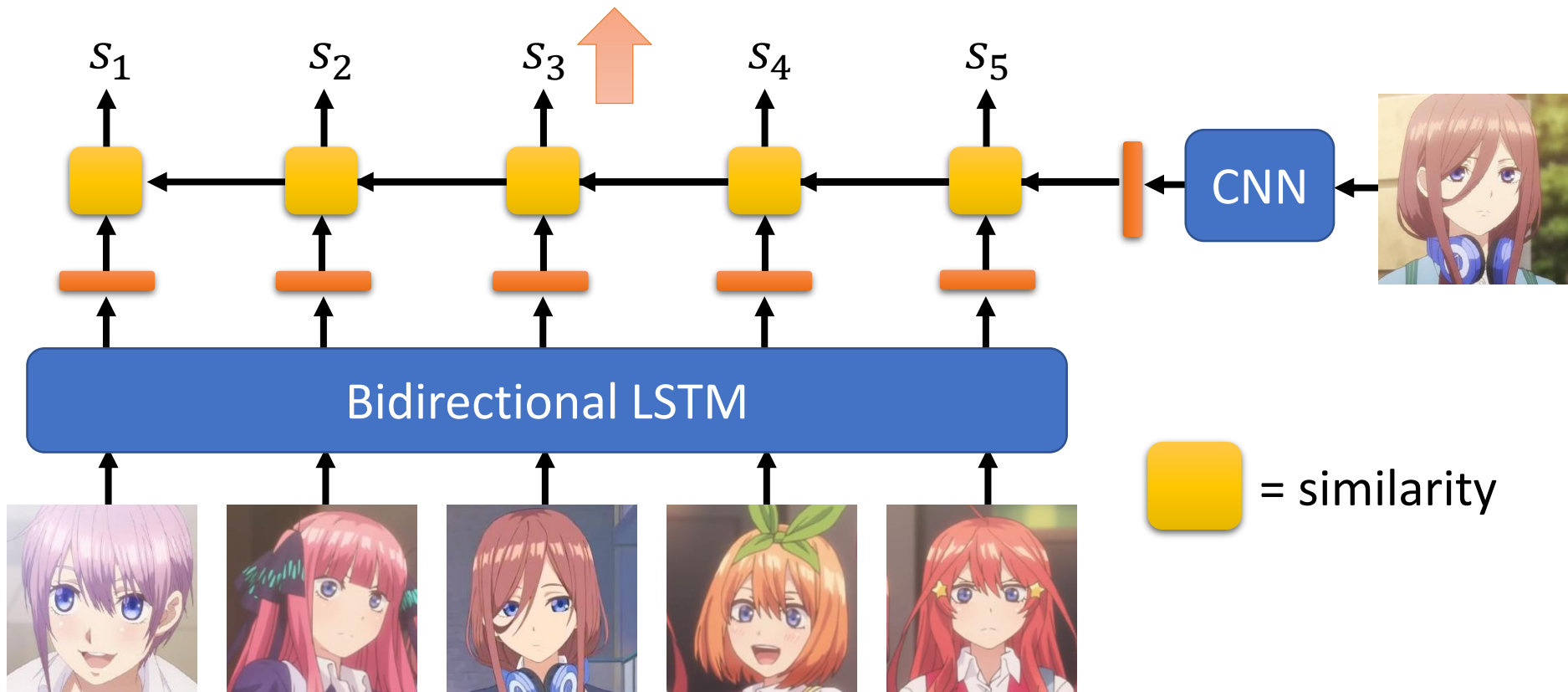


# Prototypical Network

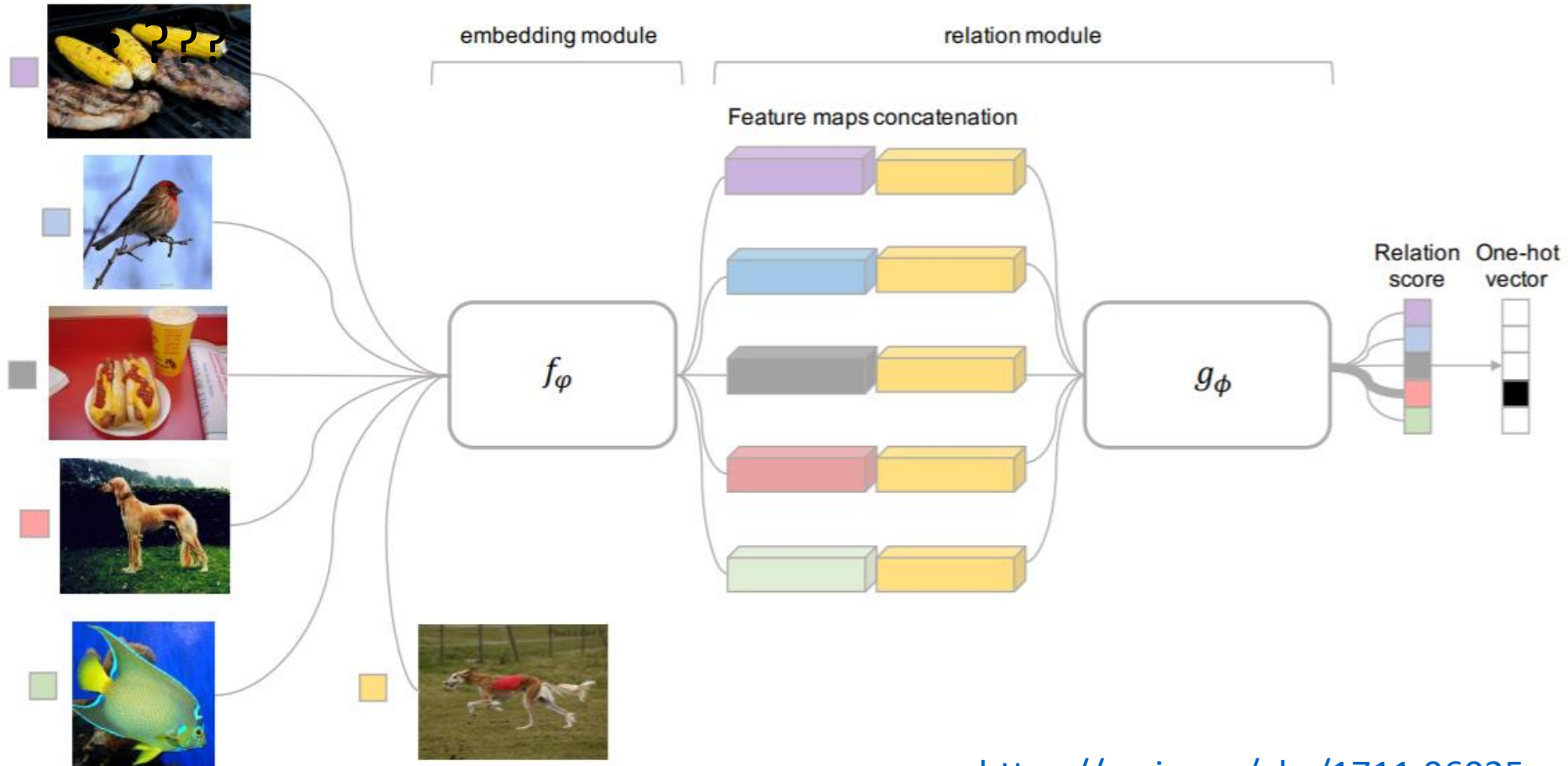


# Matching Network

Considering the relationship among the training examples



# Relation Network

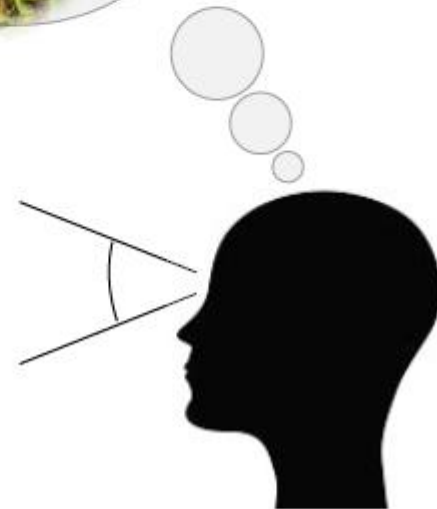


<https://arxiv.org/abs/1711.06025>

# Few-shot learning for imaginary data



blue heron



<https://arxiv.org/abs/1801.05401>

# Few-shot learning for imaginary data

