

Matching Networks for One Shot Learning

References :

Vinyals O, Blundell C, Lillicrap T, et al. Matching Networks for One Shot Learning[J]. 2016.

Download: <https://arxiv.org/abs/1606.04080>

段云志

©超星思维Nova Mind

2017/6/25 1

Outline

- Background
- Introduction
 - Introduce to One Shot Learning
 - Introduce to MANN
- Model
 - Motivation
 - Matching Networks
 - Backpropagation
- Experiments
- Summary

One-shot learning with Matching Networks

- Background
- Introduction
 - Introduce to One Shot Learning
 - Introduce to MANN
- Model
 - Motivation
 - Matching Networks
 - Backpropagation
- Experiments
- Summary

Background

- Learning from a few examples: remains challenge
- new data: models must relearn parameters
- Memory-augmented neural network has the ability to make accurate predictions after only a few samples

One-shot learning with Matching Networks

- Background
- Introduction
 - Introduce to One Shot Learning
 - Introduce to MANN
- Model
 - Motivation
 - Matching Networks
 - Backpropagation
- Experiments
- Summary

Introduction

- Introduce to One Shot Learning
 - What is one-shot learning?
 - Why do we need one-shot learning?
- Introduce to Neural Turing Machine(NTM)

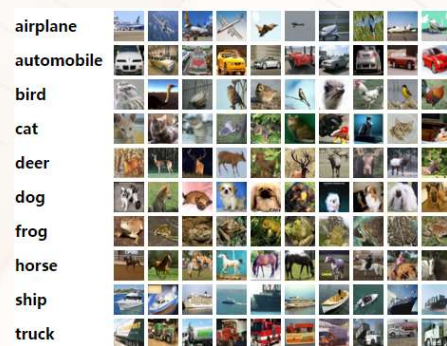
What is one-shot learning?

■ One-shot?

L :Data set

S:Support set

B:Batch



L

Sampling N labels

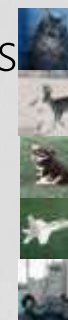
N-ways

sampling 1 example

| |
|----------|
| cat |
| deer |
| dog |
| airplane |
| truck |

Sampling k examples from each label

K-shot



S

Task: classify \hat{x} into 5 classes, {cat,deer,dog,airplane,truck}, using support set S



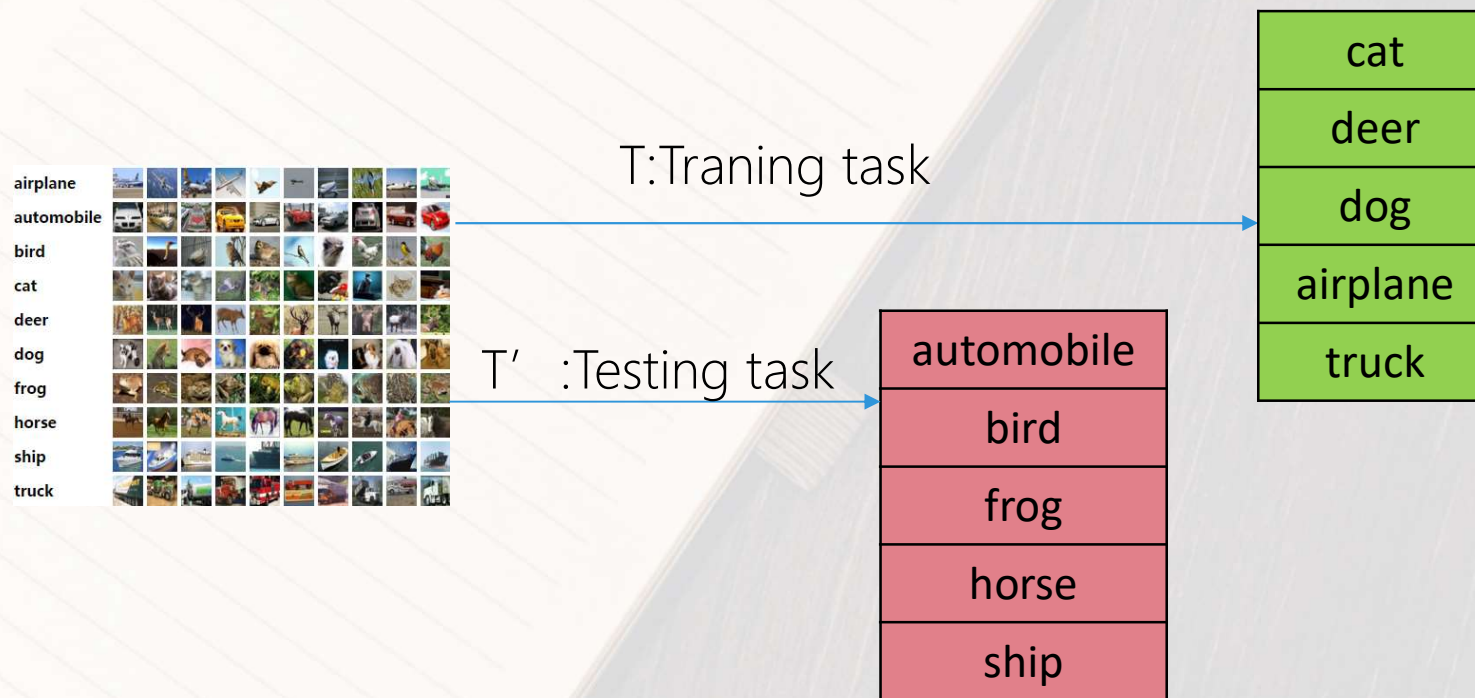
B

\hat{x}

N-ways k-shot learning

What is one-shot learning?

- Machine Learning Principle: Test and Train conditions Must Match
- Separate labels for training and testing
testing phase are not used in training phase !!!



Why do we need one-shot learning?

- Problems:
can not get enough training data
- few data for training/testing?

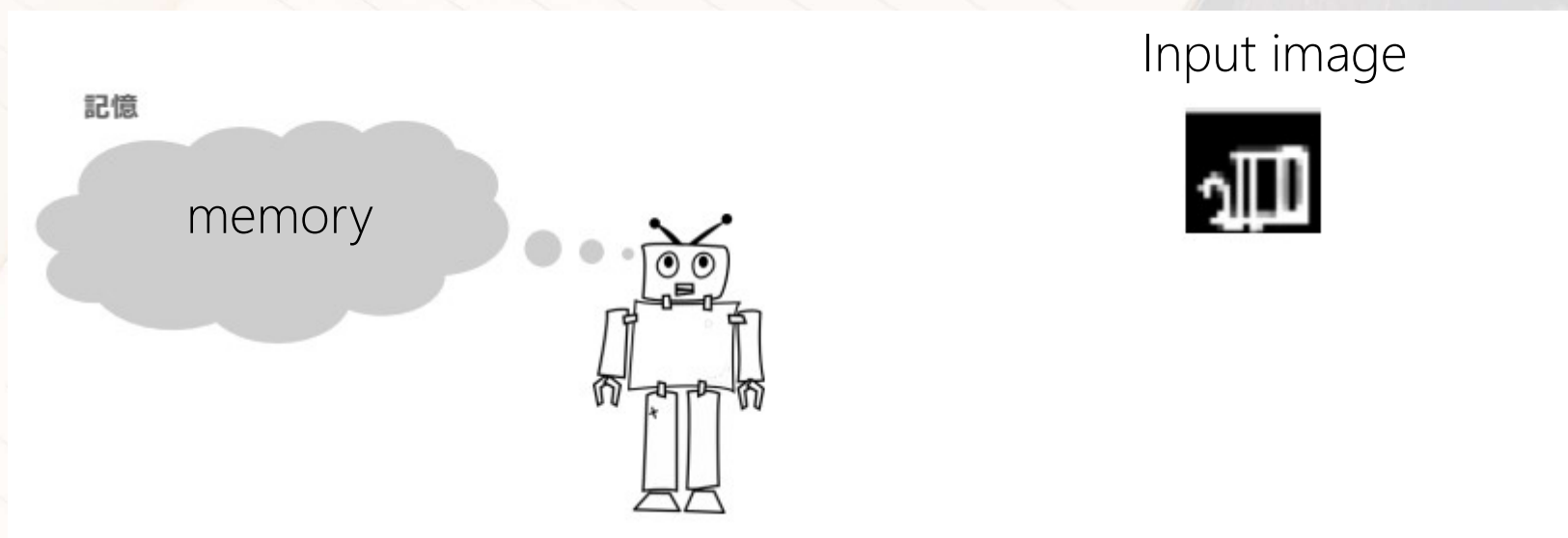


VS

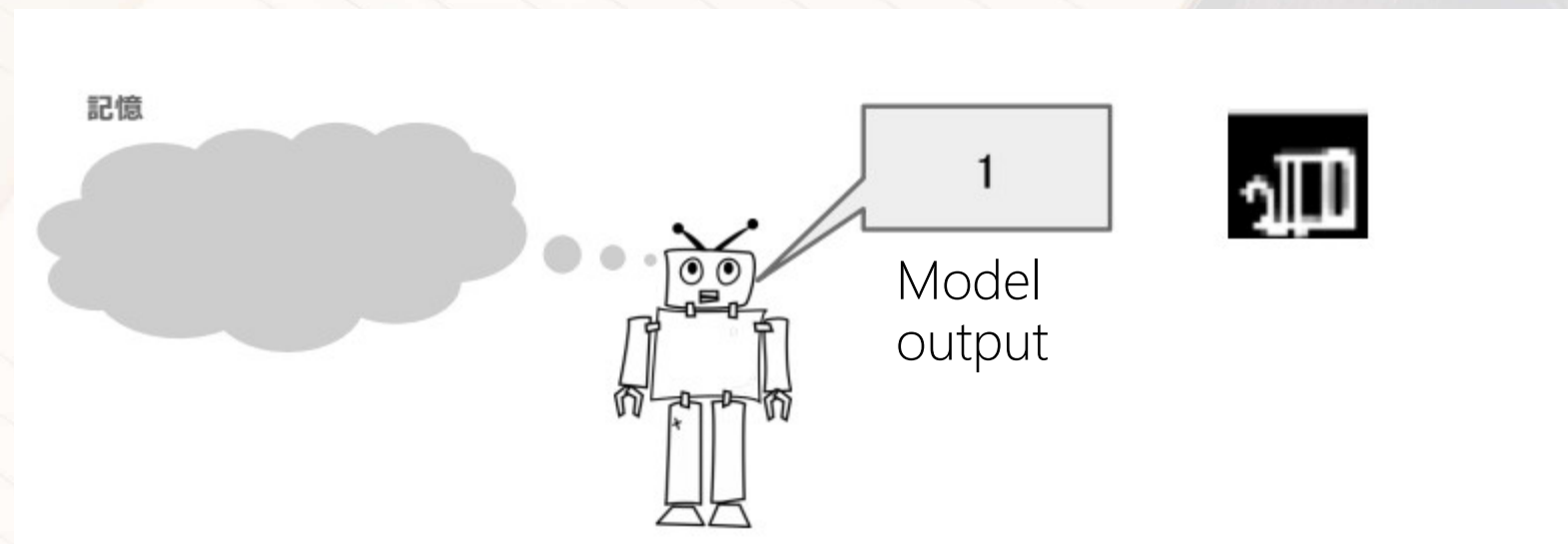


Neural Turing Machine(NTM)

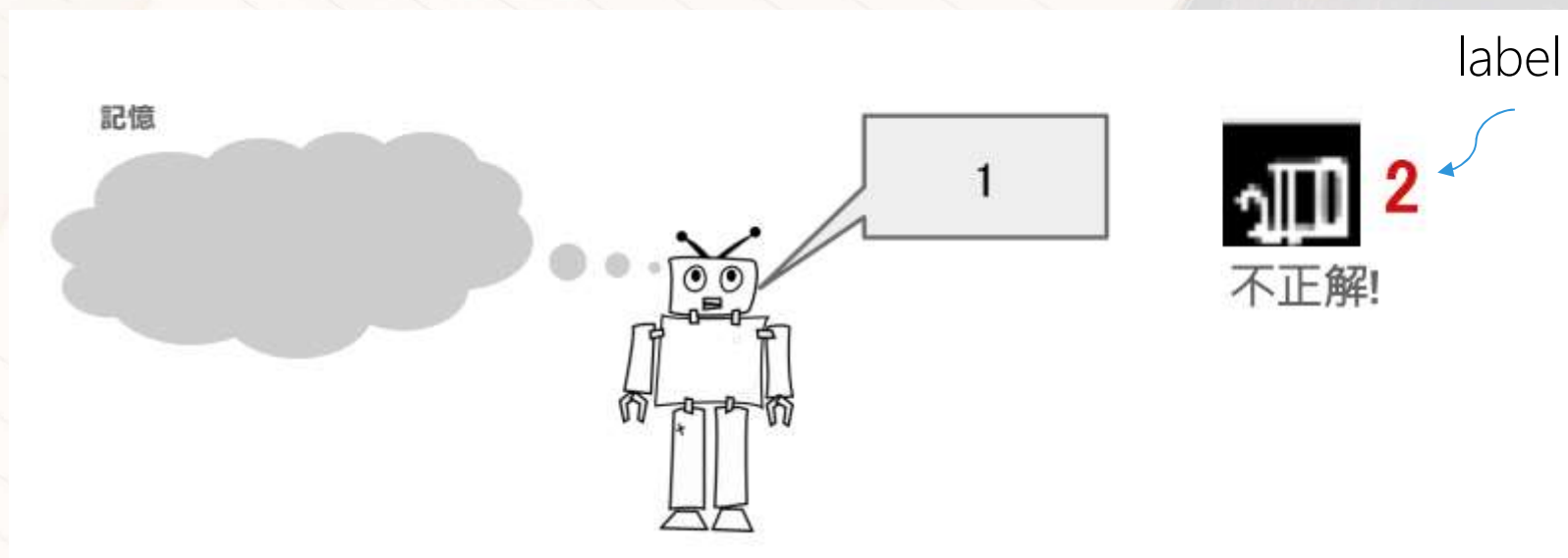
Memory will update with training



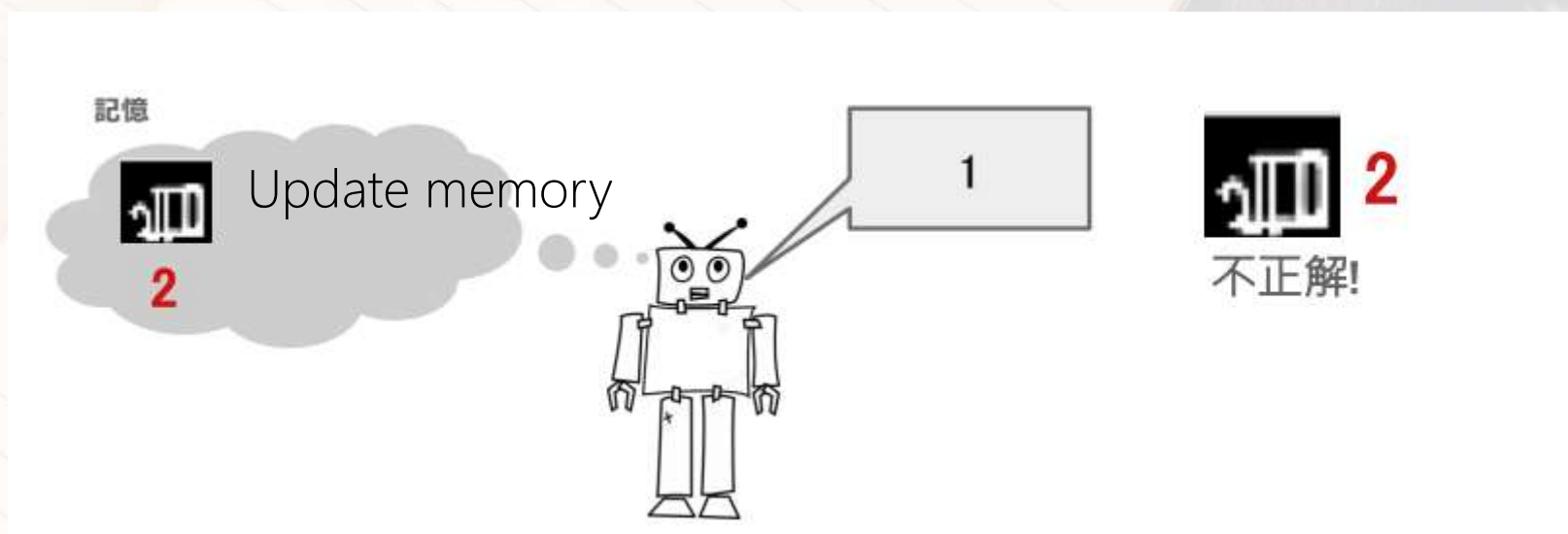
Neural Turing Machine(NTM)



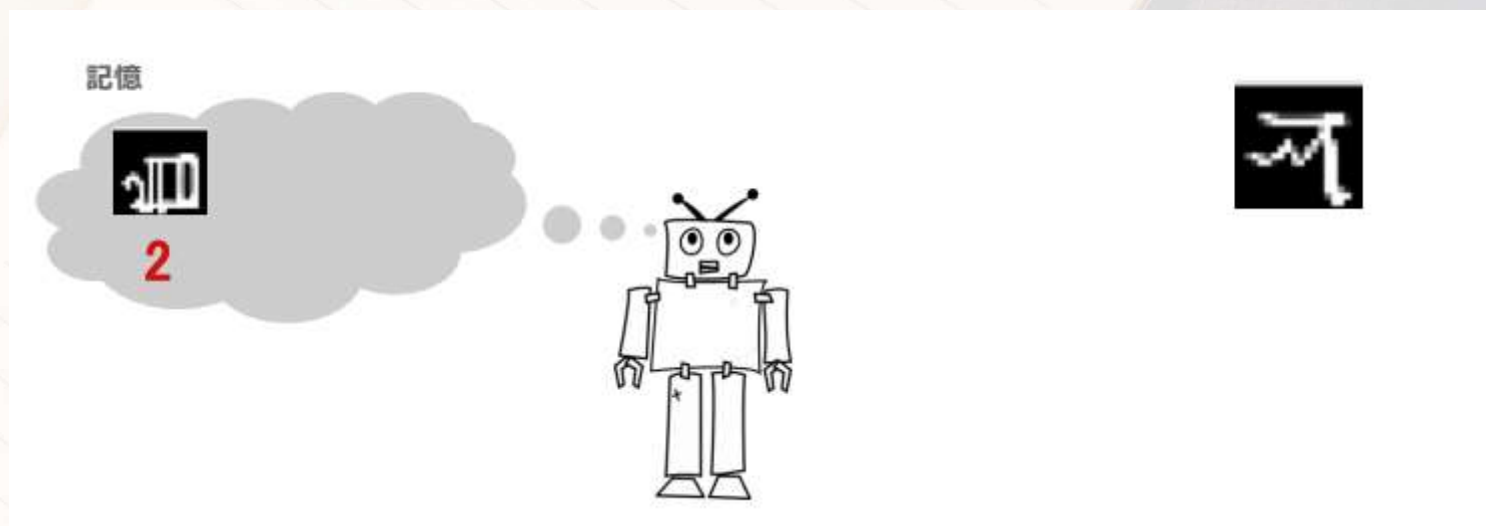
Neural Turing Machine(NTM)



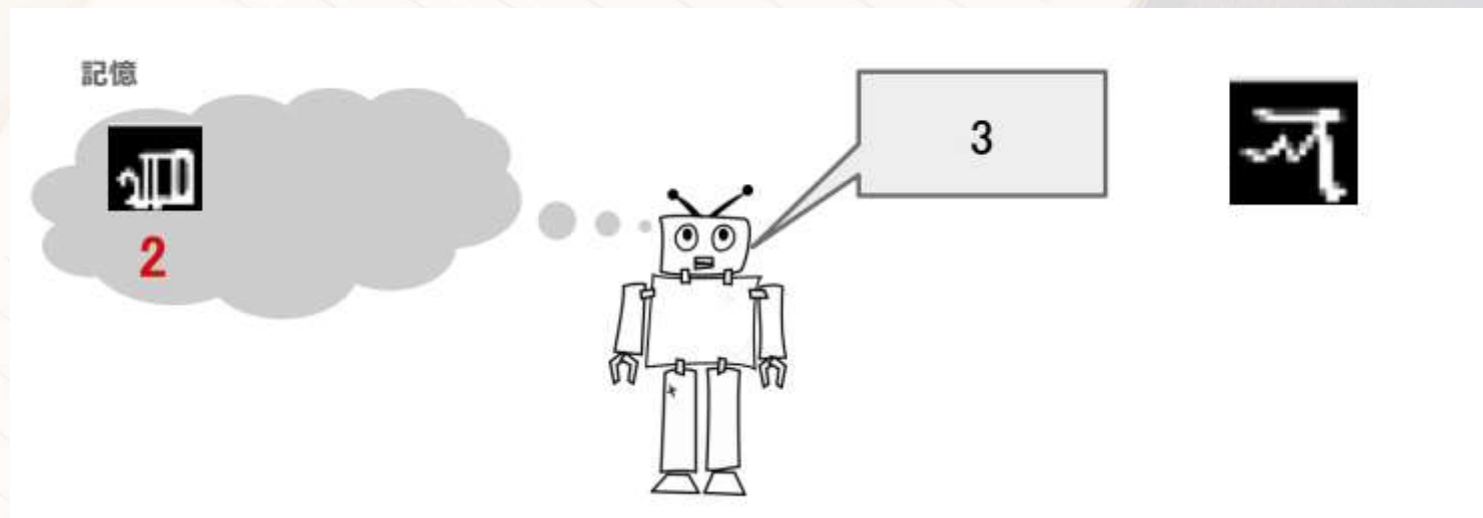
Neural Turing Machine(NTM)



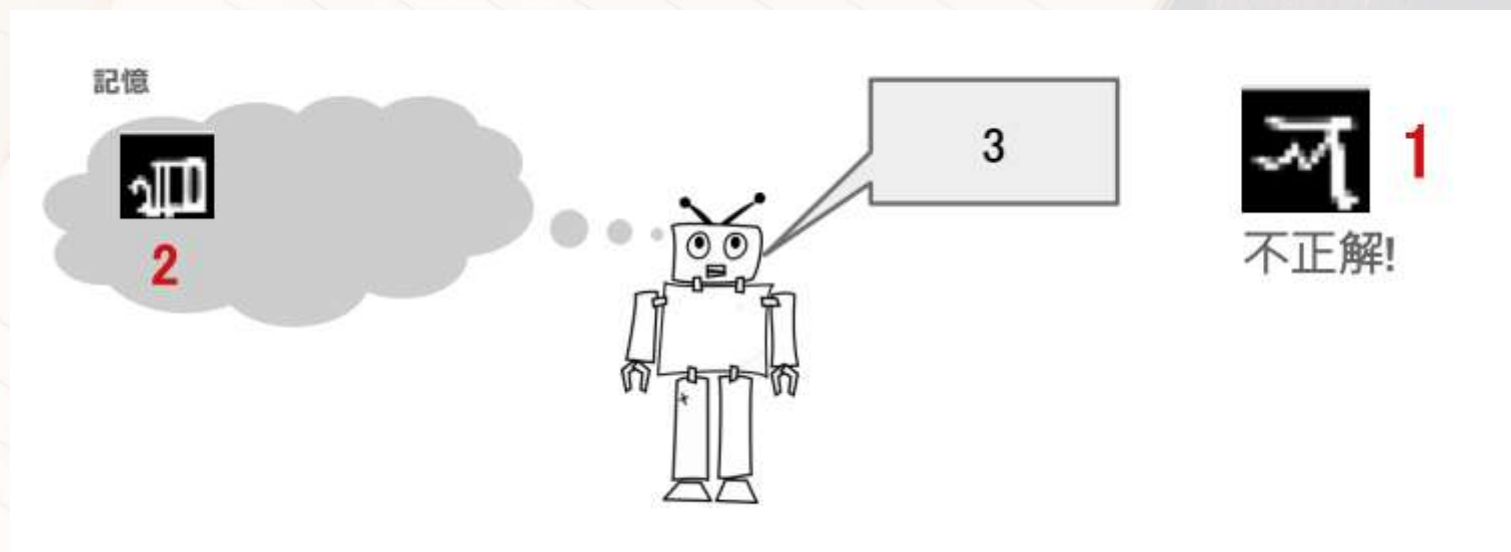
Neural Turing Machine(NTM)



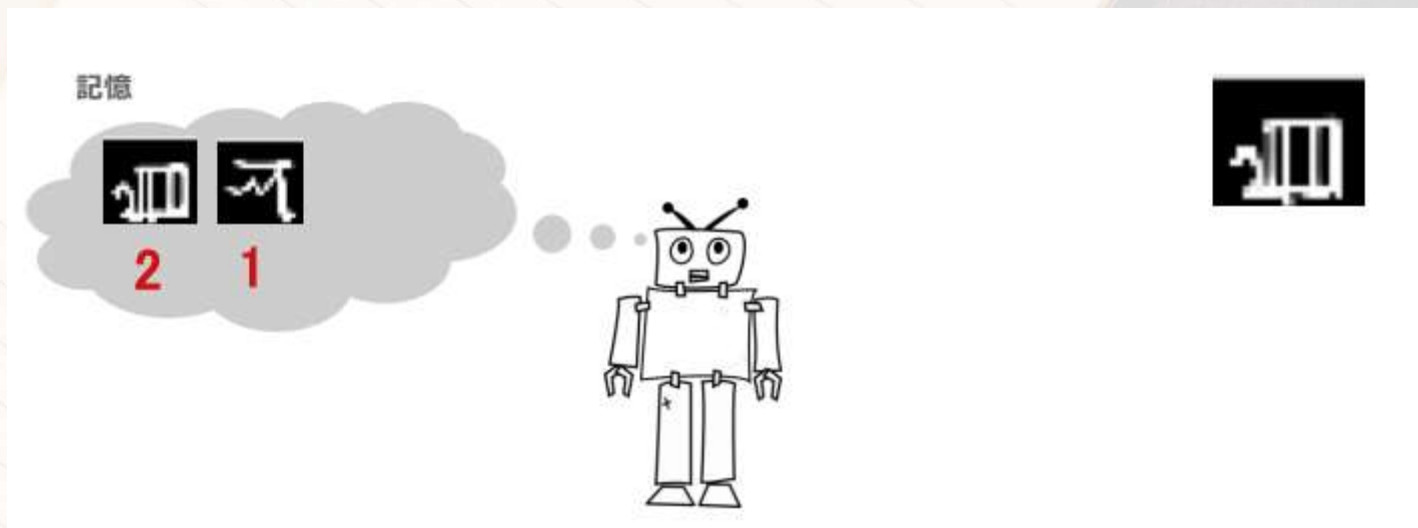
Neural Turing Machine(NTM)



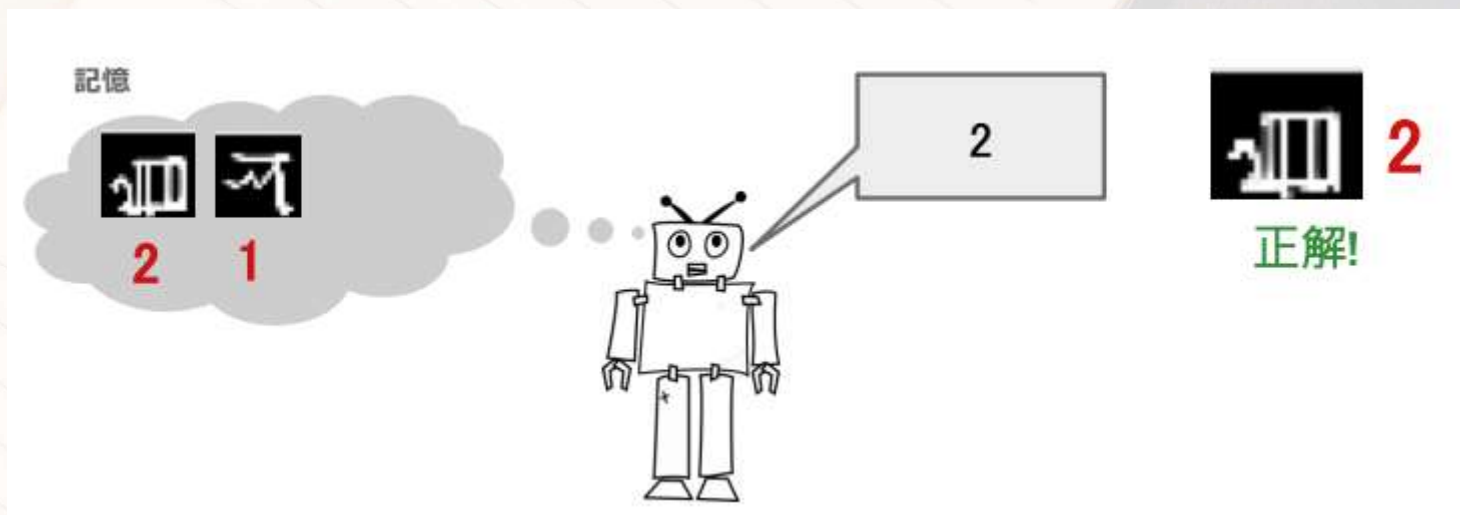
Neural Turing Machine(NTM)



Neural Turing Machine(NTM)



Neural Turing Machine(NTM)



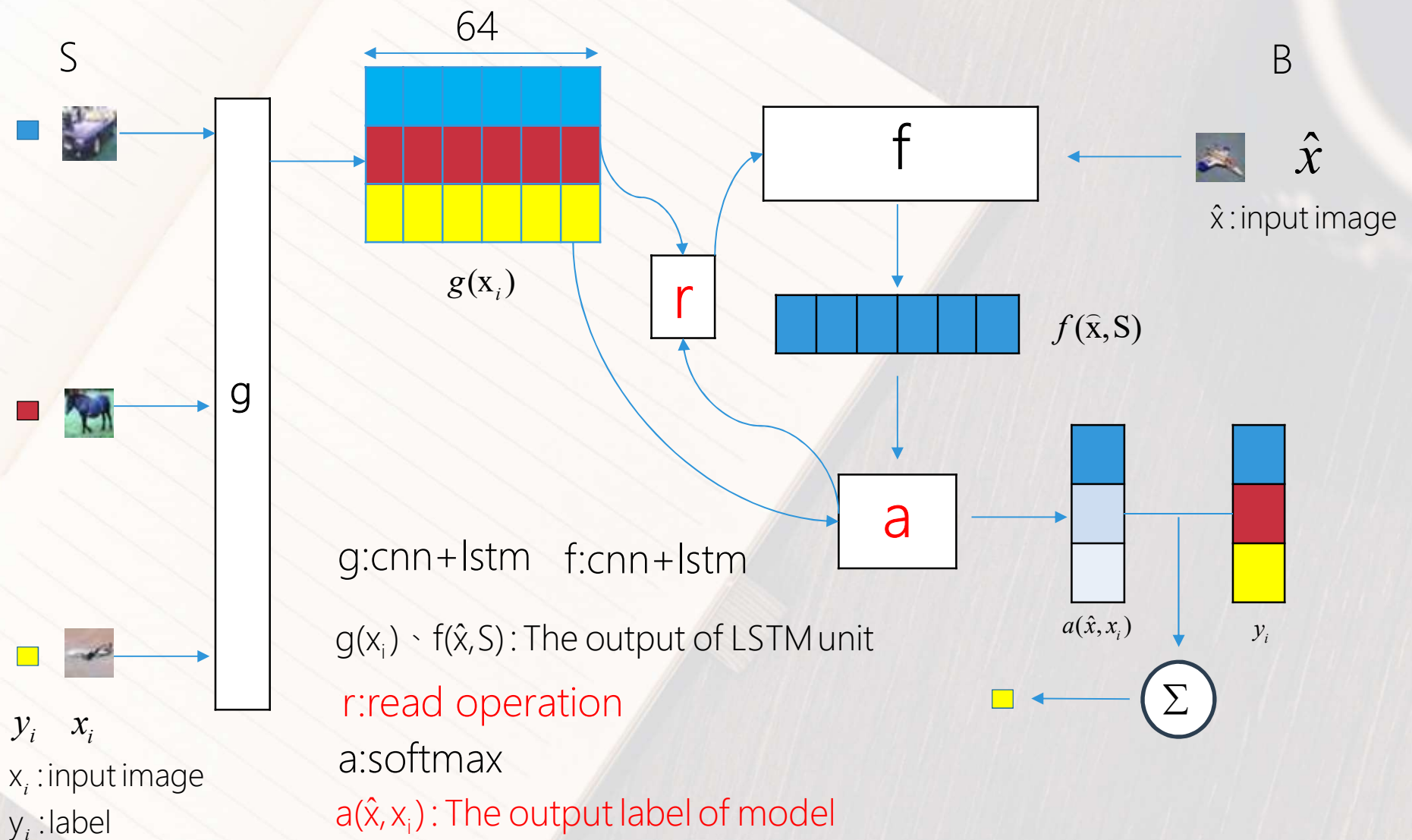
One-shot learning with Matching Networks

- Background
- Introduction
 - Introduce to One Shot Learning
 - Introduce to MANN
- **Model**
 - Motivation
 - Matching Networks
 - Backpropagation
- Experiments
- Summary

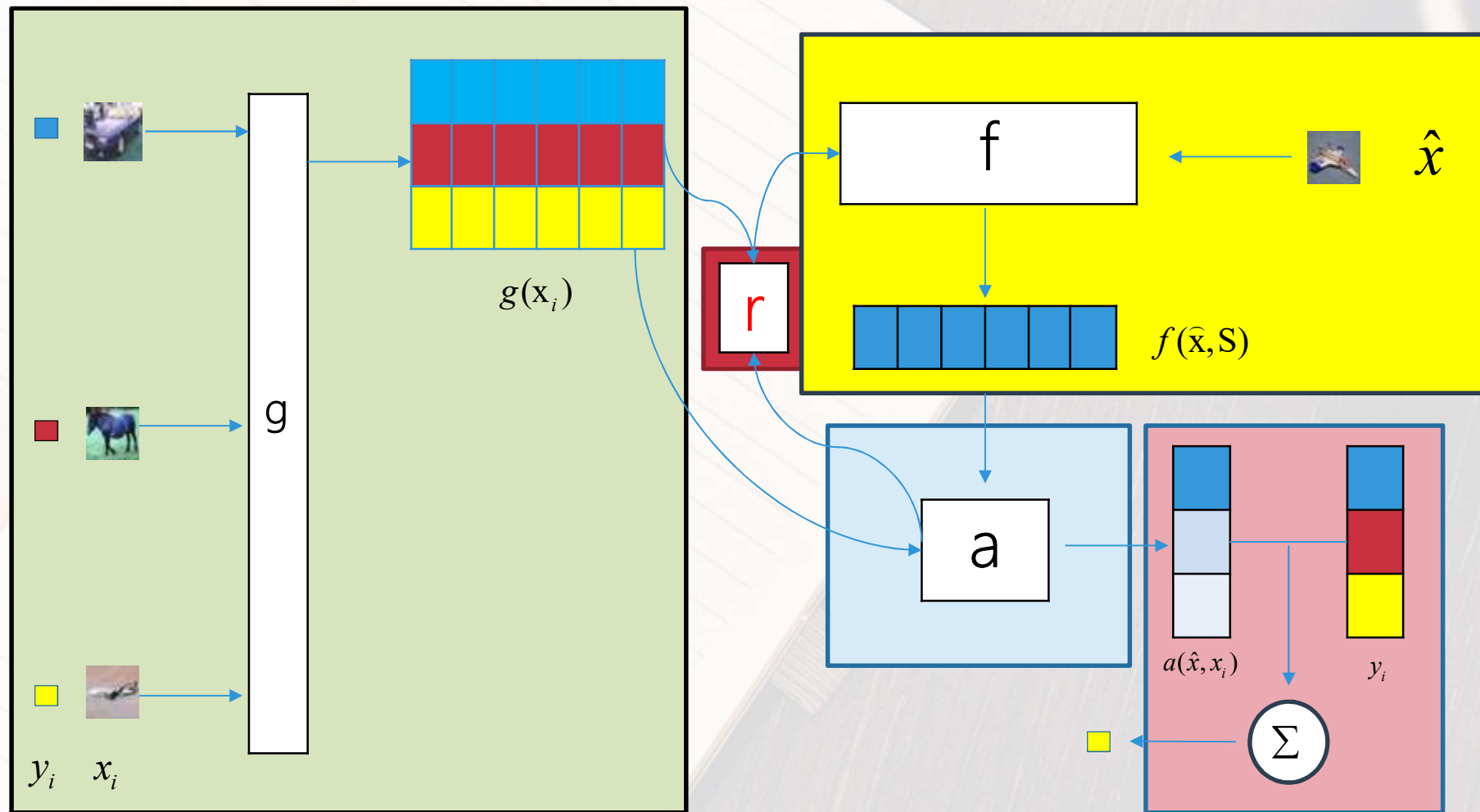
Model: Matching Networks

- Motivation:
 - Few data for training
 - Model should be updated all the time
 - Non-parametric models performance depends on the chosen metric
- Model
Matching Networks(non-parametric components)

Matching Networks

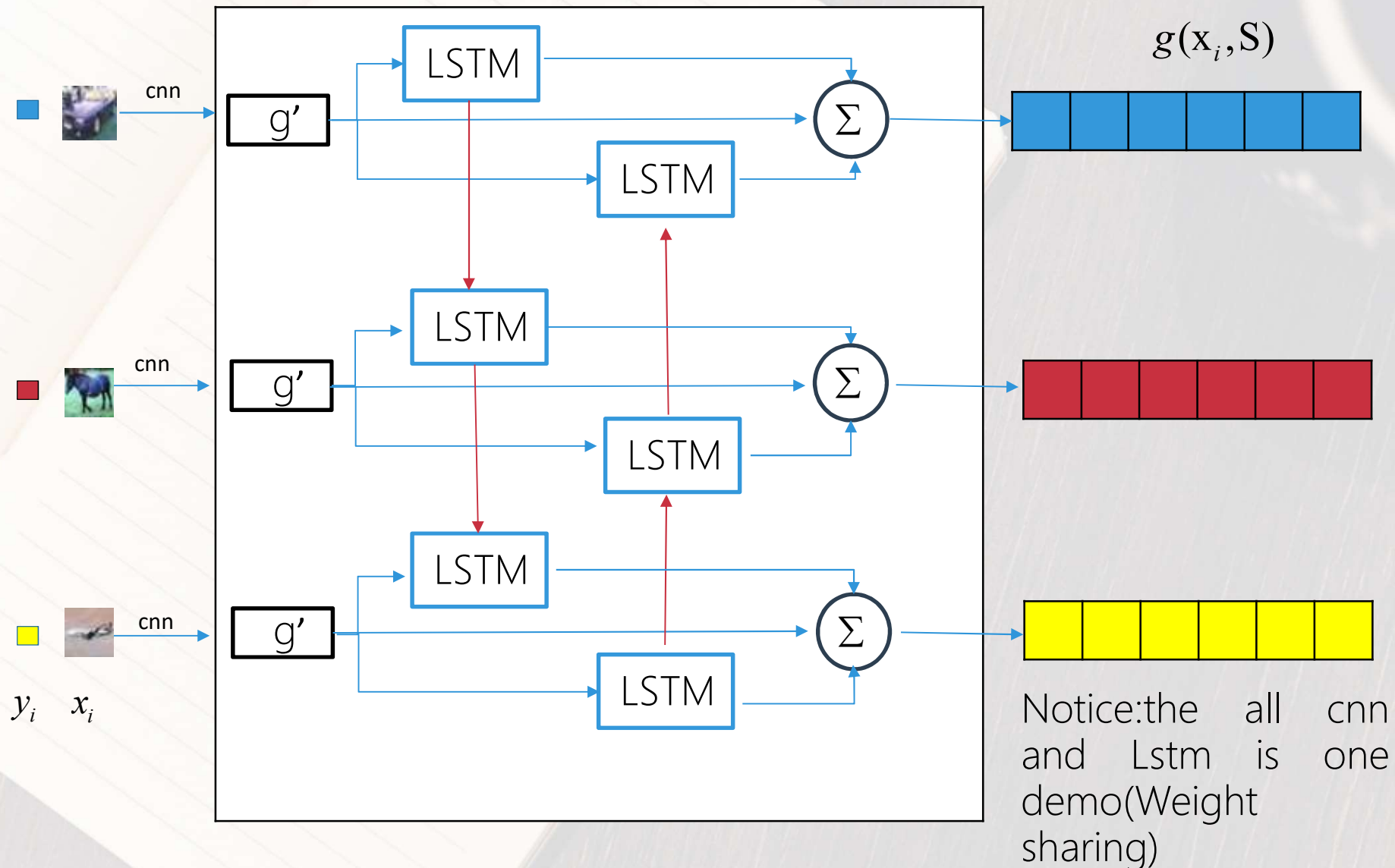


Matching Networks

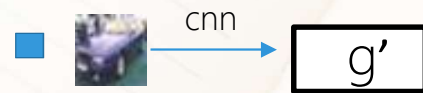


S:Support set

Matching Networks



Matching Networks



g' :Neural network

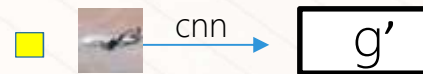
Input: image
size=28*28

Batch size:32



Model:
VGG model

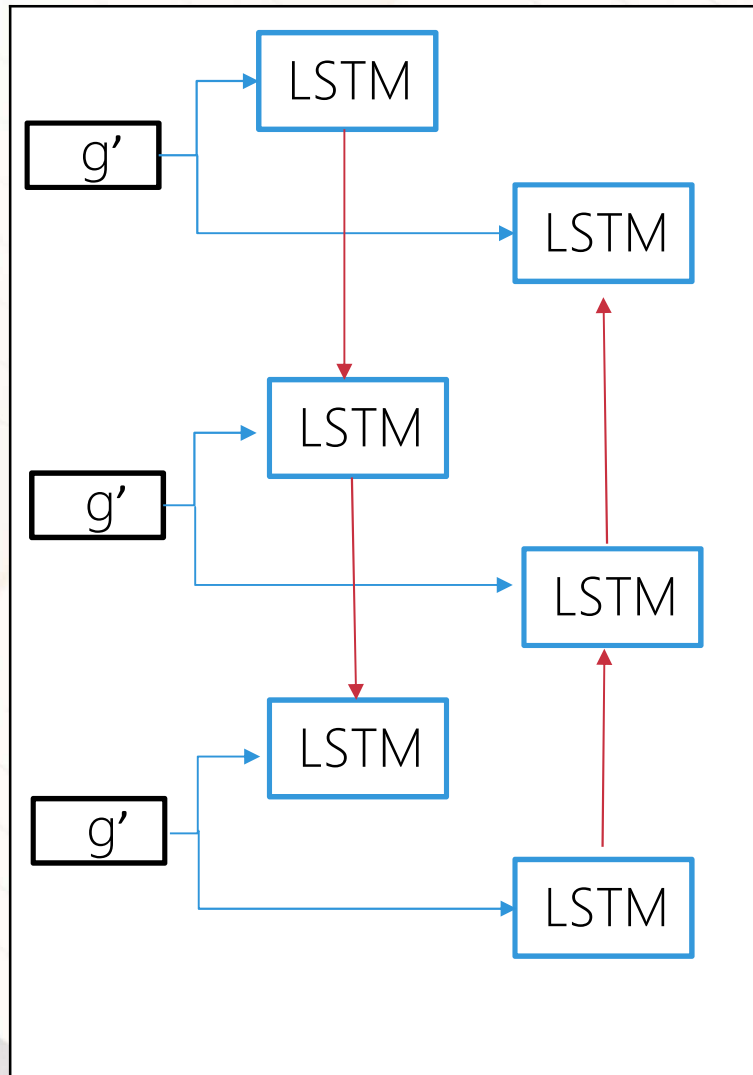
Output:
64*1 vector



y_i x_i

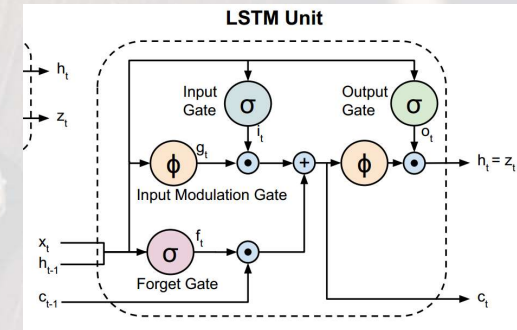
VGG model please reference:
<https://arxiv.org/abs/1409.1556>

Matching Networks



Input: g'

Output: $\vec{h}_i, \vec{c}_i, \bar{h}_i, \bar{c}_i$



$$\vec{h}_i, \vec{c}_i = LSTM(g'(x_i), \vec{h}_{i-1}, \vec{c}_{i-1})$$

$$\bar{h}_i, \bar{c}_i = LSTM(g'(x_i), \bar{h}_{i+1}, \bar{c}_{i+1})$$

LSTM DEMO : $x = g'$

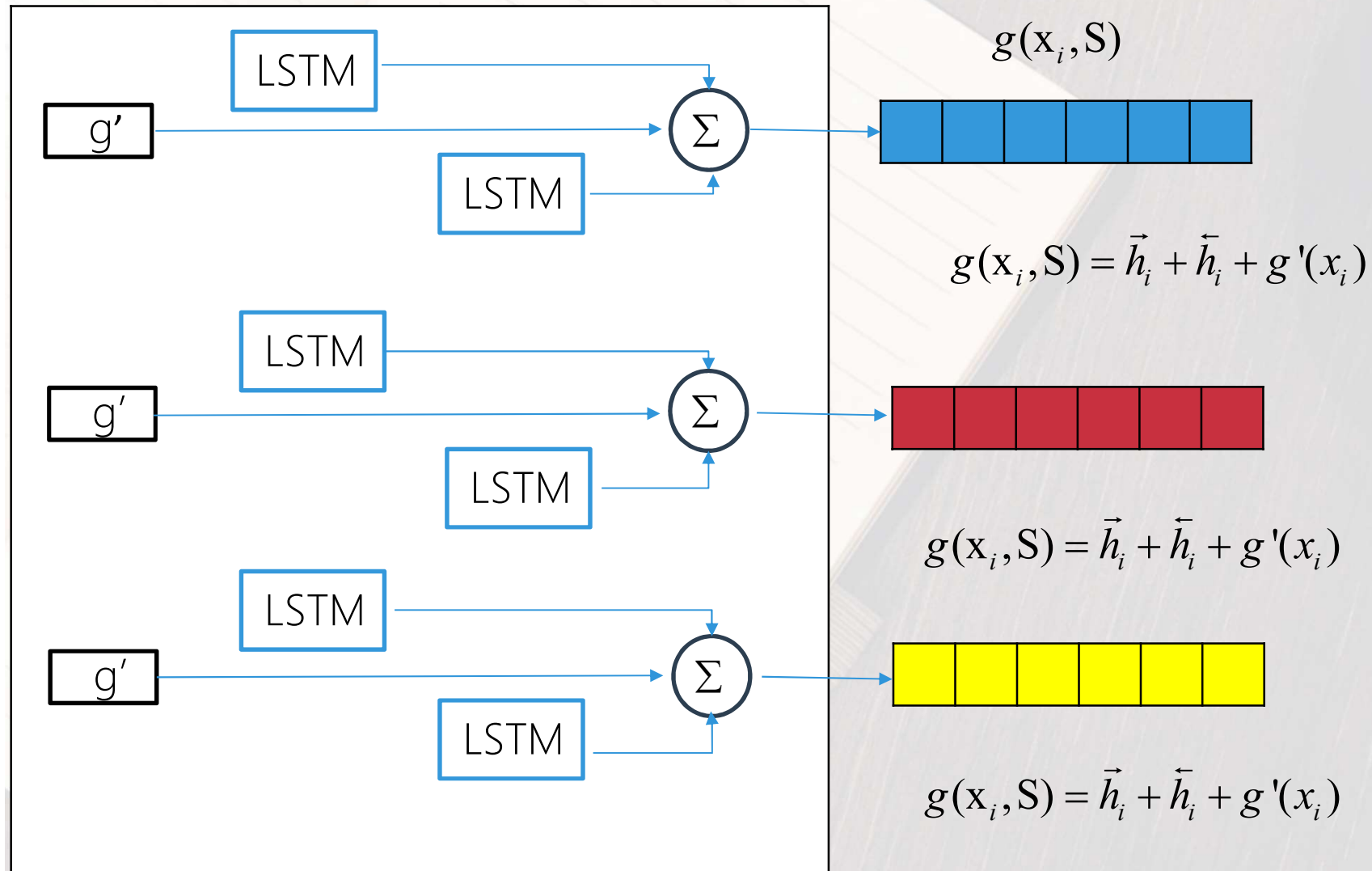
$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \quad f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \quad g_t = \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c)$$

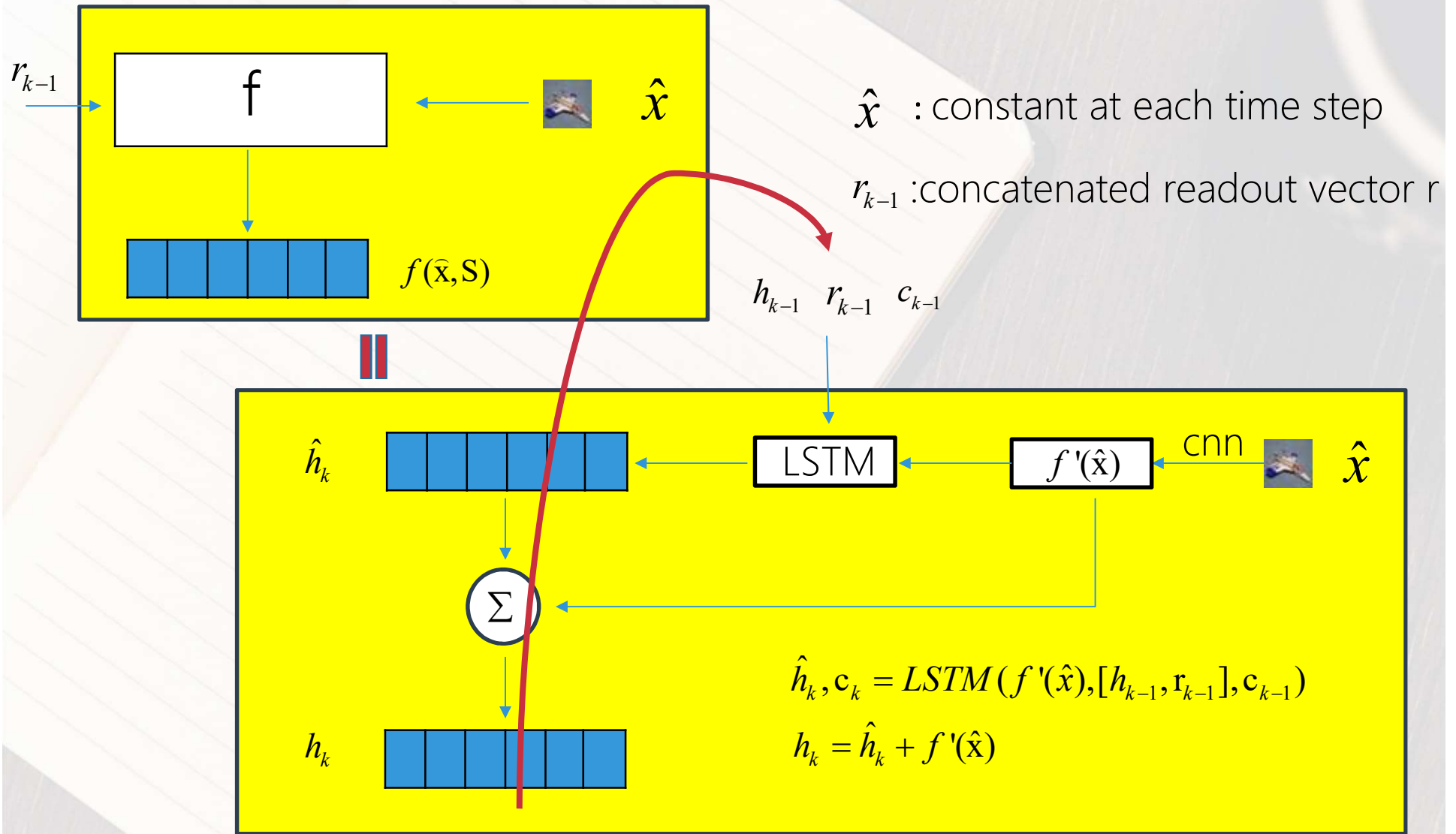
$$c_t = f_t \square c_{t-1} + i_t \square g_t \quad h_t = o_t \square \tanh(c_t)$$

Write operation

Matching Networks

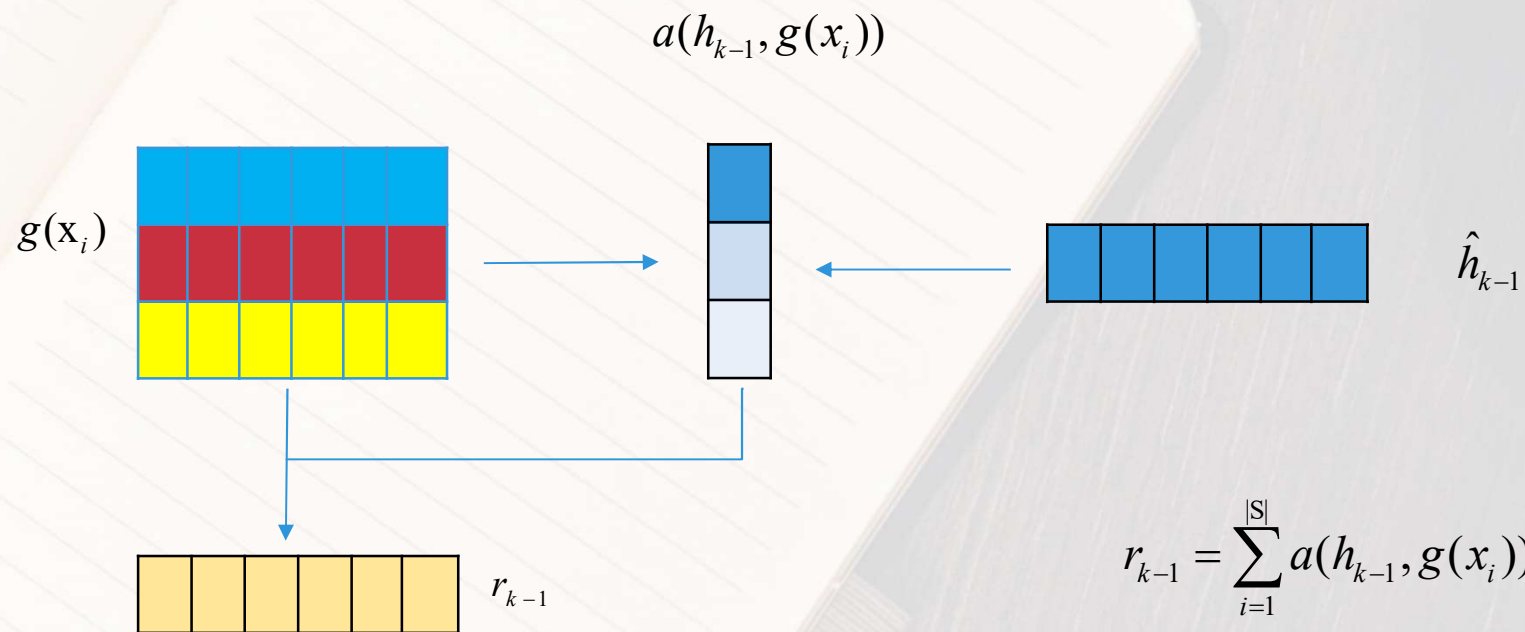


Matching Networks



Matching Networks

Calculate the relevance :

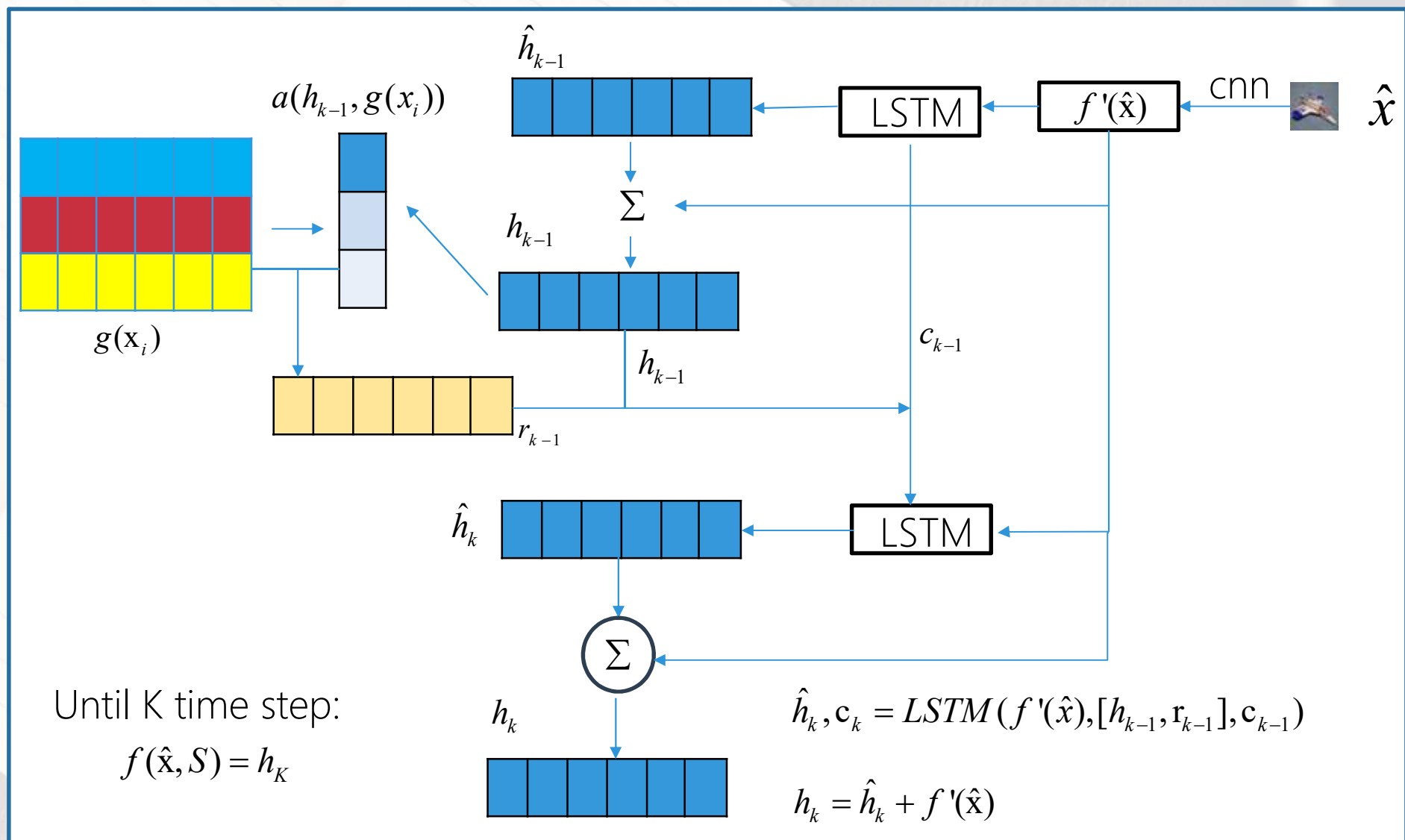


$$r_{k-1} = \sum_{i=1}^{|S|} a(h_{k-1}, g(x_i)) g(x_i)$$

$$a(f(\hat{x}), g(x_i)) = \text{soft max}(f(\hat{x}), g(x_i)) = \frac{e^{c(f(\hat{x})g(x_i))}}{\sum_{j=1}^k e^{c(f(\hat{x})g(x_j))}}$$

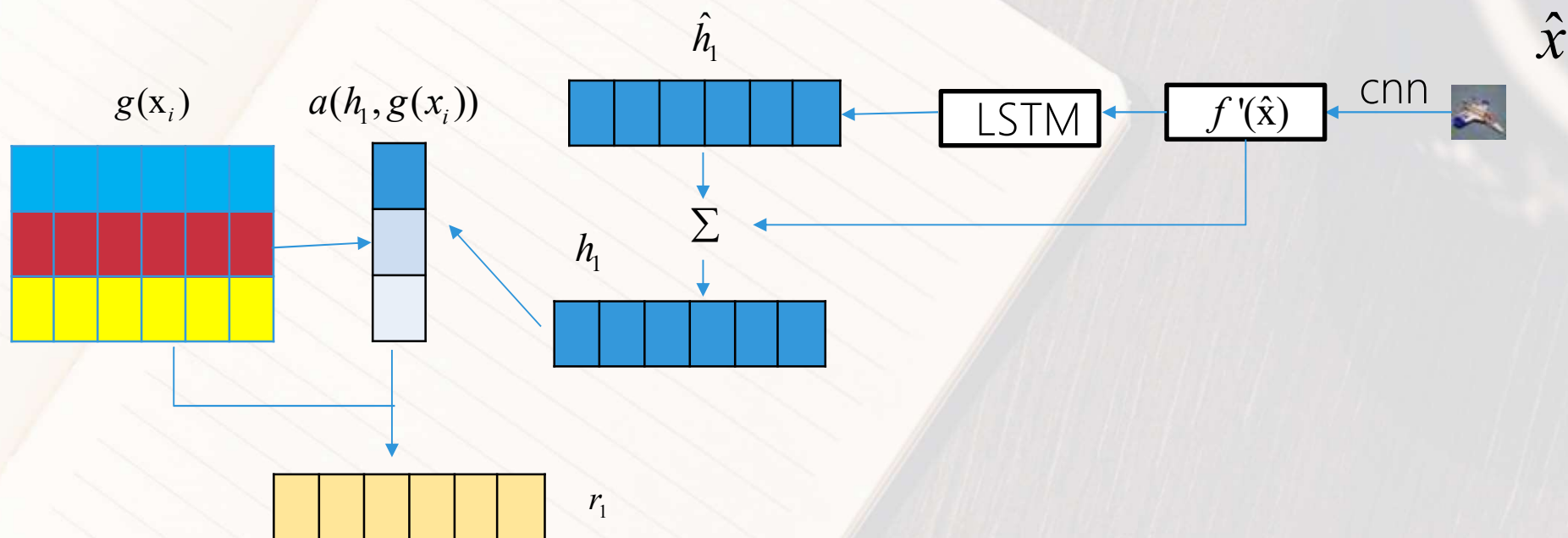
r (read) is a sum of g weighted according to the relevance to h

Matching Networks



Matching Networks

How to calculate h_1 ?

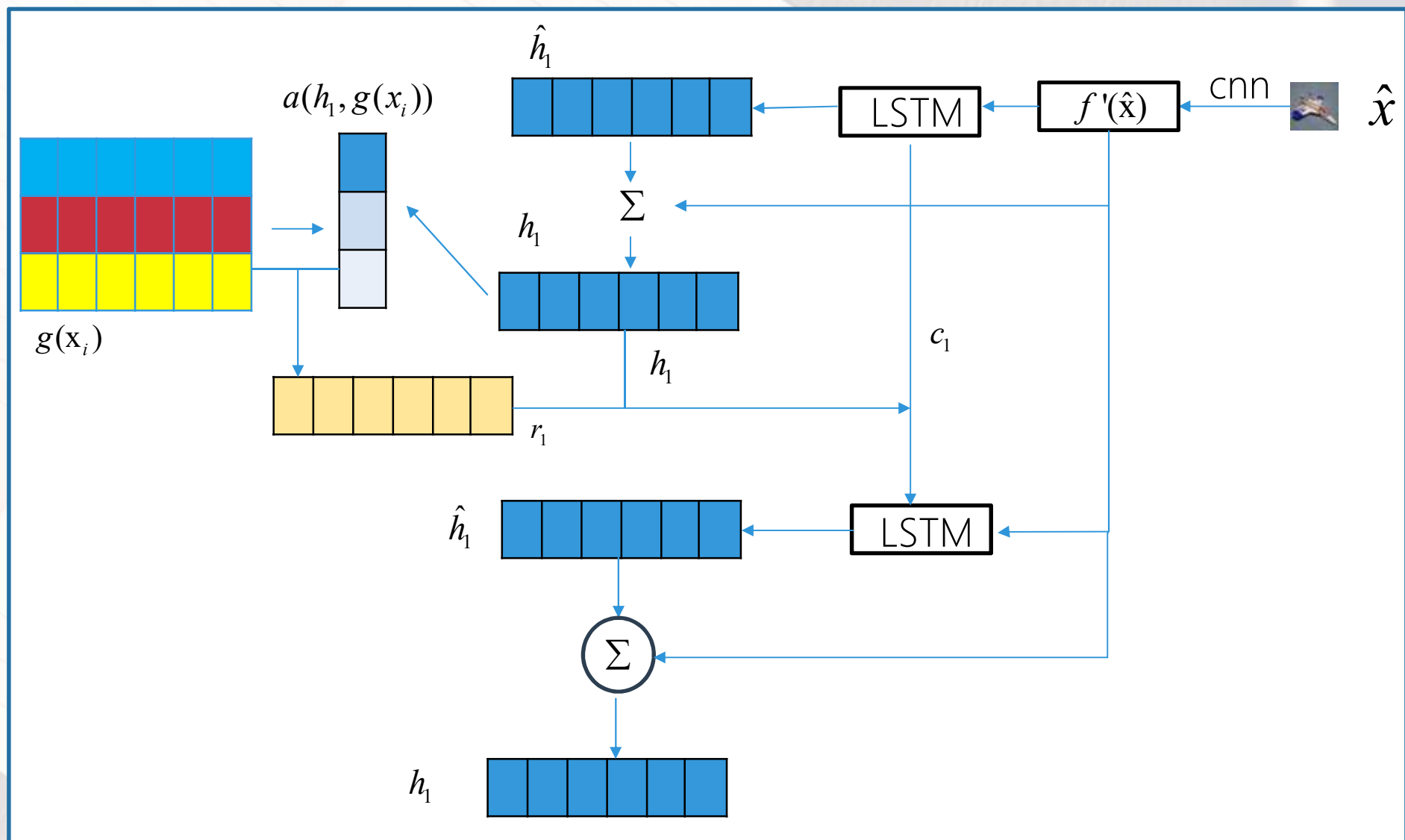


$$\hat{h}_1, c_1 = LSTM(f'(\hat{x}), [h_0, r_0], c_0)$$

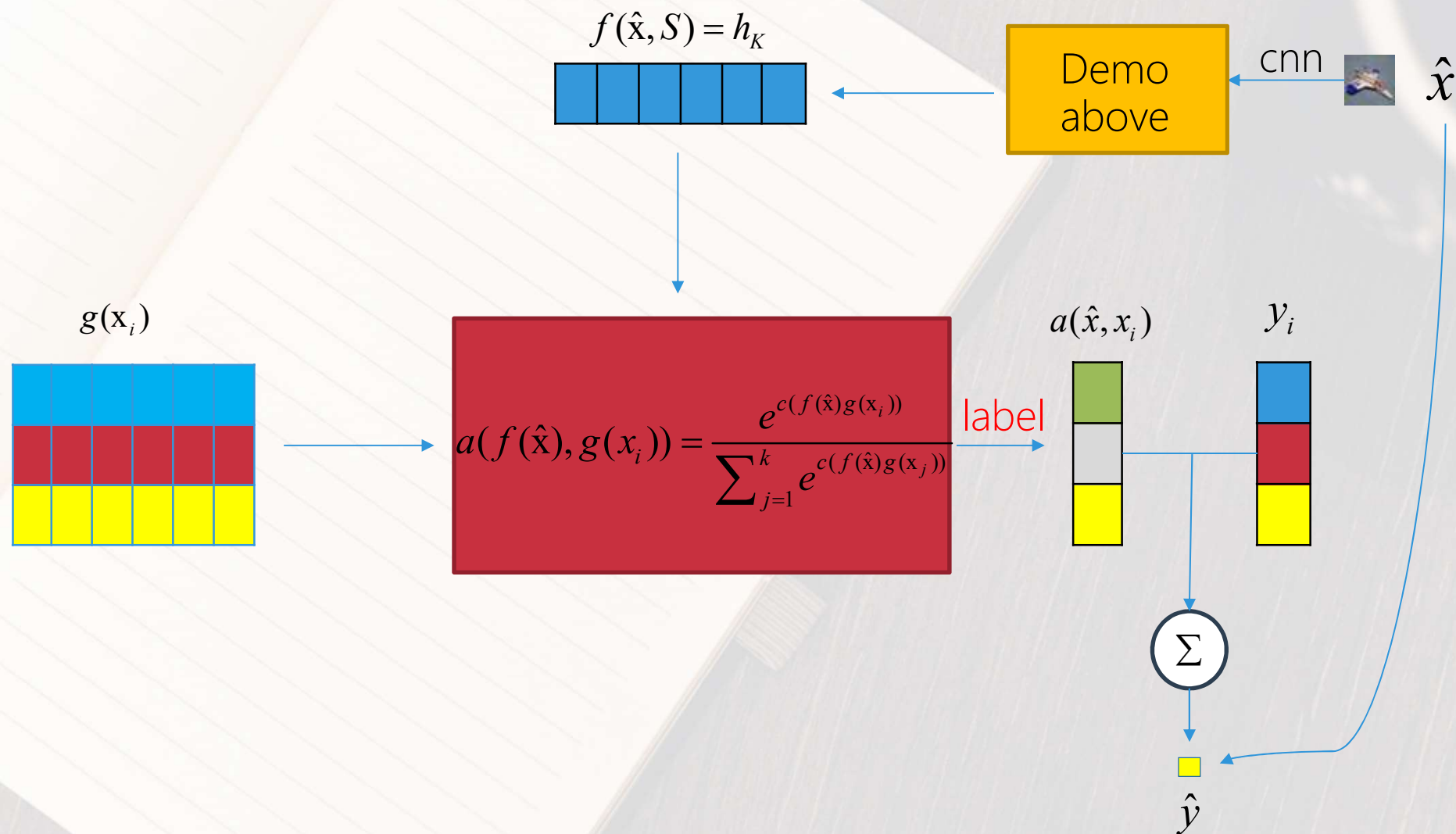
$$h_1 = \hat{h}_1 + f'(\hat{x})$$

$$\longrightarrow a(h_1, g(x_i)) \longrightarrow r_1 = \sum_{i=1}^{|S|} a(h_1, g(x_i))g(x_i)$$

Matching Networks



Matching Networks



Matching Networks

$a: [1, 8e-39, 2e-35, 0, 1.4e-35]$

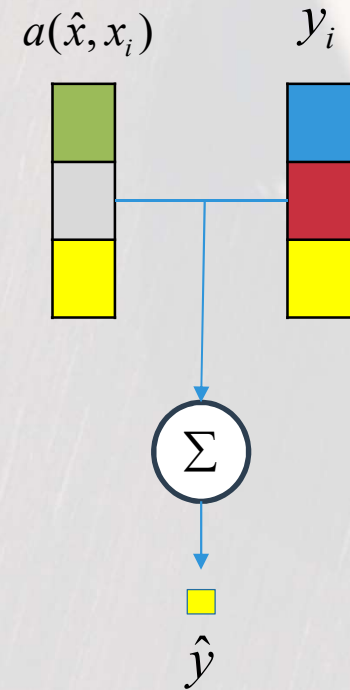
$y: [3, 1, 0, 2, 4]$ Five ways-one shot

output_label = 3

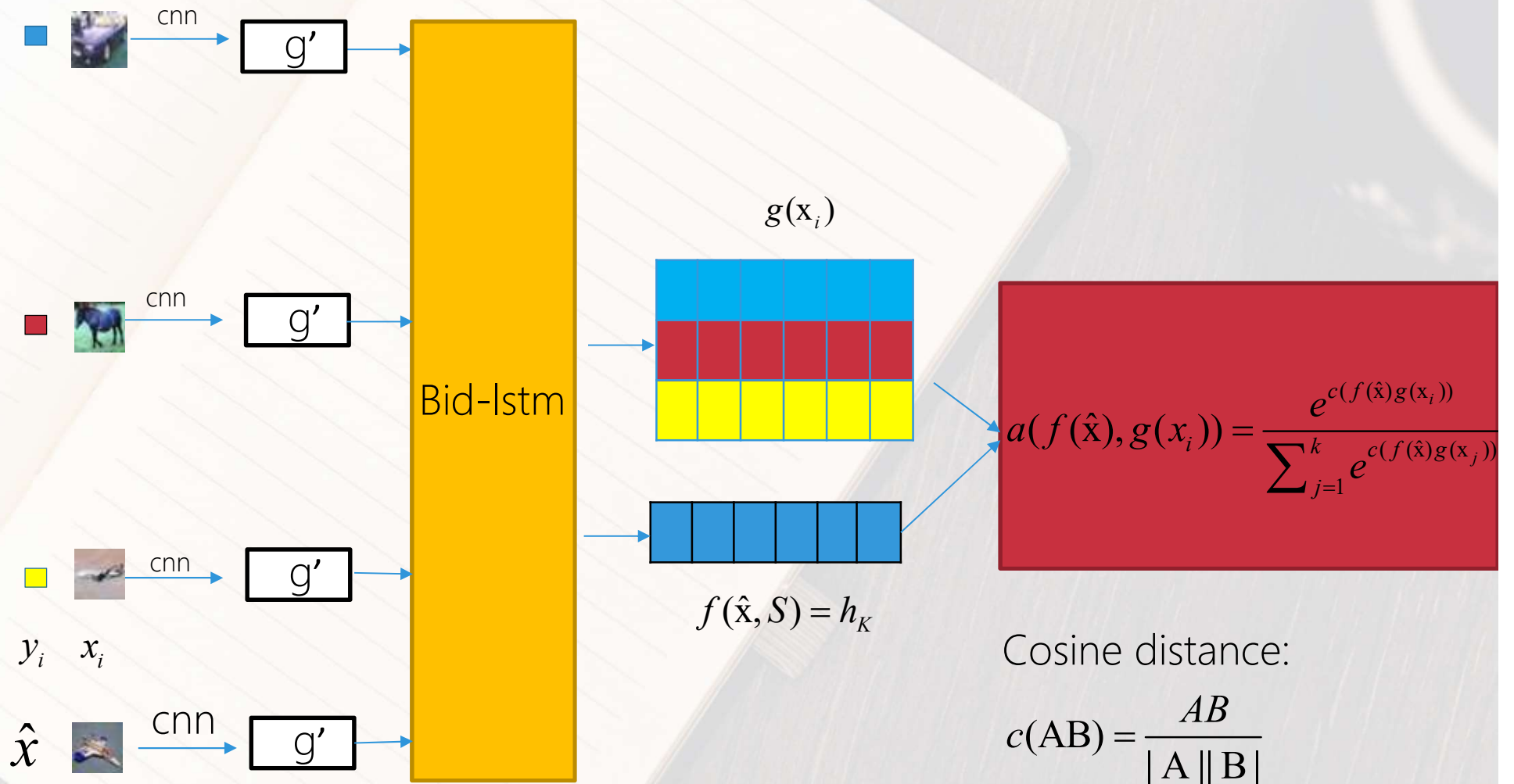
We use $a(\hat{x}, x_i)$ as the probability of memory labels y_i . The max location of $a(\hat{x}, x_i)$ as index and find the corresponding position of y_i as our output label.

$$P(\hat{y} | \hat{x}, S) = \sum_{i=1}^k a(\hat{x}, x_i) y_i$$

All of above called "episode "



Matching Networks without lstm



Backpropagation

To form an “episode” to compute gradients and update our model:

L :Data set

S :Support set (sample from L)

B :Batch (sample from L)

\hat{x} (image), \hat{y} (label) is one input of model

θ is the weight

Conditional Probability:

$$P(\hat{y} | \hat{x}, S) = \sum_{i=1}^k a(\hat{x}, x_i) y_i$$

$$\theta = \arg \max_{\theta} E_{L \sim T} [E_{S \sim L, B \sim L} [\sum_{(x,y) \in B} \log P_{\theta}(y | x, S)]]$$

We learning the mapping function P of model!

Backpropagation

$$\theta = \arg \max_{\theta} \sum_{(x,y) \in B} \log P_{\theta}(y | x, S) = -\arg \min \textit{Costfunction}$$

$$\textit{Cost} = \sum_{(x,y) \in B} \log P(\hat{y} | \hat{x}, S)$$

$$P(\hat{y} | \hat{x}, S) = \sum_{i=1}^k a(\hat{x}, x_i) y_i$$

For every $\hat{x}, \hat{y} \sim B$

$$E_{\hat{x}, \hat{y}} = \log P(\hat{y} | \hat{x}, S)$$

$$\begin{aligned} \frac{\partial E}{\partial W} \Big|_{\theta=-W} &= \frac{1}{P} \frac{\partial P}{\partial W} = \frac{1}{P} \frac{\partial}{\partial W} (a(\hat{x}, x_1) y_1 + \dots a(\hat{x}, x_k) y_k) \\ &= \frac{1}{P} \left(\frac{\partial a(\hat{x}, x_1)}{\partial W_1} y_1 + \dots \frac{\partial a(\hat{x}, x_k)}{\partial W_k} y_k \right) \end{aligned}$$

Backpropagation

Forward propagation

$$a(f(\hat{x}), g(x_i)) = \frac{e^{f(\hat{x})g(x_i)}}{\sum_{j=1}^k e^{f(\hat{x})g(x_j)}}$$

Backpropagation

$$\begin{aligned} \frac{\partial a(f(x), g(x_i))}{\partial W_i} &= \frac{\partial \frac{e^{f(\hat{x})g(x_i)}}{\sum_{j=1}^k e^{f(\hat{x})g(x_j)}}}{\partial W_i} \\ &= \frac{(e^{f(\hat{x})g(x_i)})' \sum_{j=1}^k e^{f(\hat{x})g(x_j)} - e^{f(\hat{x})g(x_i)} (\sum_{j=1}^k e^{f(\hat{x})g(x_j)})'}{(\sum_{j=1}^k e^{f(\hat{x})g(x_j)})^2} \end{aligned}$$

$$\begin{aligned} \frac{\partial e^{f(\hat{x})g(x_i)}}{\partial W} &= e^{f(\hat{x})g(x_i)} \frac{\partial (f(\hat{x})g(x_i))}{\partial W} \\ &= e^{f(\hat{x})g(x_i)} \left(\frac{\partial (f(\hat{x}))}{\partial W} g(x_i) + \frac{\partial (g(x_i))}{\partial W} f(\hat{x}) \right) \end{aligned}$$

Backpropagation

Forward propagation

$$h_k = LSTM(f'(\hat{x}), [h_{k-1}, \mathbf{r}_{k-1}], \mathbf{c}_{k-1}) + f'(\hat{x})$$

Backpropagation

$$\frac{\partial h_k}{\partial W} = \frac{\partial \hat{h}_k}{\partial W} + \frac{\partial f'(\hat{x})}{\partial W}$$

$$\frac{\partial f'(\hat{x})}{\partial W} \square \frac{\partial g'(x)}{\partial W}$$

Backpropagation

Forward propagation

$$\hat{h}_k, \mathbf{c}_k = LSTM(f'(\hat{x}), [h_{k-1}, r_{k-1}], \mathbf{c}_{k-1})$$

$$i_k = \sigma(W_{xi}f' + W_{hi}h_{k-1} + W_{ri}r_{k-1} + b_i)$$

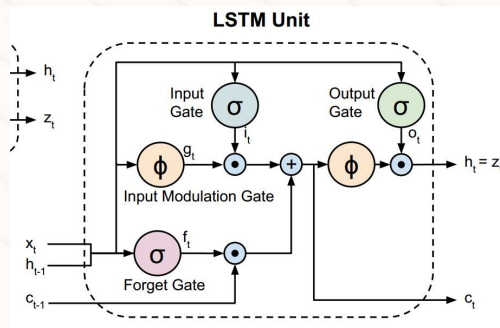
$$f_k = \sigma(W_{xf}f' + W_{hf}h_{k-1} + W_{rf}r_{k-1} + b_f)$$

$$o_k = \sigma(W_{xo}f' + W_{ho}h_{k-1} + W_{ro}r_{k-1} + b_o)$$

$$g_k = \tanh(W_{xc}f' + W_{hc}h_{k-1} + W_{rc}r_{k-1} + b_c)$$

$$c_k = f_k \square c_{k-1} + i_k \square g_k$$

$$h_k = o_k \square \tanh(c_k)$$



Backpropagation

Layers t:

$$\delta h^t = \frac{\delta LSTM}{\delta W}$$

$$\delta o^t = \delta h^t \square \tanh(c^t)$$

$$\delta c^t = \delta h^t \square o^t \square (1 - \tanh^2(c^t)) + \delta c^{t+1} \square f^{t+1}$$

$$\delta i^t = \delta c^t \square g^t \quad \delta f^t = \delta c^t \square c^{t-1}$$

$$\delta g^t = \delta c^t \square i^t \quad \delta c^{t-1} = \delta c^t \square f^t$$

$$\delta \hat{i}^t = \delta i^t \square i^t \square (1 - i^t)$$

$$\delta \hat{f}^t = \delta f^t \square f^t \square (1 - f^t)$$

$$\delta \hat{g}^t = \delta g^t \square (1 - \tanh^2(\hat{g}_t))$$

Backpropagation

Forward propagation

$$\hat{h}_k, \mathbf{c}_k = LSTM(f'(\hat{x}), [h_{k-1}, r_{k-1}], \mathbf{c}_{k-1})$$

$$\hat{i}_k = W_{xi}f' + W_{hi}h_{k-1} + W_{ri}r_{k-1} + b_i$$

$$\hat{f}_k = W_{xf}f' + W_{hf}h_{k-1} + W_{rf}r_{k-1} + b_f$$

$$\hat{o}_k = W_{xo}f' + W_{ho}h_{k-1} + W_{ro}r_{k-1} + b_o$$

$$\hat{g}_k = W_{xc}f' + W_{hc}h_{k-1} + W_{rc}r_{k-1} + b_c$$

Backpropagation

Layers t:

$$\delta \hat{o}^t = \delta o^t \square o^t \square (1 - o^t)$$

$$\tanh'(\hat{g}^t) = 1 - \tanh^2(\hat{g}^t)$$

$$\delta \mathbf{z}^t = [\delta \hat{g}^t, \delta \hat{i}^t, \delta \hat{f}^t, \delta \hat{o}^t]$$

for simple: we can make:

$$\mathbf{z}^t = \begin{bmatrix} \hat{g}^t \\ \hat{i}^t \\ \hat{f}^t \\ \hat{o}^t \end{bmatrix} = \begin{bmatrix} W_{gx} & W_{gh} & W_{gr} \\ W_{ix} & W_{ih} & W_{ir} \\ W_{fx} & W_{fh} & W_{fr} \\ W_{ox} & W_{oh} & W_{or} \end{bmatrix} \begin{bmatrix} x_t \\ h_{t-1} \\ r_{t-1} \end{bmatrix} + \begin{bmatrix} b_g \\ b_i \\ b_f \\ b_o \end{bmatrix}$$

Backpropagation

Forward propagation

$$\hat{h}_k, \mathbf{c}_k = LSTM(f'(\hat{x}), [h_{k-1}, r_{k-1}], \mathbf{c}_{k-1})$$

$$\hat{i}_k = W_{xi}f' + W_{hi}h_{k-1} + W_{ri}r_{k-1} + b_i$$

$$\hat{f}_k = W_{xf}f' + W_{hf}h_{k-1} + W_{rf}r_{k-1} + b_f$$

$$\hat{o}_k = W_{xo}f' + W_{ho}h_{k-1} + W_{ro}r_{k-1} + b_o$$

$$\hat{g}_k = W_{xc}f' + W_{hc}h_{k-1} + W_{rc}r_{k-1} + b_c$$

Backpropagation

Layers t:

$$\mathbf{z}^t = \begin{bmatrix} \hat{g}^t \\ \hat{i}^t \\ \hat{f}^t \\ \hat{o}^t \end{bmatrix} = \begin{bmatrix} W_{gx} & W_{gh} & W_{gr} \\ W_{ix} & W_{ih} & W_{ir} \\ W_{fx} & W_{fh} & W_{fr} \\ W_{ox} & W_{oh} & W_{or} \end{bmatrix} \begin{bmatrix} x_t \\ h_{t-1} \\ r_{t-1} \end{bmatrix} + \begin{bmatrix} b_g \\ b_i \\ b_f \\ b_o \end{bmatrix}$$

$$\mathbf{z}^t = \mathbf{W} \mathbf{I}^t + \mathbf{B}$$

$$\delta \mathbf{W}^t = \delta \mathbf{z}^t \cdot (\mathbf{I}^t)^T \quad \delta \mathbf{I}^t = \delta \mathbf{z}^t \cdot \mathbf{W}^t = \begin{bmatrix} \delta x^t \\ \delta h^{t-1} \\ \delta r^{t-1} \end{bmatrix}$$

$$\delta \mathbf{B}^t = \delta \mathbf{z}^t = \begin{bmatrix} \delta b_g^t \\ \delta b_i^t \\ \delta b_f^t \\ \delta b_o^t \end{bmatrix}$$

Backpropagation

Forward propagation

$$\hat{h}_k, \mathbf{c}_k = LSTM(f'(\hat{x}), [h_{k-1}, r_{k-1}], \mathbf{c}_{k-1})$$

$$\hat{i}_k = W_{xi}f' + W_{hi}h_{k-1} + W_{ri}r_{k-1} + b_i$$

$$\hat{f}_k = W_{xf}f' + W_{hf}h_{k-1} + W_{rf}r_{k-1} + b_f$$

$$\hat{o}_k = W_{xo}f' + W_{ho}h_{k-1} + W_{ro}r_{k-1} + b_o$$

$$\hat{g}_k = W_{xc}f' + W_{hc}h_{k-1} + W_{rc}r_{k-1} + b_c$$

Backpropagation

Layers t:

$$z^t = W \square I^t + B$$

$$\delta W_{gx} = \delta \hat{g}^t \cdot x^t$$

$$\delta W_{ix} = \delta \hat{i}^t \cdot x^t$$

$$\delta W_{gh} = \delta \hat{g}^t \cdot h^{t-1}$$

$$\delta W_{ih} = \delta \hat{i}^t \cdot h^{t-1}$$

$$\delta W_{gr} = \delta \hat{g}^t \cdot r^{t-1}$$

$$\delta W_{ir} = \delta \hat{i}^t \cdot r^{t-1}$$

$$\delta W_{fx} = \delta \hat{f}^t \cdot x^t$$

$$\delta W_{ox} = \delta \hat{o}^t \cdot x^t$$

$$\delta W_{fh} = \delta \hat{f}^t \cdot h^{t-1}$$

$$\delta W_{oh} = \delta \hat{o}^t \cdot h^{t-1}$$

$$\delta W_{fr} = \delta \hat{f}^t \cdot r^{t-1}$$

$$\delta W_{or} = \delta \hat{o}^t \cdot r^{t-1}$$

Backpropagation

Forward propagation

$$g(\mathbf{x}_i, \mathbf{S}) = \vec{h}_i + \tilde{h}_i + g'(x_i)$$

$$\vec{h}_i, \vec{c}_i = LSTM(g'(x_i), \vec{h}_{i-1}, \vec{c}_{i-1})$$

$$\tilde{h}_i, \tilde{c}_i = LSTM(g'(x_i), \tilde{h}_{i+1}, \tilde{c}_{i+1})$$

Backpropagation

$$\delta g = \delta \vec{h}_i + \delta \tilde{h}_i + \delta g'$$

$$\delta \vec{h}_i \square \delta \tilde{h}_i = \delta(LSTM)$$

The same as above

Backpropagation

Forward propagation

Backpropagation

cnn demo

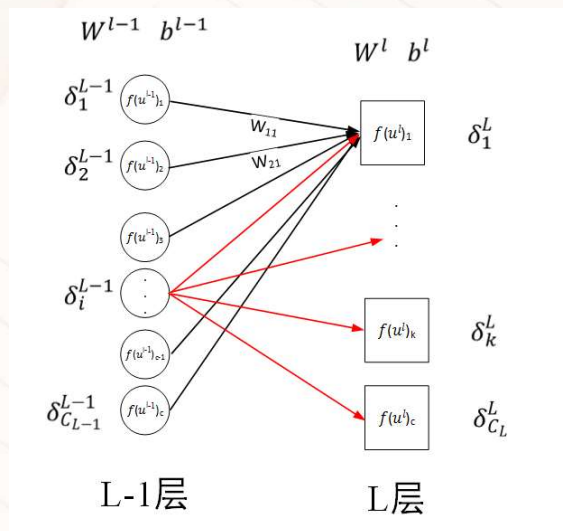
Input: g'

$$\delta_k = \delta g'$$

Backpropagation

Forward propagation

cnn demo



$$u^l = W^l f(u^{l-1}) + b^l$$

Backpropagation

$$\begin{aligned} f'(u^l) &= \left(\frac{1}{1 + e^{-u^l}} \right)' \\ &= \frac{1}{1 + e^{-u^l}} \left(1 - \frac{1}{1 + e^{-u^l}} \right) \\ &= f(u^l)(1 - f(u^l)) \end{aligned}$$

$$\delta u^l = \delta_k^l$$

Backpropagation

Forward propagation

cnn demo

Backpropagation

$$\begin{aligned}\delta_i^{l-1} &= \sum_{k=1}^{c_l} (\delta_k^{l-1}) = \sum_{k=1}^{c_l} \frac{\partial u_i^l}{\partial u_i^{l-1}} \\ &= \frac{\partial}{\partial u_i^{l-1}} \sum_{k=1}^{c_l} (Wf + b) \\ &= \sum_{k=1}^{c_l} \delta_k^l W_{ik}^l f'(u_i^{l-1})\end{aligned}$$

$$\delta^{l-1} = (W^l)^T \delta^{l.*} f'(u^{l-1})$$

Backpropagation

Forward propagation

Backpropagation

cnn demo

$$b = b - \alpha \frac{\partial}{\partial b} J(W, b) = b - \alpha * \delta^l$$

$$W = W - \alpha \frac{\partial}{\partial W} J(W, b)$$

$$\frac{\partial u}{\partial W_{ik}} = \delta_k^l * f(u_i^{l-1})$$

$$W = W - \alpha * \delta_k^l * f(u_i^{l-1})$$

One-shot learning with Matching Networks

- Background
- Introduction
 - Introduce to One Shot Learning
 - Introduce to MANN
- Model
 - Motivation
 - Matching Networks
 - Backpropagation
- Experiments
- Summary

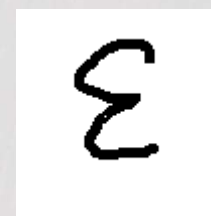
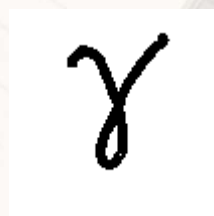
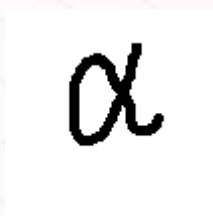
Experiments

- Data: Omniglot

Omniglot consists of 1623 characters from 50 different alphabets. Each of these was hand drawn by 20 different people. The large number of classes (characters) with relatively few data per class(20), makes this an ideal data set for testing small-scale one-shot classification.

Download:<https://github.com/brendenlake/omniglot>

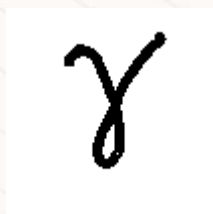
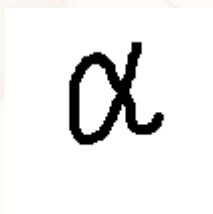
Example:



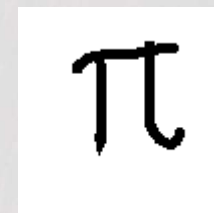
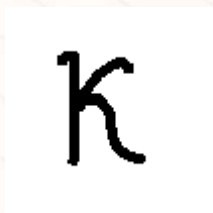
Experiments

- 5-way, 5-shot learning

Traning data set:



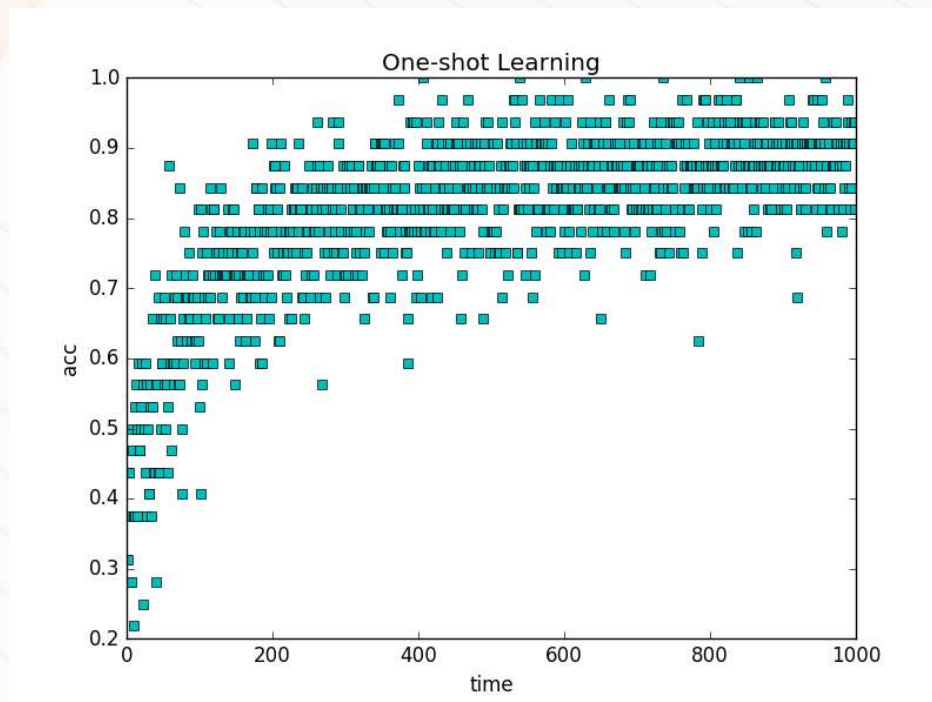
Test data set:



Experiments

■ Results

Test the relationship between accuracy and the number of bath.



Simulation environment:

- Python 3.5
- Tensorflow-1.0

Traning Input:

- Five classes image of Omniglot

Output:

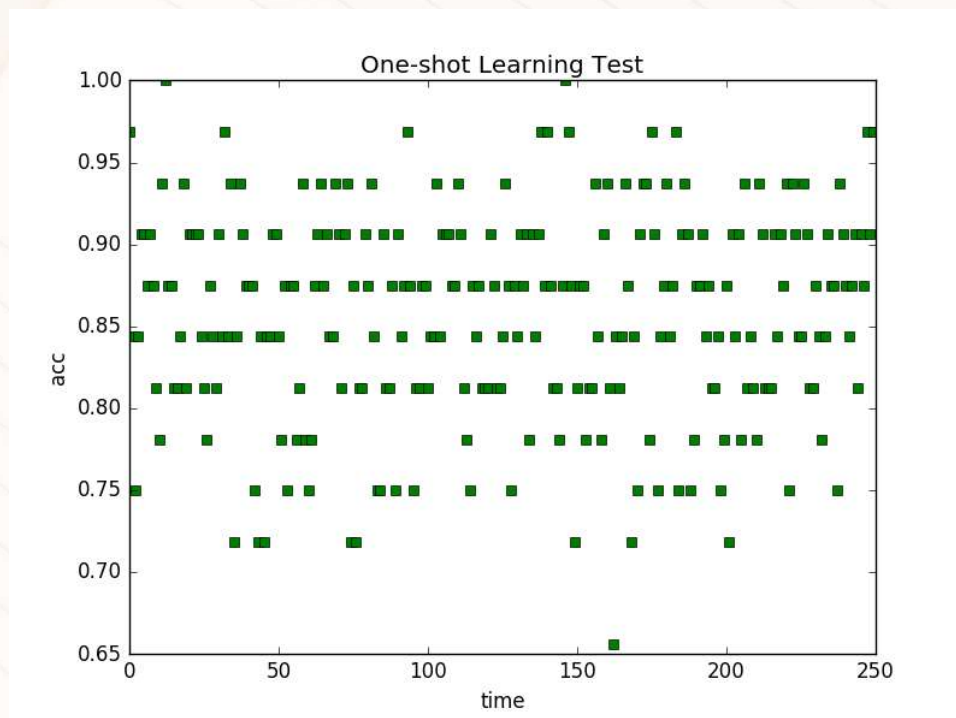
- Accuracy of test bath

Time:total_train_batches

Experiments

■ Results

Test the accuracy of one-shot model.



Simulation environment:

- Python 3.5
- Tensorflow-1.0

Test Input:

- Another five classes image of Omniglot

Output:

- Accuracy of test bath

One-shot learning with Matching Networks

- Background
- Introduction
 - Introduce to One Shot Learning
 - Introduce to MANN
- Model
 - Motivation
 - Matching Networks
 - Backpropagation
- Experiments
- Summary

Summary

- One-shot learning learns the mapping function between input image and memory
- One-shot learning search the memory information for training
- Matching Network has non-parametric structure, thus has ability to acquisition of new examples rapidly