

Modelo

Matching Network

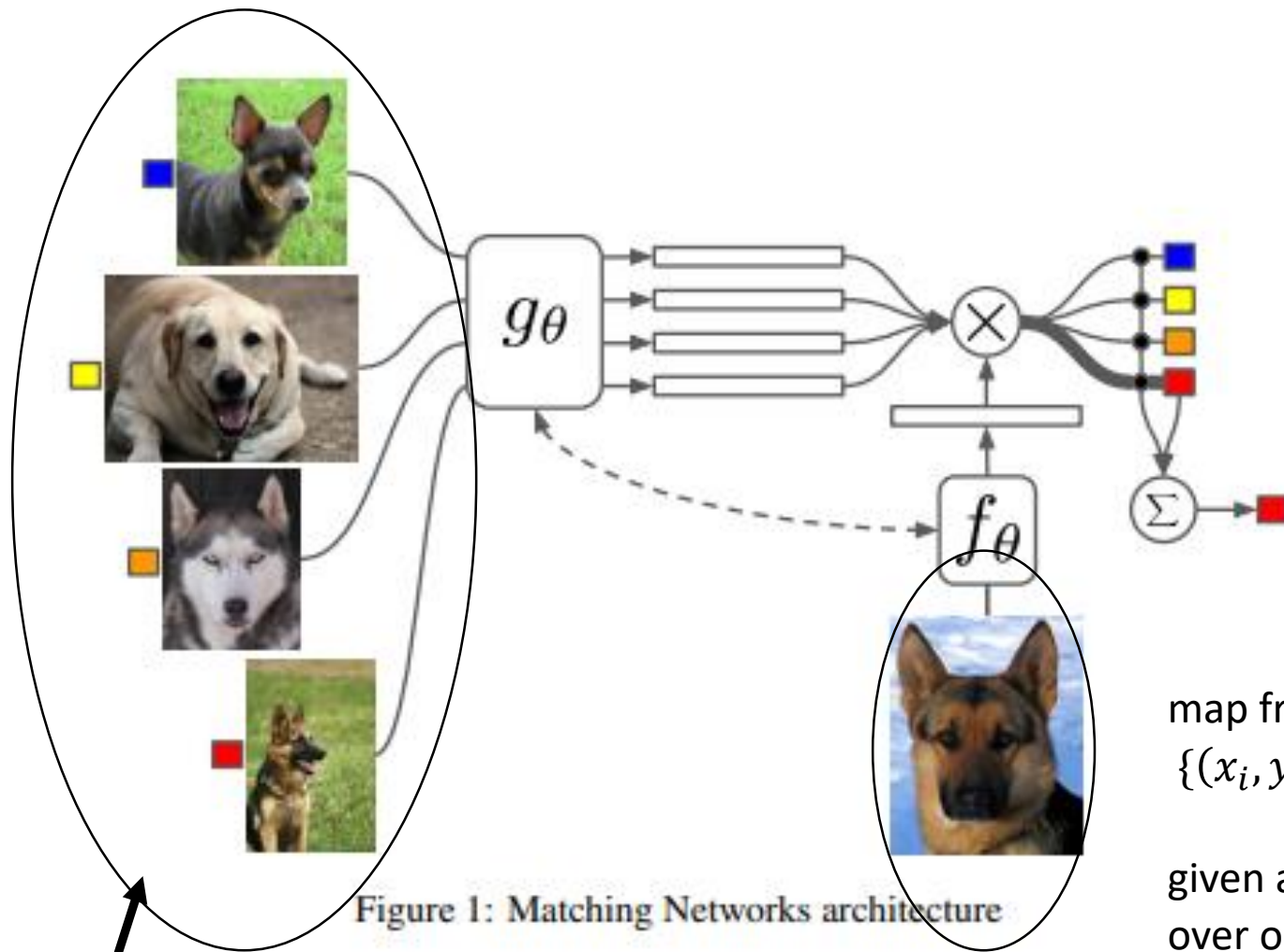


Figure 1: Matching Networks architecture

- Uses recent advances in NN with augmented memory.
- Training procedure is based on simple ML principle: test and train conditions must match (Showing only a few examples per class).

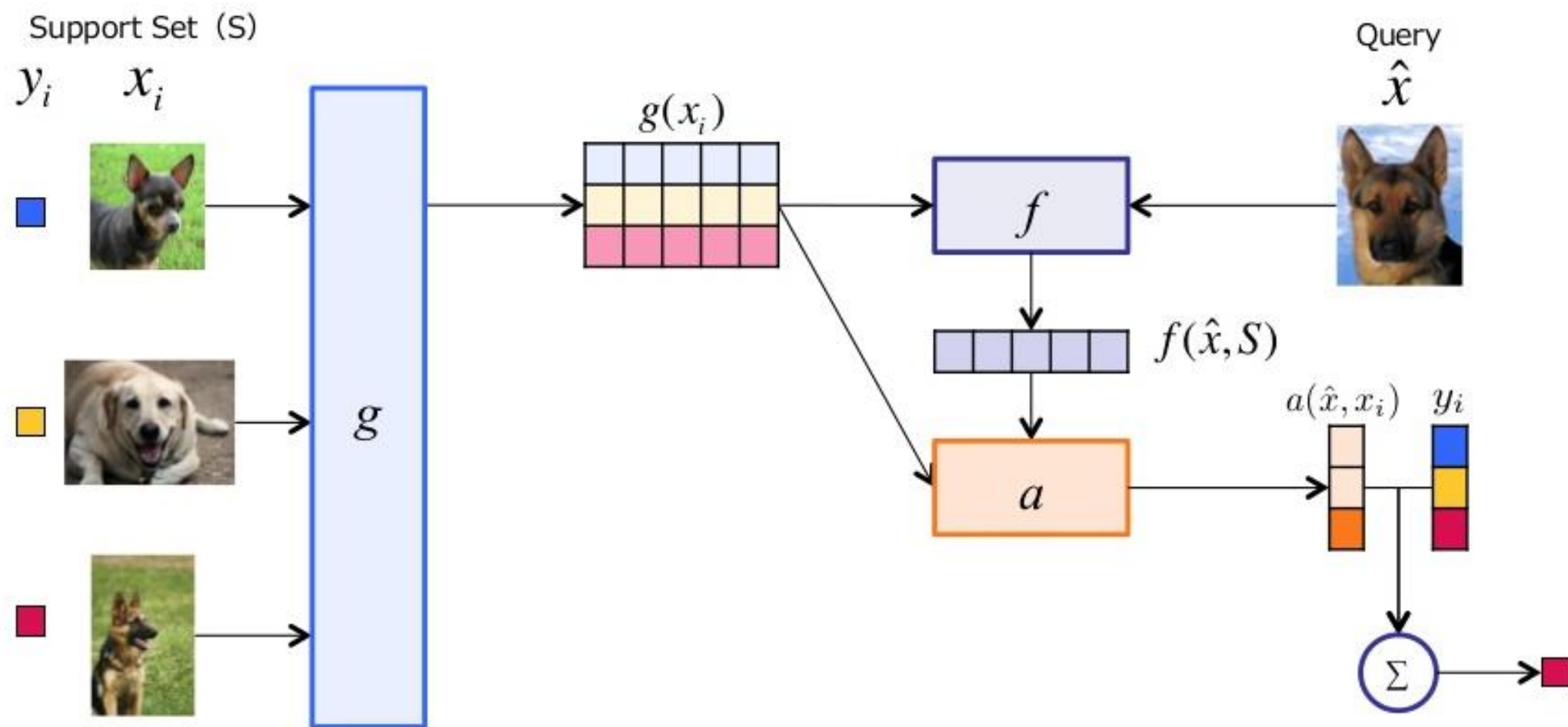
map from a S of k examples of image-label pairs $S = \{(x_i, y_i)\}_{i=1}^k$ to a classifier $C_S(\hat{x})$

given a test example (\hat{x}) defines a probability distribution over outputs (\hat{Y})

We define the mapping $S \rightarrow C_S(\hat{x})$ to be $P(\hat{Y}|\hat{x}, S)$ where P is parameterised by a neural network.

• Support set (S)

• Query set (S)



Our model in its simplest form computes \hat{y} as follows:

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

the attention mechanism **a** is a kernel on $X \times X$, then is akin to a **kernel density estimator**

describes the output for a new class as a linear combination of the labels in the support set.

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

■ The Attention Kernel

- Calculate softmax over the cosine distance between $f(\hat{x}, S)$ and $g(x_i)$
 - Similar to nearest neighbor calculation
- Train a network using cross entropy loss

embedding functions **f** and **g** being appropriate neural networks (potentially **with f = g**) to embed \hat{x} and x_i .

Loss

- This kind of loss is also related to methods such as **Neighborhood Component Analysis (NCA)**, **triplet loss** or **large margin nearest neighbor (Investigar)**.

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}},$$

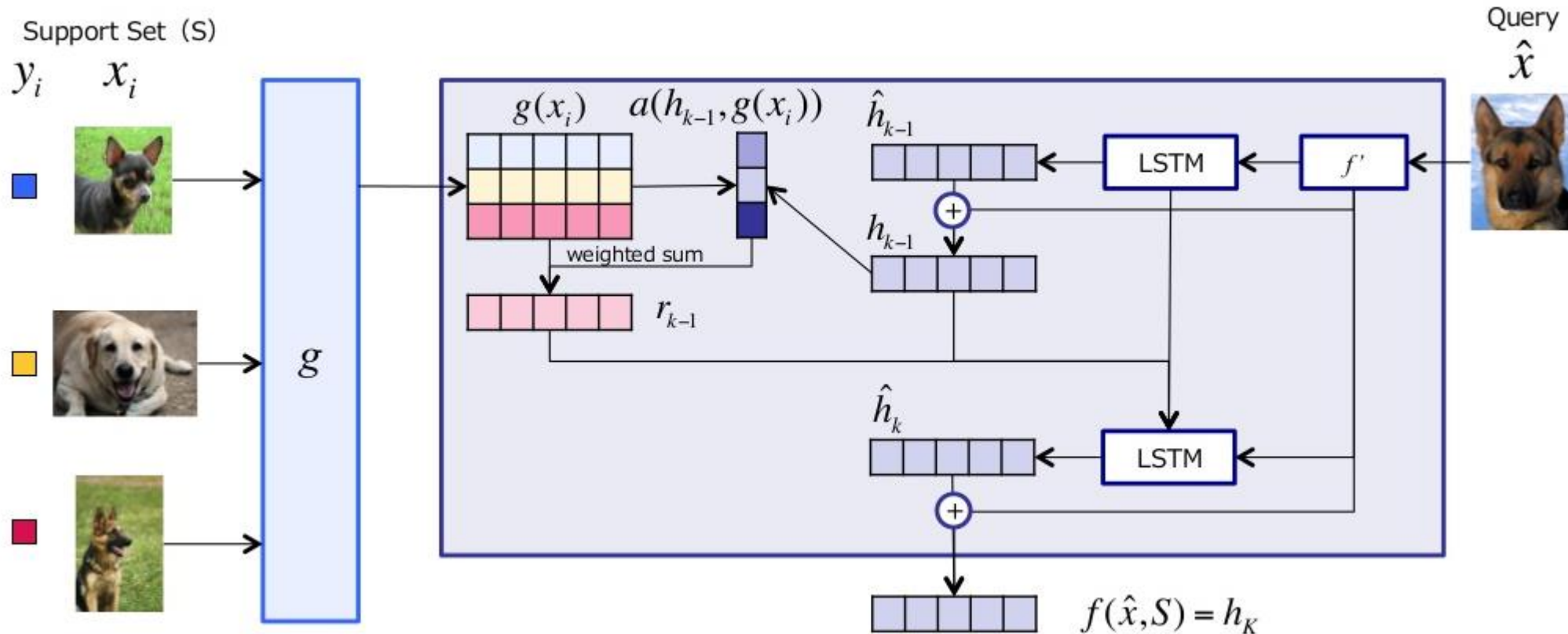
Full Context Embeddings f

$$f(\hat{x}, S) = \text{attLSTM}(f'(\hat{x}), g(S), K)$$

where f' is a neural network (e.g., VGG or Inception, as described in the main text). We define K to be the number of “processing” steps following work from [26] from their “Process” block. $g(S)$ represents the embedding function g applied to each element x_i from the set S .

K steps of “reads”, $\text{attLSTM}(f'(\hat{x}), g(S), K) = h_K$

- a “content” based attention
- the softmax $g(x_i)$.
- The read-out h_{k-1} from $g(S)$ is concatenated to h_{k-1} .



$$\hat{h}_k, c_k = \text{LSTM}(f'(\hat{x}), [h_{k-1}, r_{k-1}], c_{k-1})$$

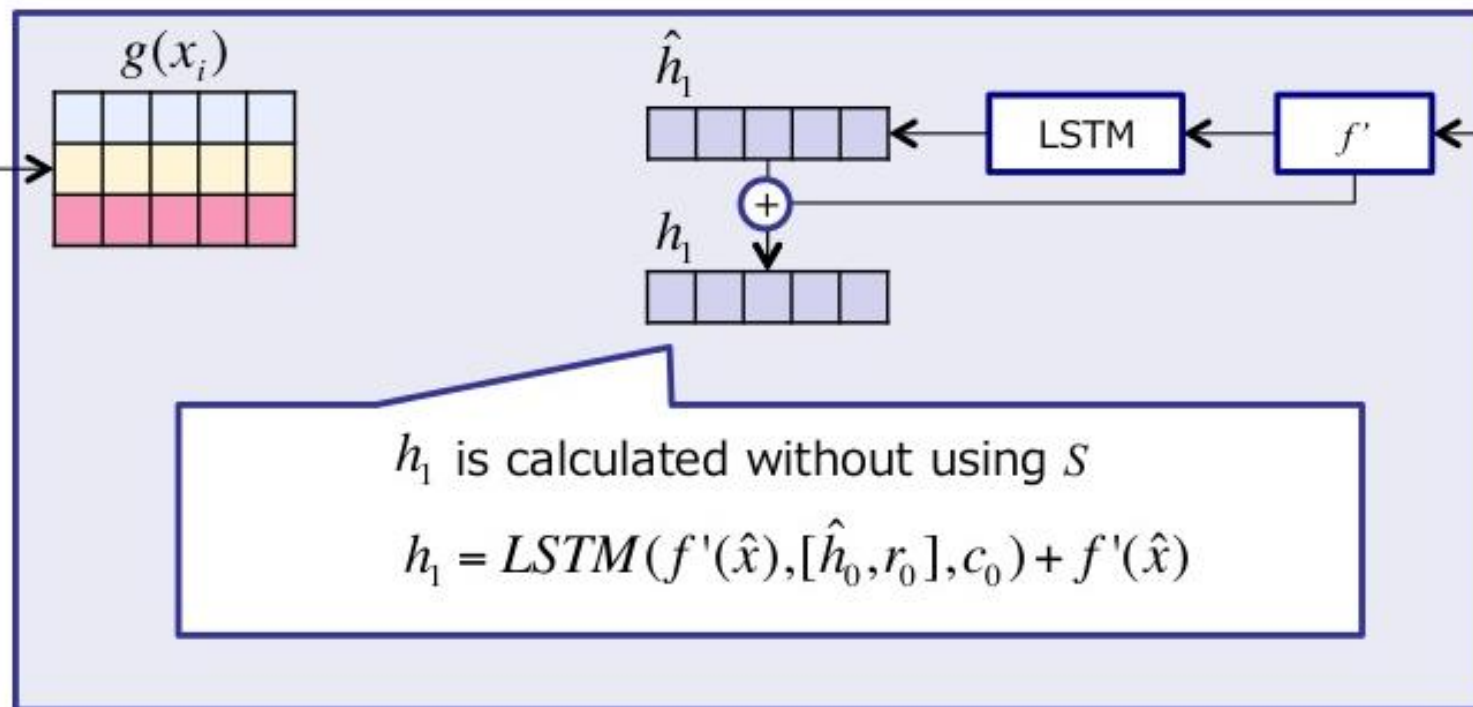
$$h_k = \hat{h}_k + f'(\hat{x})$$

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

$$a(h_{k-1}, g(x_i)) = \text{softmax}(h_{k-1}^T g(x_i))$$

Support Set (S)

y_i x_i



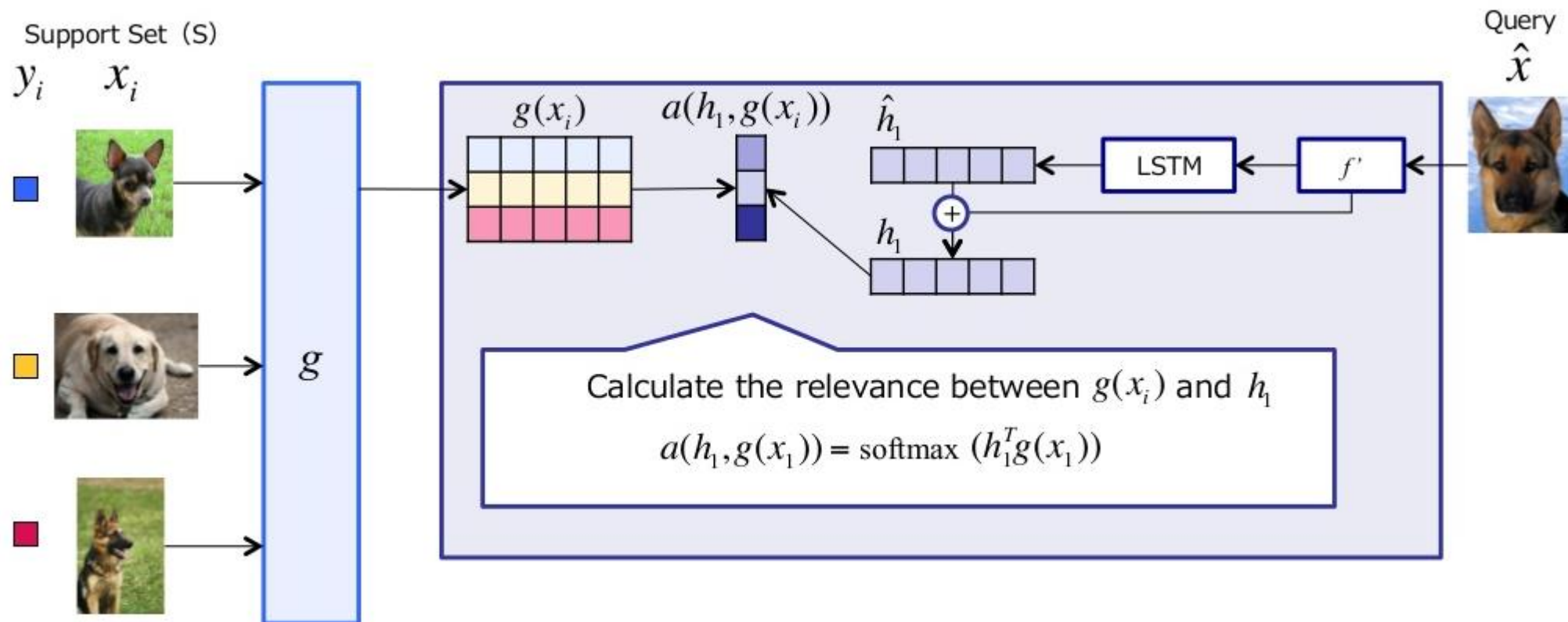
Query

\hat{x}



$$\begin{aligned} \hat{h}_k, c_k &= \text{LSTM}(f'(\hat{x}), [h_{k-1}, r_{k-1}], c_{k-1}) \\ h_k &= \hat{h}_k + f'(\hat{x}) \end{aligned}$$

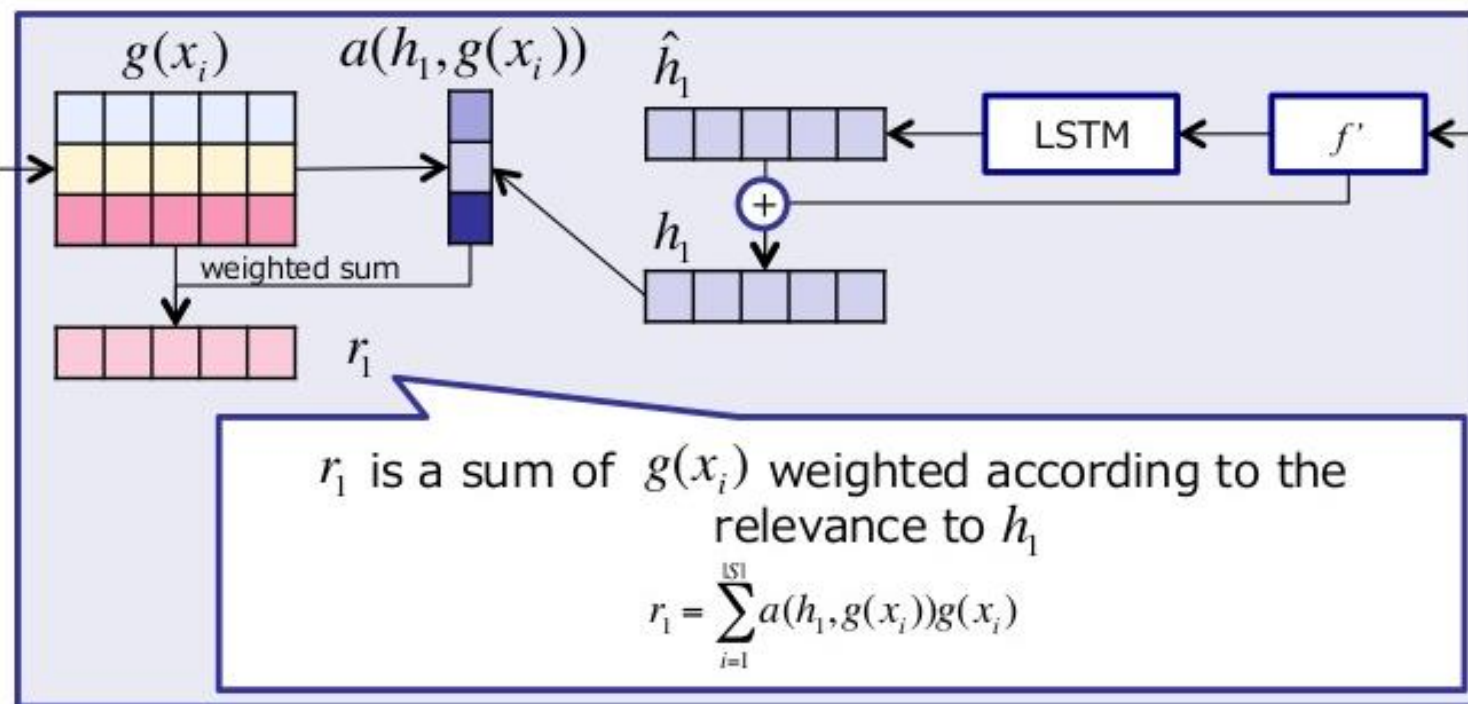
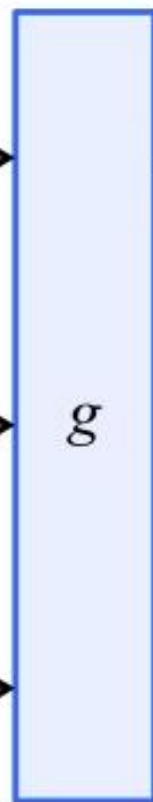
$$r_{k-1} = \sum_{i=1}^{|S|} a(h_{k-1}, g(x_i)) g(x_i)$$



$$a(h_{k-1}, g(x_i)) = \text{softmax}(h_{k-1}^T g(x_i))$$

Support Set (S)

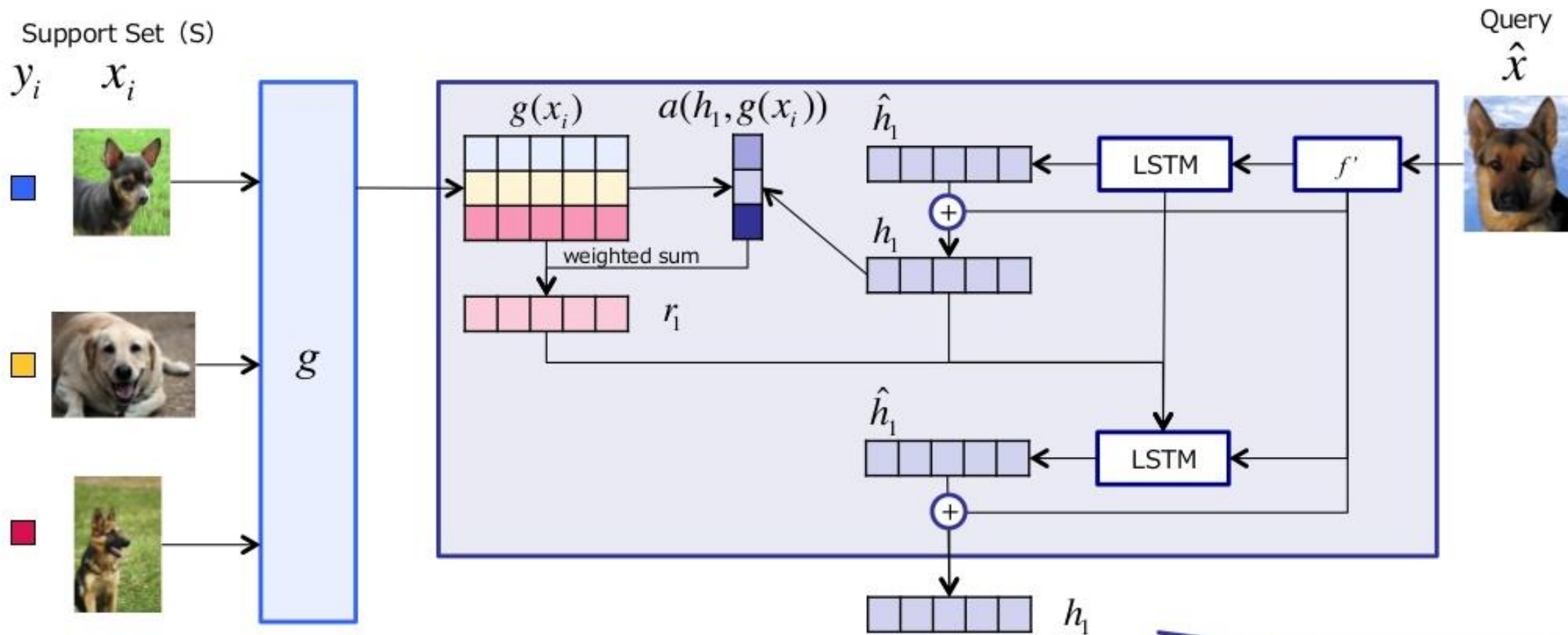
y_i x_i



Query

\hat{x}

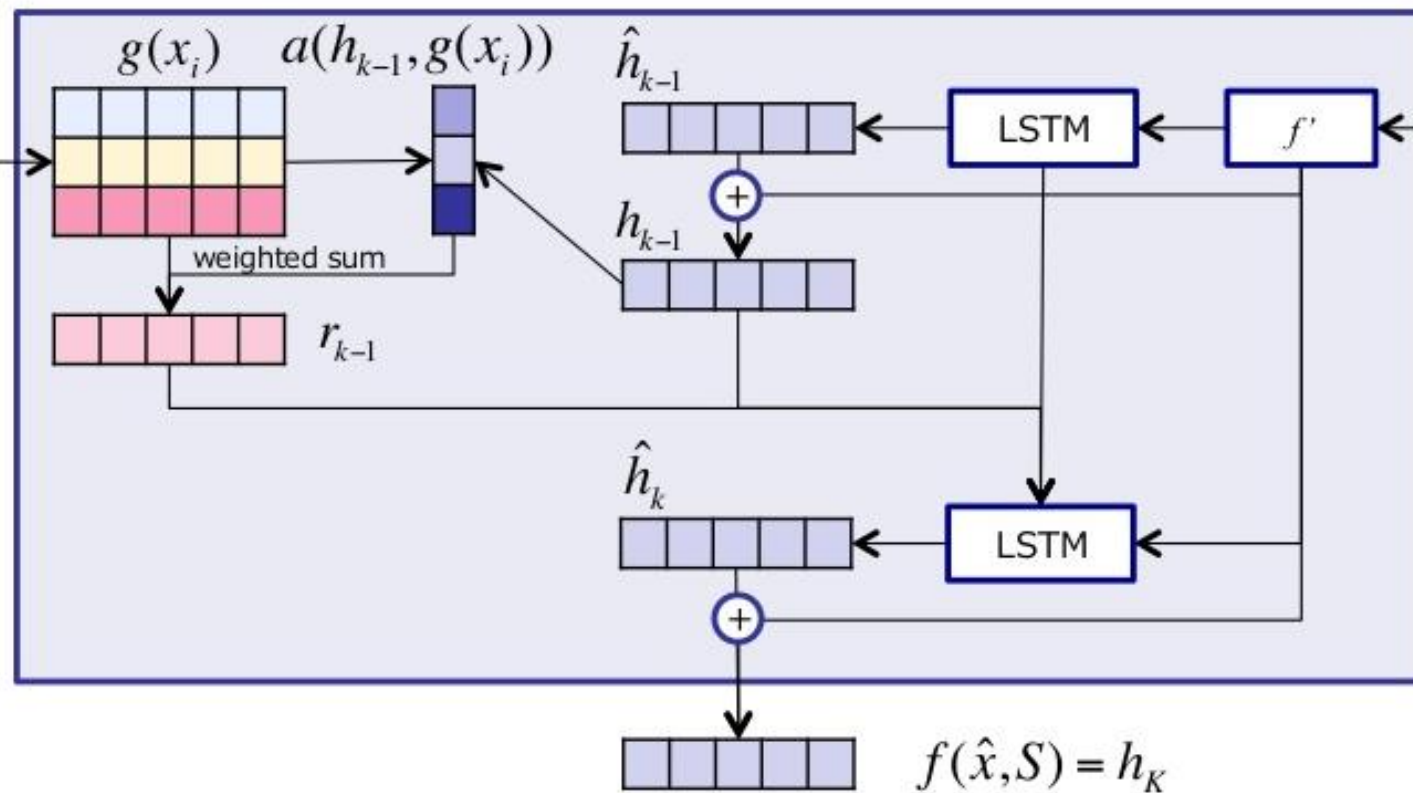
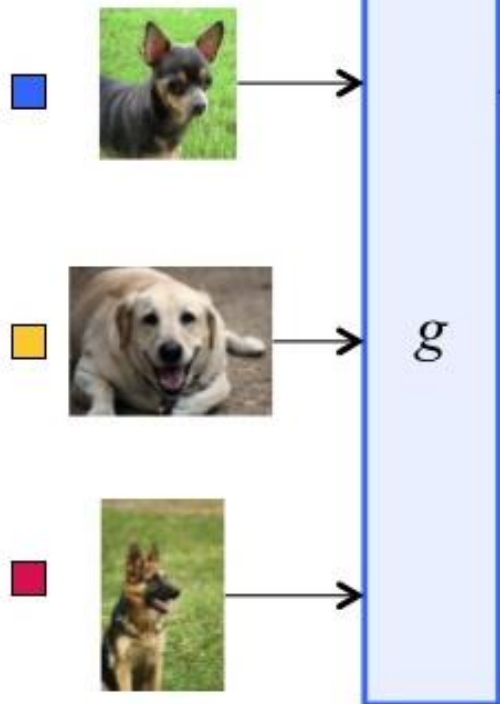




$$\begin{aligned} \hat{h}_k, c_k &= \text{LSTM}(f'(\hat{x}), [h_{k-1}, r_{k-1}], c_{k-1}) \\ h_k &= \hat{h}_k + f'(\hat{x}) \end{aligned}$$

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

h_1 is calculated using S

$$y_i \quad x_i$$


Query
 \hat{x}



Let $f(\hat{x}, S)$ be the output
after K steps

Full Context Embeddings g

- The Fully Conditional Embedding g
 - Embed x_i in consideration of S

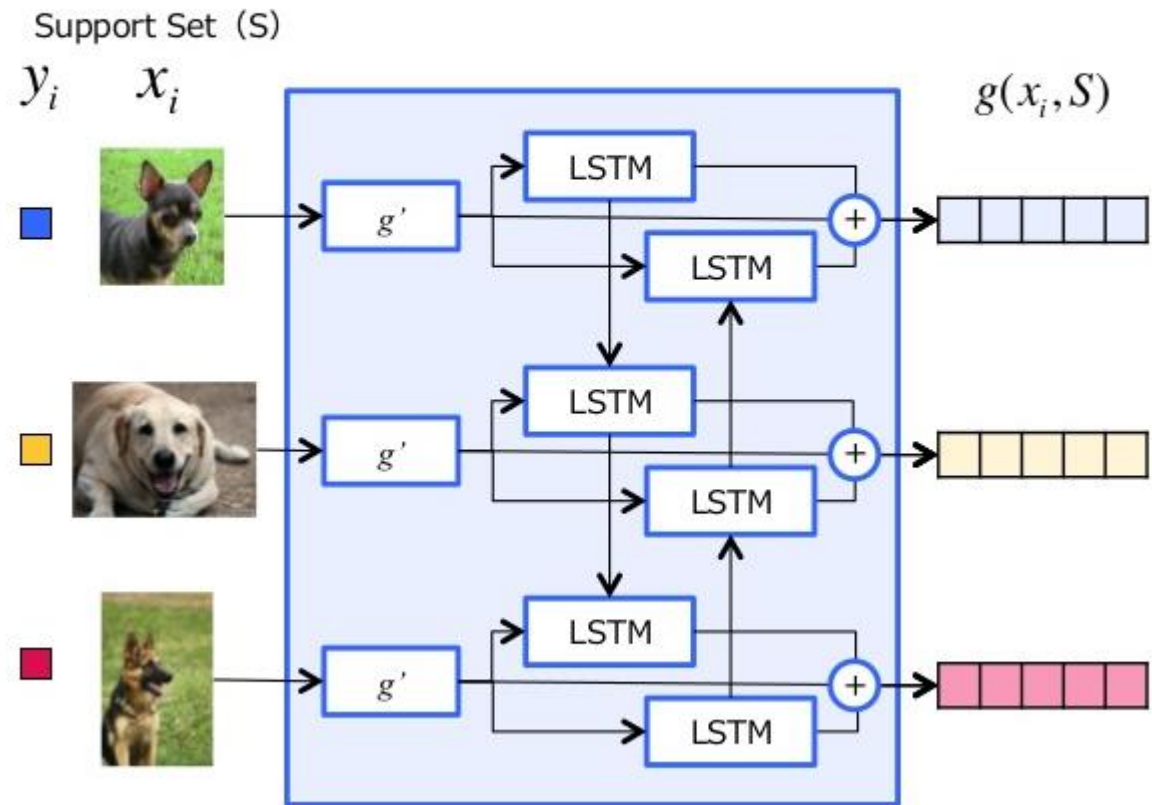
Then we define $g(x_i, S) = \vec{h}_i + \tilde{h}_i + g'(x_i)$ with:

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

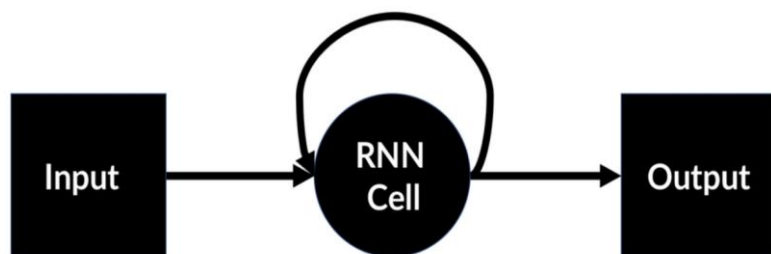
$$\tilde{h}_i, \tilde{c}_i = \text{LSTM}(g'(x_i), \tilde{h}_{i+1}, \tilde{c}_{i+1})$$

g' : neural network (e.g., VGG or Inception)

\overleftarrow{h} starts from $i = |S|$

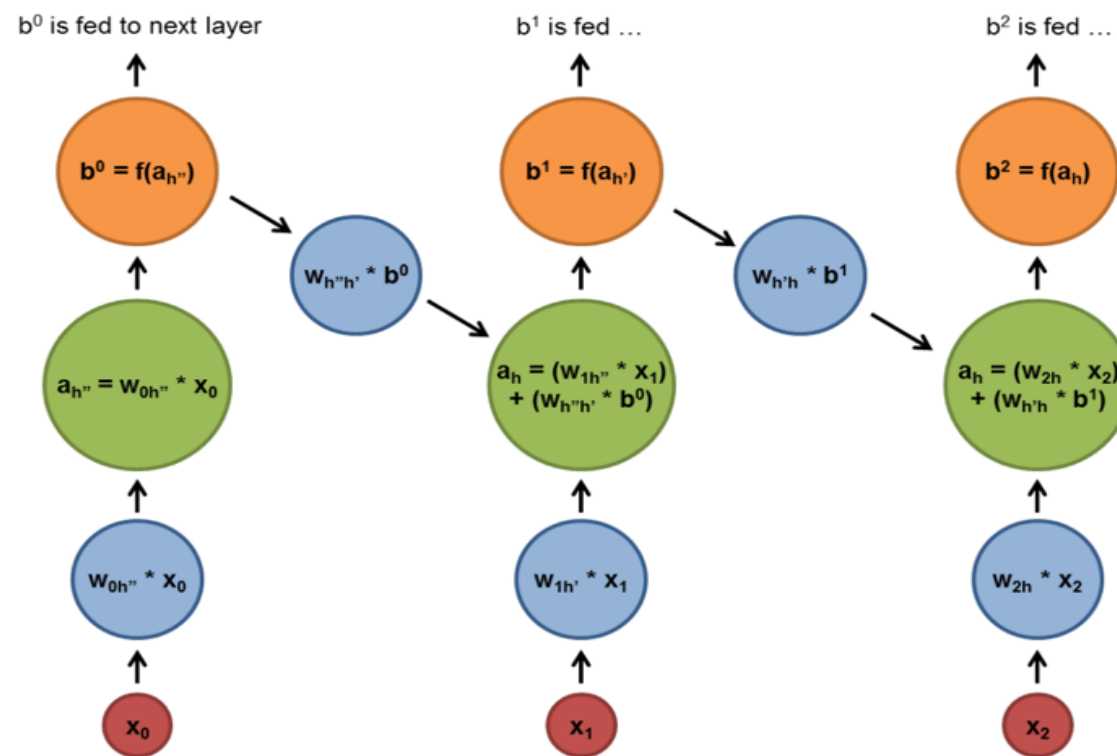


RNN



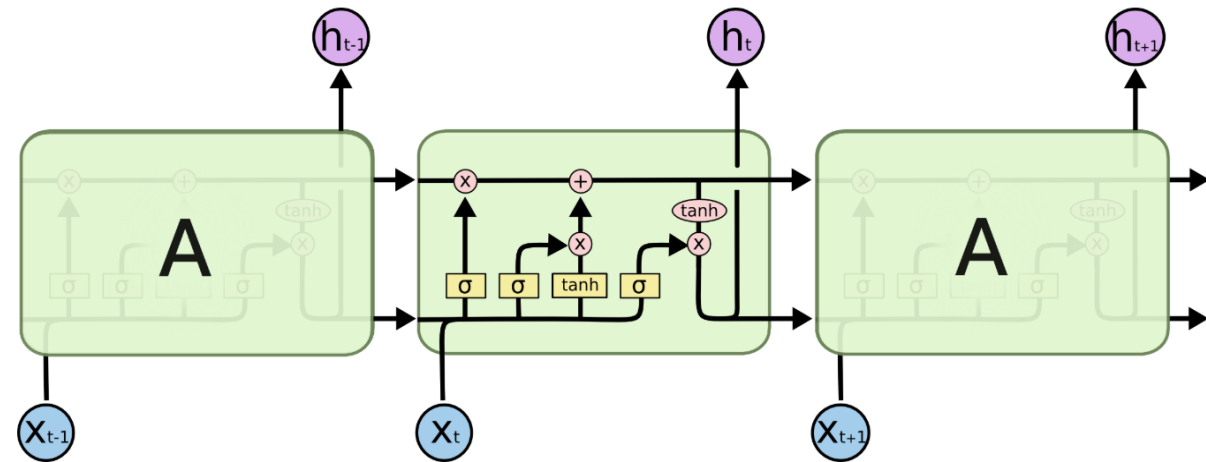
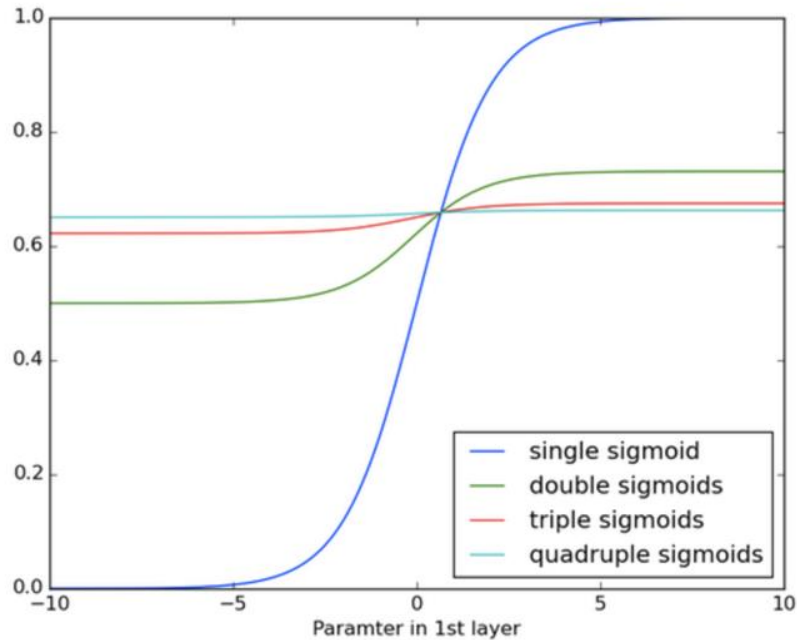
The scheme of an RNN

$$\mathbf{h}_t = \phi(W\mathbf{x}_t + U\mathbf{h}_{t-1})$$

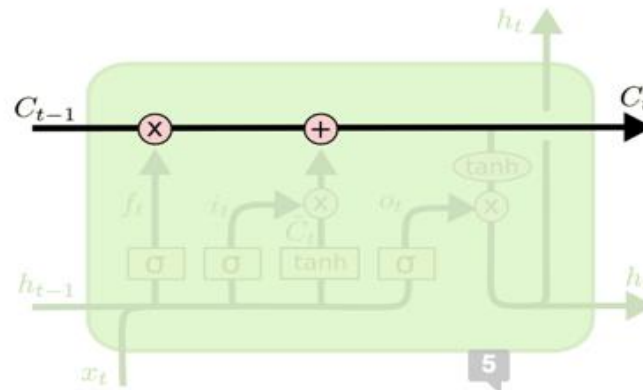


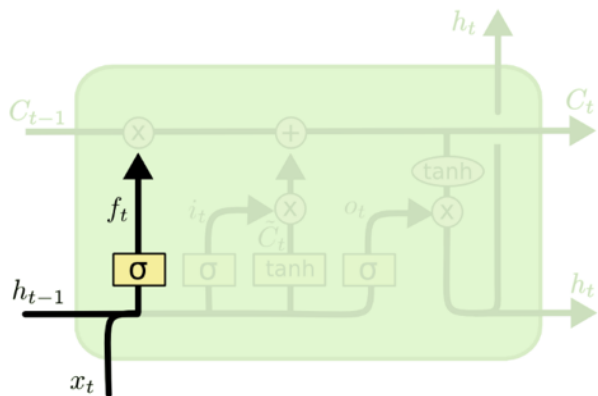
RNN LSTM

In direct response to the vanishing gradients problem of simple RNNs, the **Long Short-Term Memory (LSTM)** layer was invented. This layer performs much better at longer time series.

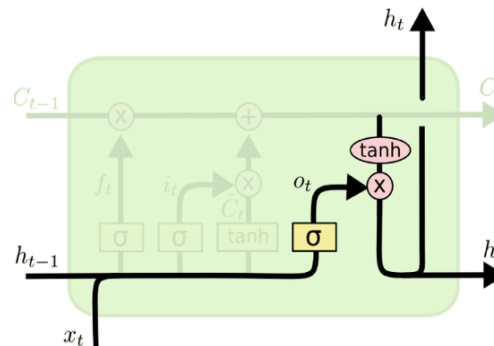


The repeating module in an LSTM contains four interacting layers.



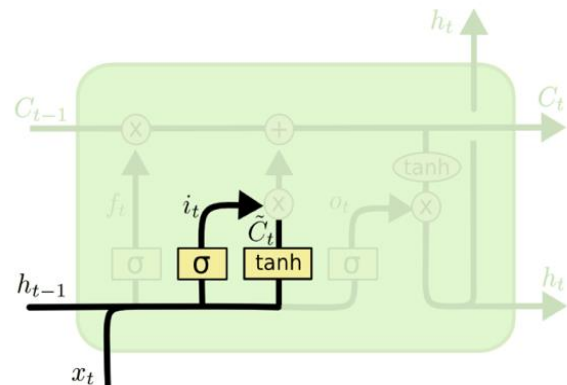


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$



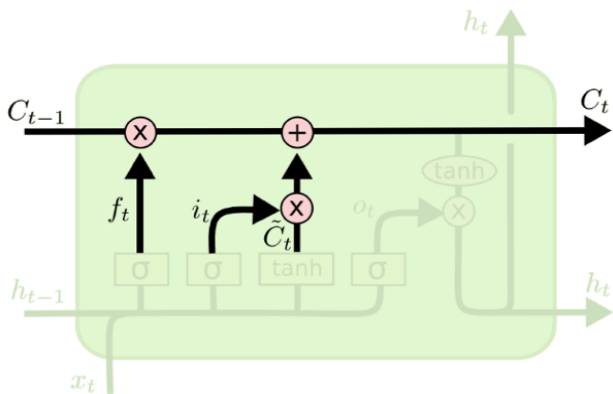
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



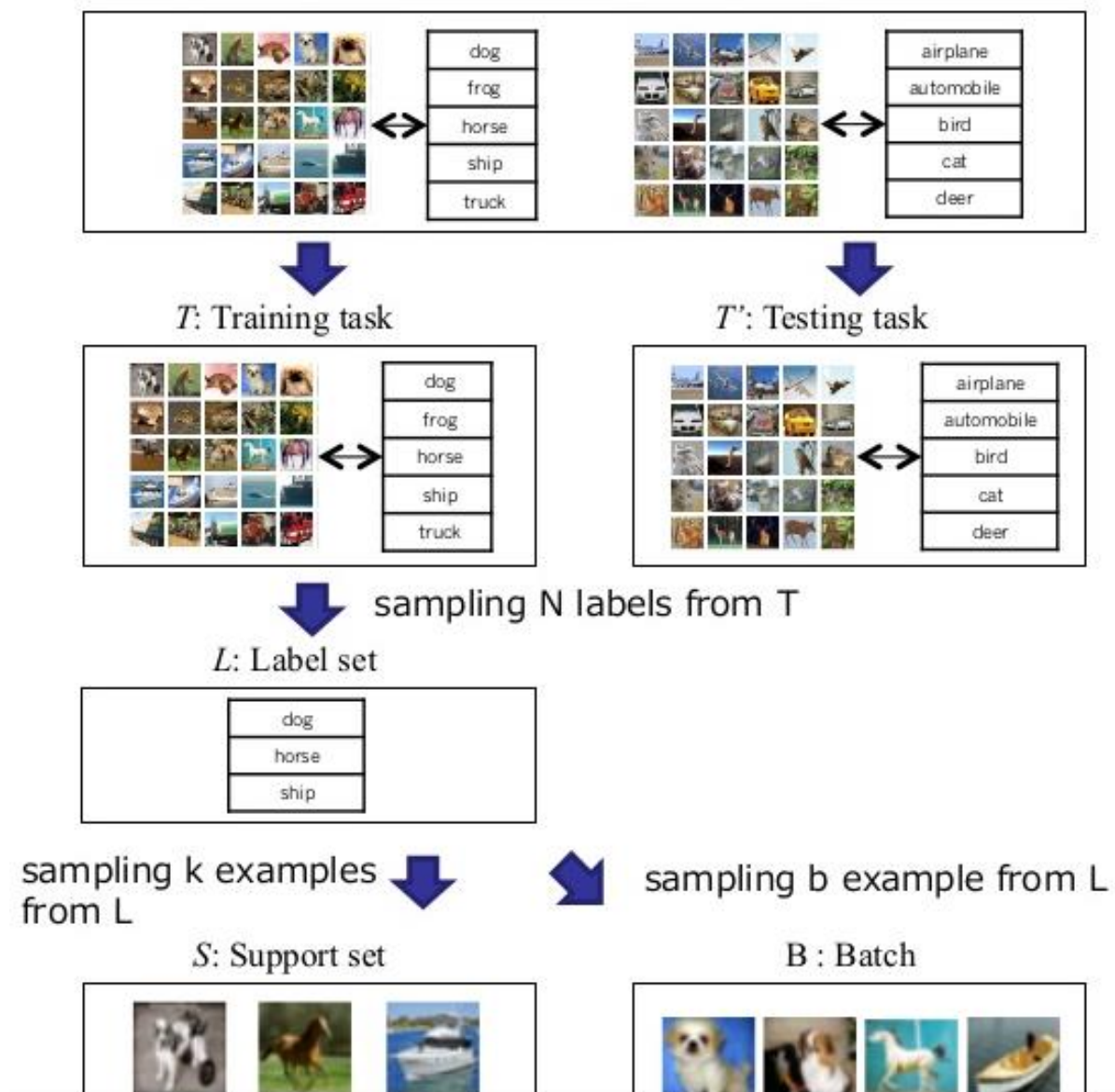
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Training Strategy

- Our model has to perform well with support sets \mathcal{S} which contain classes never seen during training.
- \mathbf{T} as distribution over possible label sets \mathbf{L} .
- \mathbf{T} to uniformly weight all data sets of up to a few unique classes (e.g., 5) with a few examples per class (e.g., up to 5)
- A label set \mathbf{L} sampled from a task \mathbf{T} , $\mathbf{L} \sim \mathbf{T}$, will typically have 5 to 25 examples.
- “episode” - compute gradients and update our model
- sample \mathbf{L} from \mathbf{T} (e.g., \mathbf{L} could be the label set **{cats, dogs}**). Use \mathbf{L} to sample the support set \mathbf{S} and a batch \mathbf{B} (i.e., both \mathbf{S} and \mathbf{B} are labelled examples of cats and dogs).
- The Matching Net is then trained to minimise the **error predicting the labels in the batch \mathbf{B} conditioned on the support set \mathbf{S} .**

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

■ Train a classifier through one-shot learning



Results

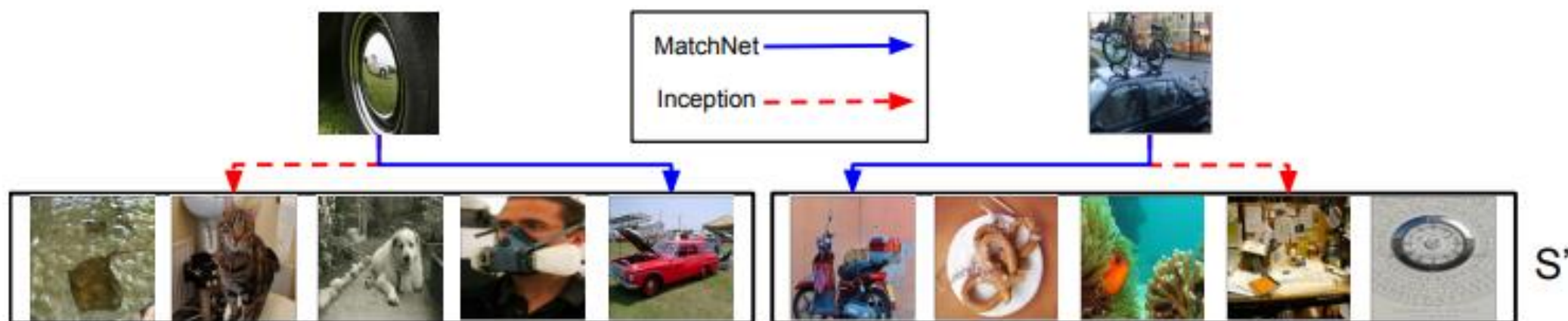


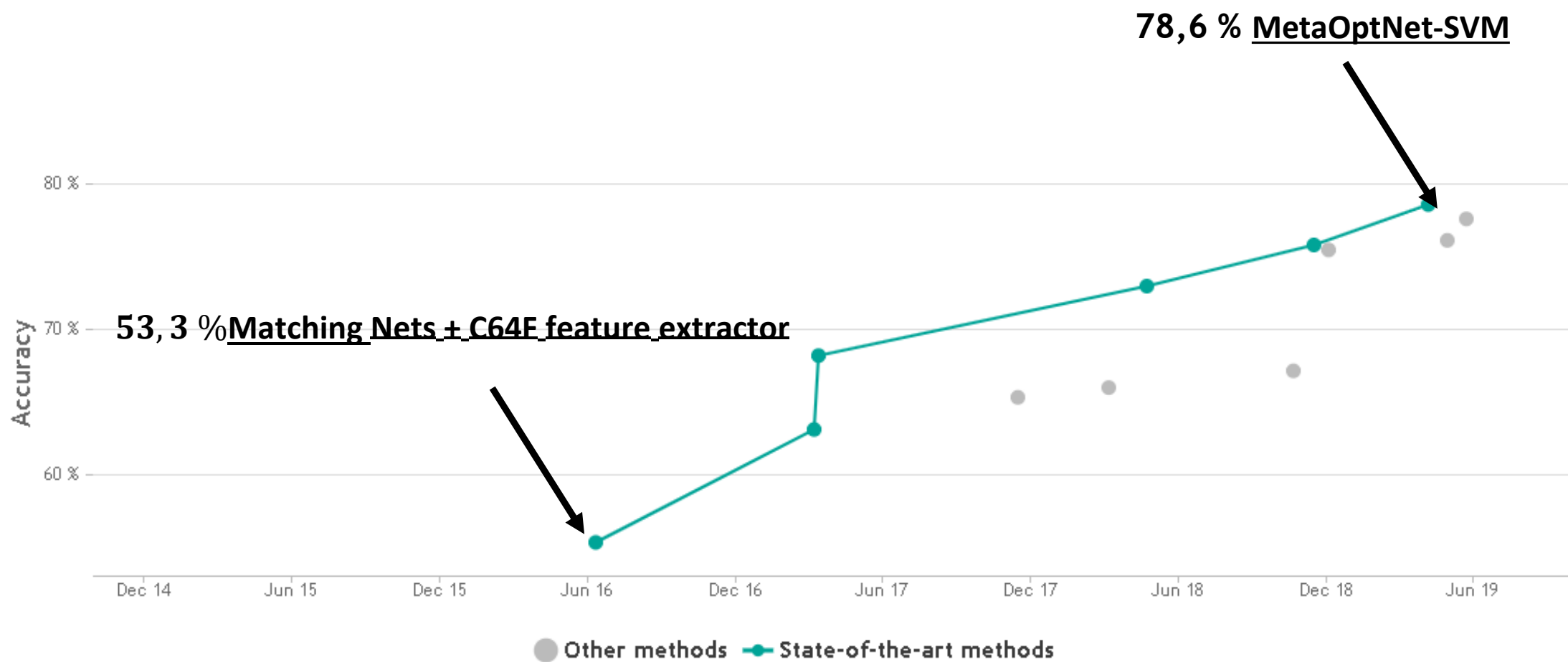
Figure 2: Example of two 5-way problem instance on ImageNet. The images in the set S' contain classes never seen during training. Our model makes far less mistakes than the Inception baseline.

Table 2: Results on *miniImageNet*.

Model	Matching Fn	Fine Tune	5-way Acc	
			1-shot	5-shot
PIXELS	Cosine	N	23.0%	26.6%
BASELINE CLASSIFIER	Cosine	N	36.6%	46.0%
BASELINE CLASSIFIER	Cosine	Y	36.2%	52.2%
BASELINE CLASSIFIER	Softmax	Y	38.4%	51.2%
MATCHING NETS (OURS)	Cosine	N	41.2%	56.2%
MATCHING NETS (OURS)	Cosine	Y	42.4%	58.0%
MATCHING NETS (OURS)	Cosine (FCE)	N	44.2%	57.0%
MATCHING NETS (OURS)	Cosine (FCE)	Y	46.6%	60.0%

Table 3: Results on full ImageNet on *rand* and *dogs* one-shot tasks. Note that $\neq L_{rand}$ and $\neq L_{dogs}$ are sets of classes which are seen during training, but are provided for completeness.

Model	Matching Fn	Fine Tune	ImageNet 5-way 1-shot Acc			
			L_{rand}	$\neq L_{rand}$	L_{dogs}	$\neq L_{dogs}$
PIXELS	Cosine	N	42.0%	42.8%	41.4%	43.0%
INCEPTION CLASSIFIER	Cosine	N	87.6%	92.6%	59.8%	90.0%
MATCHING NETS (OURS)	Cosine (FCE)	N	93.2%	97.0%	58.8%	96.4%
INCEPTION ORACLE	Softmax (Full)	Y (Full)	$\approx 99\%$	$\approx 99\%$	$\approx 99\%$	$\approx 99\%$



Codigo

- <https://github.com/AntreasAntoniou/MatchingNetworks>