

# NAS **via** RNN + RL

**Shusen Wang**

<http://wangshusen.github.io/>

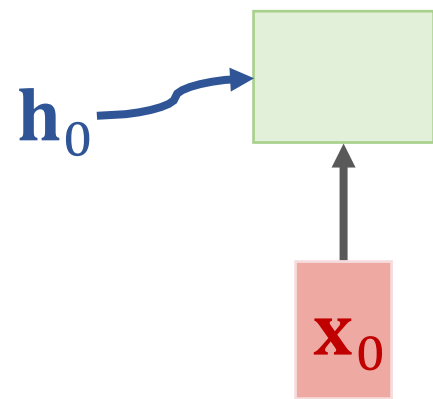
# Prerequisites

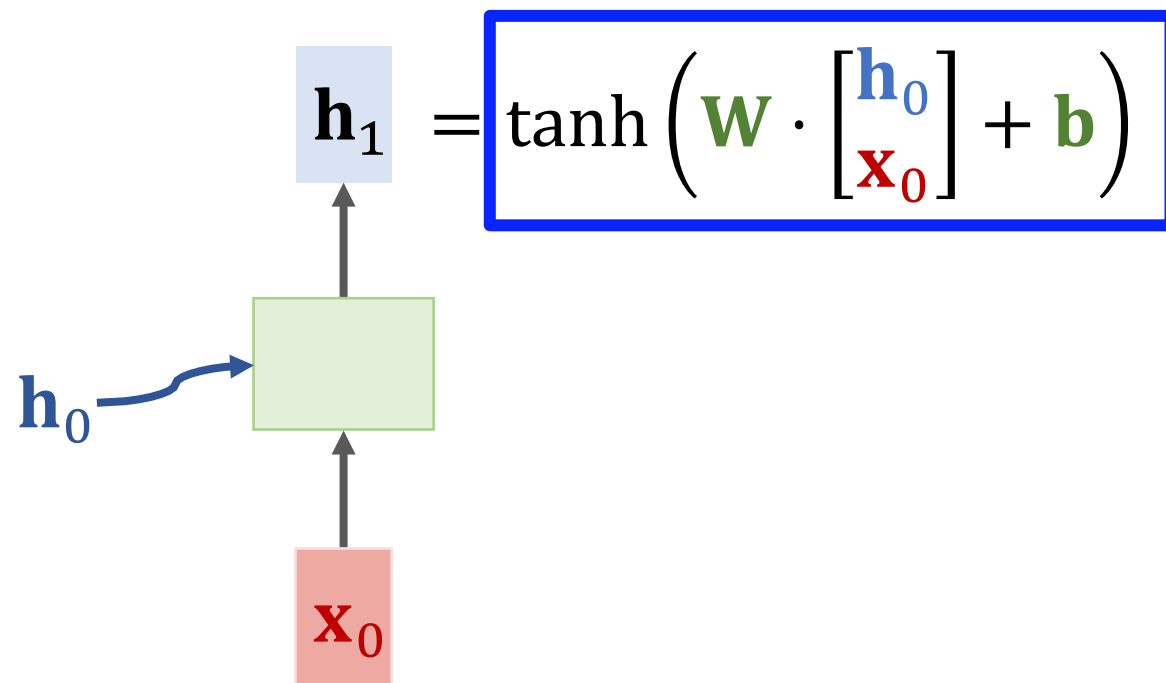
- Recurrent neural networks (RNNs).
- Policy-based reinforcement learning.

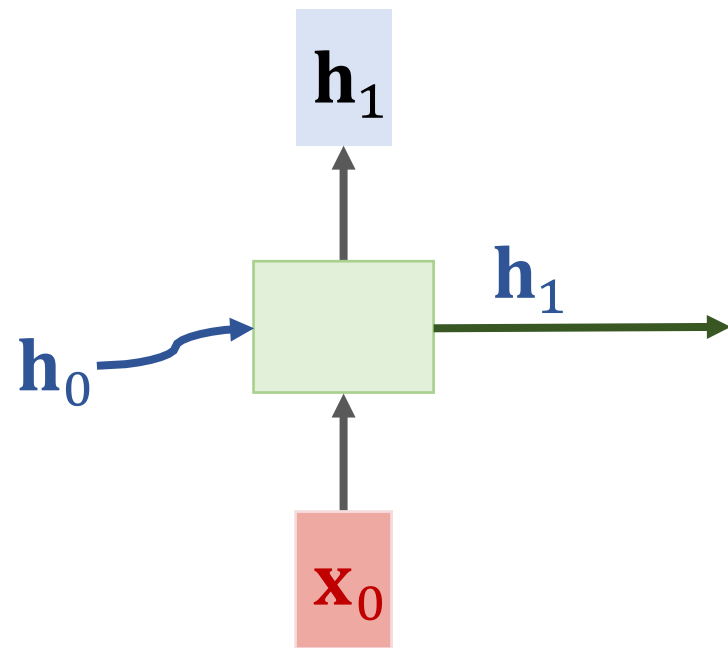
# RNN for Generating CNN Architectures

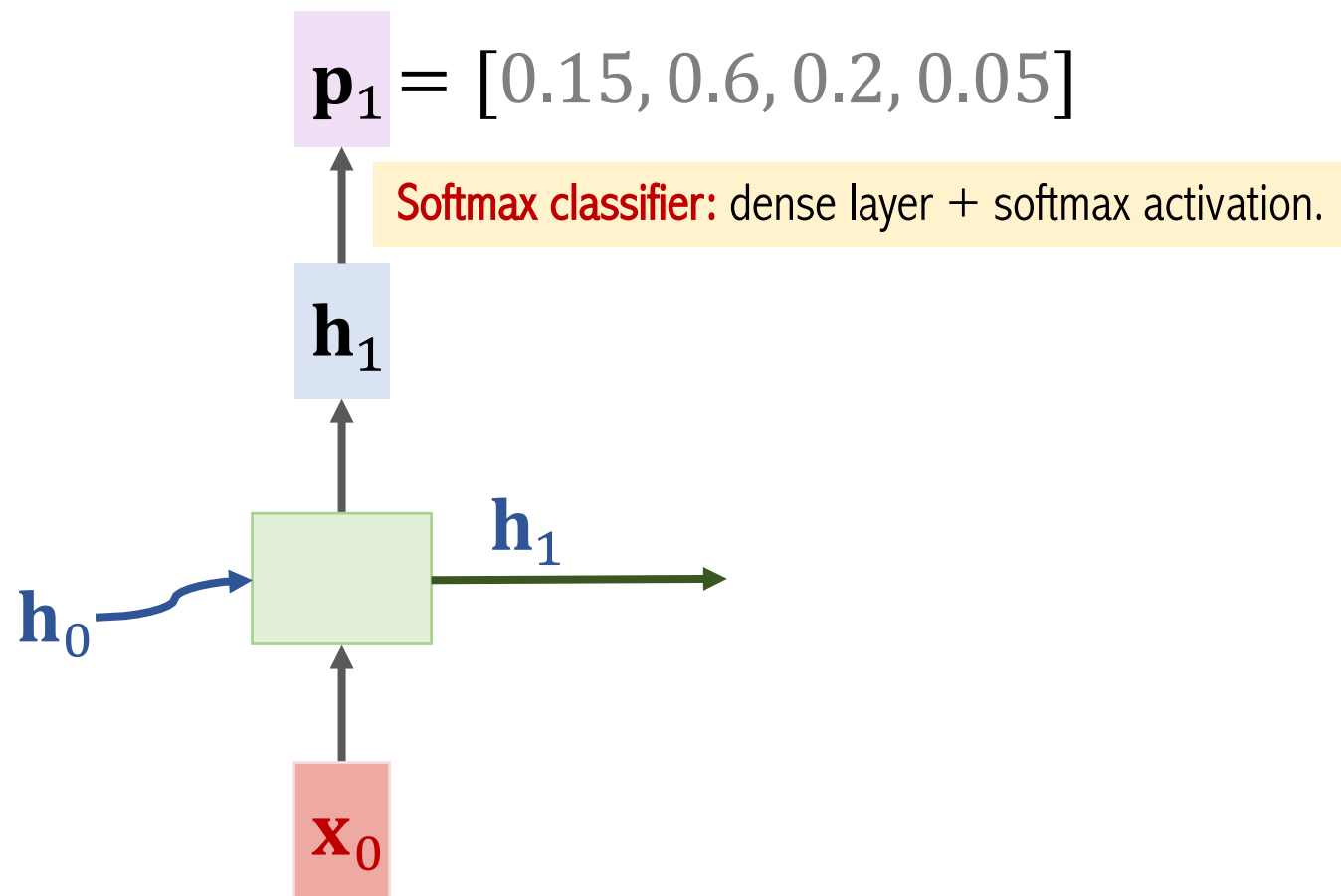
## Reference:

- Zoph & Le. [Neural architecture search with reinforcement learning](#). In *ICLR*, 2017.

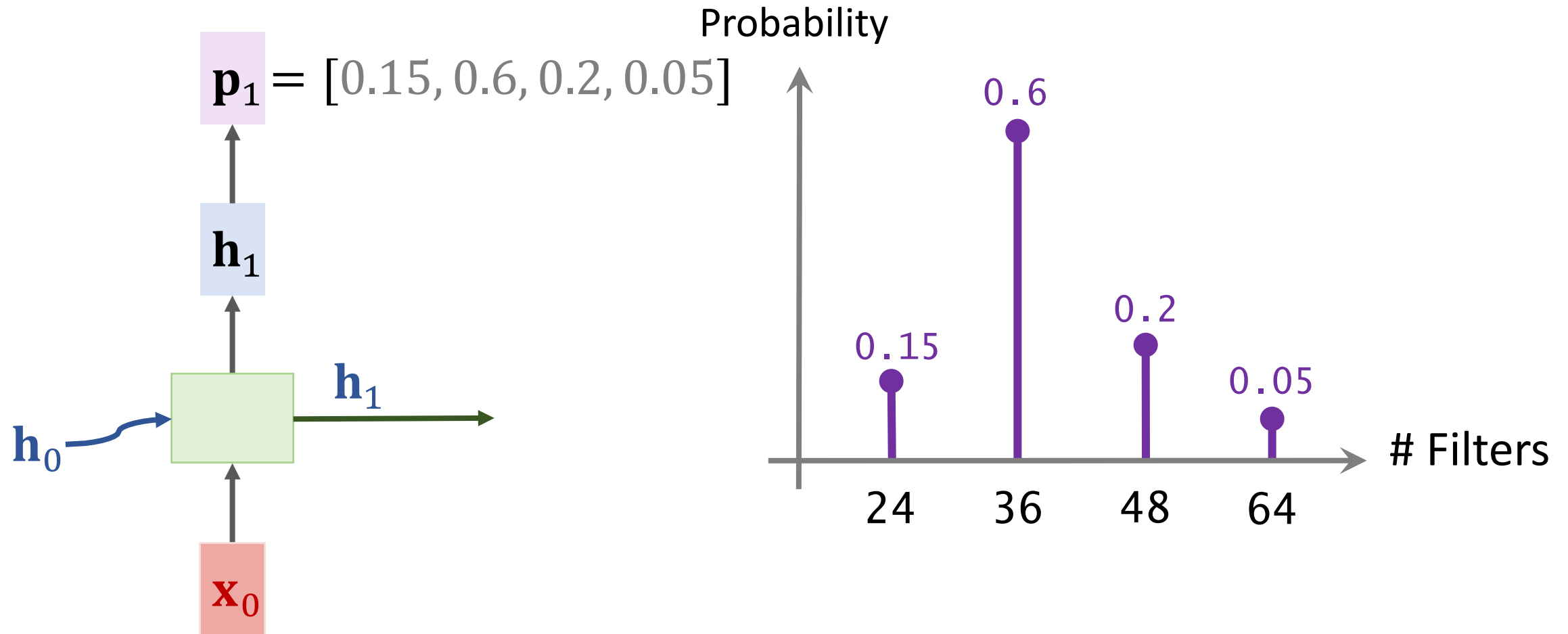






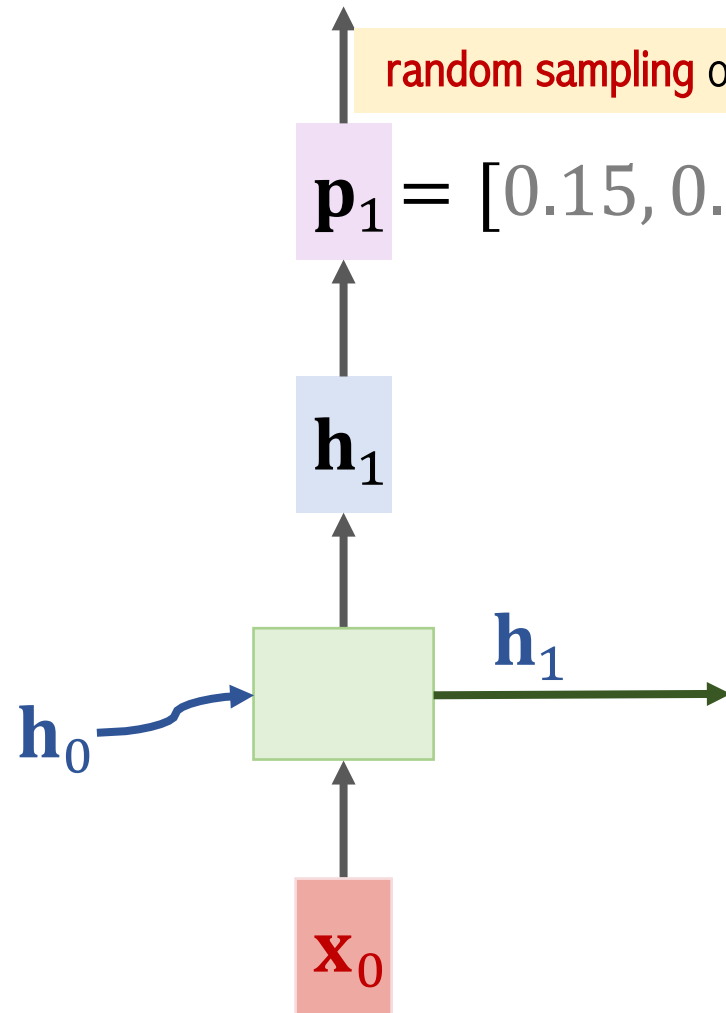


## Predict Number of Filters

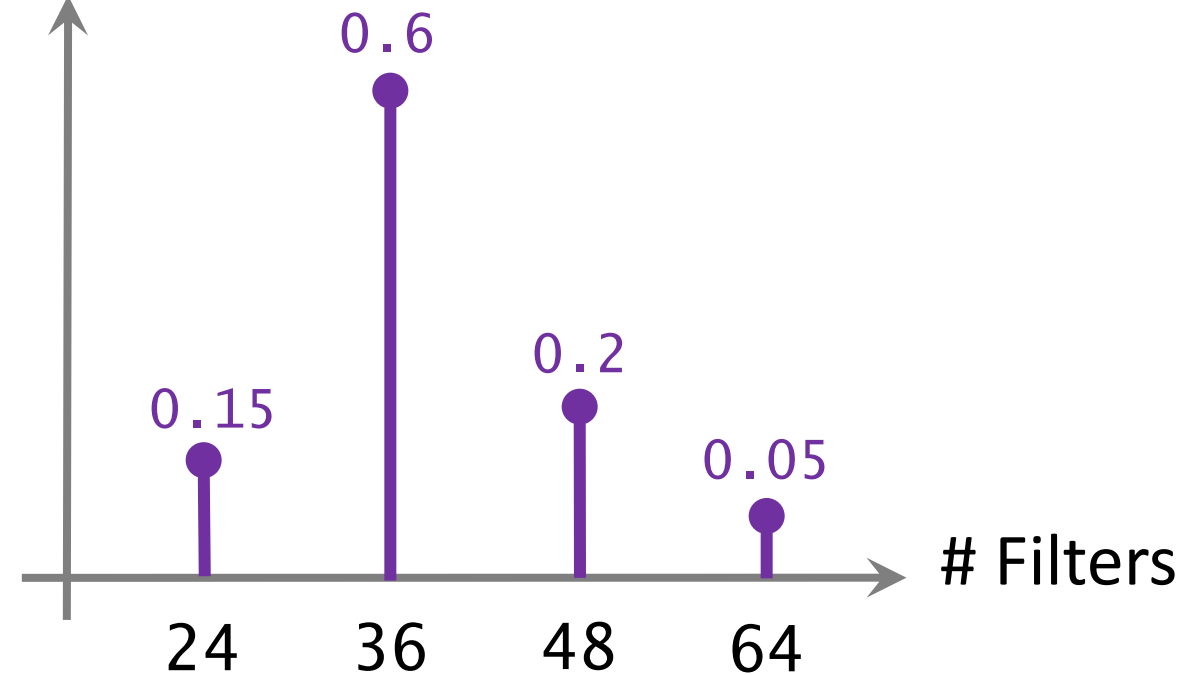




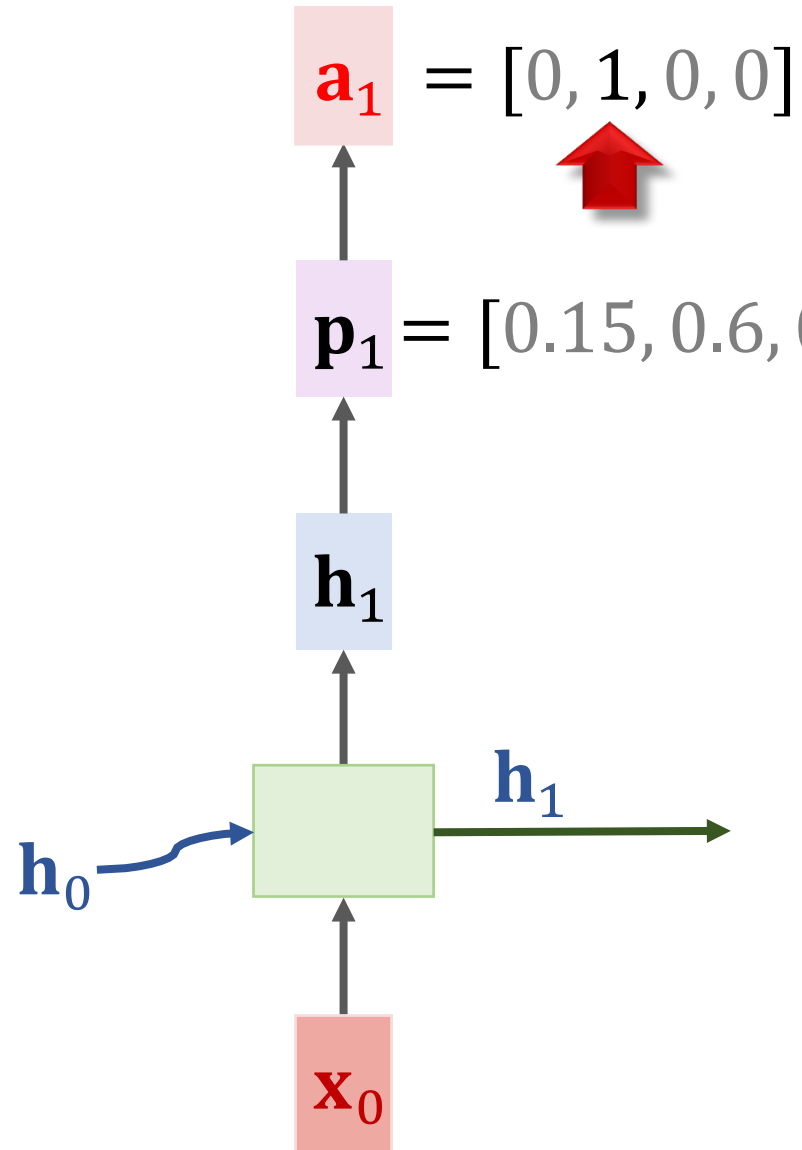
## Predict Number of Filters



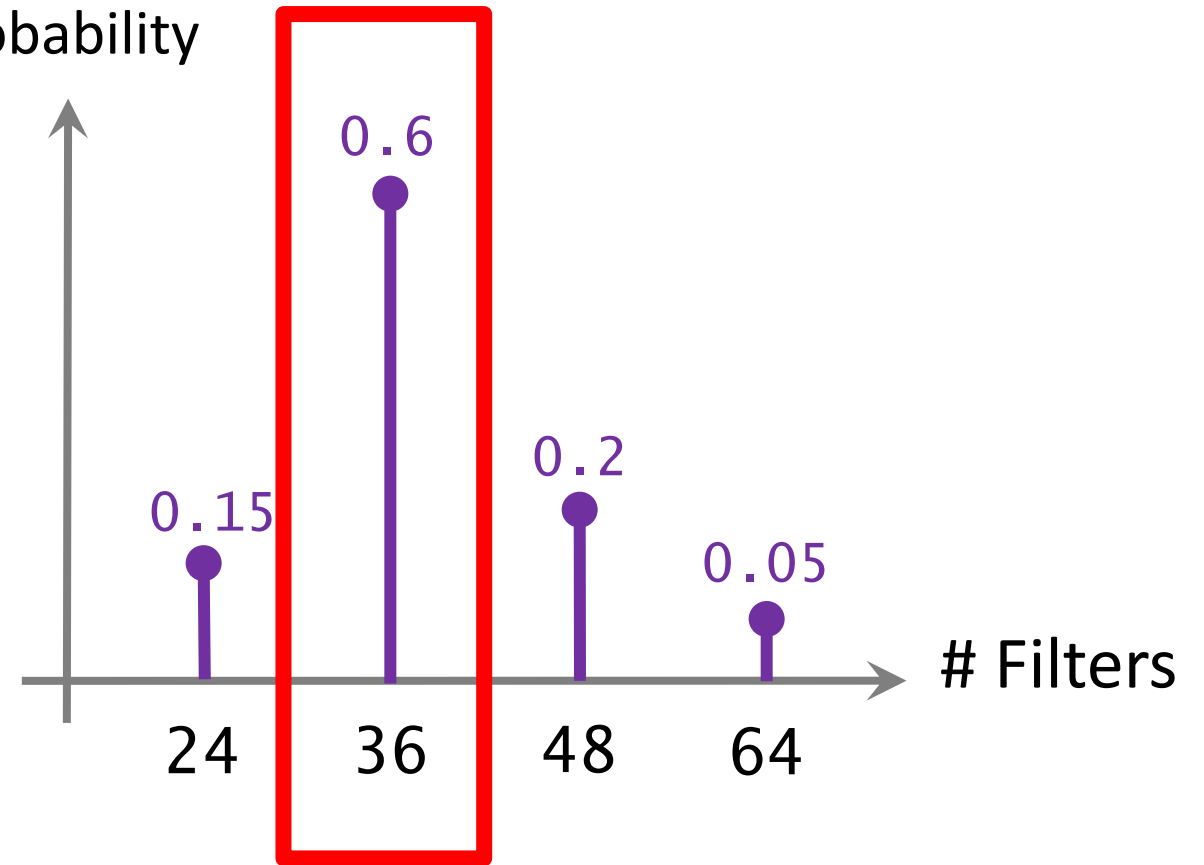
Probability

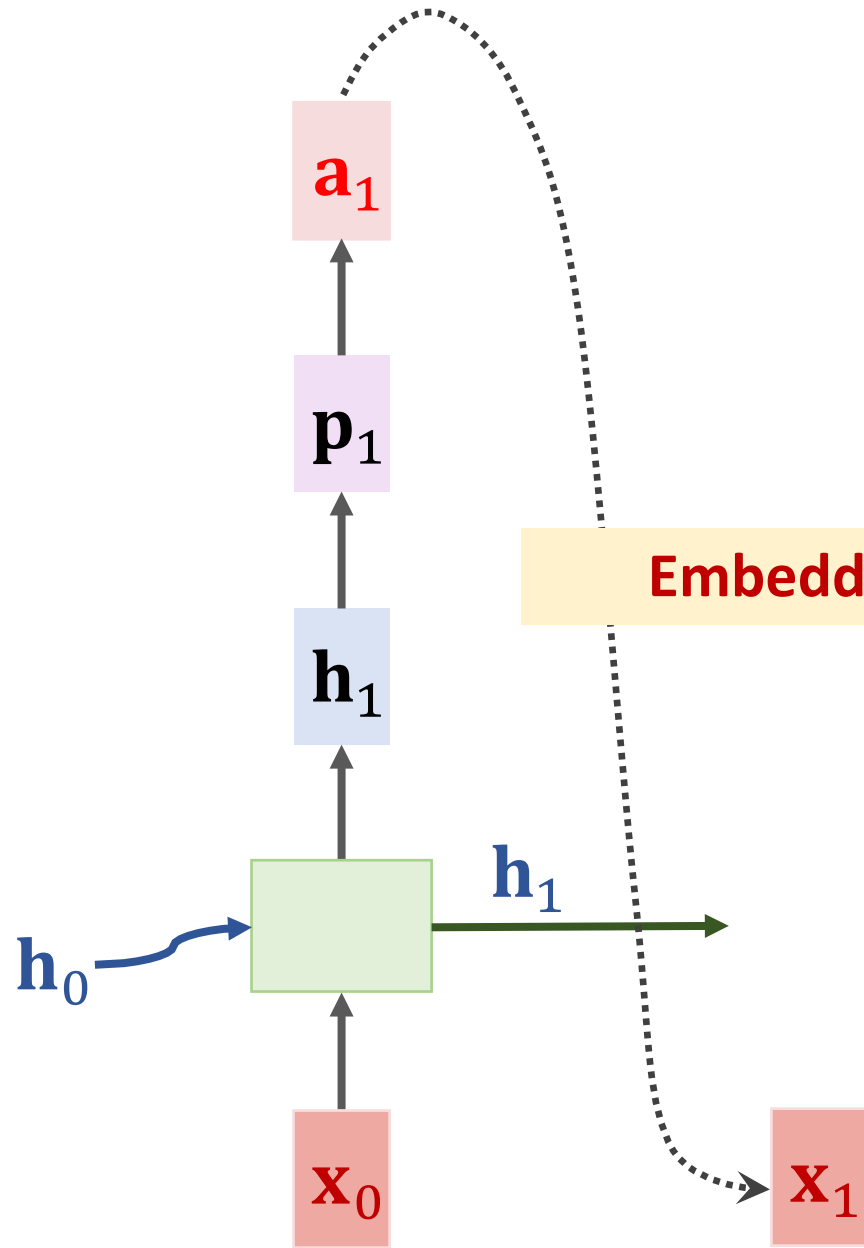


## Predict Number of Filters

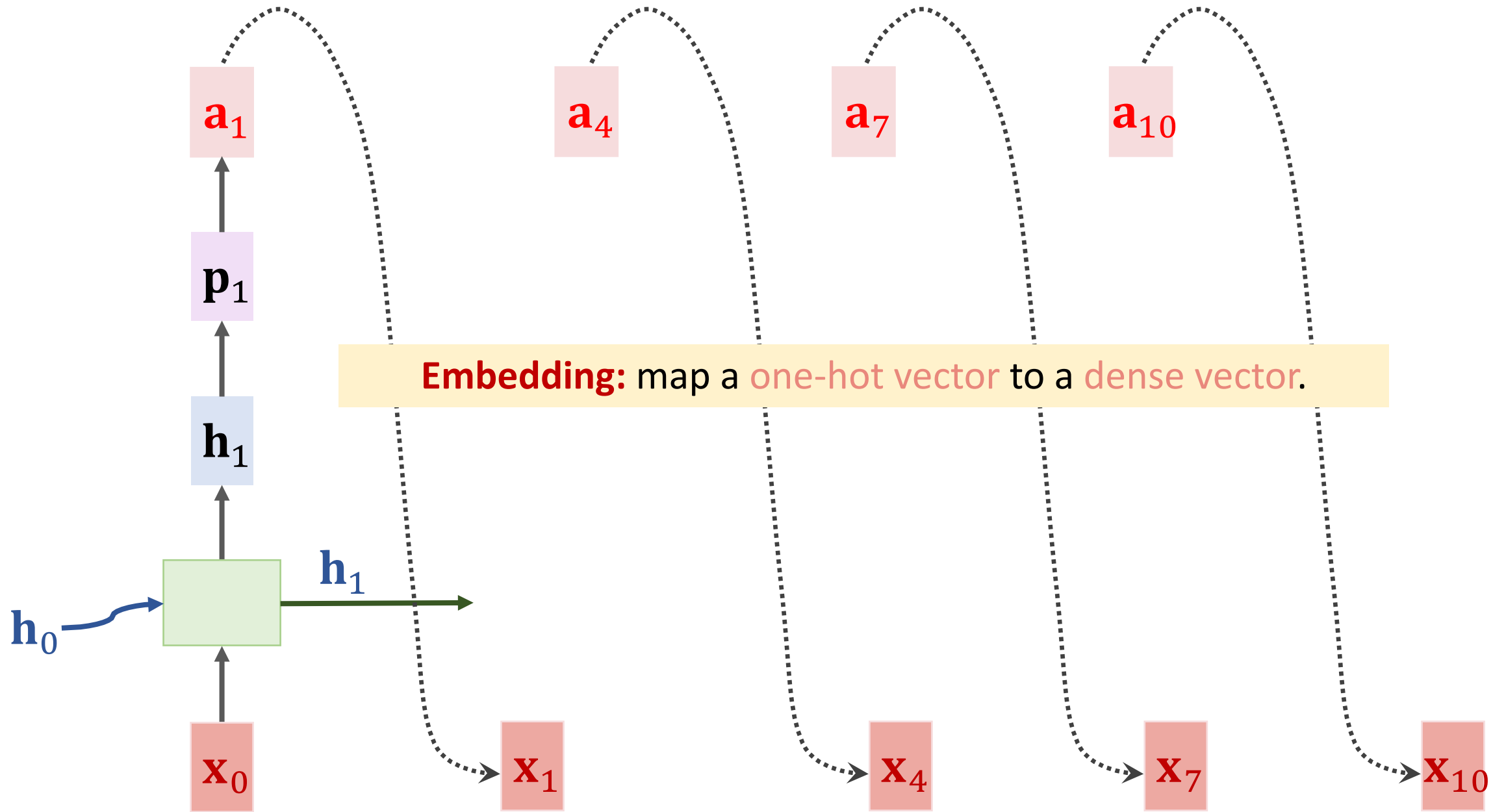


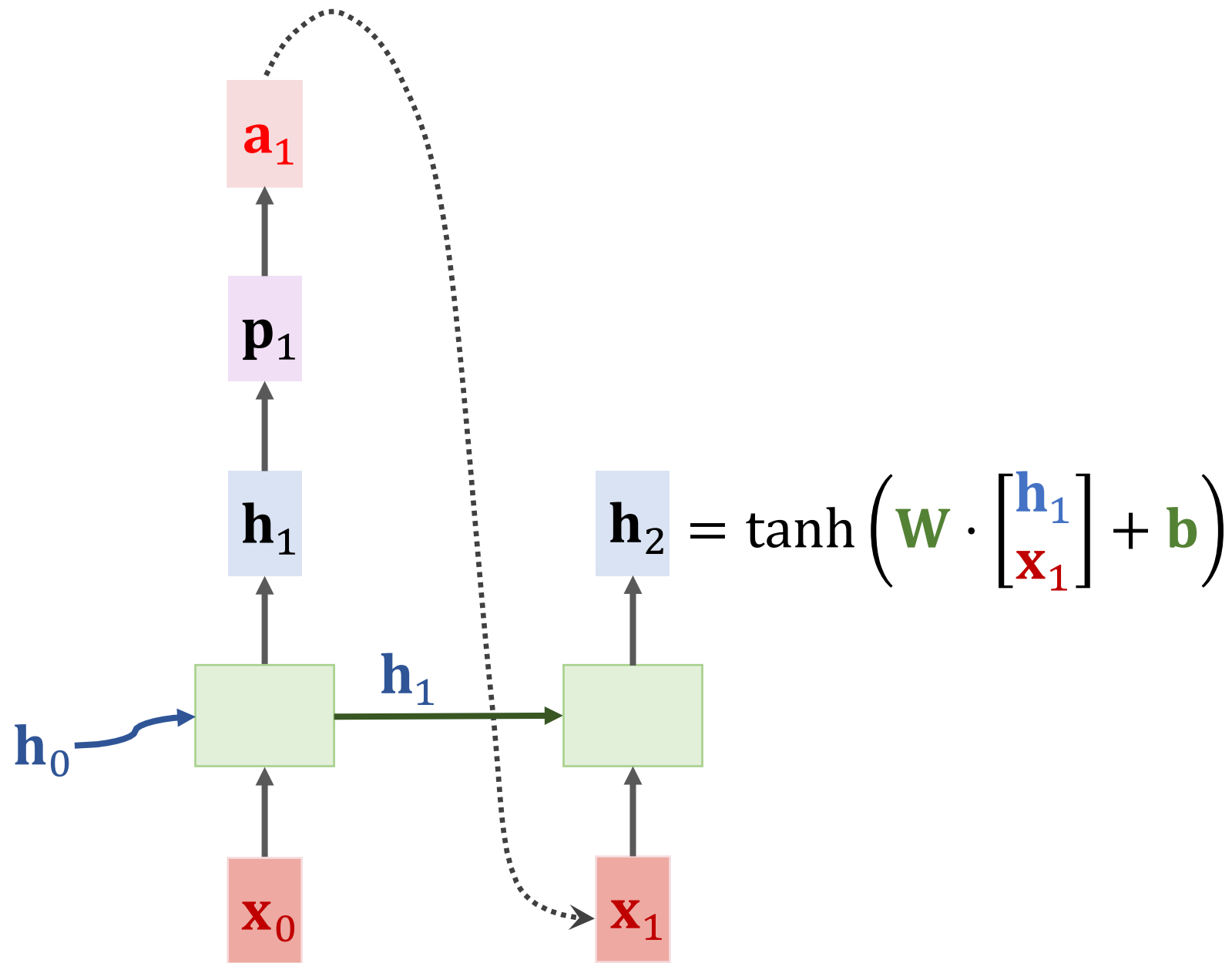
Probability

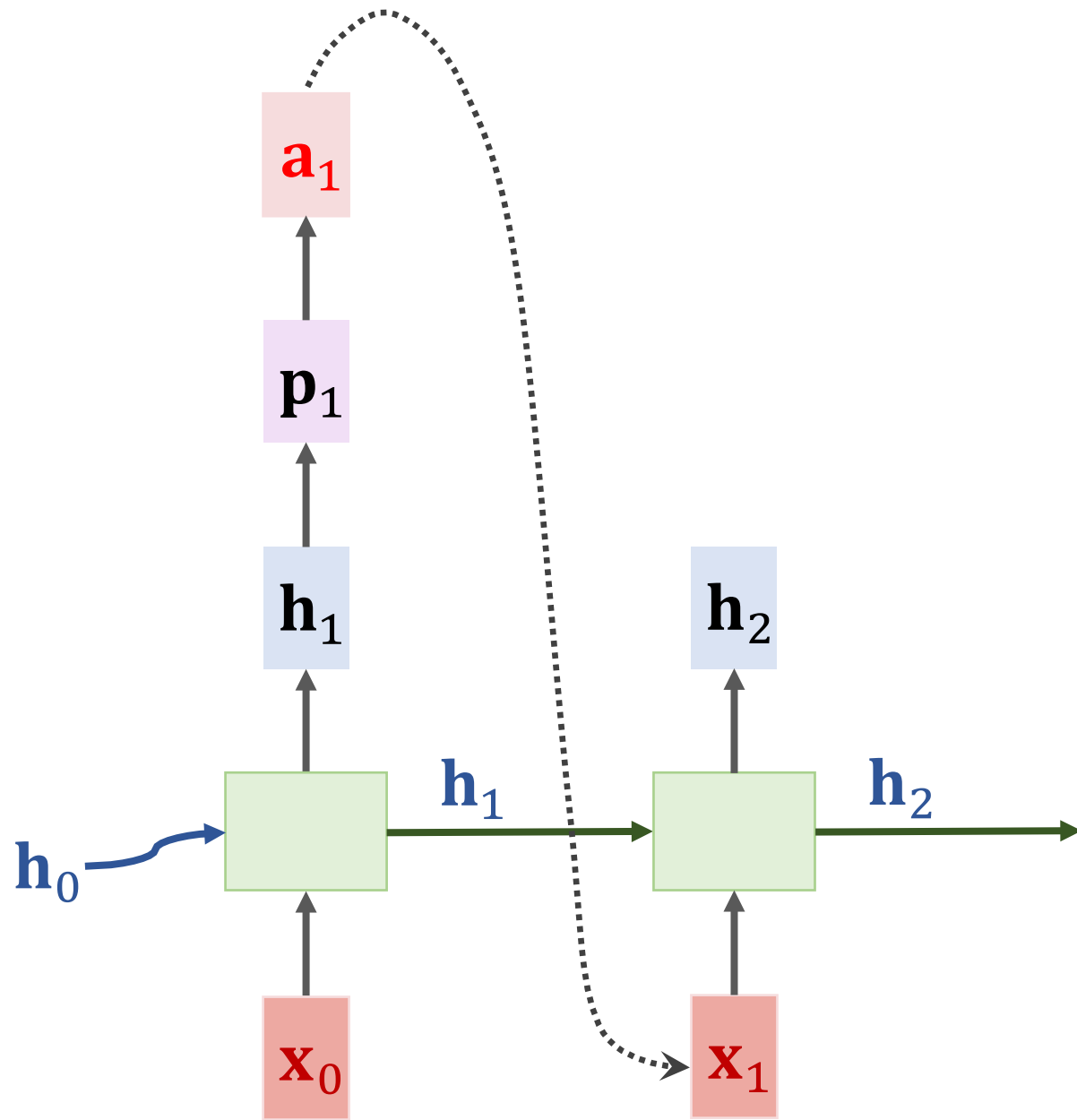


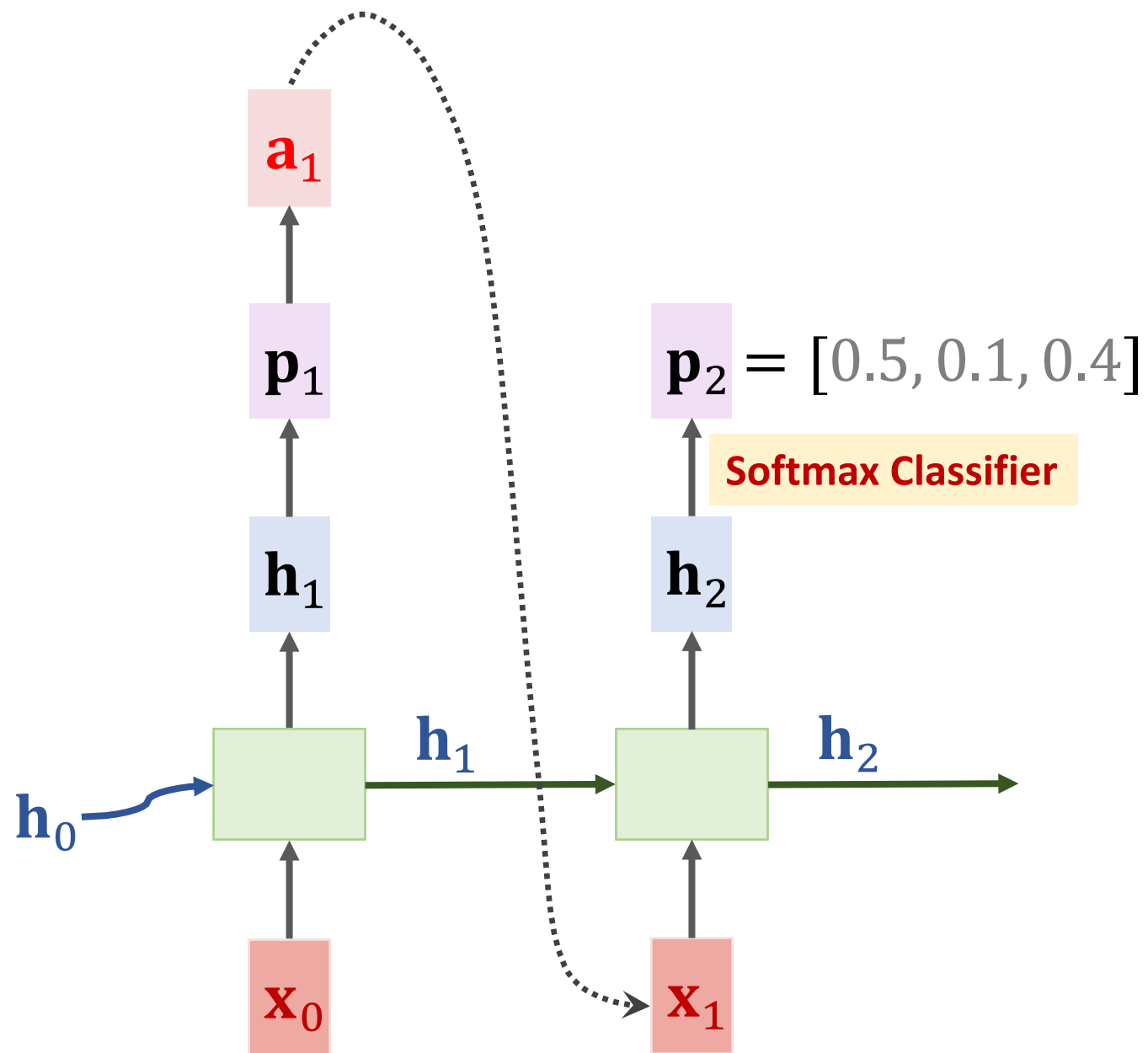


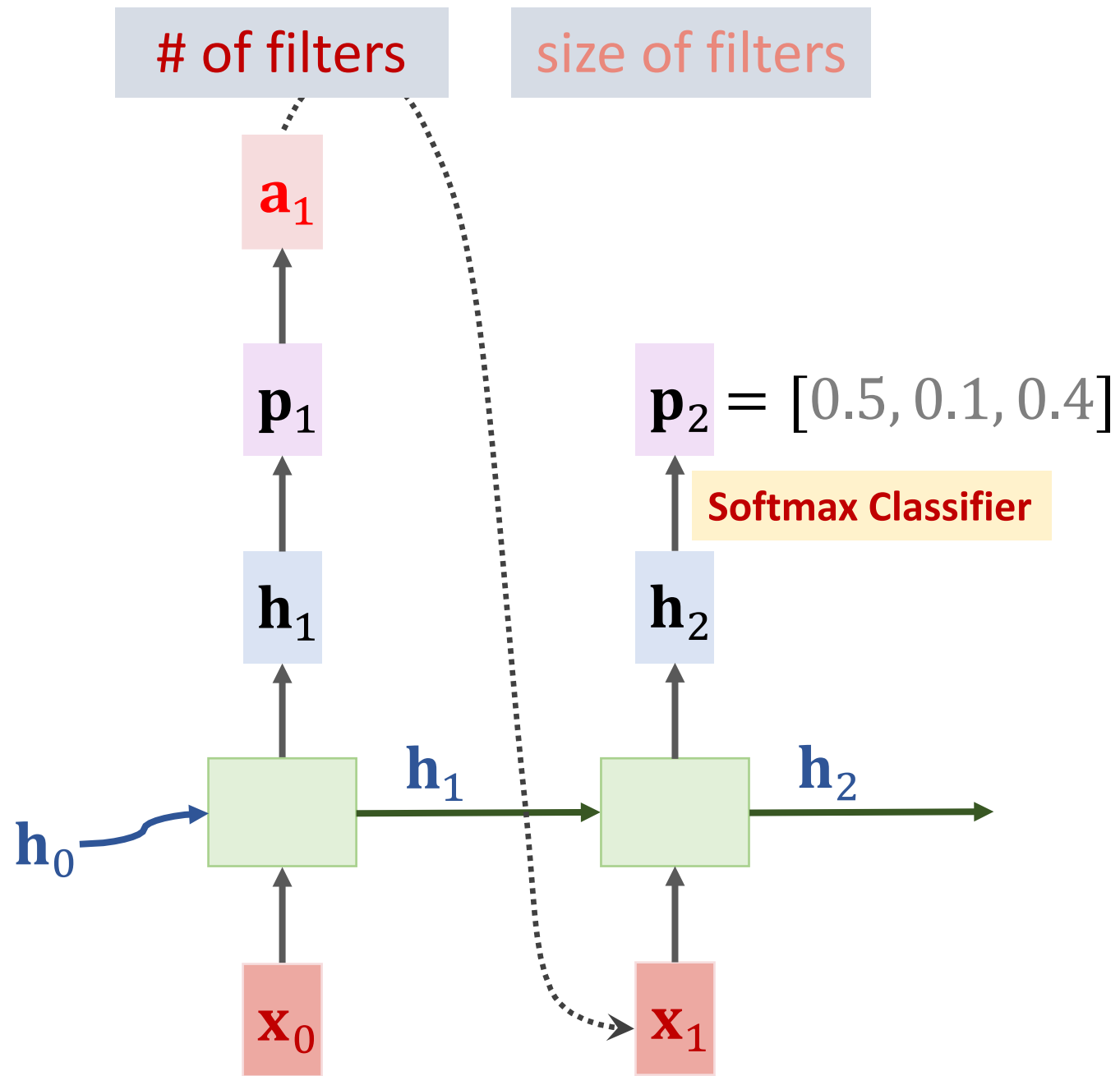
**Embedding:** map a one-hot vector to a dense vector.



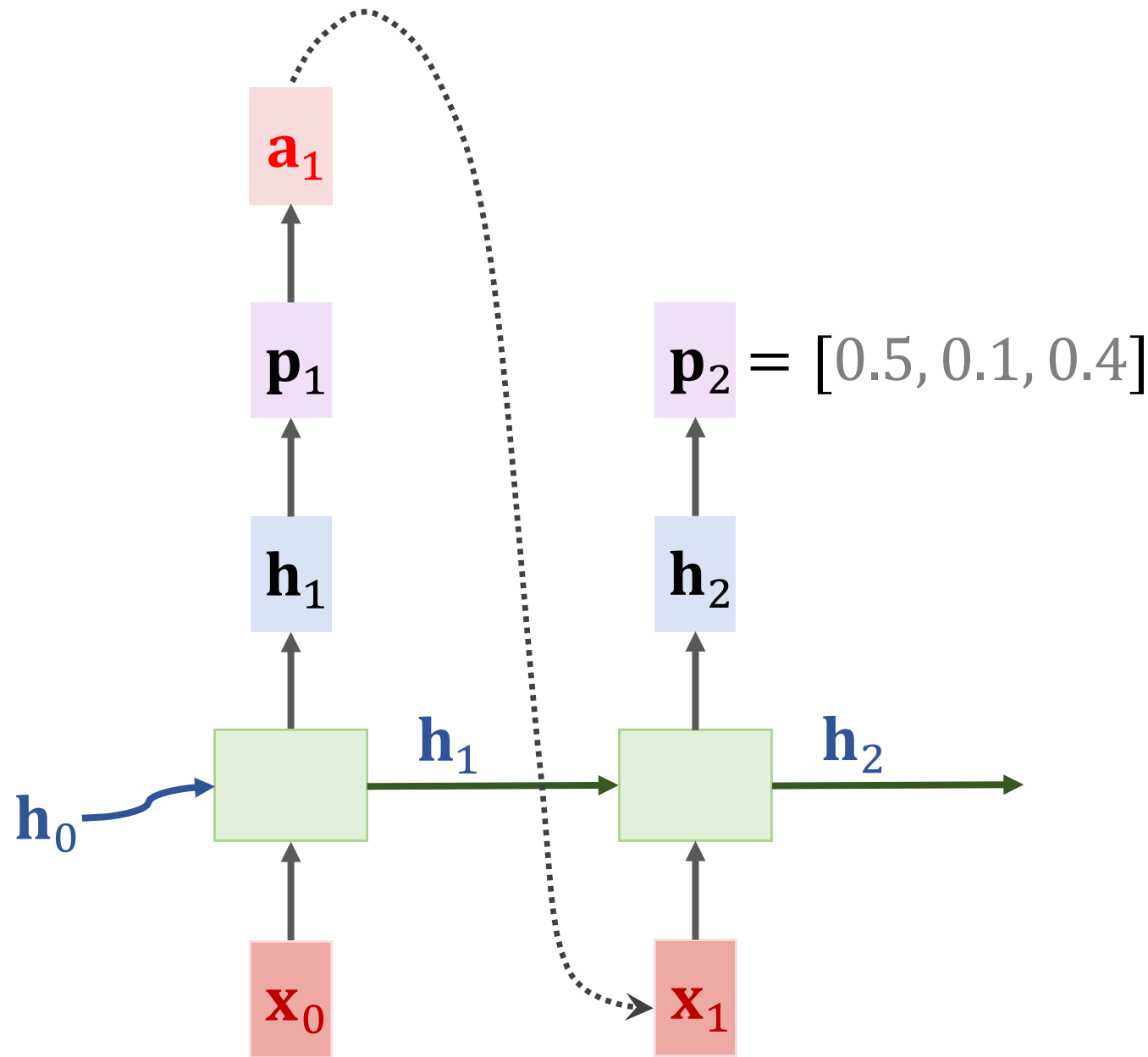






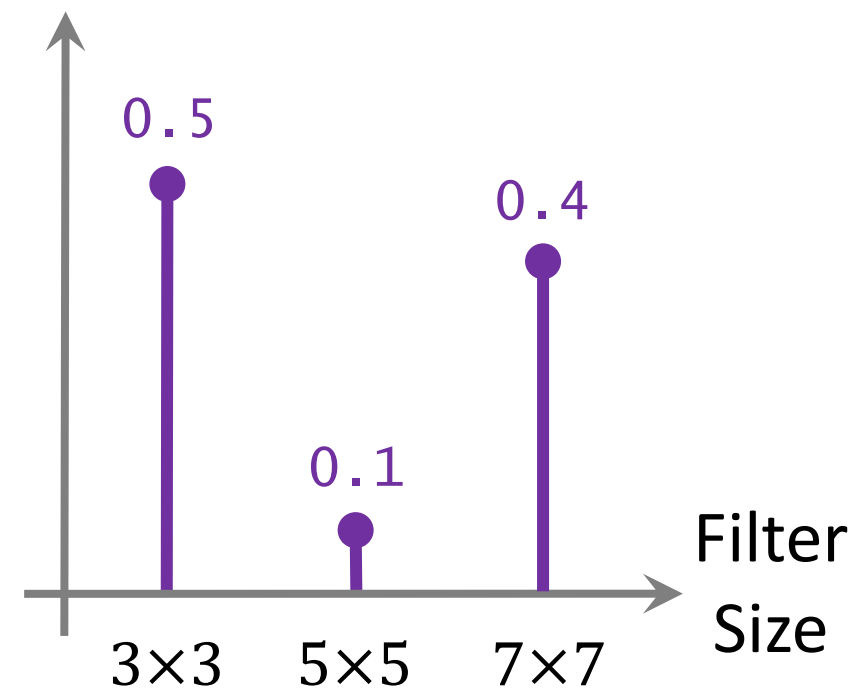


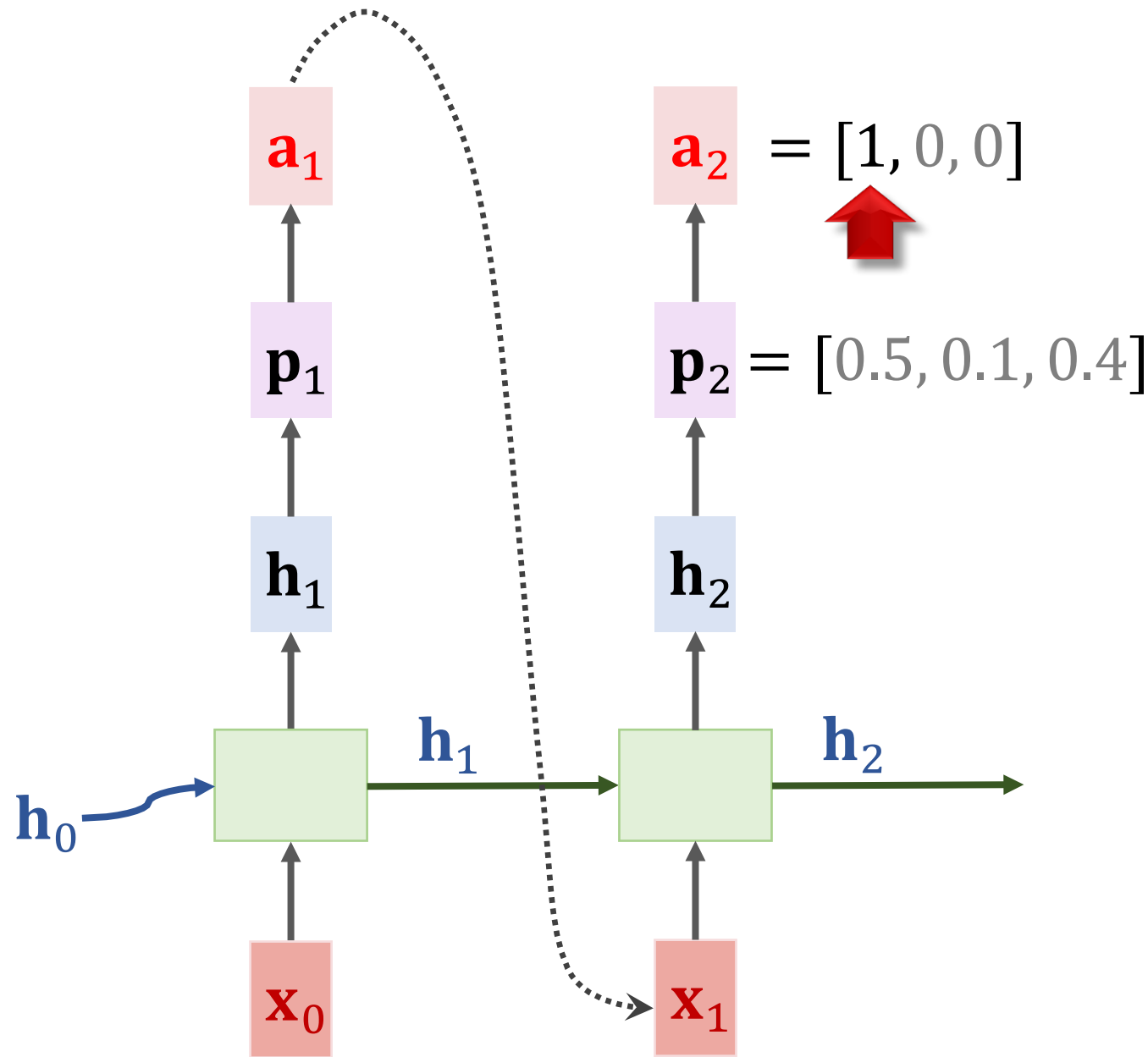




## Predict Size of Filters

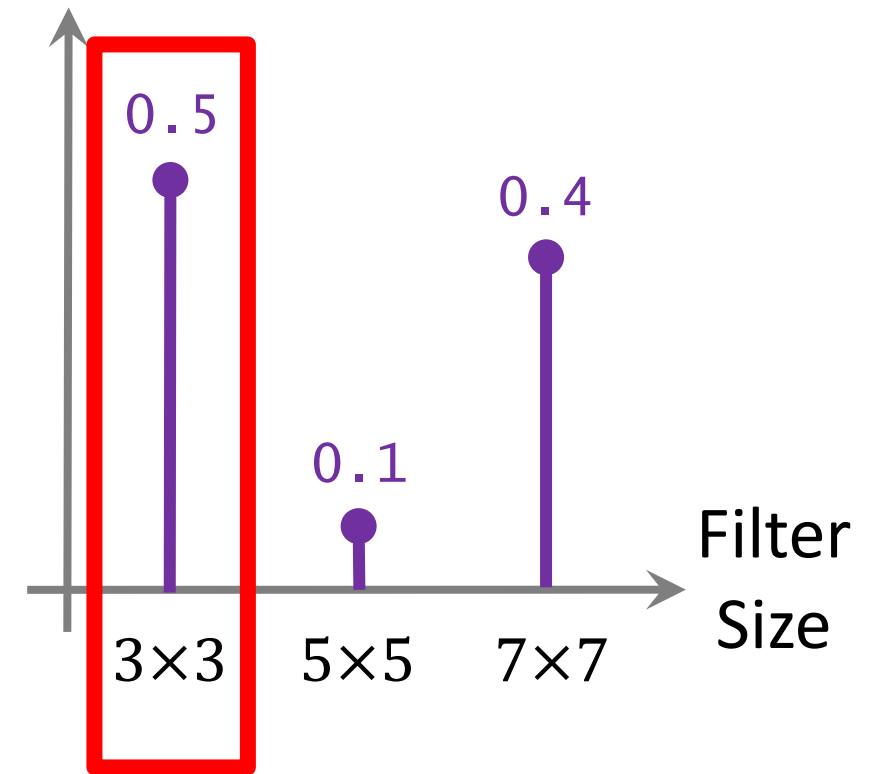
Probability

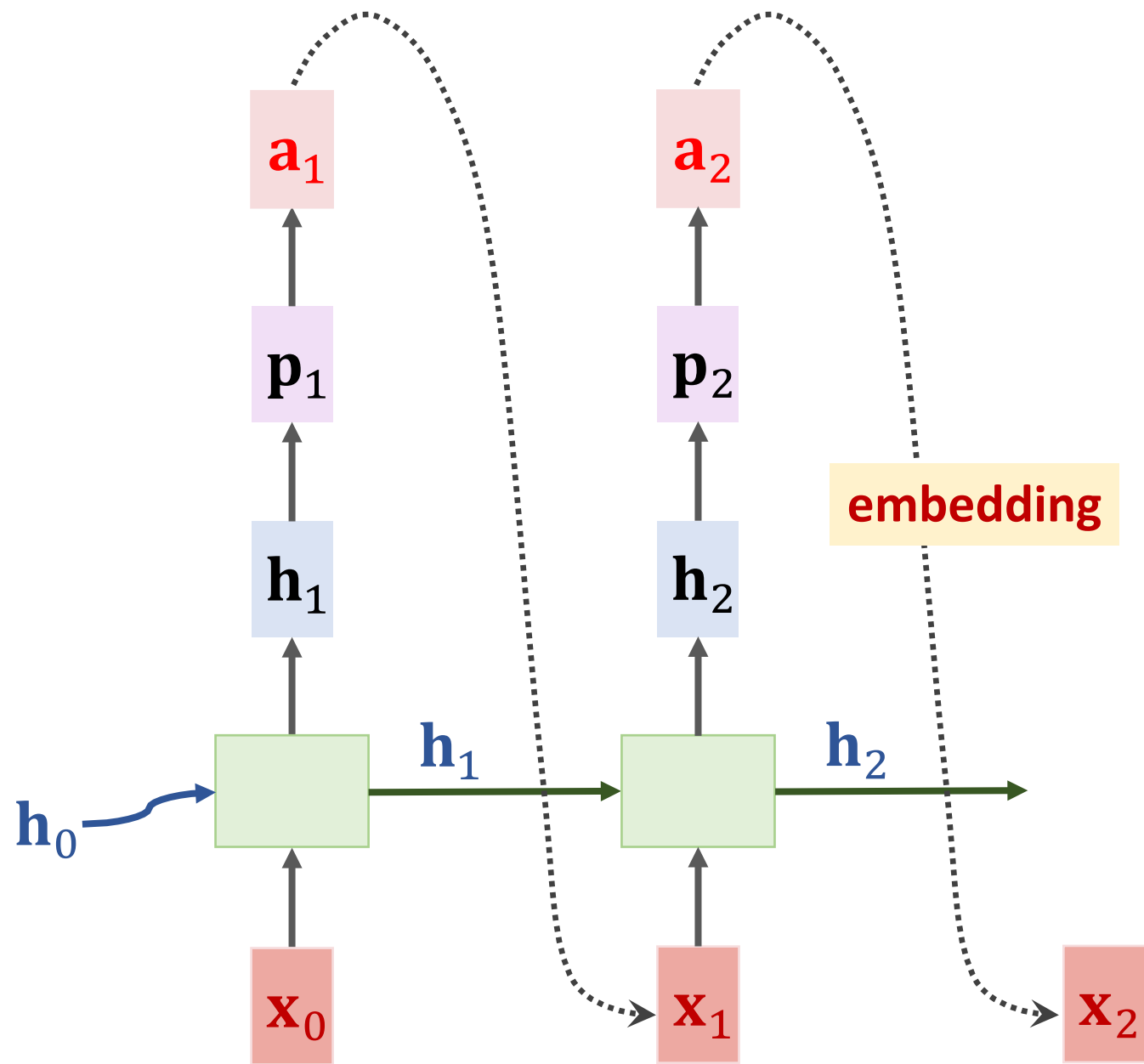


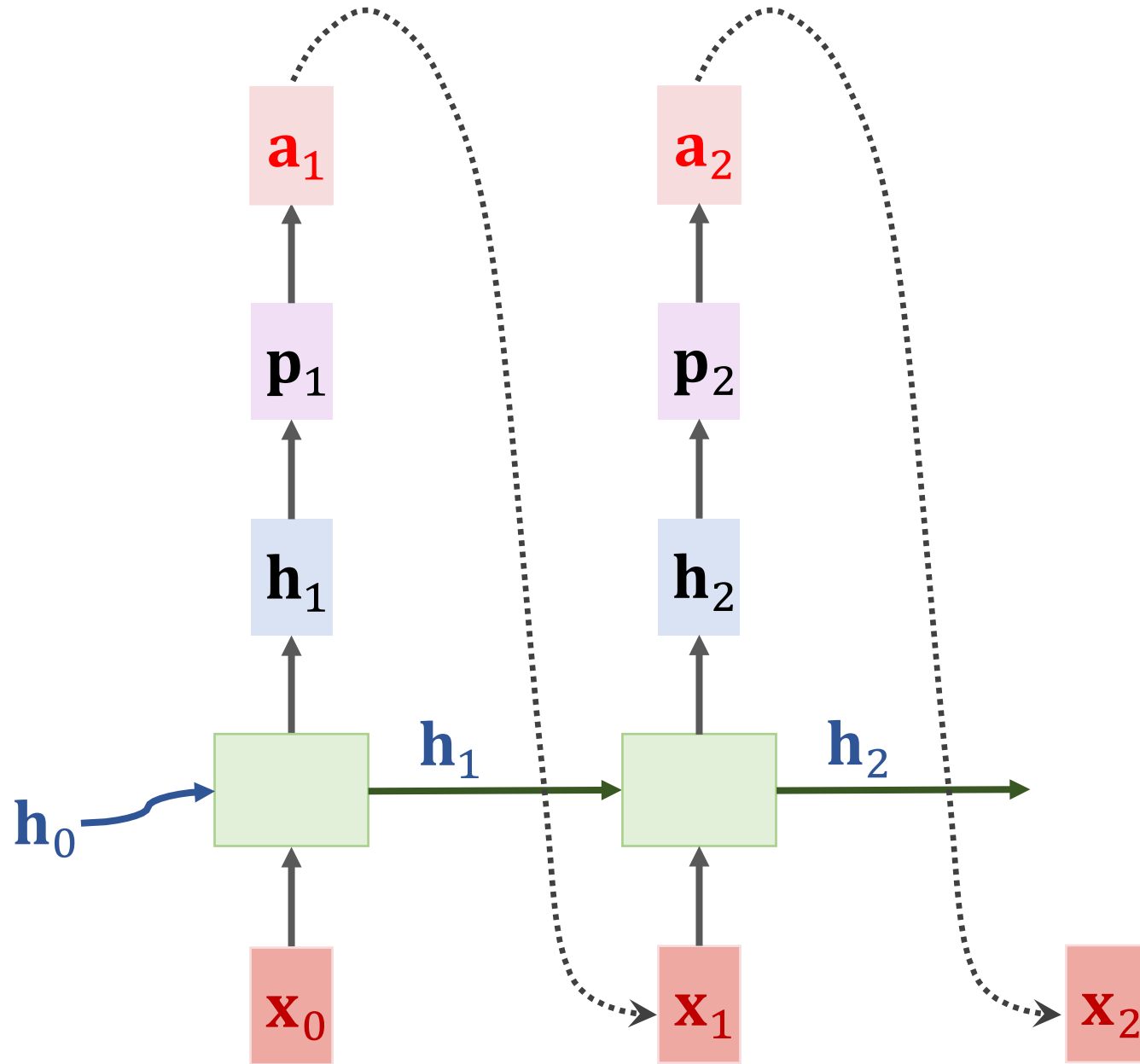


## Predict Size of Filters

Probability

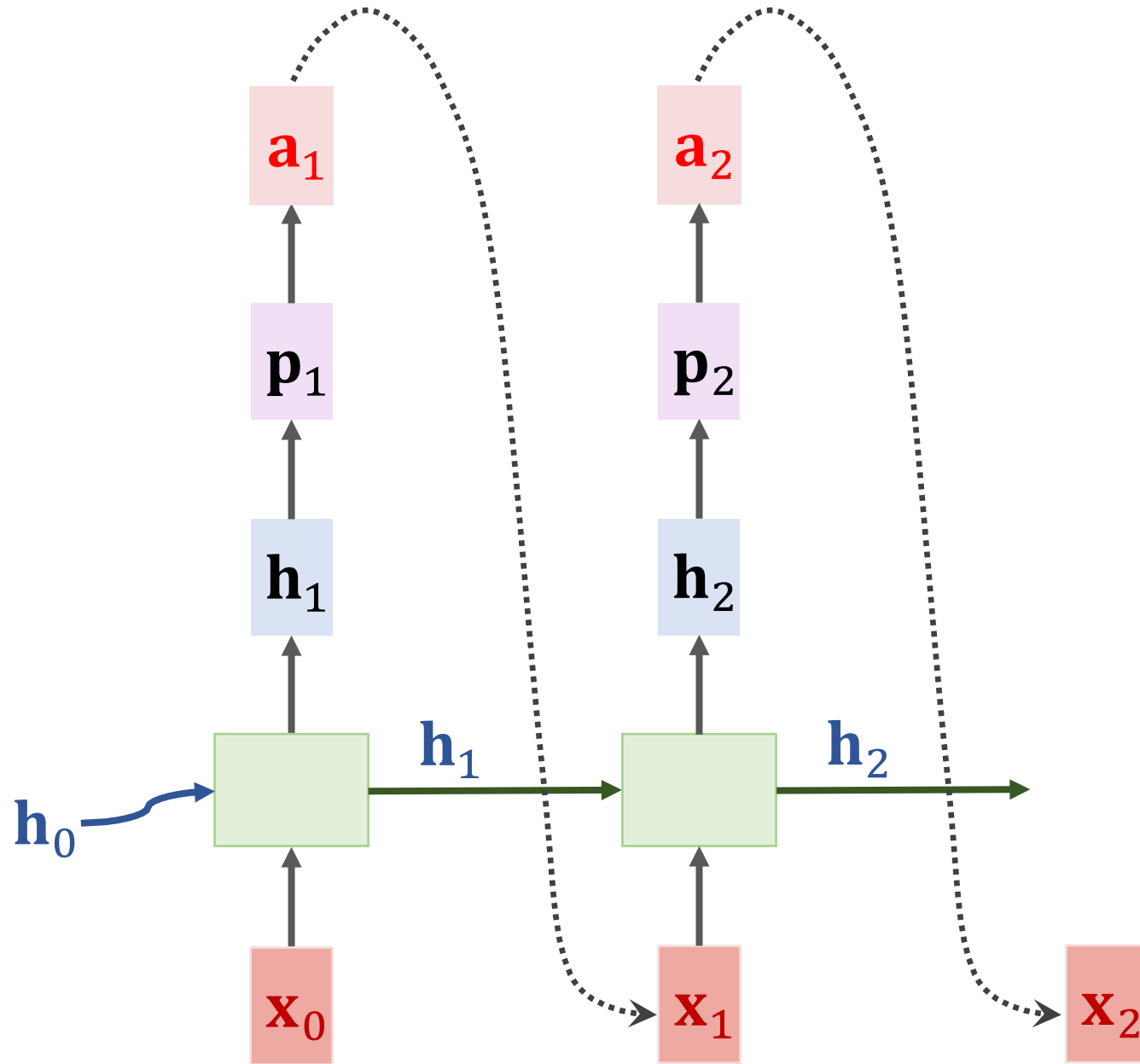






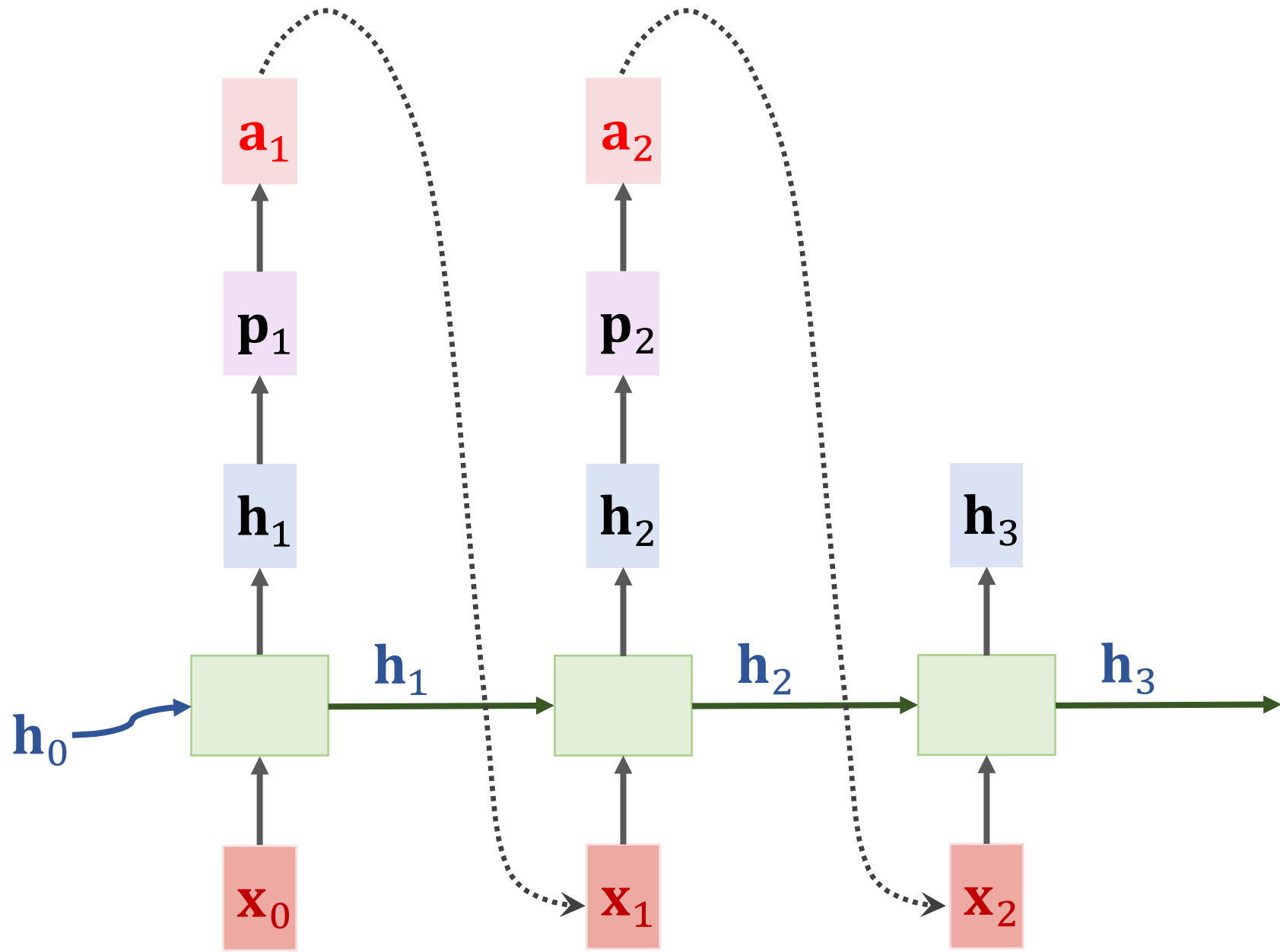
An embedding layer can be shared among the same task.

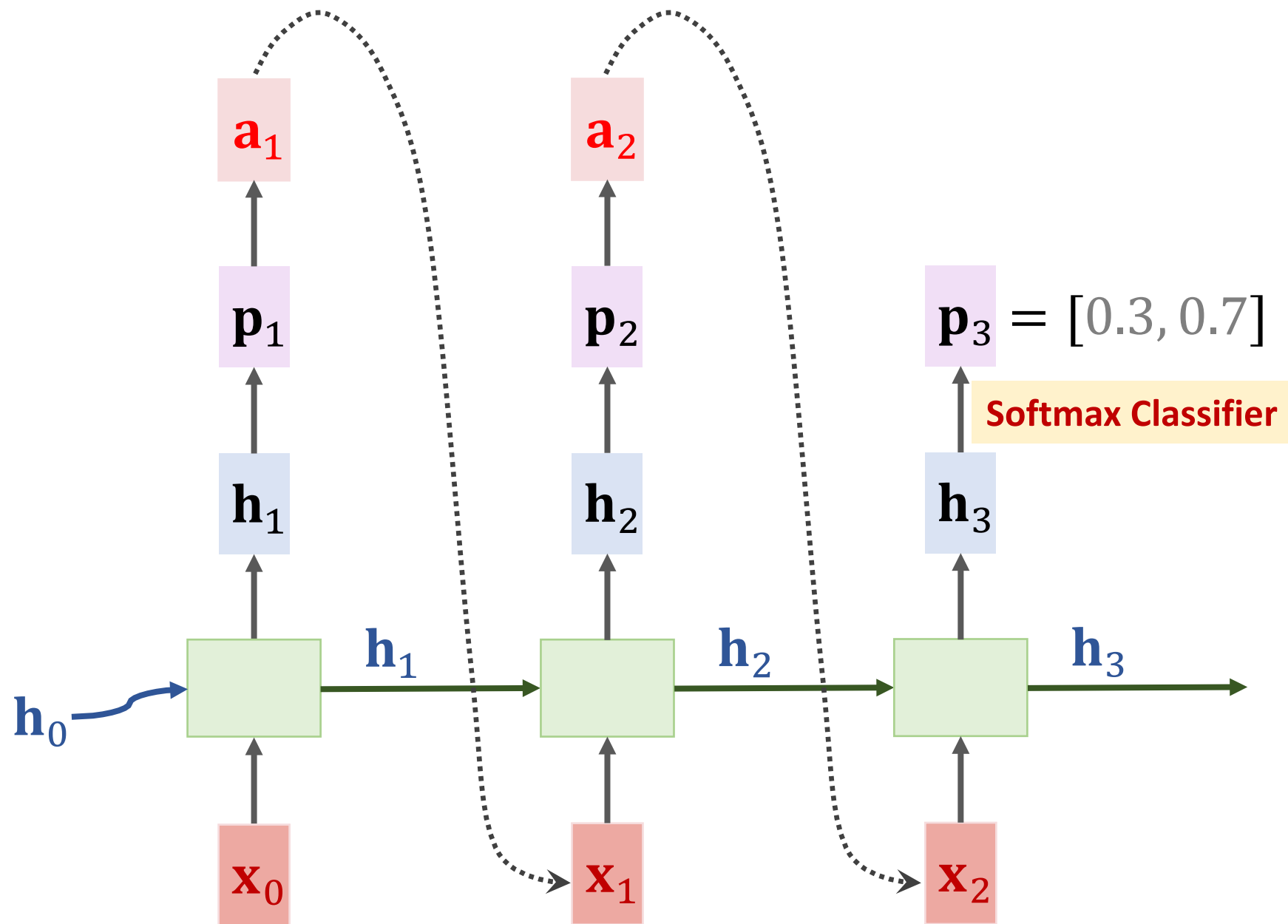
- E.g.,  $a_2$ ,  $a_5$ ,  $a_8$ ,  $a_{11}$  are for size of filters.
- They can share embedding layer.

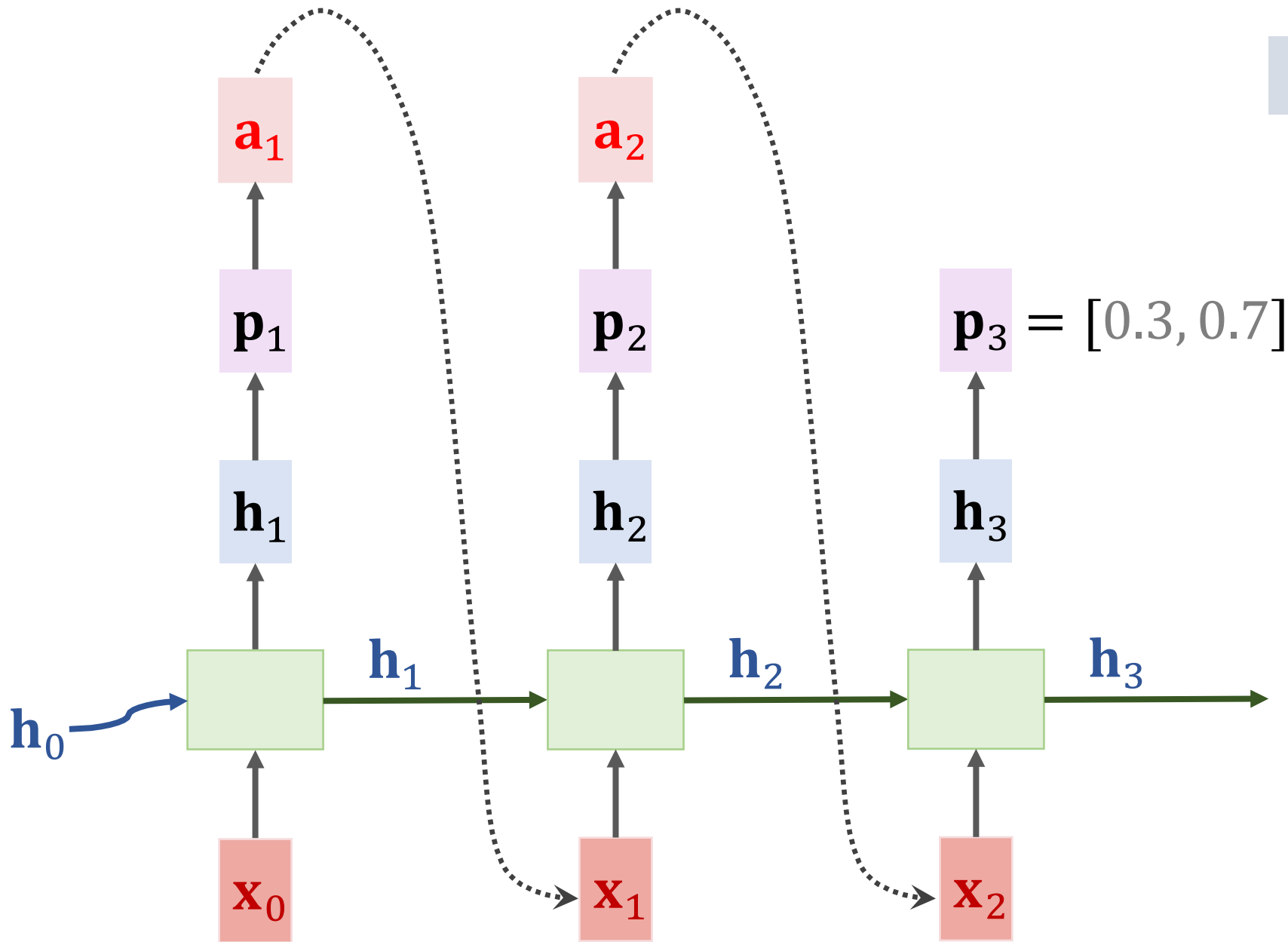


Different tasks do not share embedding layers.

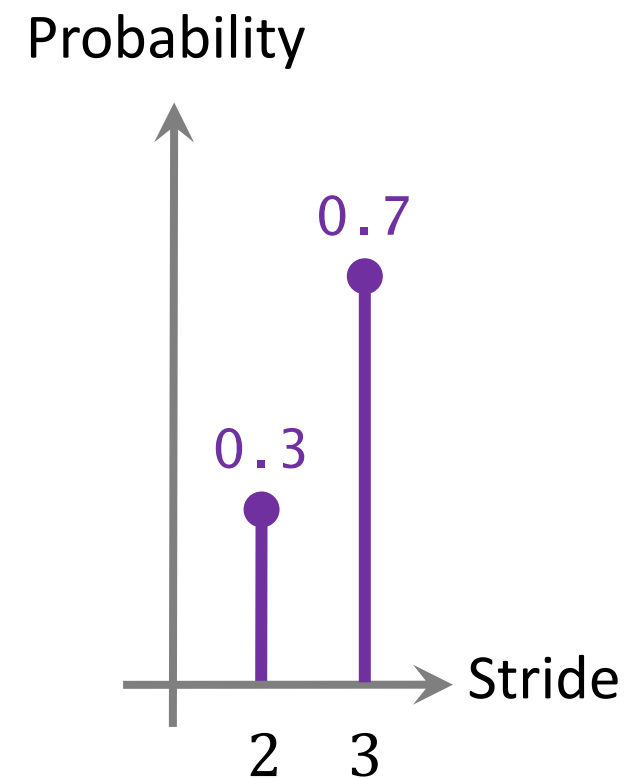
- E.g.,  $a_1$  and  $a_2$  are for different tasks.
- They cannot share embedding layer.



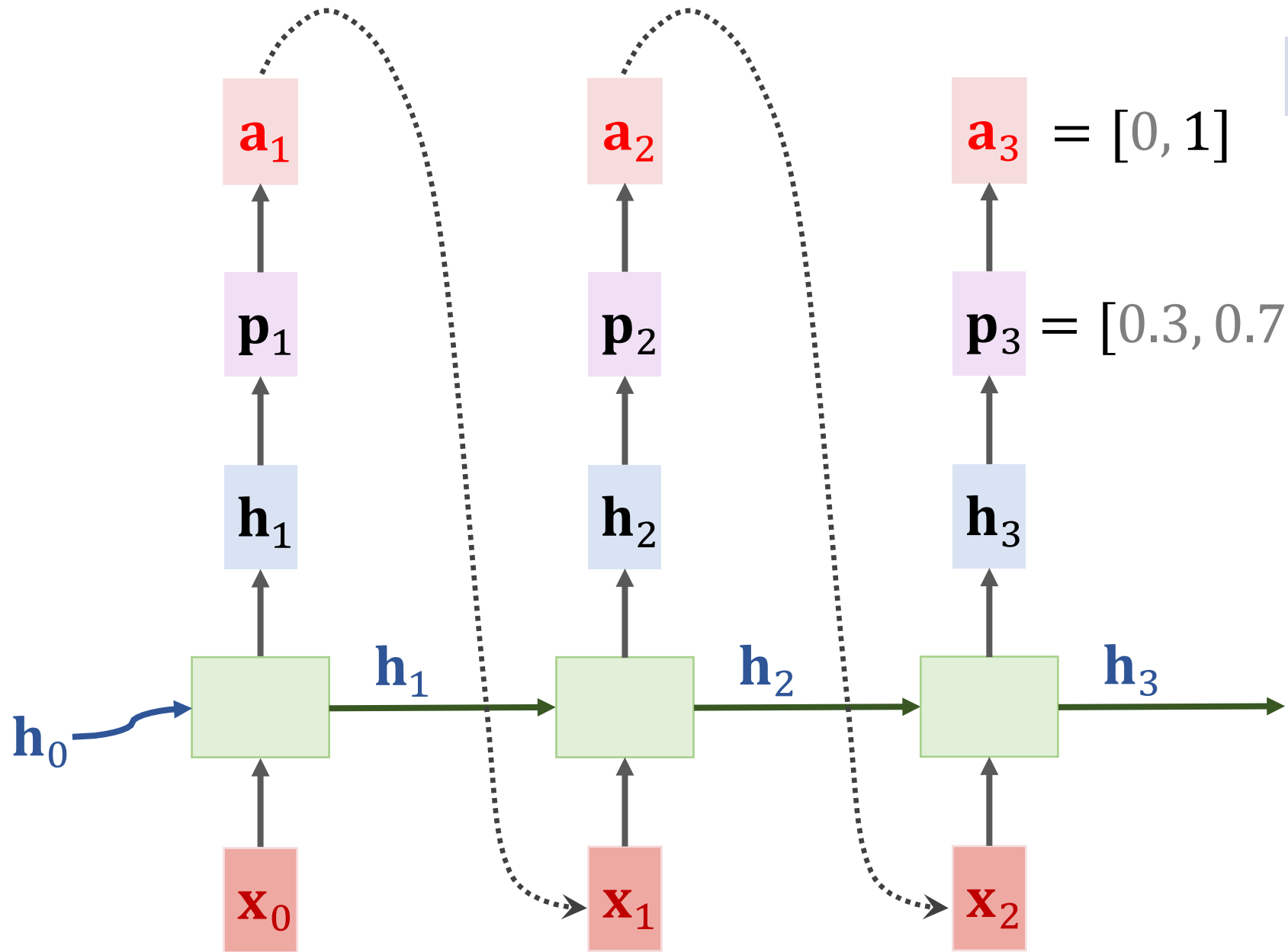




## Predict Stride

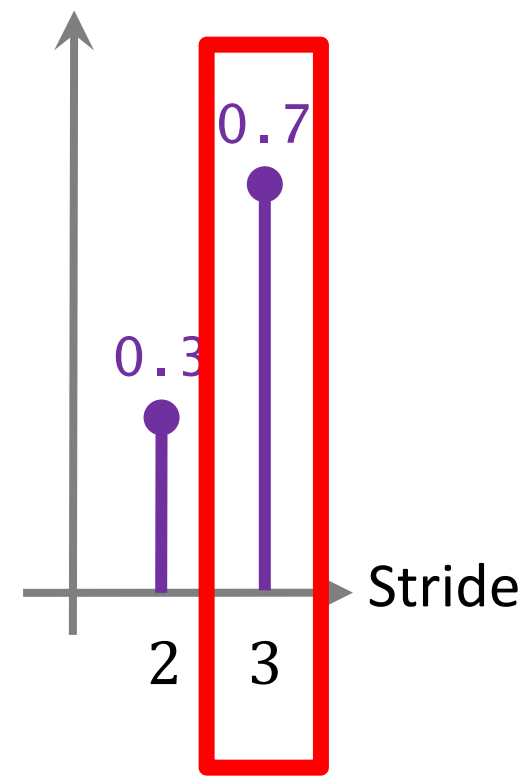


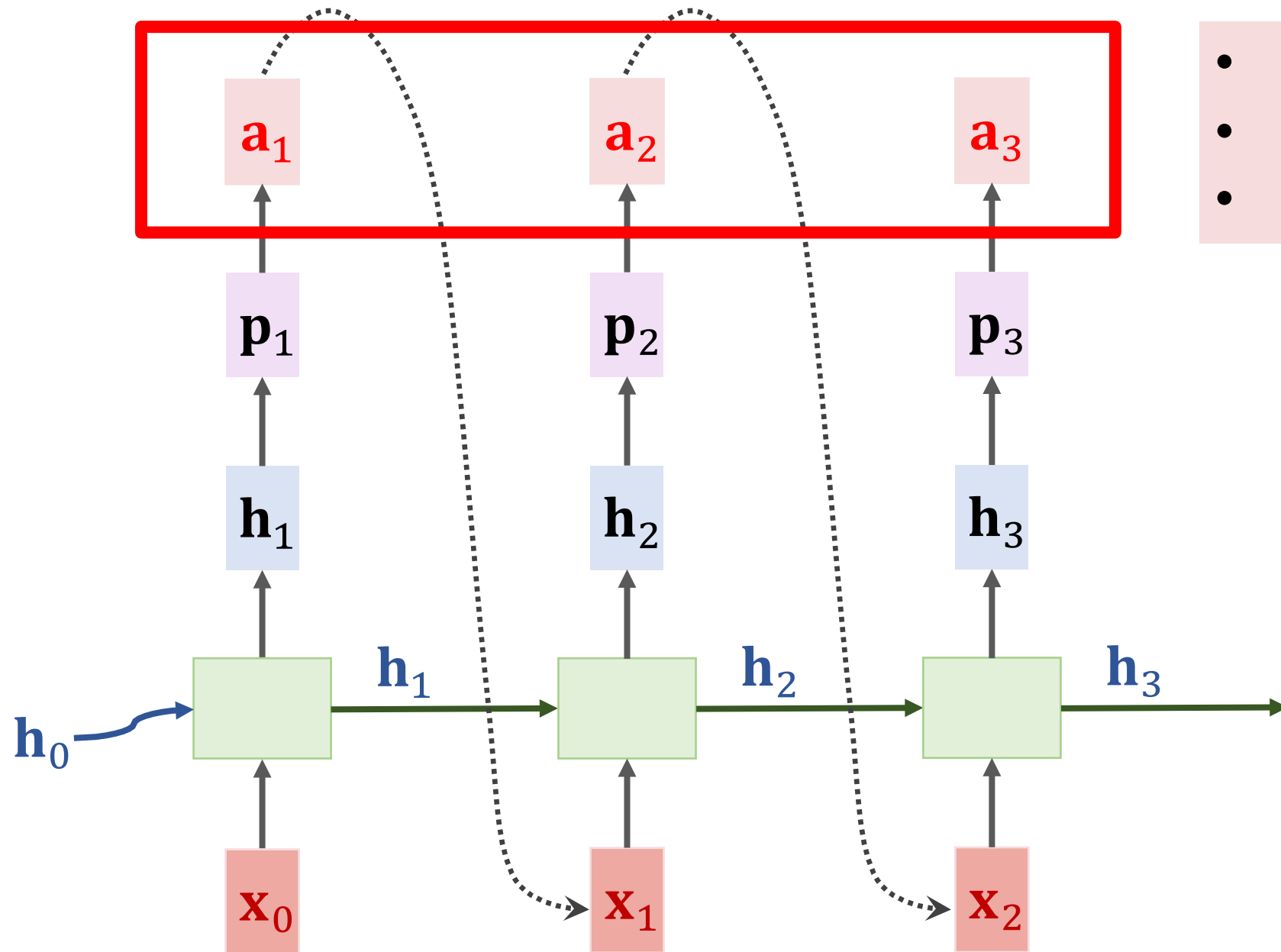




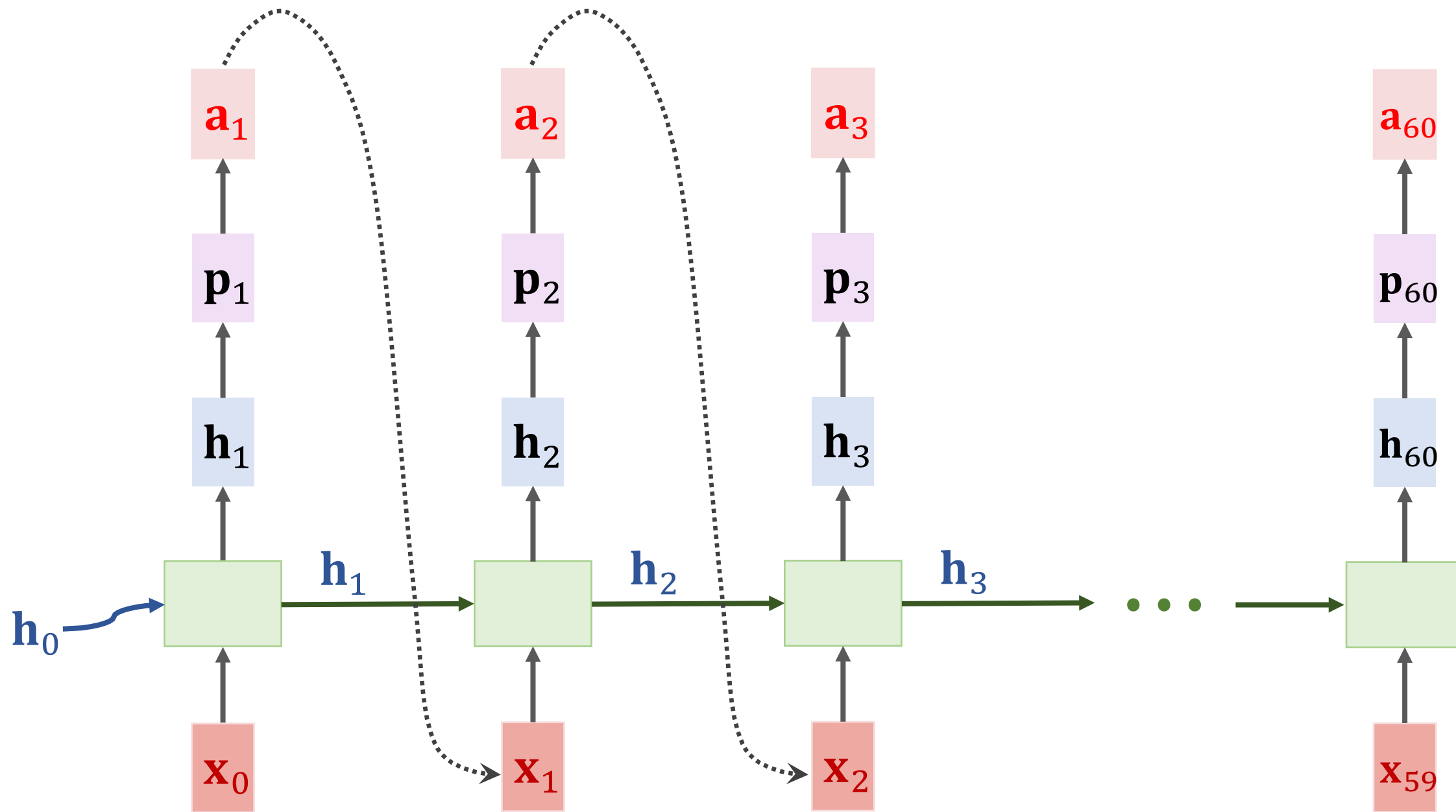
**Predict Stride**

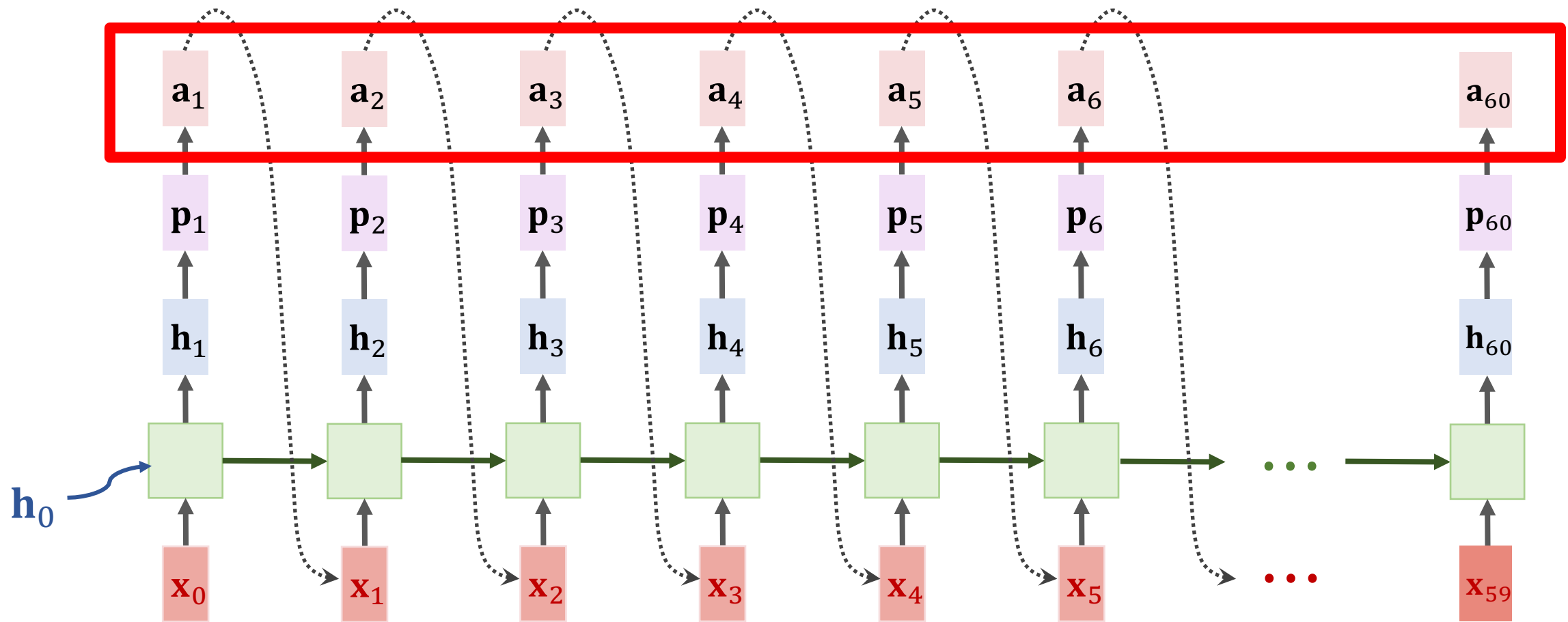
Probability

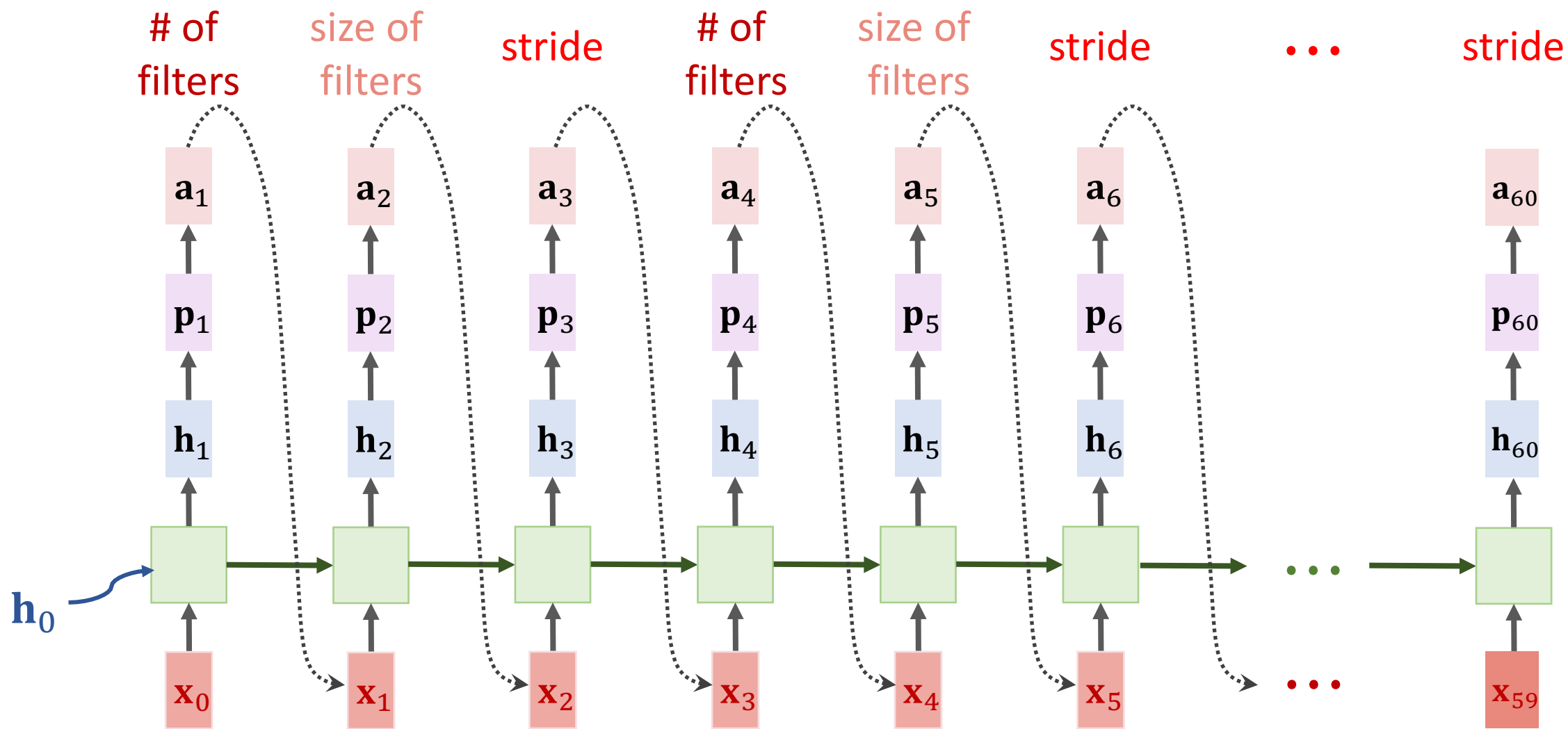


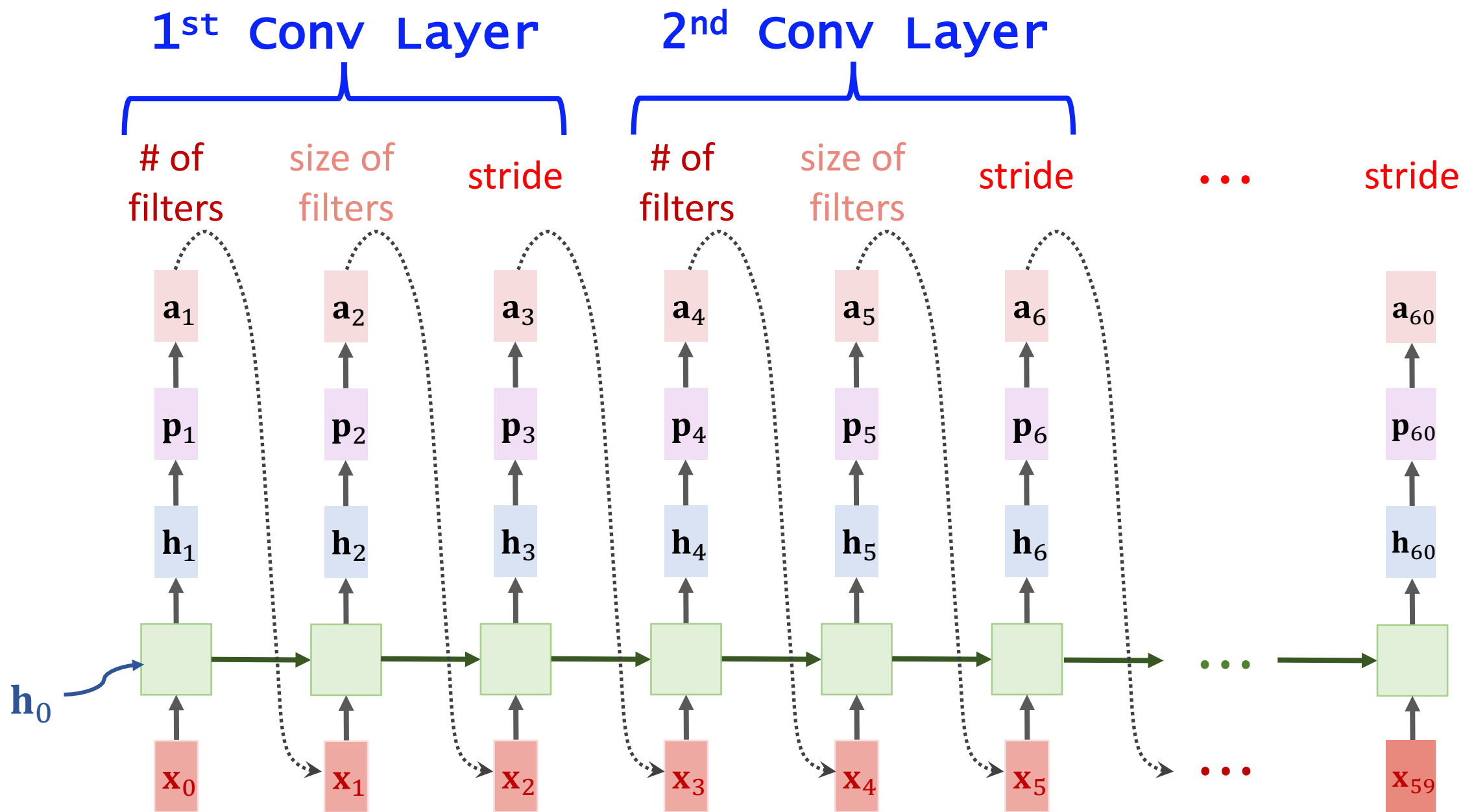


- Filter number = 36
- Filter size =  $3 \times 3$
- Stride = 3









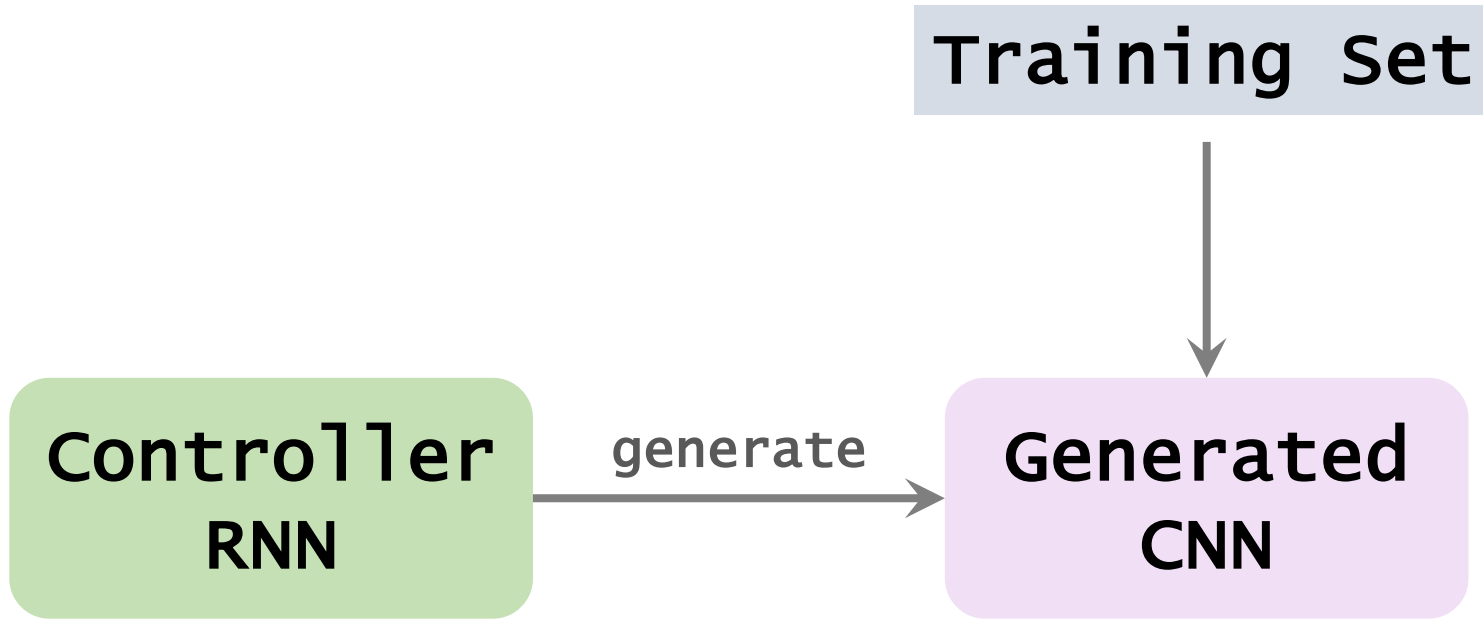
# **Training Controller RNN**

# How to train the controller RNN?

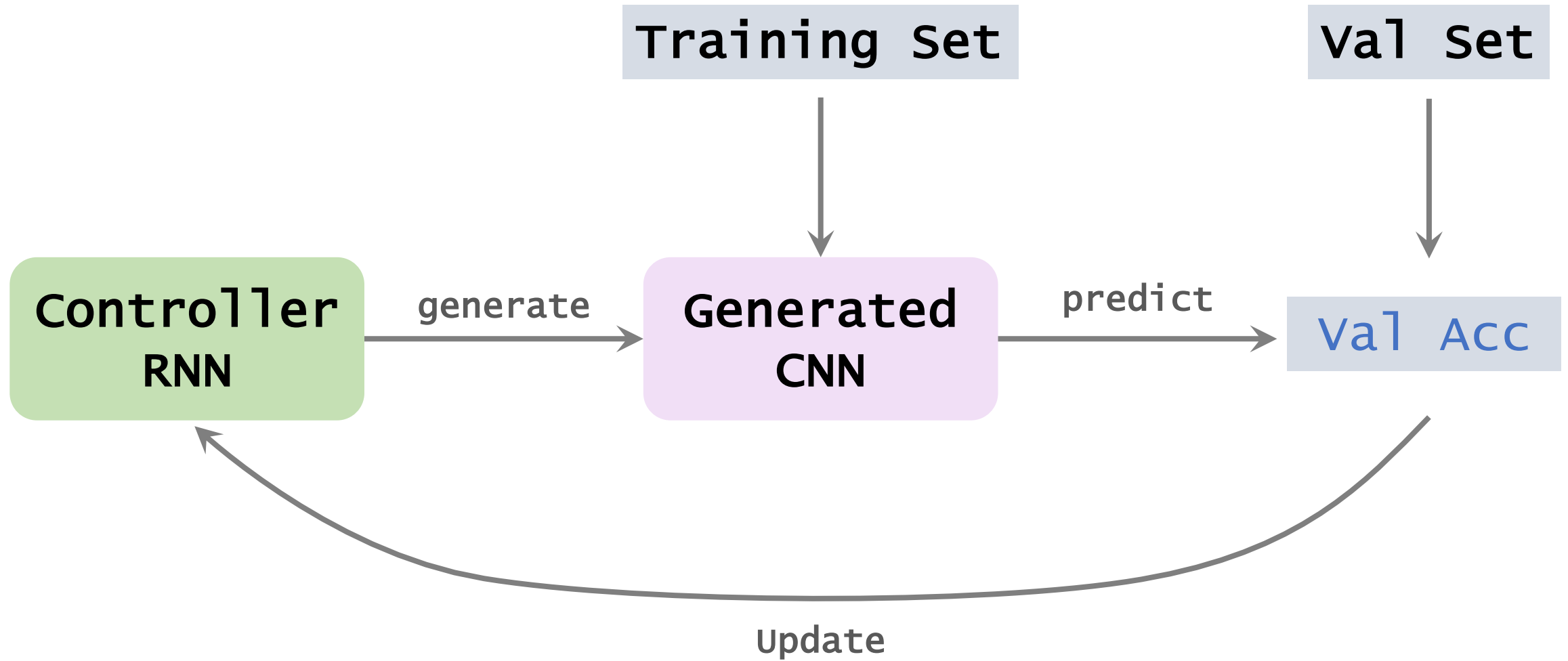
- The controller RNN outputs the hyper-parameters of a CNN.
- With the hyper-parameters at hand, instantiate a CNN.
- Train the CNN on a dataset, e.g., CIFAR-10, ImageNet, etc.
- Compute validation accuracy on a held-out dataset.
- Validation accuracy is the supervision for training the controller RNN.



# How to train the controller RNN?



# How to train the controller RNN?



# Challenges

- ➡ •  $r$ : objective function (to maximize).
- ➡ •  $\theta$ : optimization variable.

# Challenges

- $r$ : objective function (to maximize).
- $\theta$ : optimization variable.
- What if  $r$  is a differentiable function of  $\theta$ ?
- Update  $\theta$  by gradient ascent:

$$\theta \leftarrow \theta + \beta \frac{\partial r}{\partial \theta}.$$

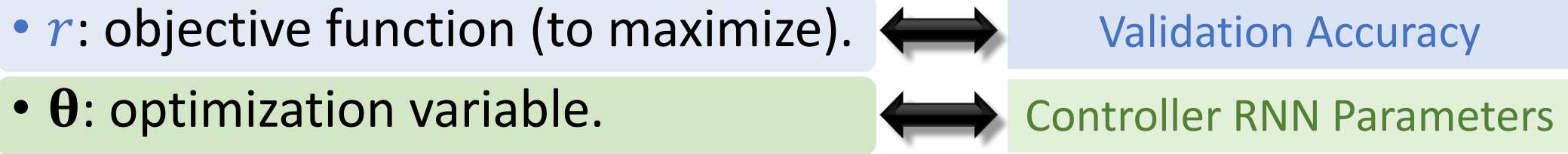
# Challenges

- $r$ : objective function (to maximize).
- $\theta$ : optimization variable.
- What if  $r$  is a differentiable function of  $\theta$ ?
- Update  $\theta$  by gradient ascent:

$$\theta \leftarrow \theta + \beta \cdot \frac{\partial r}{\partial \theta}.$$

- However, if  $r$  is **not** a differentiable function of  $\theta$ , then we cannot use the gradient to update  $\theta$ .

# Challenges

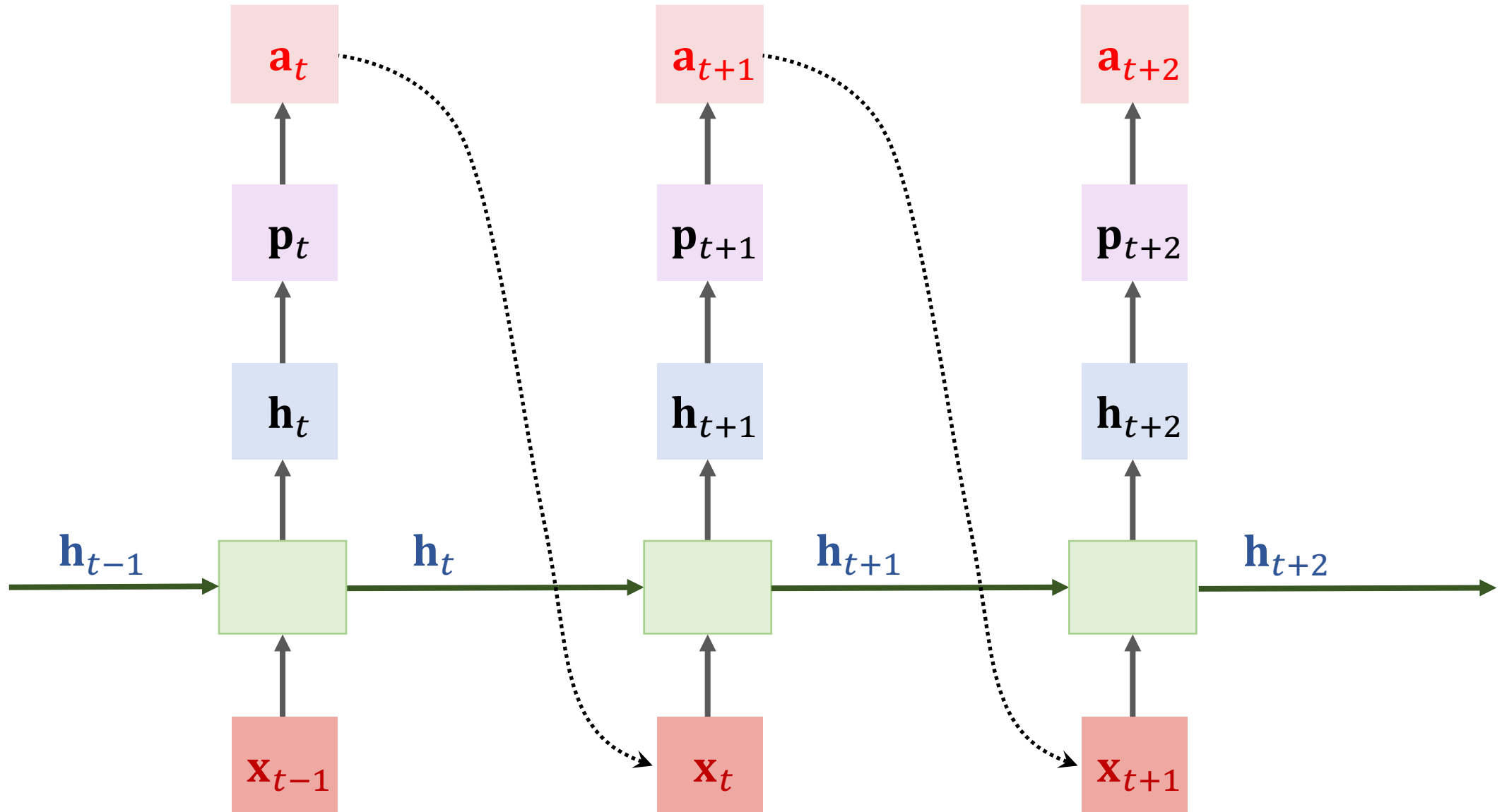


- Validation accuracy ( $r$ ) is **not** a differentiable function of the controller RNN parameters ( $\theta$ ).
- They have to use reinforcement learning.

# Reinforcement Learning

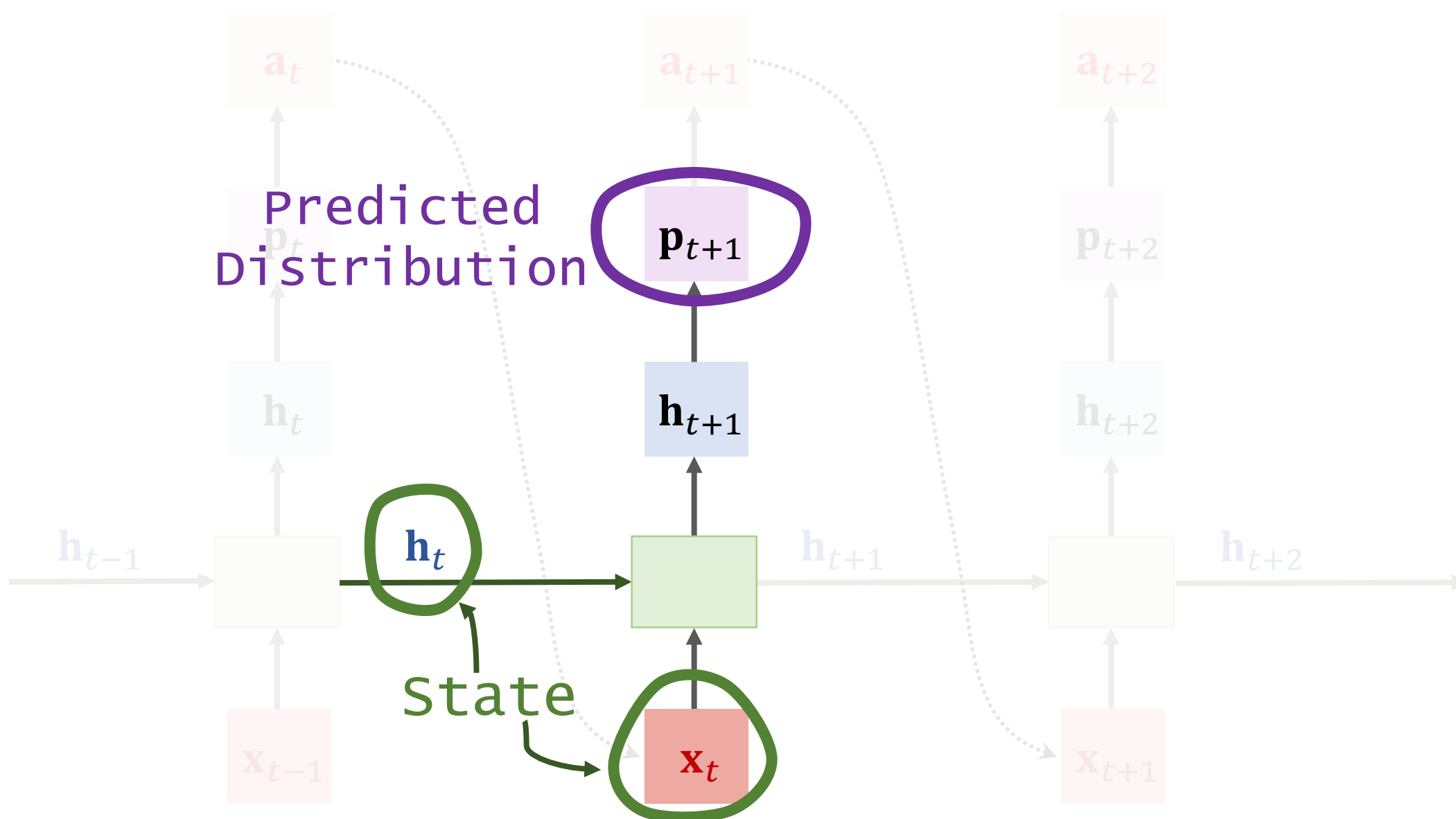
- **Objective:** Improve the controller RNN so that validation accuracies improve over time.
- **Rewards:** validation accuracies.
- **Policy function:** the controller RNN.
- Improve the policy function by **policy gradient** ascent.

# Policy Function

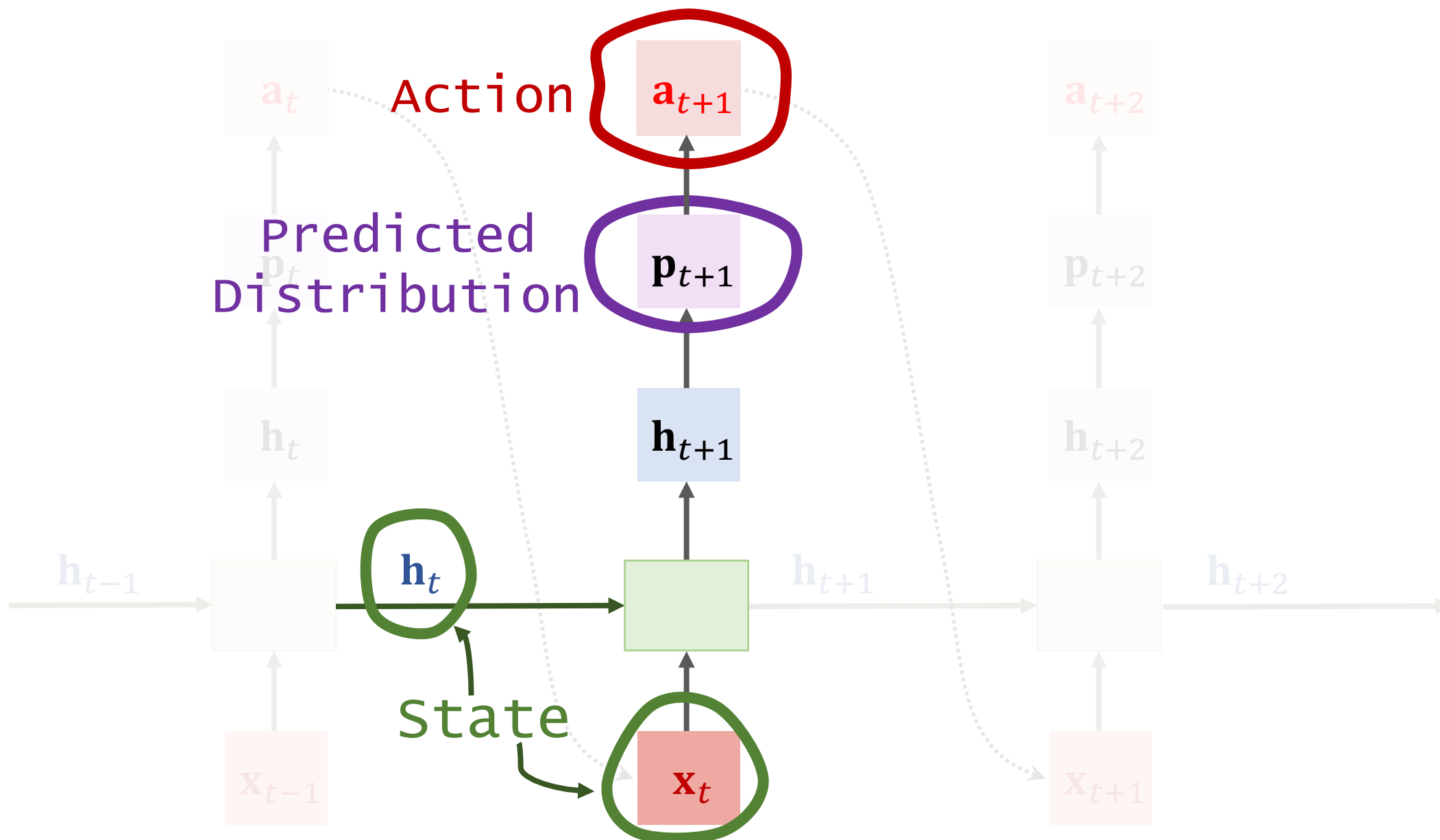




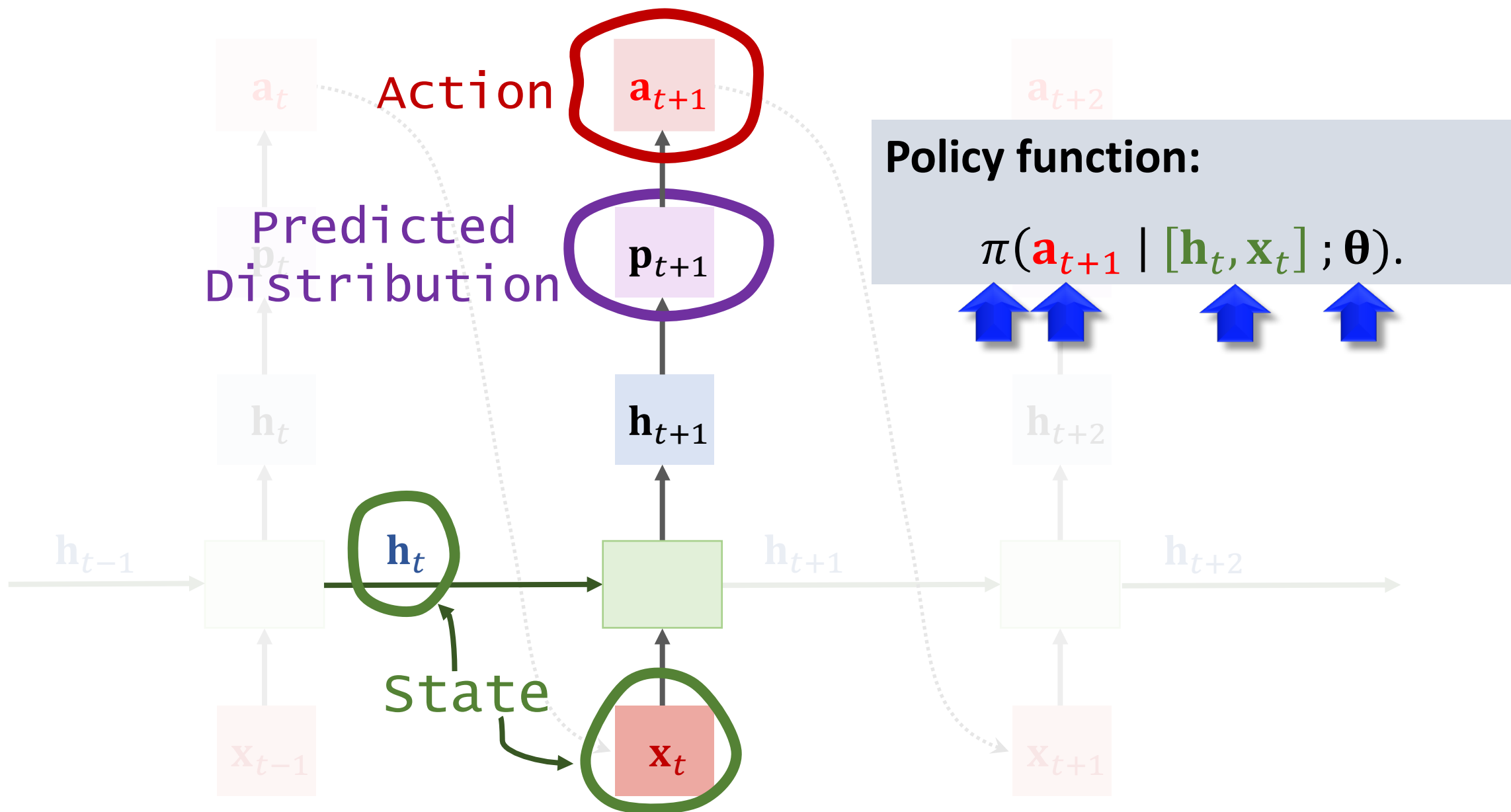
# Policy Function



# Policy Function

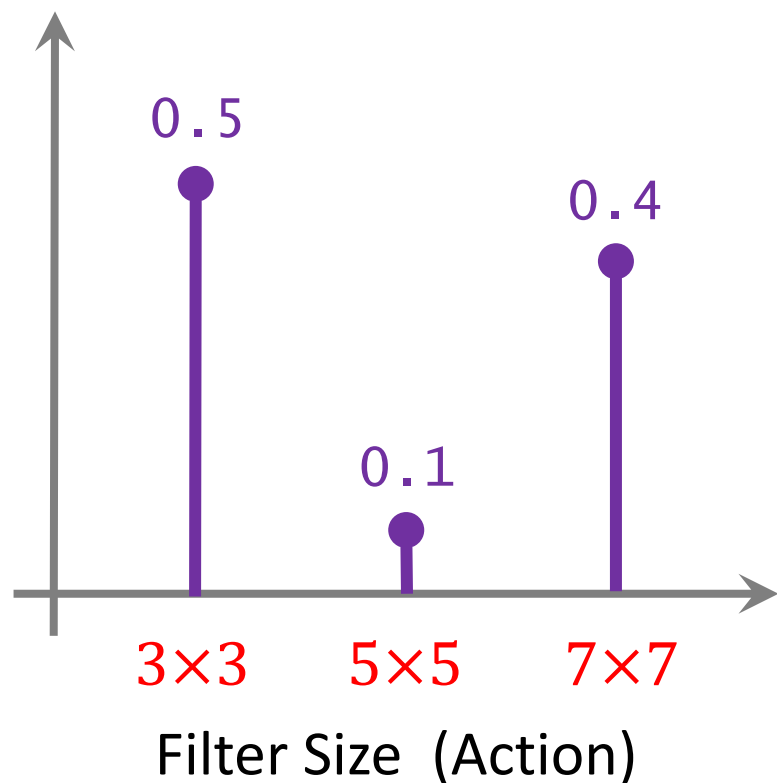


# Policy Function



# Policy Function

Probability  
Density



**Policy function:**

$$\pi(\mathbf{a}_{t+1} \mid [\mathbf{h}_t, \mathbf{x}_t]; \boldsymbol{\theta}).$$


- $\pi(\text{"3x3"} \mid [\mathbf{h}_t, \mathbf{x}_t]; \boldsymbol{\theta}) = 0.5.$
- $\pi(\text{"5x5"} \mid [\mathbf{h}_t, \mathbf{x}_t]; \boldsymbol{\theta}) = 0.1.$
- $\pi(\text{"7x7"} \mid [\mathbf{h}_t, \mathbf{x}_t]; \boldsymbol{\theta}) = 0.4.$

# Reward & Return


- Suppose the controller RNN runs 60 steps.
- The first 59 rewards are zeros:  $r_1 = r_2 = \dots = r_{59} = 0$ .
- The last reward is the validation accuracy:  $r_{60} = \text{ValAcc}$ .

# Reward & Return

- Suppose the controller RNN runs 60 steps.
- The first 59 rewards are zeros:  $r_1 = r_2 = \dots = r_{59} = 0$ .
- The last reward is the validation accuracy:  $r_{60} = \text{ValAcc}$ .
- Return (aka cumulative reward) is defined as:


$$u_t = r_t + r_{t+1} + r_{t+2} + \dots + r_{59} + r_{60}.$$


- Thus, all the returns are equal:

$$u_1 = u_2 = \dots = u_{60} = \text{ValAcc}.$$


# REINFORCE Algorithm

- Approximate policy gradients by:

$$\frac{\partial \log \pi(\mathbf{a}_{t+1} \mid [\mathbf{h}_t, \mathbf{x}_t]; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot u_t.$$


# REINFORCE Algorithm

- Approximate policy gradients by:

$$\frac{\partial \log \pi(\mathbf{a}_{t+1} \mid [\mathbf{h}_t, \mathbf{x}_t]; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot u_t.$$

- Update trainable parameters by:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \beta \cdot \sum_{t=1}^{60} \frac{\partial \log \pi(\mathbf{a}_{t+1} \mid [\mathbf{h}_t, \mathbf{x}_t]; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot u_t.$$



# Recap

- Run the controller RNN to generate the hyper-parameters of the 20 convolutional layers.
- Instantiate a CNN, train the CNN, and then obtain a validation accuracy (to be used as a **reward**.)
- REINFORCE algorithm uses the **reward** to update the policy function (i.e., the controller RNN.)
- Repeat this process thousands of times.

# NAS is expensive!

- To update the controller RNN once, we need to train a CNN from scratch. (Once is already expensive.)
- 10,000+ updates to train the RNN well → Train 10,000+ CNNs from scratch. (Extremely expensive!)
- The controller RNN itself has tuning hyper-parameters.
  - E.g., # of layers, size of  $\mathbf{x}$ , size of  $\mathbf{h}$ , etc.
  - Hyper-parameter tuning will make the overall time cost many times higher!

# Thank You!

<http://wangshusen.github.io/>