

Recurrent Neural Networks: LSTMs & GRUs

Antoine Bosselut

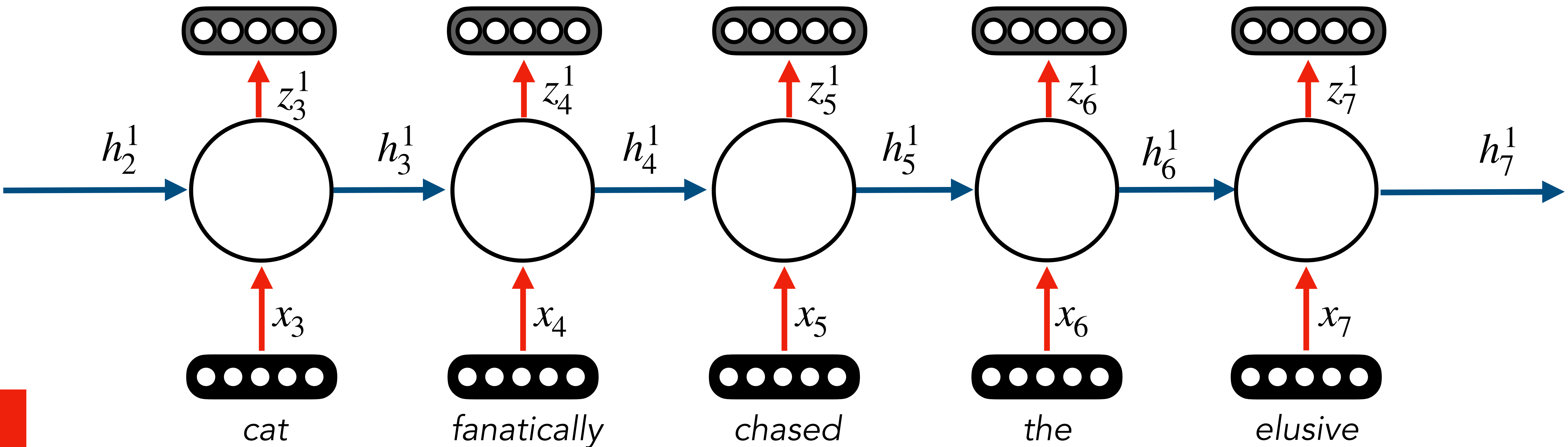
Section Outline

- **Recurrent neural networks continued:** Multiple Layers, Bidirectionality
- **Recurrent neural network variants:** LSTMs, GRUs
- **Recurrent neural networks in practice:** sequence-to-sequence models

Last Week Recap

- Early neural language models (and n-gram models) suffer from **fixed context windows**
- Recurrent neural networks can **theoretically** learn to model an **unbounded context length** using back propagation through time (BPTT)
- Practically, **vanishing gradients** can still stop many RNNs from learning **long-range dependencies**

RNNs in Practice



Classification

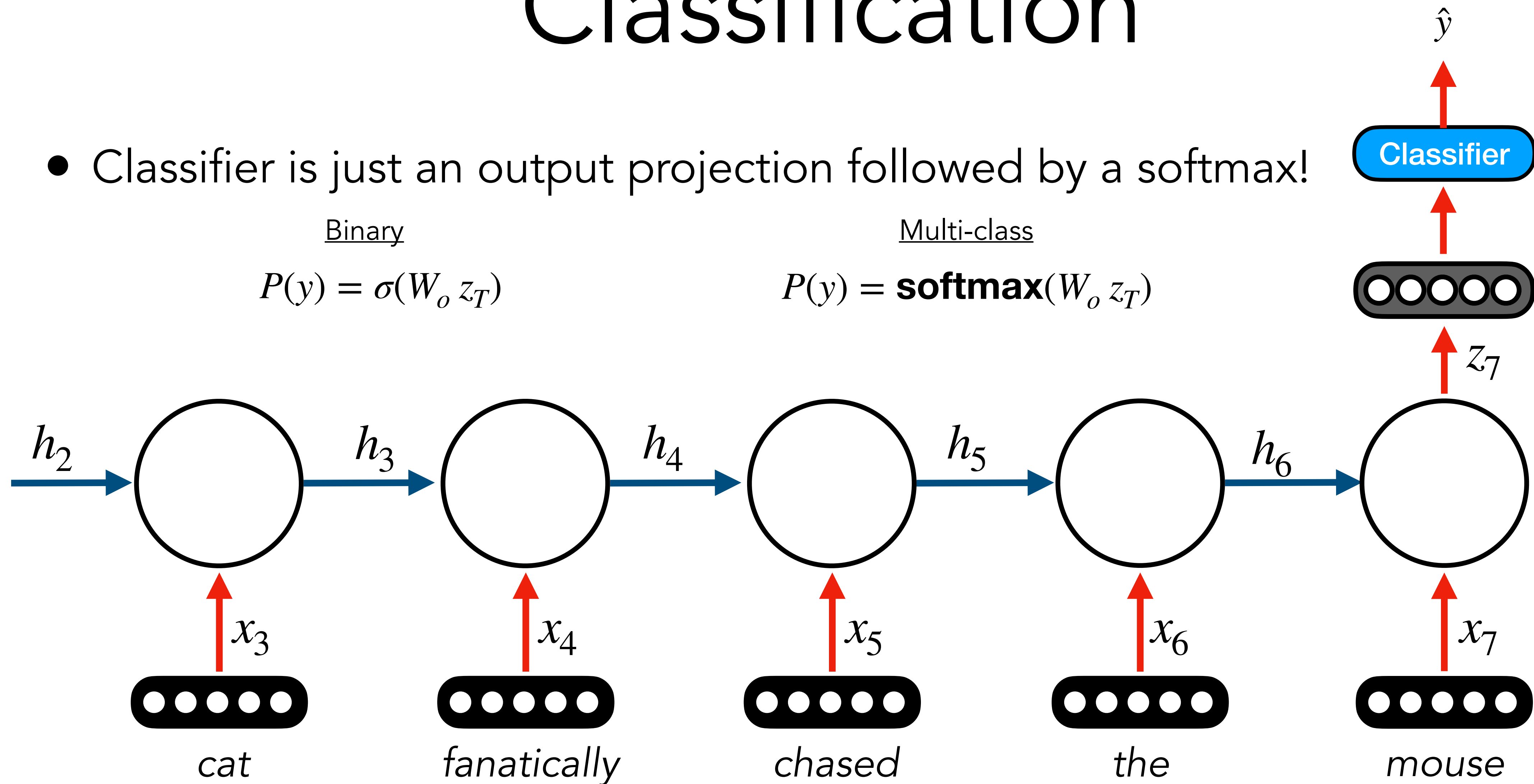
- Classifier is just an output projection followed by a softmax!

Binary

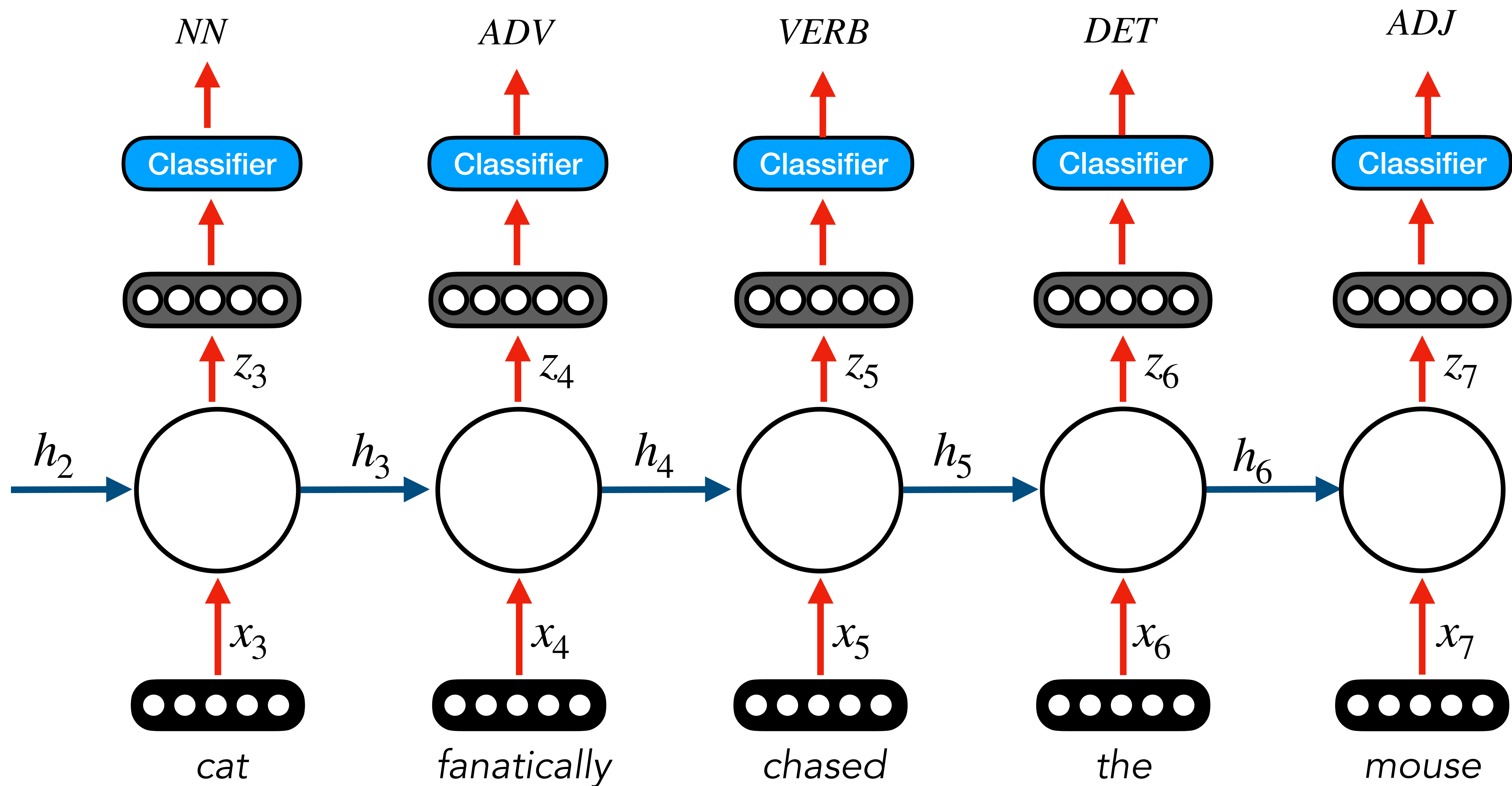
$$P(y) = \sigma(W_o z_T)$$

Multi-class

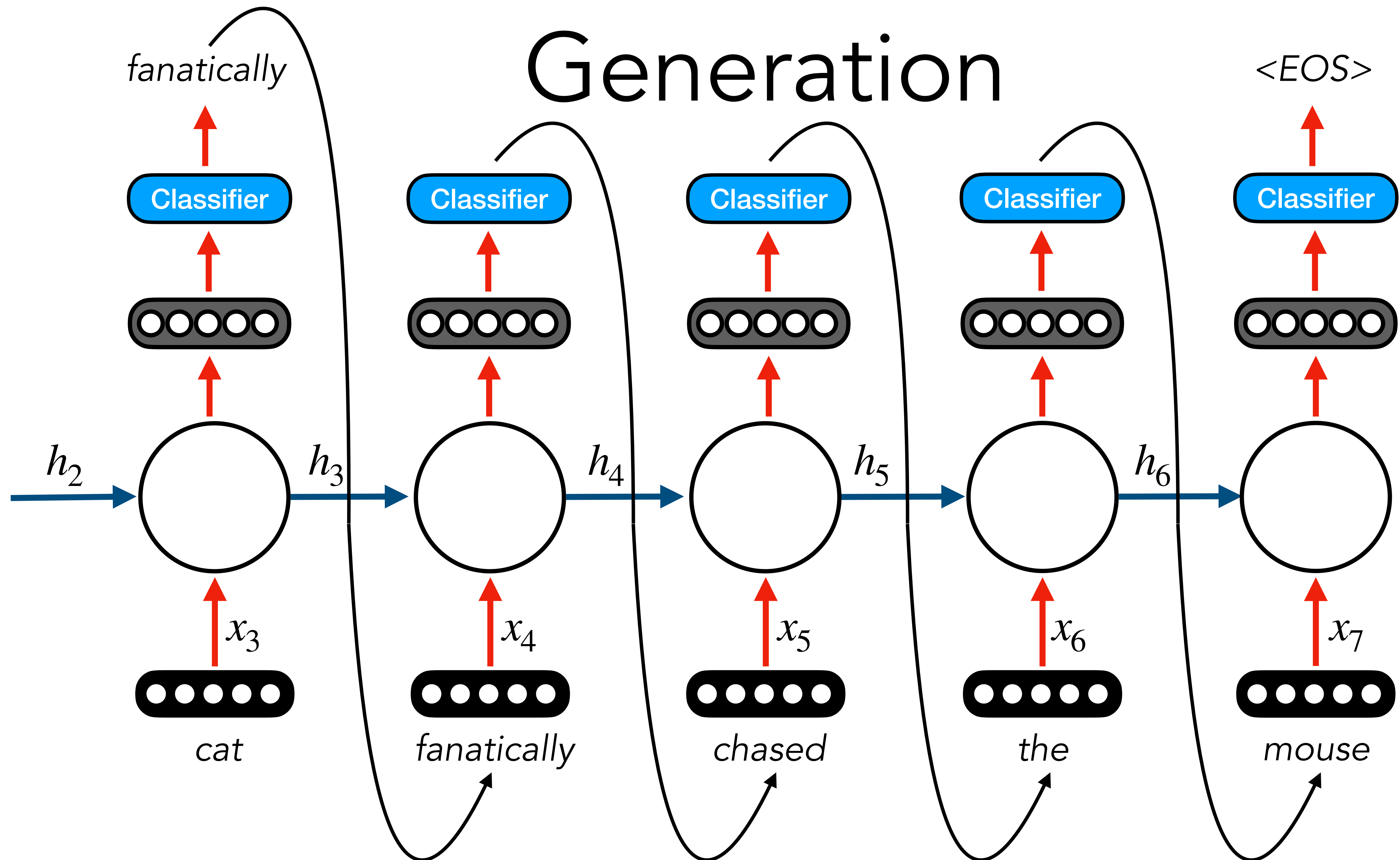
$$P(y) = \mathbf{softmax}(W_o z_T)$$



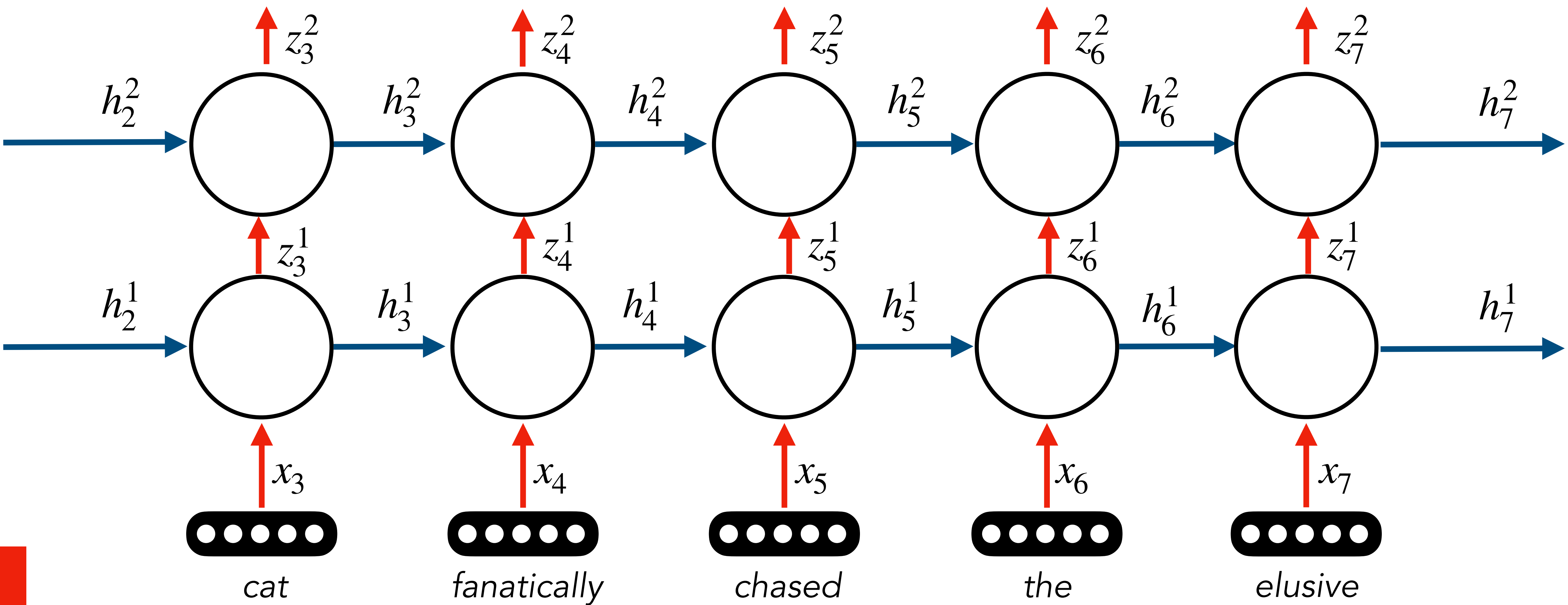
Sequence Labeling

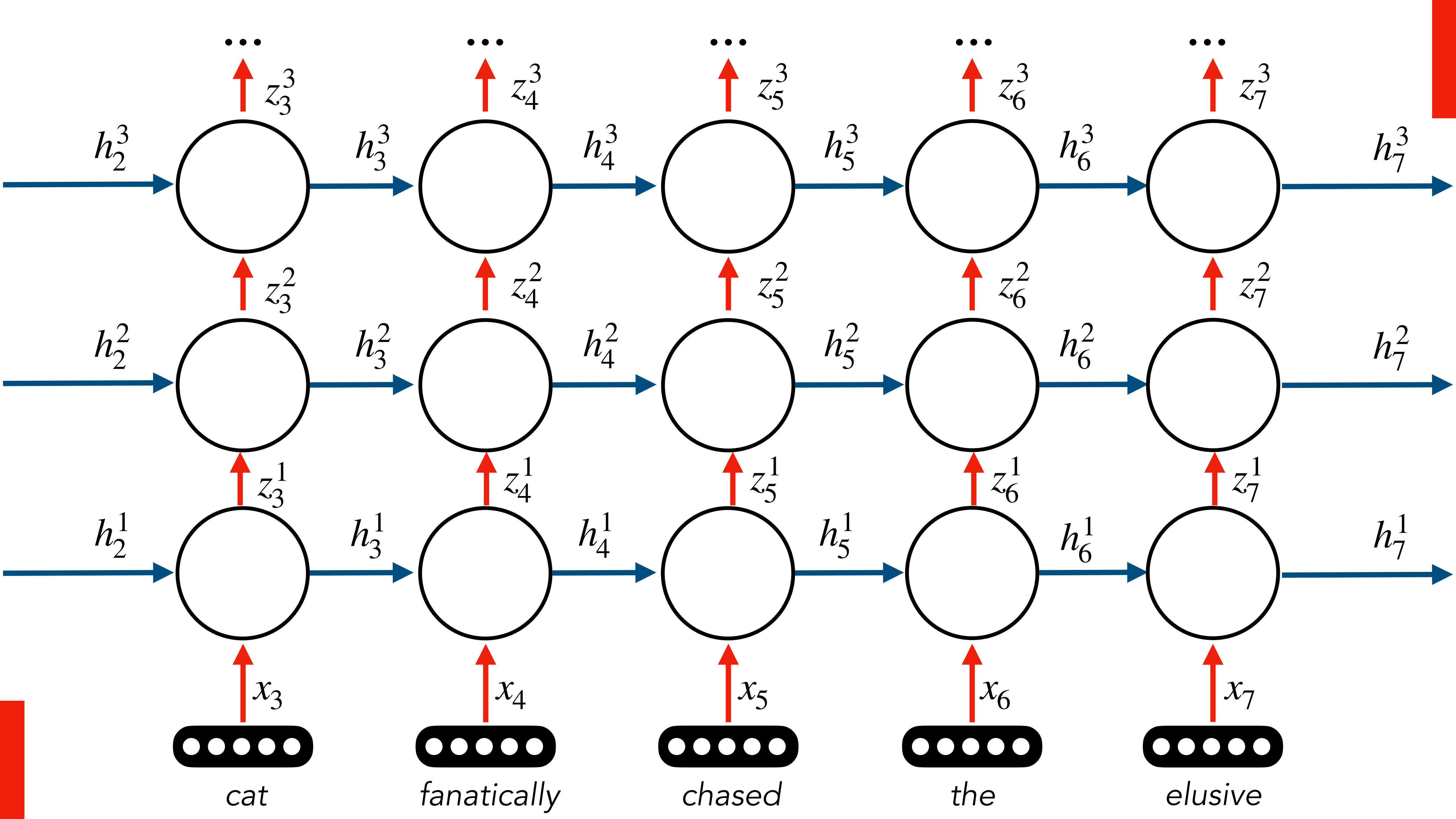


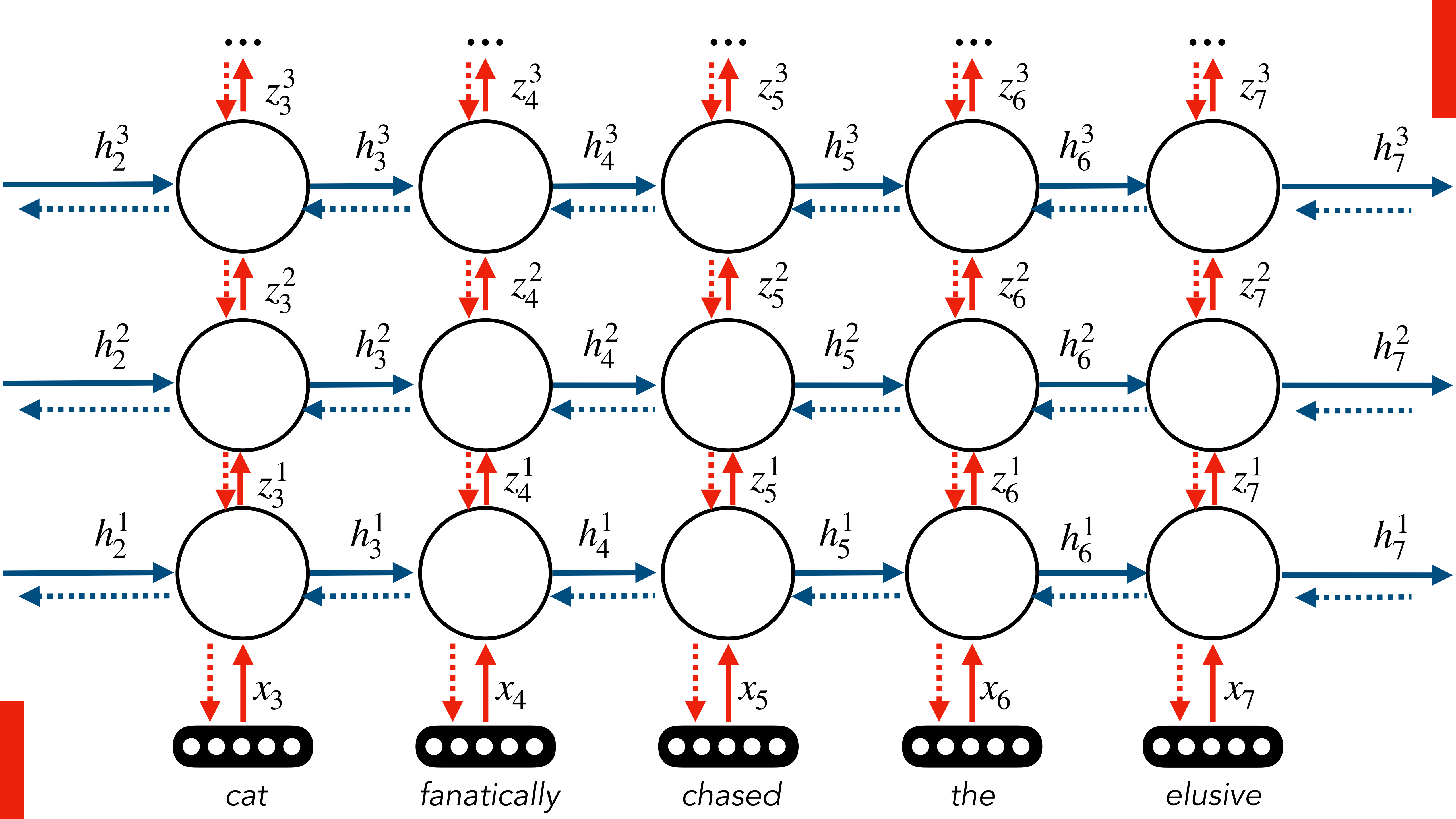
Generation

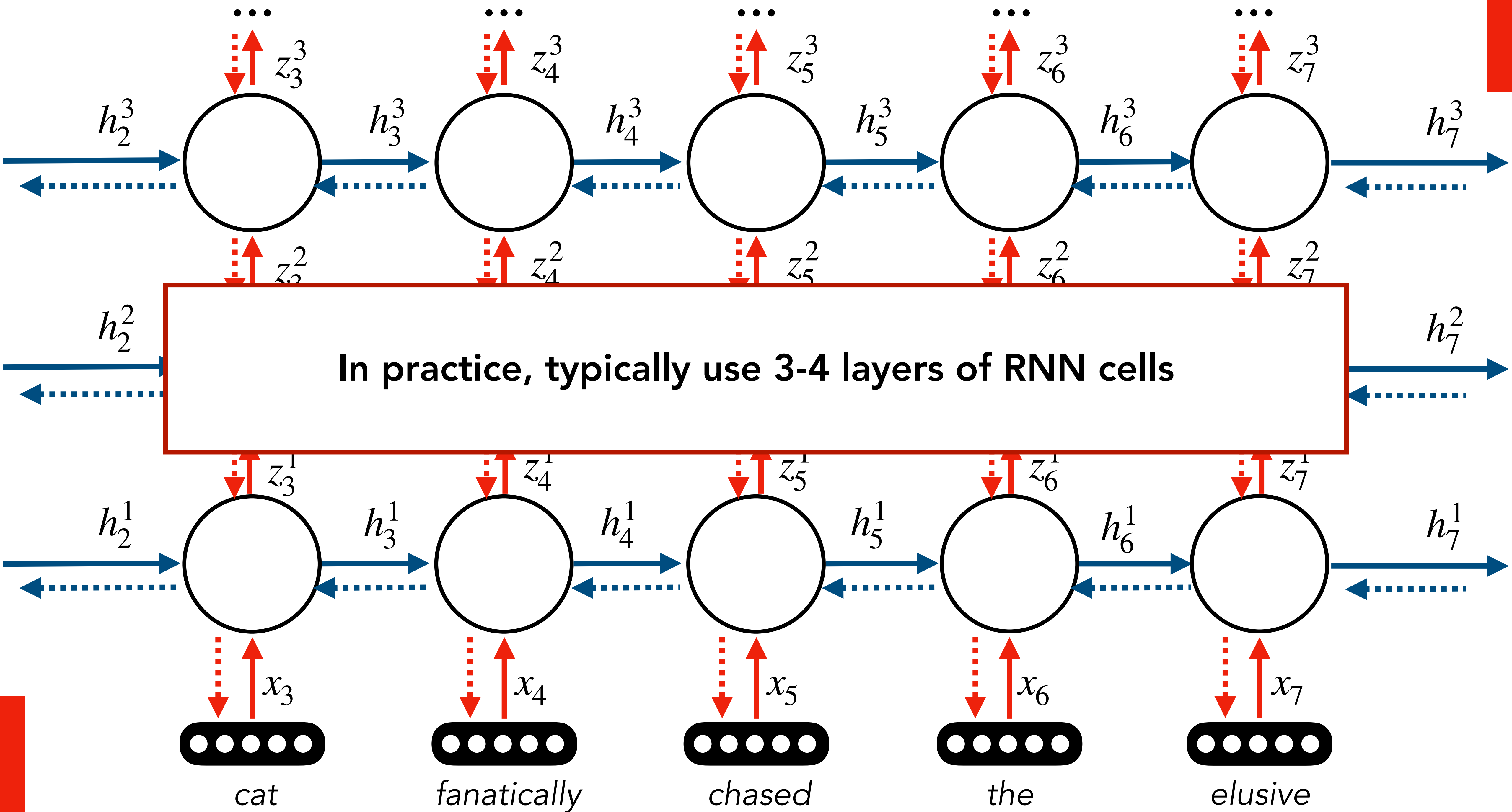


Multiple Layers

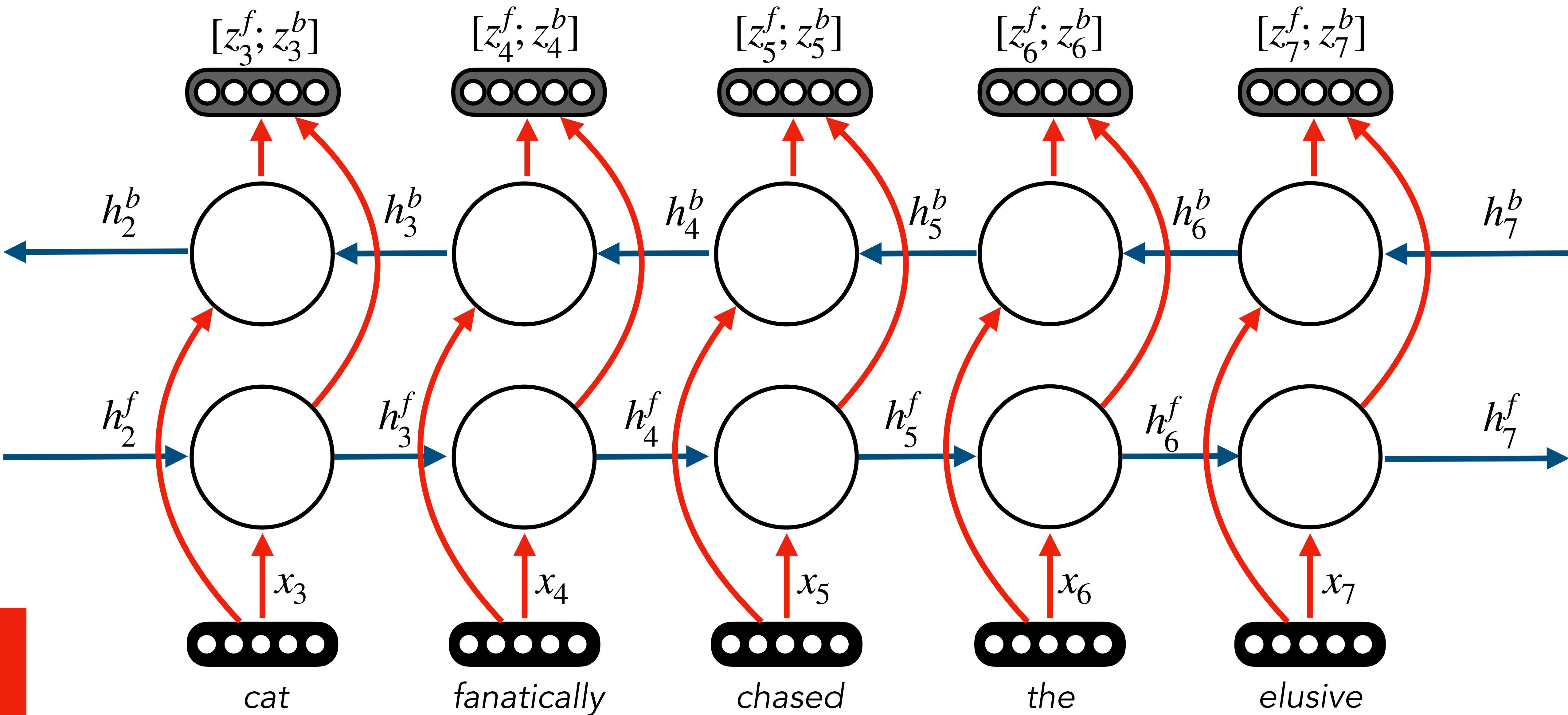




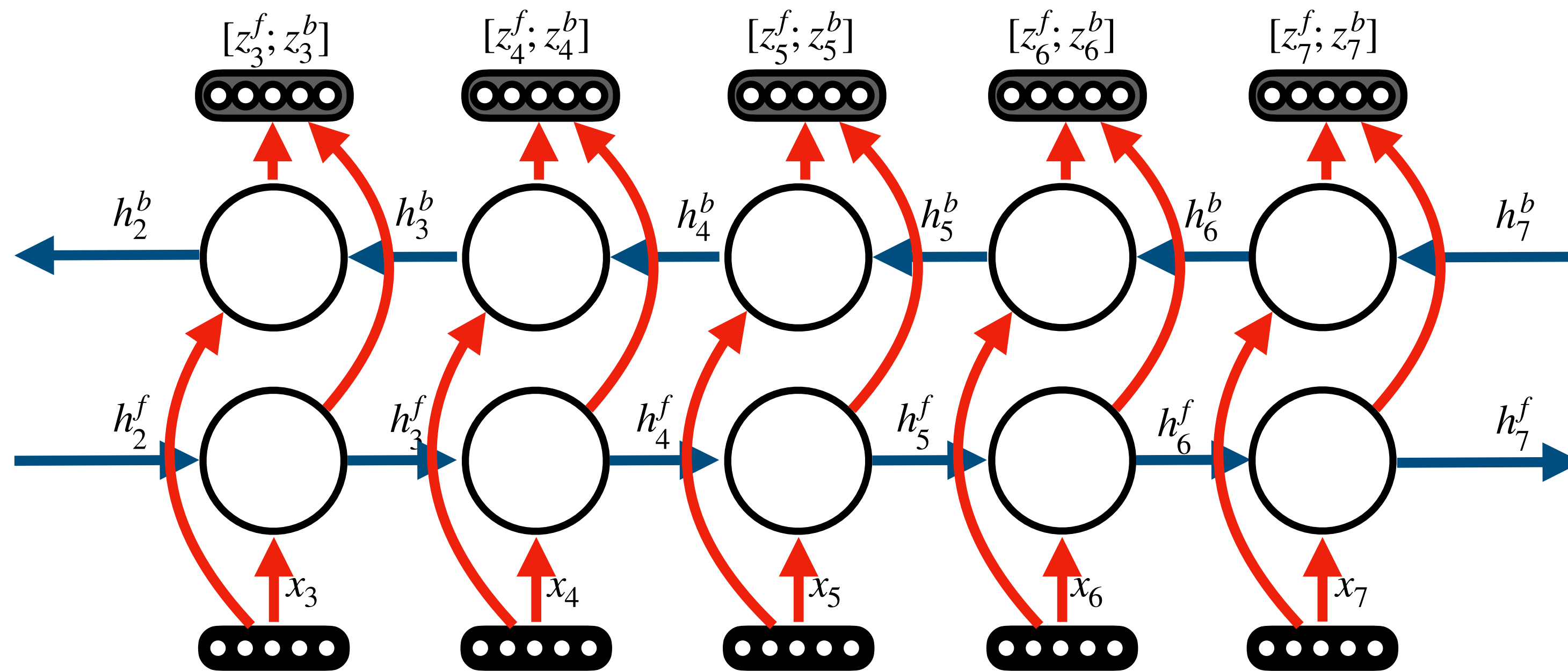




Bidirectionality



Bidirectionality



- **Concatenate** the output states for final representation at each step
 - Can also do mean / max
- Separate parameters for forward and backward RNNs
- If you can use the future text for the task, then you should use **bidirectionality**

Question

**For which of the following task types
can we use a bidirectional RNN?**

- (a) Classification**
- (b) Sequence labelling**
- (c) Text Generation**

Issue with Recurrent Models

- Multiple steps of state overwriting makes it challenging to learn long-range dependencies.

*They tuned, discussed for a moment, then struck up a lively **jig**. Everyone joined in, turning the courtyard into an even more chaotic scene, people now **dancing** in circles, **swinging** and **spinning** in circles, everyone making up their own **dance steps**. I felt my feet tapping, my body wanting to move. Aside from writing, I 've always loved dancing .*

- Nearby words should affect each other more than farther ones, but RNNs make it challenging to learn any long-range interactions

Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

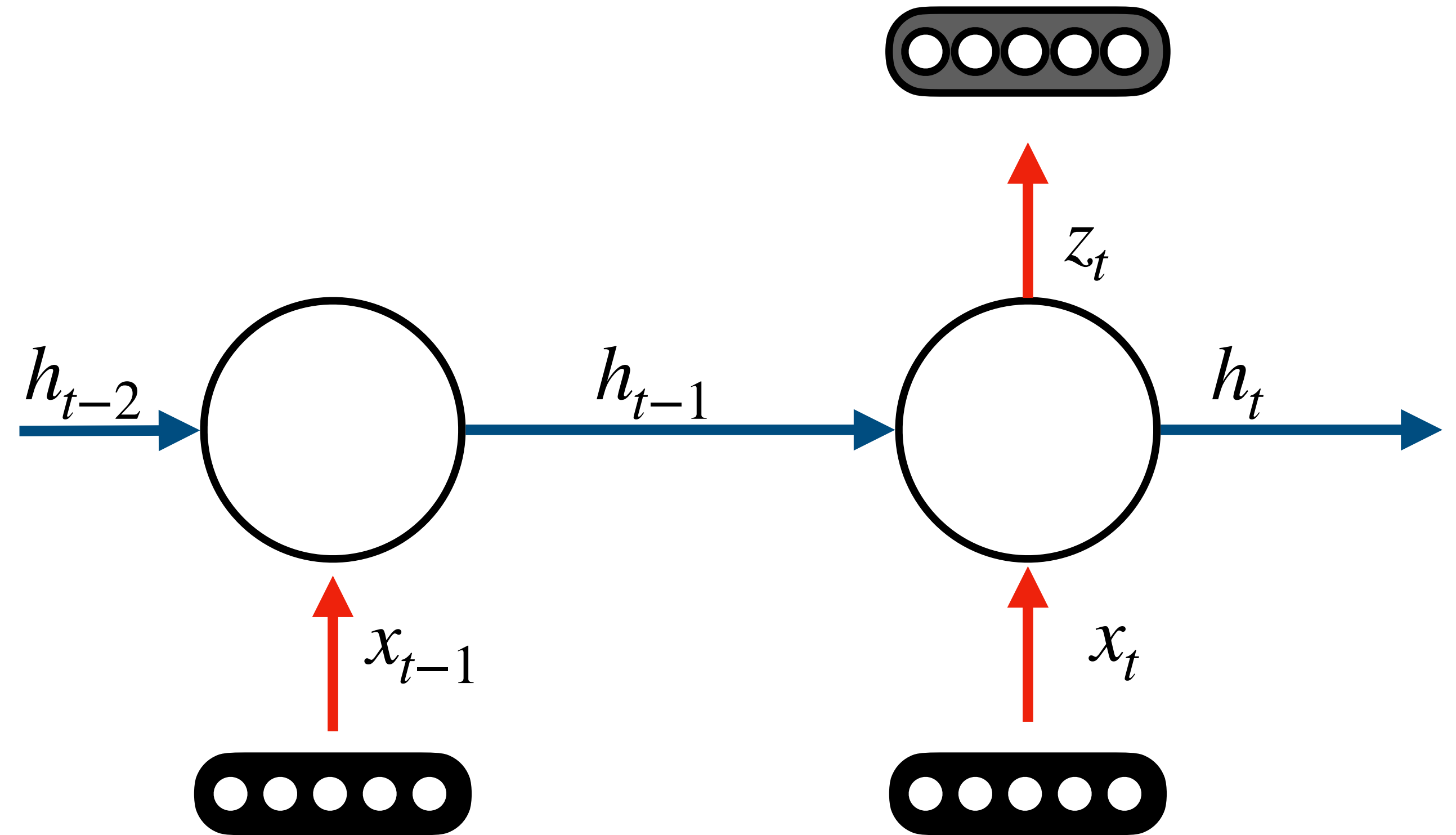
$$v_t = W_{zh}h_t + b_z \quad z_t = \sigma(v_t)$$

$$u_t = W_{hx}x_t + W_{hh}h_{t-1} + b_h \quad h_t = \sigma(u_t)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} \frac{\partial v_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} W_{zh}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} \frac{\partial u_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh}$$

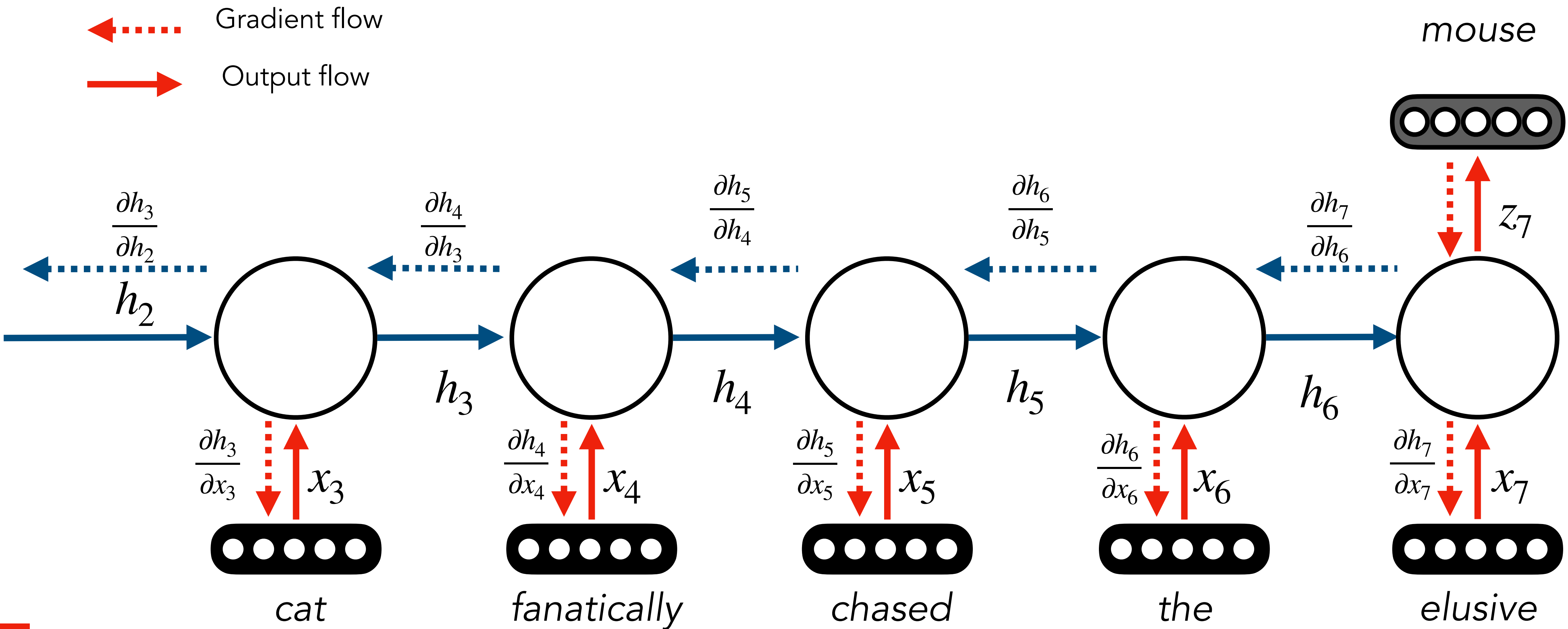
$$\frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \frac{\partial u_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$



$$\frac{\partial z_t}{\partial h_{t-2}} = \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(v_t)}{\partial v_t} W_{zh} \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh} \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$

Note that these are actually
the same matrix

Backpropagation through time



Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

$$\frac{\partial z_t}{\partial h_{t-2}} = \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(v_t)}{\partial v_t} \textcolor{yellow}{W_{zh}} \frac{\partial \sigma(u_t)}{\partial u_t} \textcolor{pink}{W_{hh}} \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \textcolor{blue}{W_{hh}}$$

$$v_t = \textcolor{yellow}{W_{zh}}h_t + b_z \quad z_t = \sigma(v_t)$$

$$u_t = W_{hx}x_t + \textcolor{pink}{W_{hh}}h_{t-1} + b_h \quad h_t = \sigma(u_t)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} \textcolor{yellow}{\frac{\partial v_t}{\partial h_t}} = \frac{\partial \sigma(v_t)}{\partial v_t} \textcolor{yellow}{W_{zh}}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} \textcolor{pink}{\frac{\partial u_t}{\partial h_{t-1}}} = \frac{\partial \sigma(u_t)}{\partial u_t} \textcolor{pink}{W_{hh}}$$

$$\frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \textcolor{blue}{\frac{\partial u_{t-1}}{\partial h_{t-2}}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \textcolor{blue}{W_{hh}}$$

Generalising this:

$$\frac{\partial h_t}{\partial h_{t-T}} = \prod_{i=t-T}^{i=t} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=t-T}^{i=t} \frac{\partial \sigma(u_i)}{\partial u_i} \textcolor{blue}{W_{hh}}$$

Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

$$\frac{\partial z_t}{\partial h_{t-2}} = \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(v_t)}{\partial v_t} \boxed{W_{zh}} \frac{\partial \sigma(u_t)}{\partial u_t} \boxed{W_{hh}} \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \boxed{W_{hh}}$$

$$v_t = \boxed{W_{zh}}h_t + b_z \quad z_t = \sigma(v_t)$$

$$u_t = W_{hx}x_t + \boxed{W_{hh}}h_{t-1} + b_h \quad h_t = \sigma(u_t)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} \boxed{\frac{\partial v_t}{\partial h_t}} = \frac{\partial \sigma(v_t)}{\partial v_t} \boxed{W_{zh}}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} \boxed{\frac{\partial u_t}{\partial h_{t-1}}} = \frac{\partial \sigma(u_t)}{\partial u_t} \boxed{W_{hh}}$$

$$\frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \boxed{\frac{\partial u_{t-1}}{\partial h_{t-2}}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \boxed{W_{hh}}$$

Generalising this:

$$\frac{\partial h_t}{\partial h_{t-T}} = \prod_{i=t-T}^{i=t} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=t-T}^{i=t} \frac{\partial \sigma(u_i)}{\partial u_i} \boxed{W_{hh}}$$

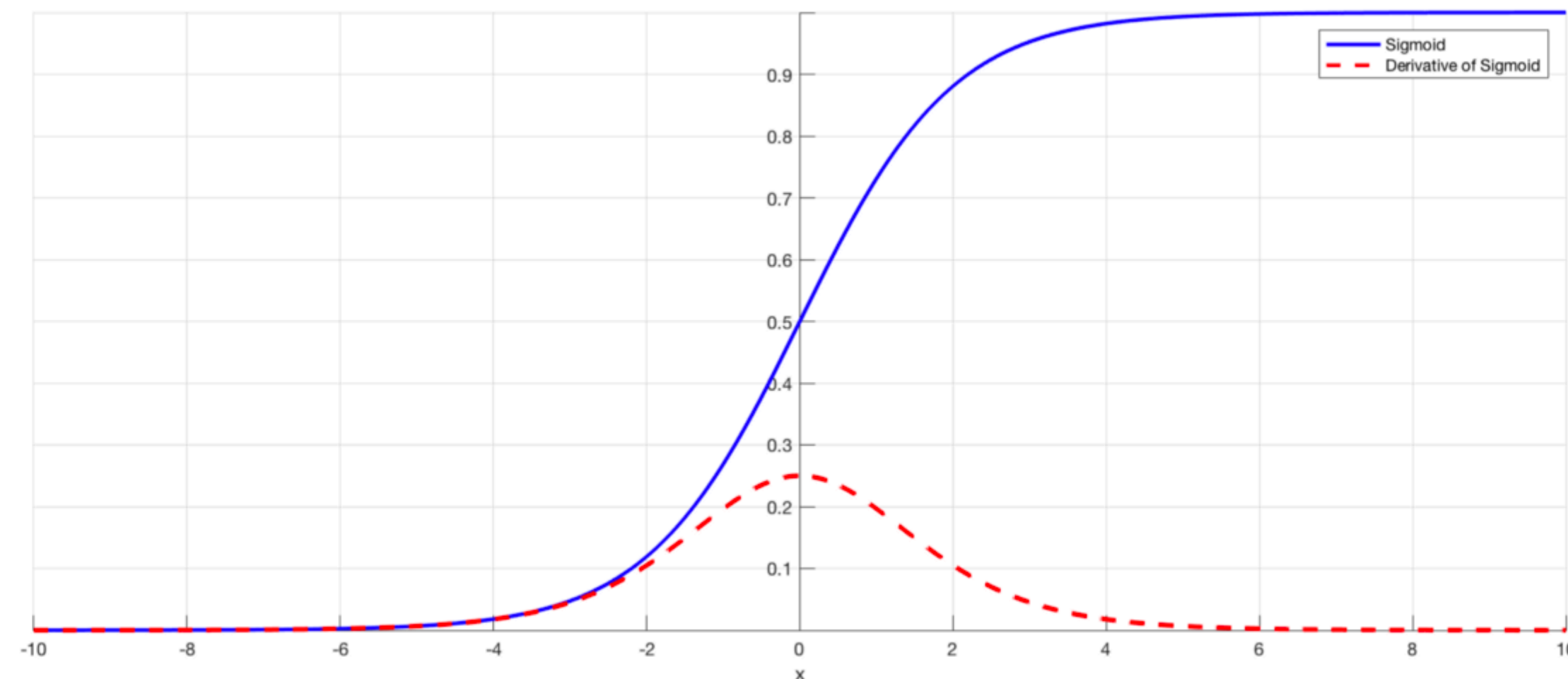
< 1 for many
activation fxns

Typically small
(Regularisation)

Vanishing Gradients

- **Learning Problem:** Long unrolled networks will crush gradients that backpropagate to earlier time steps

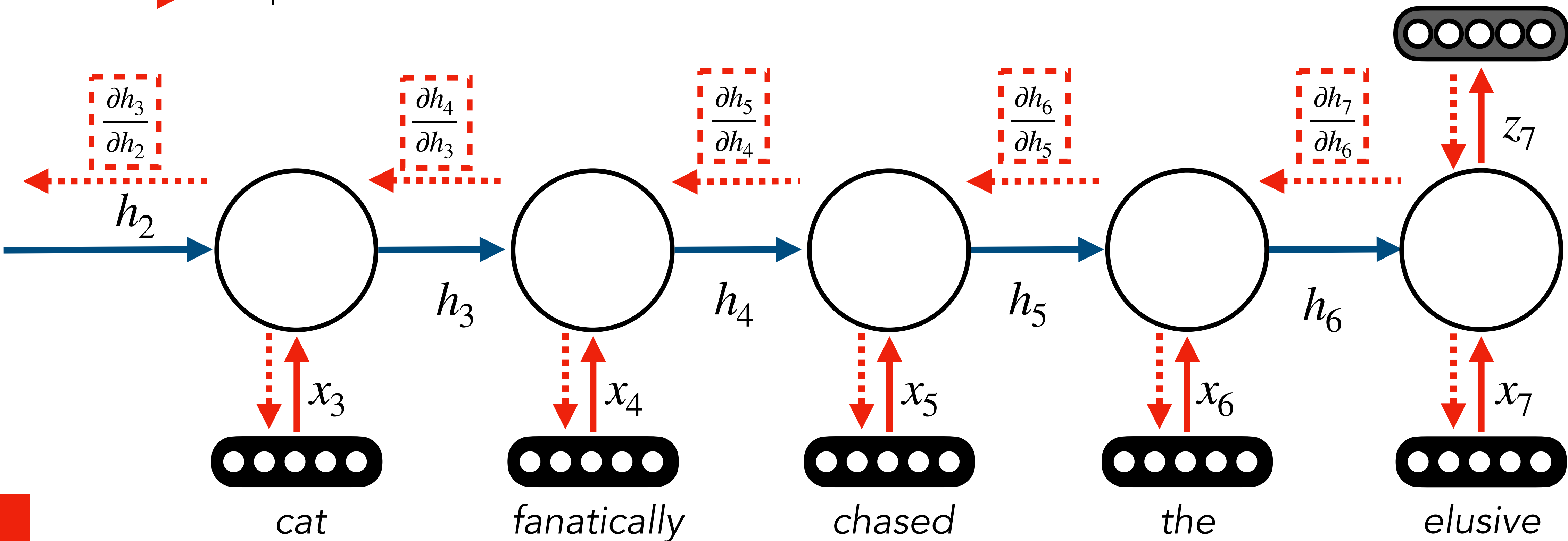
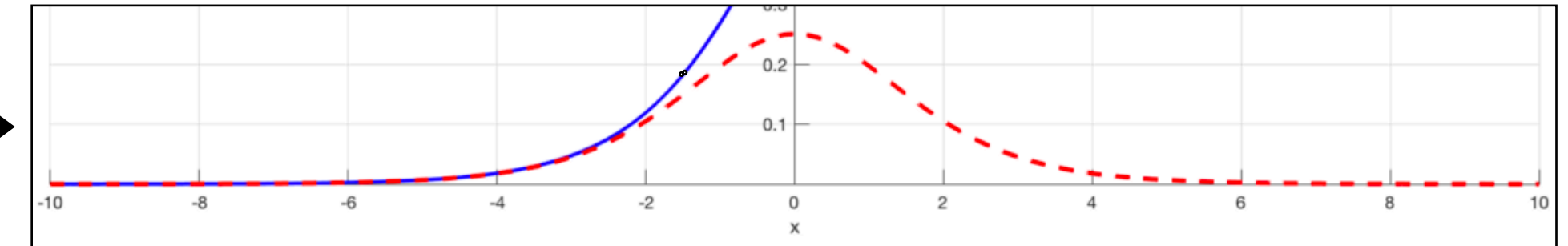
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$
$$u = W_{hx}x_t + W_{hh}h_{t-1} + b_h$$
$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = W_{hh} \frac{\partial \sigma(u)}{\partial u}$$



Vanishing Gradients

← Gradient flow
→ Output flow

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = W_{hh} \frac{\partial \sigma(u)}{\partial u}$$



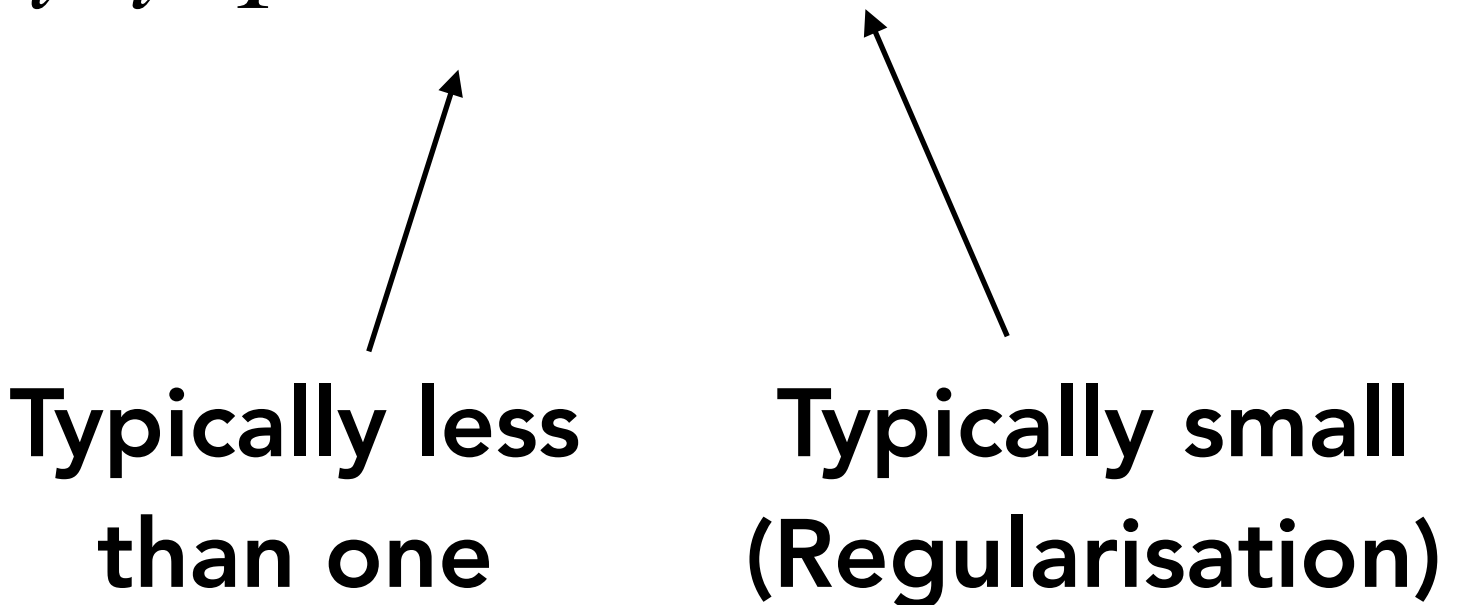
Question

How could we fix this vanishing gradient problem?

$$\frac{\partial h_t}{\partial h_{t-T}} = \prod_{i=t-T}^{i=t} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=t-T}^{i=t} \frac{\partial \sigma(u_i)}{\partial u_i} W_{hh}$$

Typically less than one

Typically small (Regularisation)



Gated Recurrent Neural Networks

- Use gates to avoid dampening gradient signal every time step

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

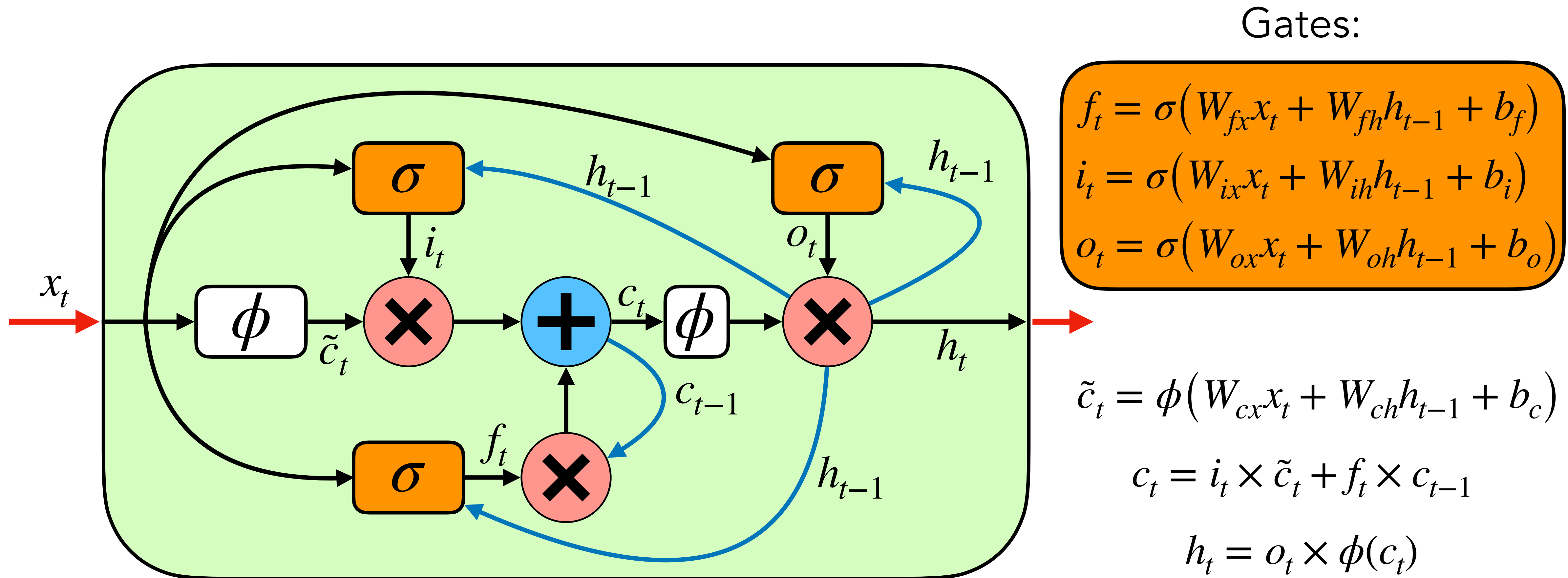
Elman Network

$$h_t = h_{t-1} \odot \mathbf{f} + \mathbf{func}(x_t)$$

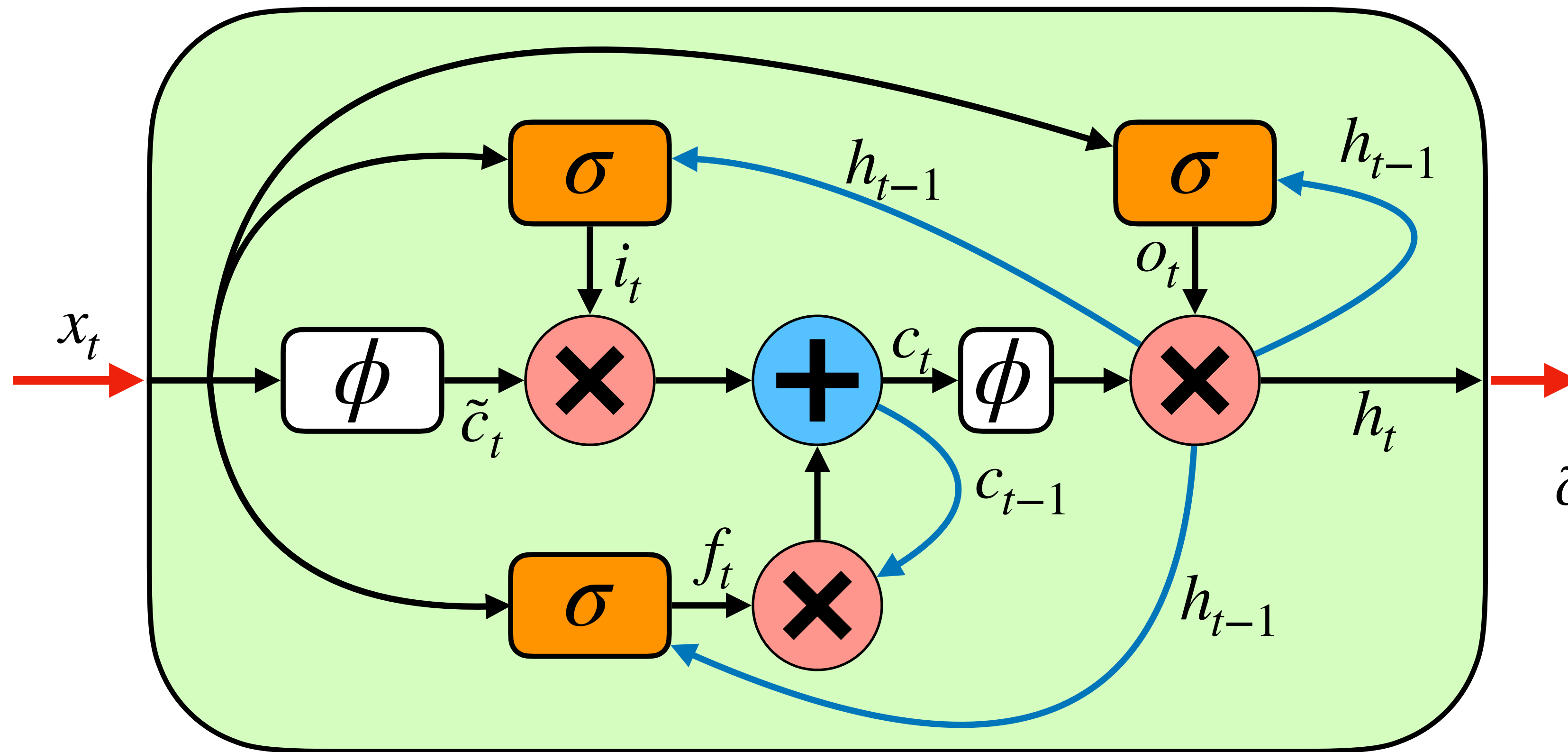
Gated Network Abstraction

- Gate value \mathbf{f} computes how much information from previous hidden state moves to the next time step $\rightarrow 0 < \mathbf{f} < 1$
- Because h_{t-1} is no longer inside the activation function, it is not automatically constrained, reducing vanishing gradients!

Long Short Term Memory (LSTM)



Long Short Term Memory (LSTM)

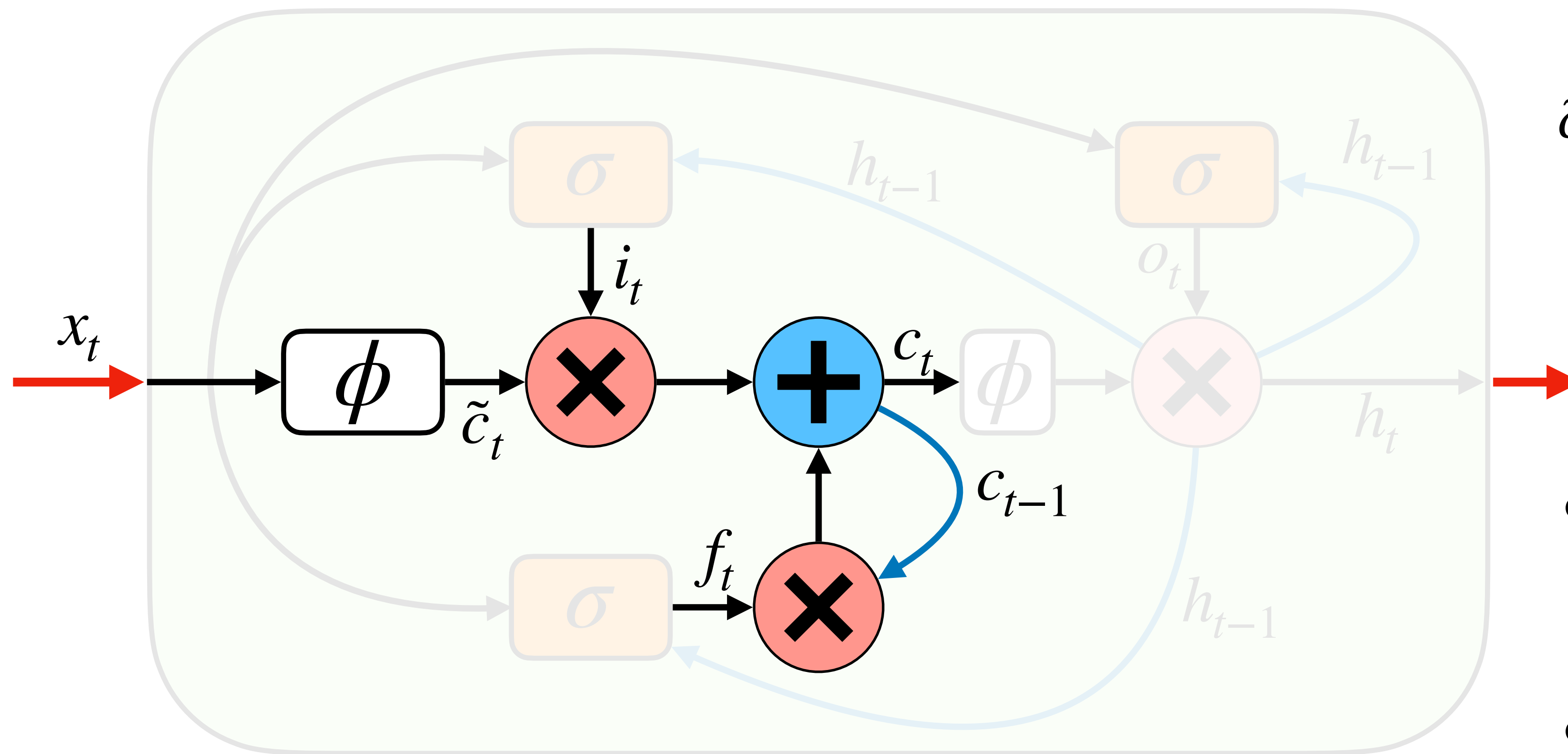


$$\tilde{c}_t = \phi(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$

$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

$$h_t = o_t \times \phi(c_t)$$

Cell State



$$\tilde{c}_t = \phi(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$

$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

- Hidden state h_{t-1} is now short-term memory
- Cell state c_t tracks longer-term dependencies

What does the cell state track?

- Can visualise the activations of cell state (i.e., dimensions of **c**) and find semantic behaviour!
- Stack Overflow example:

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

What does the cell state track?

- Can visualise the activations of cell state (i.e., dimensions of **c**) and find semantic behaviour!
- Stack Overflow example: **track indentation**

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

What does the cell state track?

- Can visualise the activations of cell state (i.e., dimensions of **c**) and find semantic behaviour!
- Stack Overflow example: track indentation

```
#ifdef CONFIG_AUDIT_SYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

- War and Peace:

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

What does the cell state track?

- Can visualise the activations of cell state (i.e., dimensions of **c**) and find semantic behaviour!
- Stack Overflow example: track indentation

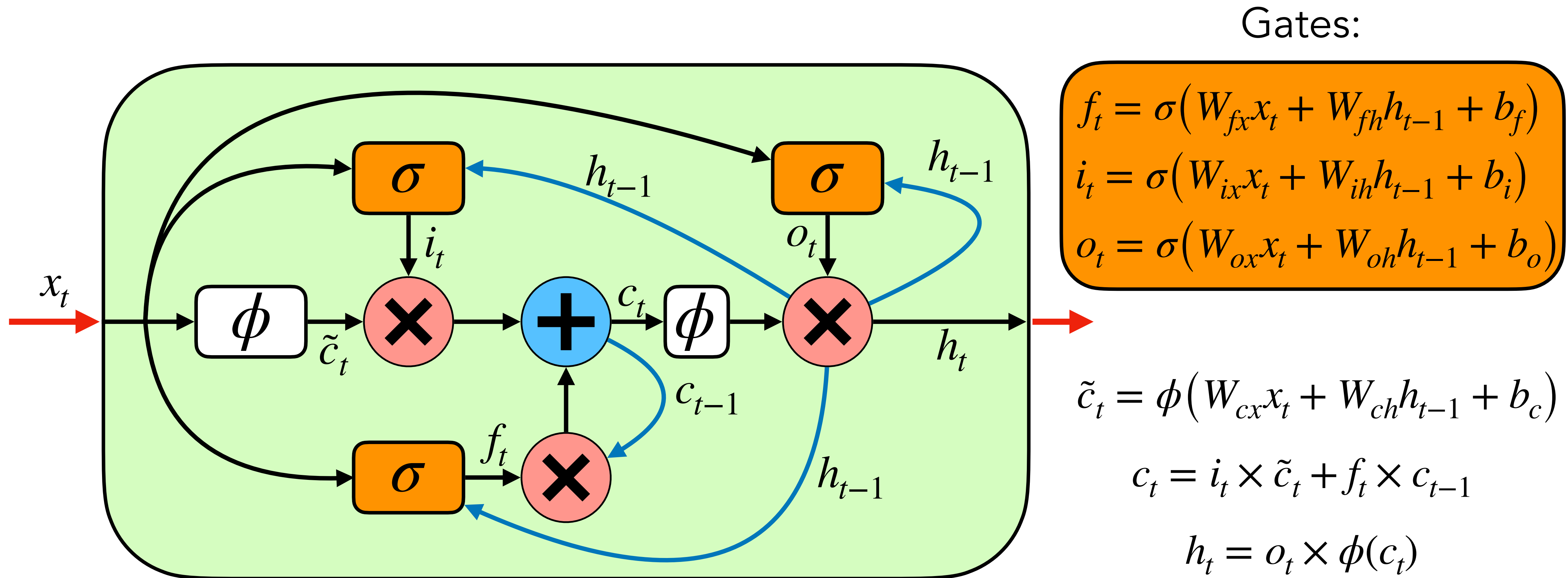
```
#ifdef CONFIG_AUDIT_SYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

- War and Peace: **are we in a quote or not?**

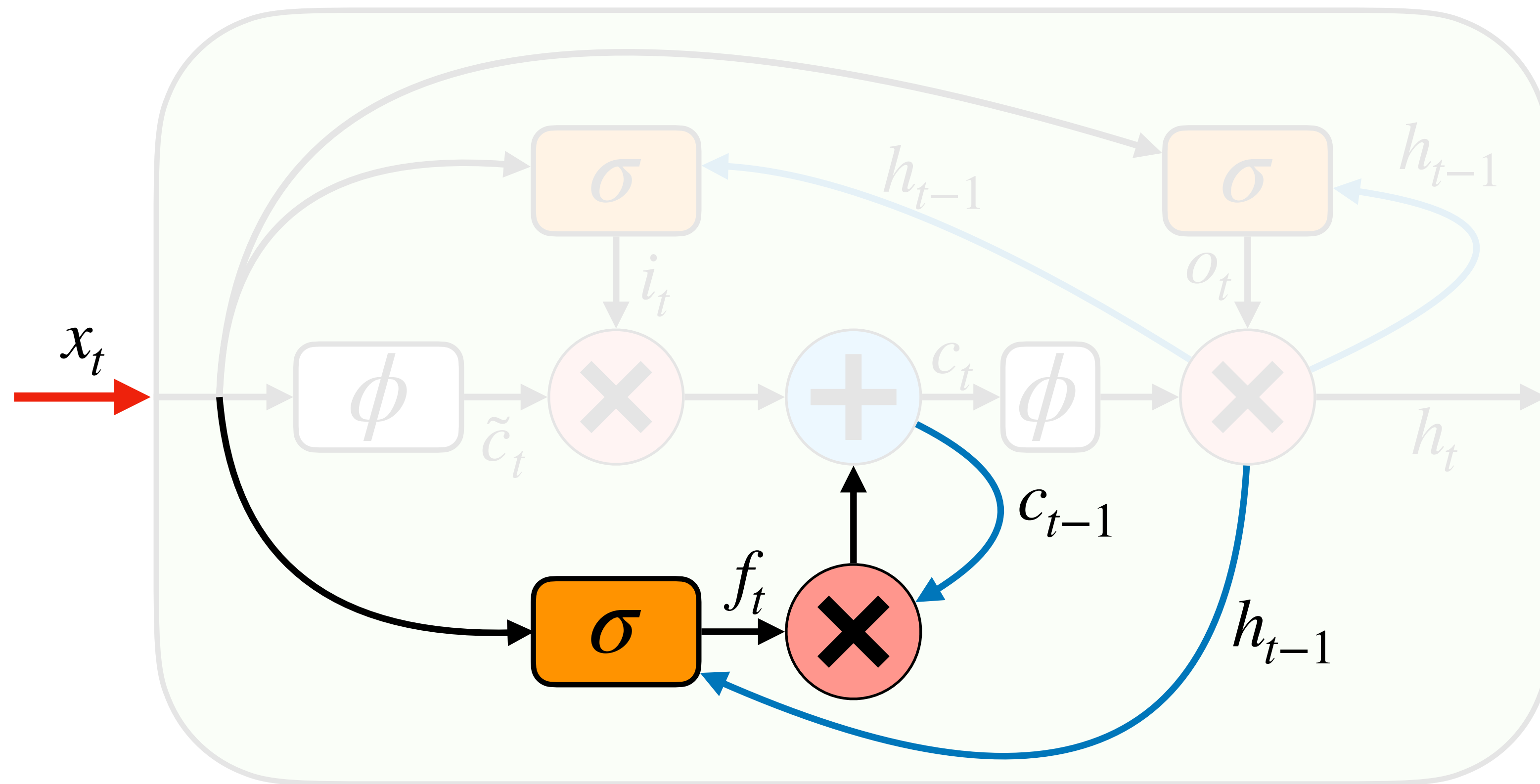
"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

Long Short Term Memory (LSTM)



Forget Gate



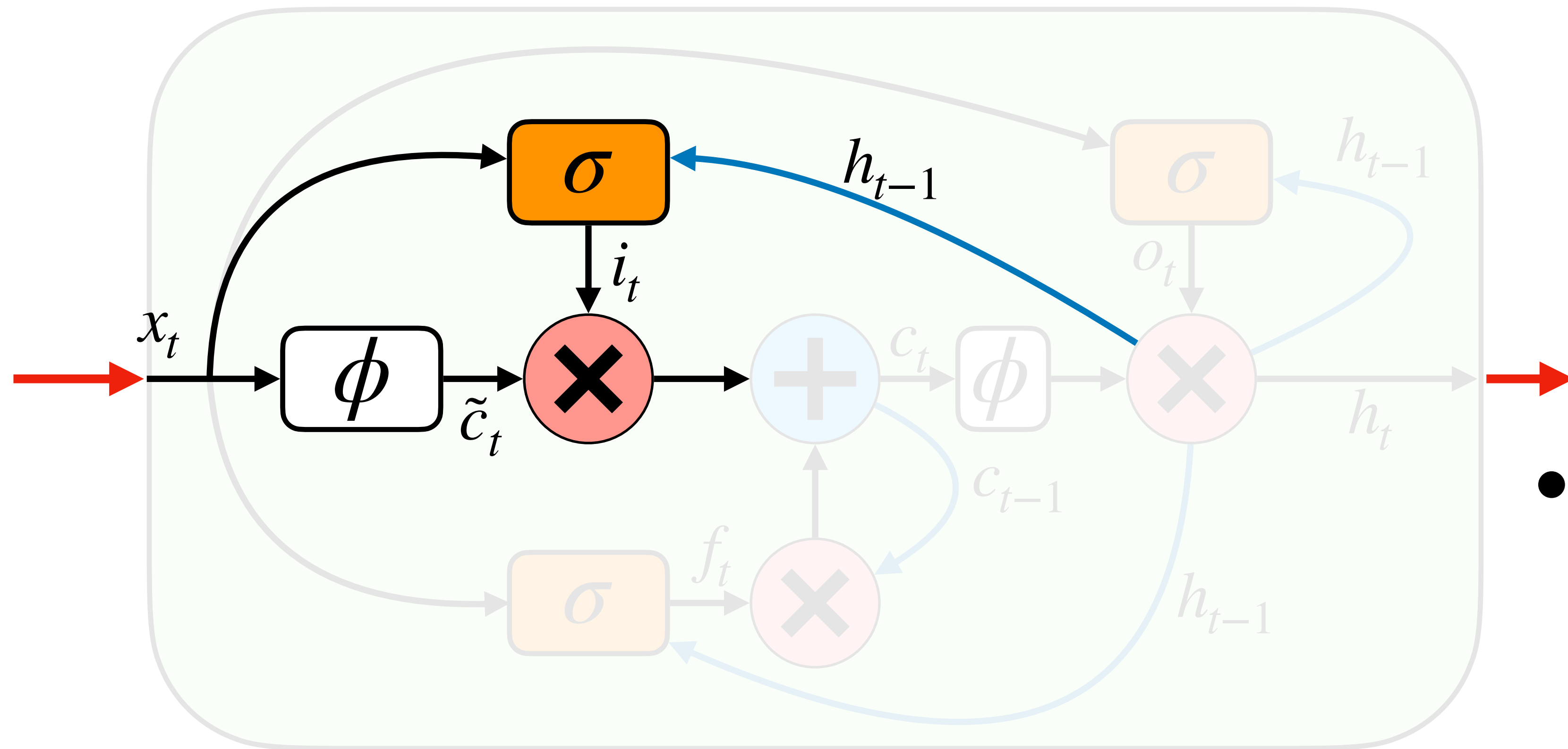
$$f_t = \sigma(W_{fx}x_t + W_{fh}h_{t-1} + b_f)$$

$$\tilde{c}_t = \phi(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$

$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

- Forget gate controls how much memory is forgotten
 - 1 -> remember the past
 - 0 -> forget everything up to now

Input Gate



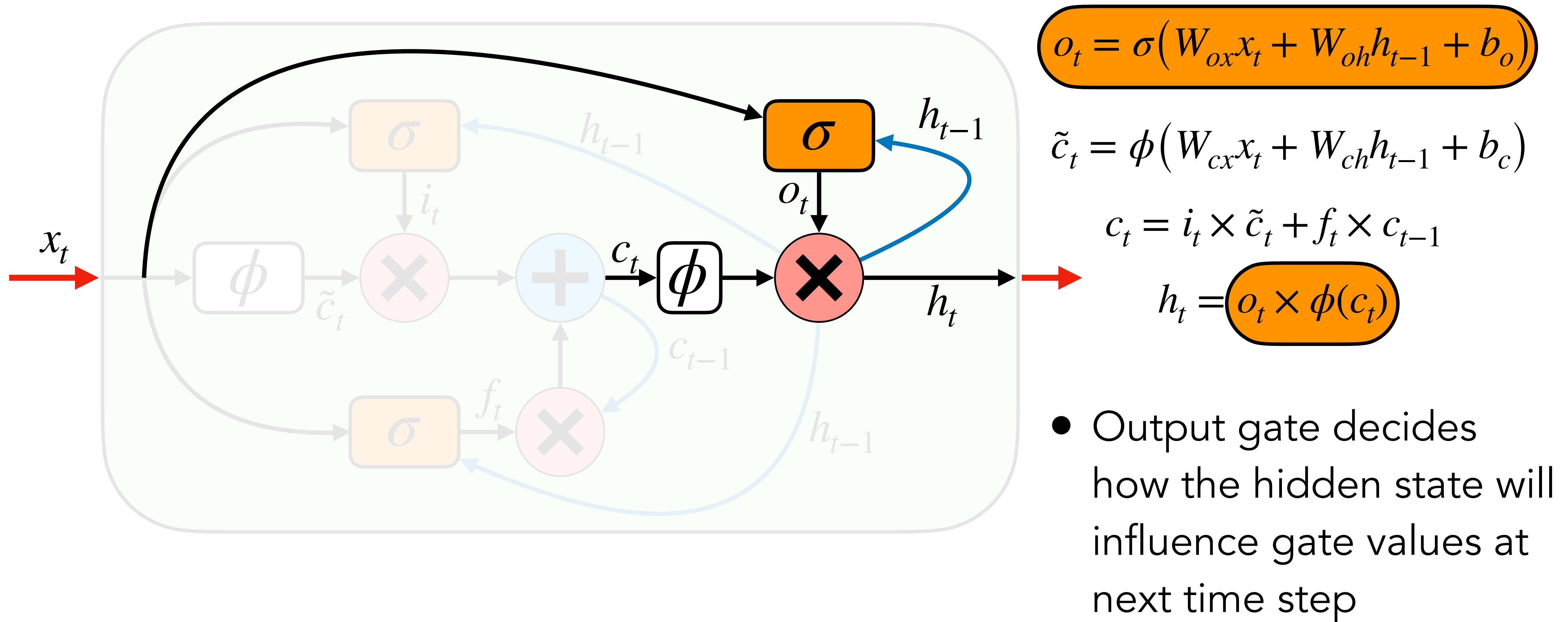
$$i_t = \sigma(W_{ix}x_t + W_{ih}h_{t-1} + b_i)$$

$$\tilde{c}_t = \phi(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$

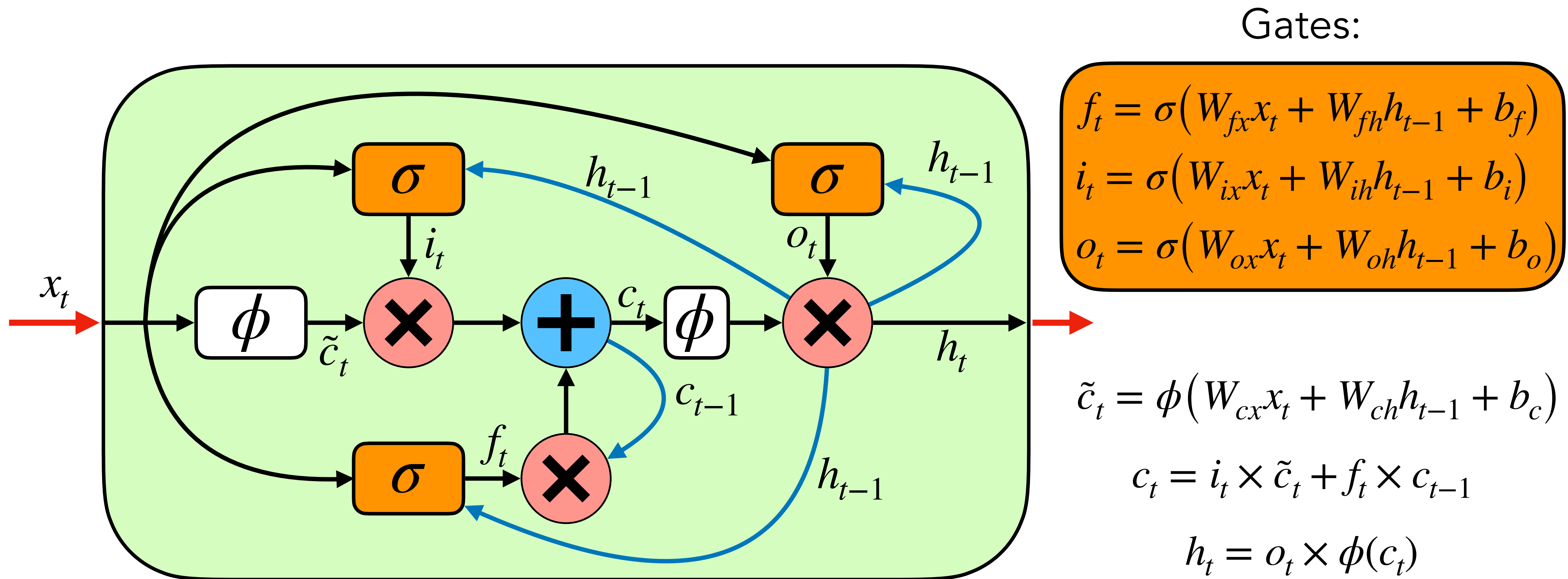
$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

- Input gate controls how new info is added to state memory
 - 0 -> ignore current time step
 - What might be ignored?

Output Gate



Long Short Term Memory (LSTM)



Questions!

- For what type of input might the model learn to make the input gate be 0 ?
- What happens if the forget gate is 0?
- What happens if both the forget gate and input gate are 0 ?
- What happens if both the forget gate and input gate are 1 ?

Gates:

$$\begin{aligned}f_t &= \sigma(W_{fx}x_t + W_{fh}h_{t-1} + b_f) \\i_t &= \sigma(W_{ix}x_t + W_{ih}h_{t-1} + b_i) \\o_t &= \sigma(W_{ox}x_t + W_{oh}h_{t-1} + b_o)\end{aligned}$$

$$\tilde{c}_t = \phi(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$

$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

$$h_t = o_t \times \phi(c_t)$$

Gated Recurrent Unit (GRU)

- Also uses gates to avoid dampening gradient signal every time step

$$h_t = (1 - \mathbf{z}) \odot h_{t-1} + \mathbf{z} \odot \mathbf{func}(x_t, h_{t-1}) \quad h_t = h_{t-1} \odot \mathbf{f} + \mathbf{func}(x_t)$$

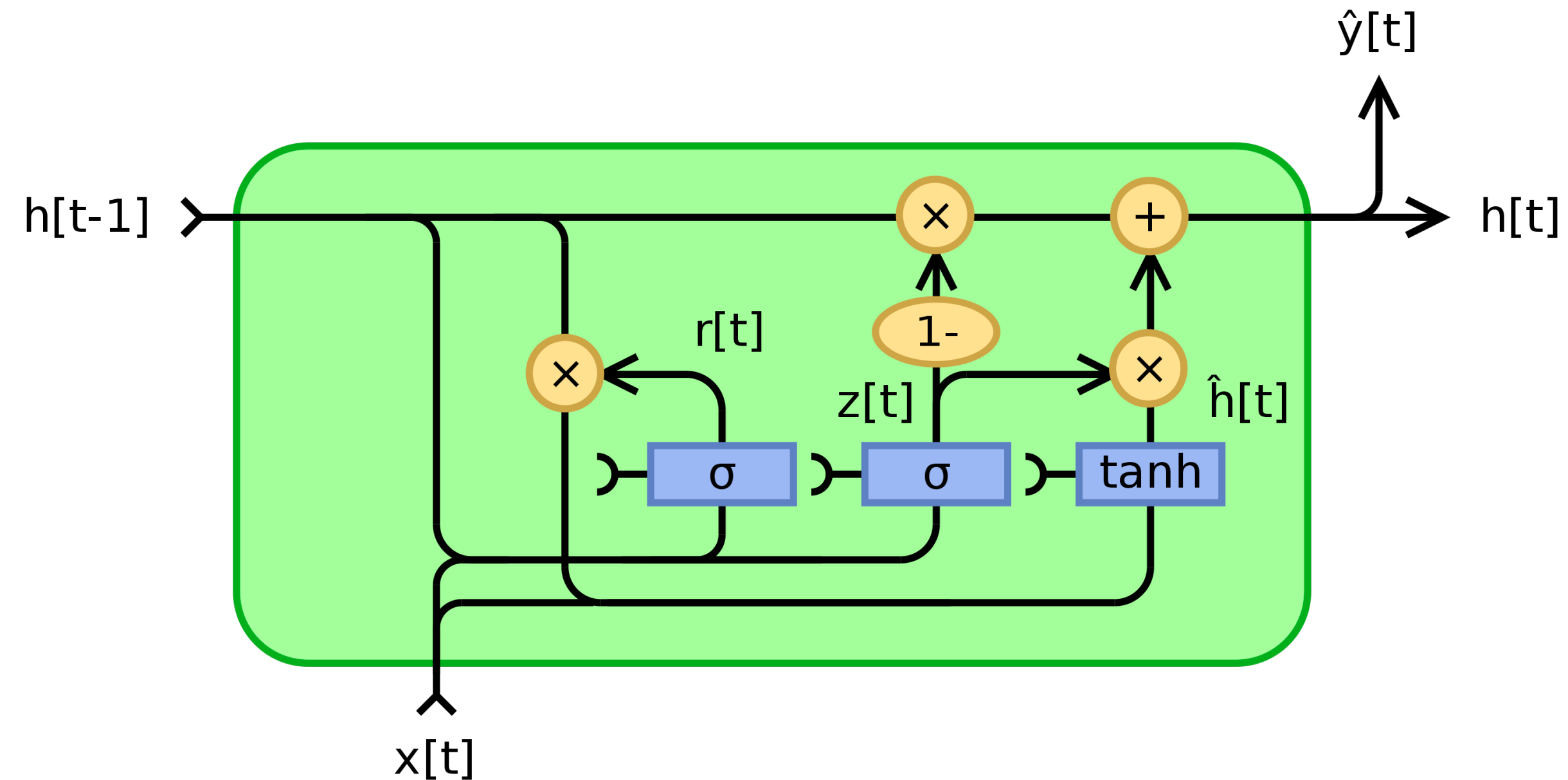
GRU

LSTM

- Works similarly to LSTM
 - Theoretically less powerful (find out why tomorrow!)
 - Typically faster to train and sometimes works better than LSTMs

Gated Recurrent Unit (GRU)

- \mathbf{z} is update gate (used to update hidden state), \mathbf{r} is reset gate (used to reset hidden state)
- The single hidden state and simpler update gate gives simpler mixing algorithm than in LSTMs



$$z_t = \sigma_g(W_z x_t + U_z h_{t-1} + b_z)$$

$$r_t = \sigma_g(W_r x_t + U_r h_{t-1} + b_r)$$

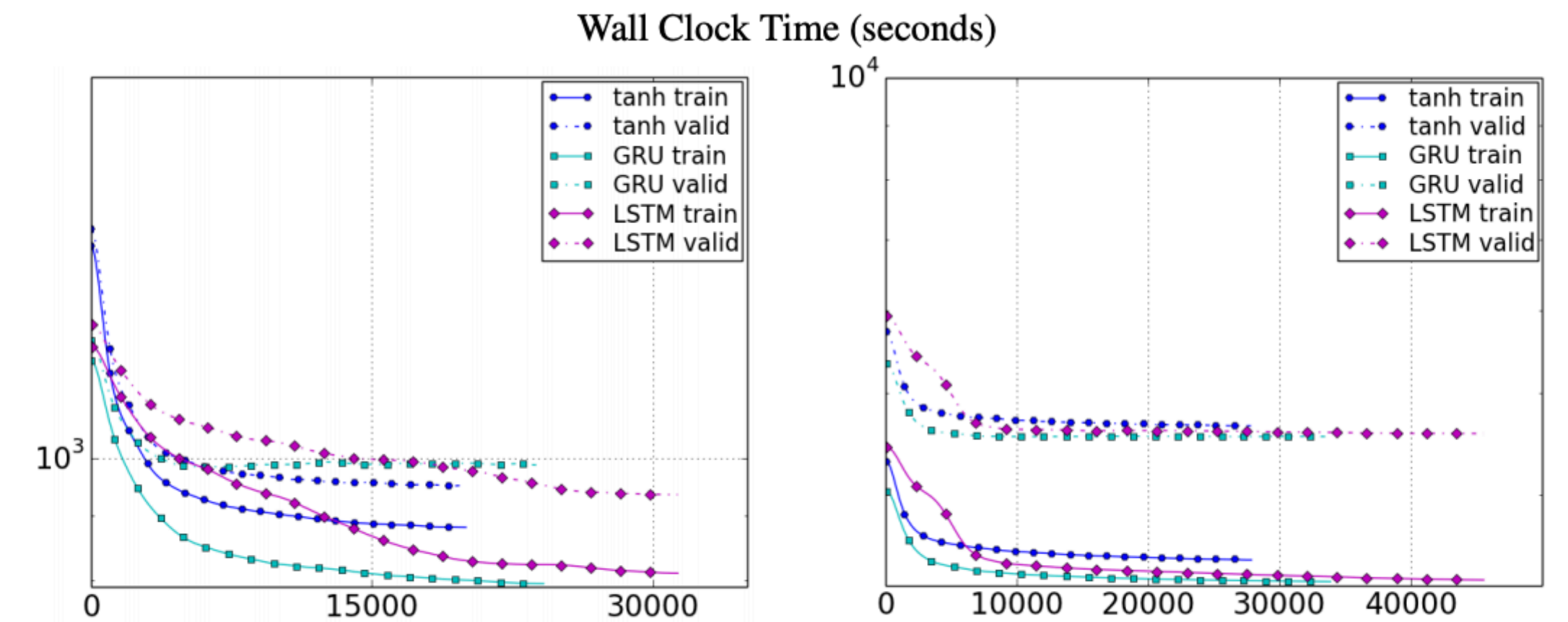
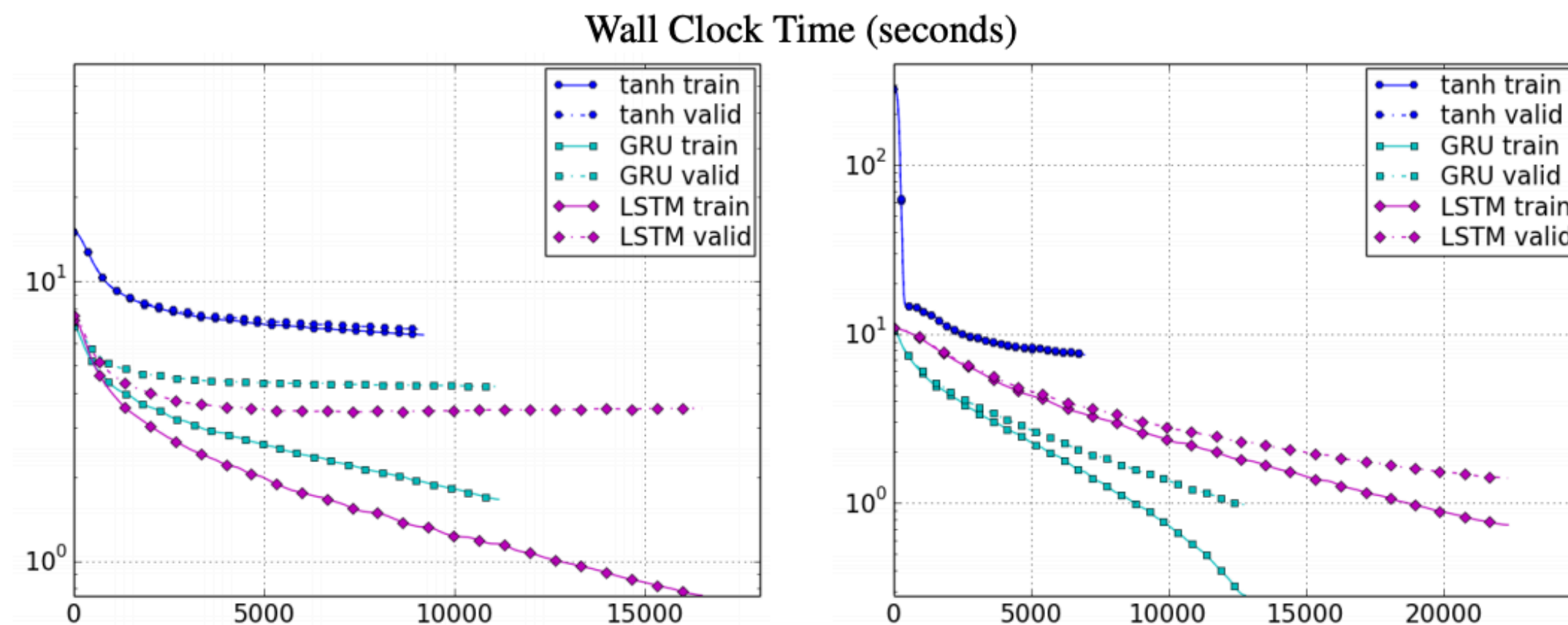
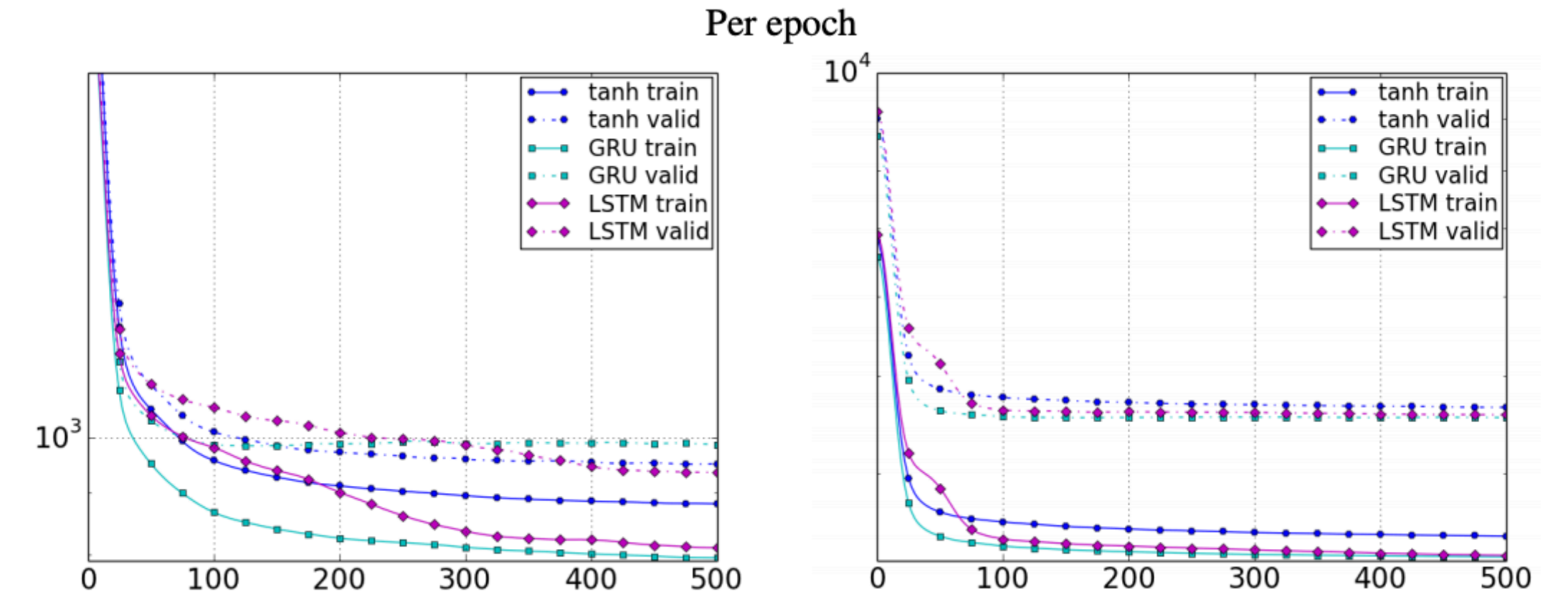
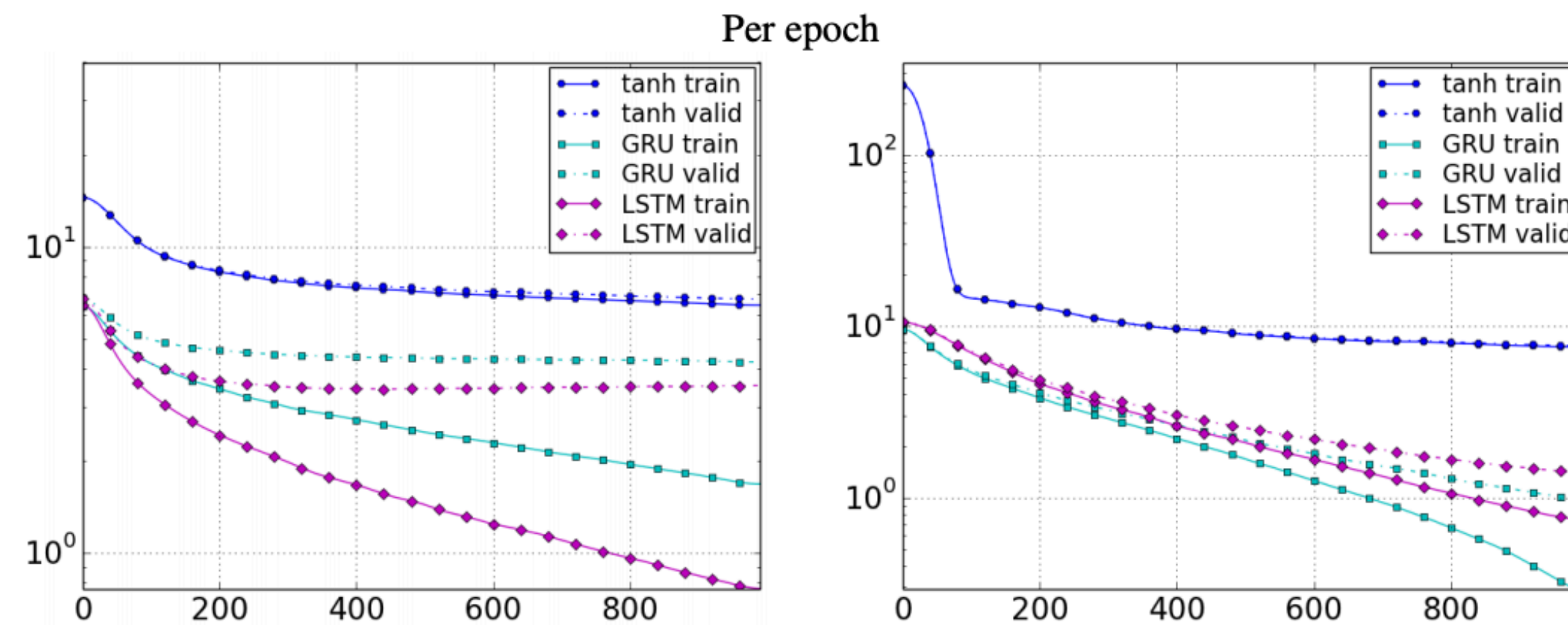
$$h_t = (1 - z_t) \circ h_{t-1} + z_t \circ \sigma_h(W_h x_t + U_h(r_t \circ h_{t-1}) + b_h)$$

Which is better?

Speech Signal Modeling

Music Modeling

Negative Loglikelihood



(a) UbiSoft Dataset A

(b) UbiSoft Dataset B

(a) Nottingham Dataset

(b) MuseData Dataset

Question

What are the advantages of using LSTMs and GRUs?

Vanishing Gradients?

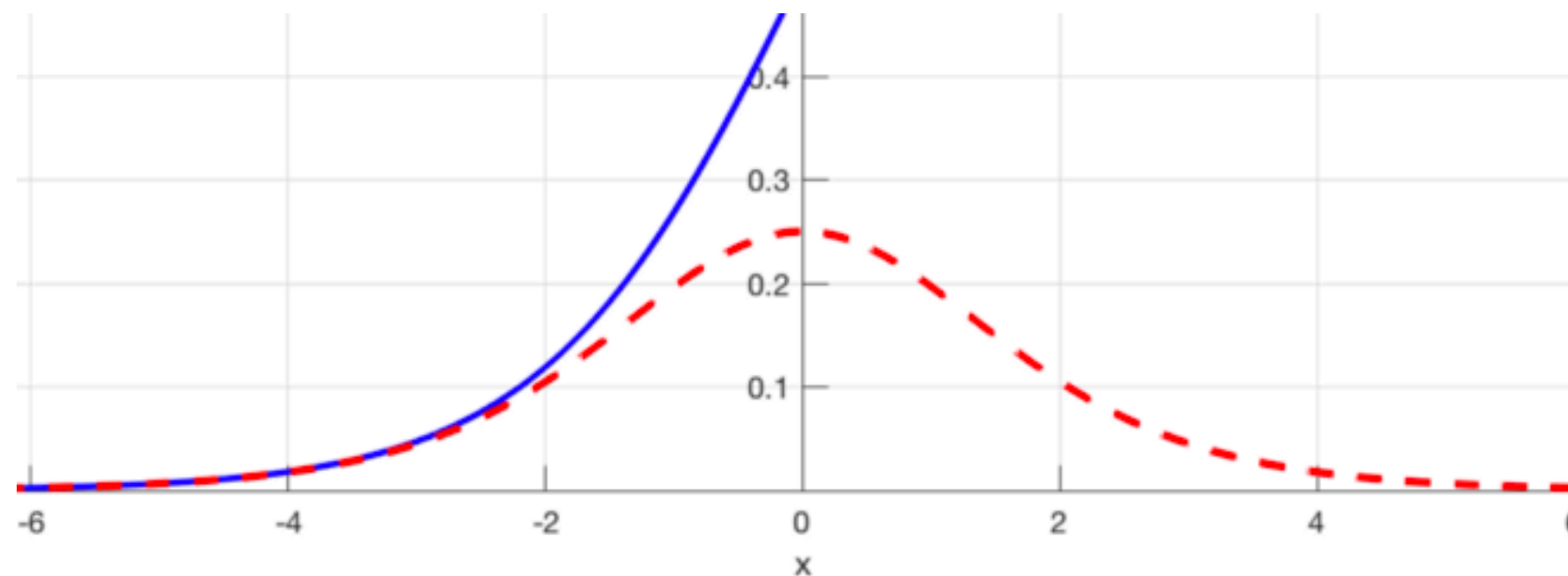
Recurrent Neural Networks

Long Short Term Memory

State maintained by hidden state feedback

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

Gradient systemically squashed by sigmoid



State maintained by cell value

$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

Gradient set by value of forget gate

$$\frac{\partial c_t}{\partial c_{t-1}} = f_t$$

Can still vanish, but only if forget gate closes!

Question

What's a disadvantages of using a LSTM or GRU?

Question

What's a disadvantages of using a LSTM or GRU?

More parameters!

$$f_t = \sigma(W_{fx}x_t + W_{fh}h_{t-1} + b_f)$$

$$i_t = \sigma(W_{ix}x_t + W_{ih}h_{t-1} + b_i)$$

$$o_t = \sigma(W_{ox}x_t + W_{oh}h_{t-1} + b_o)$$

$$\tilde{c}_t = \phi(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$$

$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

$$h_t = o_t \times \phi(c_t)$$

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

Question

Could there be better architectures than GRUs and LSTMs?

Optimal Architectures?

MUT1:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT2:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z) \\ r &= \text{sigm}(x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT3:

$$\begin{aligned} z &= \text{sigm}(W_{xz}x_t + W_{hz} \tanh(h_t) + b_z) \\ r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

Arch.	Arith.	XML	PTB
Tanh	0.29493	0.32050	0.08782
LSTM	0.89228	0.42470	0.08912
LSTM-f	0.29292	0.23356	0.08808
LSTM-i	0.75109	0.41371	0.08662
LSTM-o	0.86747	0.42117	0.08933
LSTM-b	0.90163	0.44434	0.08952
GRU	0.89565	0.45963	0.09069
MUT1	0.92135	0.47483	0.08968
MUT2	0.89735	0.47324	0.09036
MUT3	0.90728	0.46478	0.09161

Arch.	5M-tst	10M-v	20M-v	20M-tst
Tanh	4.811	4.729	4.635	4.582 (97.7)
LSTM	4.699	4.511	4.437	4.399 (81.4)
LSTM-f	4.785	4.752	4.658	4.606 (100.8)
LSTM-i	4.755	4.558	4.480	4.444 (85.1)
LSTM-o	4.708	4.496	4.447	4.411 (82.3)
LSTM-b	4.698	4.437	4.423	4.380 (79.83)
GRU	4.684	4.554	4.559	4.519 (91.7)
MUT1	4.699	4.605	4.594	4.550 (94.6)
MUT2	4.707	4.539	4.538	4.503 (90.2)
MUT3	4.692	4.523	4.530	4.494 (89.47)

Recap

- Recurrent neural networks can **theoretically** learn to model an **unbounded context length**
 - no increase in model size because weights are shared across time steps
- Practically, however, **vanishing gradients** stop vanilla RNNs from learning useful **long-range dependencies**
- LSTMs and GRUs are variants of recurrent networks that mitigate the vanishing gradient problem
 - used for for **many sequence-to-sequence tasks (up next!)**

References

- Elman, J.L. (1990). Finding Structure in Time. *Cogn. Sci.*, 14, 179-211.
- Schuster, M., & Paliwal, K. K. (1997). Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11), 2673–2681.
- Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9, 1735-1780.
- Cho, K., Merrienboer, B.V., Gülçehre, Ç., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. *Conference on Empirical Methods in Natural Language Processing*.
- Greff, K., Srivastava, R.K., Koutník, J., Steunebrink, B.R., & Schmidhuber, J. (2015). LSTM: A Search Space Odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 28, 2222-2232.