





RNN

Name	Nathan Varghese
Identity Key	nava 3000

	Level	Completed	Goal	
	Beginner	8	5722	10
	Intermediate	2		
	Advanced	0	Total Completed	
	Expert	0	10	

Recurrent Neural Network

CSCI 5722: Computer Vision

Fall 2024

Dr. Tom Yeh

Evolution to RNNs

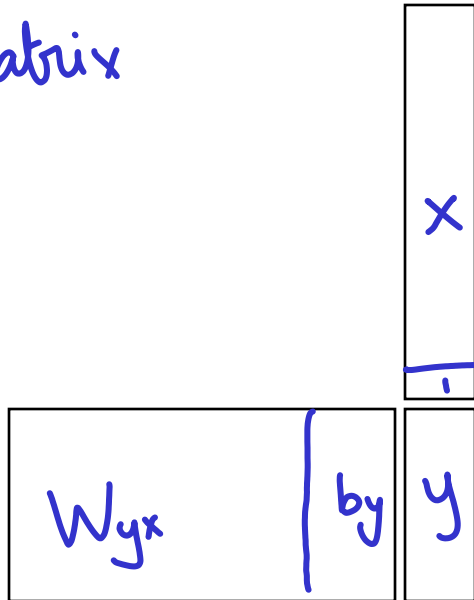
CSCI 5722 Computer Vision



University of Colorado
Boulder

Linear transformation

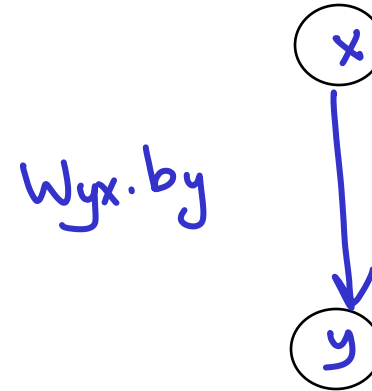
matrix



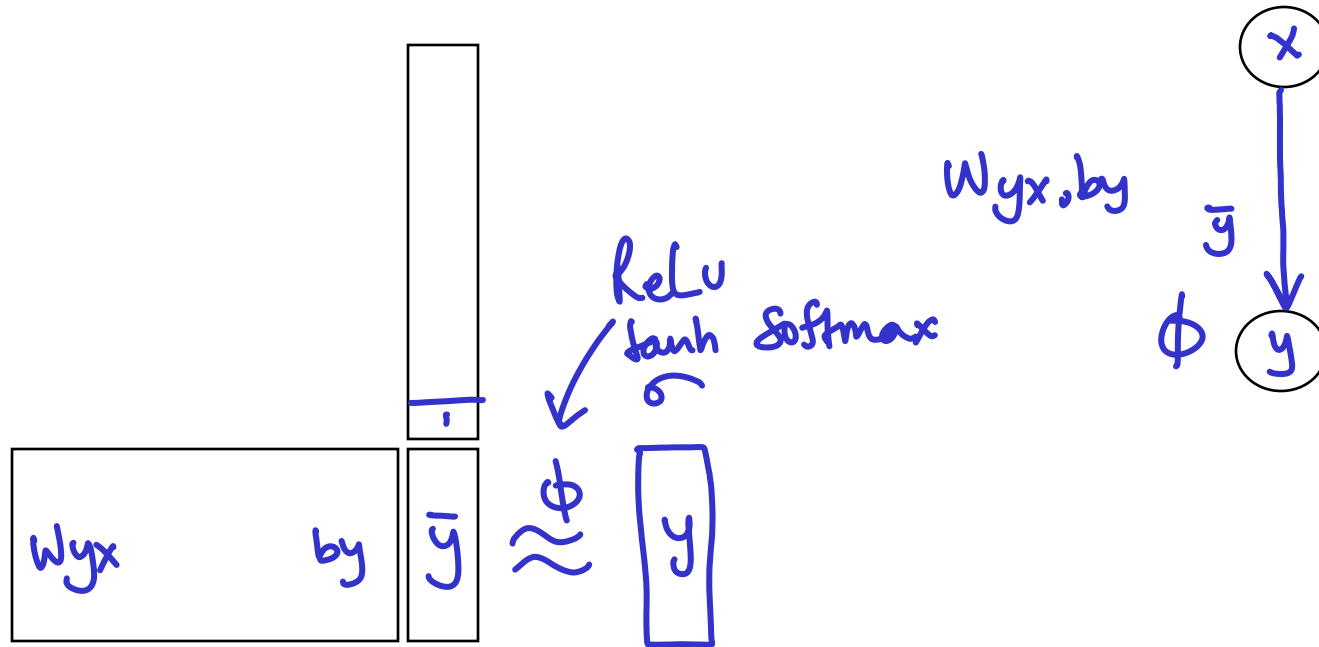
algebraic

$$y = W_{yx} * X + b_y$$

graphical

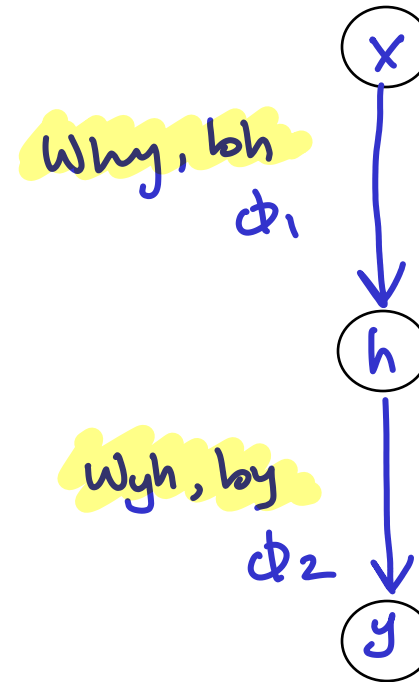
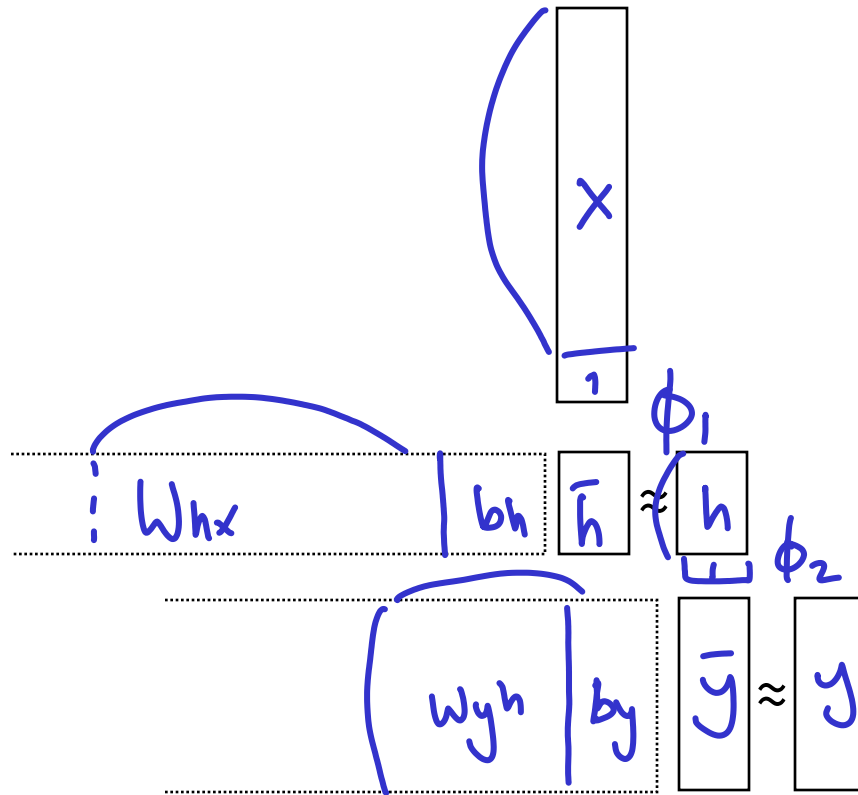


+ non-linearity



$$\underline{y} = \underline{\phi(w_{yx} \cdot X + b_y)}$$

+ hidden layer



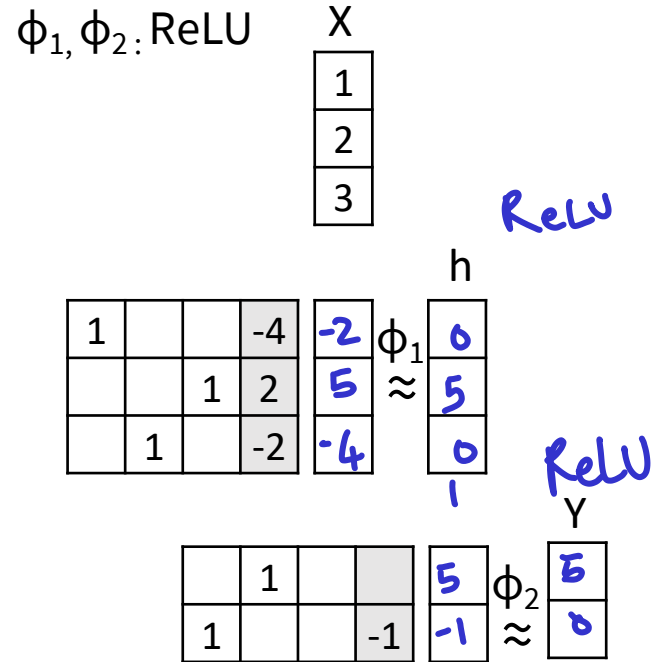
● - trainable parameters

$$h = \phi_1 (W_{hx} \cdot X + b_h)$$

$$y = \phi_2 (W_{yh} \cdot h + b_y)$$

$$= \phi_2 (W_{yh} \cdot (\phi_1 (W_{hx} \cdot X + b_h)) + b_y)$$

Calculate an NN with a hidden layer



$$Y = \phi_2(W_2 \cdot \phi_1(W_1 \cdot X + b_1) + b_2)$$

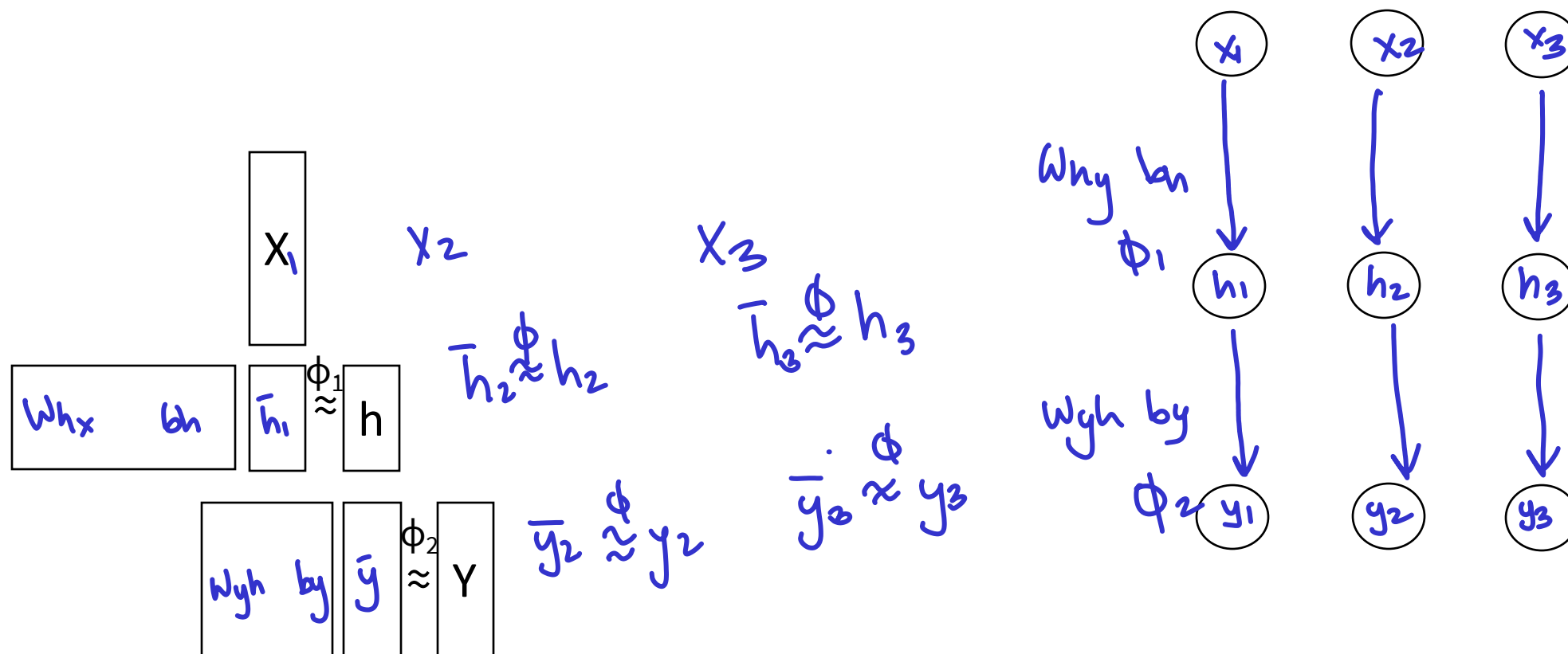
$$\text{size}(W_1) = \underline{3 \times 3}$$

$$\text{size}(b_1) = \underline{3 \times 1}$$

$$\text{size}(W_2) = \underline{2 \times 3}$$

$$\text{size}(b_2) = \underline{2 \times 1}$$

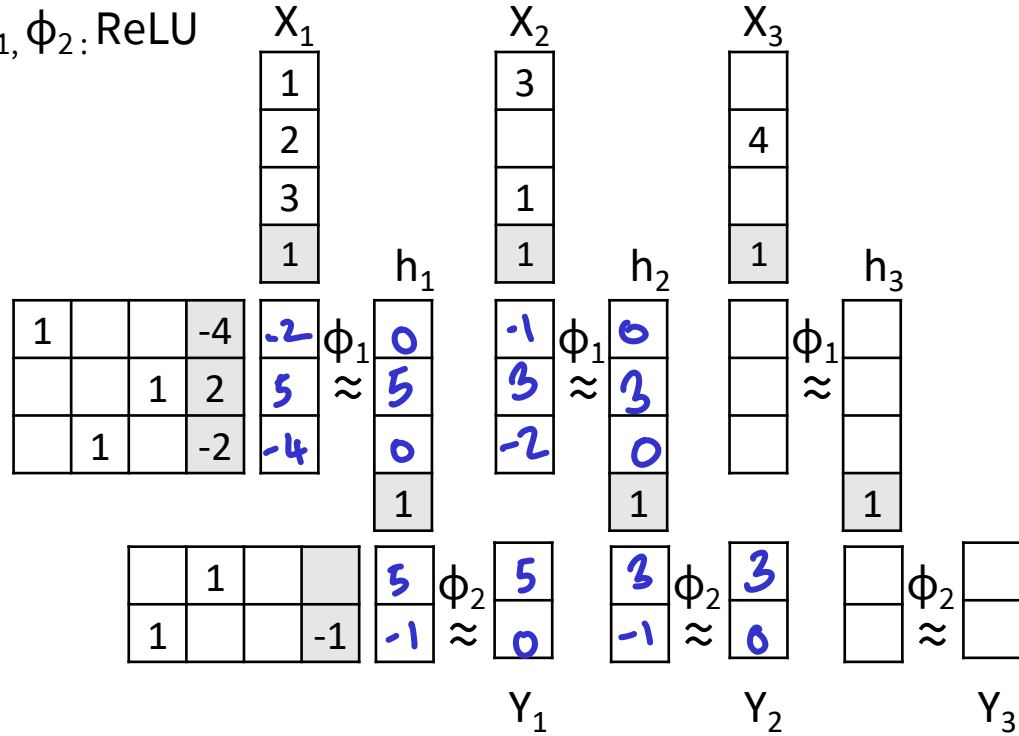
+ sequence



$$y_t = \phi_2 (w_{yh} \phi_1 (w_{hx} x_t + b_h) + b_y) \quad t=1, 2, 3 \dots$$

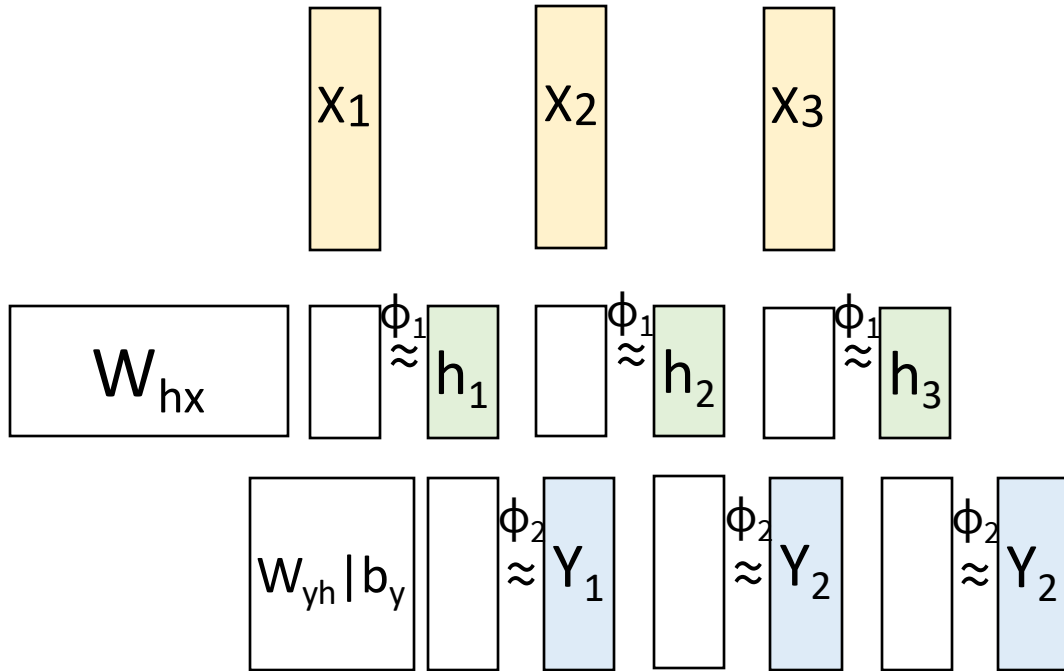
Calculate an NN for a batch of inputs

ϕ_1, ϕ_2 : ReLU

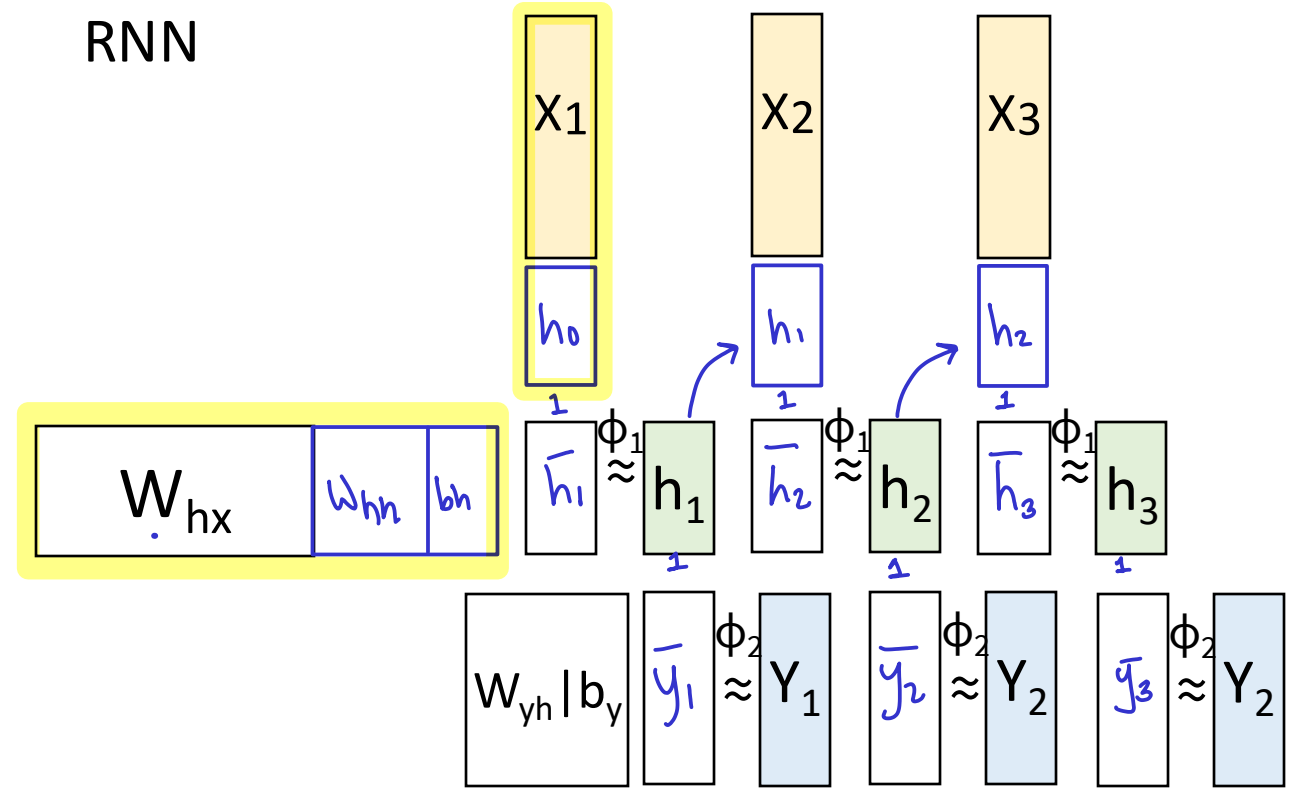


+ dependence on previous inputs

Independent NN



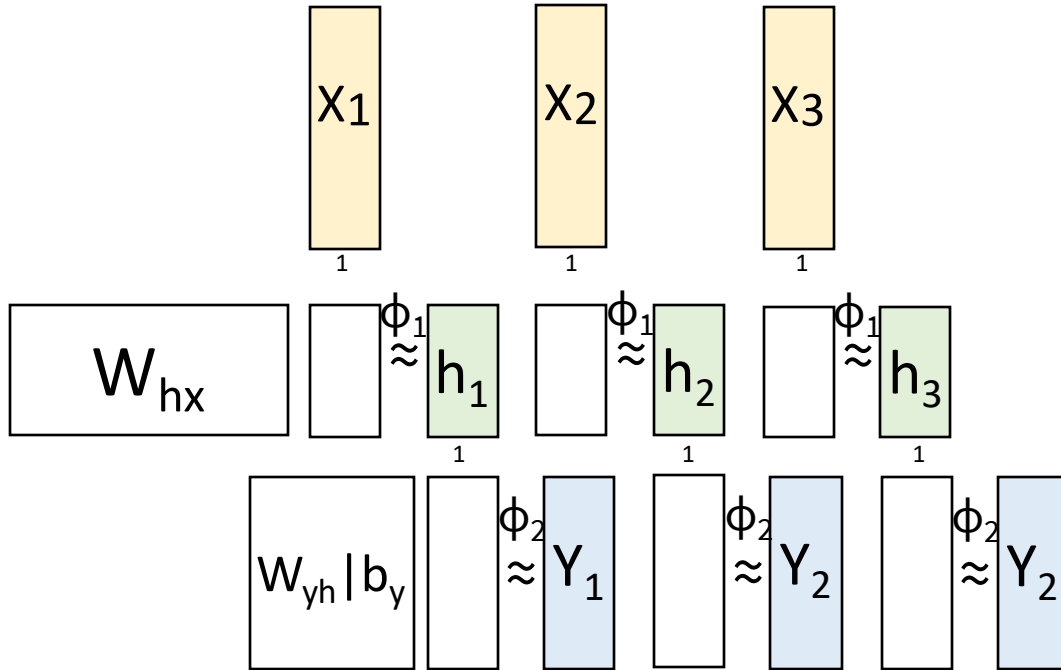
RNN



 = Multiplied together

NN vs RNN: Math

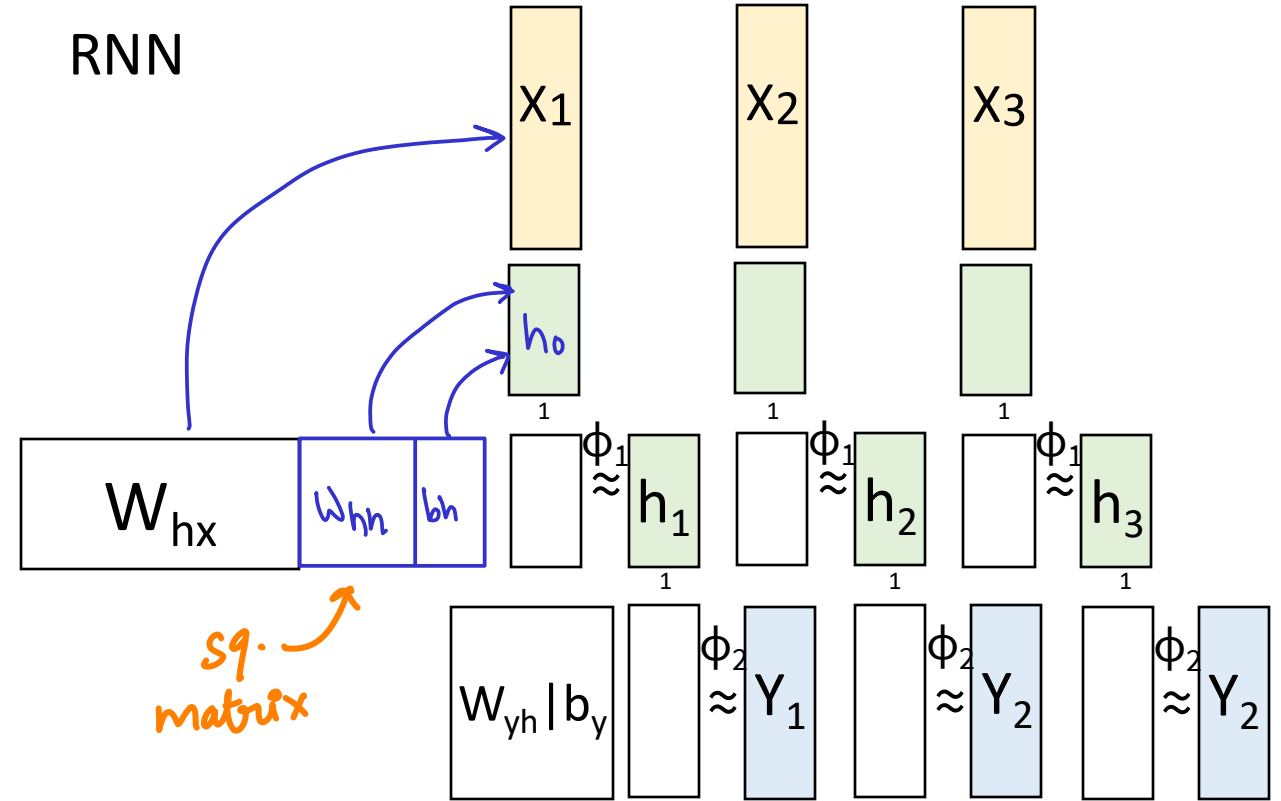
Independent NN



$$h_t = \phi_1(W_{hx} \cdot X_t + bh)$$

$$Y_t = \phi_2(W_{yh} \cdot h_t + by)$$

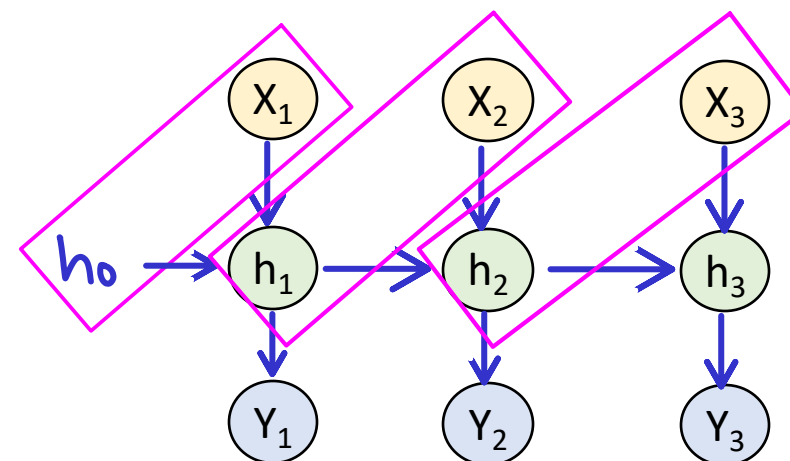
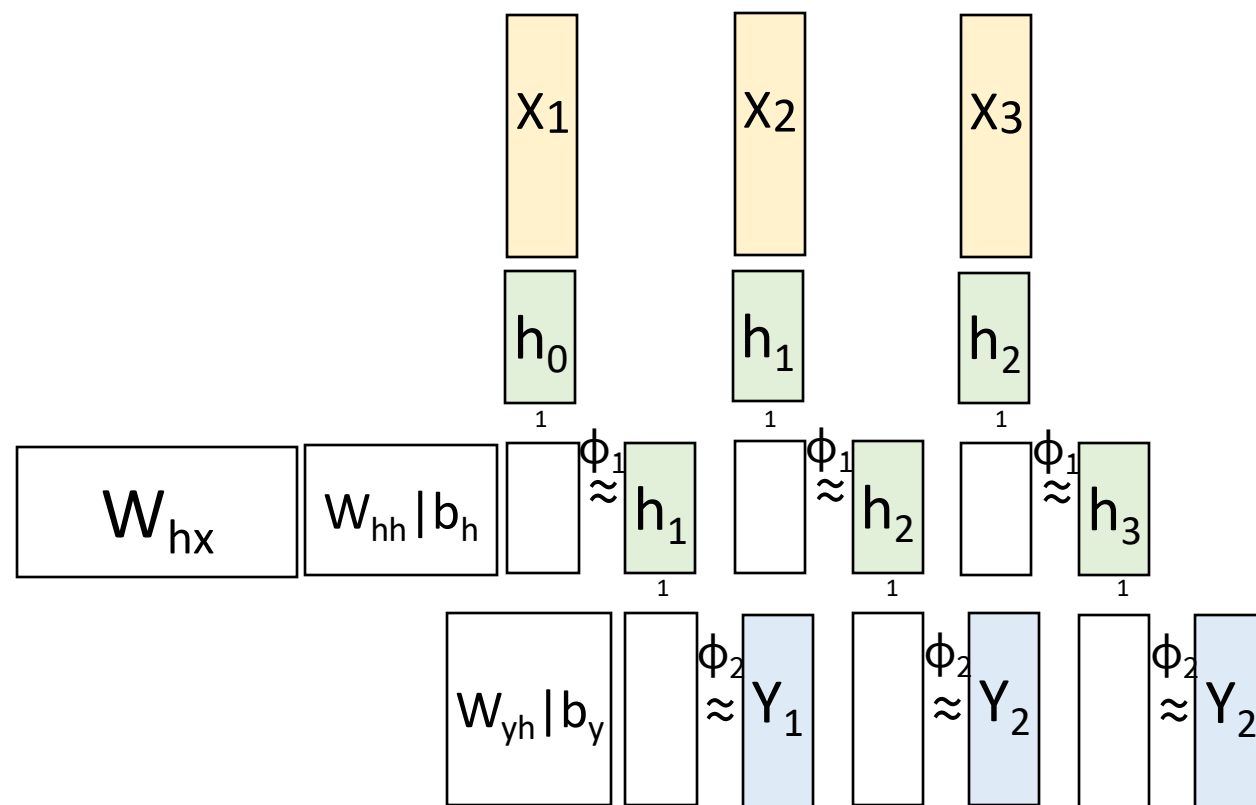
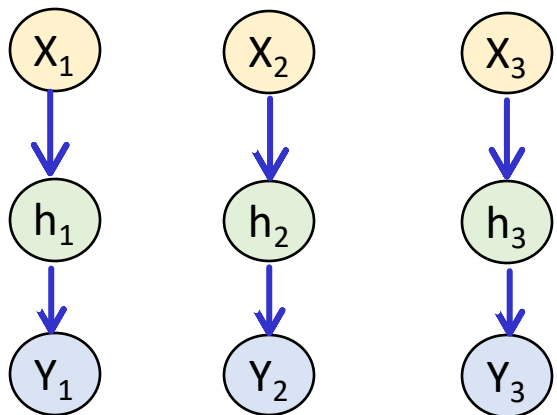
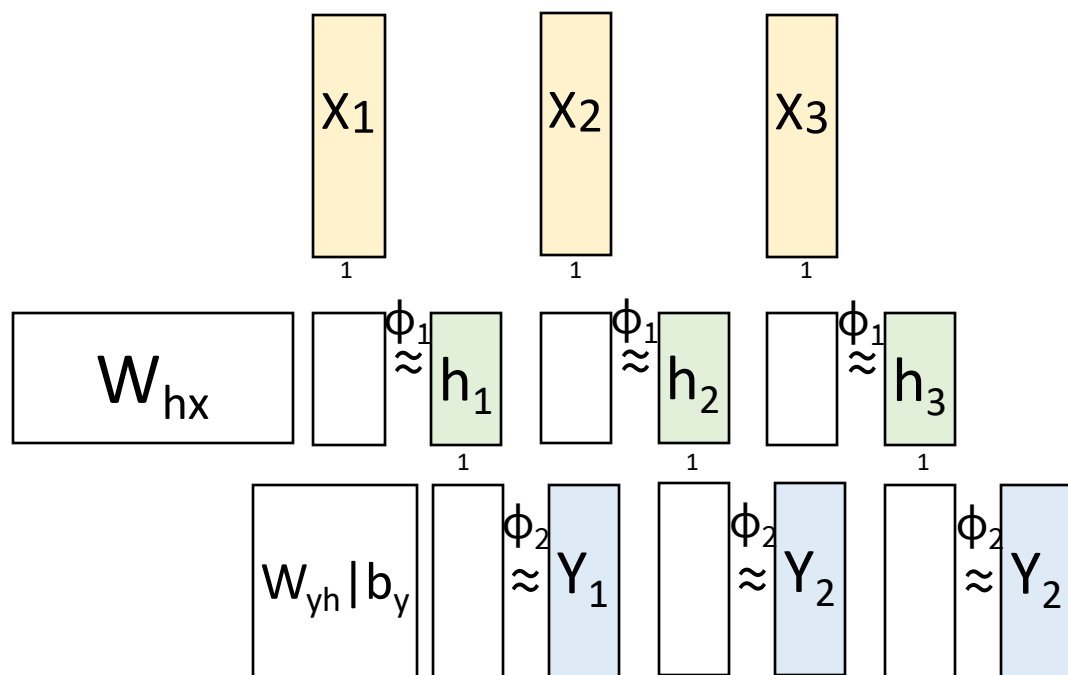
RNN



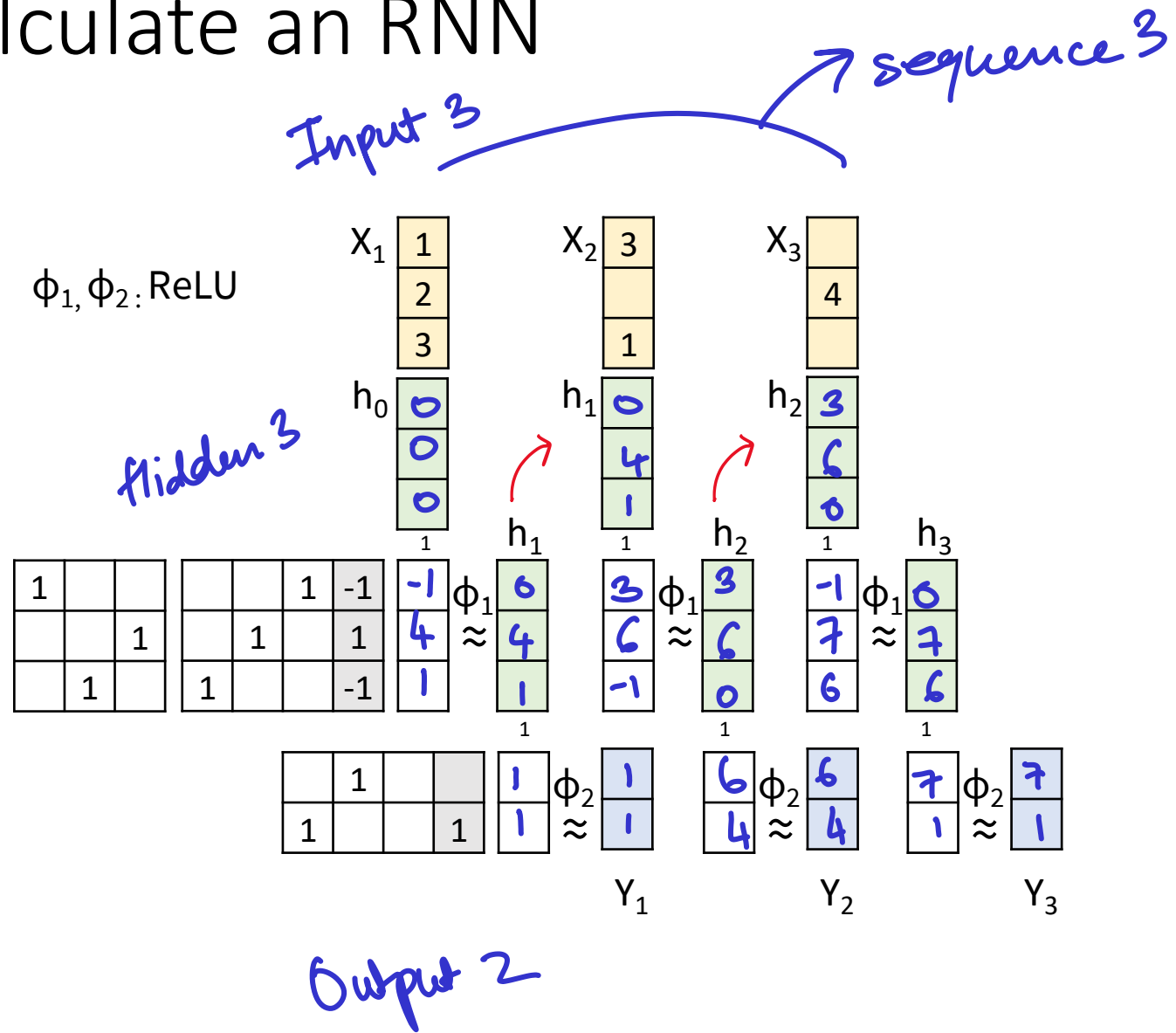
$$h_t = \phi_1(W_{hx} \cdot X_t + W_{hh} \cdot h_{t-1} + bh)$$

$$Y_t = \phi_2(W_{yh} \cdot h_t + by)$$

NN vs. RNN: Graph



Calculate an RNN



Counting Parameters

Given :

size(X) = 128

size(h) = 64

size(Y) = 96

Determine:

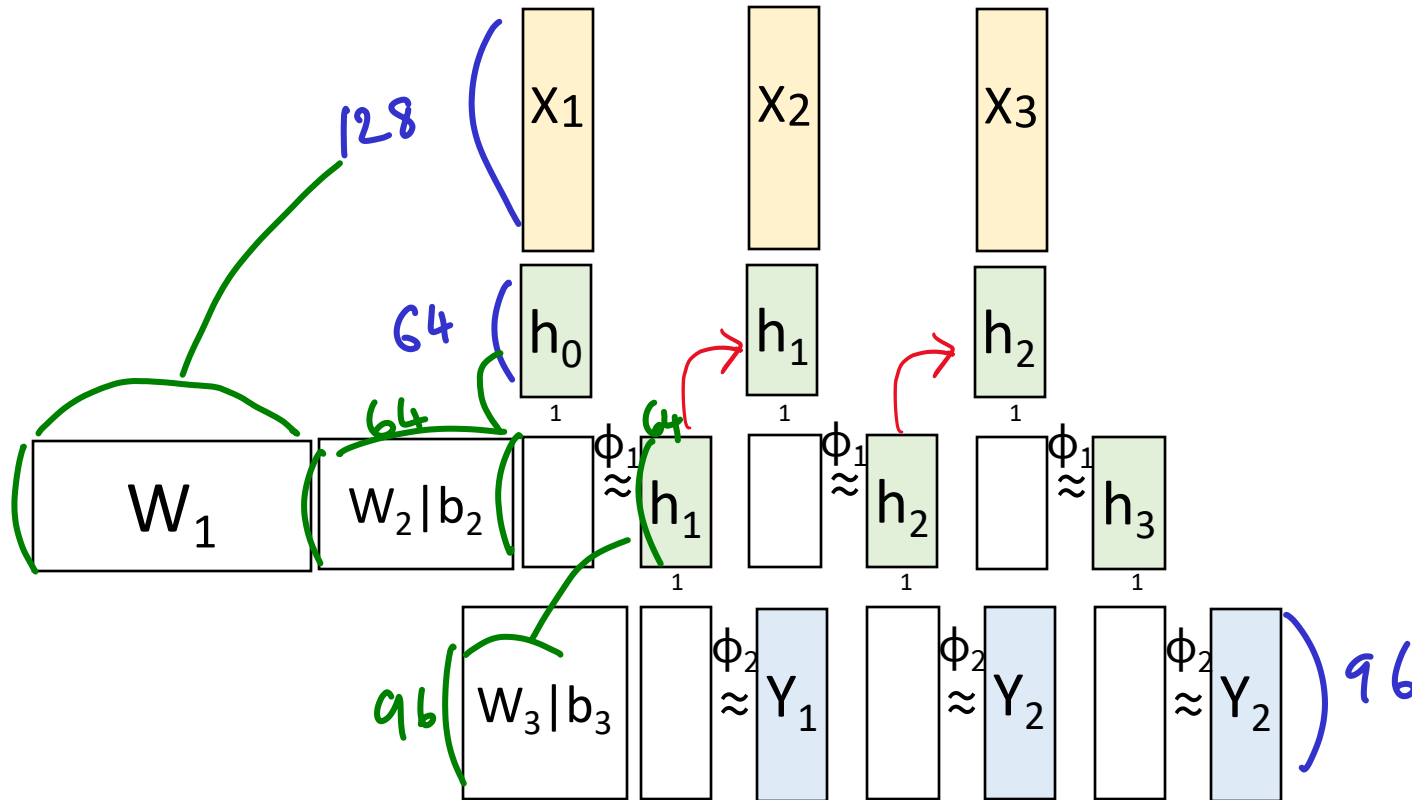
size(W₁) = 64 × 128

size(W₂) = 64 × 64

size(b₂) = 64 × 1

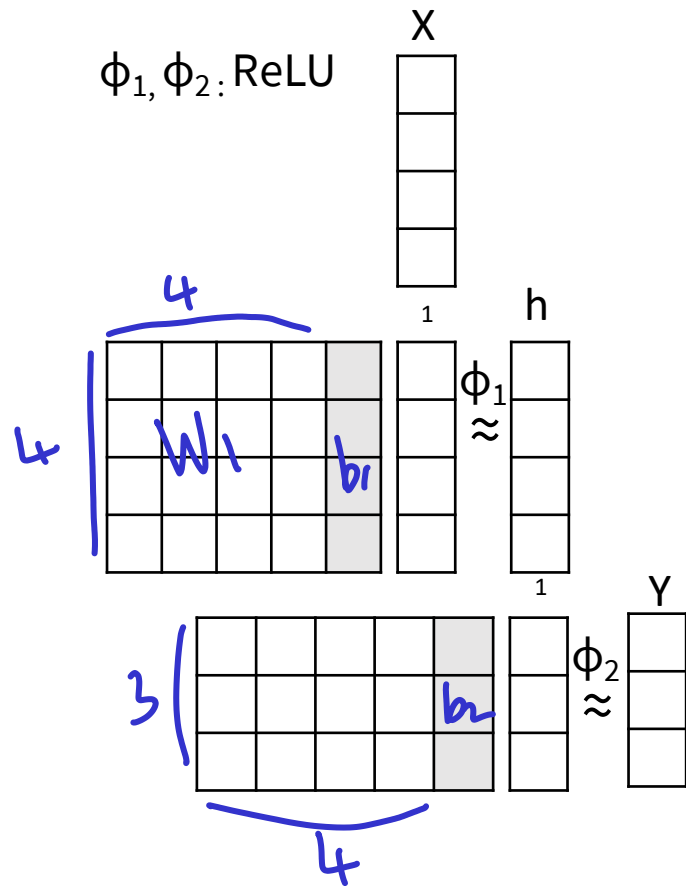
size(W₃) = 96 × 64

size(b₃) = 96 × 1





MLP Parameter Sizes



$$\text{size}(W_1) = \underline{4 \times 4}$$

$$\text{size}(b_1) = \underline{4 \times 1}$$

$$\text{size}(W_2) = \underline{3 \times 4}$$

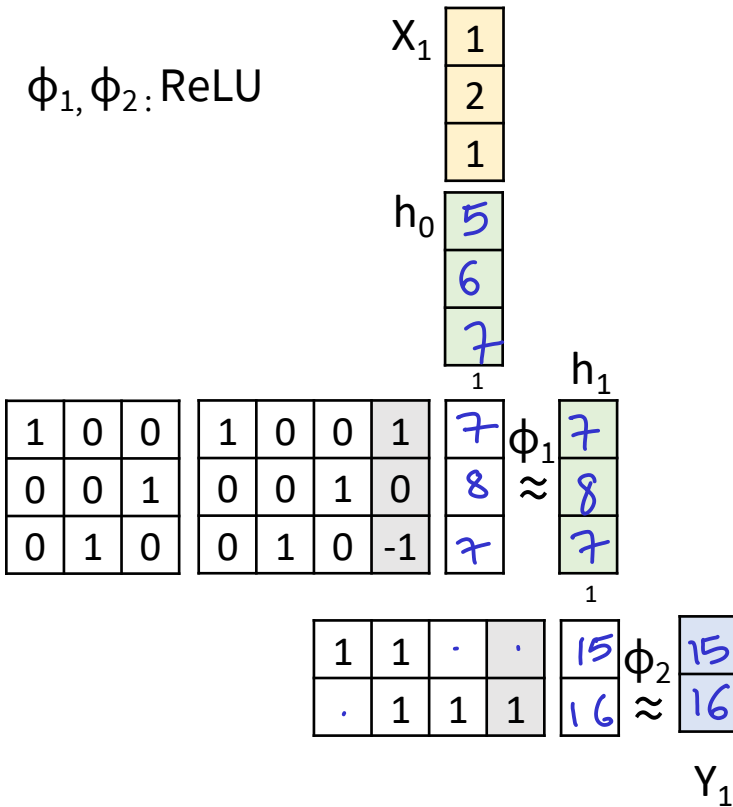
$$\text{size}(b_2) = \underline{3 \times 1}$$

$$Y = \phi_2(W_2 \cdot \phi_1(W_1 \cdot X + b_1) + b_2)$$



Calculate an RNN ($t = 1$)

This activity is standalone, not dependent on other activities.



$$1 + 5 + 1 = 7$$

$$1 + 7 = 8$$

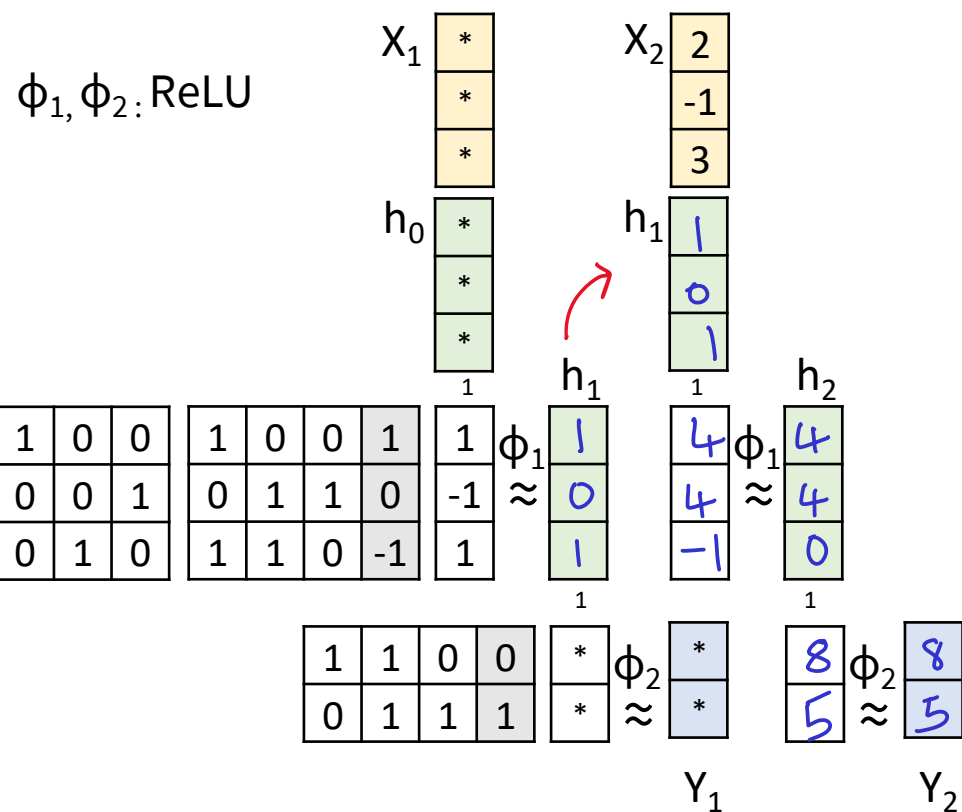
$$2 + 6 - 1 = 7$$

$$7 + 8 = 15$$

$$8 + 7 + 1 = 16$$



Calculate an RNN (t = 2)

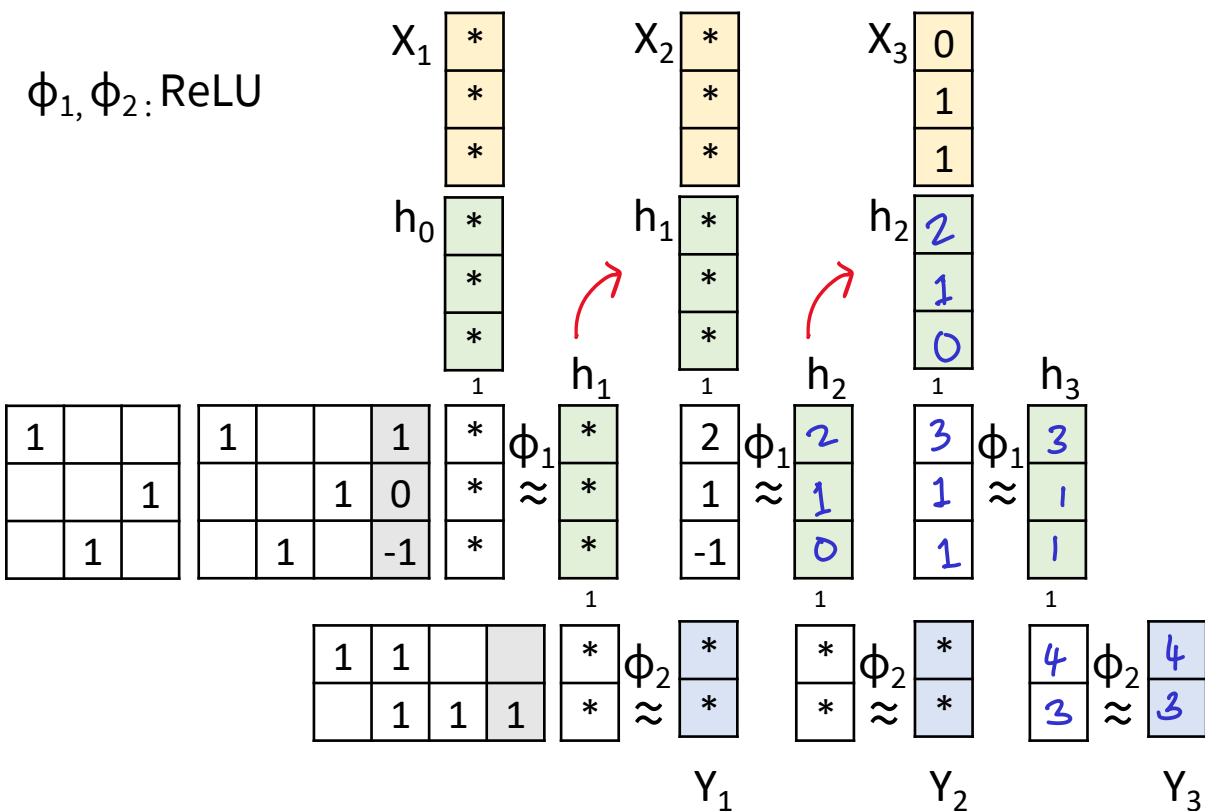


This activity is standalone, not dependent on other activities.



Calculate an RNN ($t = 3$)

This activity is standalone, not dependent on the previous one.



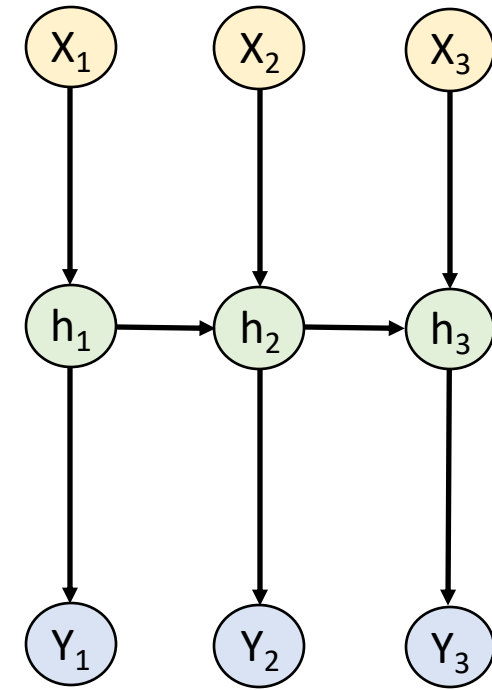
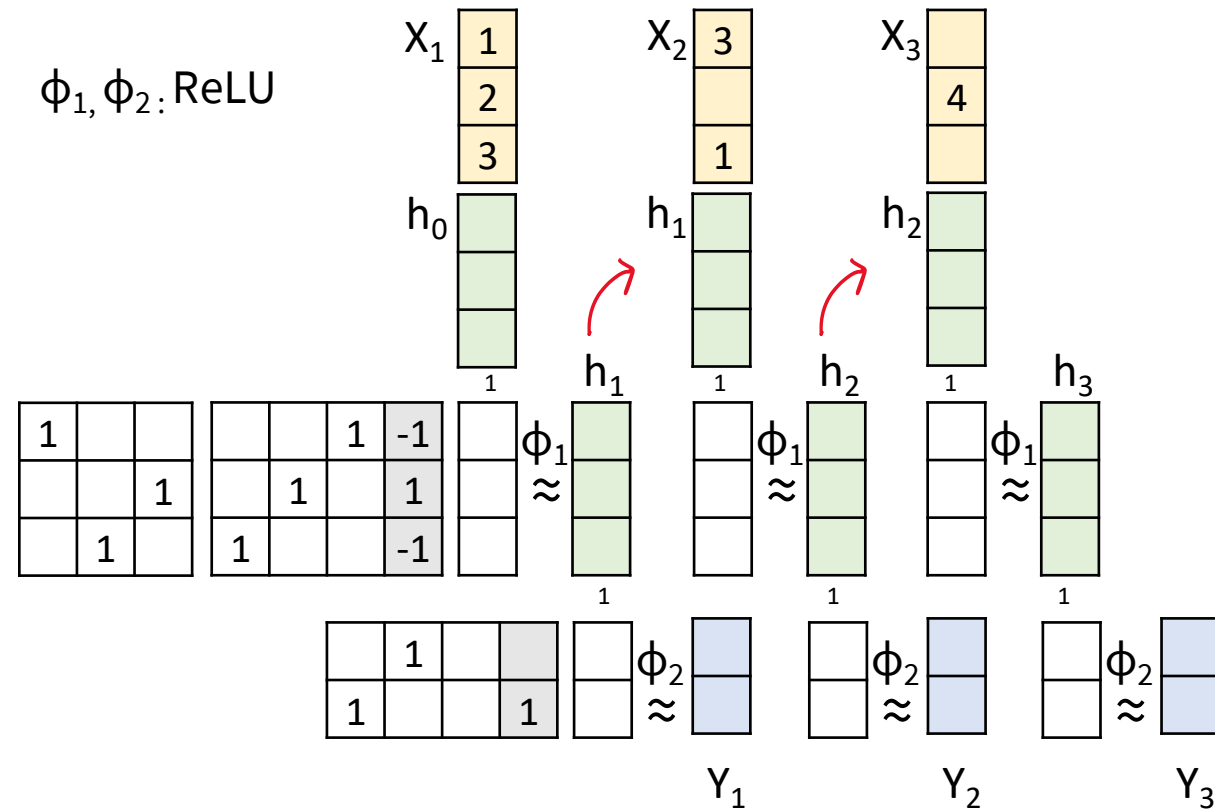
Variations of RNNs

CSCI 5722 Computer Vision

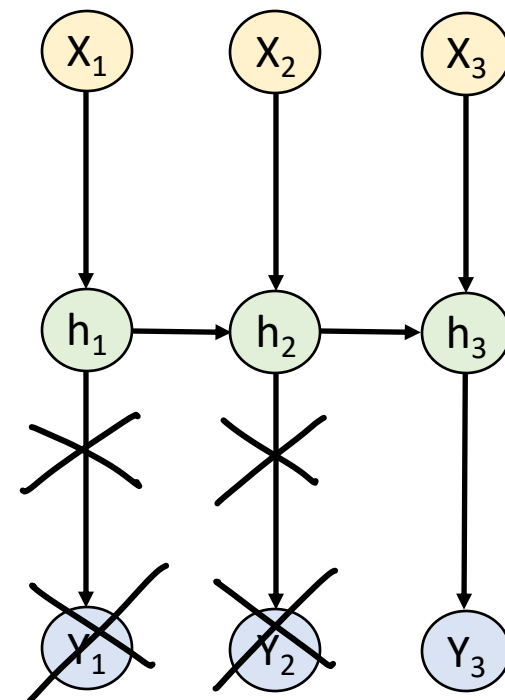
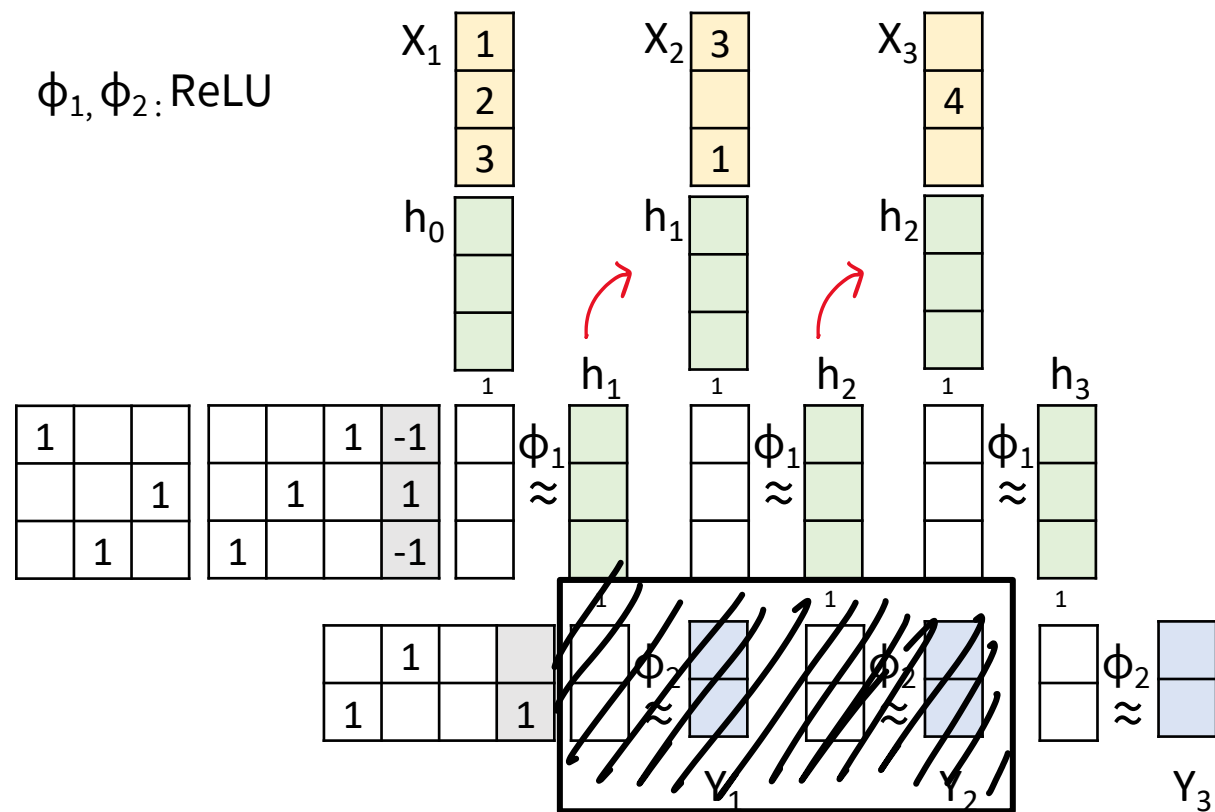


University of Colorado
Boulder

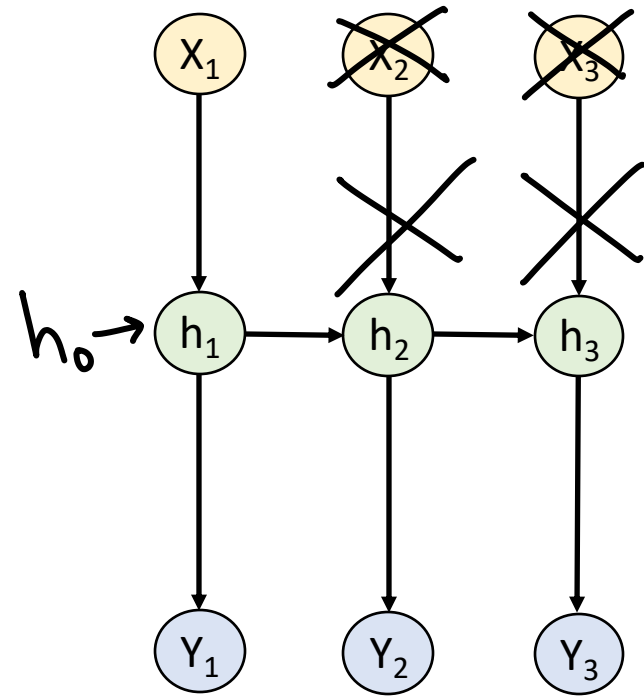
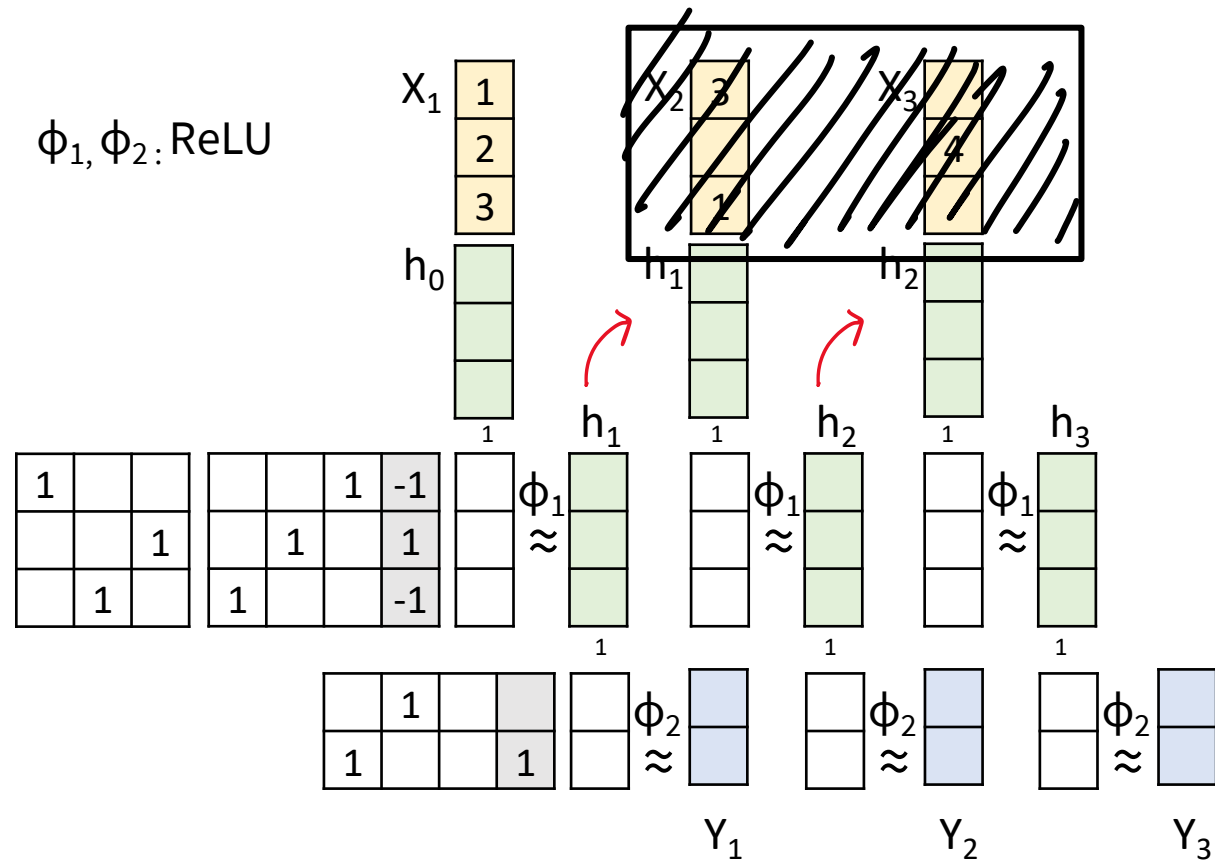
Many to Many



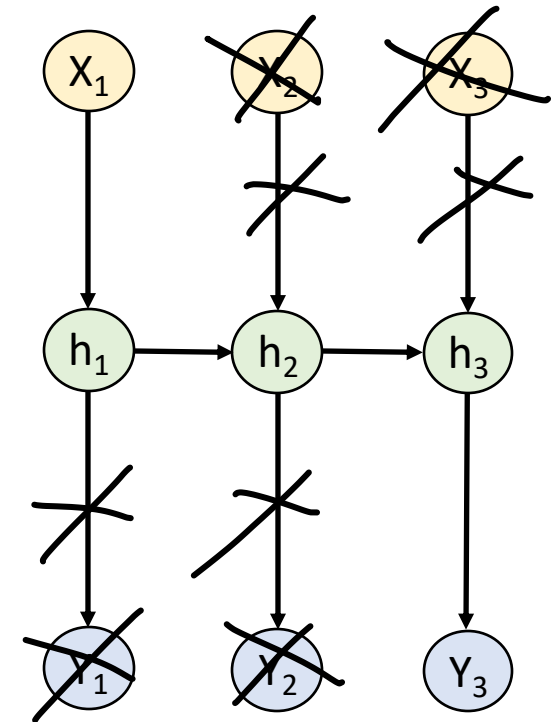
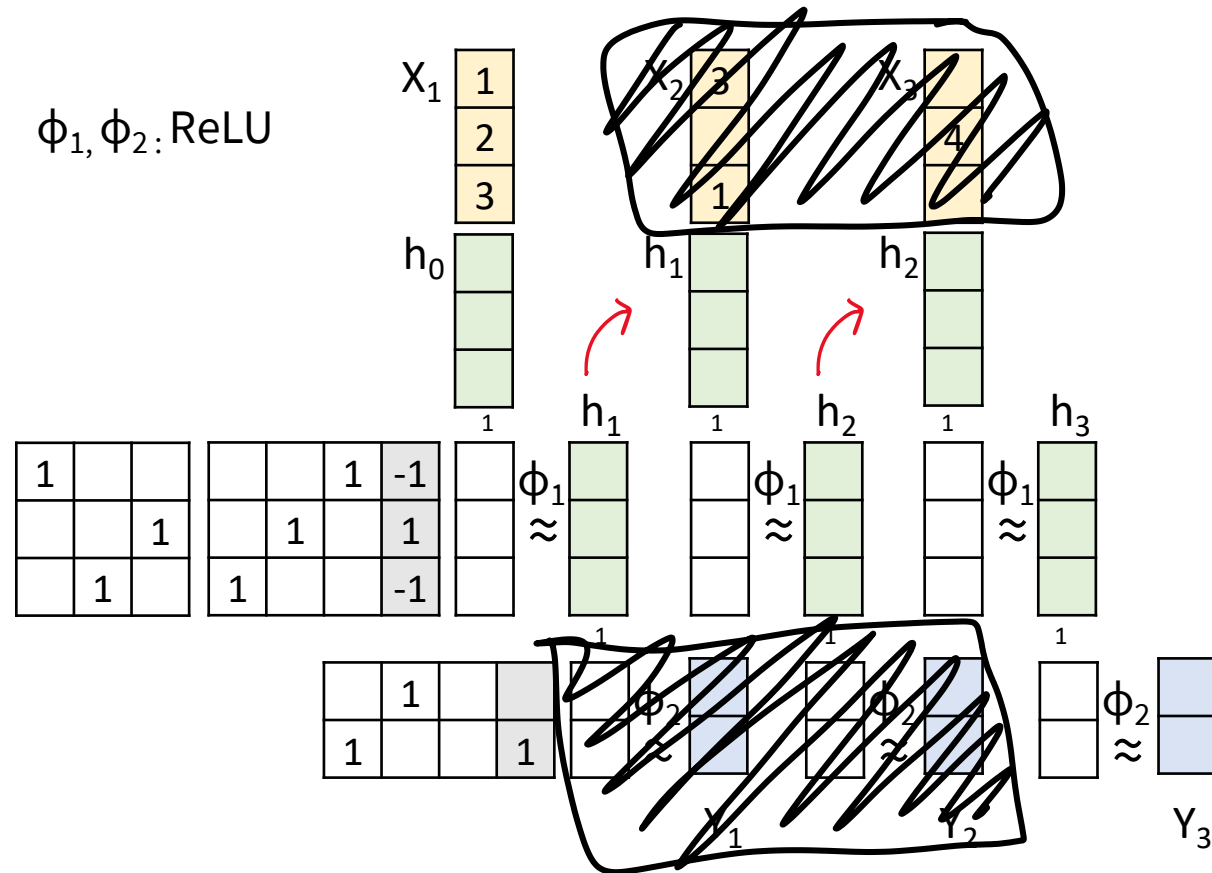
Many to One



One to Many



One to One



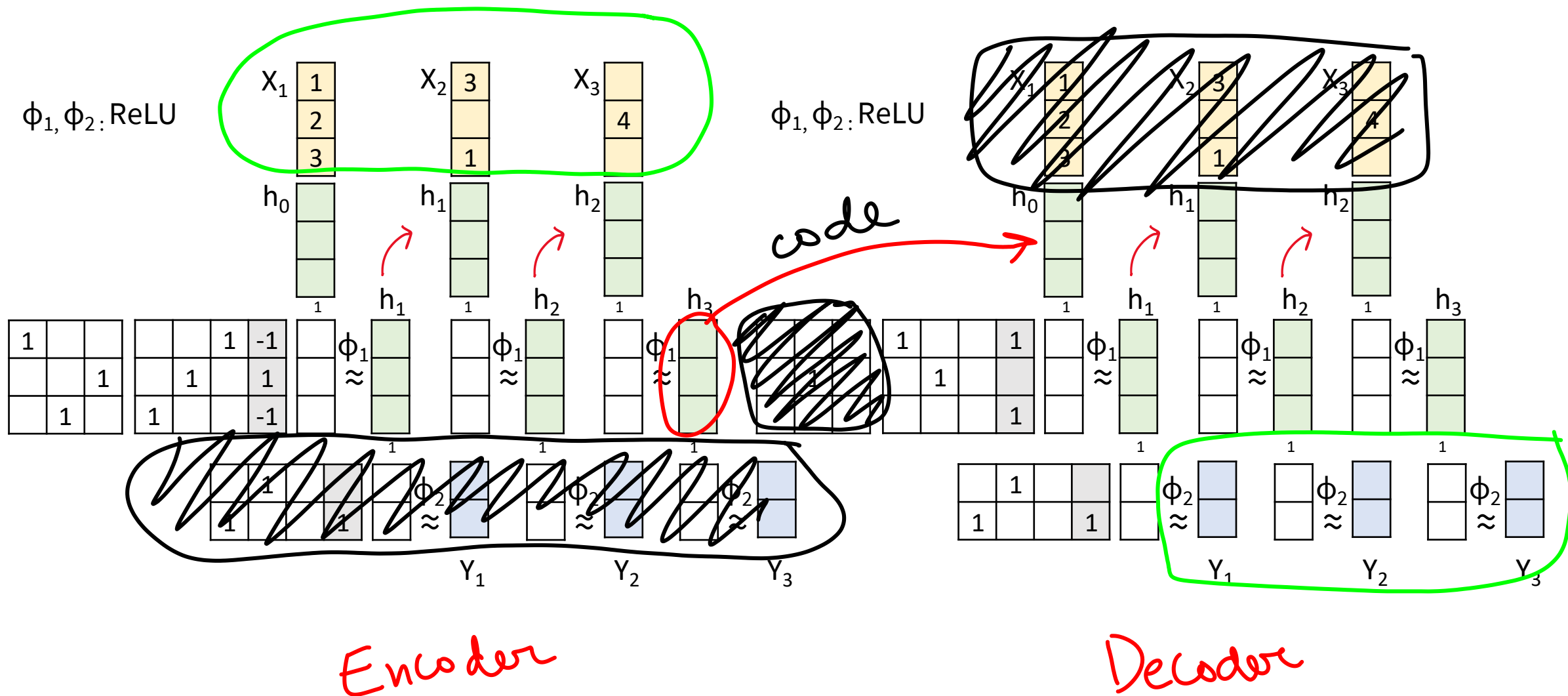
Advanced RNN Architectures

CSCI 5722 Computer Vision

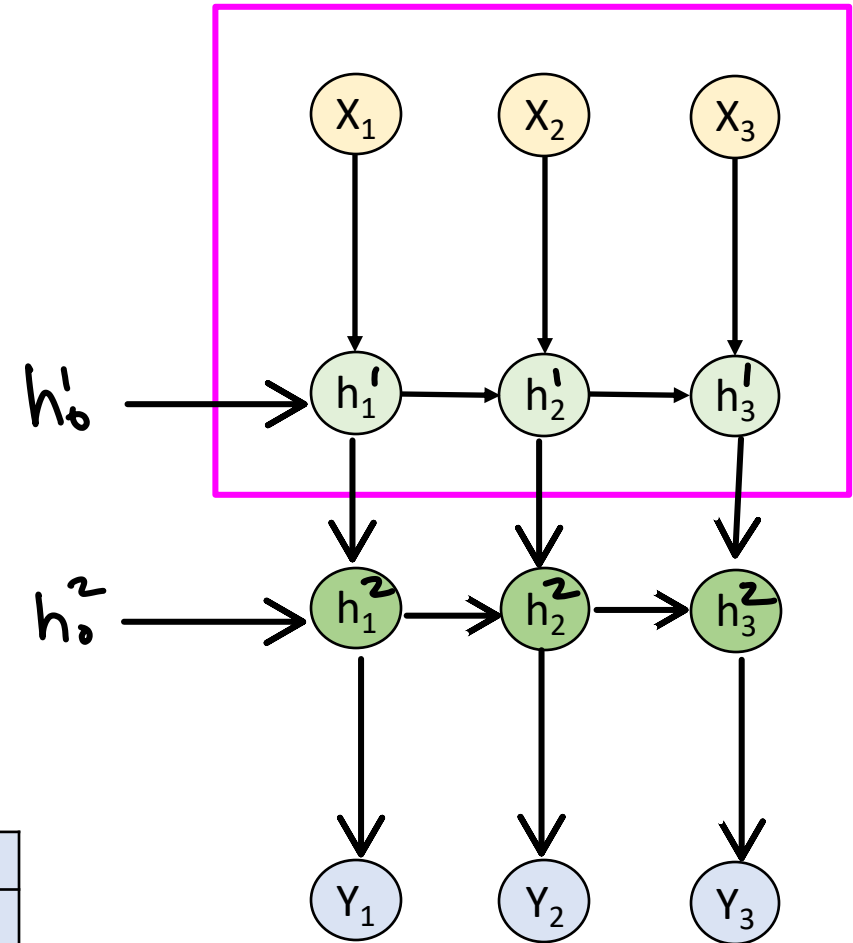
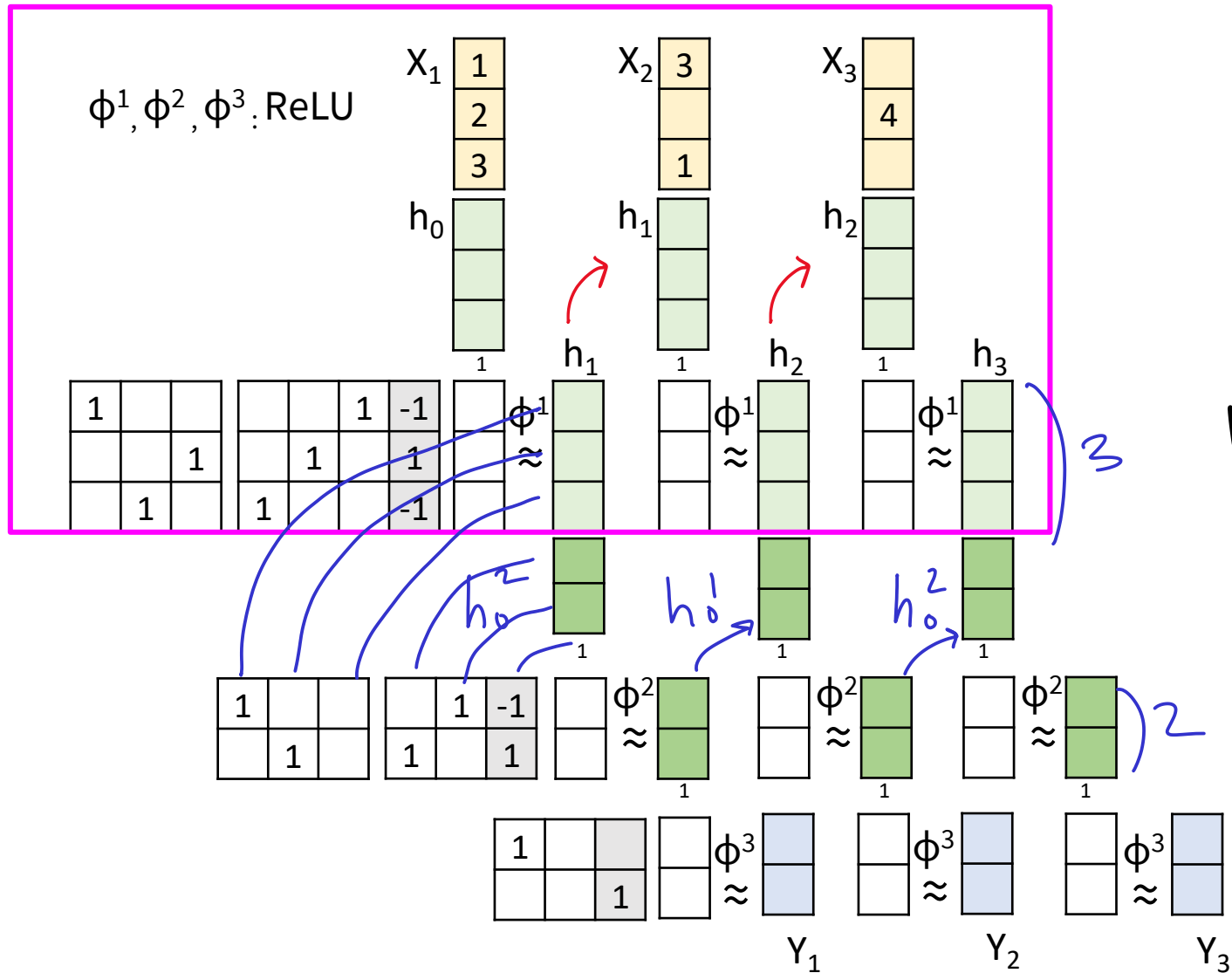


University of Colorado
Boulder

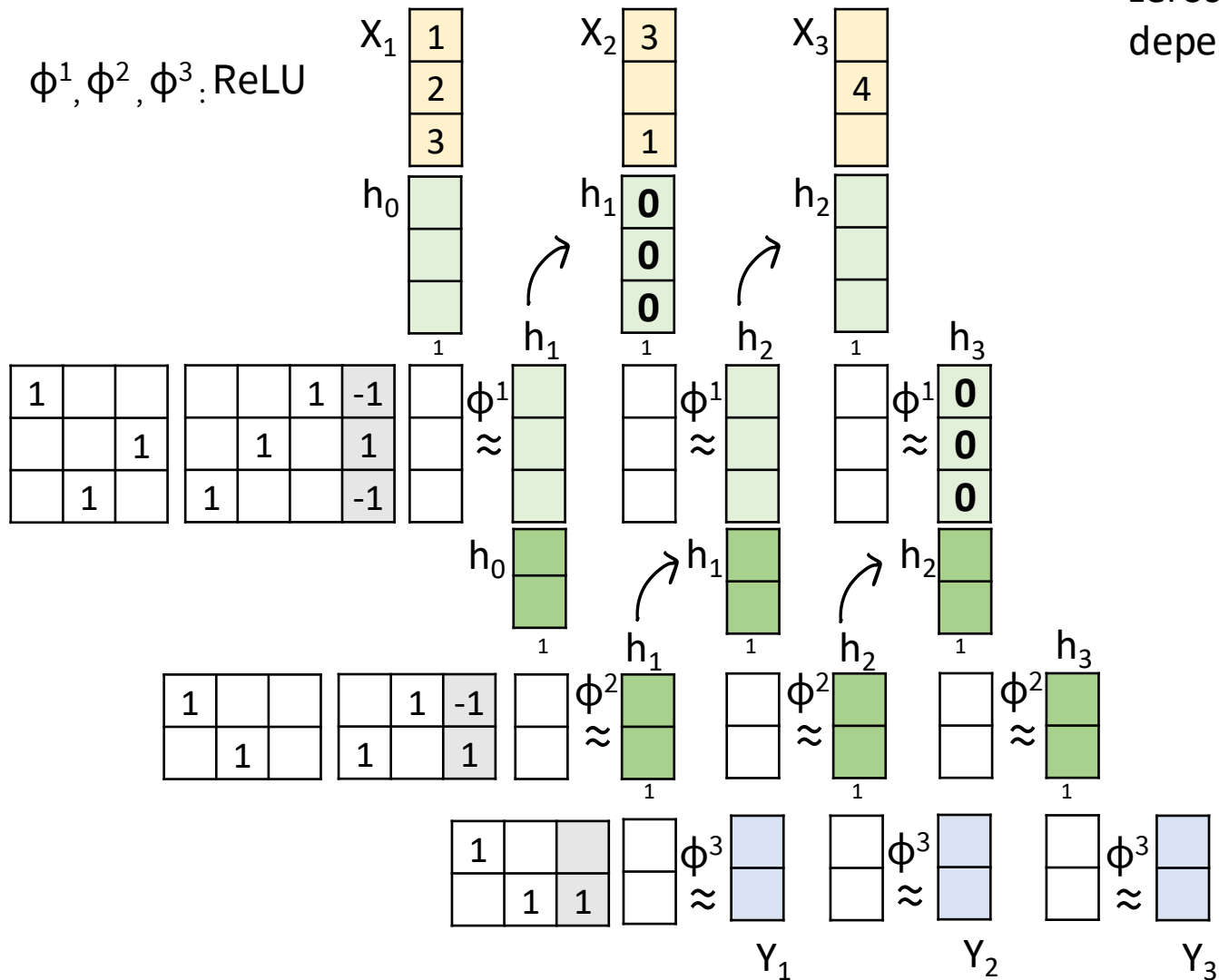
Sequence to Sequence



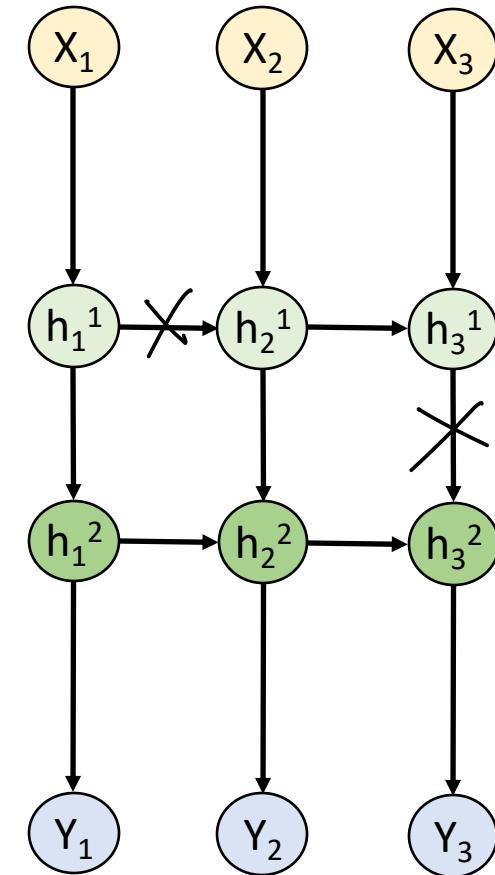
Multilayer RNN



Identify links

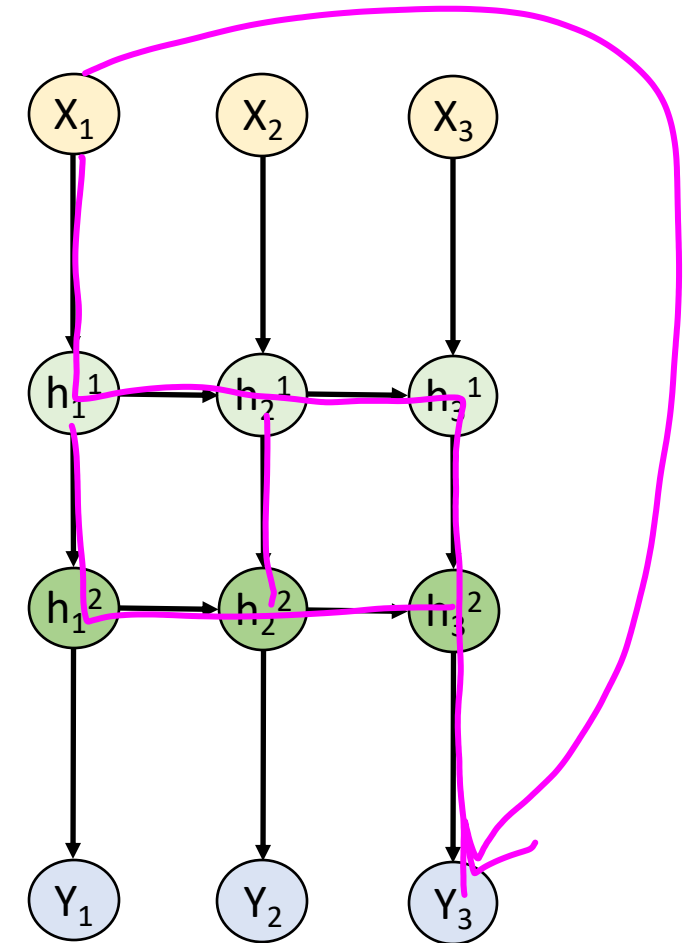
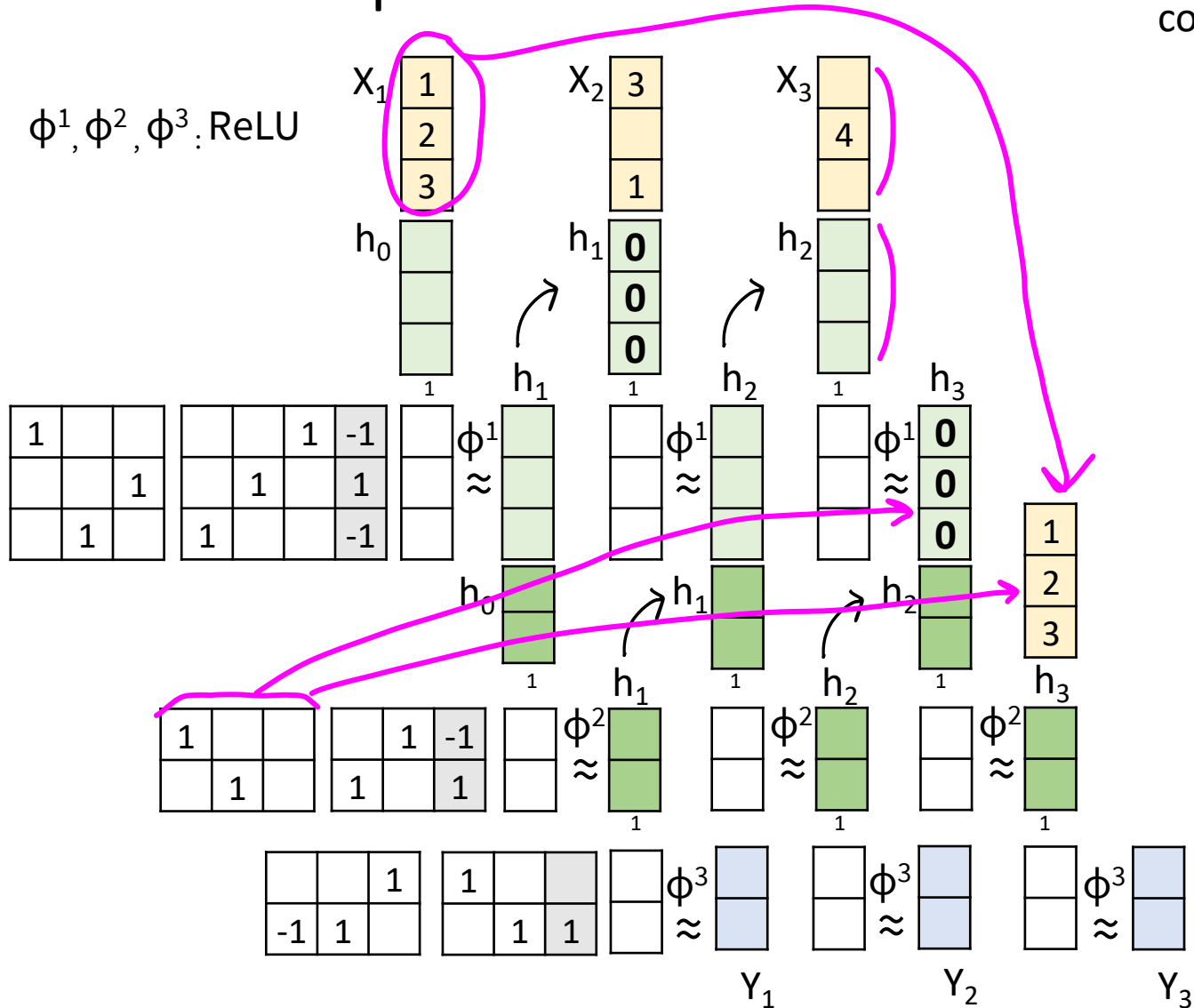


Suppose we set some hidden states to zeros on purpose. Cross out the affected dependency links.

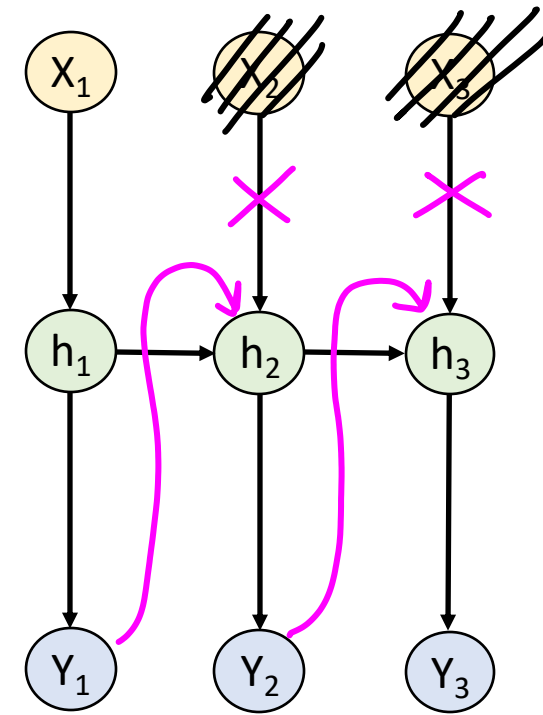
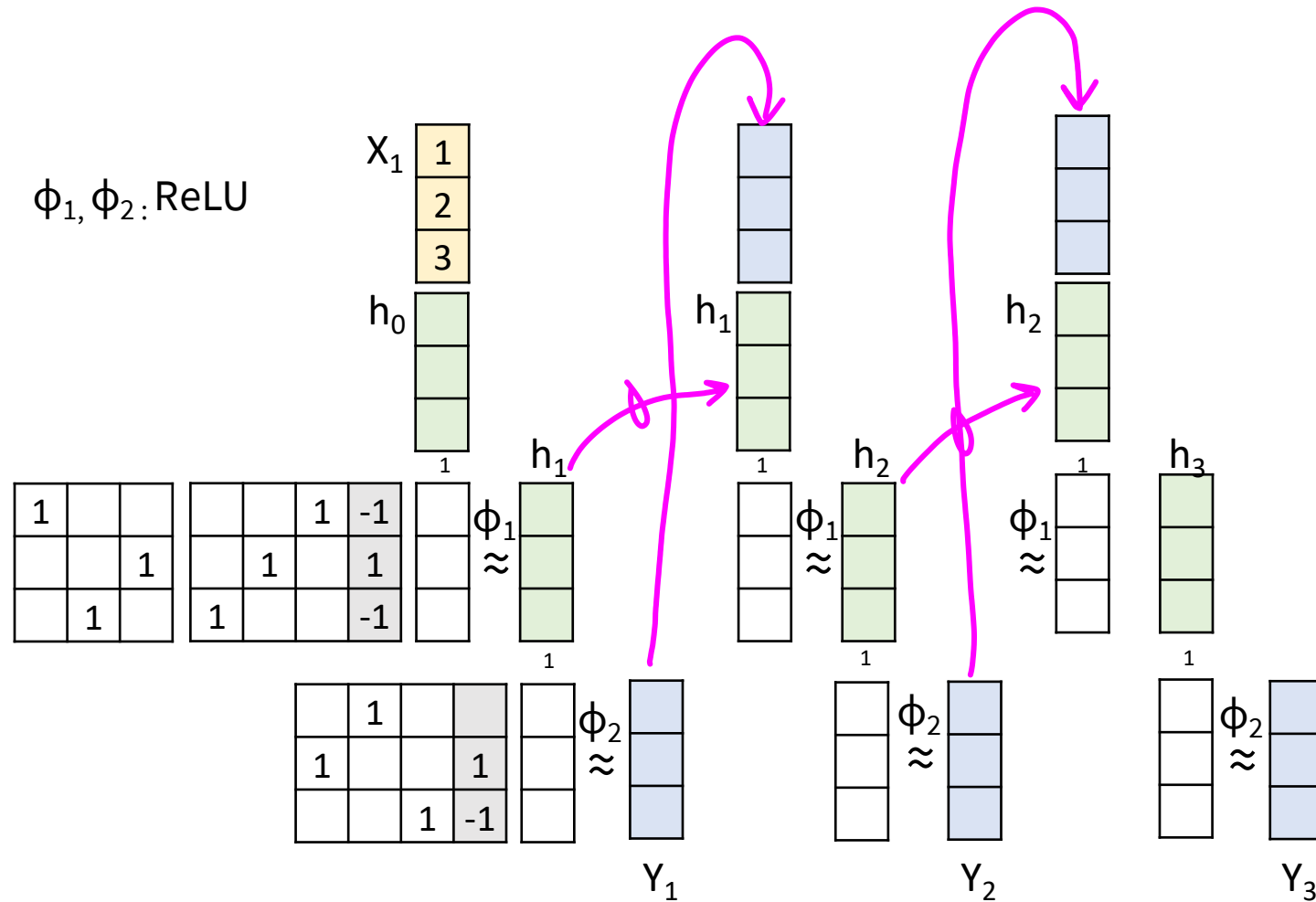


Add a skip connection

Draw a new arrow to illustrate the skip connection from X_1 .

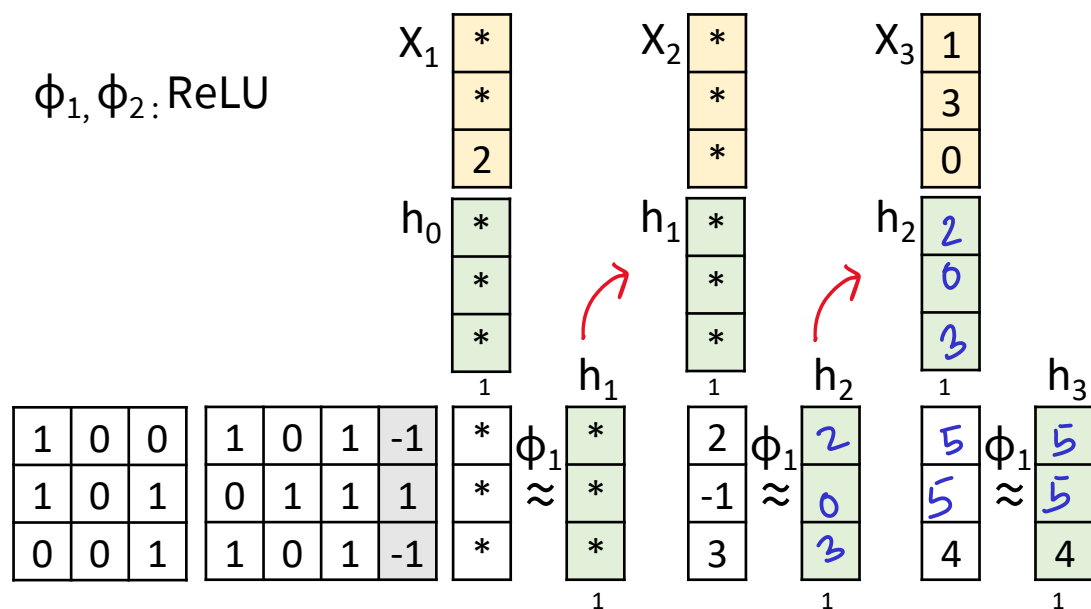


Auto-Regressive RNN





Sequence to Sequence

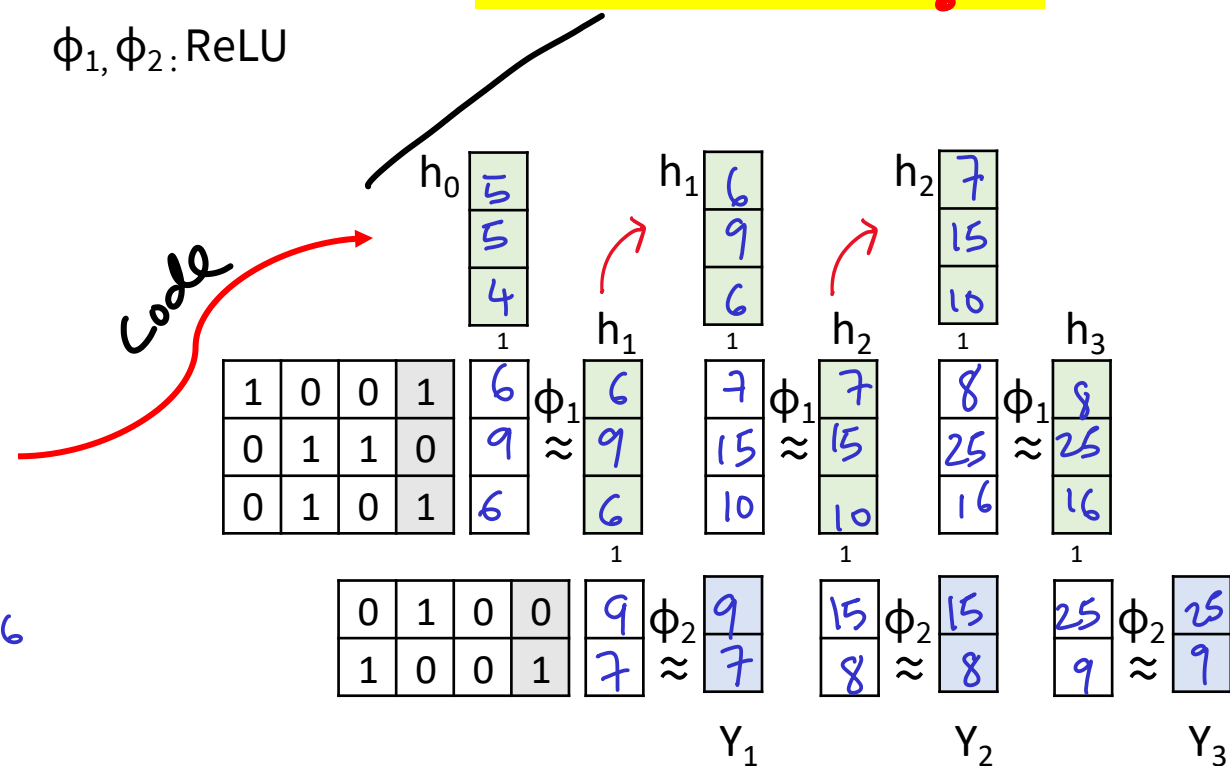
 $\phi_1, \phi_2: \text{ReLU}$ 

$$1+2+3-1=5$$
$$1+3+1=5$$

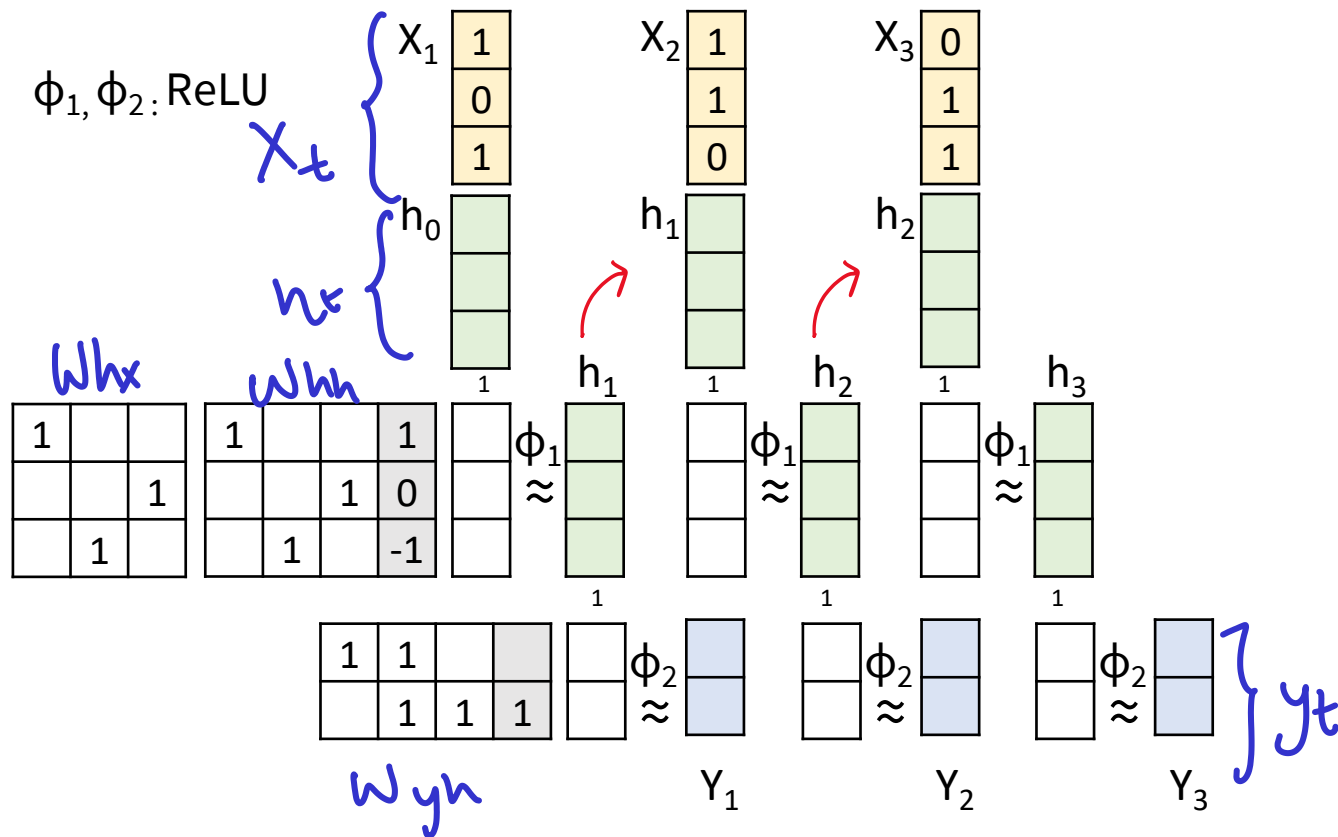
$$5 + 1 = 6$$

 $\phi_1, \phi_2: \text{ReLU}$

Remove unwanted components to make it a seq-to-seq model ?



☒ ☐ Counting Parameters (small)



$$\text{size}(X_t) = \underline{3 \times 1}$$

$$\text{size}(h_t) = \underline{3 \times 1}$$

$$\text{size}(y_t) = \underline{2 \times 1}$$

$$\text{size}(W_{hx}) = \underline{3 \times 3}$$

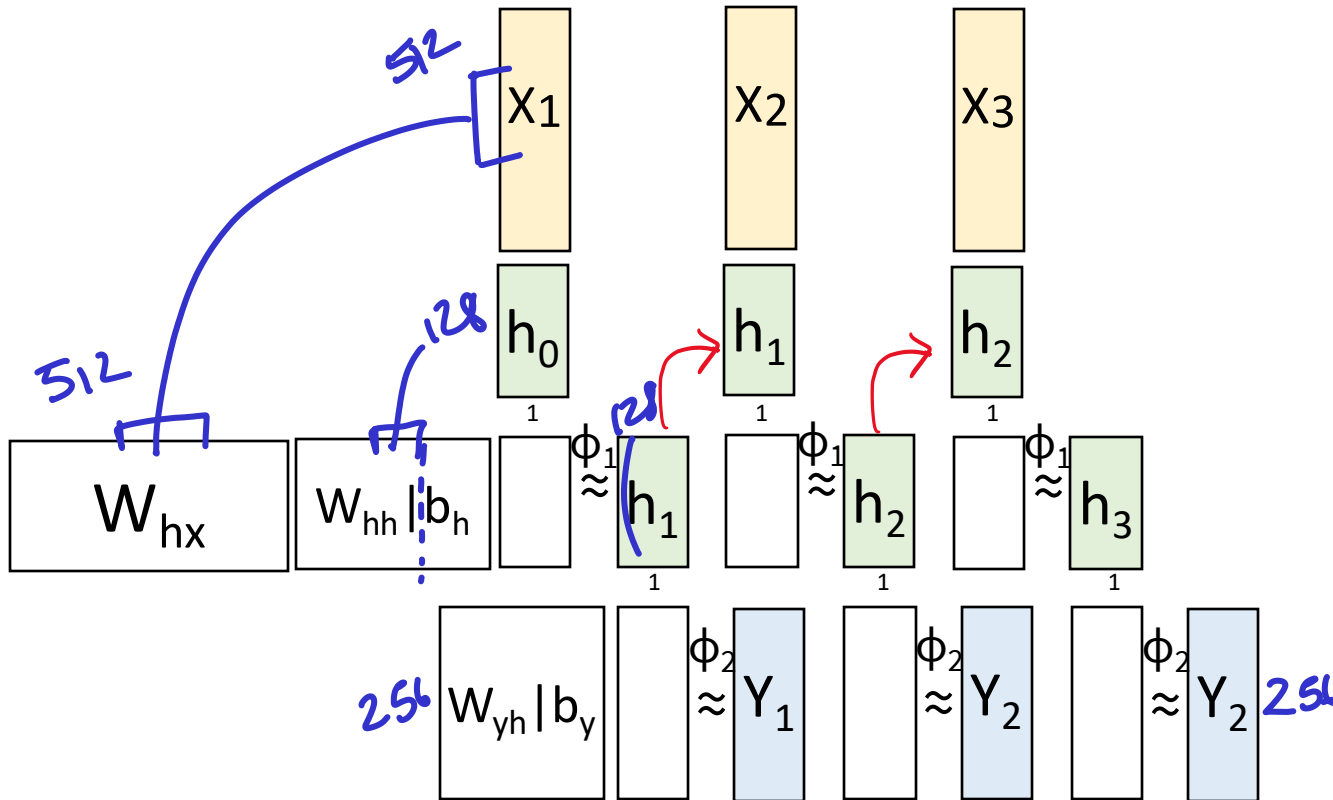
$$\text{size}(b_h) = \underline{3 \times 1}$$

$$\text{size}(W_{hh}) = \underline{3 \times 3}$$

$$\text{size}(W_{yh}) = \underline{2 \times 3}$$

$$\text{size}(b_y) = \underline{2 \times 1}$$

☒ ☐ Counting Parameters (large)



$$\text{size}(X_t) = 512 \times 1$$

$$\text{size}(h_t) = 128 \times 1$$

$$\text{size}(Y_t) = 256 \times 1$$

$$\text{size}(W_{hx}) = \underline{512 \times 128}$$

$$\text{size}(b_h) = \underline{128 \times 1}$$

$$\text{size}(W_{hh}) = \underline{128 \times 128}$$

$$\text{size}(W_{yh}) = \underline{256 \times 128}$$

$$\text{size}(b_y) = \underline{256 \times 1}$$

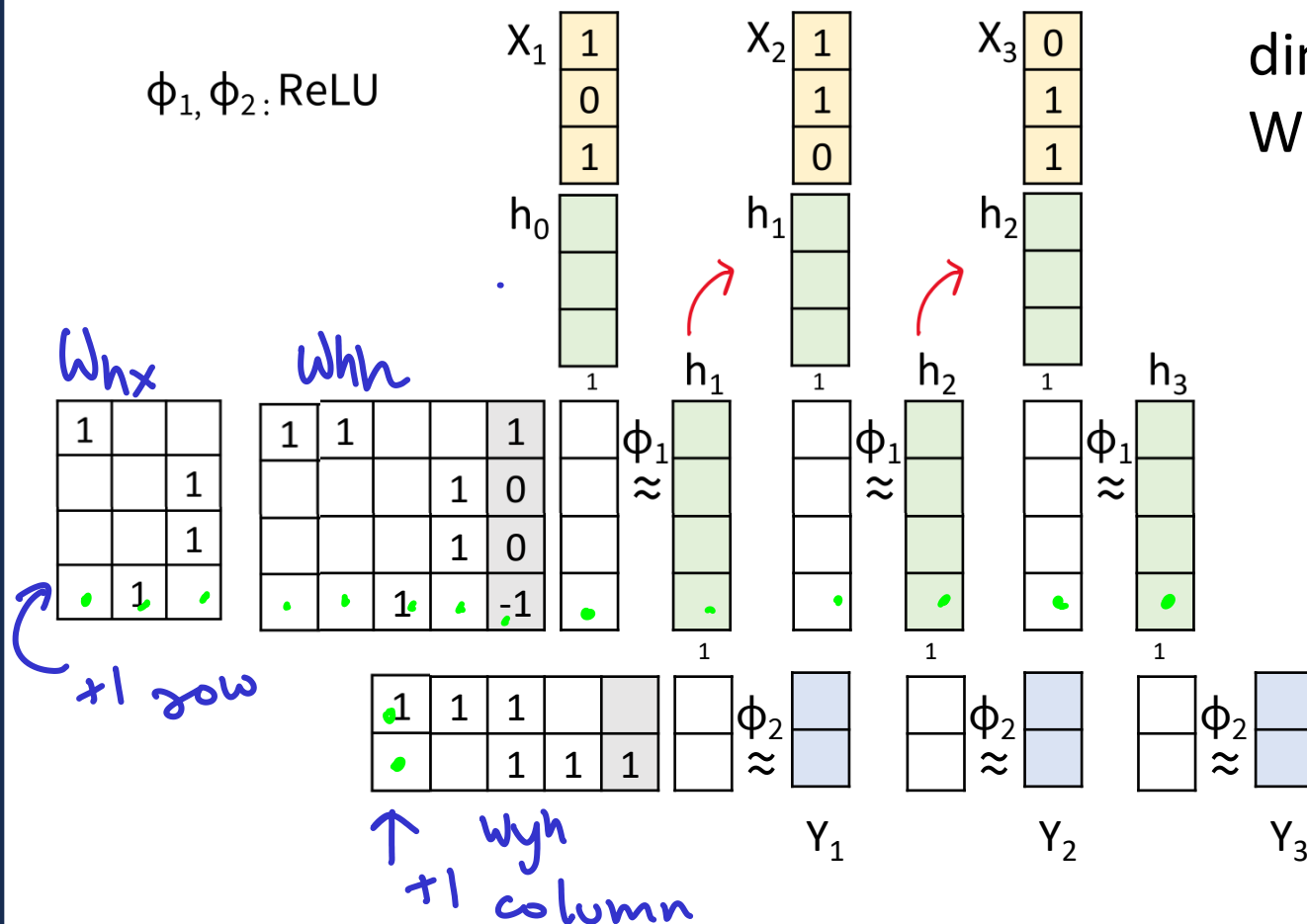


Adding Parameters

I copy pasted a screenshot

Suppose we increase the hidden state's dimension by 1.

What would be the new parameter sizes?



(Hint: You can try to draw the extra cells as visual aid)

$$\text{size}(W_{hx}) = 4 \times 3$$

$$\text{size}(b_h) = 4 \times 1$$

$$\text{size}(W_{hh}) = 4 \times 4$$

$$\text{size}(W_{yh}) = 2 \times 4$$

$$\text{size}(b_y) = 2 \times 1$$

$$\# \text{ of New Parameters} = 16$$

Modeling Probabilities

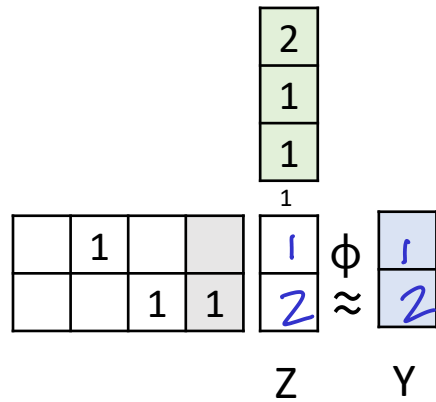
CSCI 5722 Computer Vision



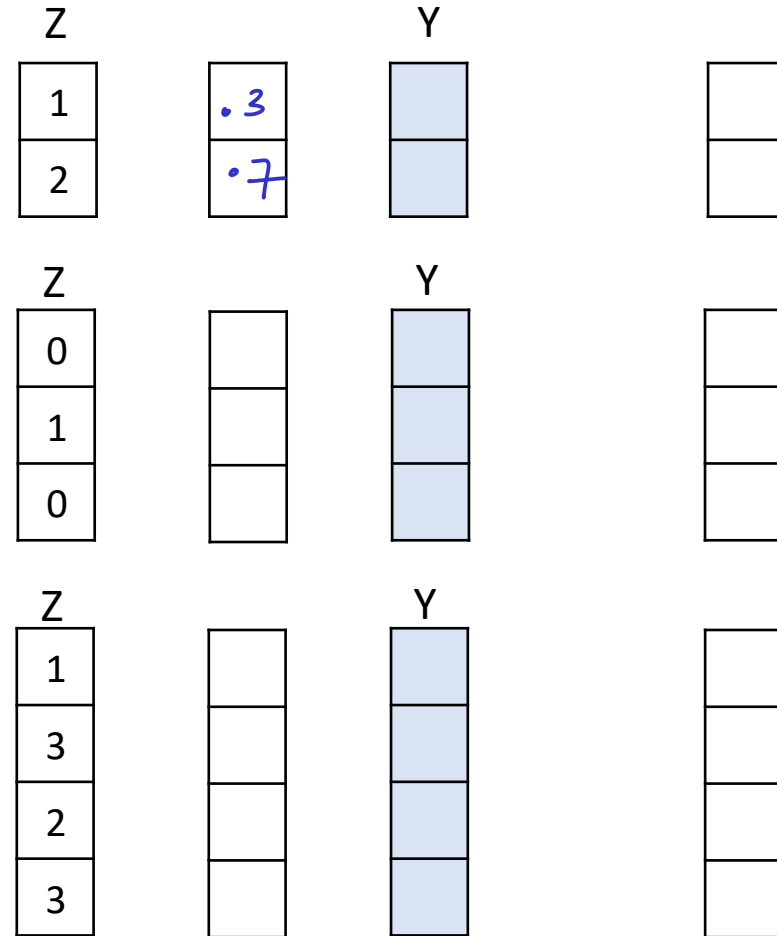
University of Colorado
Boulder

Output Values \rightarrow Probability Distribution

$\phi = \text{ReLU}$



$$Y = \phi(Z) =$$



x	e^x	round
0	1	1
1	2.71828	3
2	7.38906	7
3	20.08554	20
4	54.58715	55
5	148.41316	148

Gradient of Softmax + CE Loss

$$\frac{\partial L}{\partial Z} = \underline{\hspace{2cm}}$$

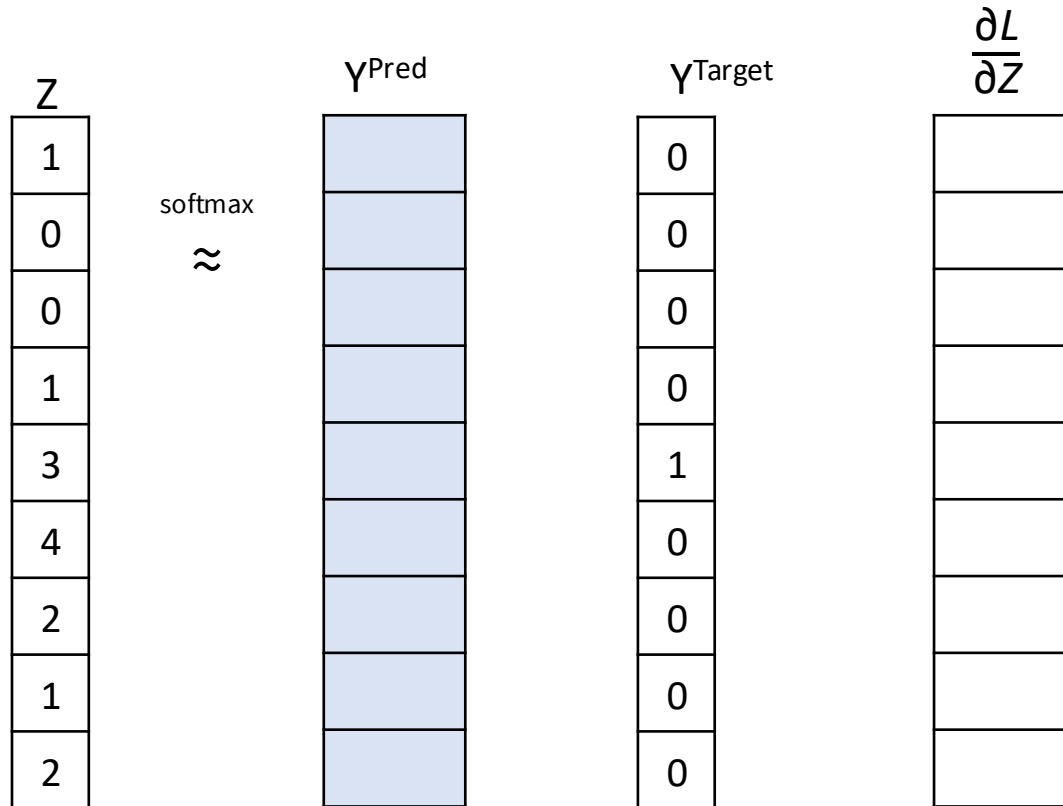
$$\text{softmax} = S(y_i) = \frac{e^{y_i}}{\sum e^{y_j}}$$

$$\text{CE loss} = -\frac{1}{N} \sum y_i \log(\hat{y}_i)$$

Z		y^{Pred}	y^{Target}	
1	softmax \approx	.3	0	
2		.7	1	
0	softmax \approx	.2	0	
1		.6	1	
0		.2	0	
1	softmax \approx	.06	0	
3		.4	1	
2		.14	0	
3		.4	0	

Calculate the Gradient of Softmax + CE Loss

x	e^x	round
0	1	1
1	2.71828	3
2	7.38906	7
3	20.08554	20
4	54.58715	55
5	148.41316	148





Calculate the Gradient of Softmax + CE Loss

x	e^x	round
0	1	1
1	2.71828	3
2	7.38906	7
3	20.08554	20
4	54.58715	55
5	148.41316	148

Z	softmax	Y _{Pred}	Y _{Target}	$\frac{\partial L}{\partial Z}$
0 ✓	≈	0.01	0	0.01
1 ✓		0.03	0	0.03
4 ✓		0.55	0	-0.55
1 ✓		0.03	0	0.03
0 ✓		0.01	1	0.01
3 ✓		0.20	0	0.20
2 ✓		0.07	0	0.07
2 ✓		0.07	0	0.07
1 ✓		0.01	0	0.02

$$CE \text{ loss} = -\frac{1}{N} \sum y_i \log(\hat{y}_i)$$

$$S(y_i) = \frac{e^{y_i}}{\sum e^{y_i}}$$

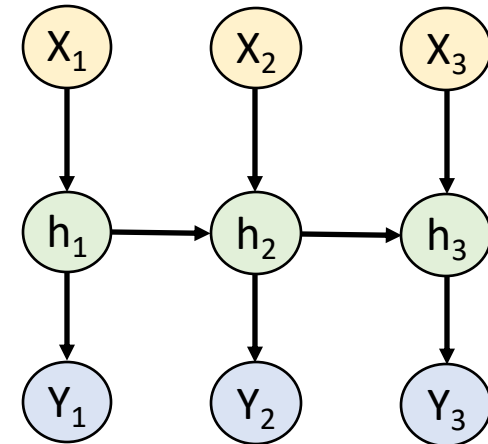
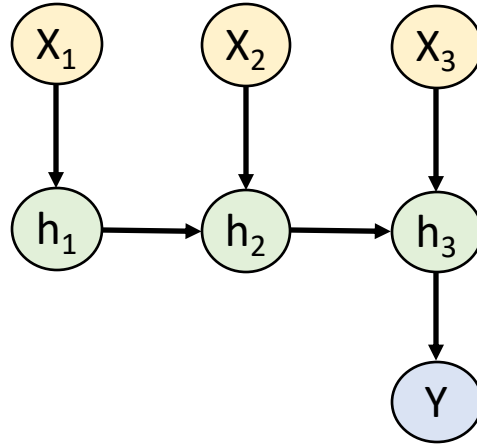
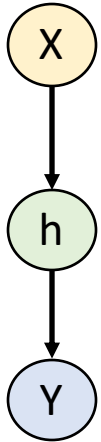
$$2e^0 + 3e^1 + 2e^2 + 1e^3 + 1e^4$$

$$= 2(1) + 3(3)$$

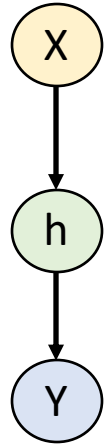
$$+ 2(7) + 20 + 55$$

$$= 100$$

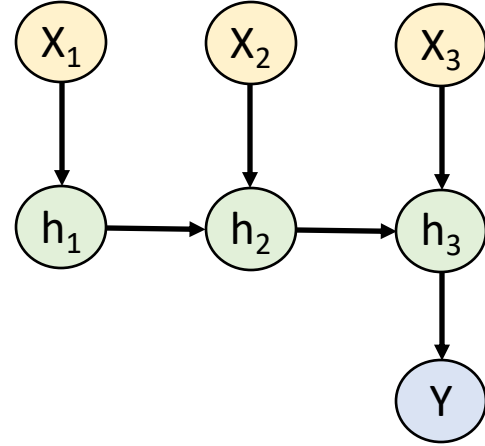
Model \rightarrow Function



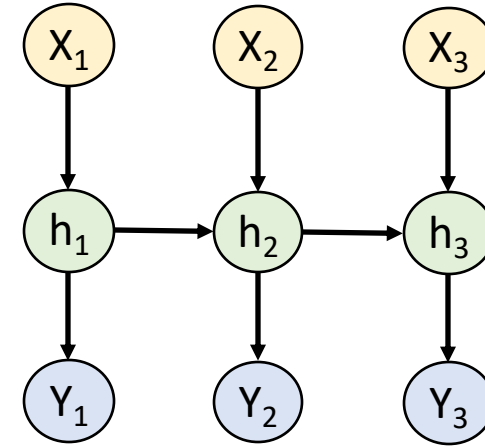
Function \rightarrow Conditional Probability Function



$$Y = f(X)$$



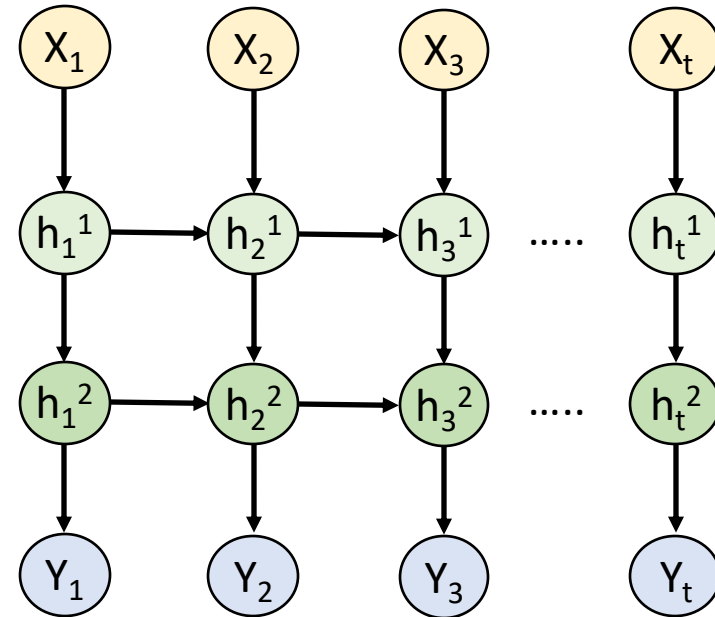
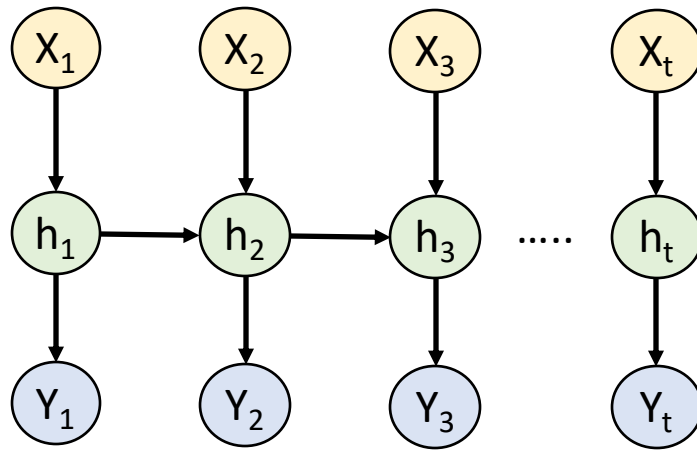
$$Y = f(X_1, X_2, X_3)$$



$$Y = f(Y_t | X_t, X_{t-1}, \dots X_1)$$

Model \rightarrow Conditional Probability Distributions

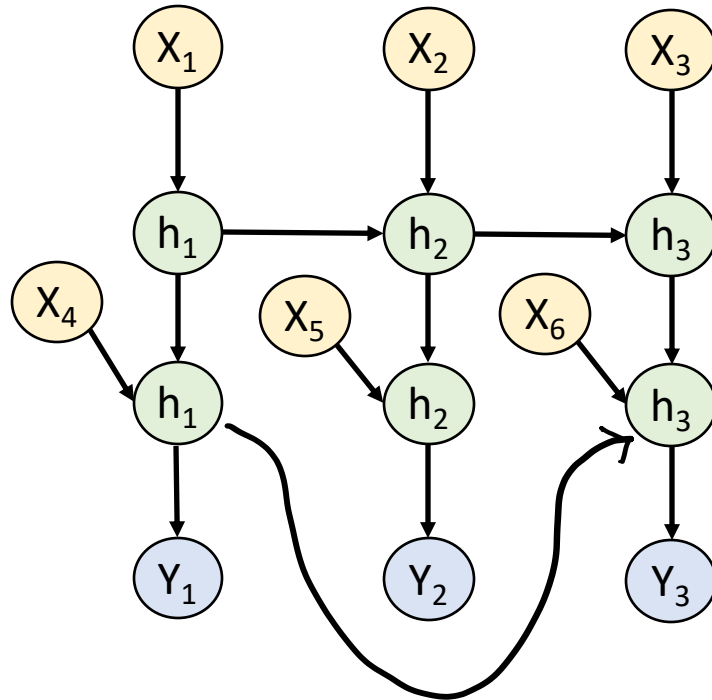
What conditional probability distributions could these two models estimate?



$P(Y_t \mid$)

$P(Y_t \mid$)

Model \rightarrow Conditional Probability Distributions



What conditional probability distributions could this model learn to estimate?

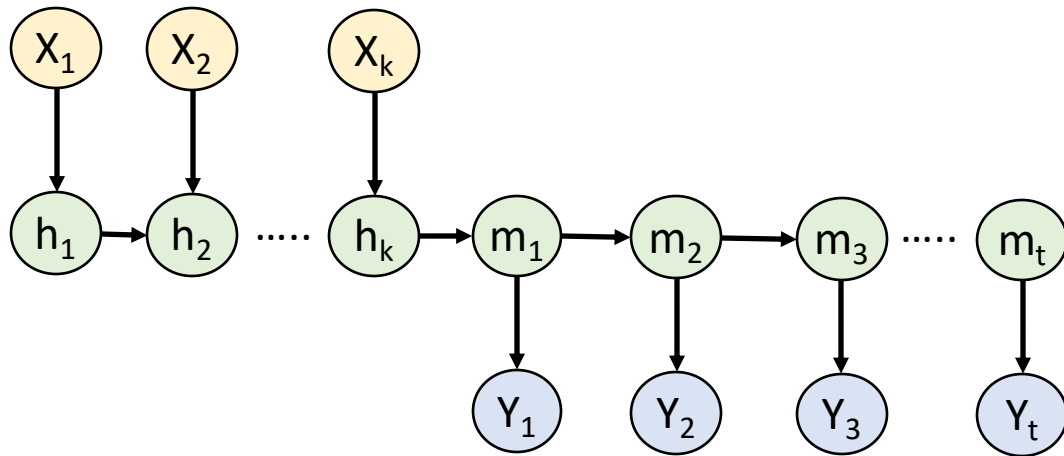
$$P(Y_1 \mid \quad \quad \quad)$$

$$P(Y_2 \mid \quad \quad \quad)$$

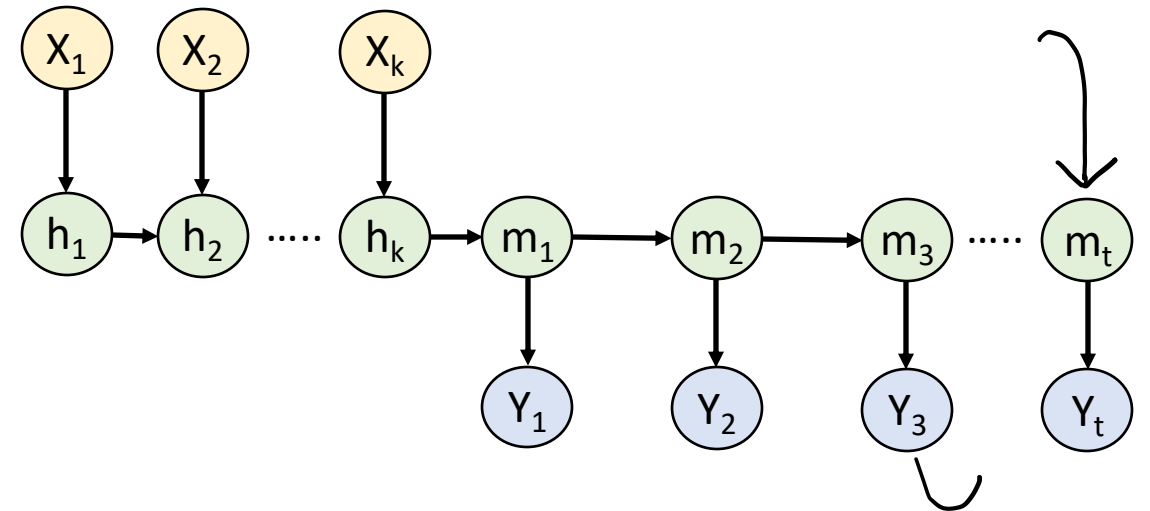
$$P(Y_3 \mid \quad \quad \quad)$$

Model \rightarrow Conditional Probability Distributions

What conditional probability distributions could these two seq-to-seq models estimate?

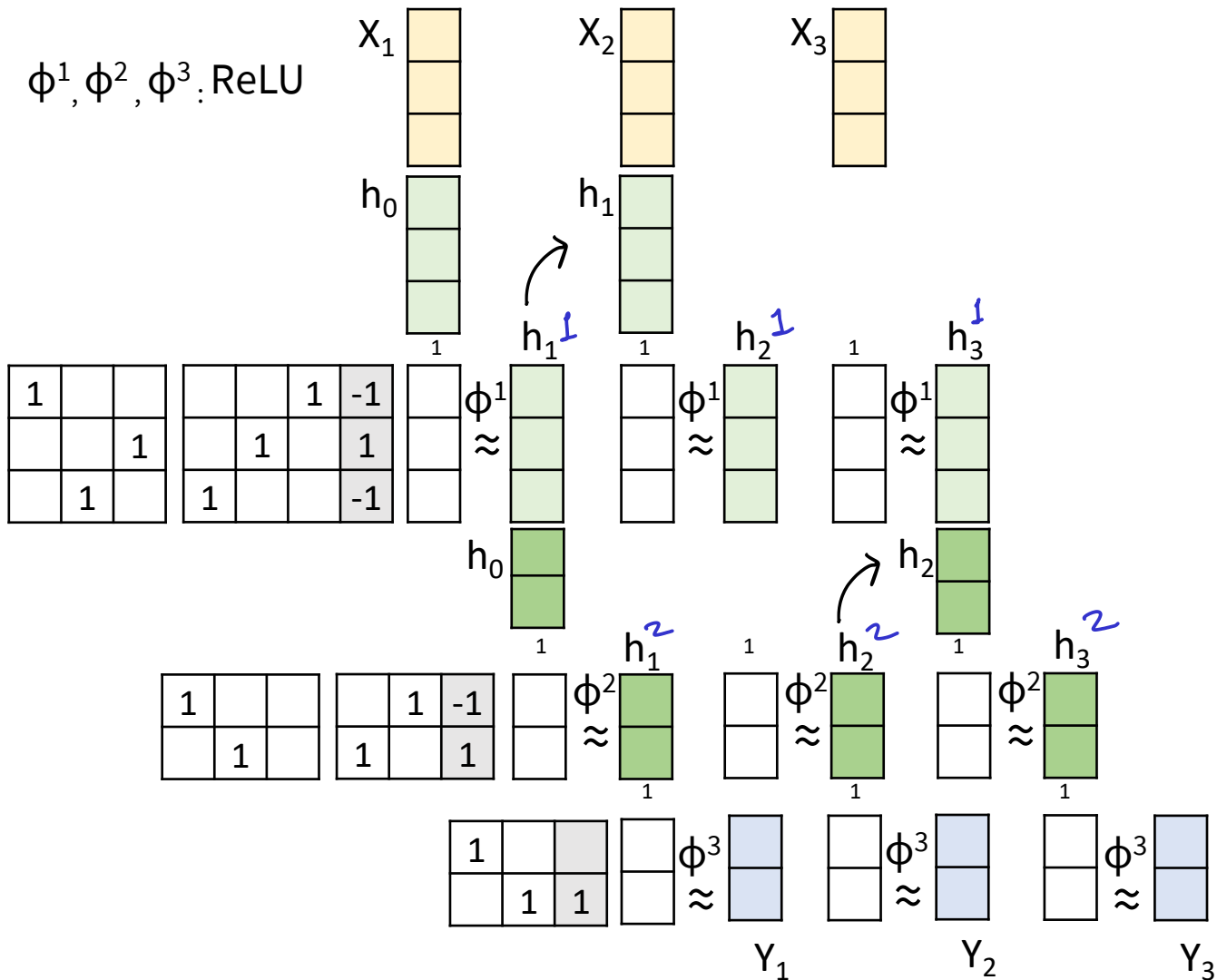


$P(Y_t \mid$)



$P(Y_t \mid$)

Identify “No Dependency” Links



Cross out the dependency links to match the matrix form.

