

AI by Hand



25 Exercise Workbook Series

1. Dot Product
2. Matrix Multiplication
3. Linear Layer
4. Activation
5. Artificial Neuron
6. Batch
7. Connection
8. Hidden Layer
9. Deep
10. Wide
11. Softmax
12. Gradient

More to come ...

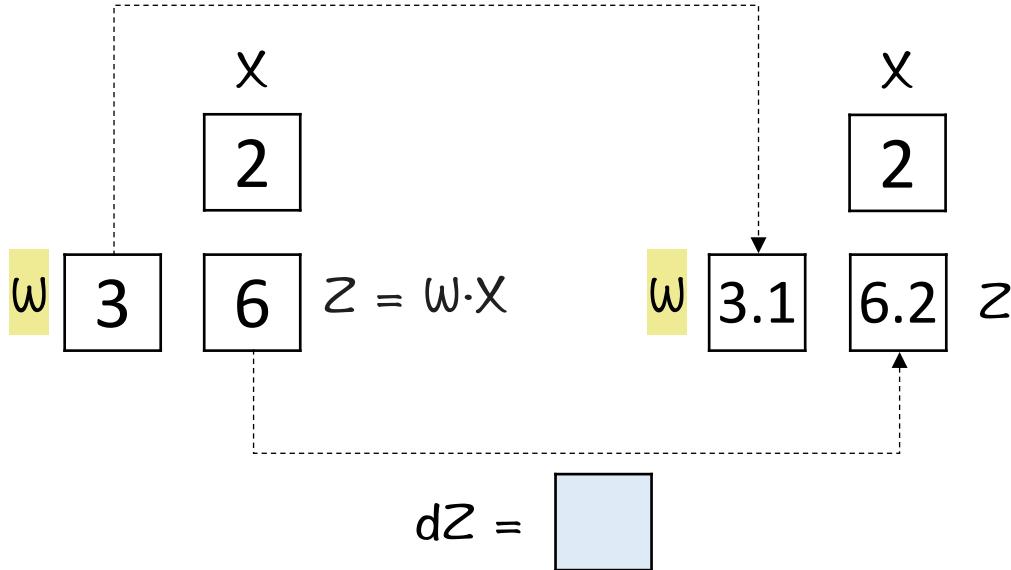
<http://by-hand.ai/workbook>



Gradient

Exercise 1

$$d\omega = \boxed{0.1}$$

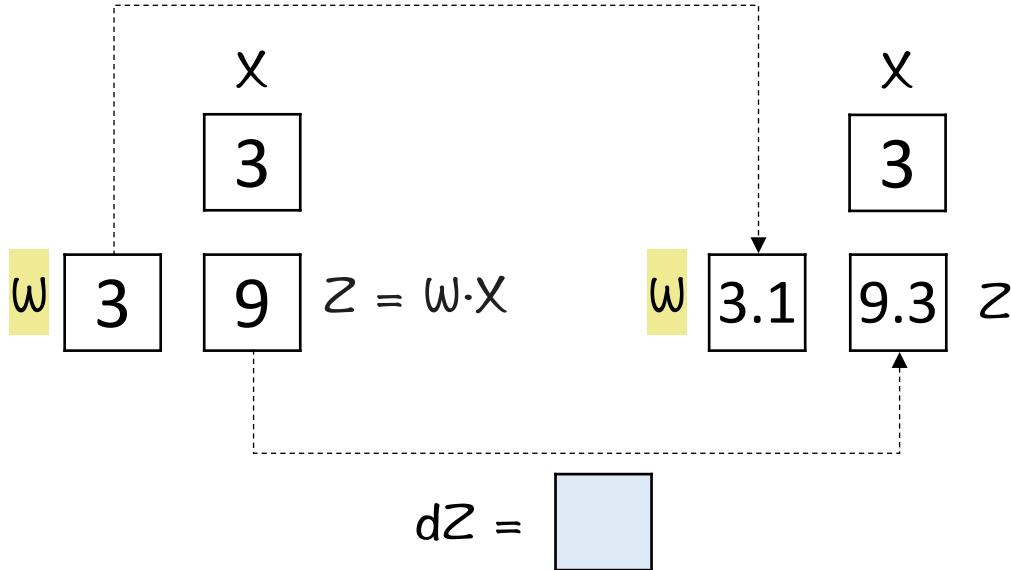


$$\frac{dz}{d\omega} = \frac{\boxed{}}{\boxed{0.1}} = \boxed{} = x$$

Gradient

Exercise 2

$$d\omega = \boxed{0.1}$$

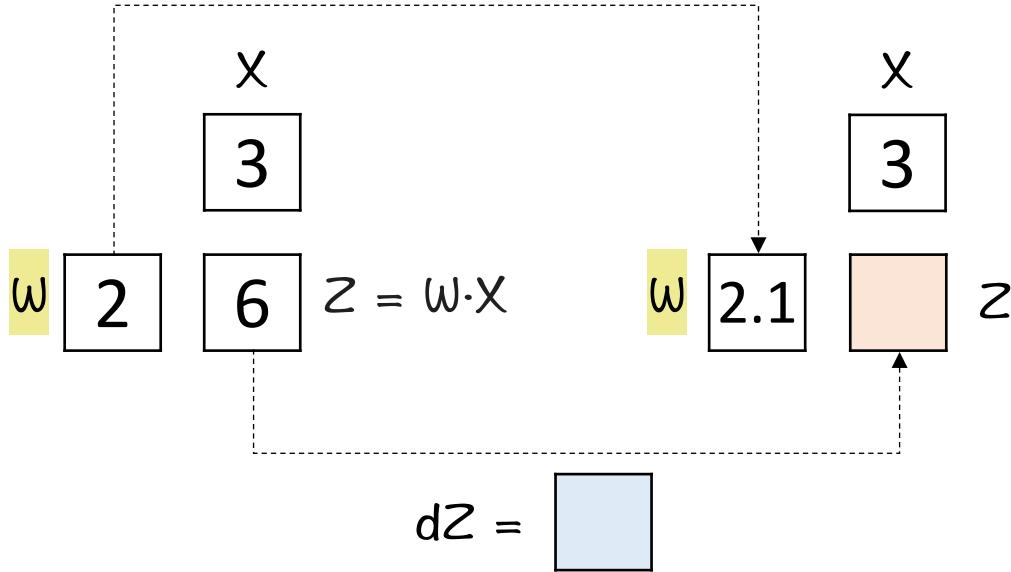


$$\frac{dZ}{d\omega} = \frac{\boxed{\quad}}{0.1} = \boxed{\quad} = X$$

Gradient

Exercise 3

$$d\omega = \boxed{0.1}$$

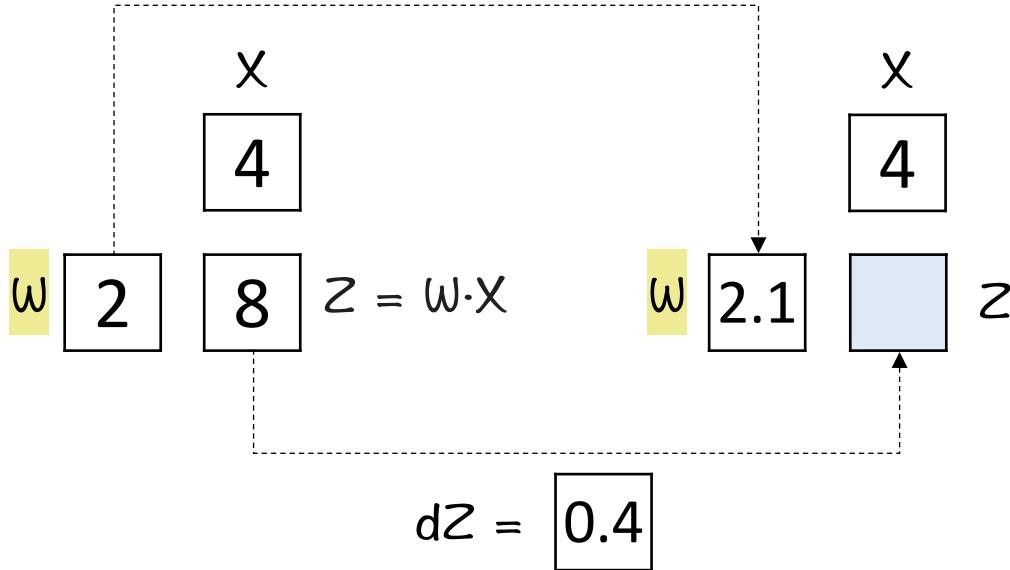


$$\frac{dz}{d\omega} = \frac{\boxed{}}{0.1} = \boxed{3} = x$$

Gradient

Exercise 4

$$d\omega = \boxed{0.1}$$

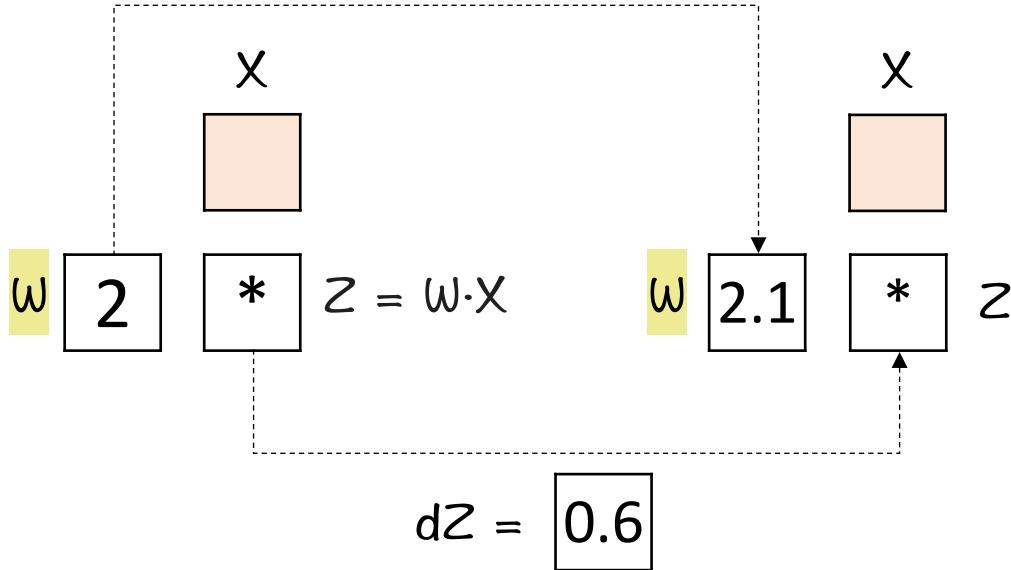


$$\frac{dZ}{d\omega} = \frac{\boxed{0.4}}{\boxed{0.1}} = \boxed{\quad} = x$$

Gradient

Exercise 5

$$d\omega = \boxed{0.1}$$

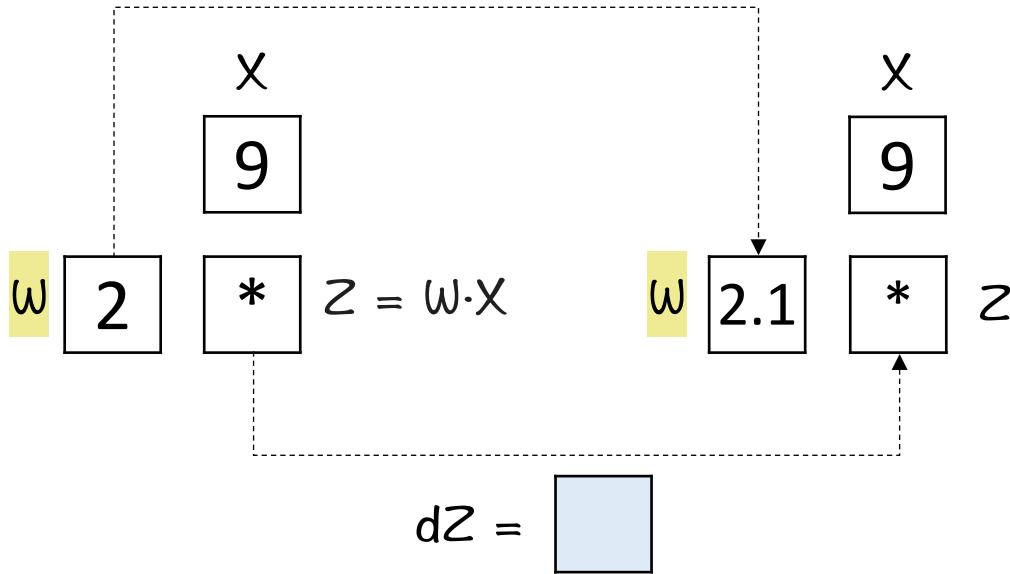


$$\frac{dz}{d\omega} = \frac{\boxed{0.6}}{\boxed{0.1}} = \boxed{} = x$$

Gradient

Exercise 6

$$d\omega = \boxed{0.1}$$

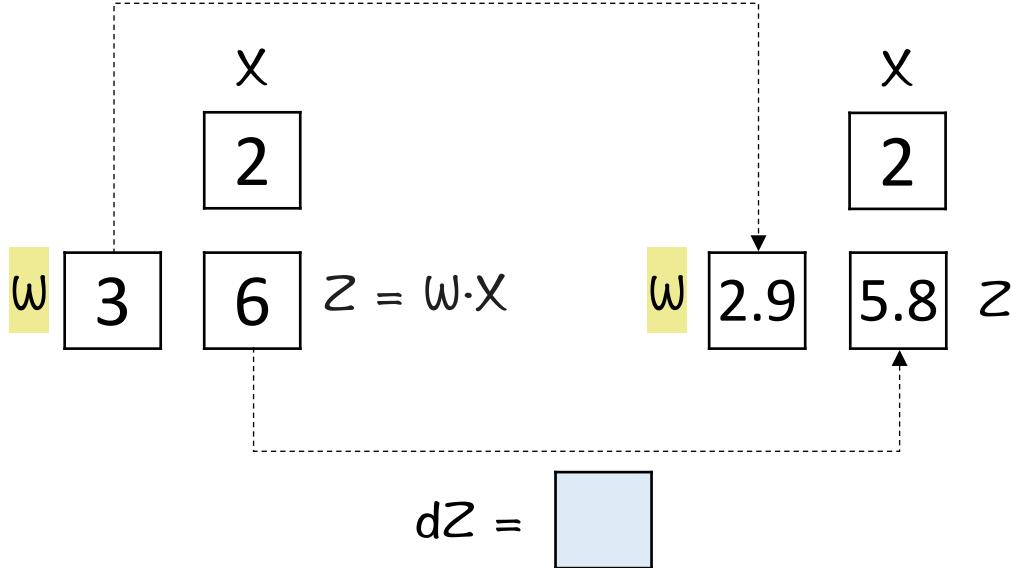


$$\frac{dZ}{d\omega} = \frac{\boxed{\quad}}{\boxed{0.1}} = \boxed{\quad} = X$$

Gradient

Exercise 7

$$d\omega = \boxed{-0.1}$$

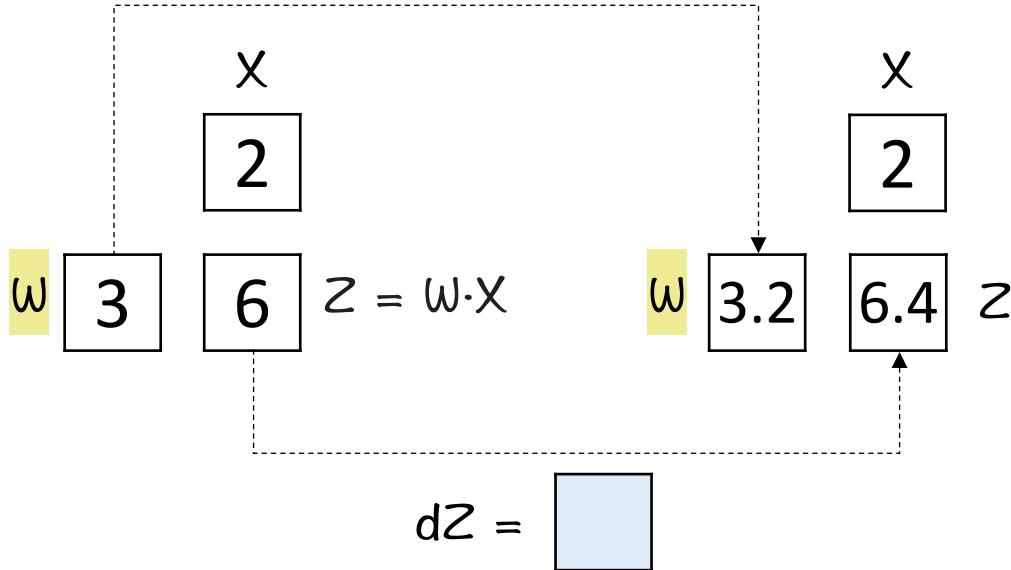


$$\frac{dz}{d\omega} = \frac{\boxed{\quad}}{\boxed{-0.1}} = \boxed{\quad} = X$$

Gradient

Exercise 8

$$d\omega = \boxed{0.2}$$

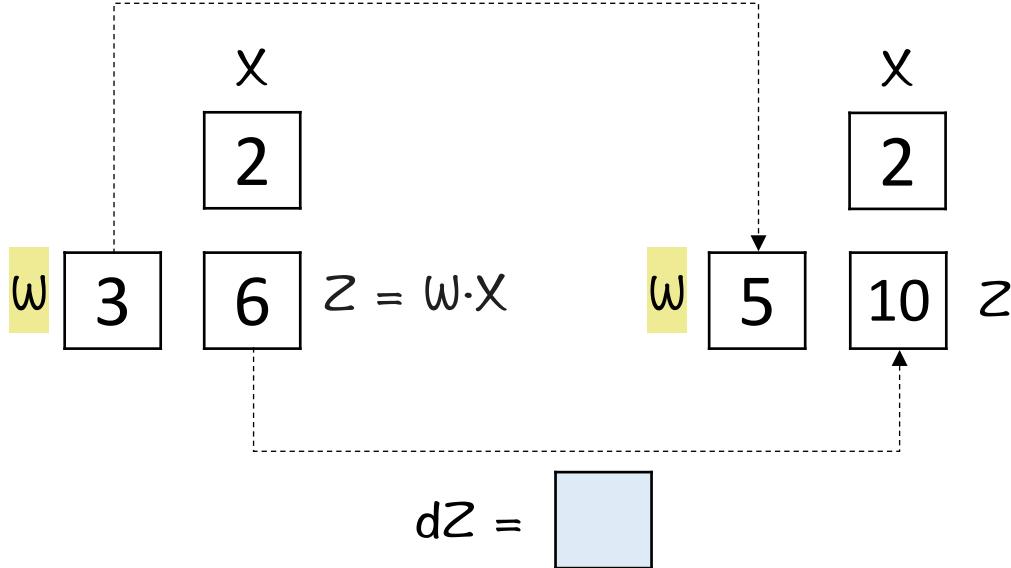


$$\frac{dZ}{d\omega} = \frac{\boxed{\quad}}{\boxed{0.2}} = \boxed{\quad} = X$$

Gradient

Exercise 9

$$d\omega = \boxed{2}$$

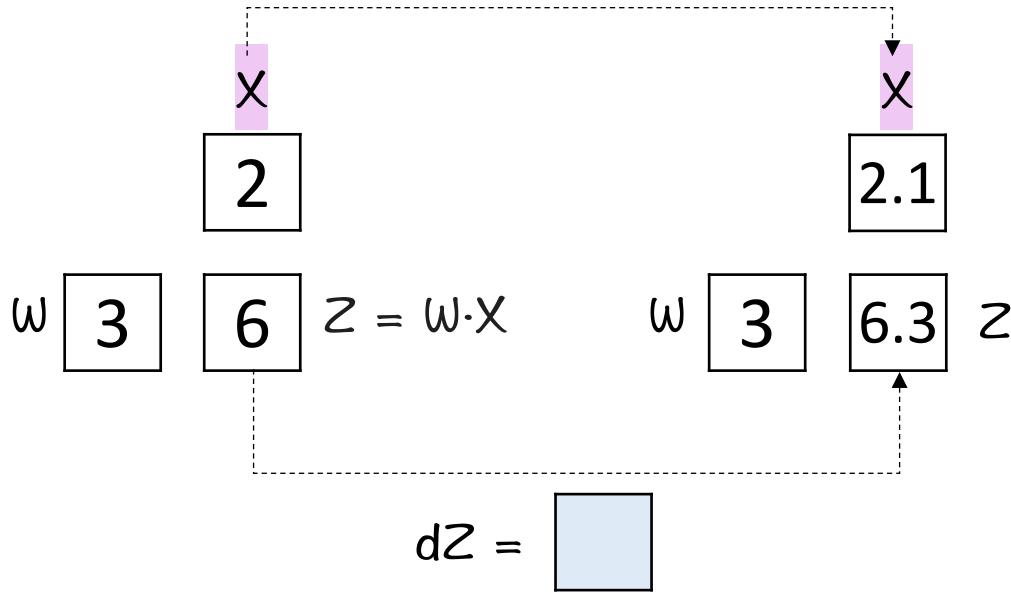


$$\frac{dZ}{d\omega} = \frac{\boxed{\quad}}{\boxed{2}} = \boxed{\quad} = X$$

Gradient

Exercise 10

$$dX = \boxed{0.1}$$

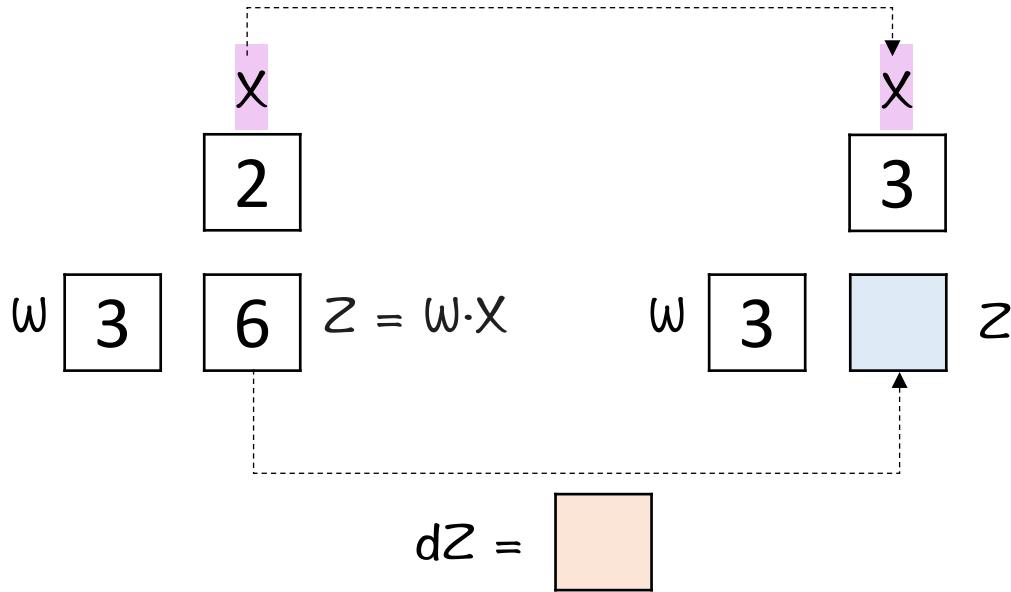


$$\frac{dZ}{dX} = \frac{\boxed{0.1}}{\boxed{0.1}} = \boxed{6} = \omega$$

Gradient

Exercise 11

$$dX = \boxed{1}$$

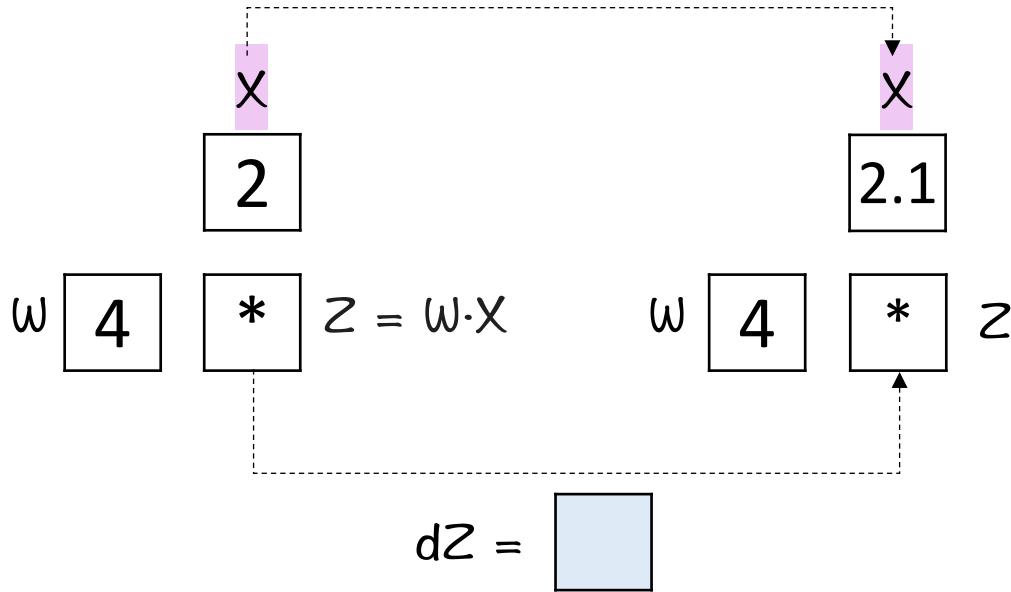


$$\frac{dZ}{dX} = \frac{\boxed{}}{\boxed{1}} = \boxed{3} = w$$

Gradient

Exercise 12

$$dX = \boxed{0.1}$$

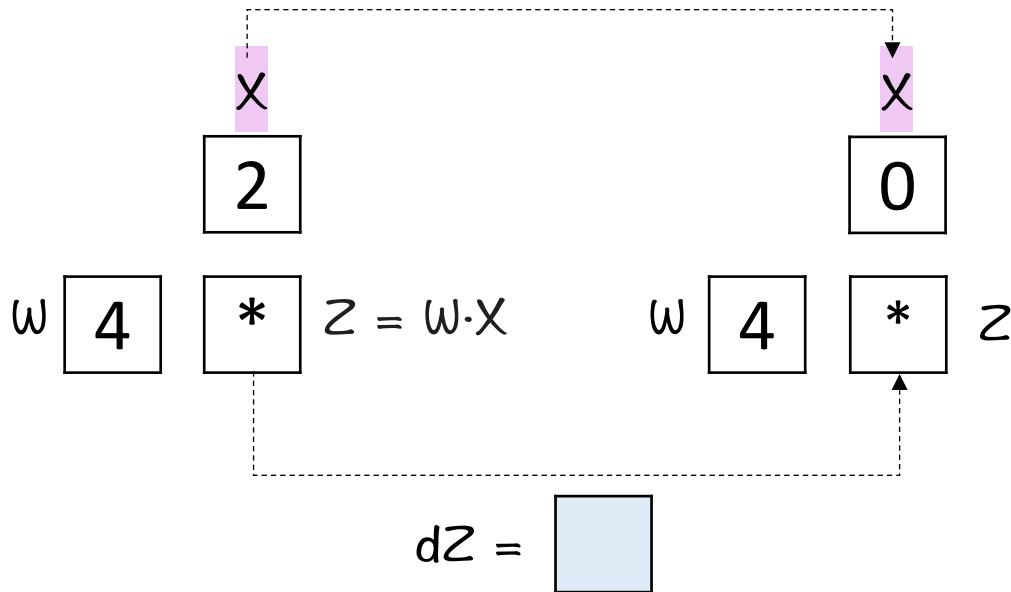


$$\frac{dZ}{dX} = \frac{\boxed{\quad}}{\boxed{0.1}} = \boxed{\quad} = \omega$$

Gradient

Exercise 13

$$dX = \boxed{-2}$$

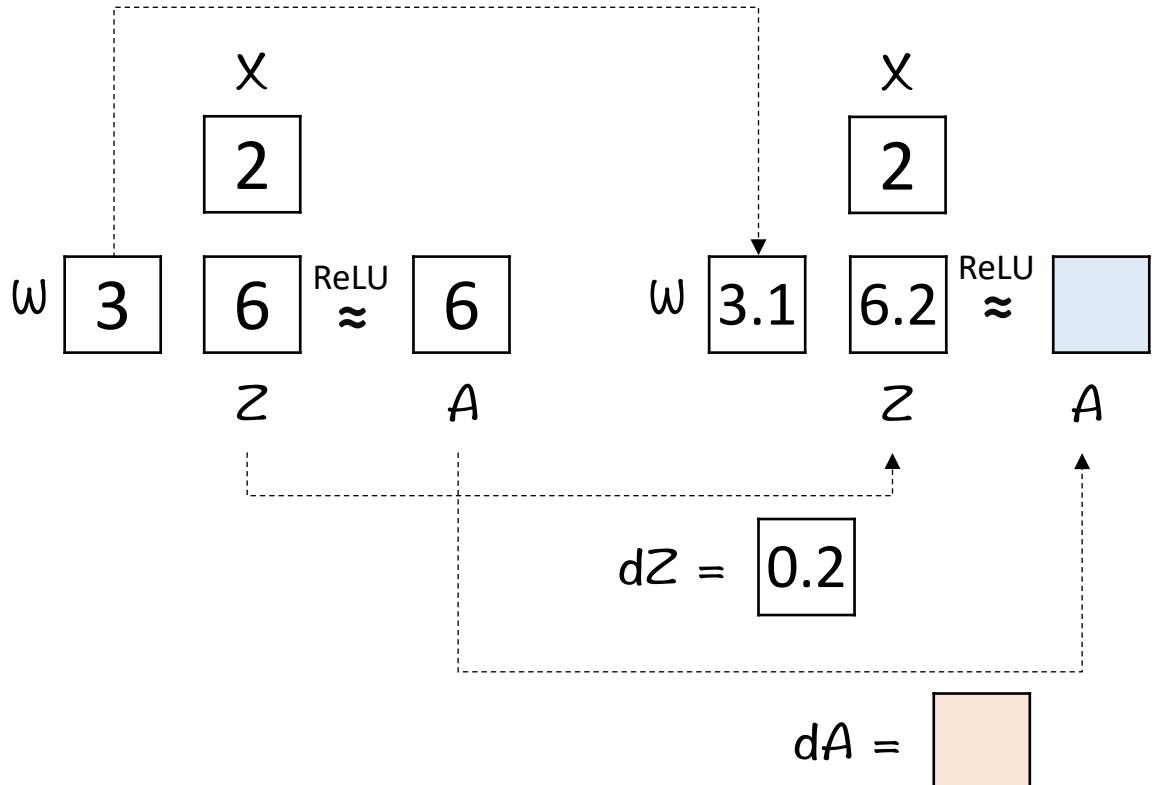


$$\frac{dZ}{dX} = -\frac{\boxed{}}{\boxed{-2}} = \boxed{} = w$$

Gradient

Exercise 14

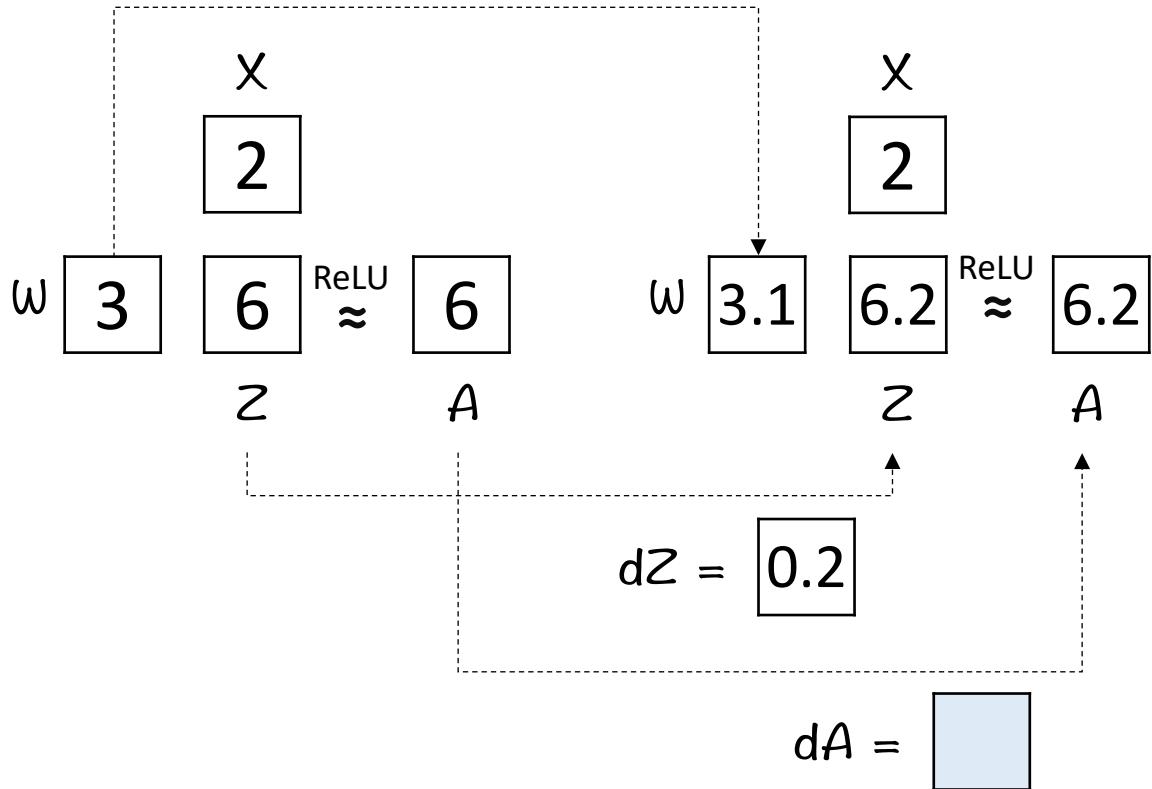
$$dW = \boxed{0.1}$$



Gradient

Exercise 15

$$dW = \boxed{0.1}$$

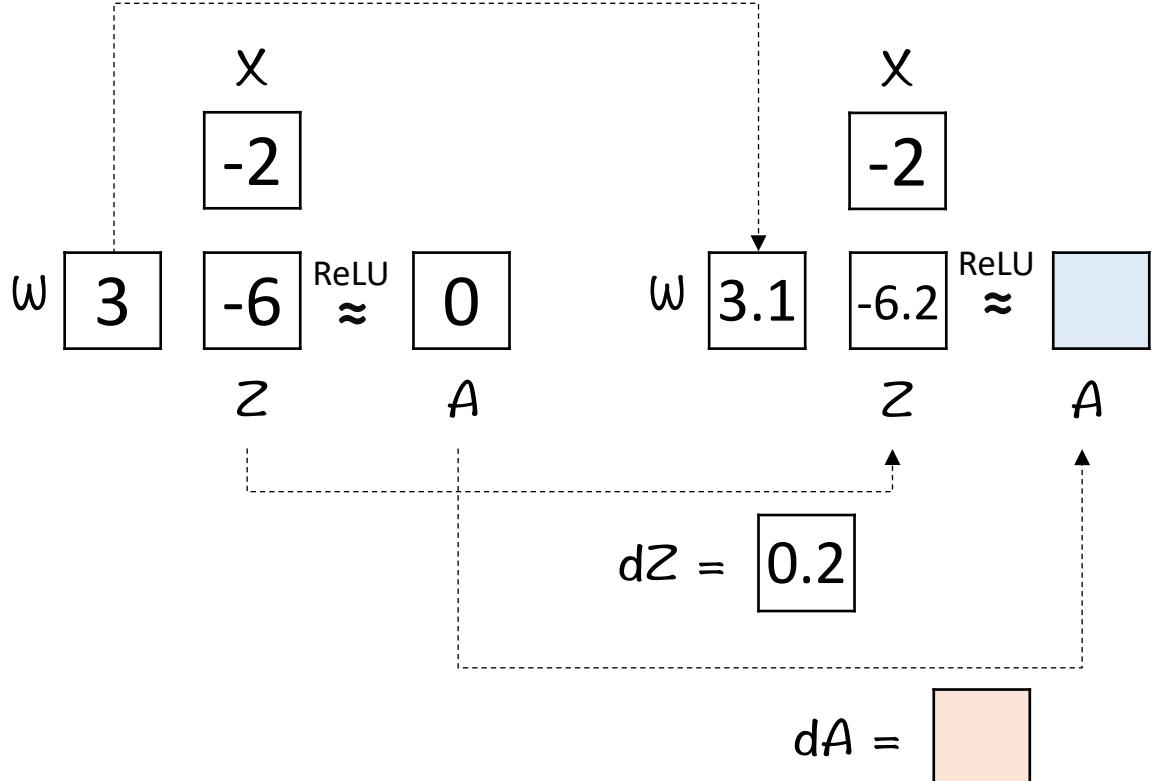


$$\frac{dA}{dZ} = \frac{\boxed{}}{\boxed{0.2}} = \boxed{}$$

Gradient

Exercise 16

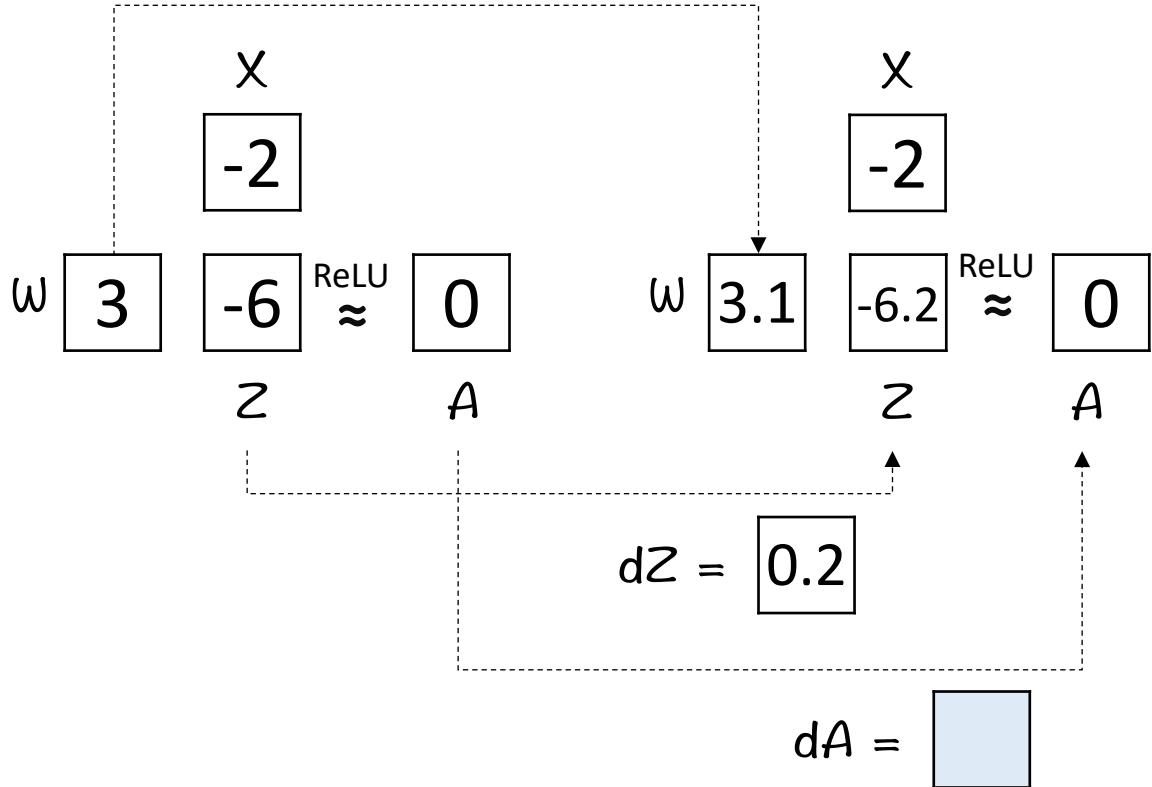
$$dW = \boxed{0.1}$$



Gradient

Exercise 17

$$dW = \boxed{0.1}$$

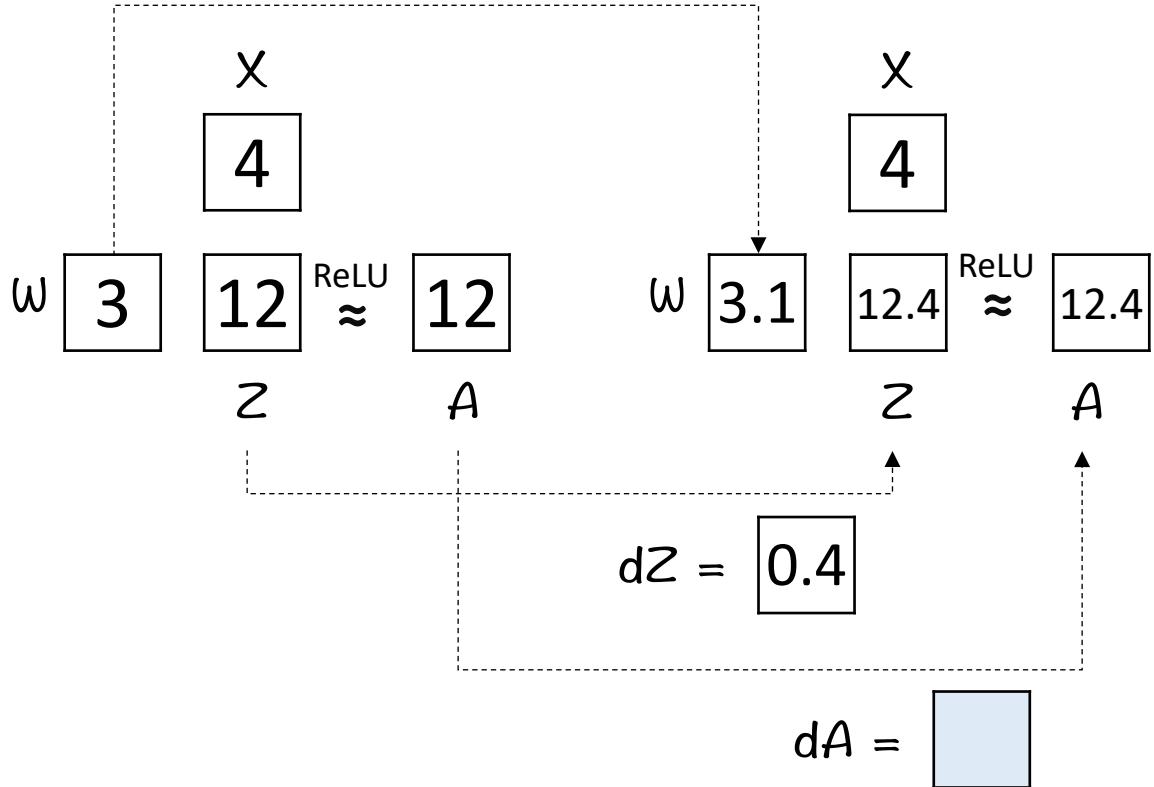


$$\frac{dA}{dZ} = \frac{\boxed{\text{?}}}{\boxed{0.2}} = \boxed{\text{?}}$$

Gradient

Exercise 18

$$dW = \boxed{0.1}$$

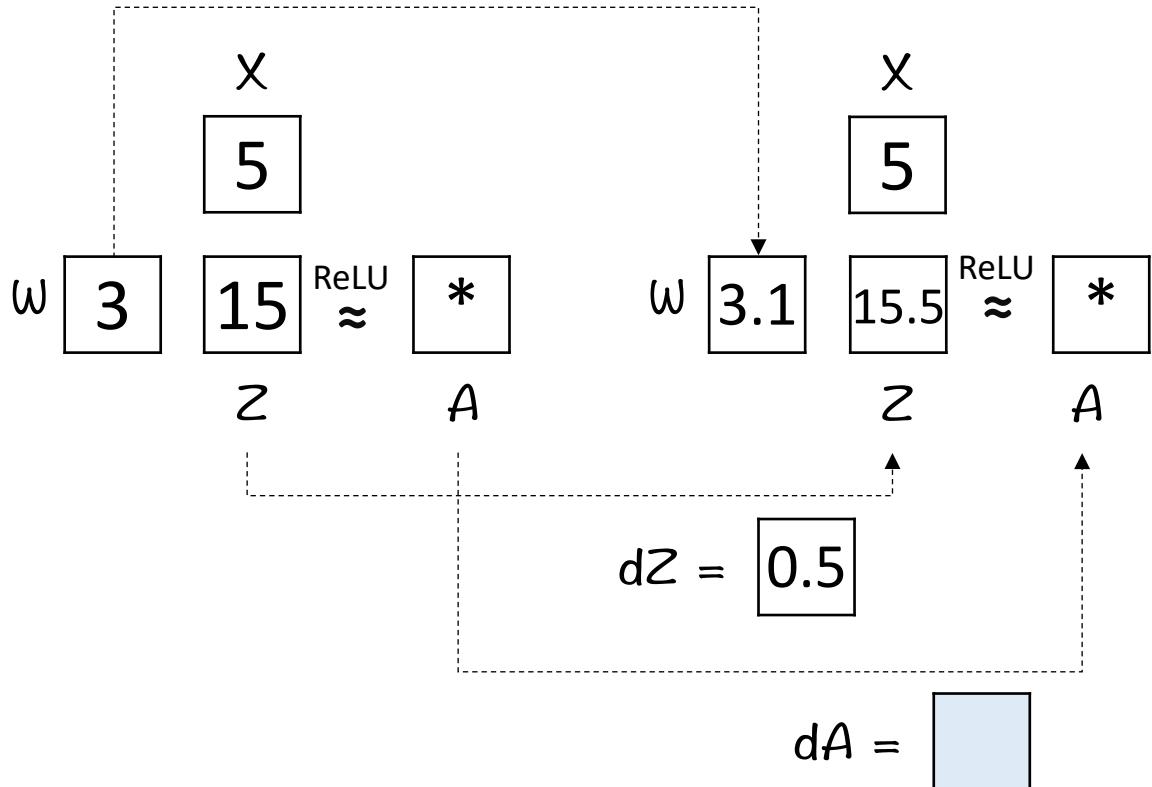


$$\frac{dA}{dZ} = \frac{\boxed{}}{\boxed{0.4}} = \boxed{}$$

Gradient

Exercise 19

$$dW = \boxed{0.1}$$

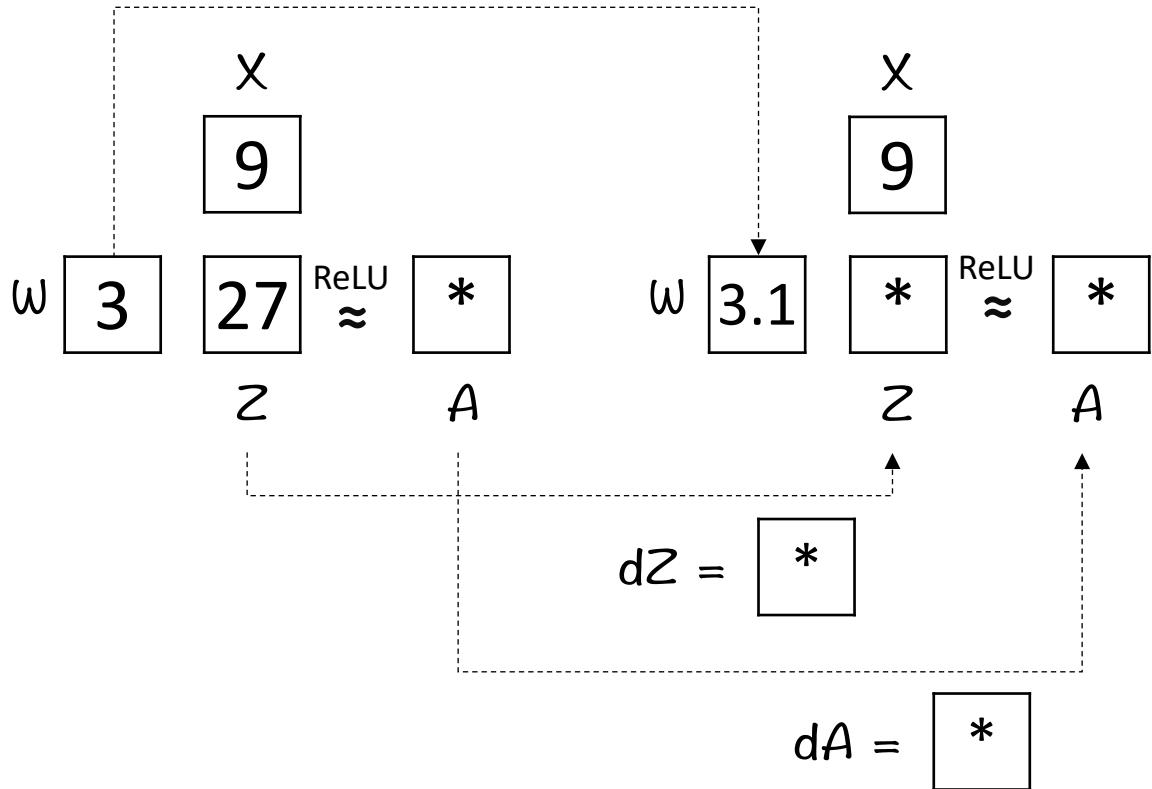


$$\frac{dA}{dZ} = \frac{\boxed{0.5}}{0.5} = \boxed{1}$$

Gradient

Exercise 20

$$dW = \boxed{0.1}$$

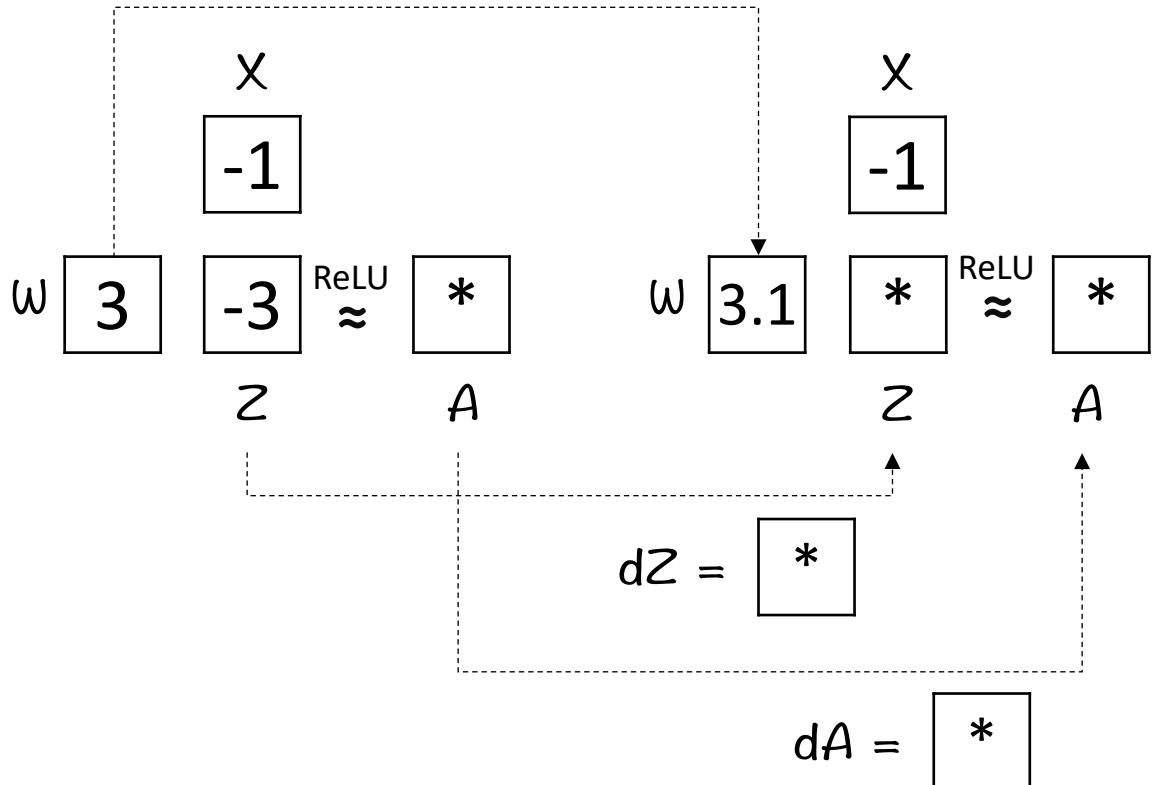


$$\frac{dA}{dZ} = \frac{\boxed{*}}{\boxed{*}} = \boxed{\text{orange box}}$$

Gradient

Exercise 21

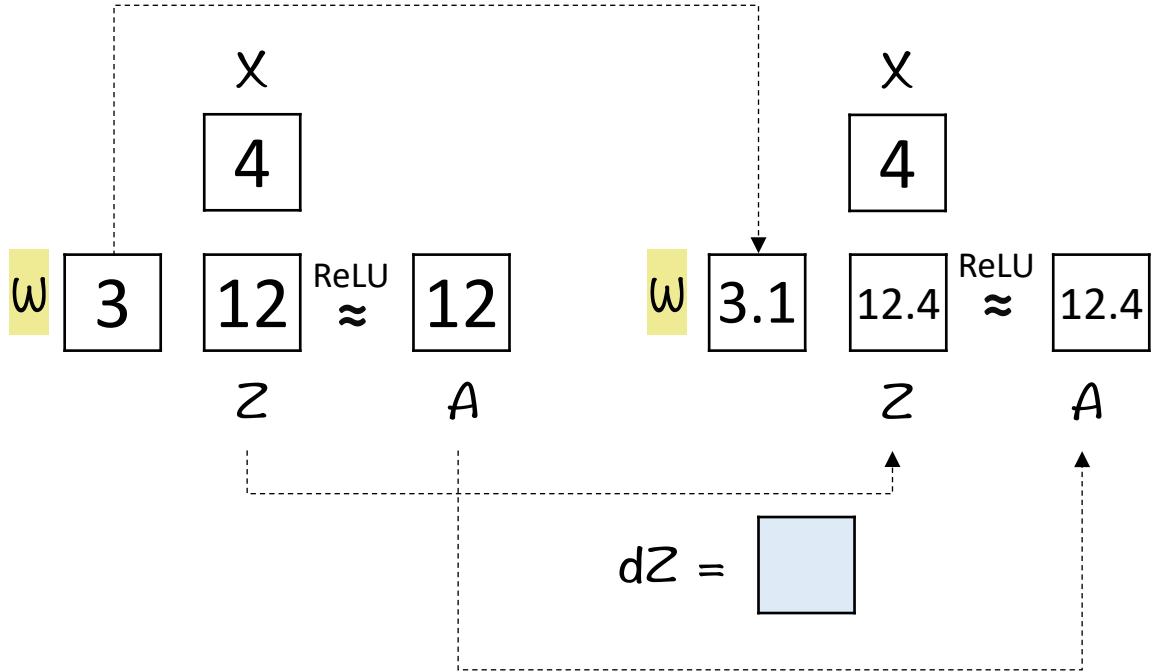
$$dW = \boxed{0.1}$$



Gradient

Exercise 22

$$dW = \boxed{0.1}$$



$$\frac{dA}{dz} = \frac{\boxed{0.4}}{\boxed{1}} = \boxed{1}$$

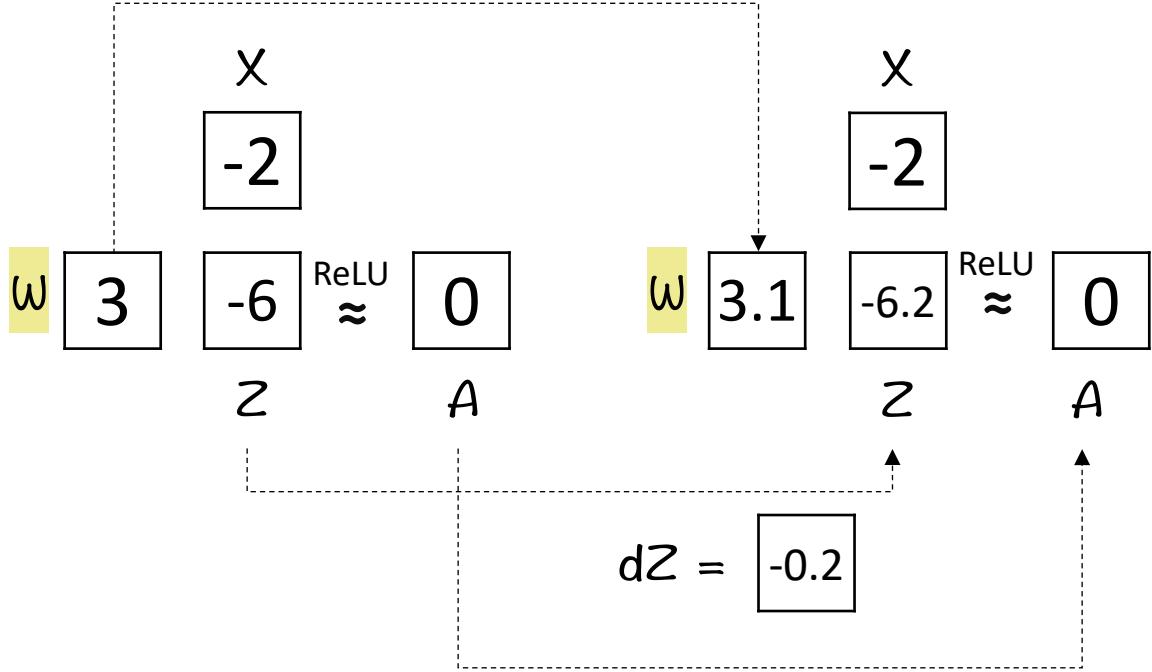
$$\frac{dZ}{dW} = \frac{\boxed{1}}{\boxed{0.1}} = \boxed{4}$$

$$\frac{dA}{dW} = \frac{\boxed{0.4}}{\boxed{0.1}} = \boxed{4}$$

Gradient

Exercise 23

$$dW = \boxed{0.1}$$



$$\frac{dA}{dZ} = \frac{\boxed{0}}{\boxed{-0.2}} = \boxed{0}$$

$$\frac{dZ}{dW} = \frac{\boxed{-0.2}}{\boxed{0.1}} = \boxed{-2}$$

$$\frac{dA}{dW} = \frac{\boxed{0}}{\boxed{0.1}} = \boxed{\quad}$$

Gradient

Exercise 24

$$dX = \boxed{0.1}$$

$$\begin{array}{ccc} & \boxed{x} & \\ & \downarrow & \\ \boxed{4} & & \boxed{4.1} \\ & \uparrow & \\ \omega \boxed{3} & \boxed{12} & \xrightarrow{\text{ReLU}} \boxed{12} & \quad \omega \boxed{3} & \boxed{12.3} & \xrightarrow{\text{ReLU}} \boxed{12.3} \\ & z & & & z & \\ & \downarrow & & & \uparrow & \\ & & & dZ = \boxed{} & & \end{array}$$

$$dA = \boxed{0.3}$$

$$\frac{dA}{dZ} = -\frac{\boxed{0.3}}{\boxed{}} = \boxed{1}$$

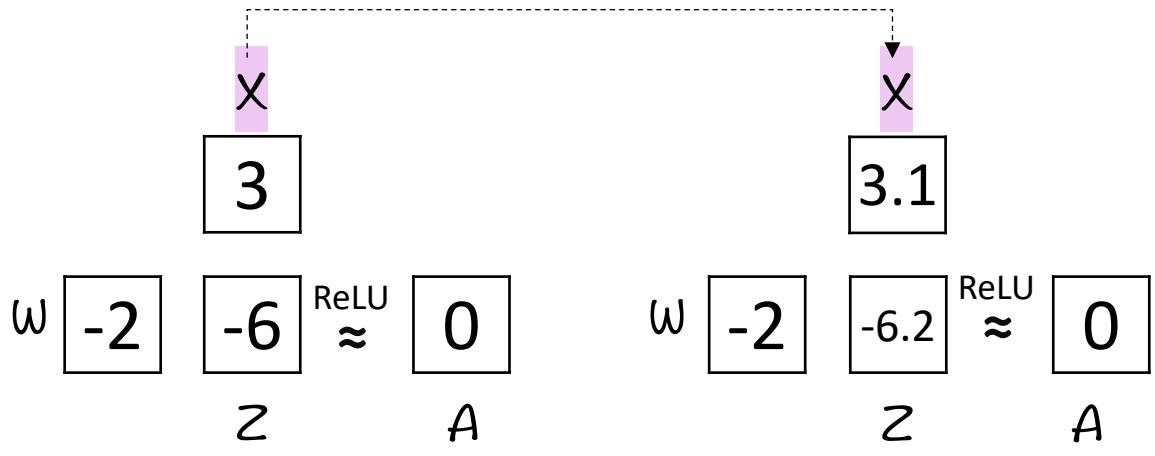
$$\frac{dA}{dX} = \frac{\boxed{0.3}}{\boxed{0.1}} = \boxed{}$$

$$\frac{dZ}{dX} = \frac{\boxed{}}{\boxed{0.1}} = \boxed{3}$$

Gradient

Exercise 25

$$dX = \boxed{0.1}$$



$$dZ = \boxed{-0.2}$$

$$dA = \boxed{\quad}$$

$$\frac{dA}{dZ} = \frac{\boxed{\quad}}{\boxed{-0.2}} = \boxed{0}$$

$$\frac{dA}{dX} = \frac{\boxed{\quad}}{\boxed{0.1}} = \boxed{\quad}$$

$$\frac{dZ}{dX} = \frac{\boxed{-0.2}}{\boxed{0.1}} = \boxed{-2}$$

Multi Layer Perceptron in pytorch

```
1 mlp_model = nn.Sequential(  
.....  
2     nn.[REDACTED]( [REDACTED, [REDACTED, bias = [REDACTED ),  
.....  
3     nn.[REDACTED](),  
.....  
4     nn.[REDACTED]( [REDACTED, [REDACTED, bias = [REDACTED ),  
.....  
5     nn.[REDACTED](),  
.....  
6     nn.[REDACTED]( [REDACTED, [REDACTED, bias = [REDACTED ),  
.....  
7     nn.[REDACTED]()  
.....  
8 )
```

The diagram illustrates the forward pass of a Multi-Layer Perceptron (MLP) with two hidden layers and one output layer. The input is a vector $\begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$ with a size of 1. It passes through a linear layer (represented by a 4×4 matrix) followed by a ReLU activation function, resulting in an intermediate vector $\begin{bmatrix} 0 \\ 3 \\ 5 \\ 3 \end{bmatrix}$. This then passes through another linear layer (represented by a 4×4 matrix) followed by a ReLU activation function, resulting in the final output vector $\begin{bmatrix} 2 \\ 1 \end{bmatrix}$ with a size of 1. Finally, a sigmoid activation function (σ) is applied to each element of the output vector, resulting in a probability vector $\begin{bmatrix} 0.95 \\ 0.50 \\ 0.12 \\ 0.99 \\ 0.01 \end{bmatrix}$.

2	1	3	1
1	-1	1	-5
1	1	0	0
0	1	1	1
1	0	1	-2
-1	3	5	3
0	3	5	3
ReLU	\approx	0	3
1	-1	1	0
0	1	-1	1
2	1	0	1
1	1	-1	1
ReLU	\approx	2	1
1	-1	2	3
-1	1	1	0
1	-2	-2	-2
2	1	0	5
-3	0	1	-5
σ	\approx	.95	.50
0	0	.12	.99
-2	5	.01	.01

Hints:

Linear Layer: { Identity | Linear | Bilinear }

Activation Function: { ReLU | Tanh | Sigmoid }

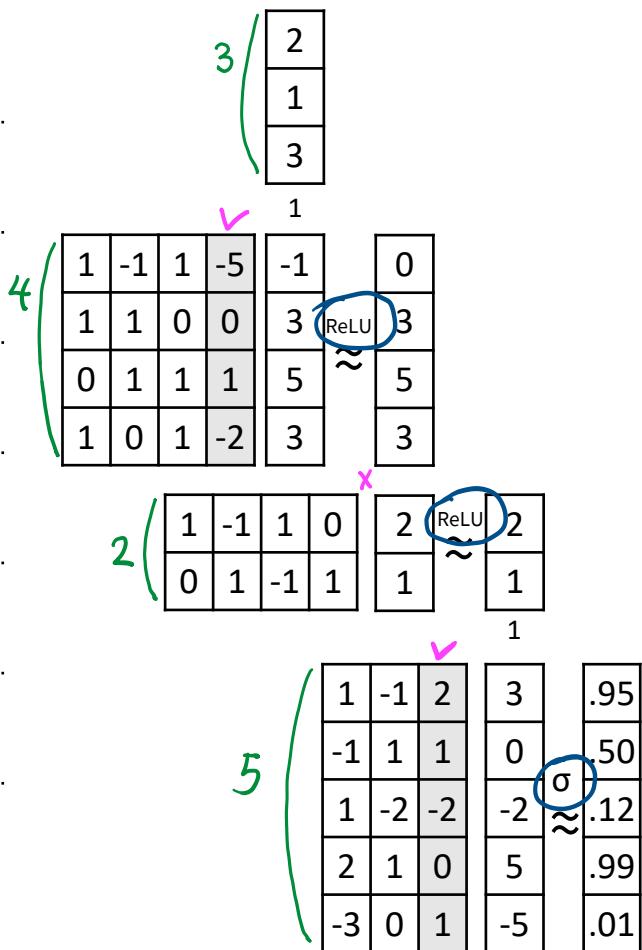
in_features: { int }

out_features: { int }

bias: { T | F }

Multi Layer Perceptron in pytorch

```
1 mlp_model = nn.Sequential(  
.....  
2     nn.Linear( 3, 4, bias = T ),  
.....  
3     nn.ReLU(),  
.....  
4     nn.Linear( 4, 2, bias = F ),  
.....  
5     nn.ReLU(),  
.....  
6     nn.Linear( 2, 5, bias = T ),  
.....  
7     nn.Sigmoid()  
.....  
8 )
```



Hints:

Linear Layer: { Identity | Linear | Bilinear }

Activation Function: { ReLU | Tanh | Sigmoid }

in_features: { int }

out_features: { int }

bias: { T | F }

Dropout

Random Sequence

.61

.39

.75

.40

.65

.42

.23

.19

.93

.42

.87

.53

.27

.69

.50

.11

.42

Training Data:

X_1	X_2
3	5
4	1

Linear

1	0	0
1	1	0
0	1	1
1	-1	0

[\approx ReLU]

1 1

Dropout
($p=0.5$)

	0	0	0
0		0	0
0	0		0
0	0	0	

1 1

Linear

1	0	0	1	0
0	1	1	0	0
1	0	-1	-1	1

[\approx ReLU]

1 1

Dropout
($p=0.33$)

	0	0
0		0
0	0	

1 1

Linear

1	-1	0	0
0	1	-1	-2

Outputs
Y

Inference

Unseen Data:

3	3
2	1

1	0	0
1	1	1
-1	1	1

[\approx ReLU]

	0	0	0
0		0	0
0	0		0
0	0	0	

1	0	1	1	0
1	1	1	0	0
1	0	-1	0	1

[\approx ReLU]

	0	0
0		0
0	0	

1	1	0	0
0	1	-1	-1

Training

-

-4	7
10	5

Targets
Y'

x 2

MSE Loss
Gradients

$$\frac{\partial L}{\partial Y}$$

Dropout

Random Sequence

.61 >.5

.39
.75
.40
.65
.42
.23

Linear

Training Data:

X_1	X_2
3	5
4	1

1 1

1	0	0
1	1	0
0	1	1
1	-1	0

[\approx ReLU]

Dropout
($p=0.5$)

$$\frac{1}{1-p} = 2$$

2	0	0	0
0	0	0	0
0	0	2	0
0	0	0	0

6	10
0	0
10	4
0	0

Linear

1	0	0	1	0
0	1	1	0	0
1	0	-1	-1	1

[\approx ReLU]

6	10
10	4
4	6
0	0

1 1

Dropout
($p=0.33$)

$$\frac{1}{1-p} = 1.5$$

1.5	0	0
0	1.5	0
0	0	0

9	15
15	6
0	0

1 1

Linear

1	-1	0	0
0	1	-1	-2

-6	9
13	4

Training



-4	7
10	5

-2	2
3	-1

x 2

-4	4
6	-2

Targets

Y'

MSE Loss
Gradients

$$\frac{\partial L}{\partial Y}$$

Inference

Unseen Data:

3	3
2	1

1 1

1	0	0
1	1	1
-1	1	1
1	-1	0

[\approx ReLU]

3	3
6	5
0	X
1	2

1 1

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

3	3
6	5
0	0
1	2

1	0	1	1	0
1	1	1	0	0
1	0	-1	0	1

[\approx ReLU]

4	5
9	8
4	4
1	1

1 1

1	0	0
0	1	0
0	0	1

4	5
9	8
4	4
1	1

1	1	0	0
0	1	-1	-1

13	13
4	3

Backpropagation

$$\begin{matrix} X \\ \boxed{2} \\ \boxed{1} \\ \boxed{3} \\ 1 \end{matrix}$$

Layer 1

$$\begin{array}{|c|c|c|c|} \hline 1 & -1 & 1 & -5 \\ \hline 1 & 1 & 0 & 0 \\ \hline 0 & 1 & 1 & 1 \\ \hline 1 & 0 & 1 & -2 \\ \hline \end{array} \quad \begin{array}{|c|} \hline -1 \\ \hline 3 \\ \hline 5 \\ \hline 3 \\ \hline 1 \end{array} \quad \text{ReLU} \approx \begin{array}{|c|} \hline 0 \\ \hline 3 \\ \hline 5 \\ \hline 3 \\ \hline 1 \end{array}$$

Layer 2

$$\begin{array}{|c|c|c|c|c|} \hline 1 & -1 & 1 & 0 & 0 \\ \hline 0 & 1 & -1 & 1 & 3 \\ \hline \end{array} \quad \begin{array}{|c|} \hline 2 \\ \hline 4 \\ \hline 1 \end{array} \quad \text{ReLU} \approx \begin{array}{|c|} \hline 2 \\ \hline 4 \\ \hline 1 \end{array}$$

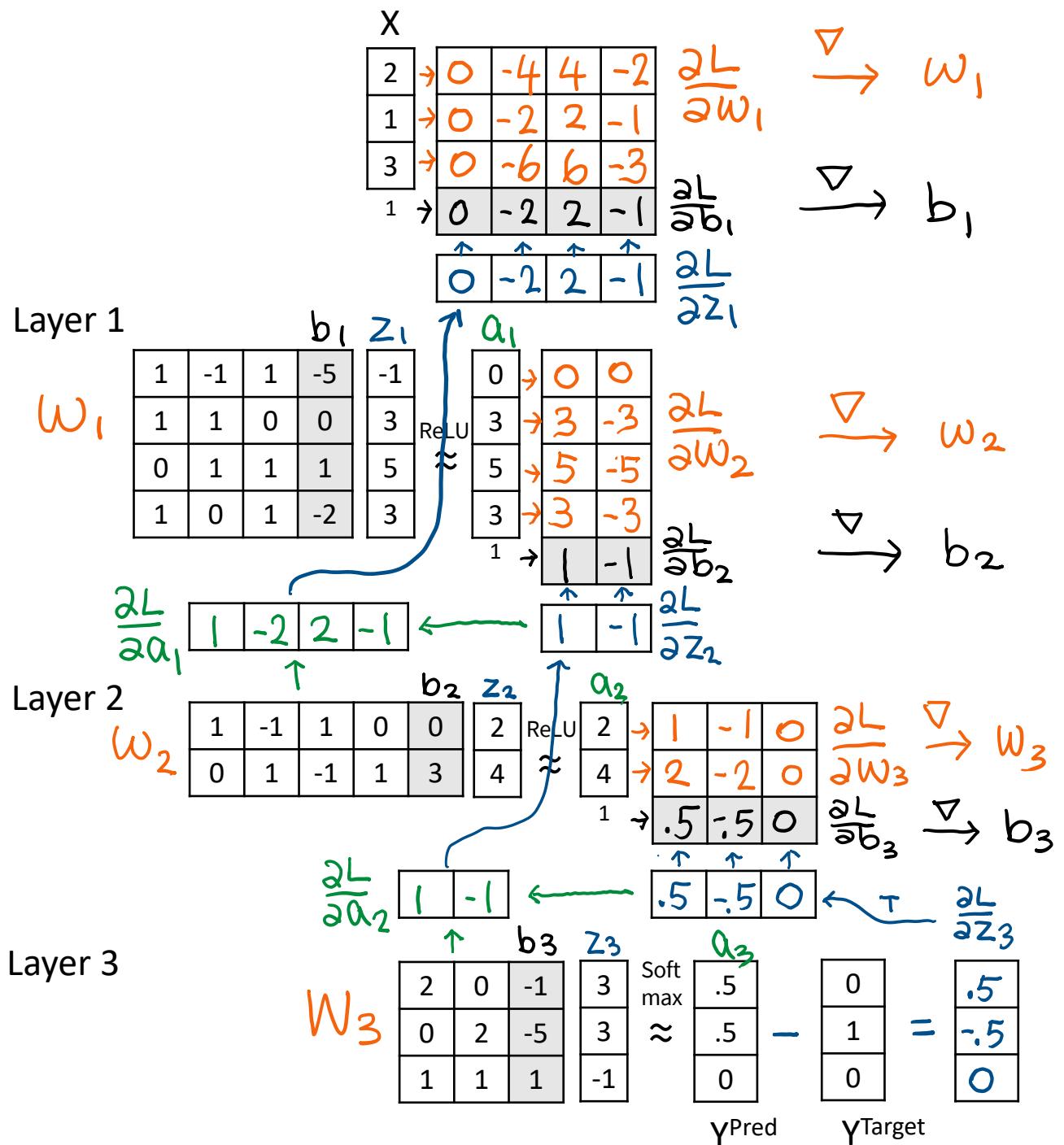
Layer 3

$$\begin{array}{|c|c|c|} \hline 2 & 0 & -1 \\ \hline 0 & 2 & -5 \\ \hline 1 & 1 & 1 \\ \hline \end{array} \quad \begin{array}{|c|} \hline 3 \\ \hline 3 \\ \hline -1 \\ \hline \end{array} \quad \text{Soft max} \approx \begin{array}{|c|} \hline .5 \\ \hline .5 \\ \hline 0 \\ \hline \end{array} \quad \begin{array}{|c|} \hline 0 \\ \hline 1 \\ \hline 0 \\ \hline \end{array}$$

γ^{Pred} γ^{Target}

L: Cross-Entropy Loss

Backpropagation



L : Cross-Entropy Loss

Batch Normalization

Mini-batch: $x_1 \ x_2 \ x_3 \ x_4$

1	0	3	0
0	3	1	1
2	1	0	2

1 1 1 1

Linear Layer

1	0	1	0
1	1	0	-1
0	2	-1	0

ReLU

\approx

Normalize

$$\mu$$

$-$

$$\sigma$$

\div

1 1 1 1

Scale & Shift

2	0	0	0
0	3	0	0
0	0	-1	1

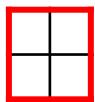
Next Layer

Batch Statistics

Σ	μ	σ^2	σ

Sum (Σ)
Mean (μ)
Variance (σ^2)
Std Dev (σ)

Trainable Parameters



Batch Normalization

Mini-batch: $x_1 \ x_2 \ x_3 \ x_4$

1	0	3	0
0	3	1	1
2	1	0	2

Linear Layer

1	0	1	0	→	3	1	3	2
1	1	0	-1	→	0	2	3	0
0	2	-1	0	→	-2	5	2	0

ReLU
≈

3	1	3	2
0	2	3	0
0	5	2	0

Batch Statistics

Σ	μ	σ^2	σ
9	2	1	1
5	1	1	1
7	2	4	2

Normalize

$$\begin{array}{r} \mu \\ - \\ \hline 2 \\ 1 \\ 2 \end{array}$$

1	-1	1	0
-1	1	2	-1
-2	3	0	-2

$$\begin{array}{r} \sigma \\ \div \\ \hline 1 \\ 1 \\ 2 \end{array}$$

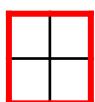
1	-1	1	0
-1	1	2	-1
-1	1	0	-1

Scale & Shift

2	0	0	0	→	2	-2	2	0
0	3	0	0	→	-3	3	6	-3
0	0	-1	1	→	2	0	1	2

Sum (Σ)
Mean (μ)
Variance (σ^2)
Std Dev (σ)

Trainable Parameters



Next Layer

Linear

y_i	x_i	x'_j	x'_1	x'_2
1	x_1		2	1
-1	x_2		4	1
1	x_3		3	1
-1	x_4			
1	x_5			
-1	x_6			

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i \quad [2 \ 0 \ 1 \ 1 \ 1 \ 1] \quad b$$

$$a_i y_i \quad [\quad \quad \quad \quad \quad \quad 2]$$

$$\text{sign} \quad [\quad \quad]$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

1	2	1
0	4	3
0	3	0
0	3	1
5	0	3
1	4	3

L2 distance $\| \cdot \|^2$

$$x_i \quad (x - x')^2 \quad \sum \sqrt{-\gamma \| \cdot \|^2_{(-0.1)}} \quad e^{\square}$$

x'_j	2
4	
3	

$-\gamma \ \cdot \ ^2$	e^{\square}
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

$$a_i \quad [20 \ 10 \ 10 \ 10 \ 0 \ 10] \quad b$$

$$a_i y_i \quad [\quad \quad \quad \quad \quad \quad -1]$$

$$\text{sign} \quad [\quad \quad]$$

Linear

x'_j	x'_1	x'_2
	2	1
	4	1
	3	1

SVM

y_i	x_i	x'_j	Kernel Matrix
1	x_1	1 2 1	
-1	x_2	0 4 3	
1	x_3	0 3 0	
-1	x_4	0 3 1	
1	x_5	5 0 3	
-1	x_6	1 4 3	

$$K(x_i, x'_j)$$

a_i	2 0 1 1 1 1	b
$a_i y_i$		2
		sign

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x_i	$(x - x')^2$	$\sum \sqrt{-\gamma x - x' ^2}$	x'_j	e^\square
1 2 1			2	
0 4 3			4	
0 3 0			3	
0 3 1				
5 0 3				
1 4 3				

$-\gamma x - x' ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1

a_i	20 10 10 10 0 10	b
$a_i y_i$		

$a_i y_i$		sign

Linear

x'_j	x'_1	x'_2
	2	1
	4	1
	3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13
-1	x_2 0 4 3	25
1	x_3 0 3 0	12
-1	x_4 0 3 1	15
1	x_5 5 0 3	19
-1	x_6 1 4 3	27

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i \quad [2 \ 0 \ 1 \ 1 \ 1 \ 1] \quad b$$

$$a_i y_i \quad [\quad \quad \quad \quad \quad \quad 2]$$

$$\text{sign} \quad [\quad \quad]$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

L2 distance $\| \cdot \|^2$

$$X_i \quad (x - x')^2 \quad \sum \sqrt{-\gamma \| \cdot \|^2_{(-0.1)}} \quad e^{\square}$$

x'_j	2	4	3
	2	4	3
	4	16	9
	3	9	1
	1	1	1

$-\gamma \ \cdot \ ^2$	e^{\square}
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

1	2	1
0	4	3
0	3	0
0	3	1
5	0	3
1	4	3

$$a_i \quad [20 \ 10 \ 10 \ 10 \ 0 \ 10] \quad b$$

$$a_i y_i \quad [\quad \quad \quad \quad \quad \quad -1]$$

$$\text{sign} \quad [\quad \quad]$$

3

Linear

x'_j	x'_1	x'_2
	2	1
	4	1
	3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13
-1	x_2 0 4 3	25
1	x_3 0 3 0	12
-1	x_4 0 3 1	15
1	x_5 5 0 3	19
-1	x_6 1 4 3	27

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i [2 \ 0 \ 1 \ 1 \ 1 \ 1] b$$

$$a_i y_i [2 \ 0 \ 1 \ -1 \ 1 \ -1 \ 2]$$

$$\text{sign} [\quad \quad]$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

L2 distance $\| \cdot \|^2$

$$x_i (x - x')^2 \sum \sqrt{-\gamma \| \cdot \|^2_{(-0.1)}} e^\gamma$$

x'_j	2	4	3
	2	4	3
	3	1	2

$-\gamma \ \cdot\ ^2$	e^γ
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

1	2	1
0	4	3
0	3	0
0	3	1
5	0	3
1	4	3

$$a_i [20 \ 10 \ 10 \ 10 \ 0 \ 10] b$$

$$a_i y_i [\quad \quad \quad \quad \quad \quad -1 \quad \quad]$$

$$\text{sign} [\quad \quad]$$

Linear

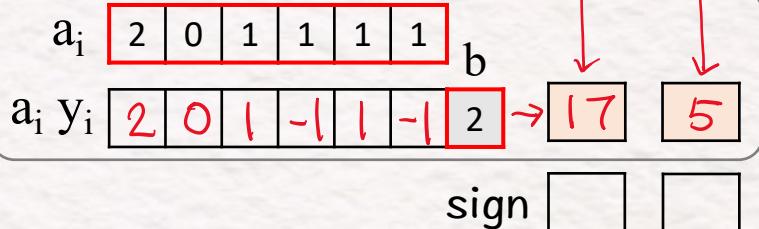
x'_j	x'_1	x'_2
	2	1
	4	1
	3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13
-1	x_2 0 4 3	25
1	x_3 0 3 0	12
-1	x_4 0 3 1	15
1	x_5 5 0 3	19
-1	x_6 1 4 3	27

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

$$b + \sum a_i y_i K$$

RBF

L2 distance $||^2$

$$x_i \quad (x - x')^2 \quad \sum \sqrt{-\gamma ||^2_{(-0.1)}} e^\square$$

x'_j	2	4	3
	2	4	3
	4	3	2
	3	2	4

$-\gamma ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

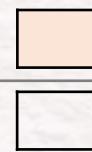
1
1
1

1	2	1
0	4	3
0	3	0
0	3	1
5	0	3
1	4	3

$$a_i [20 10 10 10 0 10] b$$

$$a_i y_i [] [] [-1] []$$

$$\text{sign} []$$



Linear

SVM

y_i	x_i	x'_j	x'_1	x'_2
1	x_1		2	1
-1	x_2		4	1
1	x_3		3	1
-1	x_4			
1	x_5			
-1	x_6			

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i [2 \ 0 \ 1 \ 1 \ 1 \ 1] b$$

$$a_i y_i [2 \ 0 \ 1 \ -1 \ 1 \ -1 \ 2] \rightarrow [17 \ 5]$$

$$\text{sign} \quad + \quad +$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

L2 distance $\| \cdot \|^2$

$$x_i \quad (x - x')^2 \quad \sum \sqrt{-\gamma \| \cdot \|^2_{(-0.1)}} e^\square$$

$$x'_j \quad \begin{matrix} 2 \\ 4 \\ 3 \end{matrix}$$

$-\gamma \ \cdot \ ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

$$\begin{matrix} 1 \\ 1 \\ 1 \end{matrix}$$

$$\begin{matrix} 1 & 2 & 1 \\ 0 & 4 & 3 \\ 0 & 3 & 0 \\ 0 & 3 & 1 \\ 5 & 0 & 3 \\ 1 & 4 & 3 \end{matrix}$$

$$\begin{matrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{matrix}$$

$$\begin{matrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{matrix}$$

$$\begin{matrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{matrix}$$

$$\begin{matrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{matrix}$$

$$\begin{matrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{matrix}$$

$$\begin{matrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{matrix}$$

$$a_i [20 \ 10 \ 10 \ 10 \ 0 \ 10] b$$

$$a_i y_i [\quad \quad \quad \quad \quad \quad -1 \quad \quad]$$

$$\text{sign} \quad \boxed{\quad}$$

$$\boxed{\quad}$$

Linear

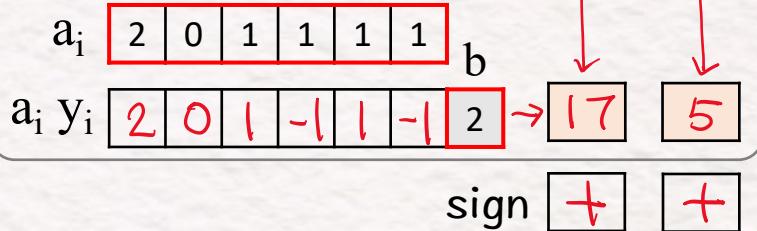
x'_1	x'_2
2	1
4	1
3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13
-1	x_2 0 4 3	25
1	x_3 0 3 0	12
-1	x_4 0 3 1	15
1	x_5 5 0 3	19
-1	x_6 1 4 3	27

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x'_j	2	4	3
x'_j	2	4	3
$(x - x')^2$	$\sum \sqrt{-\gamma x - x' ^2}$	$e^{-\gamma x - x' ^2}$	
x_i	1 2 1	1 4 4	4 0 0
	0 4 3	4 0 0	0 1 0
	0 3 0	0 1 0	4 1 4
	0 3 1	4 1 4	9 1 6 0
	5 0 3	9 1 6 0	1 0 0
	1 4 3	1 0 0	

$-\gamma x ^2$	$e^{-\gamma x ^2}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

$$a_i [20 10 10 10 0 10] b$$

$$a_i y_i [] [] [] [] [-1] []$$

$$\text{sign } []$$

Linear

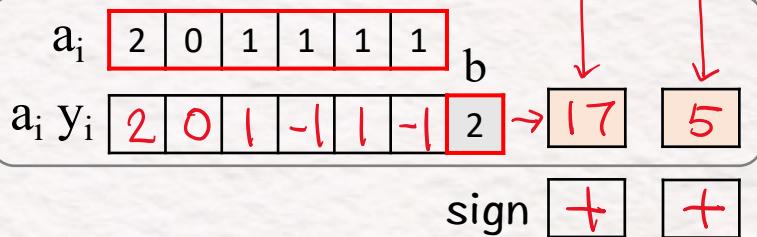
x'_1	x'_2
2	1
4	1
3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13 2
-1	x_2 0 4 3	25 7
1	x_3 0 3 0	12 3
-1	x_4 0 3 1	15 4
1	x_5 5 0 3	19 8
-1	x_6 1 4 3	27 8

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x'_j	2	4	3
x'_j	2	4	3
$L_2 \text{ distance } x - x' ^2$	9	4	1
$(x - x')^2$	4	1	9
$\sum \sqrt{-\gamma x - x' ^2}$	1	25	1
$e^{-\gamma x - x' ^2}$	1	1	1

$-\gamma x - x' ^2$	$e^{-\gamma x - x' ^2}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

$$a_i [20 10 10 10 0 10] b$$

$$a_i y_i [] [] [] [] [-1] []$$

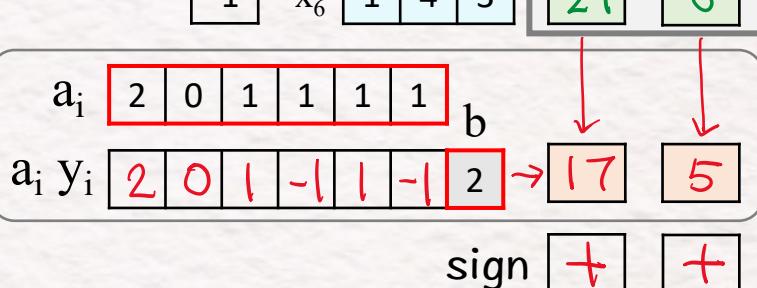
$$\text{sign} []$$

Linear

y_i	x_i	x'_j	x'_1	x'_2
1	x_1		2	1
-1	x_2		4	1
1	x_3		3	1
-1	x_4			
1	x_5			
-1	x_6			

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

$$b + \sum a_i y_i K$$

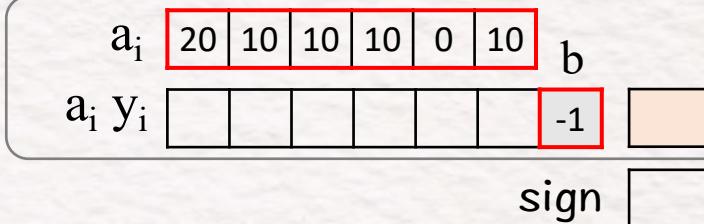
RBF

x_i	x'_j	x'_1	x'_2	x'_3
1 2 1		2		
0 4 3		4		
0 3 0		3		
0 3 1				
5 0 3				
1 4 3				

$L_2 \text{ distance } ||x - x'||^2$
 $(x - x')^2 \sum \sqrt{-\gamma ||(-0.1)||^2} e^{-\gamma ||x - x'||^2}$

$-\gamma x ^2$	$e^{-\gamma}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1



Linear

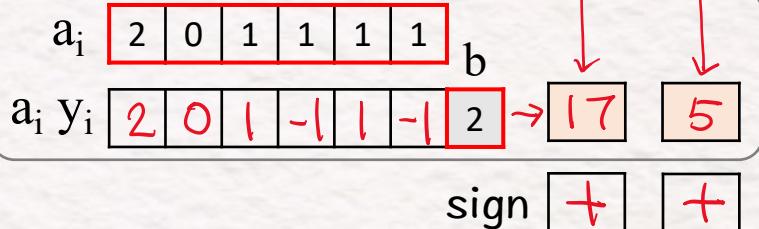
x'_1	x'_2
2	1
4	1
3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13 2
-1	x_2 0 4 3	25 7
1	x_3 0 3 0	12 3
-1	x_4 0 3 1	15 4
1	x_5 5 0 3	19 8
-1	x_6 1 4 3	27 8

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

RBF

x'_j	2	4	3
x'_j	2	4	3
$L_2 \text{ distance } x - x' ^2$	9	3	-0.3
$(x - x')^2$	4	2	-0.2
$\sum \sqrt{-\gamma x - x' ^2}$	1	1	-0.1
e^\square	9	3	-0.3
	25	5	-0.5
	1	1	-0.1

$-\gamma x ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1



Linear

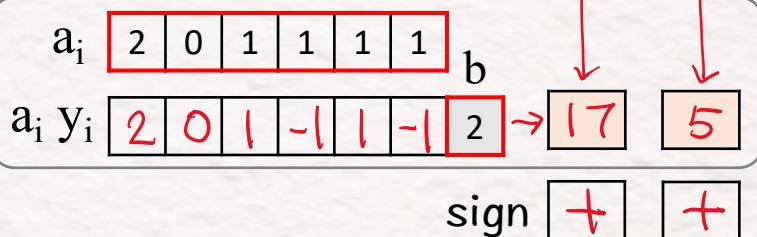
x'_1	x'_2
2	1
4	1
3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13
-1	x_2 0 4 3	25
1	x_3 0 3 0	12
-1	x_4 0 3 1	15
1	x_5 5 0 3	19
-1	x_6 1 4 3	27

Kernel Matrix

$$K(x_i, x'_j)$$



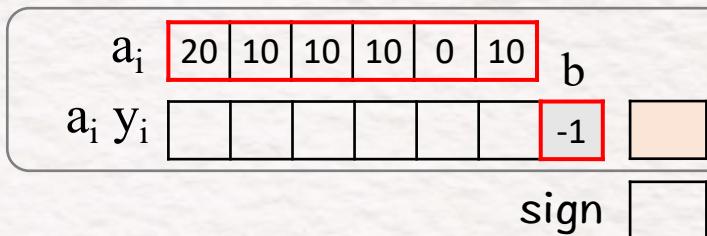
Decision Boundary

RBF

x'_j	2	4	3
x'_j	2	4	3
$(x - x')^2$	$\sum \sqrt{-\gamma x - x' ^2}$	$e^{-\gamma x - x' ^2}$	
1 2 1	9 3 -0.3	0.7	
0 4 3	4 2 -0.2	0.8	
0 3 0	1 1 -0.1	0.9	
0 3 1	9 3 -0.3	0.7	
5 0 3	25 5 -0.5	0.6	
1 4 3	1 1 -0.1	0.9	

$-\gamma x ^2$	$e^{-\gamma x ^2}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1



Linear

x'_1	x'_2
2	1
4	1
3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13
-1	x_2 0 4 3	25
1	x_3 0 3 0	12
-1	x_4 0 3 1	15
1	x_5 5 0 3	19
-1	x_6 1 4 3	27

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i \quad [2 \ 0 \ 1 \ 1 \ 1 \ 1] \quad b$$

$$a_i y_i \quad [2 \ 0 \ 1 \ -1 \ 1 \ -1 \ 2] \rightarrow [17 \ 5]$$

$$\text{sign} \quad [+] \quad [+]$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x'_j	2
4	4
3	3
0	0

$-\gamma x'_j ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

$$a_i \quad [20 \ 10 \ 10 \ 10 \ 0 \ 10] \quad b$$

$$a_i y_i \quad [\quad \quad \quad \quad \quad -1 \quad]$$

$$\text{sign} \quad [\quad]$$

1
1
1

Linear

y_i	x_i	x'_j	x'_1	x'_2
1	x_1	1 2 1	2	1
-1	x_2	0 4 3	4	1
1	x_3	0 3 0	3	1
-1	x_4	0 3 1	13	2
1	x_5	5 0 3	25	7
-1	x_6	1 4 3	12	3

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i \quad [2 \ 0 \ 1 \ 1 \ 1 \ 1] \quad b$$

$$a_i y_i \quad [2 \ 0 \ 1 \ -1 \ 1 \ -1 \ 2] \rightarrow [17 \ 5]$$

$$\text{sign} \quad [+] \quad [+]$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x_i	x'_j	2	4	3
1 2 1				
0 4 3				
0 3 0				
0 3 1				
5 0 3				
1 4 3				

L2 distance $\| \cdot \|^2$

$$(x - x')^2$$

$$\sum \sqrt{-\gamma \| \cdot \|^2}$$

$-\gamma \ \cdot \ ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

$$a_i \quad [20 \ 10 \ 10 \ 10 \ 0 \ 10] \quad b$$

$$a_i y_i \quad [\quad \quad \quad \quad \quad -1 \quad]$$

$$\text{sign} \quad [\quad]$$

Linear

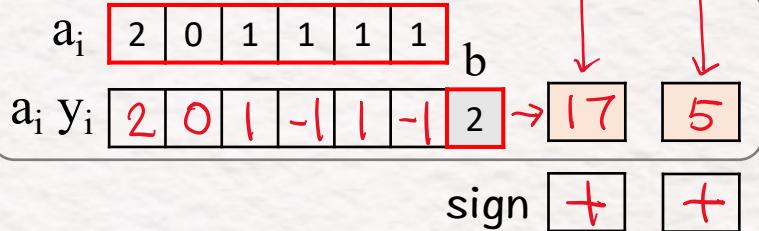
x'_j	x'_1	x'_2
	2	1
	4	1
	3	1

SVM

y_i	x_i	x'_j
1	x_1 1 2 1	13 2
-1	x_2 0 4 3	25 7
1	x_3 0 3 0	12 3
-1	x_4 0 3 1	15 4
1	x_5 5 0 3	19 8
-1	x_6 1 4 3	27 8

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x'_j	2	4	3
	2	4	3
	4	2	-0.3
	1	1	-0.2
	9	3	-0.1
	4	2	0.7
	1	1	0.8
	25	5	0.9
	1	1	0.7
	9	3	0.6
	25	5	0.9

$-\gamma x ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1
1

0	1	0	1	1	
1	9	4	14	3	
1	4	1	6	2	
1	4	0	5	2	
16	1	4	21	4	
0	9	4	13	3	

$$a_i [20 10 10 10 0 10] b$$

$$a_i y_i [] [] [] [] [-1] []$$

$$\text{sign } []$$

Linear

y_i	x_i	x'_j	x'_1	x'_2
1	x_1	1 2 1	2	1
-1	x_2	0 4 3	4	1
1	x_3	0 3 0	3	1
-1	x_4	0 3 1	13	2
1	x_5	5 0 3	25	7
-1	x_6	1 4 3	12	3

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i [2 \ 0 \ 1 \ 1 \ 1 \ 1] b$$

$$a_i y_i [2 \ 0 \ 1 \ -1 \ 1 \ -1 \ 2] \rightarrow [17 \ 5]$$

$$\text{sign} \quad + \quad +$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x_i	x'_j	2	4	3
1 2 1				
0 4 3				
0 3 0				
0 3 1				
5 0 3				
1 4 3				

L2 distance $\| \cdot \|^2$

$(x - x')^2 \sum \sqrt{-\gamma \| \cdot \|^2_{(-0.1)}} e^{\frac{-}{}}$

$-\gamma \ \cdot \ ^2$	$e^{\frac{-}{}}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1

a_i	20	10	10	10	0	10	b
$a_i y_i$						-1	

sign

0	1	0	1	1	-0.1
1	9	4	14	3	-0.3
1	4	1	6	2	-0.2
1	4	0	5	2	-0.2
16	1	4	21	4	-0.4
0	9	4	13	3	-0.3

Linear

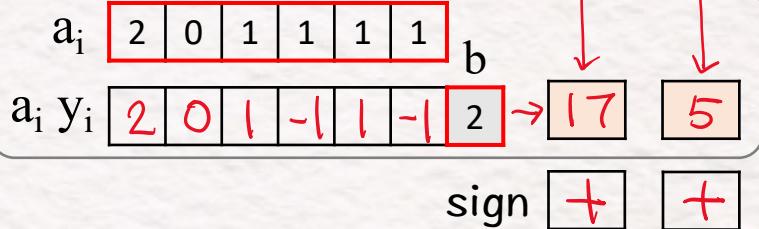
x'_1	x'_2
2	1
4	1
3	1

SVM

y_i	X_i	x'_j
1	$x_1 \begin{matrix} 1 & 2 & 1 \end{matrix}$	13
-1	$x_2 \begin{matrix} 0 & 4 & 3 \end{matrix}$	25
1	$x_3 \begin{matrix} 0 & 3 & 0 \end{matrix}$	12
-1	$x_4 \begin{matrix} 0 & 3 & 1 \end{matrix}$	15
1	$x_5 \begin{matrix} 5 & 0 & 3 \end{matrix}$	19
-1	$x_6 \begin{matrix} 1 & 4 & 3 \end{matrix}$	27

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

$$b + \sum a_i y_i K$$

RBF

X_i	x'_j	$ x - x' ^2$	$e^{-\gamma x - x' ^2}$
$\begin{matrix} 1 & 2 & 1 \end{matrix}$	$\begin{matrix} 2 \\ 4 \\ 3 \end{matrix}$	$(1-2)^2 + (2-4)^2 + (1-3)^2 = 10$	0.7
$\begin{matrix} 0 & 4 & 3 \end{matrix}$	$\begin{matrix} 4 \\ 0 \\ 0 \end{matrix}$	$(0-4)^2 + (4-0)^2 + (3-0)^2 = 41$	0.8
$\begin{matrix} 0 & 3 & 0 \end{matrix}$	$\begin{matrix} 4 \\ 1 \\ 0 \end{matrix}$	$(0-4)^2 + (3-1)^2 + (0-0)^2 = 20$	0.9
$\begin{matrix} 0 & 3 & 1 \end{matrix}$	$\begin{matrix} 4 \\ 1 \\ 4 \end{matrix}$	$(0-4)^2 + (3-1)^2 + (1-4)^2 = 26$	0.7
$\begin{matrix} 5 & 0 & 3 \end{matrix}$	$\begin{matrix} 9 \\ 16 \\ 0 \end{matrix}$	$(5-9)^2 + (0-16)^2 + (3-0)^2 = 226$	0.6
$\begin{matrix} 1 & 4 & 3 \end{matrix}$	$\begin{matrix} 1 \\ 0 \\ 0 \end{matrix}$	$(1-1)^2 + (4-0)^2 + (3-0)^2 = 25$	0.9

$-\gamma x - x' ^2$	$e^{-\gamma x - x' ^2}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1

0	1	0	1	1	-0.1
1	9	4	14	3	-0.3
1	4	1	6	2	-0.2
1	4	0	5	2	-0.2
16	1	4	21	4	-0.4
0	9	4	13	3	-0.3

$$a_i \begin{bmatrix} 20 & 10 & 10 & 10 & 0 & 10 \end{bmatrix} \quad b$$

$$a_i y_i \quad \begin{bmatrix} & & & & -1 & \end{bmatrix} \quad \begin{bmatrix} & & & & & \end{bmatrix}$$

$$\text{sign} \quad \begin{bmatrix} & \end{bmatrix}$$

Linear

SVM

y_i	x_i	x'_j	x'_1	x'_2
1	x_1		2	1
-1	x_2		4	1
1	x_3		3	1
-1	x_4			
1	x_5			
-1	x_6			

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i [2 \ 0 \ 1 \ 1 \ 1 \ 1] b$$

$$a_i y_i [2 \ 0 \ 1 \ -1 \ 1 \ -1 \ 2] \rightarrow [17 \ 5]$$

sign $+$ $+$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x_i	x'_j	2	4	3
1 2 1				
0 4 3				
0 3 0				
0 3 1				
5 0 3				
1 4 3				

L2 distance $\| \cdot \|^2$

$$(x - x')^2$$

$$\sum \sqrt{-\gamma \| \cdot \|^2}$$

$$e^{-\gamma \| \cdot \|^2}$$

$-\gamma \ \cdot \ ^2$	$e^{-\gamma \ \cdot \ ^2}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1

$$a_i [20 \ 10 \ 10 \ 10 \ 0 \ 10] b$$

$$a_i y_i [20 \ -10 \ 10 \ -10 \ 0 \ -10 \ -1] \rightarrow []$$

sign $[]$

0.9
0.7
0.8
0.8
0.67
0.9

Linear

SVM

y_i	x_i	x'_j	x'_1	x'_2
1	x_1		2	1
-1	x_2		4	1
1	x_3		3	1
-1	x_4			
1	x_5			
-1	x_6			

Kernel Matrix

$$K(x_i, x'_j)$$

$$a_i [2 \ 0 \ 1 \ 1 \ 1 \ 1] b$$

$$a_i y_i [2 \ 0 \ 1 \ -1 \ 1 \ -1 \ 2] \rightarrow [17 \ 5]$$

$$\text{sign} \quad + \quad +$$

Decision Boundary

$$b + \sum a_i y_i K$$

RBF

x_i	x'_j	x'_1	x'_2	x'_3
1 2 1		2		
0 4 3		4		
0 3 0		3		
0 3 1				
5 0 3				
1 4 3				

$-\gamma x ^2$	e^\square
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1

$$a_i [20 \ 10 \ 10 \ 10 \ 0 \ 10] b$$

$$a_i y_i [20 \ -10 \ 10 \ -10 \ 0 \ -10 \ -1] \rightarrow [-2]$$

$$\text{sign} \quad \boxed{}$$

0	1	0	1	1	-0.1	0.9
1	9	4	14	3	-0.3	0.7
1	4	1	6	2	-0.2	0.8
1	4	0	5	2	-0.2	0.8
16	1	4	21	4	-0.4	0.67
0	9	4	13	3	-0.3	0.9

1

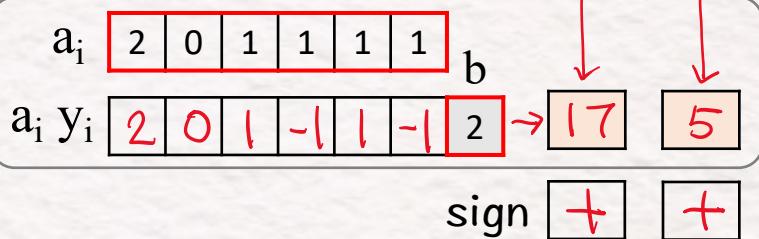
Linear

SVM

y_i	x_i	x'_j
1	x_1	x'_1
-1	x_2	x'_2
1	x_3	
-1	x_4	
1	x_5	
-1	x_6	

Kernel Matrix

$$K(x_i, x'_j)$$



Decision Boundary

$$b + \sum a_i y_i K$$

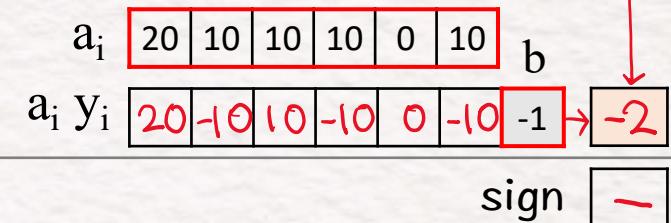
sign $+$ $+$

RBF

x_i	x'_j
$L_2 \text{ distance } x - x' ^2$	
$(x - x')^2$	
$\sum \sqrt{-\gamma x - x' ^2}$	
	$e^{-\gamma x - x' ^2}$

$-\gamma x - x' ^2$	$e^{-\gamma x - x' ^2}$
-0.5	0.6
-0.4	0.67
-0.3	0.7
-0.2	0.8
-0.1	0.9

1
1



0	1	0	1	1	-0.1	0.9
1	9	4	14	3	-0.3	0.7
1	4	1	6	2	-0.2	0.8
1	4	0	5	2	-0.2	0.8
16	1	4	21	4	-0.4	0.67
0	9	4	13	3	-0.3	0.9

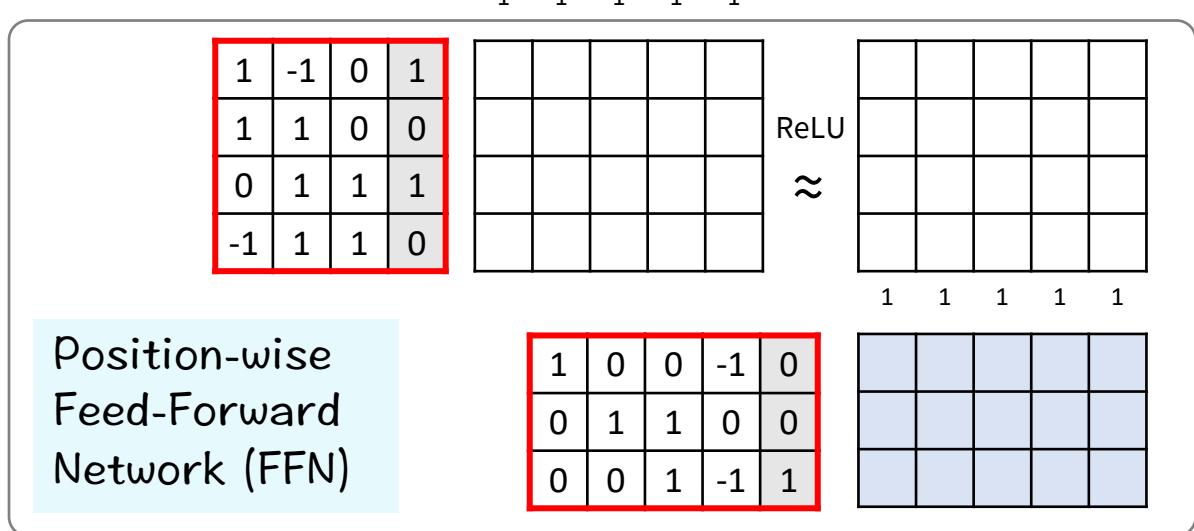
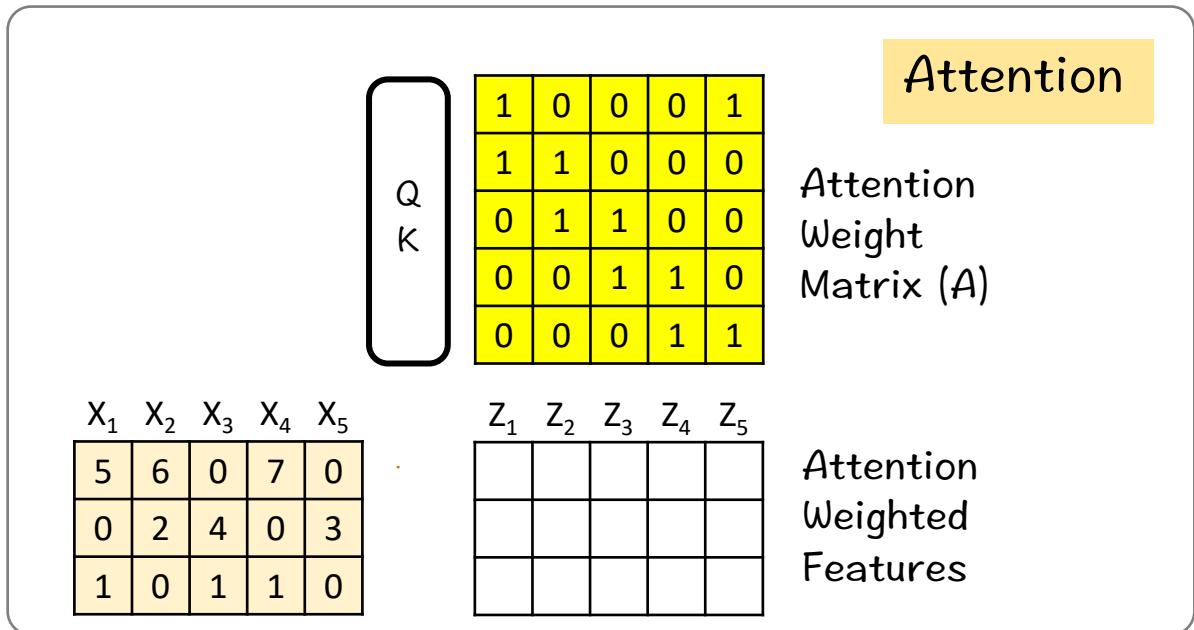
1
1

1

+

Transformer

Features from the
Previous Block

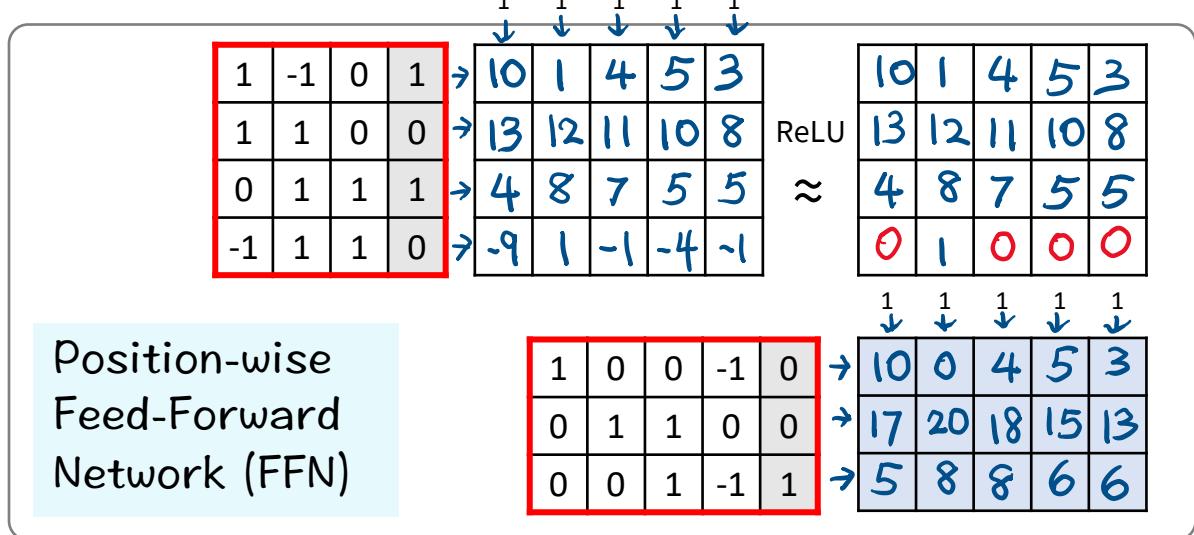
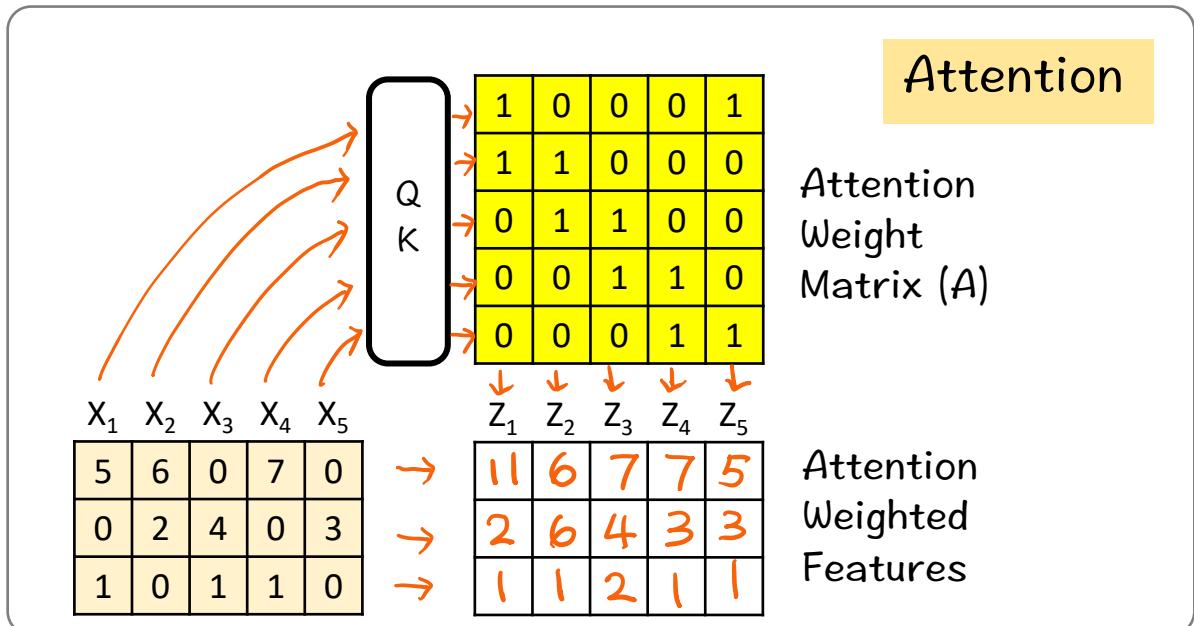


Next Block

Transformer

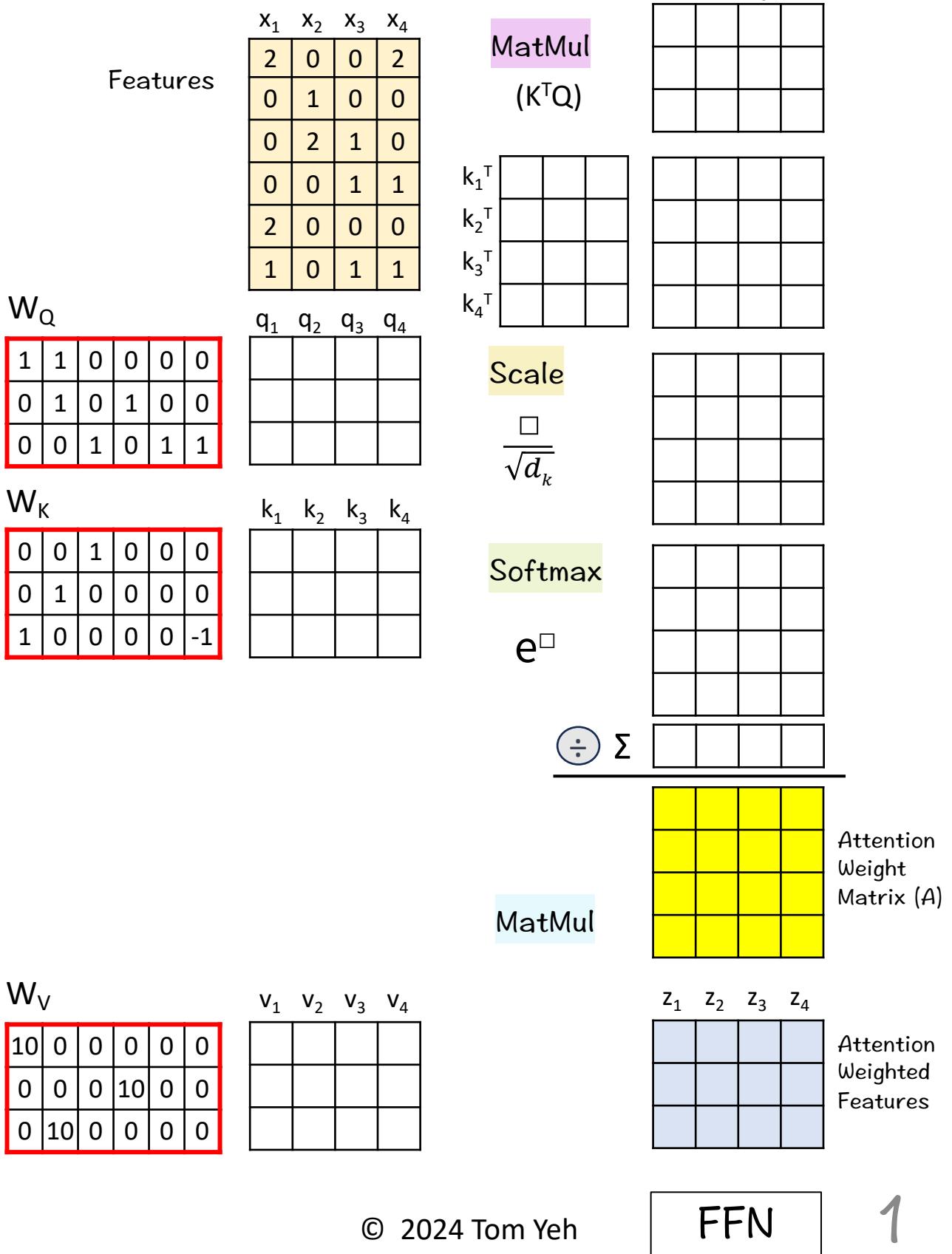
Features from the
Previous Block

↓ ↓ ↓ ↓ ↓

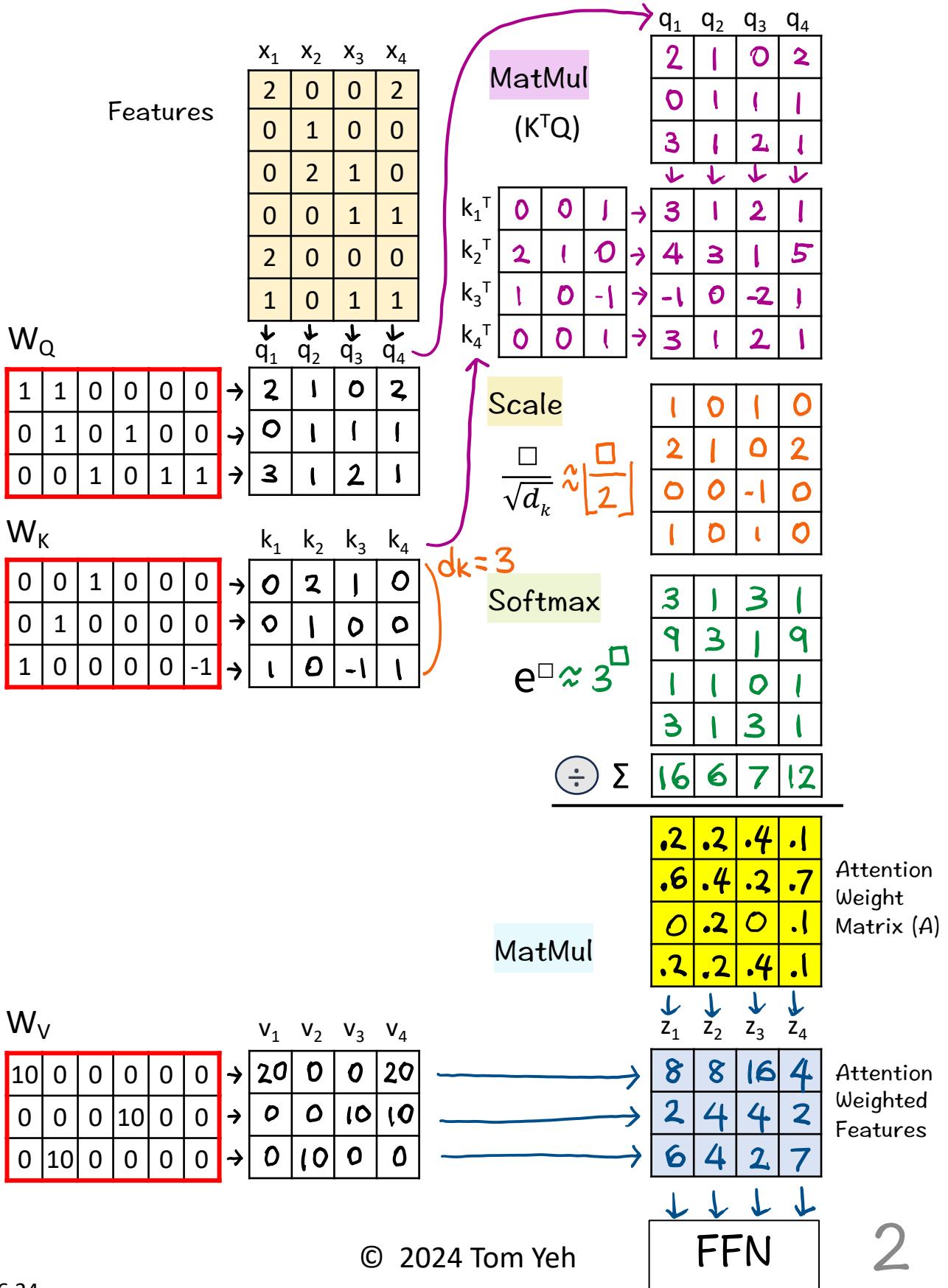


Next Block

Self Attention



Self Attention

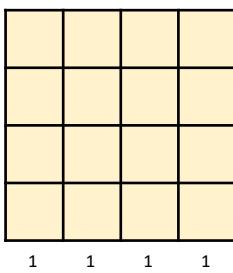


SORA's Diffusion Transformer

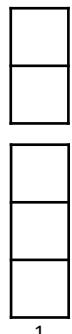
Training Video

$\begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix}$	$\begin{matrix} 2 & 0 \\ 1 & 0 \end{matrix}$	$\begin{matrix} 0 & 1 \\ 3 & 0 \end{matrix}$	$\begin{matrix} 0 & 1 \\ 4 & 0 \end{matrix}$
--	--	--	--

Spacetime
Patches
(Pixels)



Diffusion
Step $t = 3$



Prompt
“sora is sky”

Text
Encoder

0
1
-1

Visual Encoder

1	0	-1	0	0
0	1	0	1	1

[ReLU]

1	0	-1	0	0
0	1	0	1	1

Sampled
Noise



0	2	1	-1
-1	0	-2	1

Noised
Latent

1	0	-1	0
0	1	0	1

Predicted
Noise



1	0	-1	0
0	1	0	1

Noise-free
Latent

1	0	-1	0
0	1	0	1

Visual Decoder

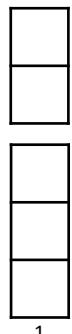
1	0	1
0	1	0
1	1	0
-1	1	0

[ReLU]

1	0	-1	0
0	1	0	1

Diffusion
Step $t = 3$

1	0	0	1	0	0
0	1	0	-1	1	0
1	1	0	0	1	0
0	2	0	1	0	2



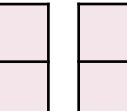
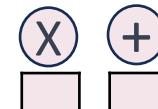
Self-Attention

1	0	0	0
1	1	0	0

1	1	0	0
0	1	1	0

1	1	1	1
---	---	---	---

Adaptive
Layer Norm



Pointwise
FFN

-1	1	-2
0	1	-5

1	1	1	1
---	---	---	---

Train

Sampled
Noise

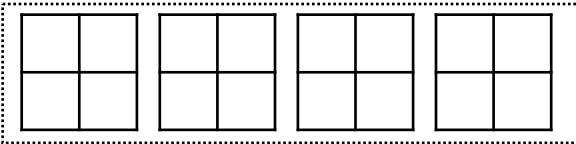


0	2	1	-1
-1	0	-2	1

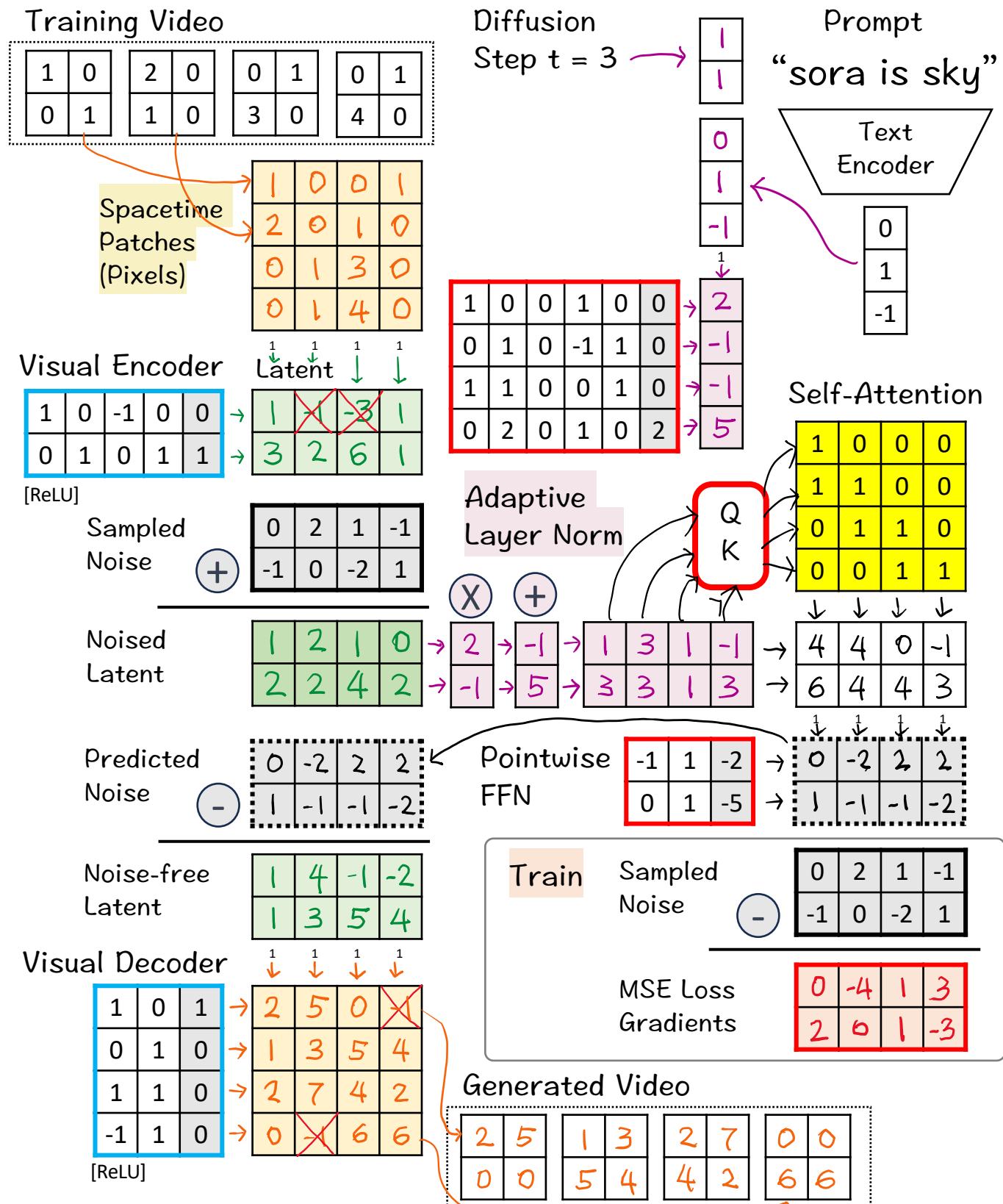
MSE Loss
Gradients

1	1	1	1
---	---	---	---

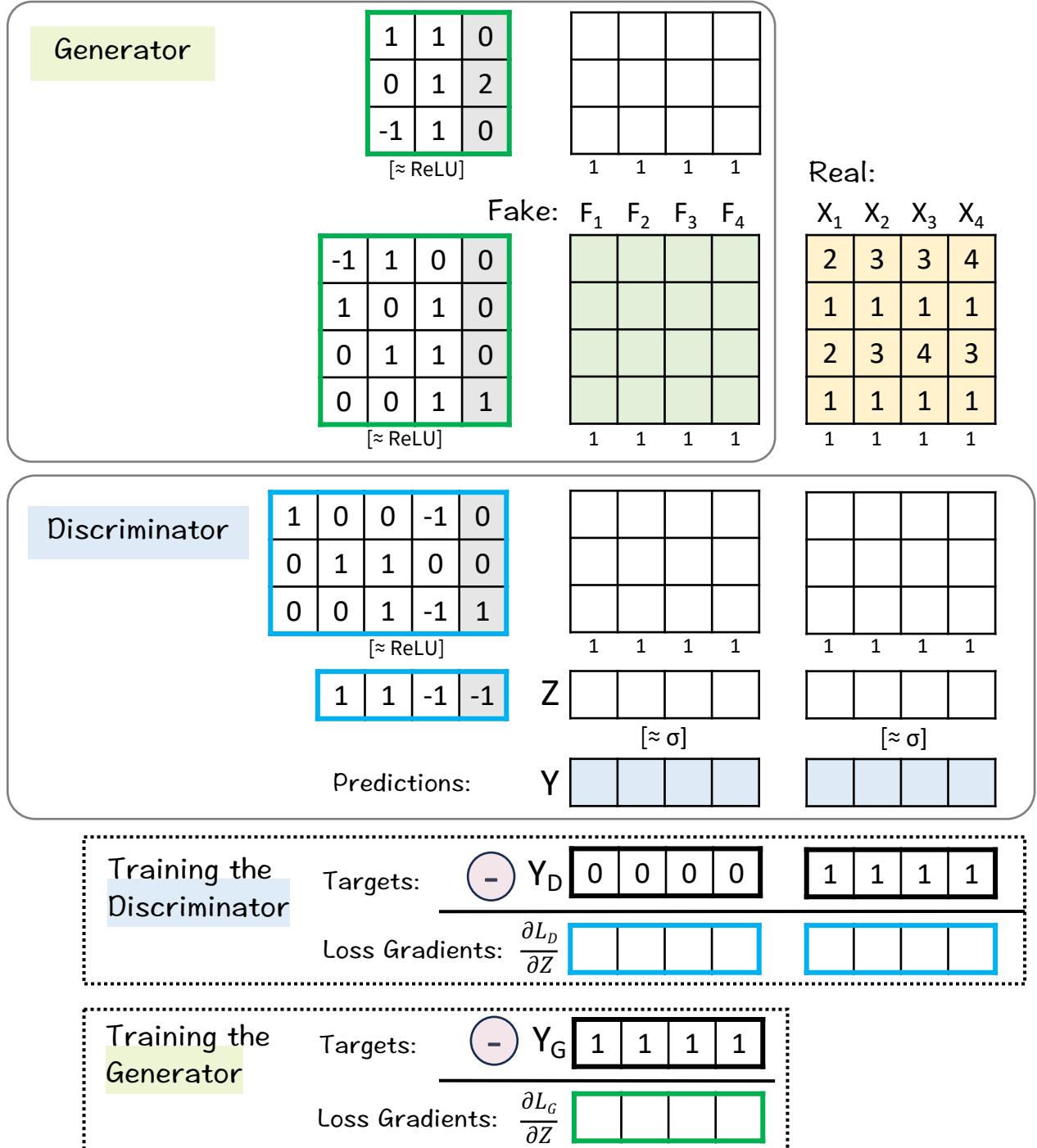
Generated Video



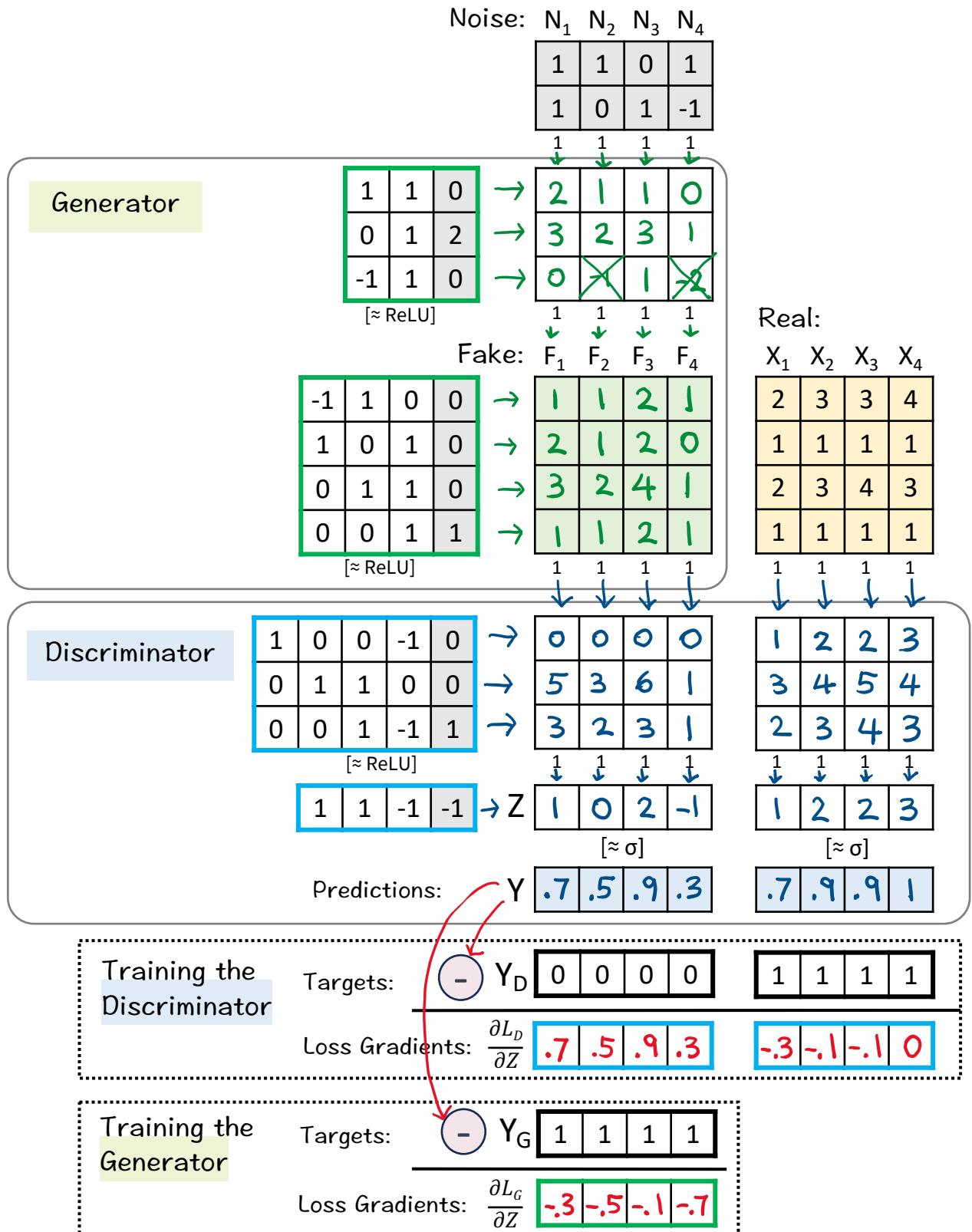
SORA's Diffusion Transformer



Generative Adversarial Network (GAN)

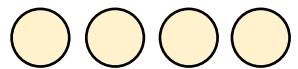


Generative Adversarial Network (GAN)



Autoencoder

X_1	X_2	X_3	X_4
1	1	2	1
2	1	2	0
3	2	4	1
1	1	2	1
1	1	1	1



Encoder

1	0	0	1	0
0	1	1	0	0
-1	0	1	0	-1

[$\approx \text{ReLU}$]

1 1 1 1

1	0	1	0
-1	1	0	0

[$\approx \text{ReLU}$]

1 1 1 1

Decoder

1	0	0
0	1	1
1	-1	0

[$\approx \text{ReLU}$]

1 1 1 1

1	0	-1	0
1	-1	0	0
0	0	1	1
0	1	1	-3

1 1 1 1

Reconstruction
Loss



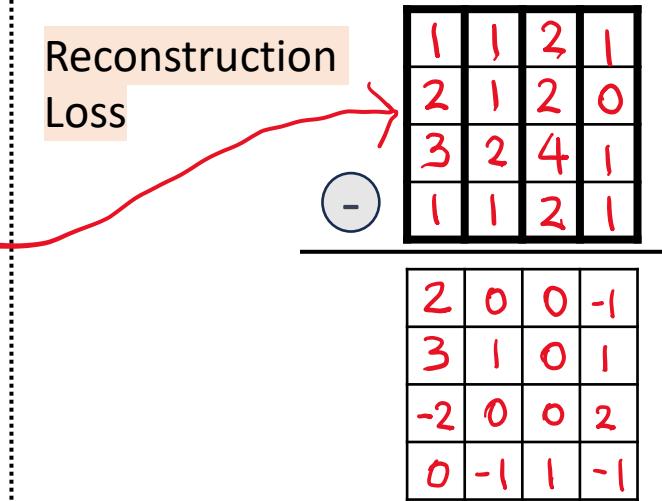
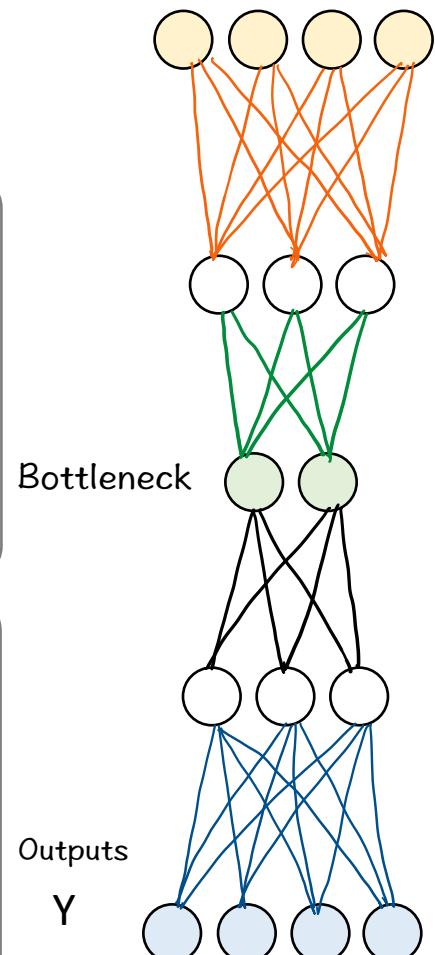
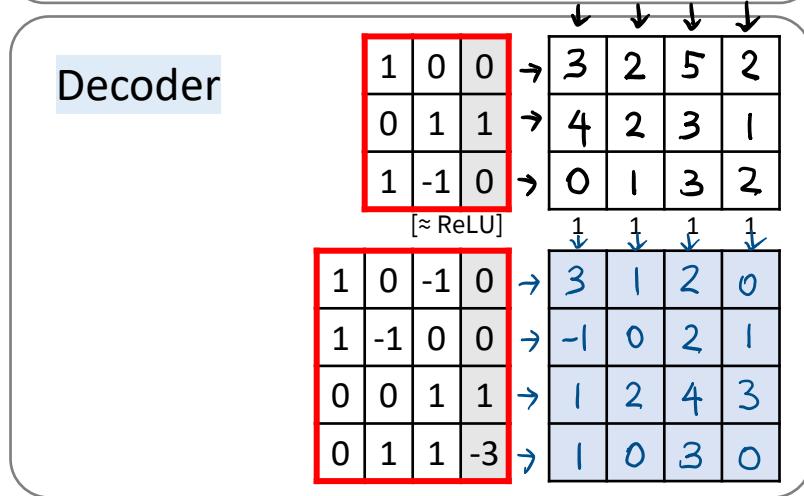
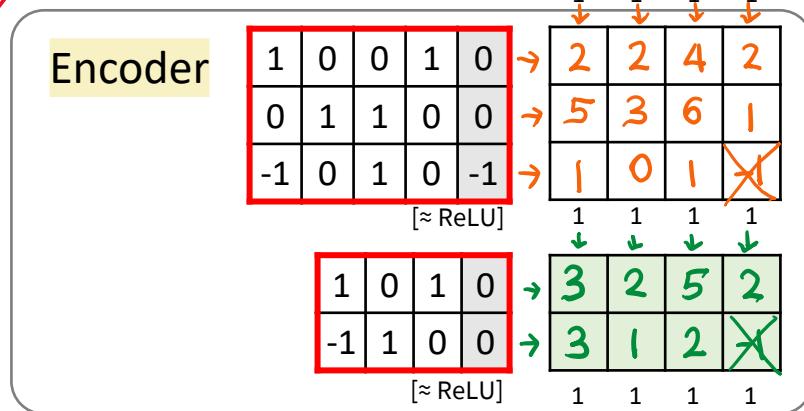
Targets

Y'

MSE Loss
Gradients $\frac{\partial L}{\partial Y}$

x 2

Autoencoder



Variational Auto Encoder

X_1	X_2	X_3
1	1	2
2	1	2
3	2	4
1	1	2
1	1	1

1	0	0	1	0
0	1	1	0	0
-1	0	1	0	-1

[$\approx \text{ReLU}$]

1	0	1	-3
-1	1	0	0

1	1	0	0
0	1	1	-2

Encoder

KL Divergence
Loss Gradients
(SGVB)

$$= \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \mu}$$

$$- \frac{1}{\sigma} = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \sigma}$$

Reparameterization
Trick

0.2	1	0.5
0.5	-1	0.2

$$\mu + \epsilon \odot \sigma$$

ϵ
$\epsilon \odot \sigma$

$\epsilon \odot \sigma$
$\epsilon \odot \sigma$

1 1 1

Decoder

1	0	0
0	1	1
1	-1	0

[$\approx \text{ReLU}$]

1	0	-1	0
1	-1	0	0
0	0	1	1
0	1	1	-3

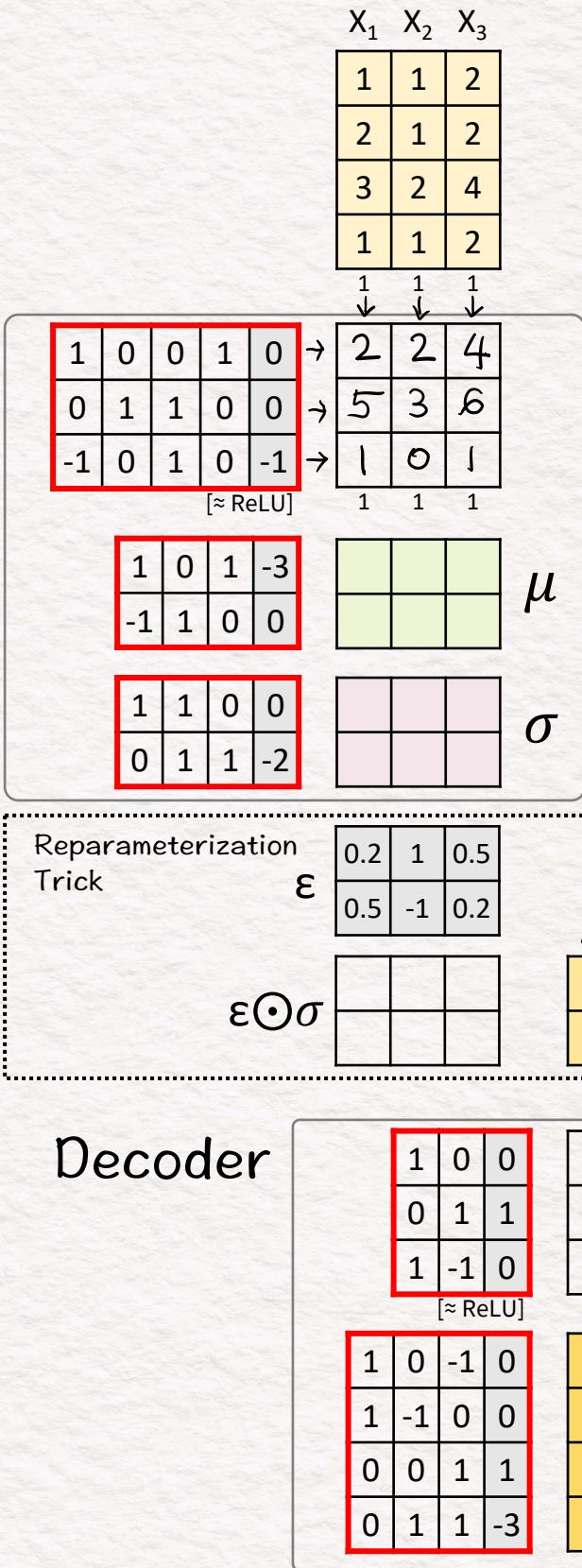
x_1	x_2	x_3
1	1	2
2	1	2
3	2	4
1	1	2

Reconstruction
Loss Gradients

$$(\frac{1}{2} \text{MSE})$$

$$= \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial Y}$$

Variational Auto Encoder



Encoder

KL Divergence
Loss Gradients
(SGVB)

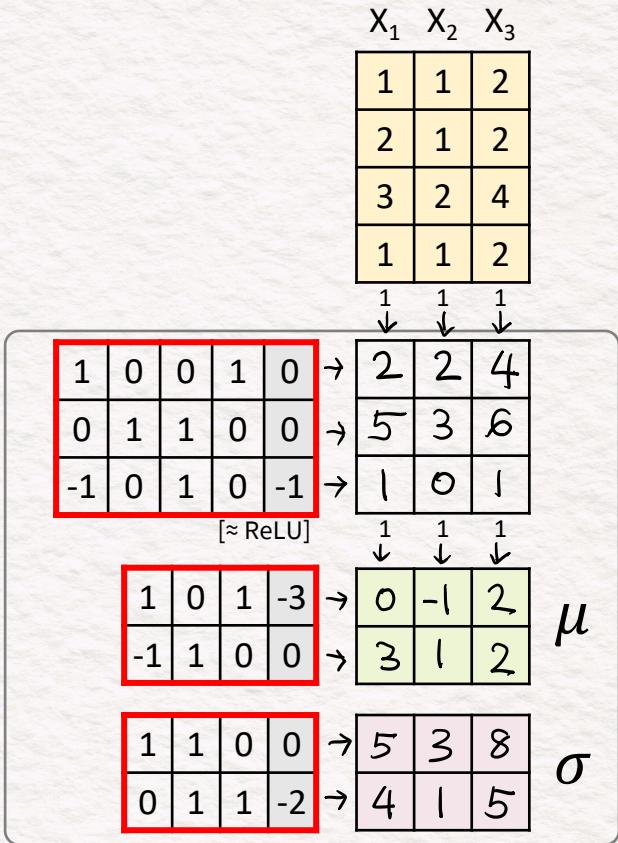
$$= \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \mu}$$

$$- \frac{1}{\sigma} = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \sigma}$$

Reconstruction
Loss Gradients

$$\left(\frac{1}{2} \text{MSE} \right) = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial Y}$$

Variational Auto Encoder



Encoder

KL Divergence
Loss Gradients
(SGVB)

$$= \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \mu}$$

$$- \frac{1}{\sigma} = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \sigma}$$

Reparameterization
Trick

$$\epsilon \begin{array}{|c|c|c|} \hline 0.2 & 1 & 0.5 \\ \hline 0.5 & -1 & 0.2 \\ \hline \end{array}$$

$$\mu + \epsilon \odot \sigma$$

$$\epsilon \odot \sigma \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

1 1 1

Decoder

$$\begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline 0 & 1 & 1 \\ \hline 1 & -1 & 0 \\ \hline \end{array}$$

$\xrightarrow{[\approx \text{ReLU}]}$

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

1 1 1

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 & 0 \\ \hline 1 & -1 & 0 & 0 \\ \hline 0 & 0 & 1 & 1 \\ \hline 0 & 1 & 1 & -3 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

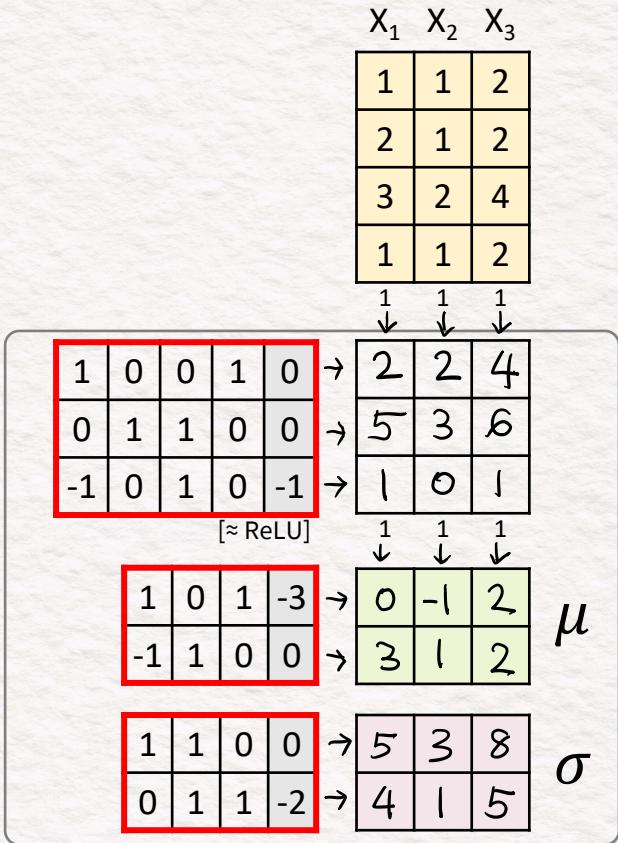
$x_1 \quad x_2 \quad x_3$

Reconstruction
Loss Gradients

$$(\frac{1}{2} \text{MSE})$$

$$= \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial Y}$$

Variational Auto Encoder



Encoder

KL Divergence
Loss Gradients
(SGVB)

$$= \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \mu}$$

$$- \frac{1}{\sigma} = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial \sigma}$$

Reparameterization
Trick

$$\epsilon \begin{array}{|c|c|c|} \hline 0.2 & 1 & 0.5 \\ \hline 0.5 & -1 & 0.2 \\ \hline \end{array}$$

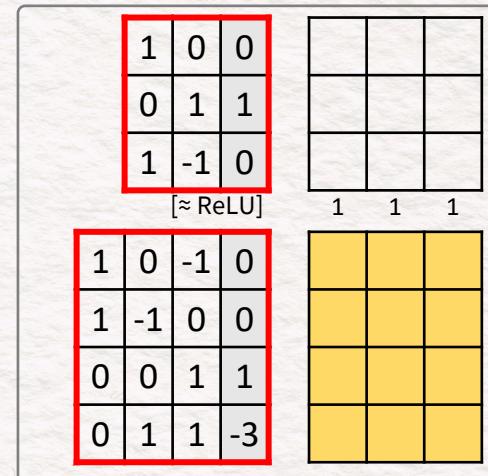
$$\mu + \epsilon \odot \sigma$$

$$\epsilon \odot \sigma \begin{array}{|c|c|c|} \hline 1 & 3 & 4 \\ \hline 2 & -1 & 1 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

1 1 1

Decoder

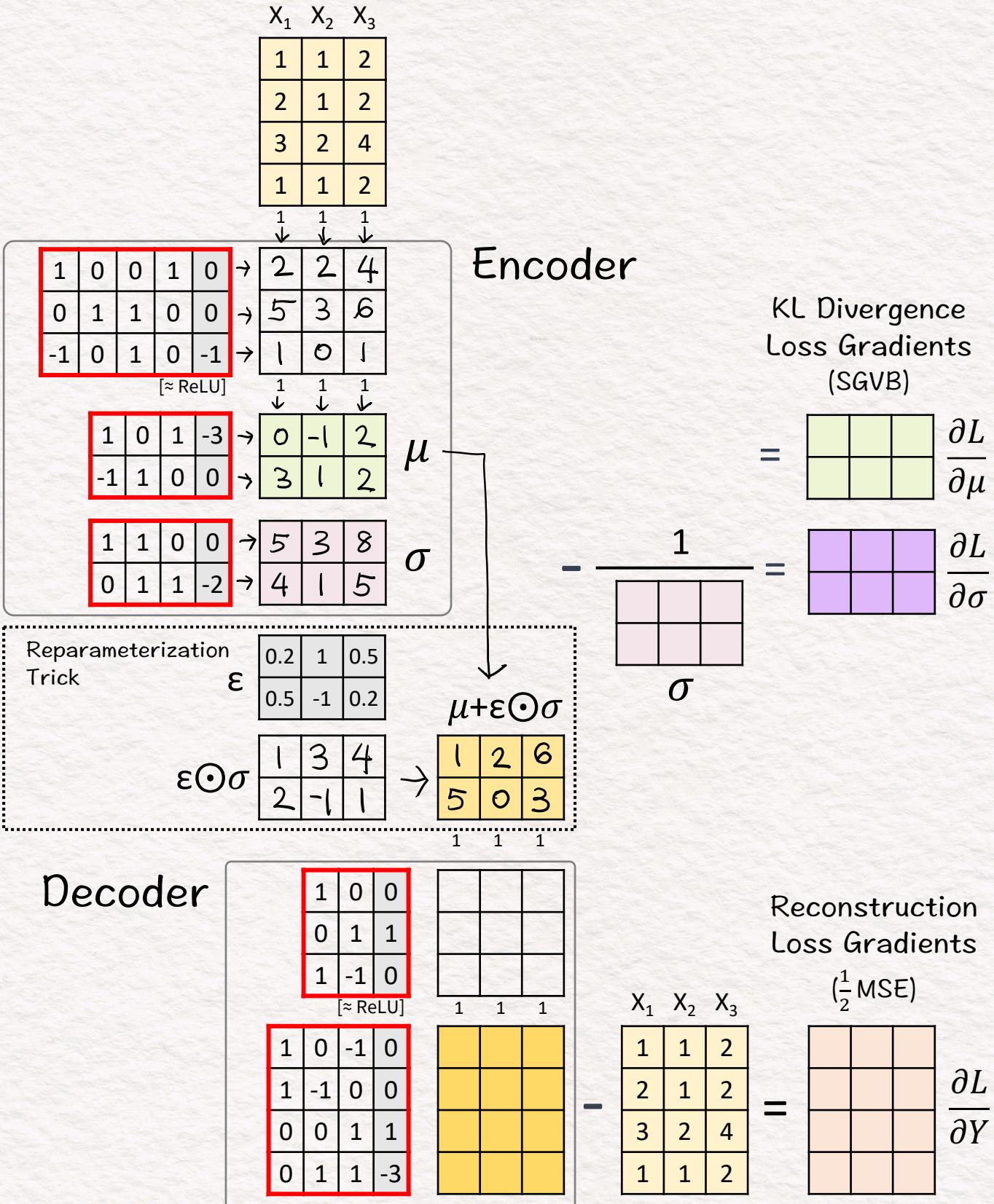


Reconstruction
Loss Gradients

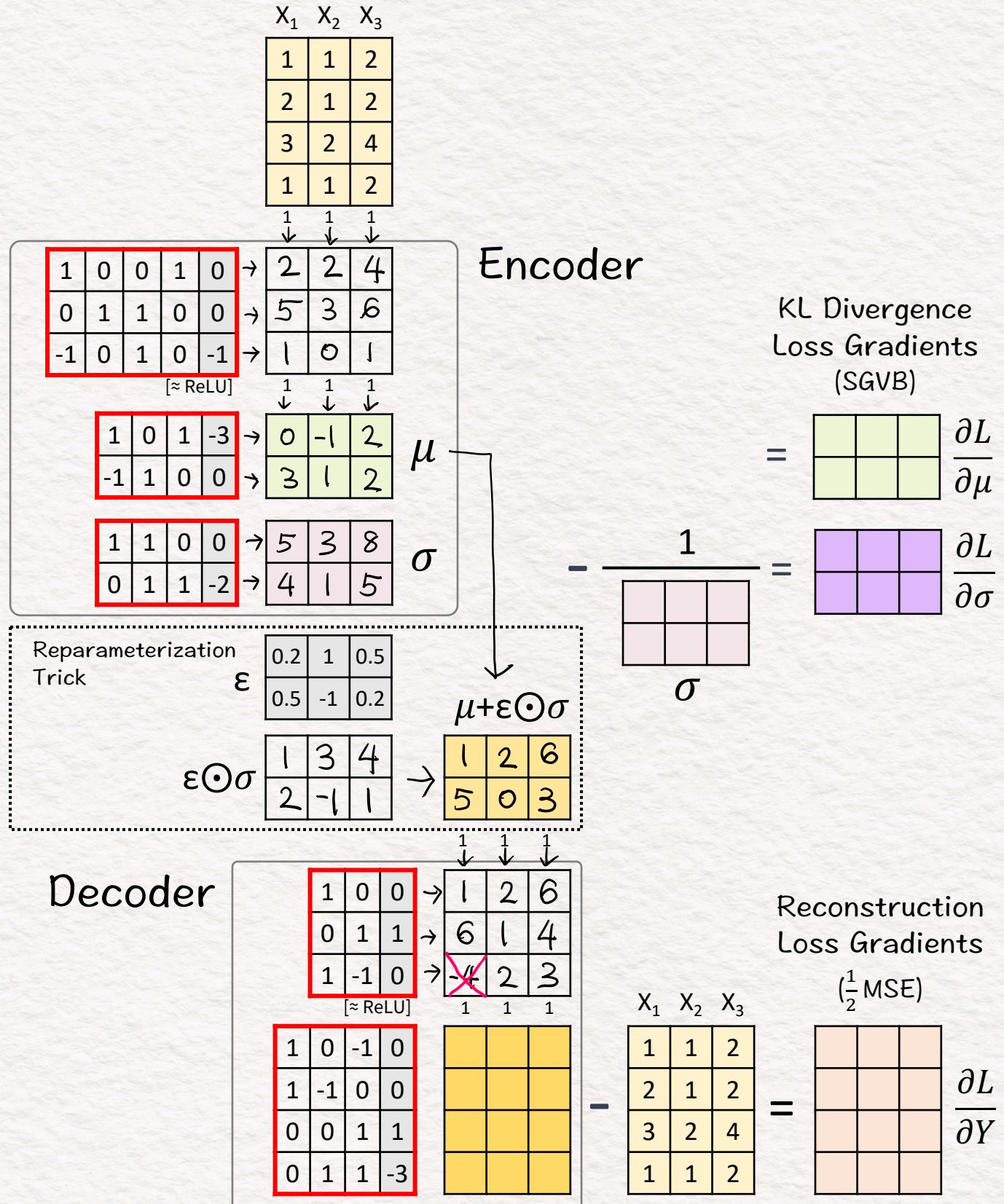
$$(\frac{1}{2} \text{MSE})$$

$$= \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \frac{\partial L}{\partial Y}$$

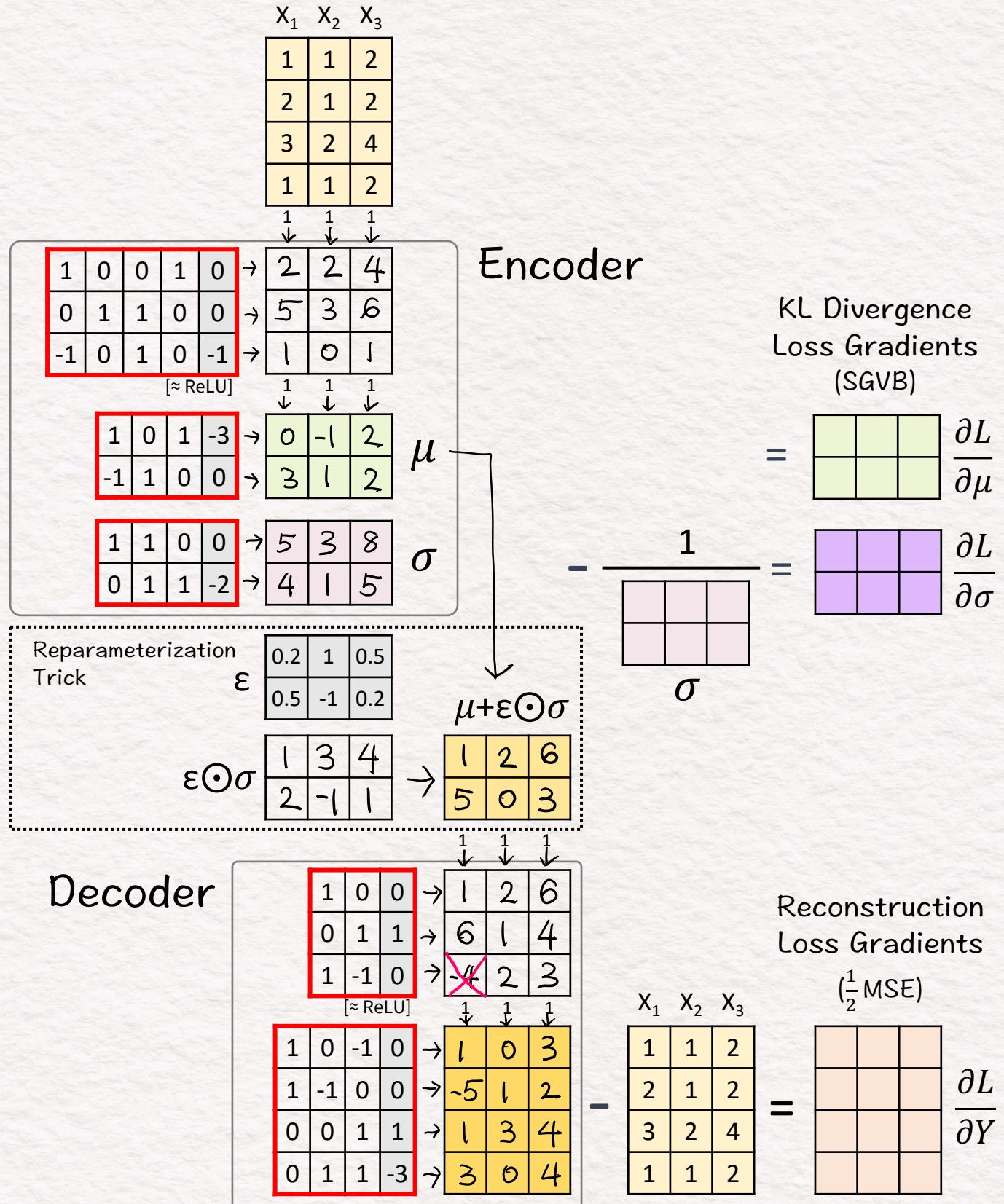
Variational Auto Encoder



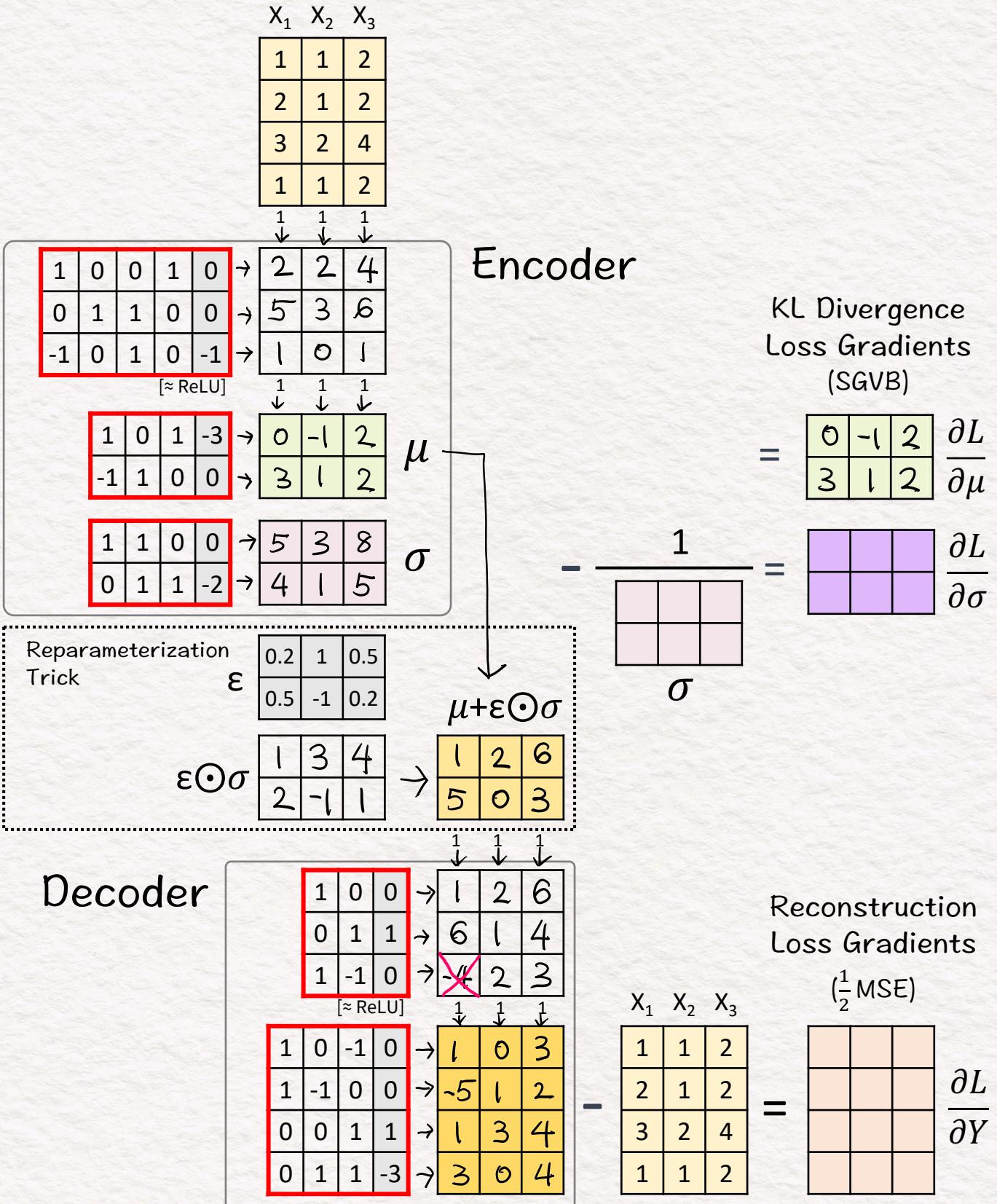
Variational Auto Encoder



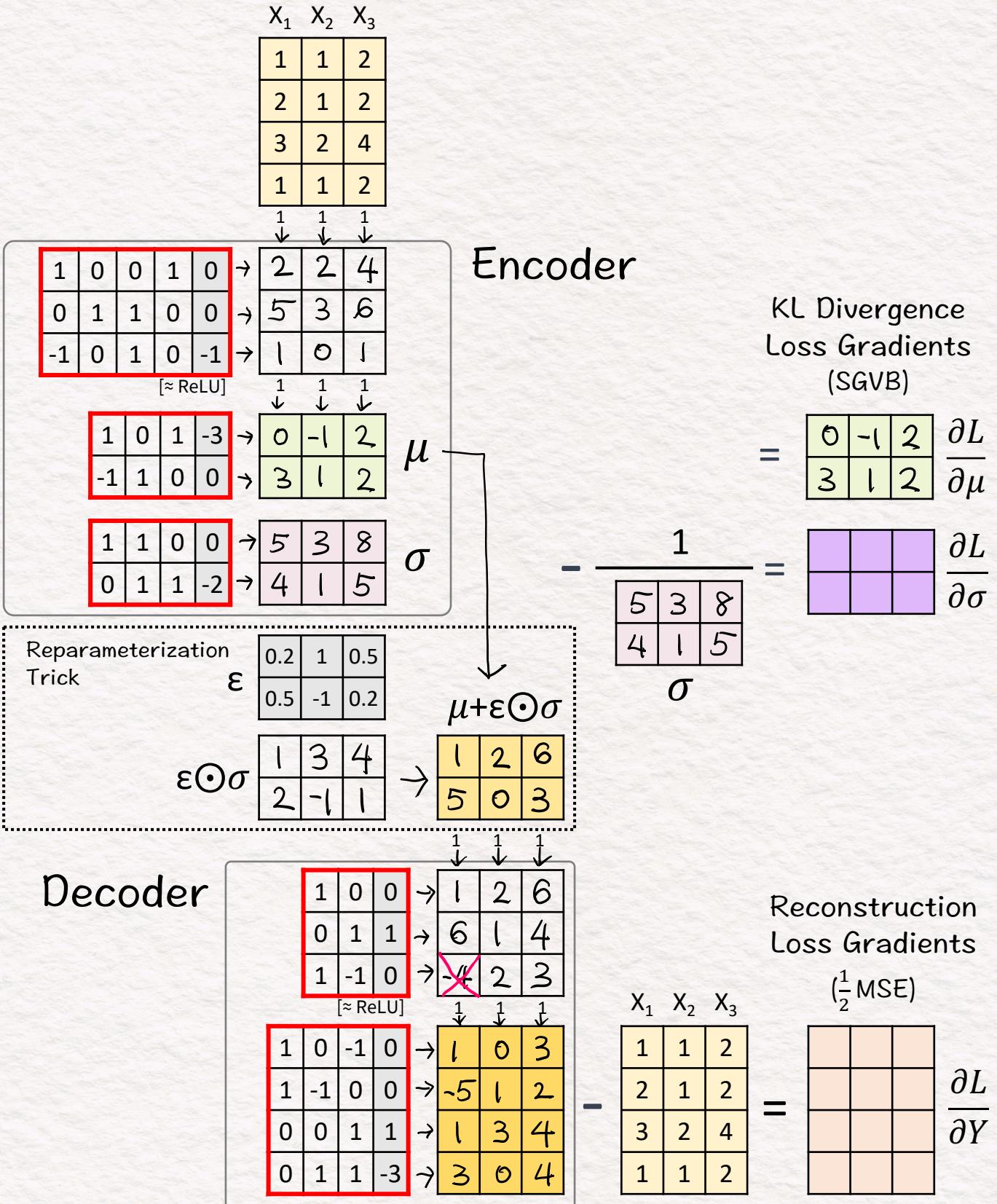
Variational Auto Encoder



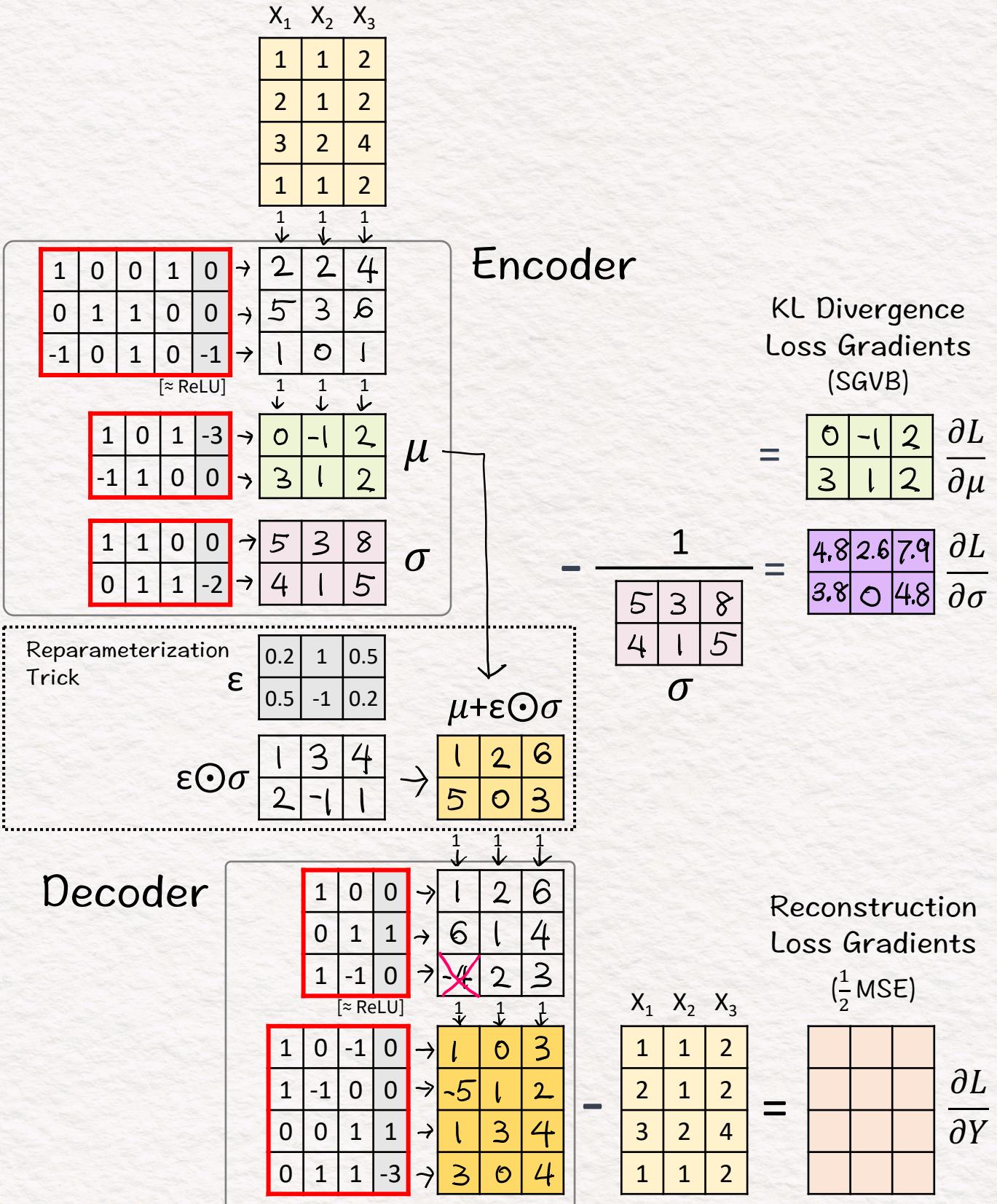
Variational Auto Encoder



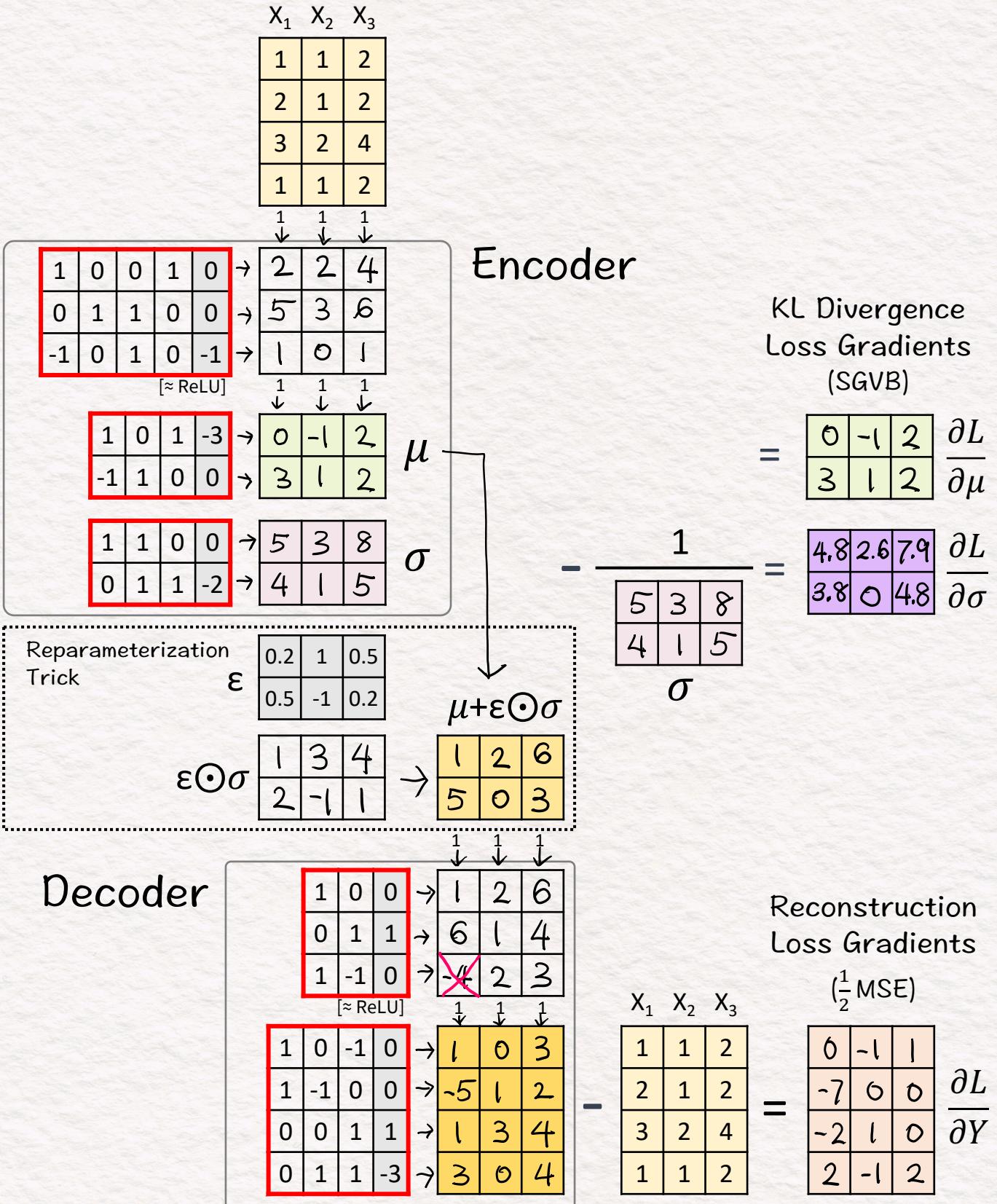
Variational Auto Encoder



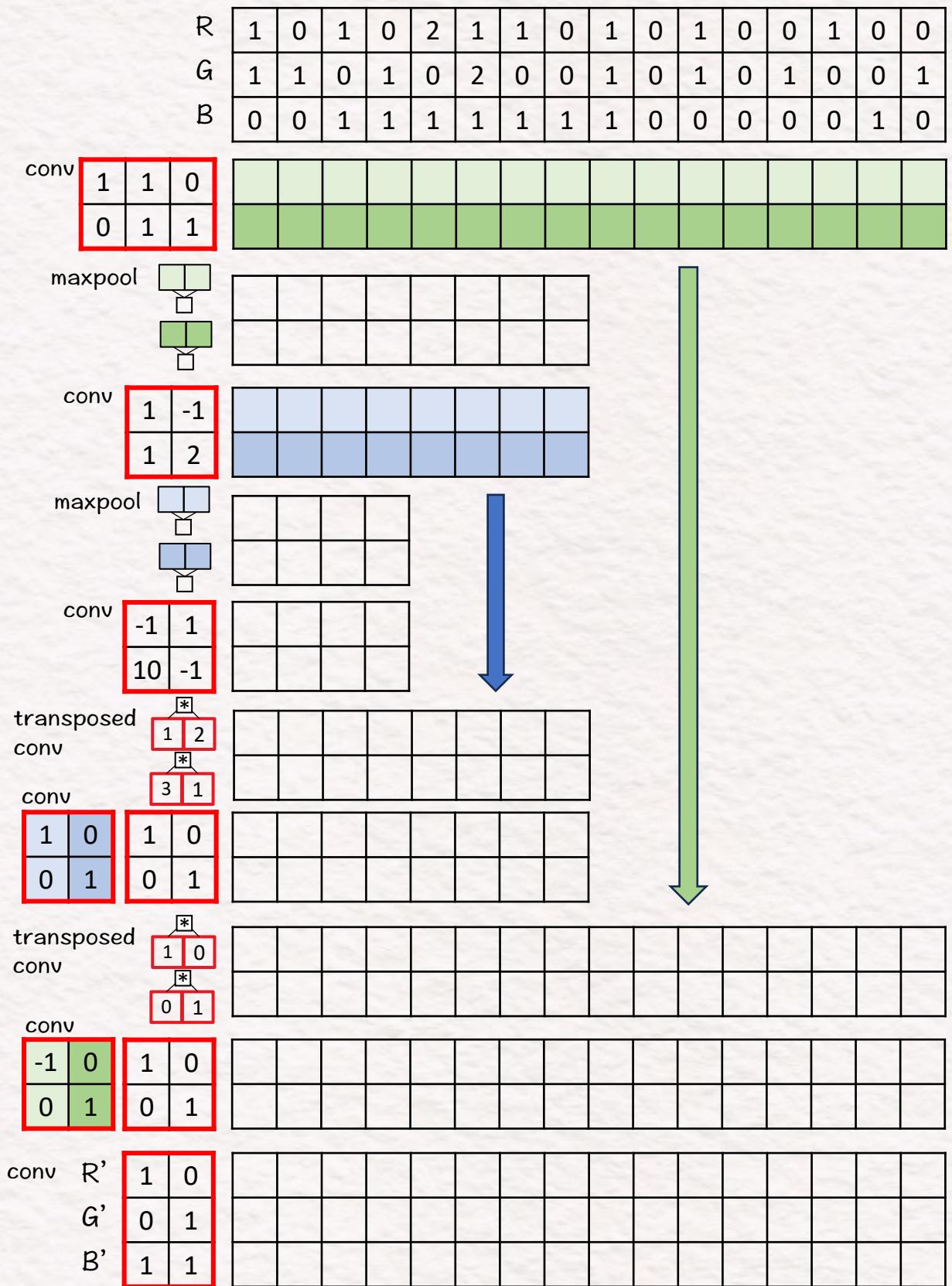
Variational Auto Encoder



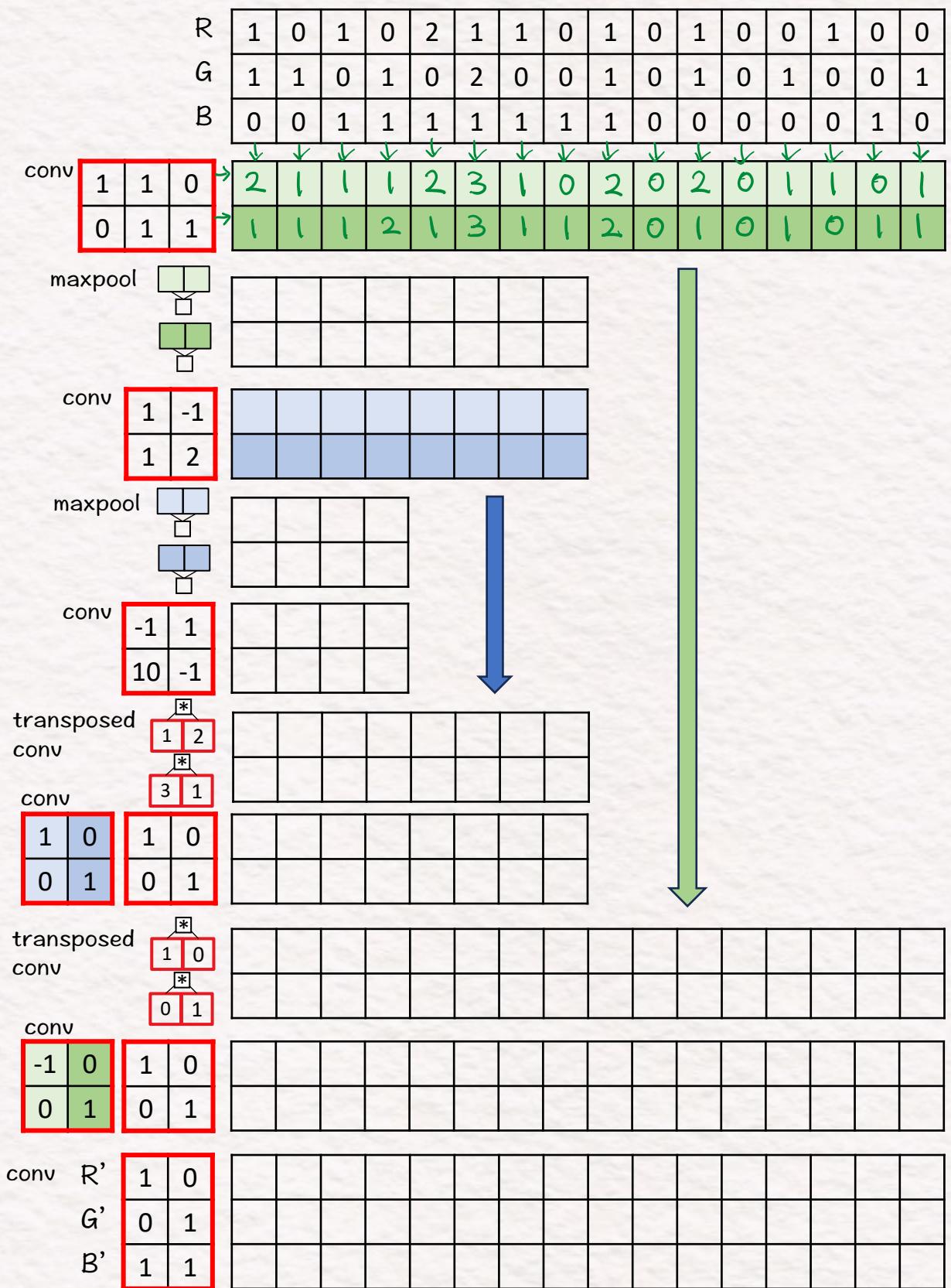
Variational Auto Encoder



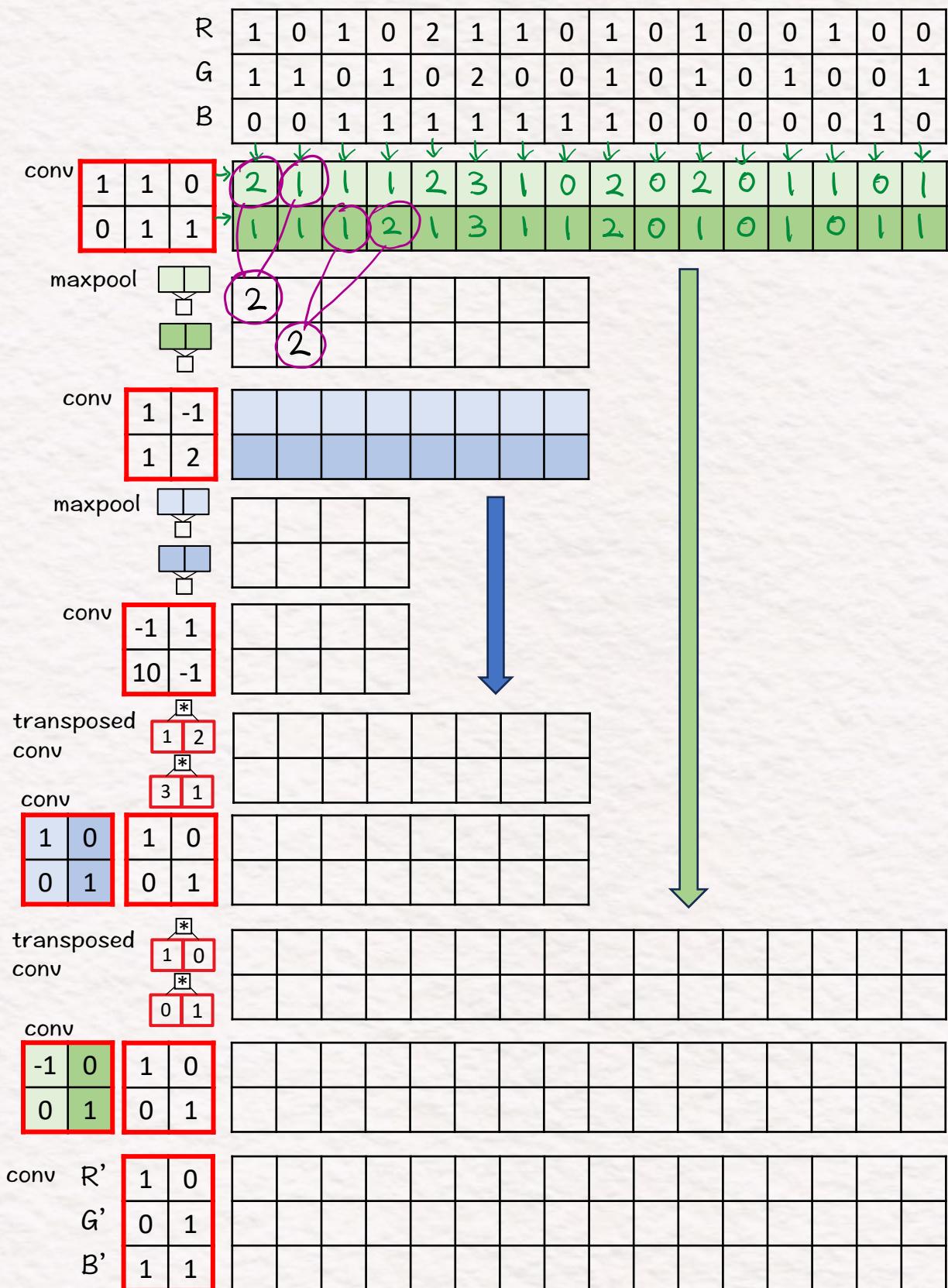
UNet



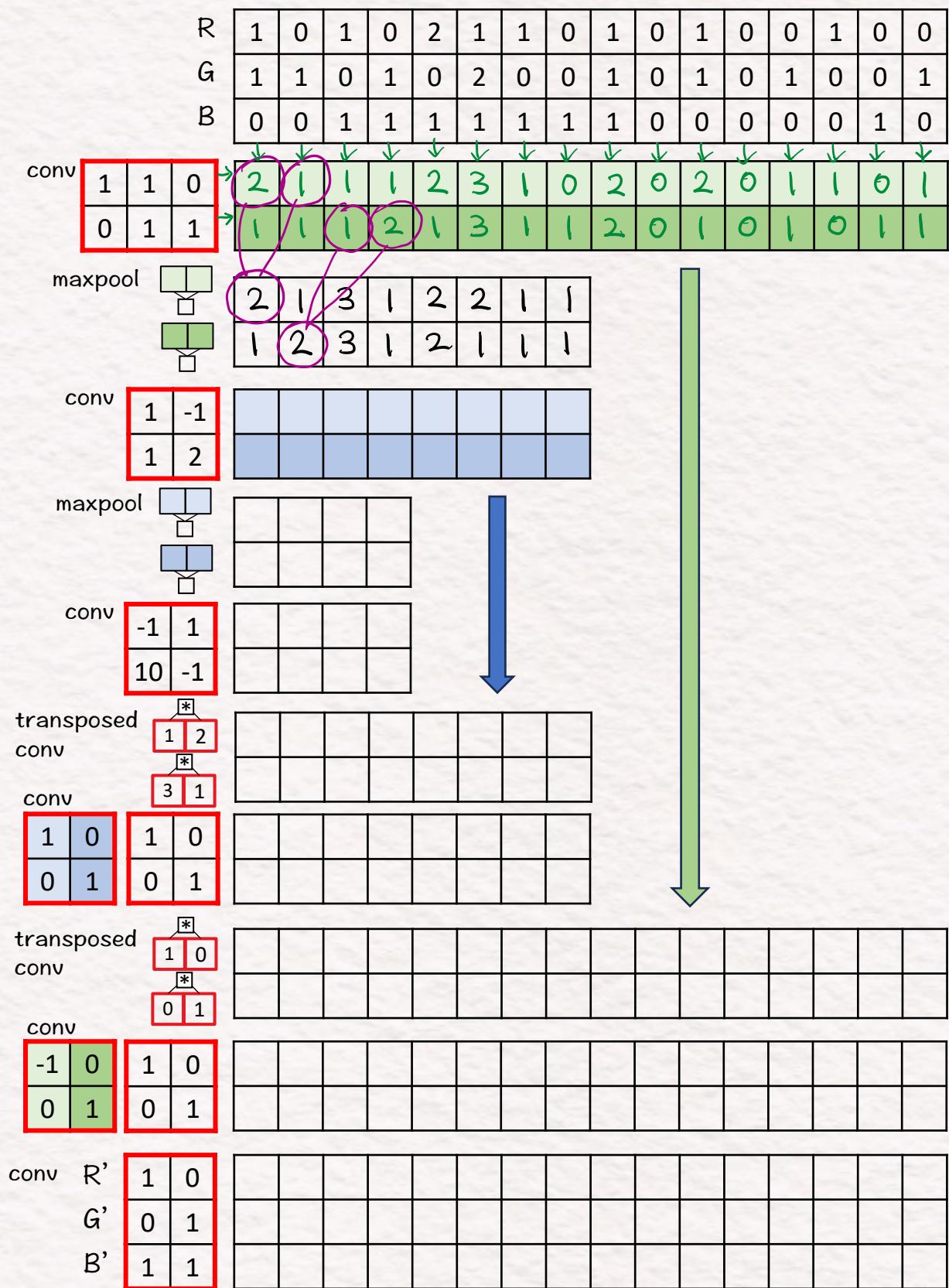
UNet



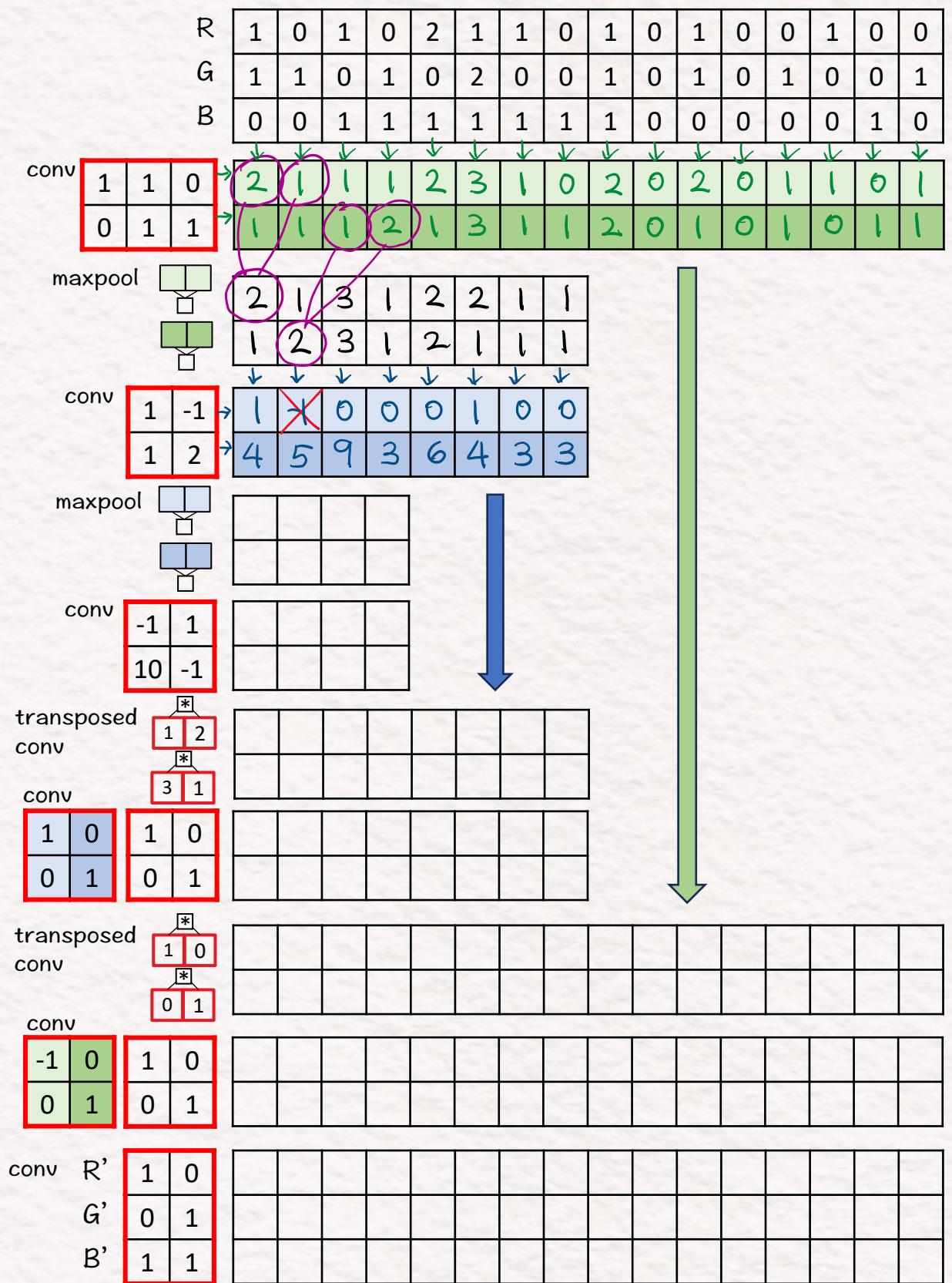
UNet



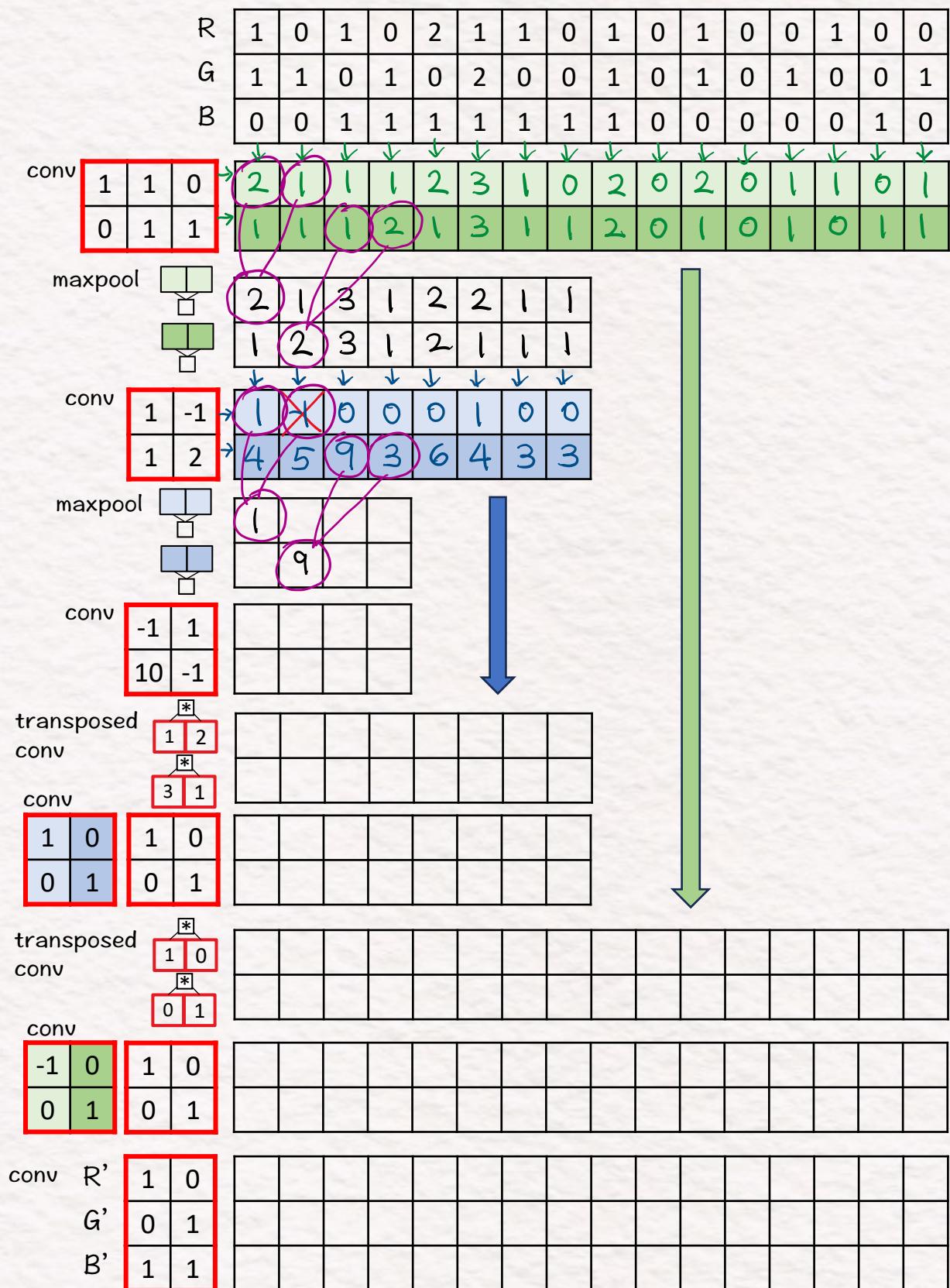
UNet



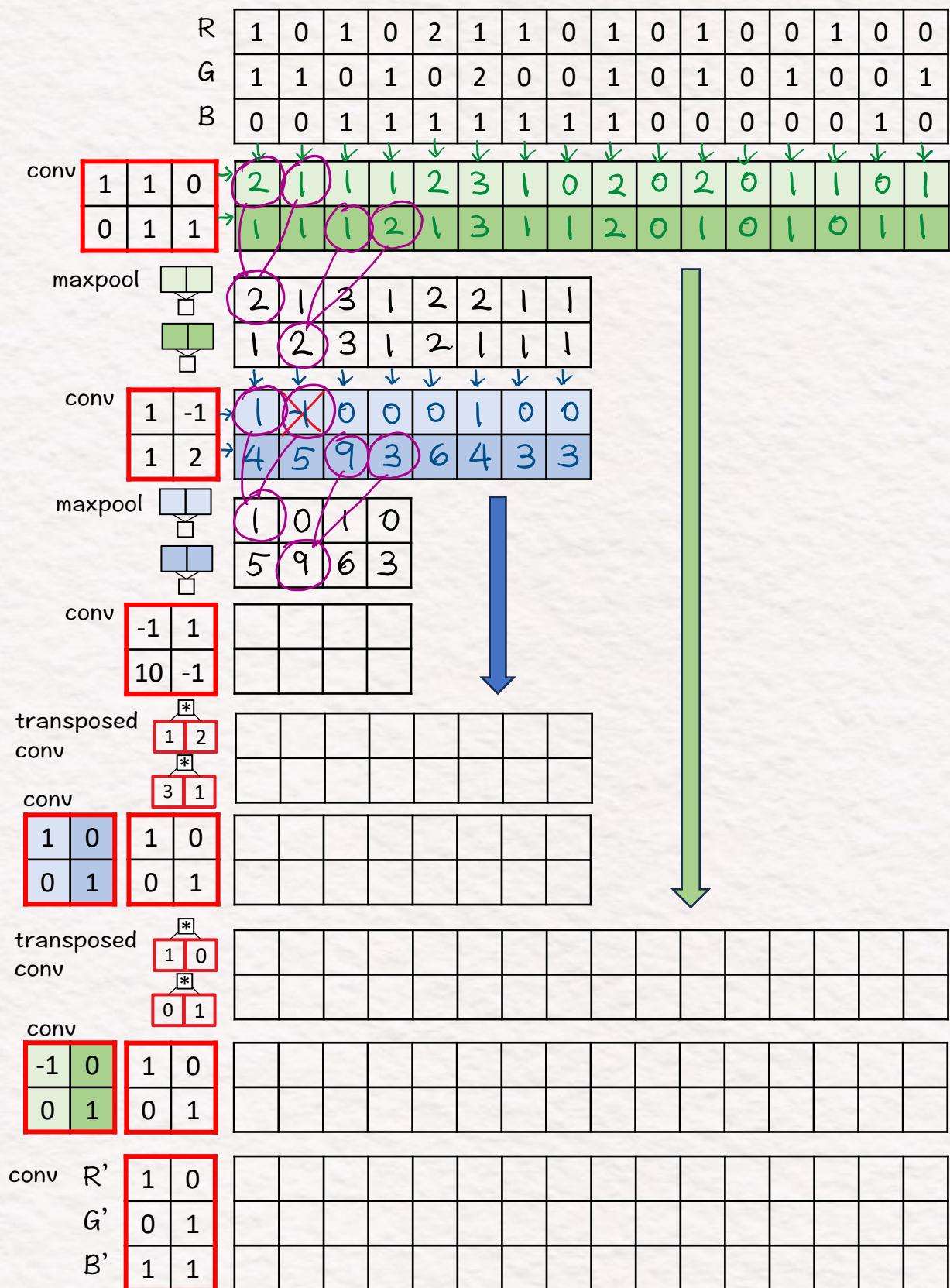
UNet



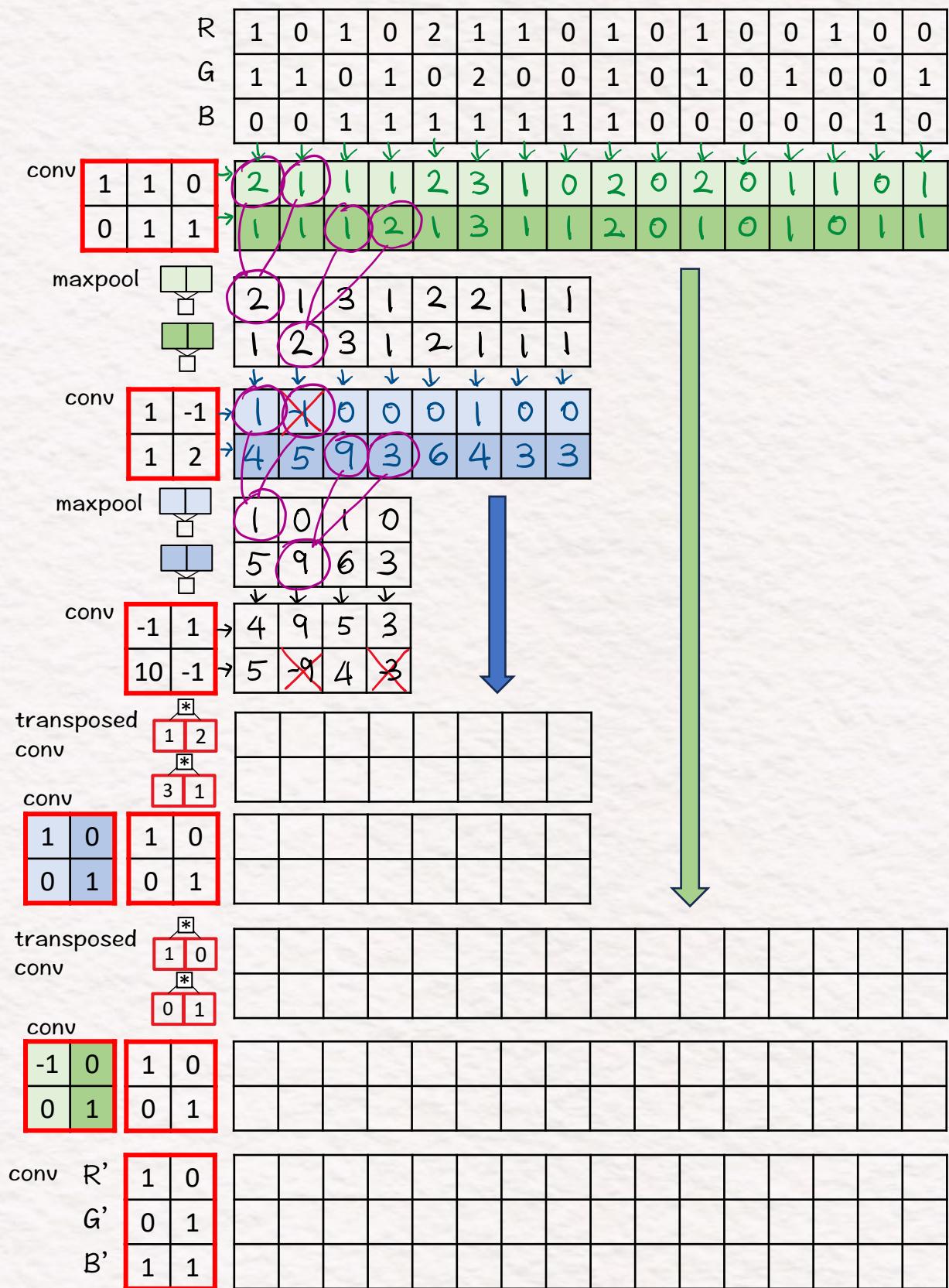
UNet



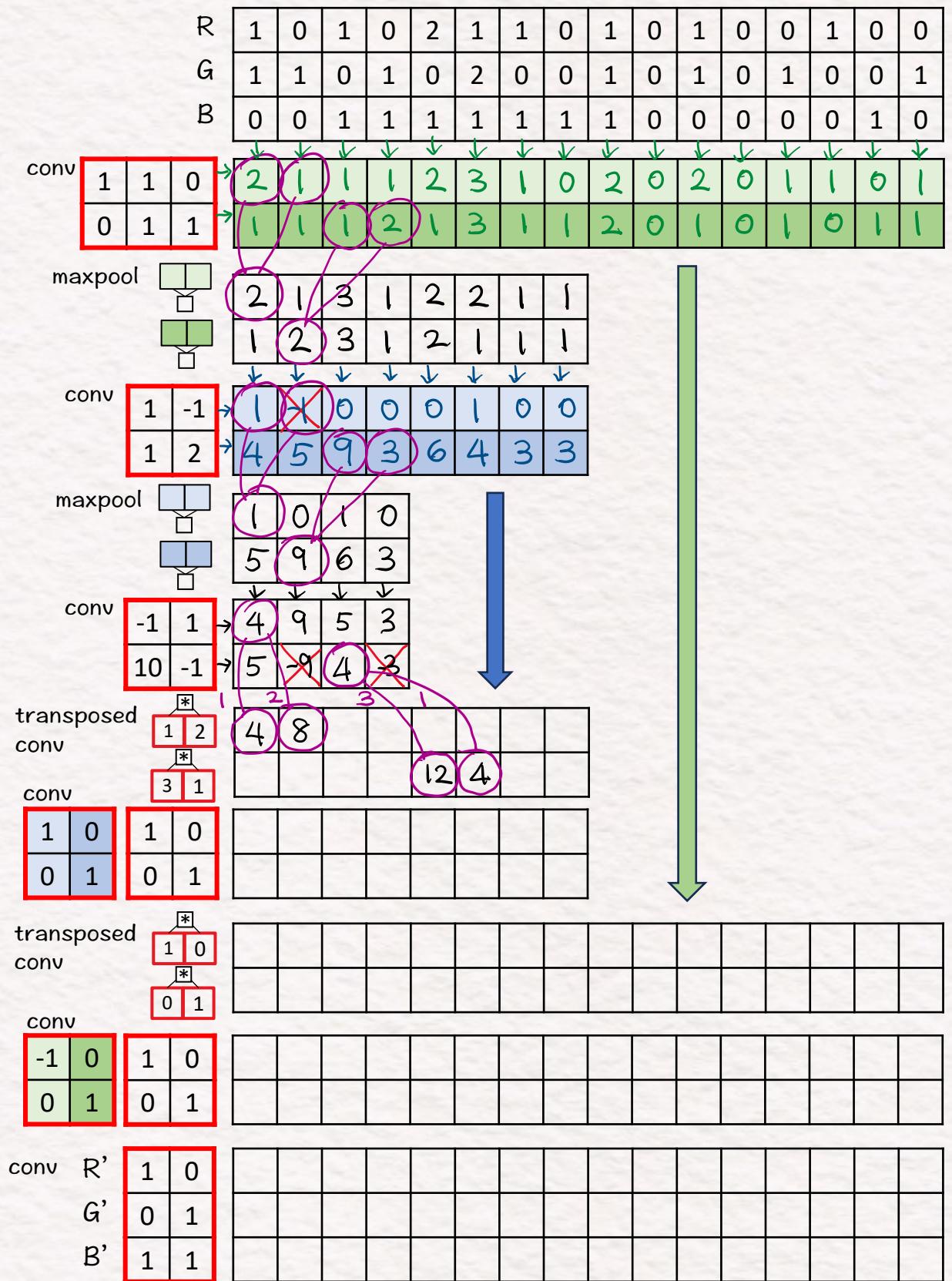
UNet



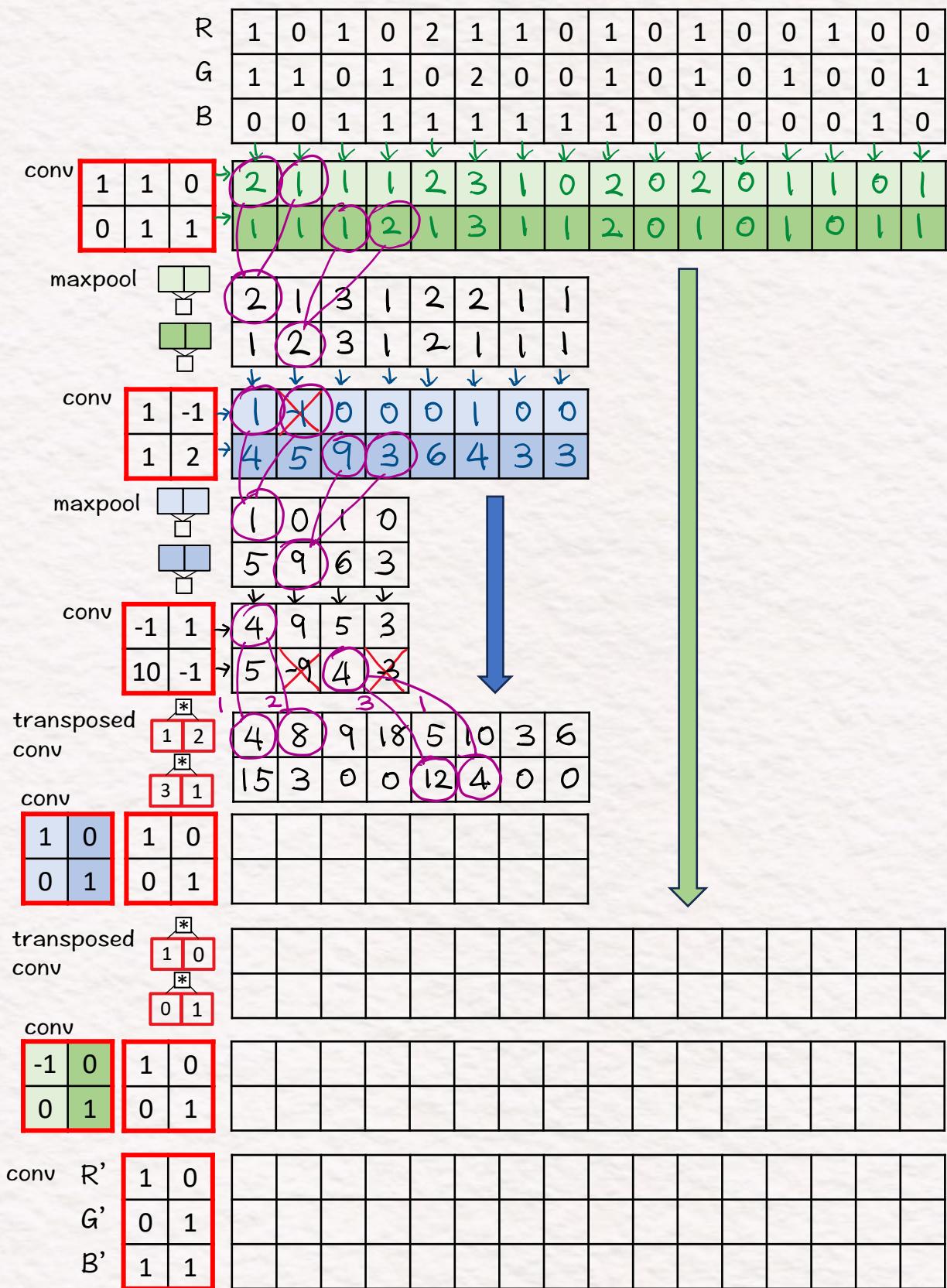
UNet



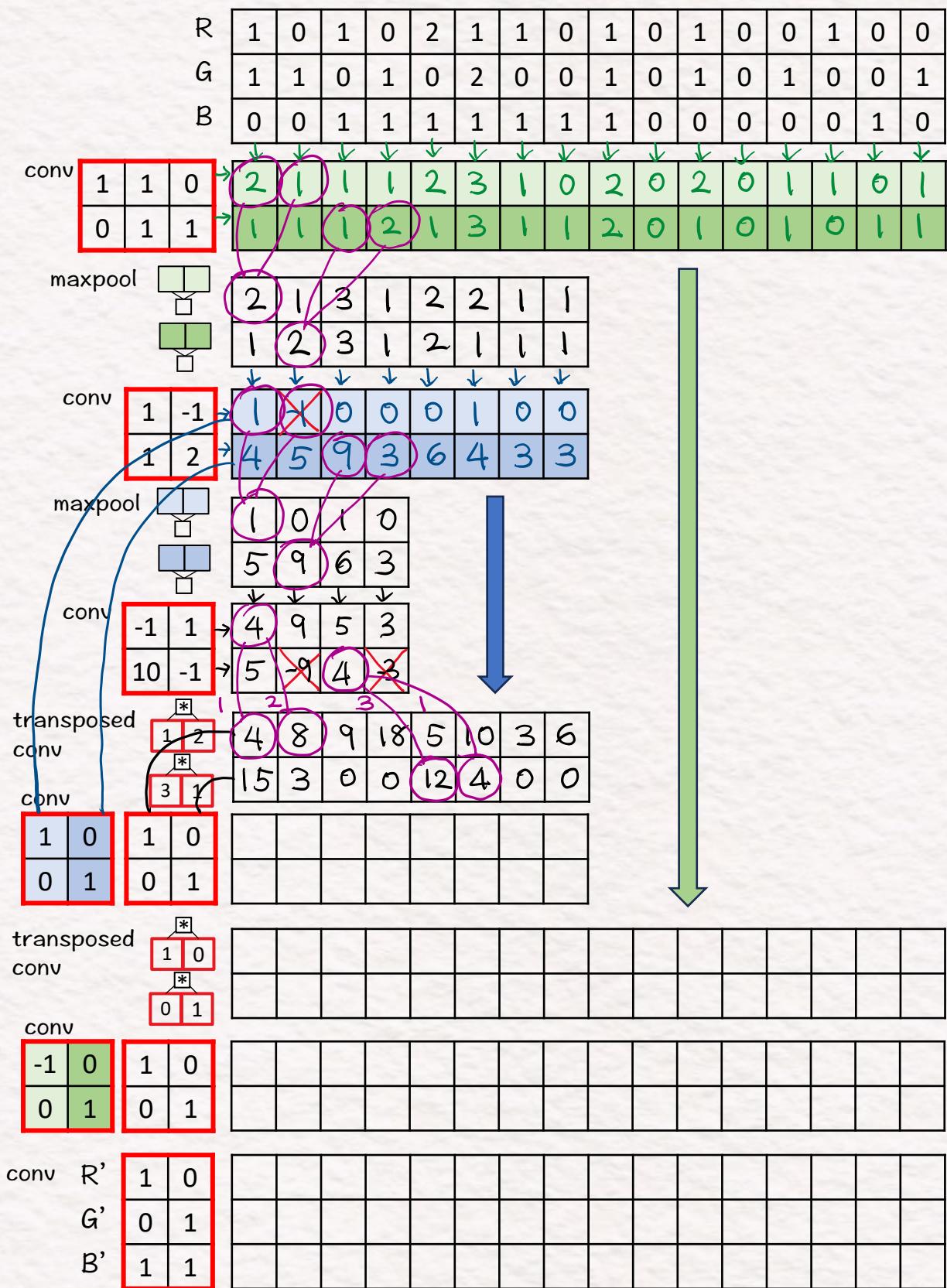
UNet



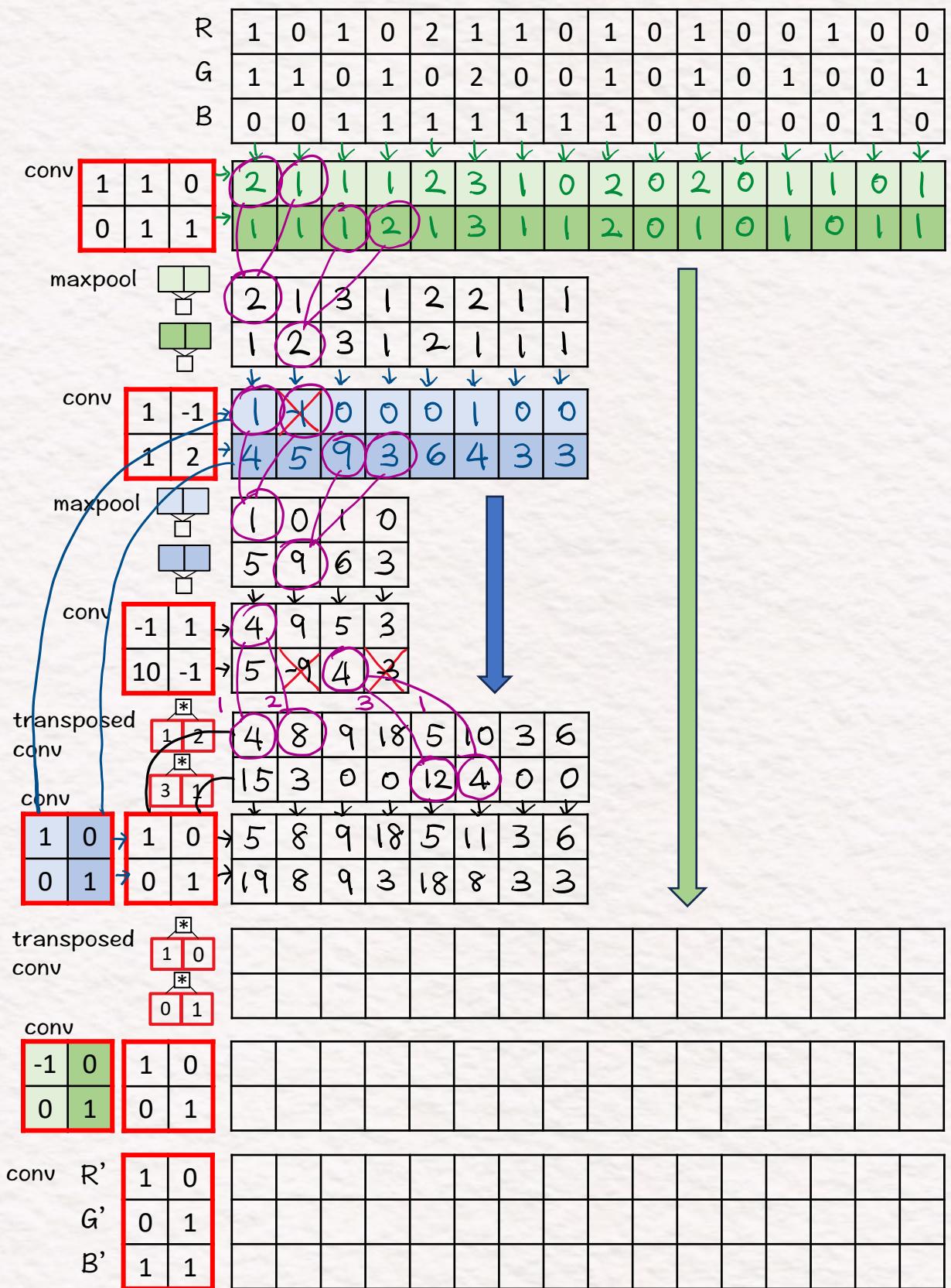
UNet



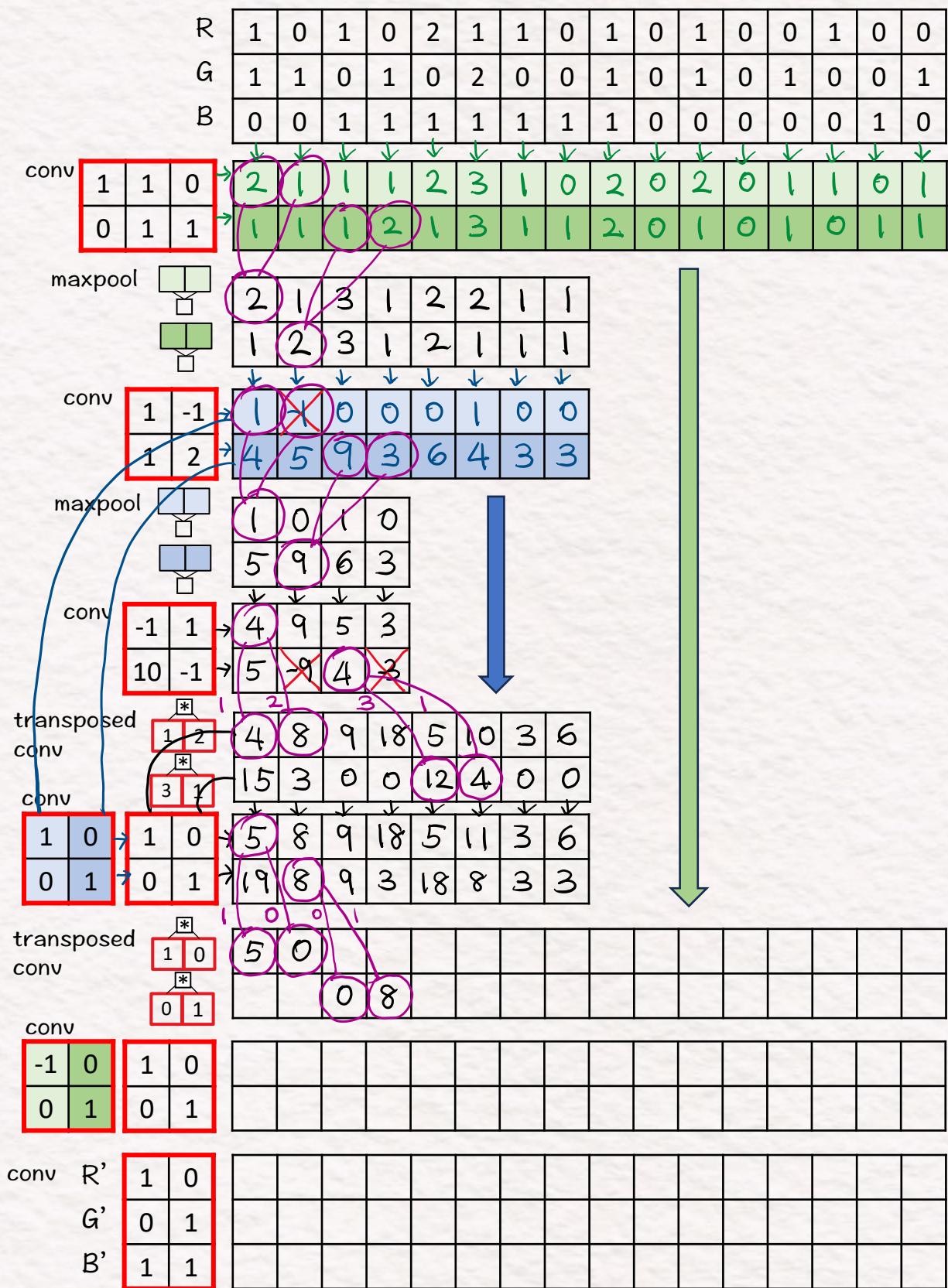
UNet



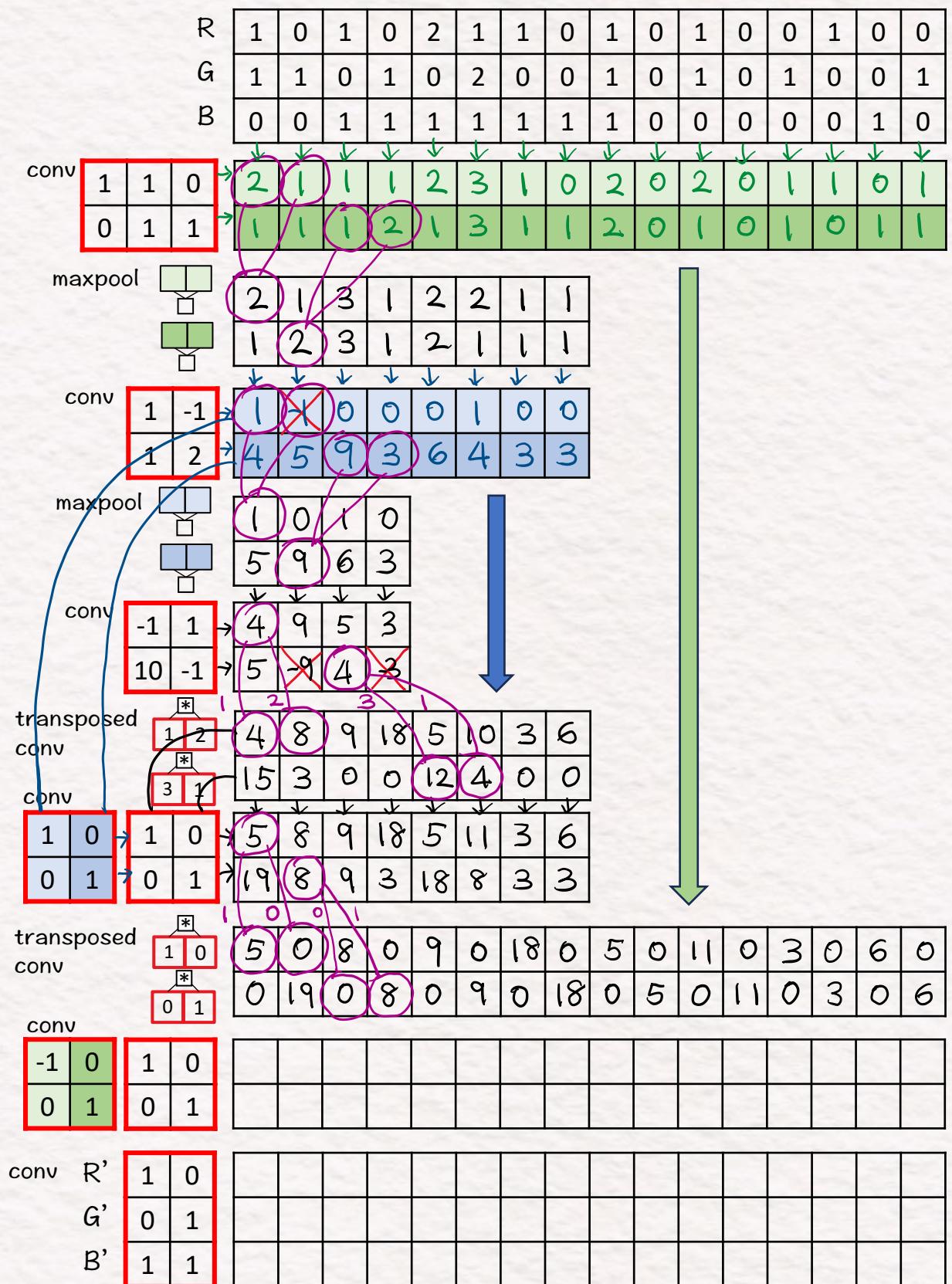
UNet



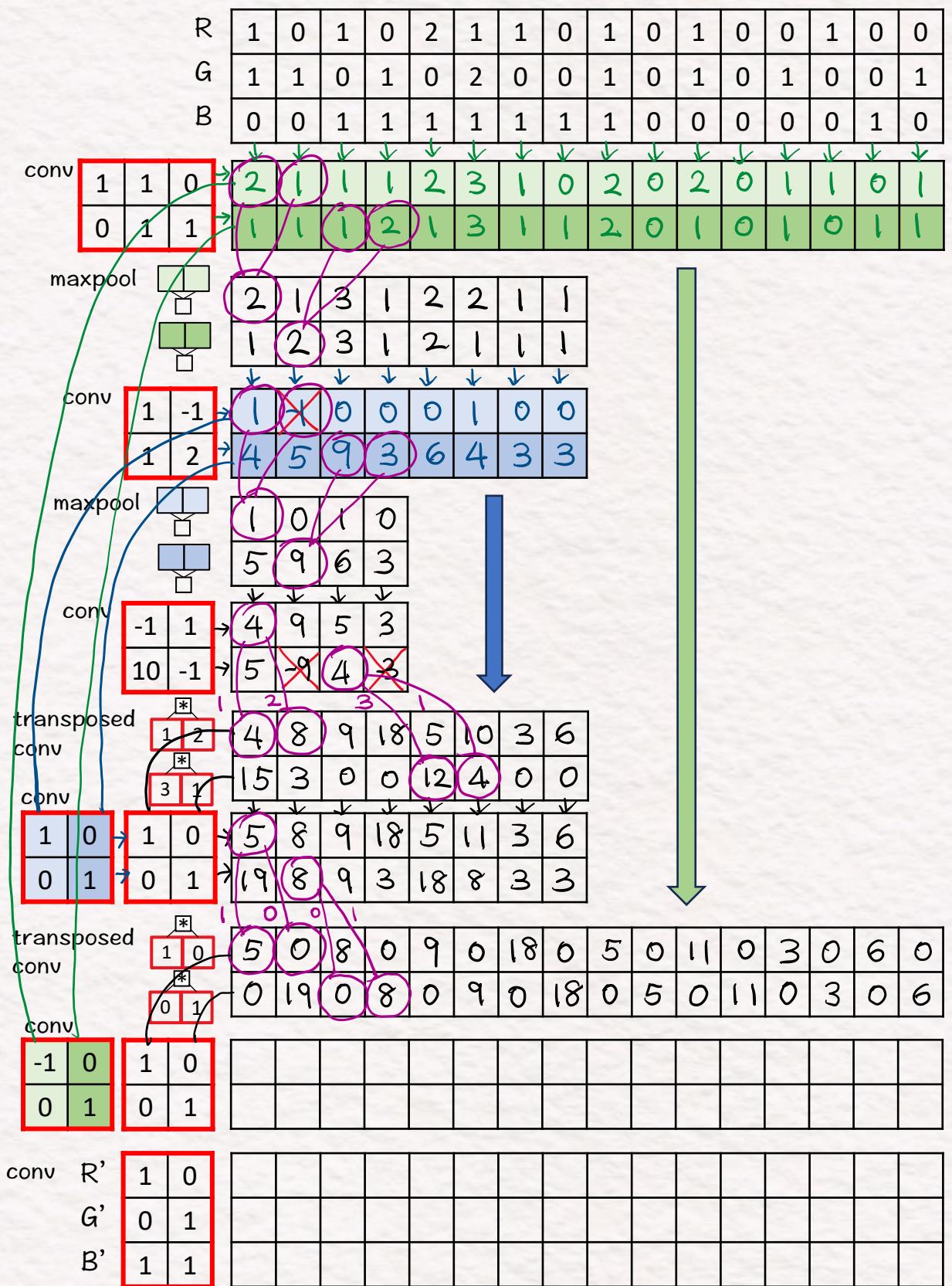
UNet



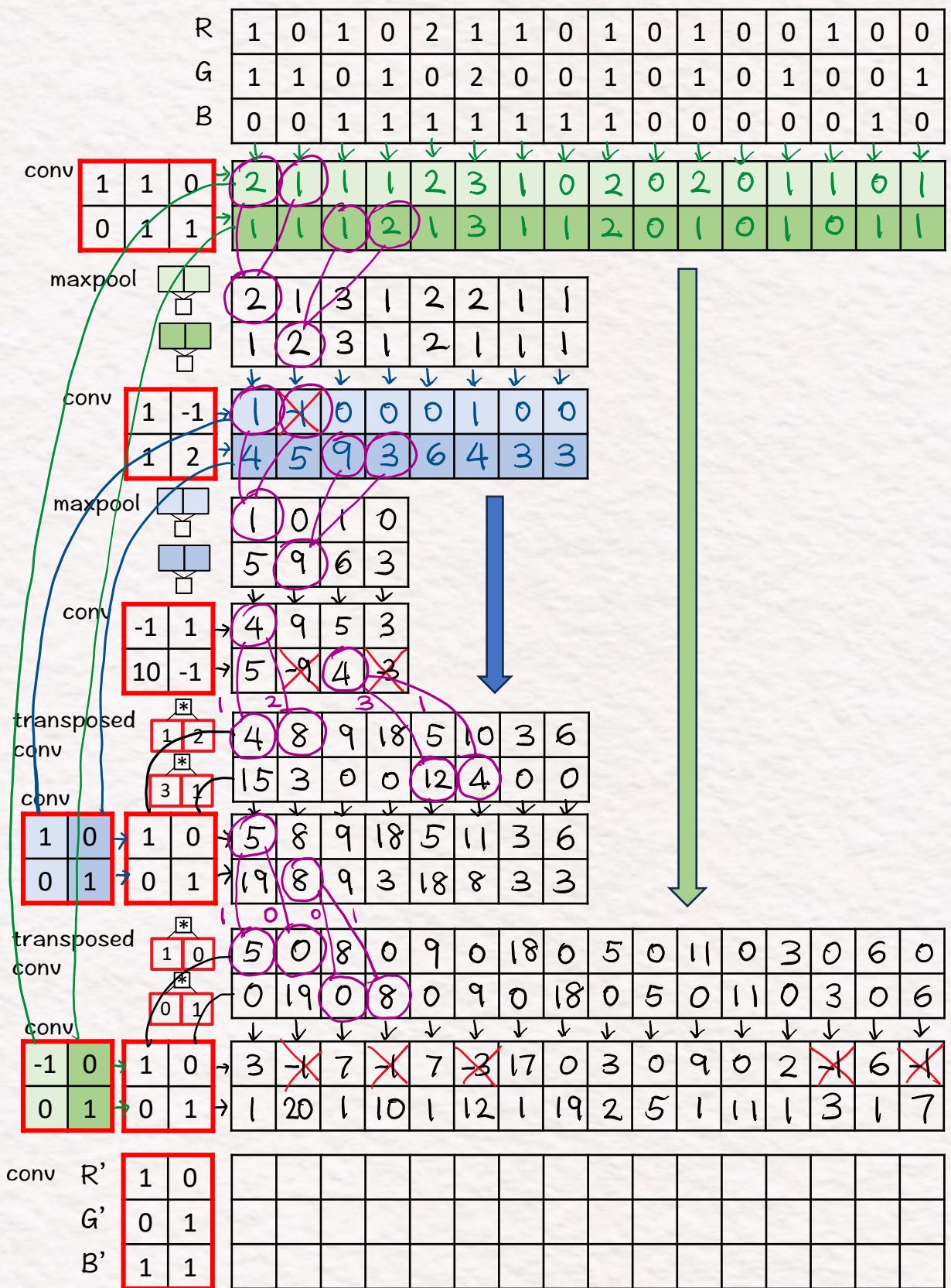
UNet



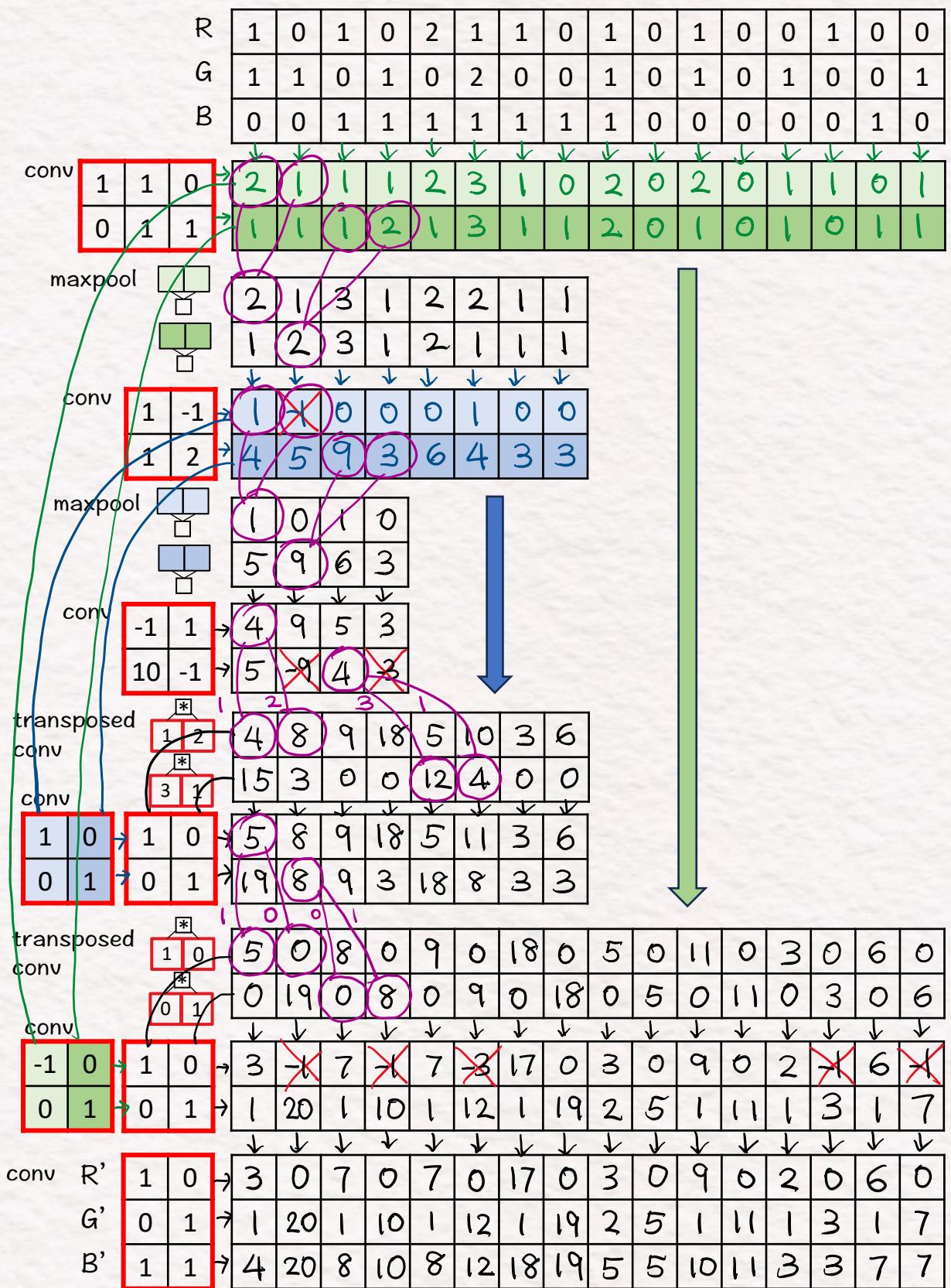
UNet



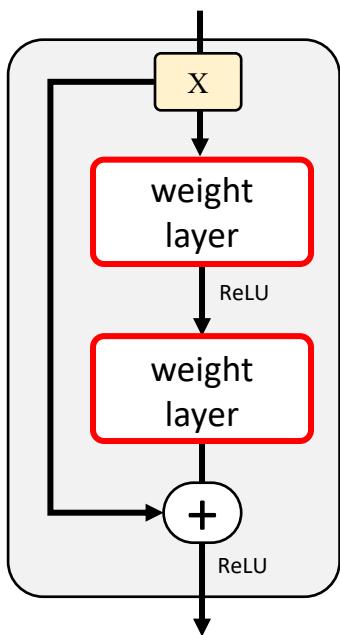
UNet



UNet

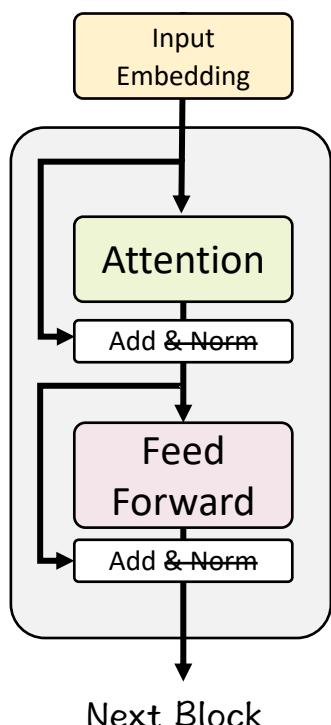


Residual Network



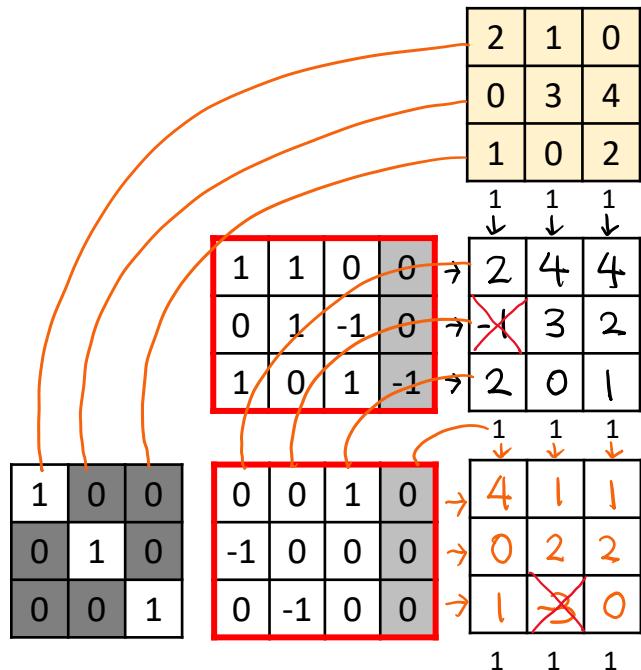
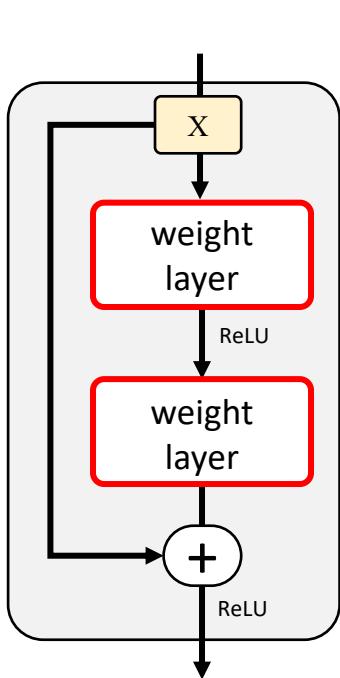
	$\begin{bmatrix} 2 & 1 & 0 \\ 0 & 3 & 4 \\ 1 & 0 & 2 \end{bmatrix}$	$\begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 1 & 0 & 1 & -1 \end{bmatrix}$	$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & 1 & 1 \end{bmatrix}$
	$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & 1 & 1 \end{bmatrix}$
	$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & 1 & 1 \end{bmatrix}$

Transformer's Encoder Block

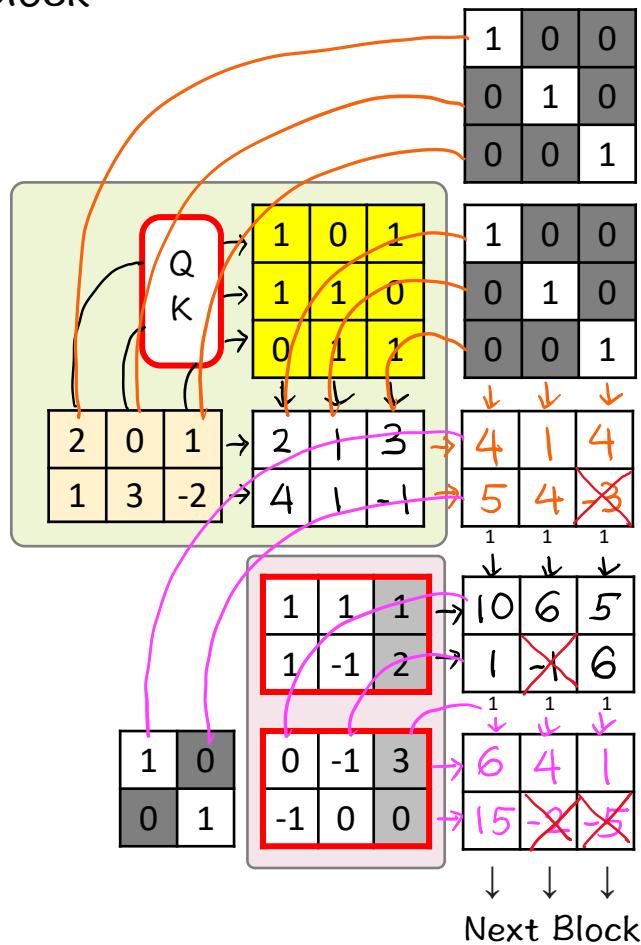
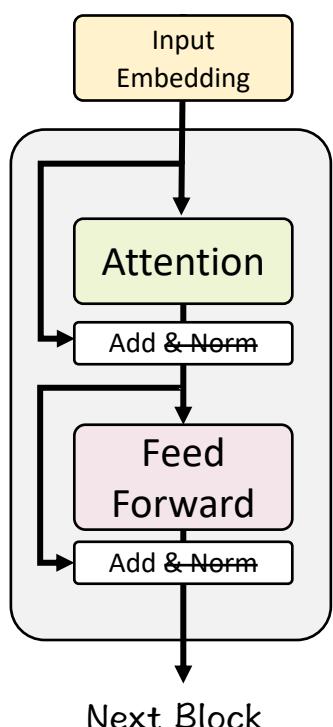


$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$	$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & 1 & 1 \end{bmatrix}$
$\begin{bmatrix} 2 & 0 & 1 \\ 1 & 3 & -2 \end{bmatrix}$	$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & 1 & 1 \end{bmatrix}$	$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & 1 & 1 \end{bmatrix}$
$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 2 \\ -1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & 1 & 1 \end{bmatrix}$
		$\downarrow \downarrow \downarrow$
		Next Block

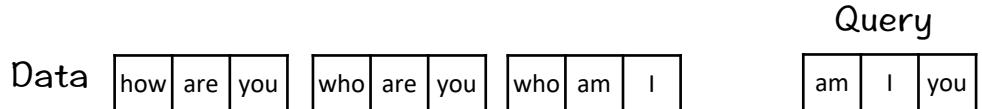
Residual Network



Transformer's Encoder Block

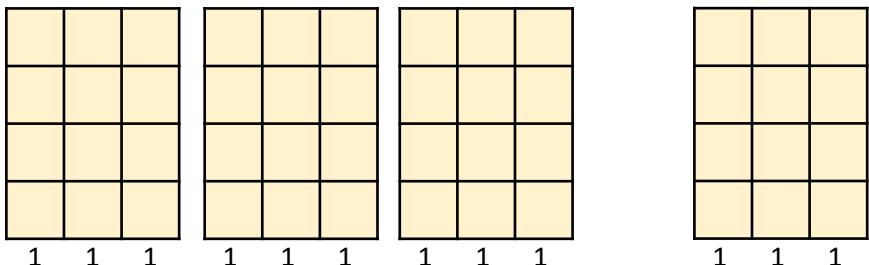


Vector Database



Word Embeddings

a	an	the	how	why	who	what	are	is	am	be	was	you	we	I	they	she	he	she	me	him	her
0	-1	0	1	0	1	0	0	-1	1	0	0	0	3	1	0	-1	0	0	0	-1	0
2	0	2	0	0	0	-1	1	0	0	0	2	1	0	2	0	2	0	0	2	0	0
-1	0	-1	1	2	0	0	1	0	1	-1	0	0	-1	0	3	0	0	-1	0	2	-1
0	1	0	0	1	0	1	0	1	0	1	-2	0	0	0	1	0	1	0	1	0	1

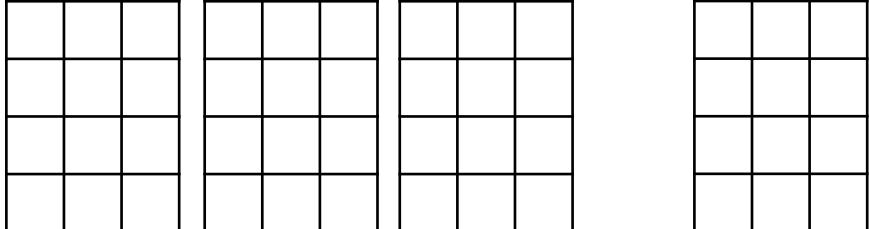


Text Embeddings

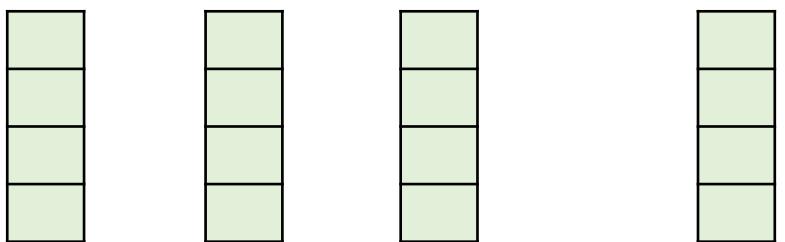
Encoder

1	1	0	0	0
0	1	0	1	0
1	0	1	0	-1
1	-1	0	0	0

Linear & ReLU



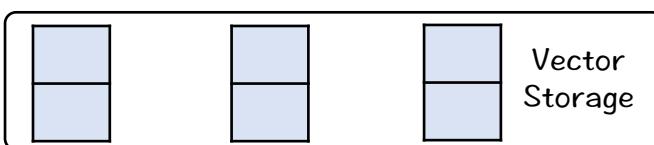
Mean Pooling



Indexing

Projection

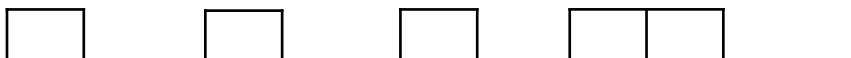
1	1	0	0
0	0	1	1



Vector Storage

Retrieval

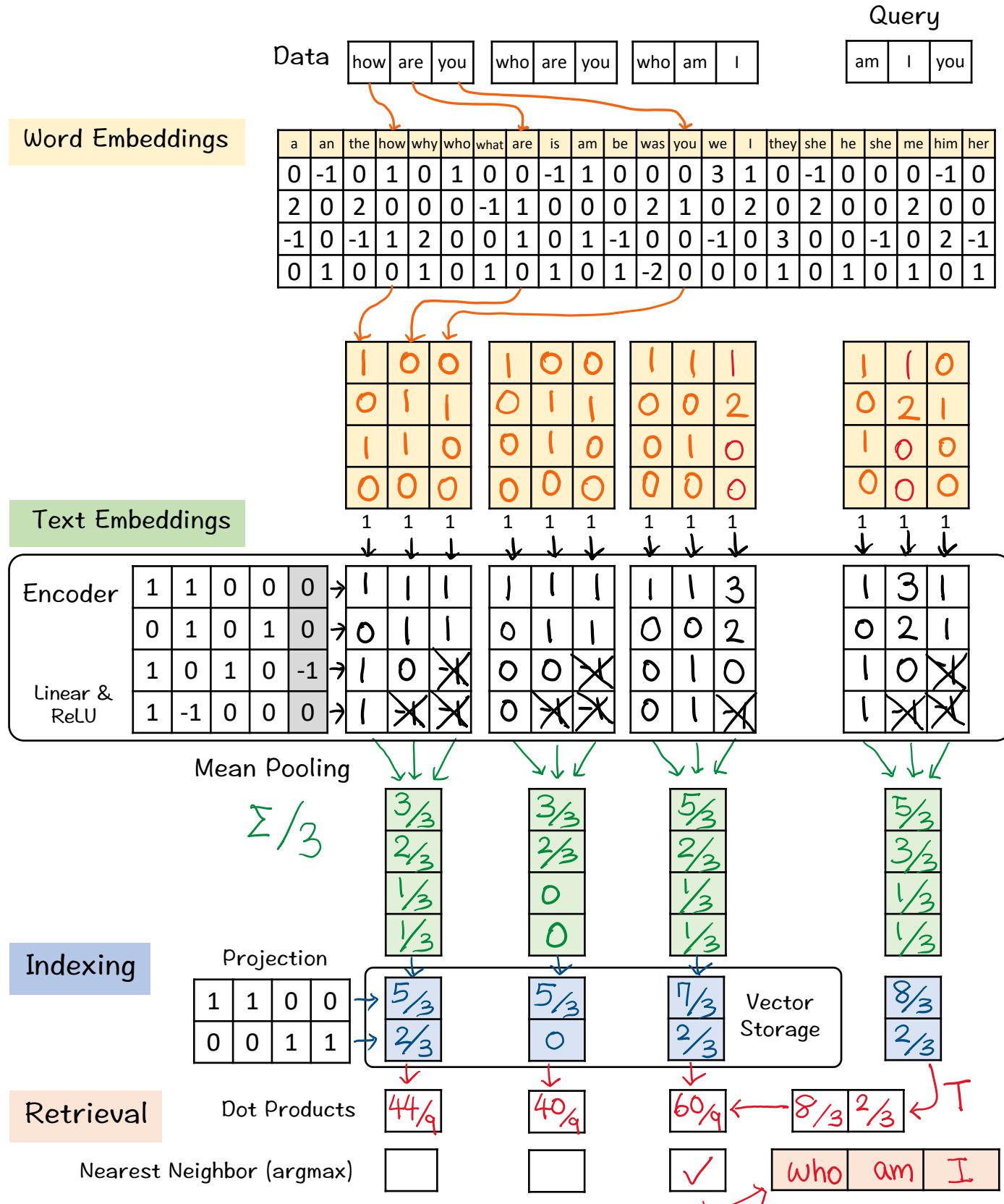
Dot Products



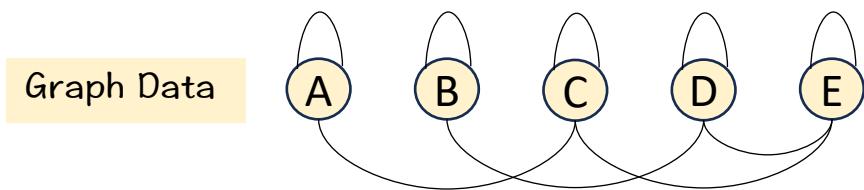
Nearest Neighbor (argmax)



Vector Database



Graph Convolutional Network



Graph
Convolutional
Network

	A	B	C	D	E
A	2	0	1	0	1
B	1	1	0	0	0
C	0	0	-1	1	1
D	0	3	0	1	0
E	1	1	1	1	1

1	1	0	0	0
0	1	0	-1	0
1	0	0	1	-1

[ReLU]

Messages

	A	B	C	D	E
A					
B					
C					
D					
E					

Adjacency
Matrix

	A	B	C	D	E

0	1	1	0
1	-1	0	-2
1	0	0	0

[ReLU]

Messages

1 1 1 1 1

Fully
Connected
Network

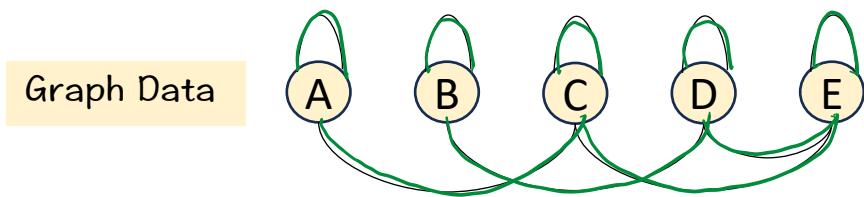
1	0	0	-2
0	1	0	-2
0	0	1	-5
1	-1	0	0
1	0	-1	0

[ReLU]

1	1	1	1	1	-9

σ

Graph Convolutional Network



Graph
Convolutional
Network

A	B	C	D	E
2	0	1	0	1
1	1	0	0	0
0	0	-1	1	1
0	3	0	1	0

1	1	0	0	0
0	1	0	-1	0
1	0	0	1	-1

[ReLU]

3	1	1	0	1
0	1	-2	0	0
1	2	0	0	0

Messages

A	B	C	D	E
1				
	1			
		1		
			1	
				1

4	1	5	2	1
0	1	0	1	0
1	2	1	2	0

0	1	1	0	
1	-1	0	-2	
1	0	0	0	

[ReLU]

1	3	1	3	0
2	-2	3	-1	-2
4	1	5	2	1

Messages

A	B	C	D	E
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1

1 1 1 1 1

Fully
Connected
Network

1	0	0	-2	
0	1	0	-2	
0	0	1	-5	
1	-1	0	0	
1	0	-1	0	

[ReLU]

1	1	1	1	1	-9
---	---	---	---	---	----

σ 0 1 1 1 0.5

Recurrent Neural Network (RNN)

Input Sequence

X	3	4	5	6
---	---	---	---	---

Parameters

$$A \begin{array}{|c|c|} \hline 1 & -1 \\ \hline 1 & 1 \\ \hline \end{array} \quad B \begin{array}{|c|c|} \hline 1 \\ \hline 2 \\ \hline \end{array} \quad C \begin{array}{|c|c|} \hline -1 & 1 \\ \hline \end{array}$$

Activation Function

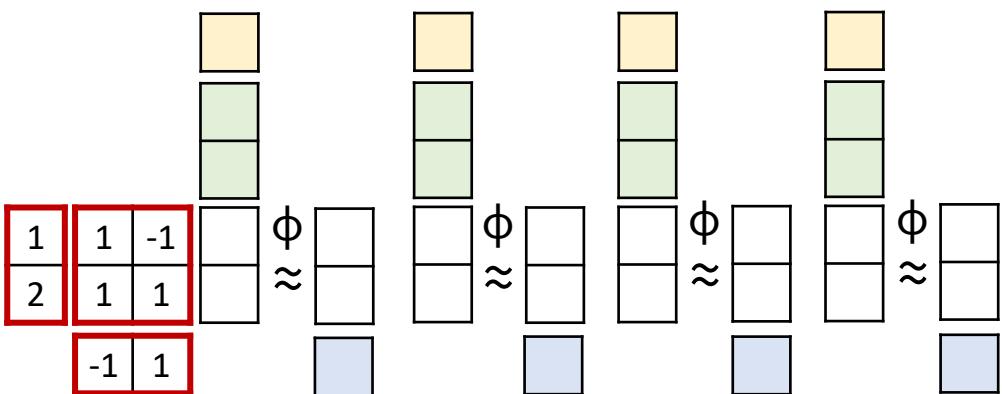
ϕ : ReLU

Hidden States

H_0	0
	0

Output Sequence

Y				



Recurrent Neural Network (RNN)

Input Sequence

X	3	4	5	6
---	---	---	---	---

Parameters

$$A \begin{array}{|c|c|} \hline 1 & -1 \\ \hline 1 & 1 \\ \hline \end{array} \quad B \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline \end{array} \quad C \begin{array}{|c|c|} \hline -1 & 1 \\ \hline \end{array}$$

Activation Function

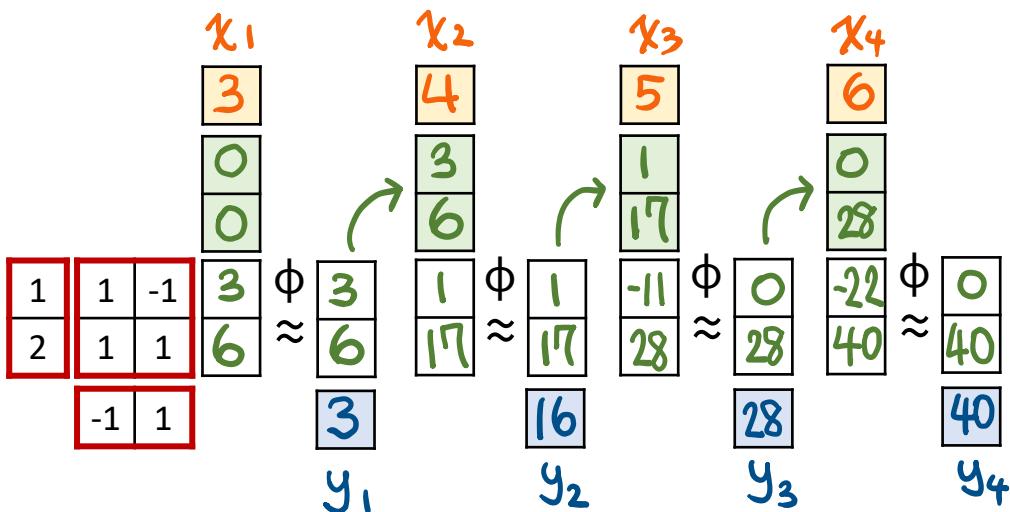
ϕ : ReLU

Hidden States

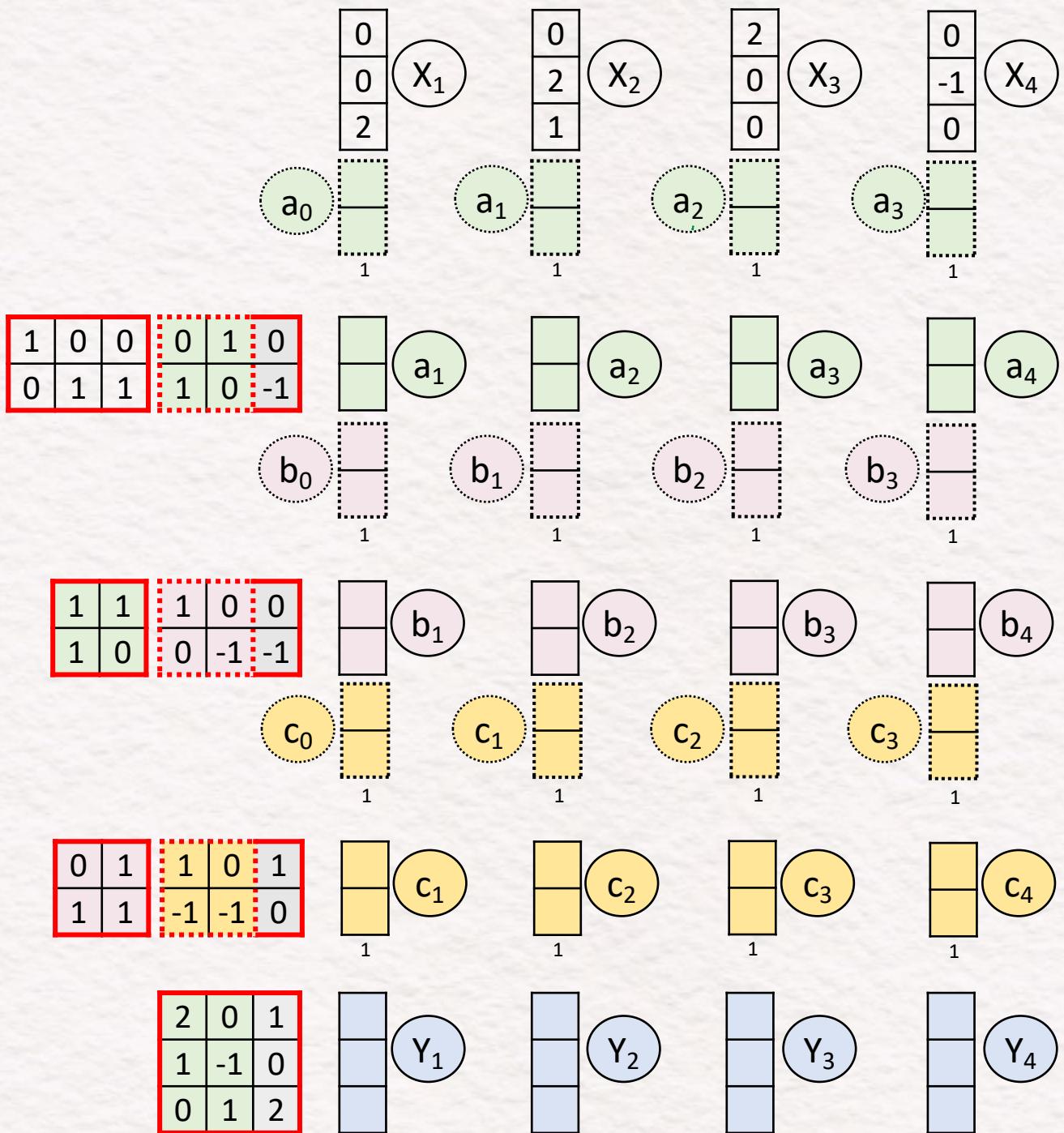
$$H_0$$

Output Sequence

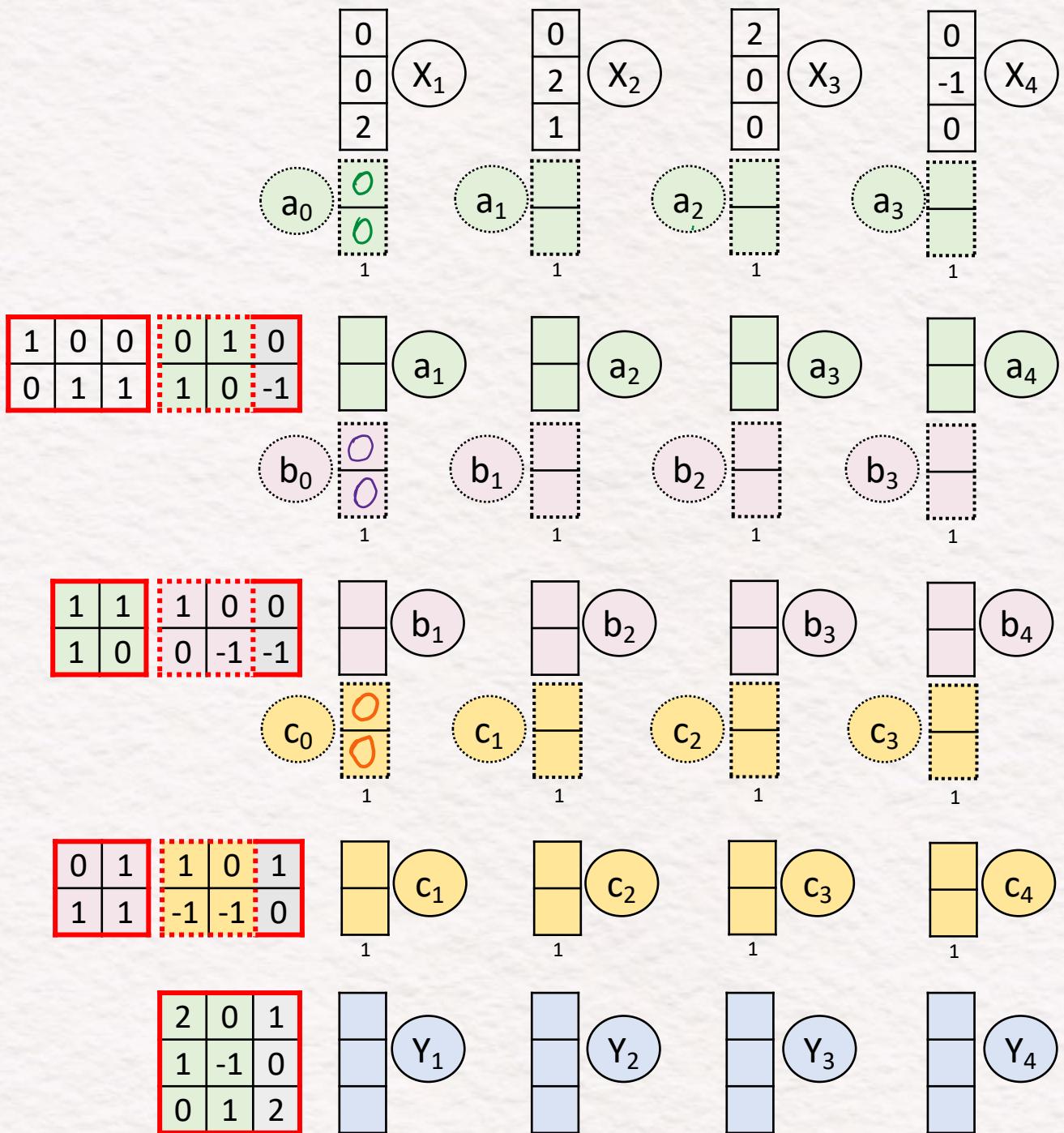
y | | | |



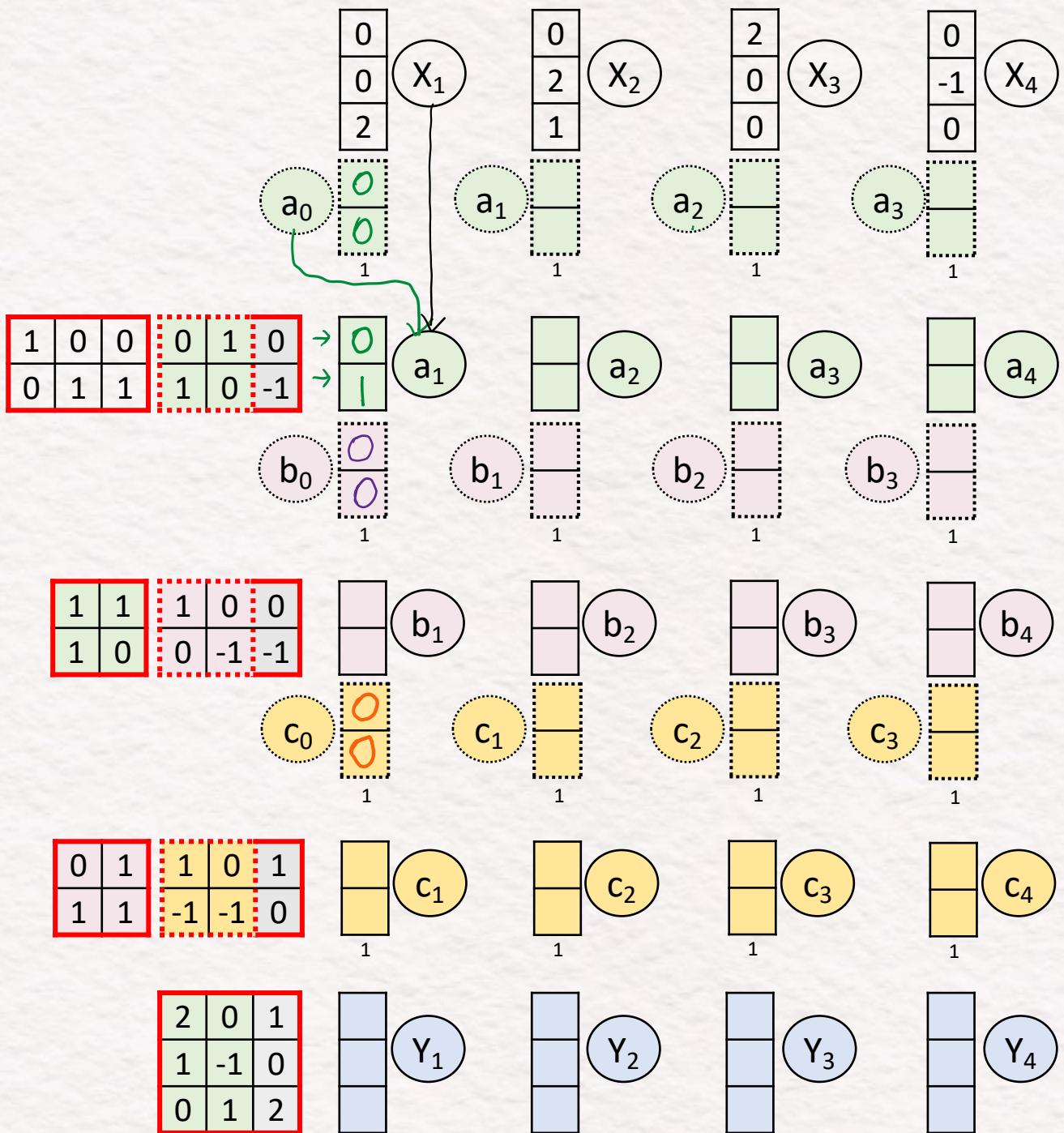
Deep RNN



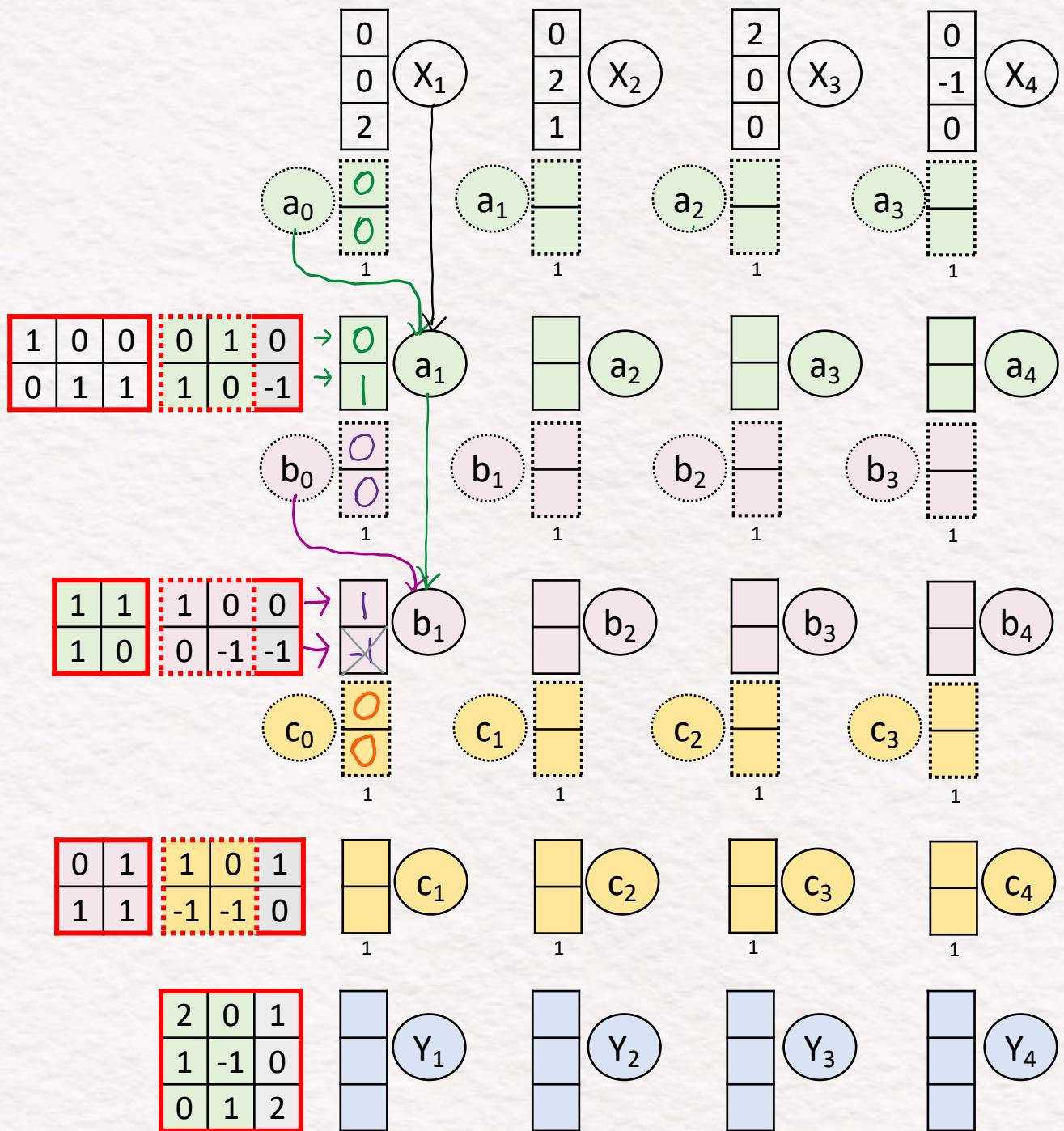
Deep RNN



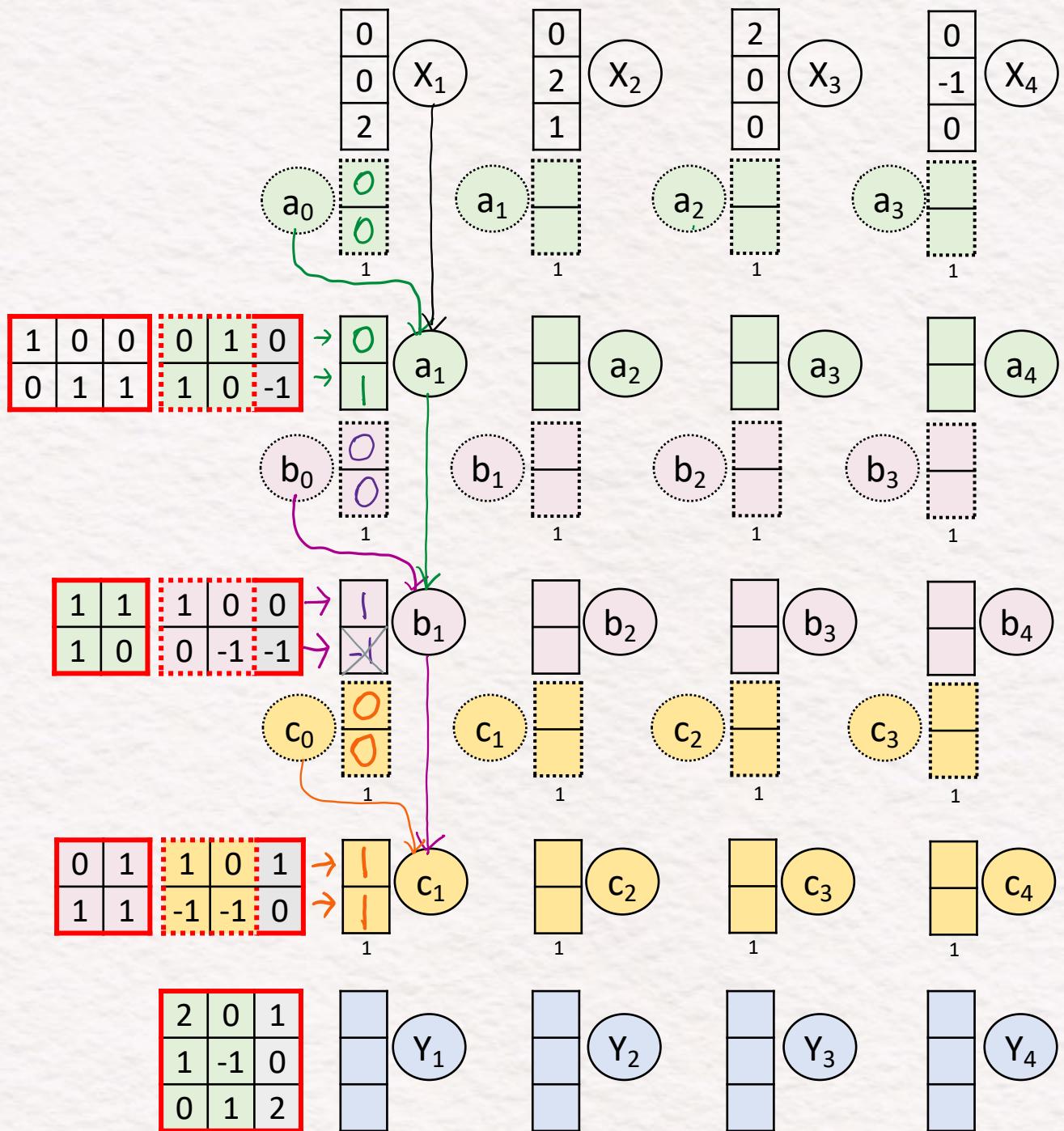
Deep RNN



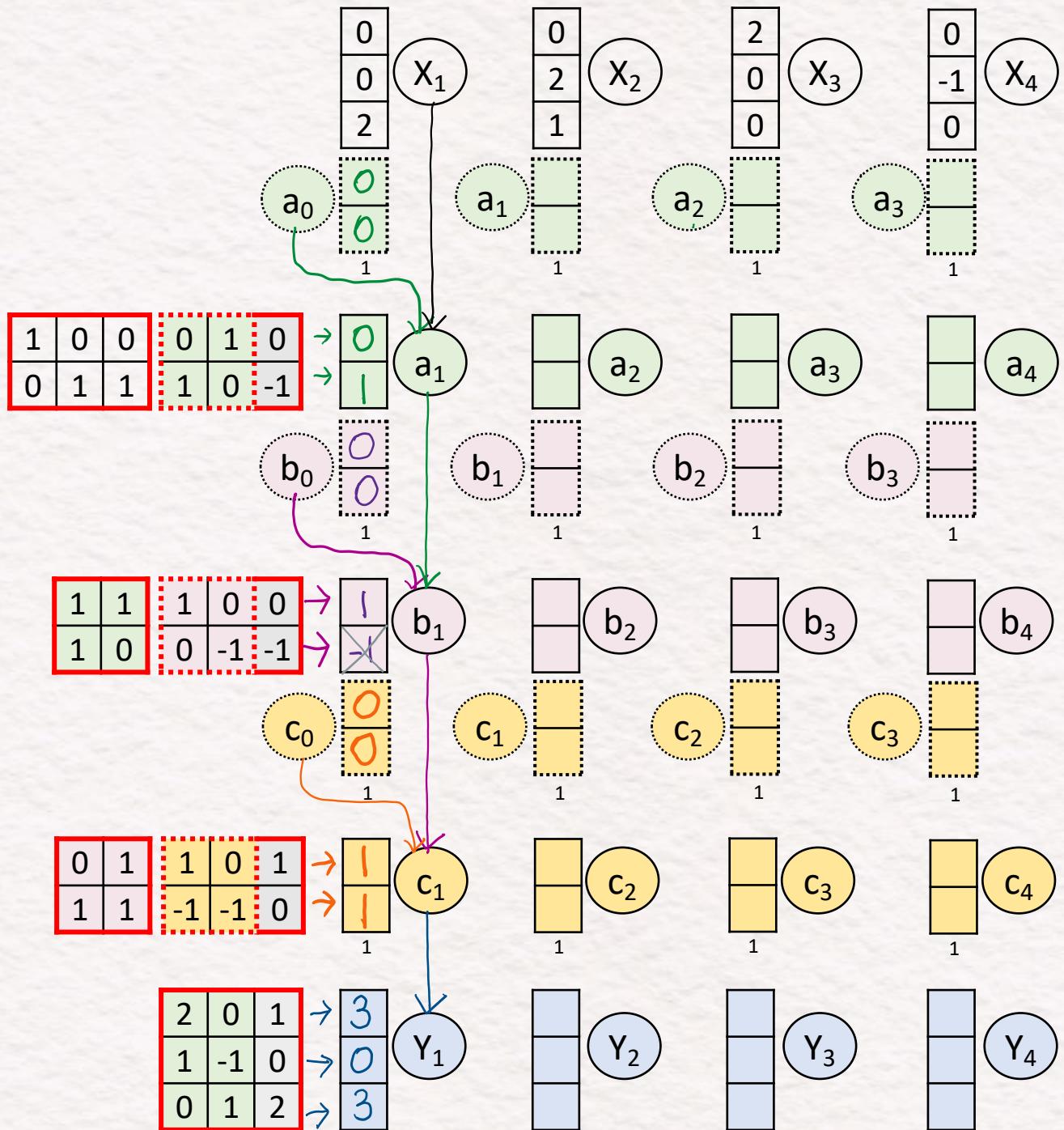
Deep RNN



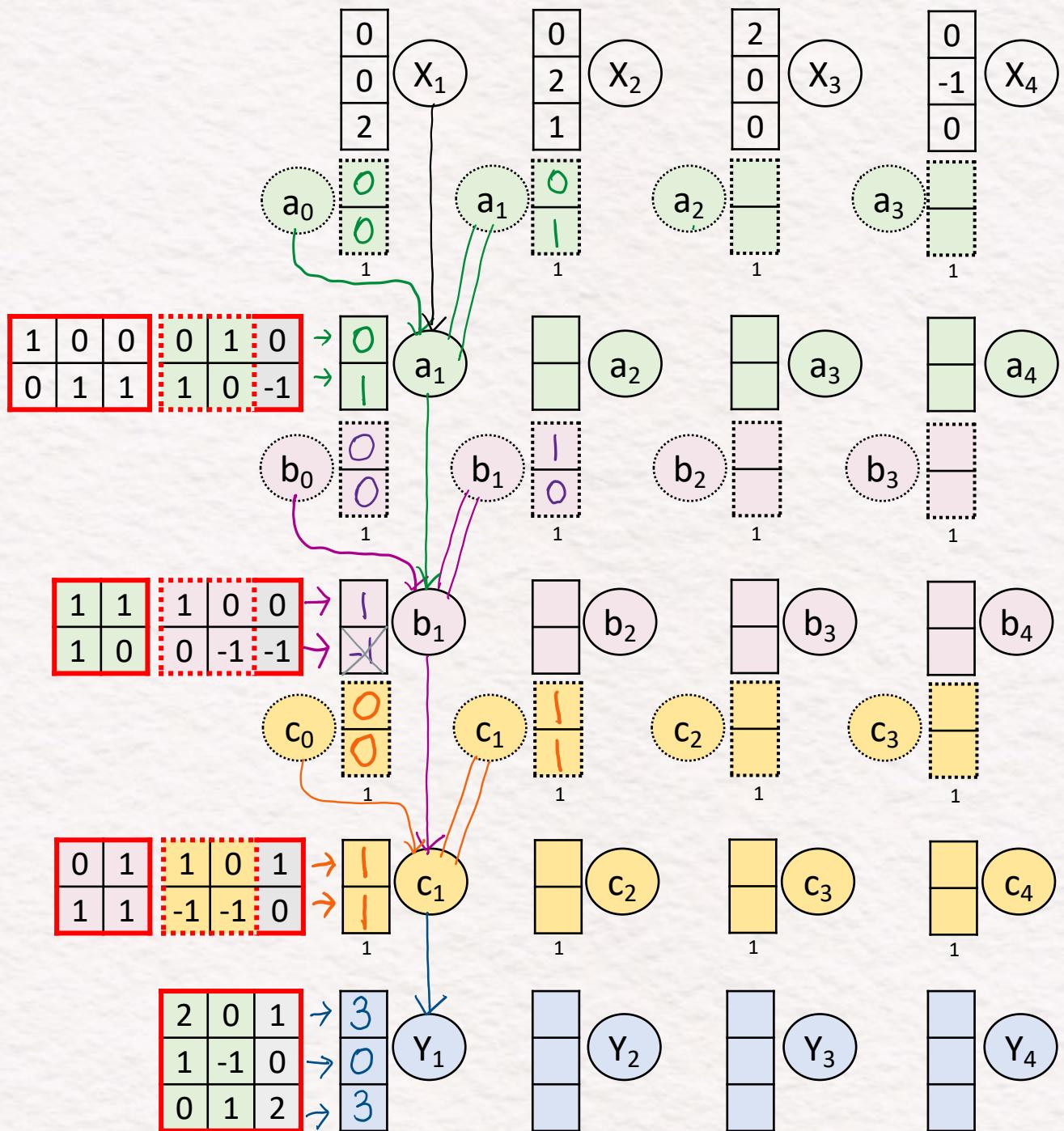
Deep RNN



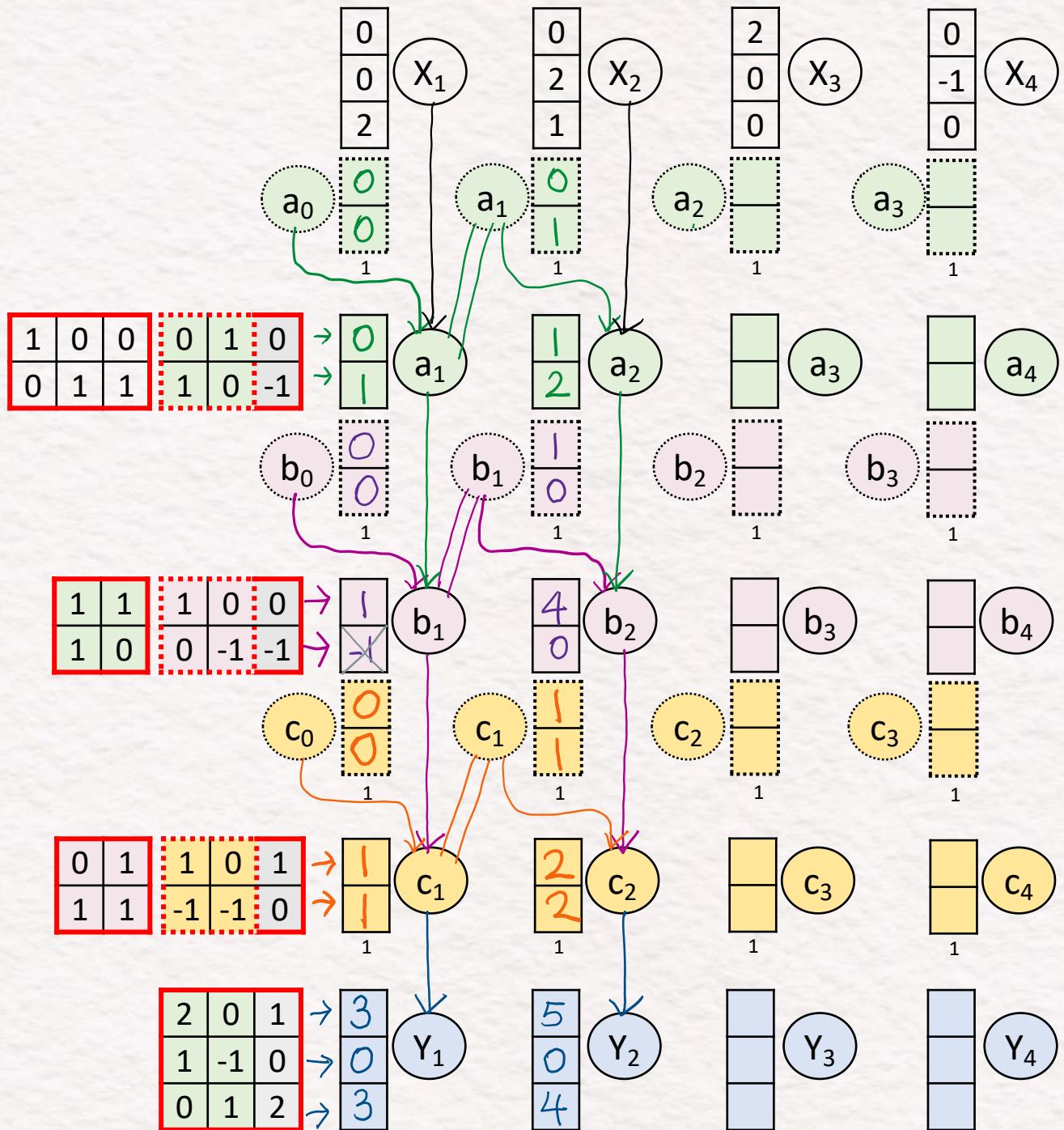
Deep RNN



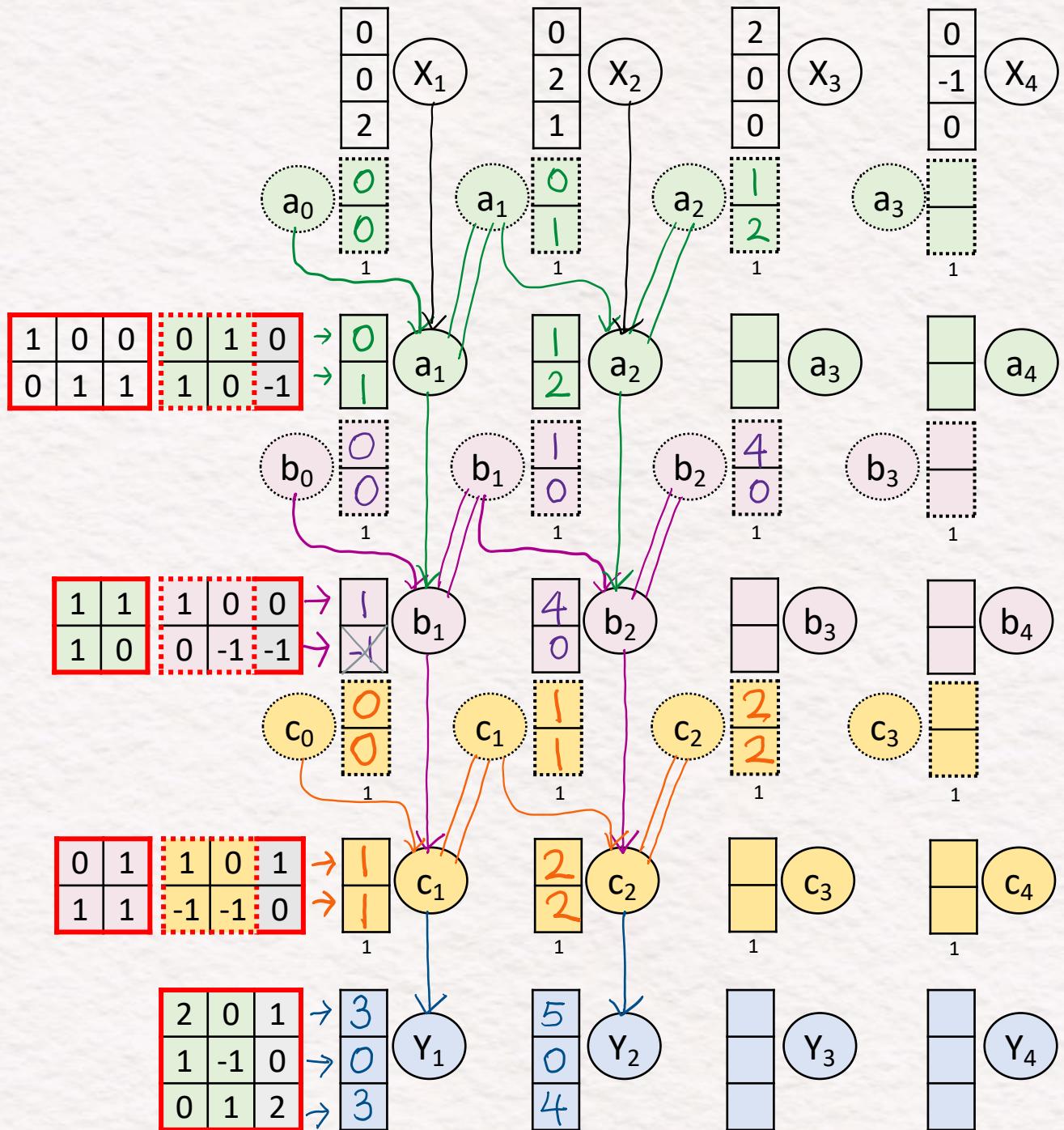
Deep RNN



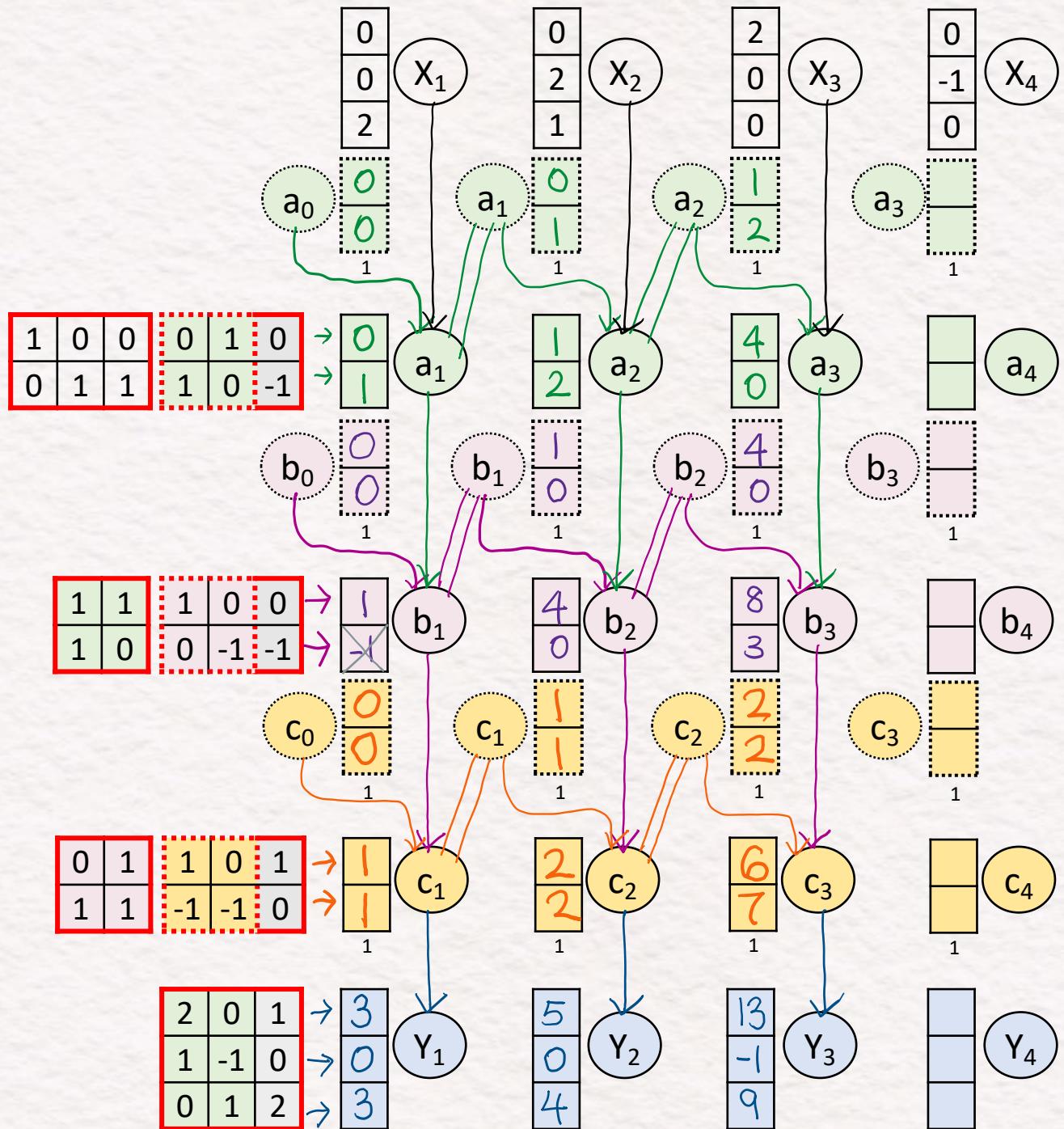
Deep RNN



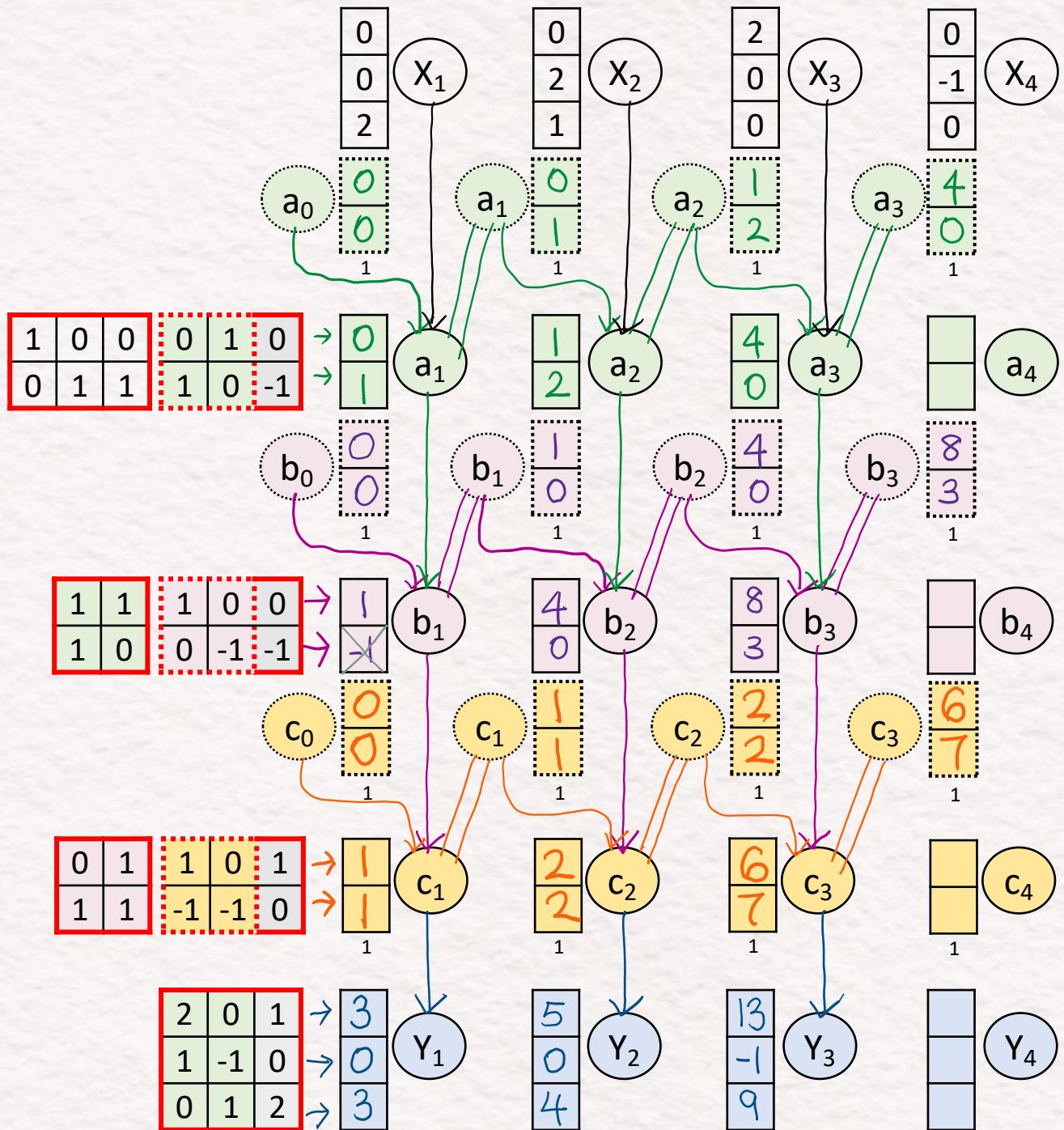
Deep RNN



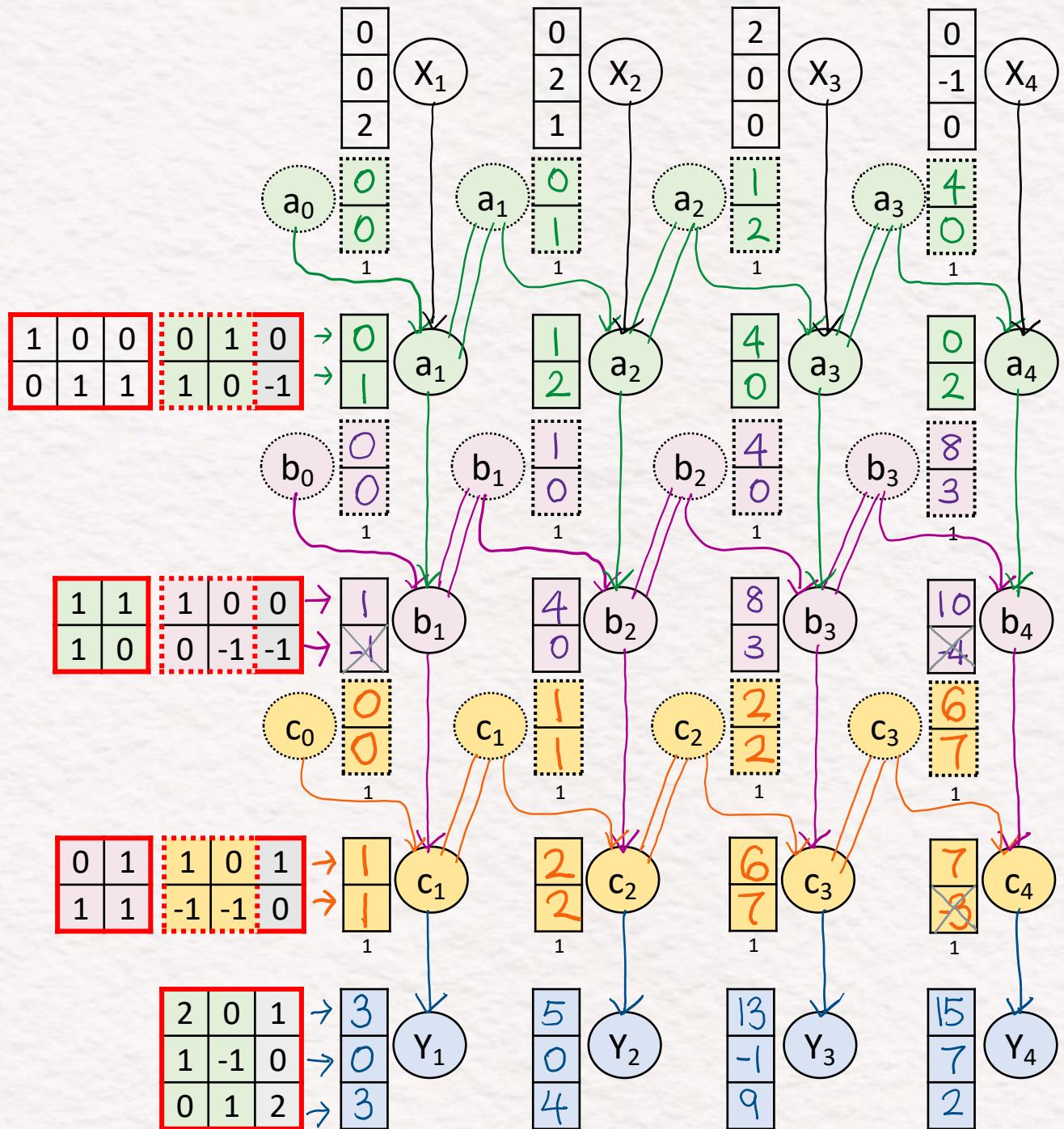
Deep RNN



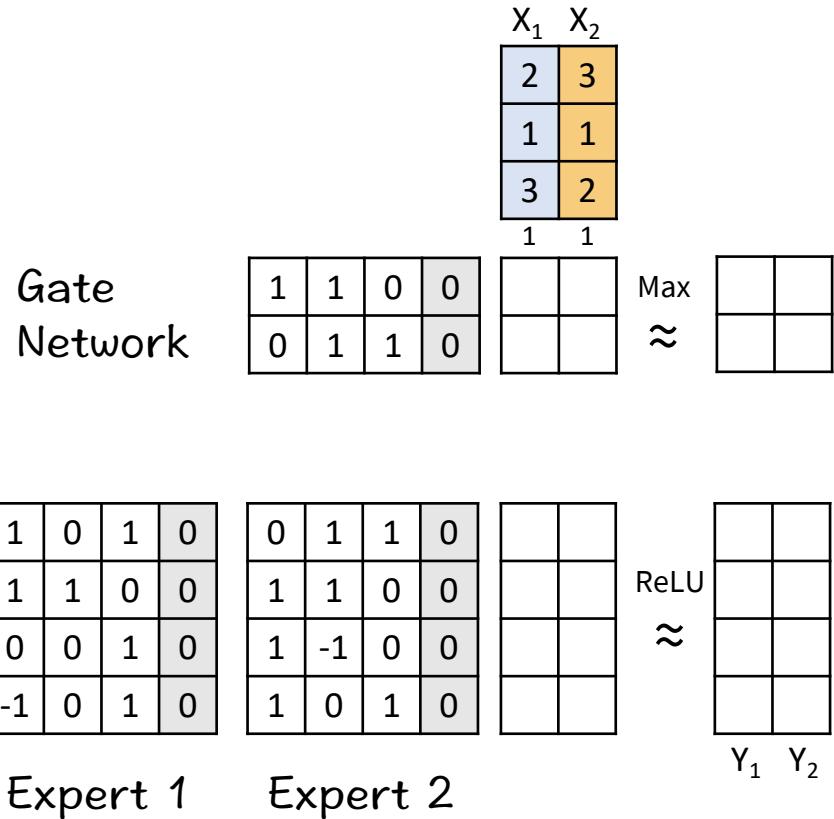
Deep RNN



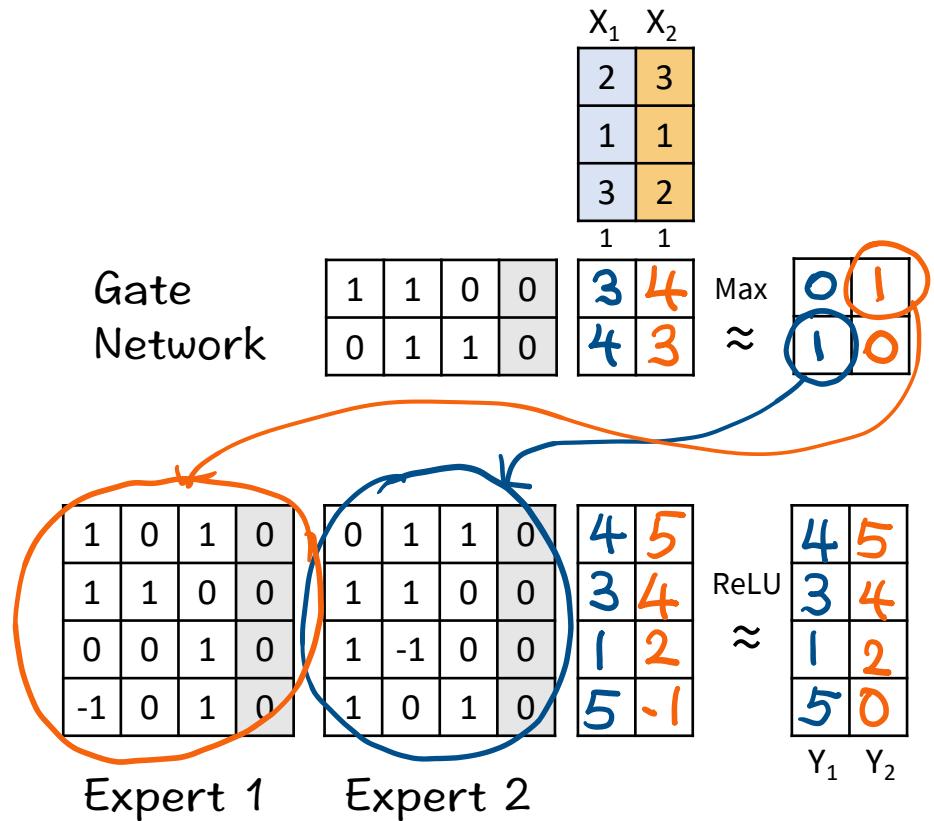
Deep RNN



Mixture of Experts



Mixture of Experts

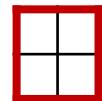


Mamba's S6 Model

Input Sequence

3	4	5	6
---	---	---	---

Parameters

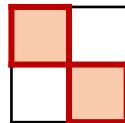


Output Sequence



Selective

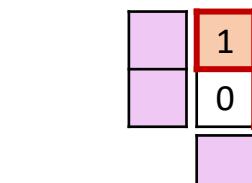
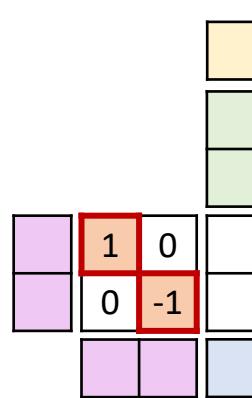
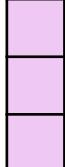
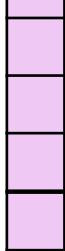
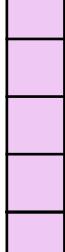
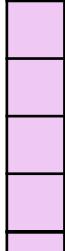
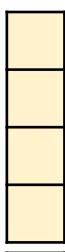
Structured



State-Space



1	-1	0	0
0	-1	0	1
1	0	-1	0
1	0	0	-1
1	0	-1	0
0	1	0	-1
1	-1	0	0
0	0	-1	1
-1	0	0	0
1	0	0	0
0	0	-1	0
0	1	0	0
1	-1	0	0
0	0	-1	1
1	0	0	0
0	-1	1	0



Scan

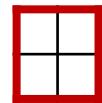


Mamba's S6 Model

Input Sequence

3	4	5	6
---	---	---	---

Parameters

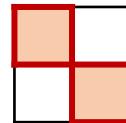


Output Sequence

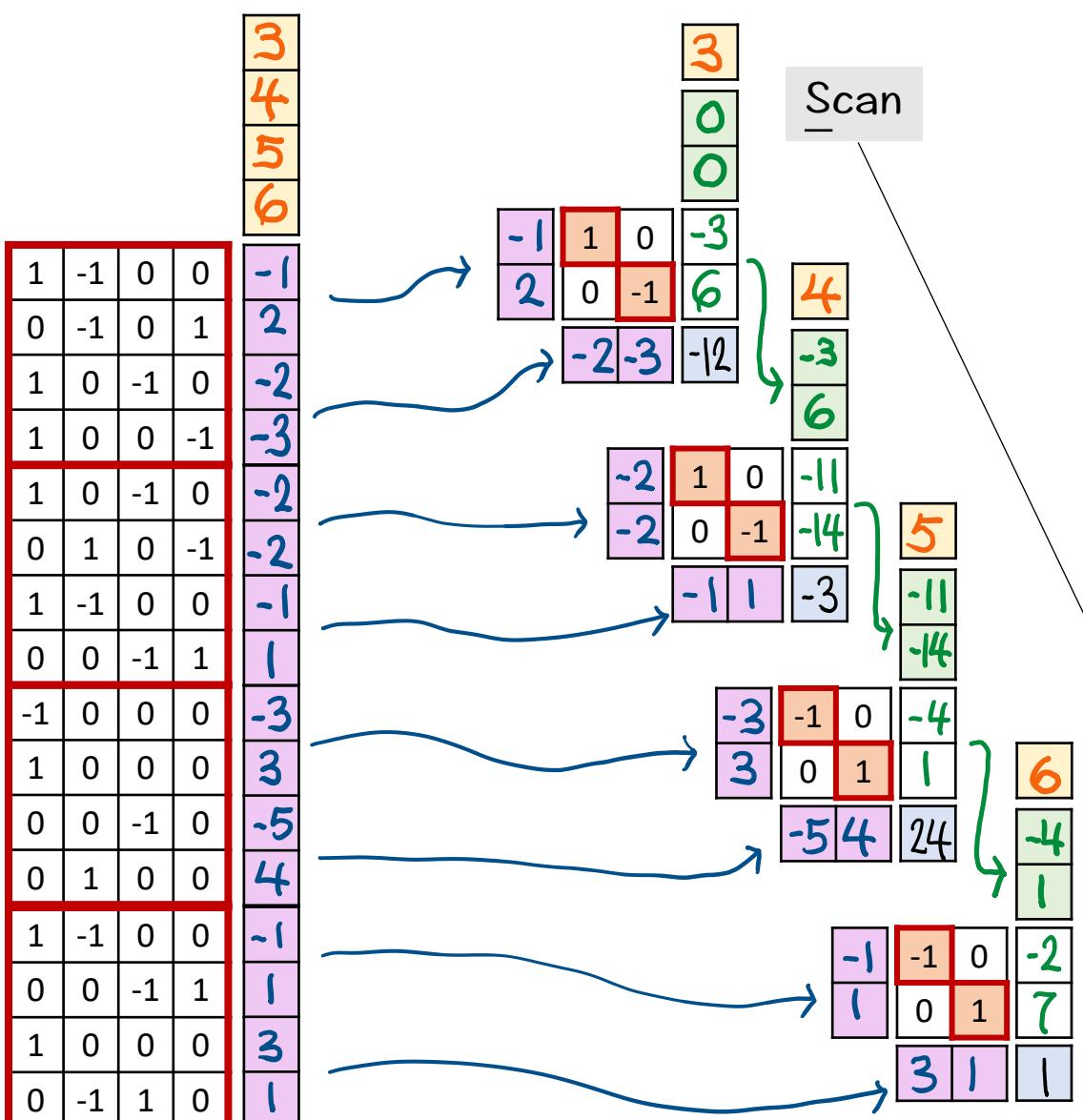
-1	2	-3	24	1
----	---	----	----	---

Selective

Structured



State-Space



Beginner's Guide to RAG

AI by Hand 

Prof. Tom Yeh

Hosted by:



University of Colorado
Boulder



Download Slides

<https://by-hand.ai/rag>

Roadmap

RAG Equation

Q/A

Prompt

Load

Advanced

Split

Index

Retrieval

Embed

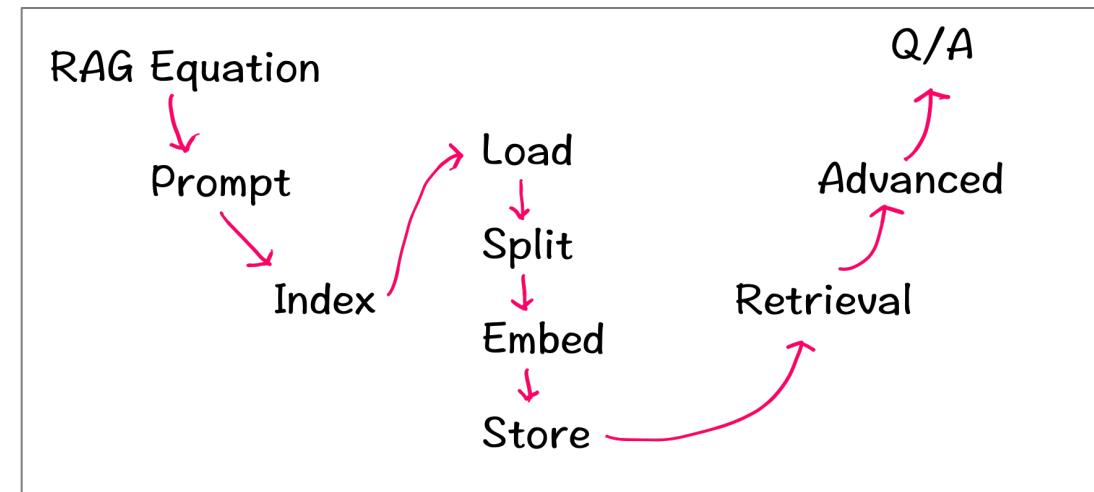
Store

RAG Equation

Beginner's Guide to RAG - AI by Hand 🖌



University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

RAG = Query + Prompt + Context + LLM

How much does
Tom make?

Search

Database

Amy, 100, CEO
Tom, 50, Prof
Sam, 120, CTO

[INST] <<SYS>> You are a [REDACTED], respectful, [REDACTED] assistant. Always answer as helpfully as possible, while being safe. Your answers should [REDACTED] include any harmful, [REDACTED], racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially [REDACTED] and positive in nature. <</SYS>>

Please answer the following question by using information from the provided context information!
CONTEXT_INFORMATION:

QUESTION:

[/INST]

LLM

RAG = Query + ~~Prompt~~ + Context + LLM

How much does
Tom make?

Search

[INST]

Please answer the following question by using
information from the provided context information!

CONTEXT_INFORMATION:

QUESTION:

[/INST]

Database

Amy, 100, CEO
Tom, 50, Prof
Sam, 120, CTO

LLM

~~RAG~~ = Query + ~~Prompt~~ + ~~Context~~ + LLM

How much does
Tom make?

Search

[INST]
QUESTION:

[/INST]

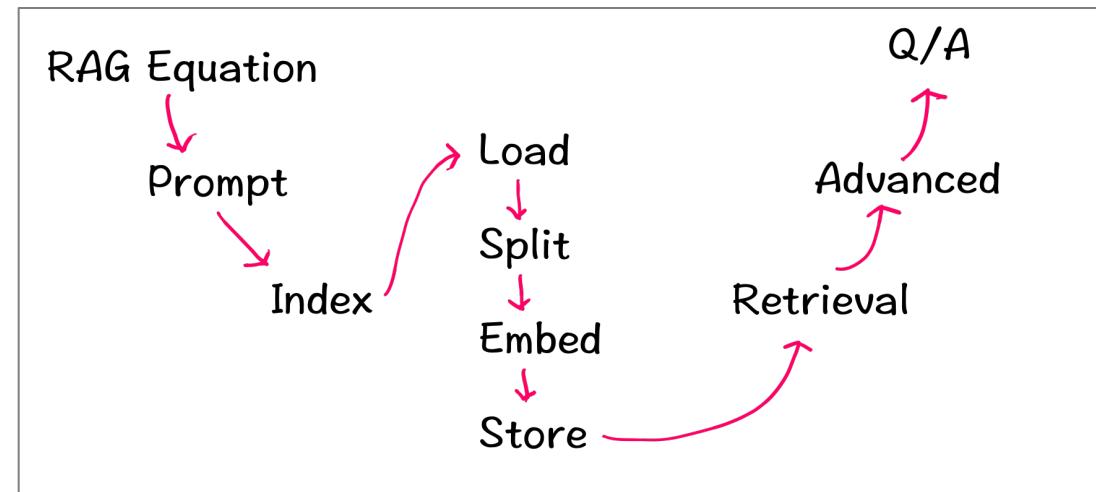
LLM

Prompt

Beginner's Guide to RAG - AI by Hand 🖊



University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

Summarize

Context information from [REDACTED] sources is below.

{context_str}

Given the information from multiple sources and not
[REDACTED], answer the query.

Query: {query_str}

Answer:

Single Choice

Some choices are given below. It is provided in a numbered list (1 to {num_chunks}), where each item in the list corresponds to a summary.

{context_list}

Using only the choices above and not prior knowledge, return the
 that is most relevant to the question:

{query_str}

Provide choice in the following format: 'ANSWER: <number>' and
 why this choice was selected in relation to the question.

Refine

The original query is as follows: {query_str}

We have provided an existing answer: {existing_answer}

We have the opportunity to [REDACTED] the existing answer
with some more context below.

{context_str}

Given the new context, refine the original answer to better answer
the query. If the context isn't useful, return the [REDACTED] answer.

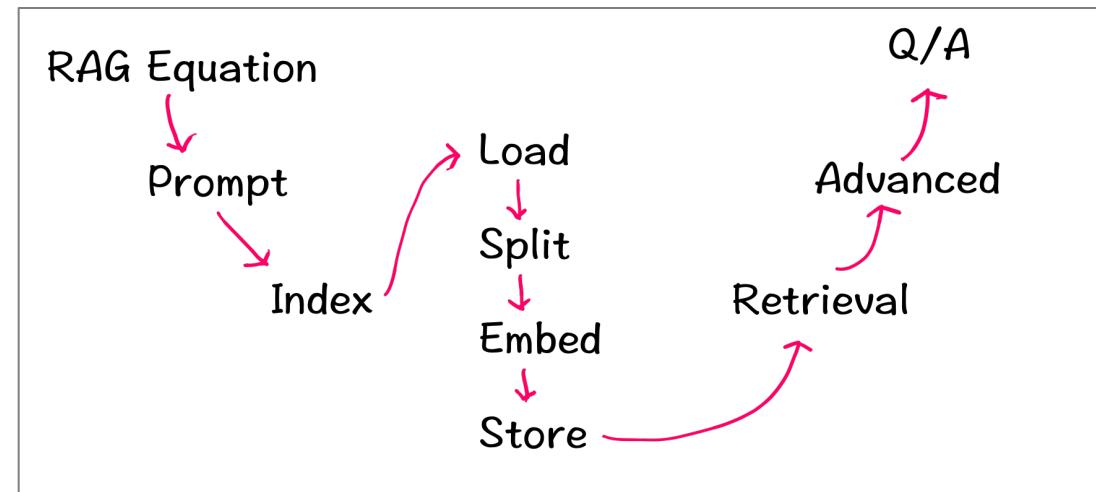
Refined Answer:

Index

Beginner's Guide to RAG - AI by Hand ✎



University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

Load

<small>id: 1</small>	<small>id: 2</small>	<small>id: 2</small>
<small>date: 5/10</small>	<small>date: 5/15</small>	<small>date: 5/17</small>
大人吃大包子，小孩吃小包子	每天都吃一個蘋果，可以保持健康	高山，大海，森林，動物，大自然

Split

	每天都	高山大
	吃一個	大海森
	蘋果可	林動物
	以保持	大自然

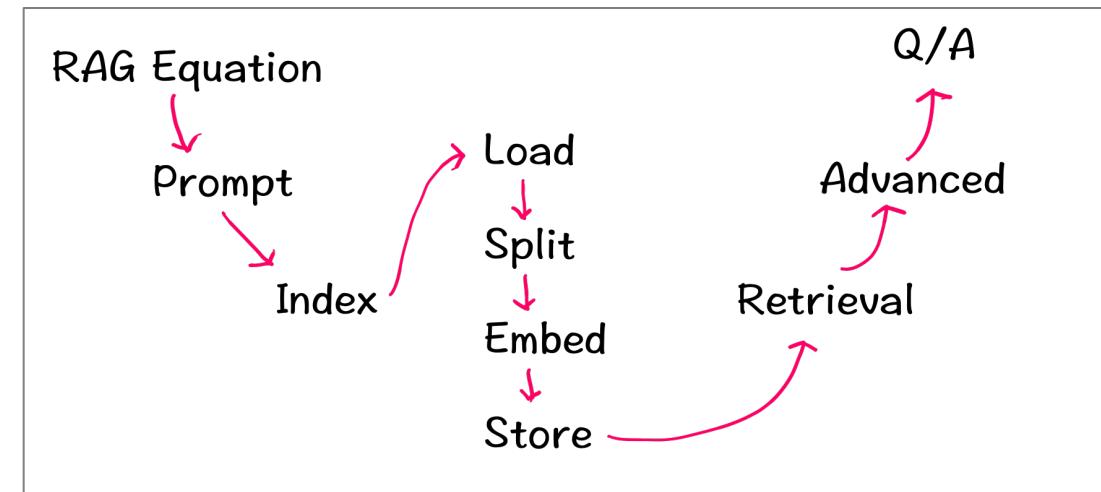
Embed

1	1	0	
1	0	1	
0	0	1	
0	1	0	

1	1	-2	
1	1	1	
0	0	1	
0	1	1	

Store

Chunk	Embed	Meta
		{id: ___, date: _____}
		{id: ___, date: _____}



Load

Beginner's Guide to RAG - AI by Hand ✎



University of Colorado
Boulder

Download: <https://by-hand.ai/rag>

LlamaHub Data Loaders

LlamaHub
Powered by Llamaindex

Github

CSVReader	Loaders	DocxReader	Verified Loaders	EpubReader	Loaders
llama-index	53 · 480868 · 5 days ago	thejessezheng	53 · 480868 · 5 days ago	haowjy	53 · 480868 · 5 days ago
FlatReader	Loaders	HTMLTagReader	Loaders	HWPReader	Loaders
llama-index	53 · 480868 · 5 days ago	llama-index	53 · 480868 · 5 days ago	sangwogenip	53 · 480868 · 5 days ago
IPYNBReader	Loaders	ImageCaptionReader	Loaders	ImageReader	Loaders
FarisHijazi	53 · 480868 · 5 days ago	FarisHijazi	53 · 480868 · 5 days ago	ravi03071991	53 · 480868 · 5 days ago
ImageTabularChartReader	Loaders	ImageVisionLLMReader	Loaders	MarkdownReader	Loaders
jon-chuang	53 · 480868 · 5 days ago	FarisHijazi	53 · 480868 · 5 days ago	hursh-desai	53 · 480868 · 5 days ago
MboxReader	Loaders	PDFReader	Loaders	PagedCSVReader	Verified Loaders
llama-index	53 · 480868 · 5 days ago	ravi03071991	53 · 480868 · 5 days ago	thejessezheng	53 · 480868 · 5 days ago



Llamaindex

Which loader to use to read CSV data?

LlamaHub
Powered by Llamaindex

Github

CSVReader	Loaders
llama-index	⭐ 53 · 📁 480868 · 5 days ago
DocxReader	Verified Loaders
thejessezhang	⭐ 53 · 📁 480868 · 5 days ago
EpubReader	Loaders
haowjy	⭐ 53 · 📁 480868 · 5 days ago
FlatReader	Loaders
llama-index	⭐ 53 · 📁 480868 · 5 days ago
HTMLTagReader	Loaders
llama-index	⭐ 53 · 📁 480868 · 5 days ago
HWPReader	Loaders
sangwogenip	⭐ 53 · 📁 480868 · 5 days ago
IPYNBReader	Loaders
FarisHijazi	⭐ 53 · 📁 480868 · 5 days ago
ImageCaptionReader	Loaders
FarisHijazi	⭐ 53 · 📁 480868 · 5 days ago
ImageReader	Loaders
ravi03071991	⭐ 53 · 📁 480868 · 5 days ago
ImageTabularChartReader	Loaders
jon-chuang	⭐ 53 · 📁 480868 · 5 days ago
ImageVisionLLMReader	Loaders
FarisHijazi	⭐ 53 · 📁 480868 · 5 days ago
MarkdownReader	Loaders
hursh-desai	⭐ 53 · 📁 480868 · 5 days ago
MboxReader	Loaders
llama-index	⭐ 53 · 📁 480868 · 5 days ago
PDFReader	Loaders
ravi03071991	⭐ 53 · 📁 480868 · 5 days ago
PagedCSVReader	Verified Loaders
thejessezhang	⭐ 53 · 📁 480868 · 5 days ago



Llamaindex

Which loader to use to read Word documents?

LlamaHub
Powered by Llamaindex

Github

CSVReader	Loaders
llama-index	53 · 480868 · 5 days ago
DocxReader	Verified Loaders
thejessezhang	53 · 480868 · 5 days ago
EpubReader	Loaders
haowjy	53 · 480868 · 5 days ago
FlatReader	Loaders
llama-index	53 · 480868 · 5 days ago
HTMLTagReader	Loaders
llama-index	53 · 480868 · 5 days ago
HWPReader	Loaders
sangwogenip	53 · 480868 · 5 days ago
IPYNBReader	Loaders
FarisHijazi	53 · 480868 · 5 days ago
ImageCaptionReader	Loaders
FarisHijazi	53 · 480868 · 5 days ago
ImageReader	Loaders
ravi03071991	53 · 480868 · 5 days ago
ImageTabularChartReader	Loaders
jon-chuang	53 · 480868 · 5 days ago
ImageVisionLLMReader	Loaders
FarisHijazi	53 · 480868 · 5 days ago
MarkdownReader	Loaders
hursh-desai	53 · 480868 · 5 days ago
MboxReader	Loaders
llama-index	53 · 480868 · 5 days ago
PDFReader	Loaders
ravi03071991	53 · 480868 · 5 days ago
PagedCSVReader	Verified Loaders
thejessezhang	53 · 480868 · 5 days ago



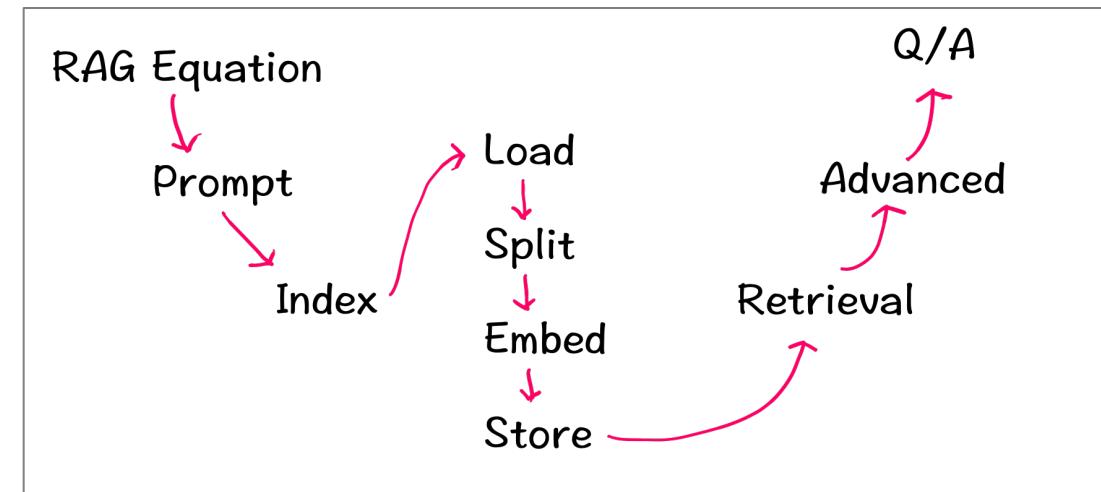
Llamaindex

Split

Beginner's Guide to RAG - AI by Hand ✎



University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

Split by Tokens

```
TokenTextSplitter(  
    chunk_size=4,  
    chunk_overlap=0  
)
```

How are you doing today? I hope you're doing well. Whether your day is busy or calm, I hope you find a moment of joy. Life moves fast, with many ups and downs. It's important to stop and enjoy the little things. Maybe today brought a good chat, a tasty meal, or a beautiful sunset. Cherish these happy moments. Each day is a chance to make good memories and connect with others. Take time to appreciate your surroundings and the people in your life. Even on challenging days, there are opportunities to learn and grow. Let's embrace each experience and look forward to tomorrow with optimism.

Split by Tokens

```
TokenTextSplitter(  
    chunk_size=4,  
    chunk_overlap=1  
)
```

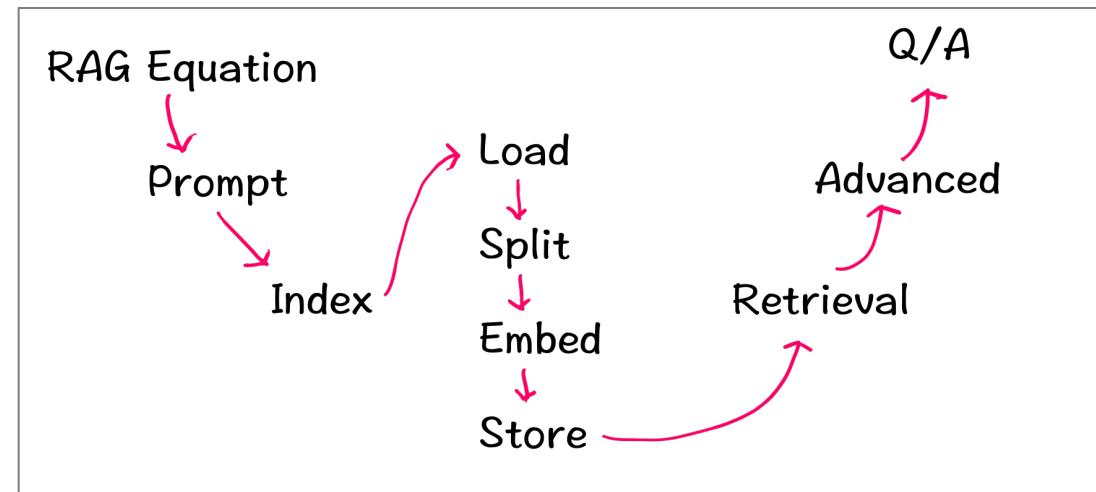
How are you doing today? I hope you're doing well. Whether your day is busy or calm, I hope you find a moment of joy. Life moves fast, with many ups and downs. It's important to stop and enjoy the little things. Maybe today brought a good chat, a tasty meal, or a beautiful sunset. Cherish these happy moments. Each day is a chance to make good memories and connect with others. Take time to appreciate your surroundings and the people in your life. Even on challenging days, there are opportunities to learn and grow. Let's embrace each experience and look forward to tomorrow with optimism.

Embed

Beginner's Guide to RAG - AI by Hand 🖊



University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

Massive Text Embedding Benchmark

<https://huggingface.co/spaces/mteb/leaderboard>

Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
1	voyage-large-2-instruct			1024	16000	68.28	81.49	53.35
2	SFR-Embedding-Mistral	7111	26.49	4096	32768	67.56	78.33	51.67
3	gte-Qwen1.5-7B-instruct	7099	26.45	4096	32768	67.34	79.6	55.83
4	voyage-lite-02-instruct	1220	4.54	1024	4000	67.13	79.25	52.42
5	GritLM-7B	7242	26.98	4096	32768	66.76	79.46	50.61
6	e5-mistral-7b-instruct	7111	26.49	4096	32768	66.63	78.47	50.26
7	google-gecko.text-embedding-f	1200	4.47	768	2048	66.31	81.17	47.48
8	GritLM-8x7B	46703	173.98	4096	32768	65.66	78.53	50.14
9	gte-large-en-v1.5	434	1.62	1024	8192	65.39	77.75	47.96
10	LLM2Vec-Meta-Llama-3-supervis	7505	27.96	4096	8192	65.01	75.92	46.45

Which models have 7B parameters?

Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
1	<u>voyage-large-2-instruct</u>			1024	16000	68.28	81.49	53.35
2	<u>SFR-Embedding-Mistral</u>	7111	26.49	4096	32768	67.56	78.33	51.67
3	<u>gte-Qwen1.5-7B-instruct</u>	7099	26.45	4096	32768	67.34	79.6	55.83
4	<u>voyage-lite-02-instruct</u>	1220	4.54	1024	4000	67.13	79.25	52.42
5	<u>GritLM-7B</u>	7242	26.98	4096	32768	66.76	79.46	50.61
6	<u>e5-mistral-7b-instruct</u>	7111	26.49	4096	32768	66.63	78.47	50.26
7	<u>google-gecko.text-embedding-p</u>	1200	4.47	768	2048	66.31	81.17	47.48
8	<u>GritLM-8x7B</u>	46703	173.98	4096	32768	65.66	78.53	50.14
9	<u>gte-large-en-v1.5</u>	434	1.62	1024	8192	65.39	77.75	47.96
10	<u>LLM2Vec-Meta-Llama-3-supervis</u>	7505	27.96	4096	8192	65.01	75.92	46.45

Which model has the most parameters?

Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
1	<u>voyage-large-2-instruct</u>			1024	16000	68.28	81.49	53.35
2	<u>SFR-Embedding-Mistral</u>	7111	26.49	4096	32768	67.56	78.33	51.67
3	<u>gte-Qwen1.5-7B-instruct</u>	7099	26.45	4096	32768	67.34	79.6	55.83
4	<u>voyage-lite-02-instruct</u>	1220	4.54	1024	4000	67.13	79.25	52.42
5	<u>GritLM-7B</u>	7242	26.98	4096	32768	66.76	79.46	50.61
6	<u>e5-mistral-7b-instruct</u>	7111	26.49	4096	32768	66.63	78.47	50.26
7	<u>google-gecko.text-embedding-p</u>	1200	4.47	768	2048	66.31	81.17	47.48
8	<u>GritLM-8x7B</u>	46703	173.98	4096	32768	65.66	78.53	50.14
9	<u>gte-large-en-v1.5</u>	434	1.62	1024	8192	65.39	77.75	47.96
10	<u>LLM2Vec-Meta-Llama-3-supervis</u>	7505	27.96	4096	8192	65.01	75.92	46.45

Which model supports the longest context window?

Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
1	<u>voyage-large-2-instruct</u>			1024	16000	68.28	81.49	53.35
2	<u>SFR-Embedding-Mistral</u>	7111	26.49	4096	32768	67.56	78.33	51.67
3	<u>gte-Qwen1.5-7B-instruct</u>	7099	26.45	4096	32768	67.34	79.6	55.83
4	<u>voyage-lite-02-instruct</u>	1220	4.54	1024	4000	67.13	79.25	52.42
5	<u>GritLM-7B</u>	7242	26.98	4096	32768	66.76	79.46	50.61
6	<u>e5-mistral-7b-instruct</u>	7111	26.49	4096	32768	66.63	78.47	50.26
7	<u>google-gecko.text-embedding-p</u>	1200	4.47	768	2048	66.31	81.17	47.48
8	<u>GritLM-8x7B</u>	46703	173.98	4096	32768	65.66	78.53	50.14
9	<u>gte-large-en-v1.5</u>	434	1.62	1024	8192	65.39	77.75	47.96
10	<u>LLM2Vec-Meta-Llama-3-supervis</u>	7505	27.96	4096	8192	65.01	75.92	46.45

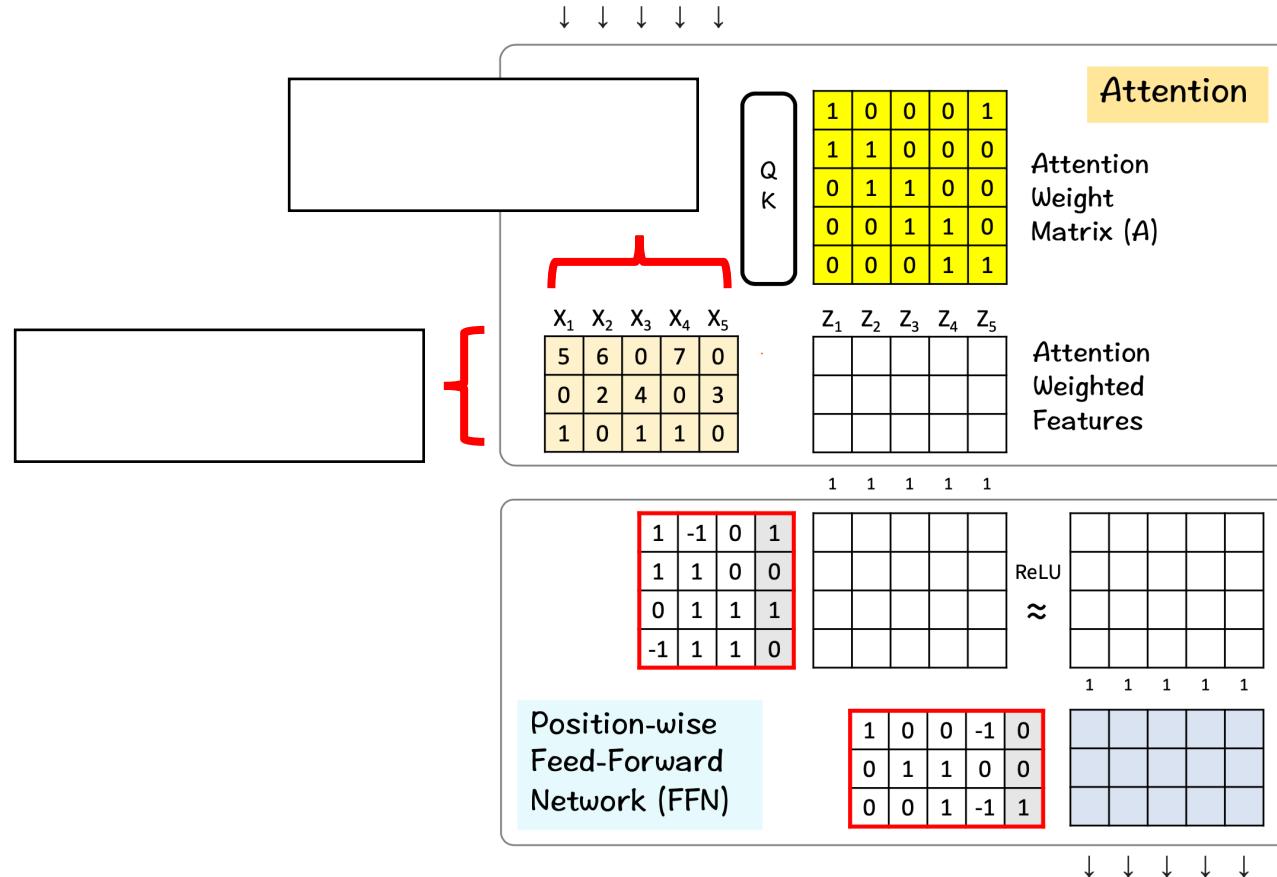
Which model has the best classification performance?

Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
1	<u>voyage-large-2-instruct</u>			1024	16000	68.28	81.49	53.35
2	<u>SFR-Embedding-Mistral</u>	7111	26.49	4096	32768	67.56	78.33	51.67
3	<u>gte-Qwen1.5-7B-instruct</u>	7099	26.45	4096	32768	67.34	79.6	55.83
4	<u>voyage-lite-02-instruct</u>	1220	4.54	1024	4000	67.13	79.25	52.42
5	<u>GritLM-7B</u>	7242	26.98	4096	32768	66.76	79.46	50.61
6	<u>e5-mistral-7b-instruct</u>	7111	26.49	4096	32768	66.63	78.47	50.26
7	<u>google-gecko.text-embedding-p</u>	1200	4.47	768	2048	66.31	81.17	47.48
8	<u>GritLM-8x7B</u>	46703	173.98	4096	32768	65.66	78.53	50.14
9	<u>gte-large-en-v1.5</u>	434	1.62	1024	8192	65.39	77.75	47.96
10	<u>LLM2Vec-Meta-Llama-3-supervis</u>	7505	27.96	4096	8192	65.01	75.92	46.45

Which model has the best clustering performance?

Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
1	<u>voyage-large-2-instruct</u>			1024	16000	68.28	81.49	53.35
2	<u>SFR-Embedding-Mistral</u>	7111	26.49	4096	32768	67.56	78.33	51.67
3	<u>gte-Qwen1.5-7B-instruct</u>	7099	26.45	4096	32768	67.34	79.6	55.83
4	<u>voyage-lite-02-instruct</u>	1220	4.54	1024	4000	67.13	79.25	52.42
5	<u>GritLM-7B</u>	7242	26.98	4096	32768	66.76	79.46	50.61
6	<u>e5-mistral-7b-instruct</u>	7111	26.49	4096	32768	66.63	78.47	50.26
7	<u>google-gecko.text-embedding-p</u>	1200	4.47	768	2048	66.31	81.17	47.48
8	<u>GritLM-8x7B</u>	46703	173.98	4096	32768	65.66	78.53	50.14
9	<u>gte-large-en-v1.5</u>	434	1.62	1024	8192	65.39	77.75	47.96
10	<u>LLM2Vec-Meta-Llama-3-supervis</u>	7505	27.96	4096	8192	65.01	75.92	46.45

Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)
text-embedding-3-large			3072	8191	64.59	75.45	49.01

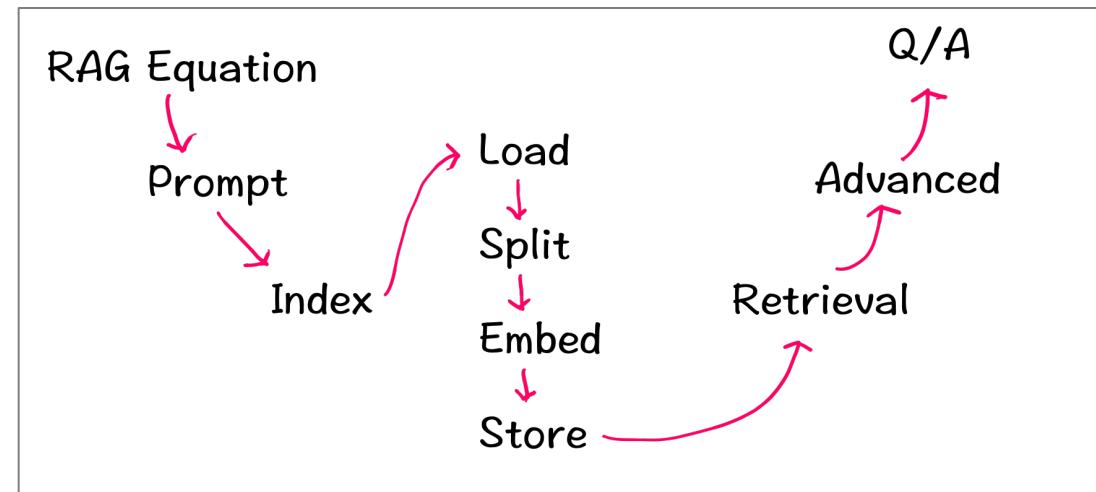


Store

Beginner's Guide to RAG - AI by Hand 🖊



University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

Choose a Vector Database

<https://superlinked.com/vector-db-comparison>

Vector DB Comparison

by Superlinked | Last Updated : Today

Search

Get insights

Give us a star

Settings

Vendor	Ops										Stats	
	Multi-Tenant	Disk Index	Ephemeral Index	Sharding	Document Size	Vector Dims	Int8 Quantiza...	Binary Quantiza...	GitHub ⭐	Doc P...		
Activeloo...	✓ ⓘ ⓘ	✓ ⓘ ⓘ	-	✖ ⓘ	Unlimited	Unlimited	-	-	7.7K	-		
Anari AI	-	-	✖	-	Unlimited	-	-	-	-	-		
Apache C...	✓ ⓘ ⓘ	✓ ⓘ ⓘ	-	-	-	Unlimited	-	-	8.5K	20		
Apache S...	✓	-	✖	✓	Unlimited	Unlimited	-	-	1.0K	27		
ApertureDB	-	-	-	-	-	Unlimited	-	-	-	-		
Azure AI S...	✓ ⓘ	✖	✖	✓ ⓘ	① 16000000	3072	✓ ⓘ	-	-	-		
Chroma	✖	✖	✓	-	-	Unlimited	-	-	12.4K	-		
ClickHouse	✓ ⓘ ⓘ	✓ ⓘ ⓘ	✖	✓ ⓘ	Unlimited	Unlimited	-	-	34.3K	-		
CrateDB	✓ ⓘ ⓘ	✓ ⓘ ⓘ	-	✓ ⓘ	-	2048	-	-	4.0K	18		
DataStax ...	✓	✓ ⓘ ⓘ	✓	✓	250000000	8192	-	-	-	-		
Elasticsea...	✓	✓	-	✓	100000000	4096	✓ ⓘ	-	67.7K	84		
Epsilla	✓ ⓘ ⓘ	-	-	✓	Unlimited	Unlimited	-	-	872	-		
GCP Verte...	-	-	-	-	-	Unlimited	-	-	-	-		
KDB.AI	-	✖	-	-	Unlimited	Unlimited	-	-	-	-		
LanceDB	✖	✓ ⓘ ⓘ	-	-	-	2048	-	-	2.9K	-		



Which database can fit 3072 dimensions?

Vector DB Comparison

by Superlinked | Last Updated : Today

Search [Get insights](#) [Give us a star](#) [Settings](#)

Vendor	Ops										Stats	
	Multi-Tenant	Disk Index	Ephemeral Index	Sharding	Document Size	Vector Dims	Int8 Quantiza...	Binary Quantiza...	GitHub ⭐	Doc P...		
Activeloo...	✓ ⓘ 🔗	✓ ⓘ 🔗	-	⚠ 🔗	Unlimited	Unlimited	-	-	7.7K	-		
Anari AI	-	-	✗	-	Unlimited	-	-	-	-	-		
Apache C...	✓ ⓘ 🔗	✓ ⓘ 🔗	-	-	-	Unlimited	-	-	8.5K	20		
Apache S...	✓	-	✗	✓	Unlimited	Unlimited	-	-	1.0K	27		
ApertureDB	-	-	-	-	-	Unlimited	-	-	-	-		
Azure AI S...	✓ 🔗	✗	✗	✓ 🔗	① 16000000	3072	✓ 🔗	-	-	-		
Chroma	✗	✗	✓	-	-	Unlimited	-	-	12.4K	-		
ClickHouse	✓ ⓘ 🔗	✓ ⓘ 🔗	✗	✓ ⓘ	Unlimited	Unlimited	-	-	34.3K	-		
CrateDB	✓ ⓘ 🔗	✓ ⓘ 🔗	-	✓ 🔗	-	2048	-	-	4.0K	18		
DataStax ...	✓	✓ ⓘ 🔗	✓	✓	250000000	8192	-	-	-	-		
Elasticsea...	✓	✓	-	✓	100000000	4096	✓ 🔗	-	67.7K	84		
Epsilla	✓ ⓘ 🔗	-	-	✓	Unlimited	Unlimited	-	-	872	-		
GCP Verte...	-	-	-	-	-	Unlimited	-	-	-	-		
KDB.AI	-	⚠	-	-	Unlimited	Unlimited	-	-	-	-		
LanceDB	✗	✓ ⓘ 🔗	-	-	-	2048	-	-	2.9K	-		

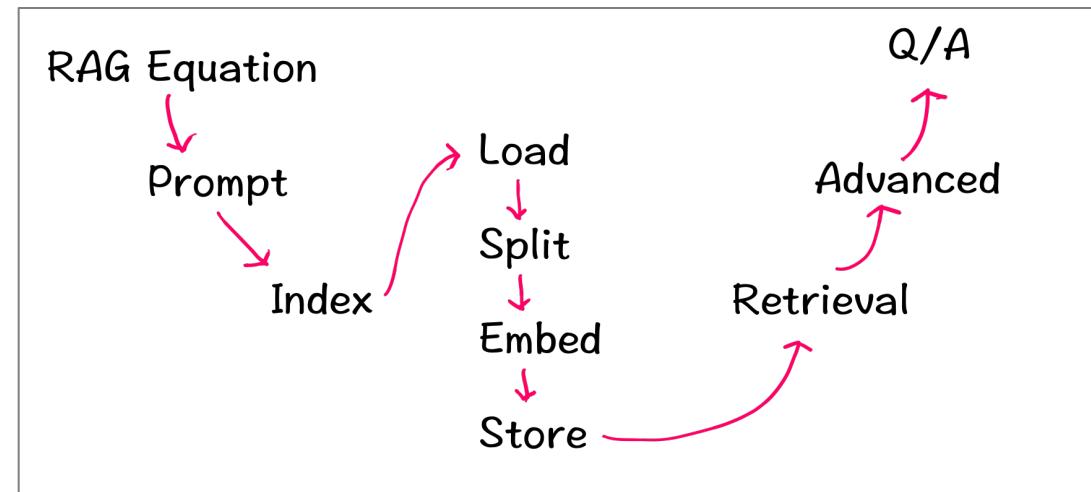


Retrieval

Beginner's Guide to RAG - AI by Hand 🖊



University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

SQL – Retrieve by Similarity

_____ comment, emb<__>[0, 2, 1, 1] **AS** score

FROM posts

ORDER BY _____ ASC | DESC ;

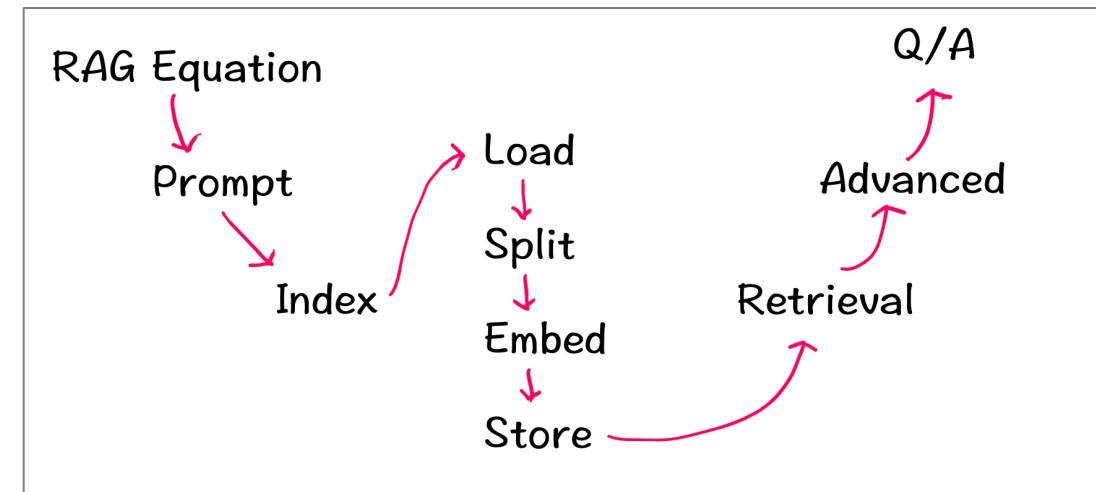


High-level API: Superlinked

```
query = Query(post_index)
    ._____(post)
    ._____(relevance_space.text, Param("_____"))

app.query(query, _____ = "how are you?")
```





Advanced RAGs

Beginner's Guide to RAG - AI by Hand ✎

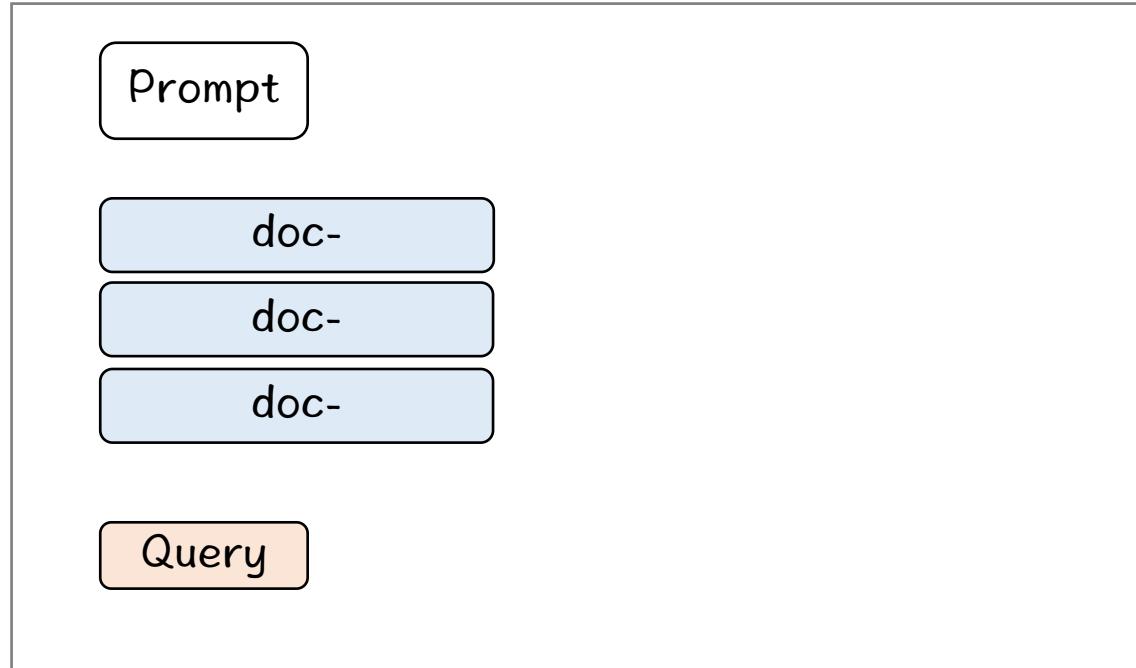


University of Colorado
Boulder

Download: <https://by-hand.ai/rag>

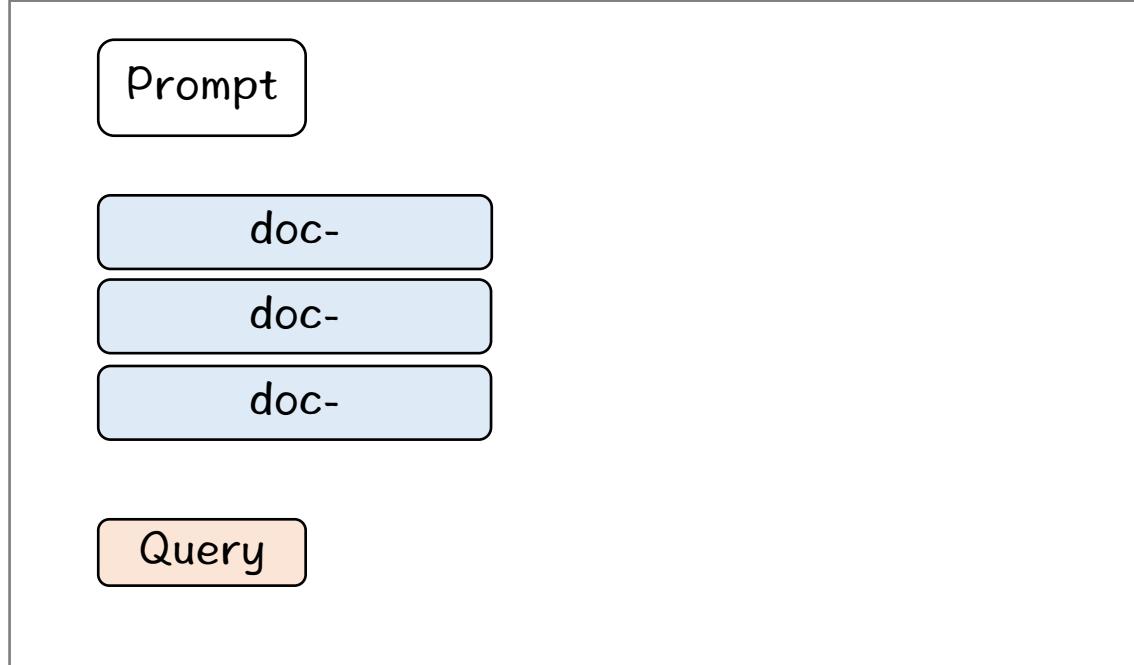
RAG

DB



LLM

Multiple Embedding Spaces DB



LLM

SuperlinkedS Vector Compute

Search

Output

```
df = result._____()  
_____ (df)
```

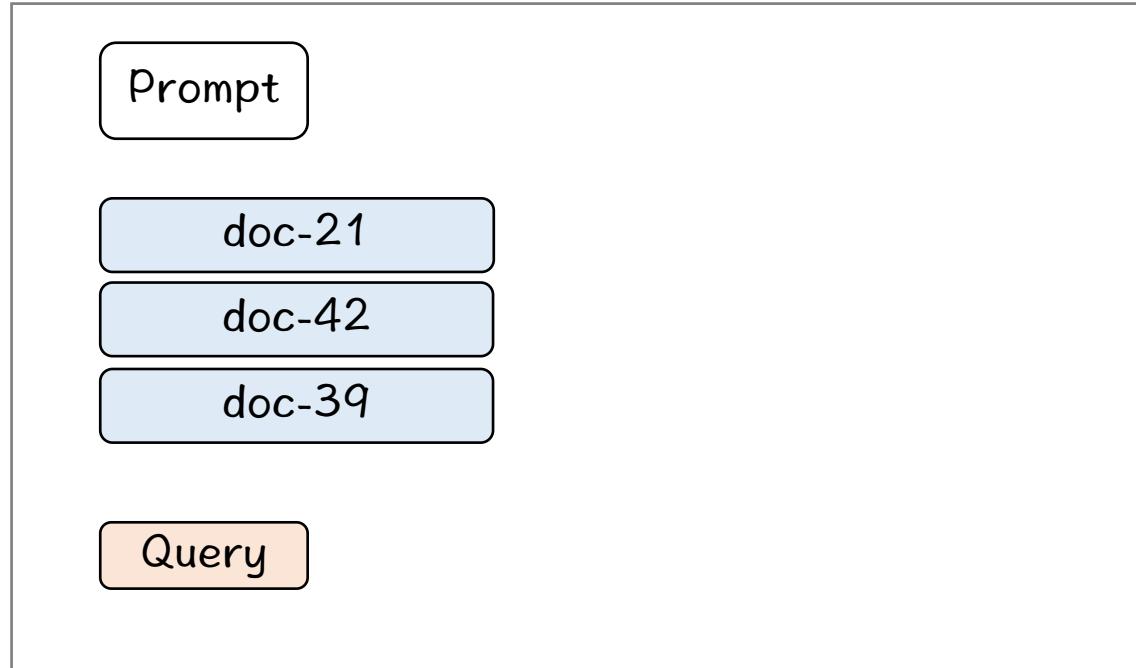
```
result = app.query(  
    simple_query,  
    query_text=[REDACTED],  
    description_weight=1,  
    headline_weight=1,  
    recency_weight=0,
```

Memory

ID	Description	Headline	Recency									
1	<table border="1"><tr><td>0</td><td>1</td><td>1</td><td>0</td></tr></table>	0	1	1	0	<table border="1"><tr><td>0</td><td>1</td><td>0</td></tr></table>	0	1	0	<table border="1"><tr><td>1</td><td>0</td></tr></table>	1	0
0	1	1	0									
0	1	0										
1	0											
2	<table border="1"><tr><td>1</td><td>1</td><td>0</td><td>0</td></tr></table>	1	1	0	0	<table border="1"><tr><td>1</td><td>0</td><td>0</td></tr></table>	1	0	0	<table border="1"><tr><td>1</td><td>0</td></tr></table>	1	0
1	1	0	0									
1	0	0										
1	0											
3	<table border="1"><tr><td>0</td><td>1</td><td>1</td><td>1</td></tr></table>	0	1	1	1	<table border="1"><tr><td>0</td><td>0</td><td>1</td></tr></table>	0	0	1	<table border="1"><tr><td>0</td><td>1</td></tr></table>	0	1
0	1	1	1									
0	0	1										
0	1											

RankGPT

DB



LLM

RankGPT

Search Query: {query}.

Rank the {num} passages above based on their [REDACTED] to the search query. The passages should be listed in [REDACTED] order using identifiers.

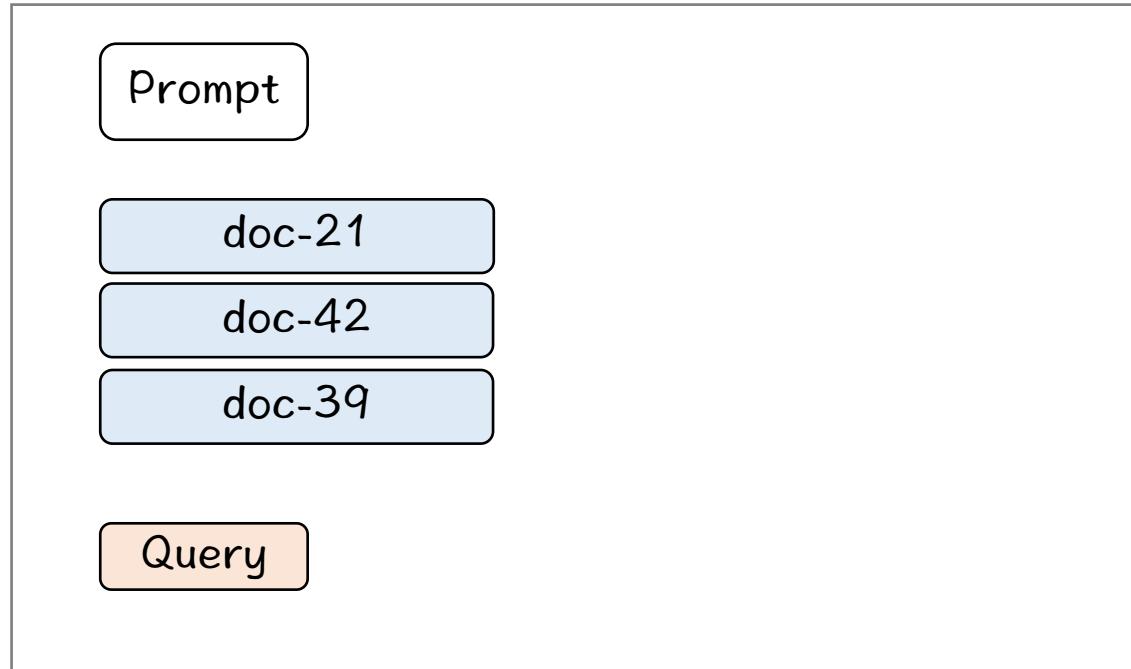
The most relevant passages should be listed first.

The output format should be [] > [], e.g., [1] > [2].

Only response the ranking results,

do not [REDACTED].

Multi-Query Retrieval DB



LLM

Multi-Query Retrieval

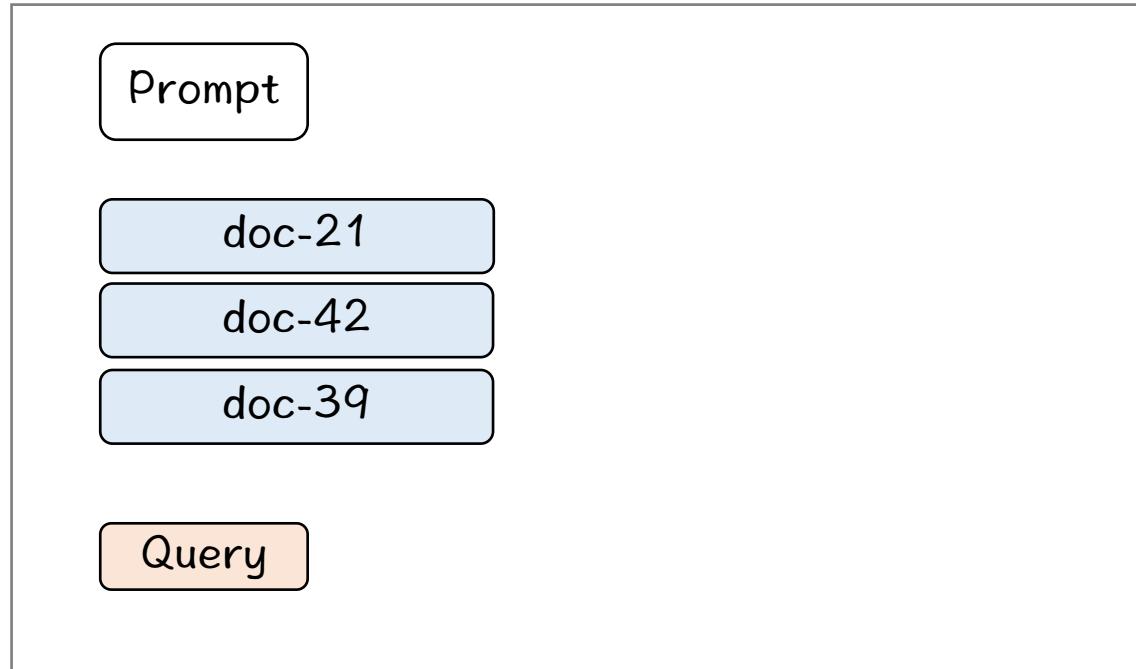
You are an AI language model assistant. Your task is to generate _____ different versions of the given user question to retrieve relevant documents from a vector database.

By generating multiple _____ on the user question, your goal is to help the user overcome some of the limitations of _____ search.

Provide these alternative questions seperated by newlines.

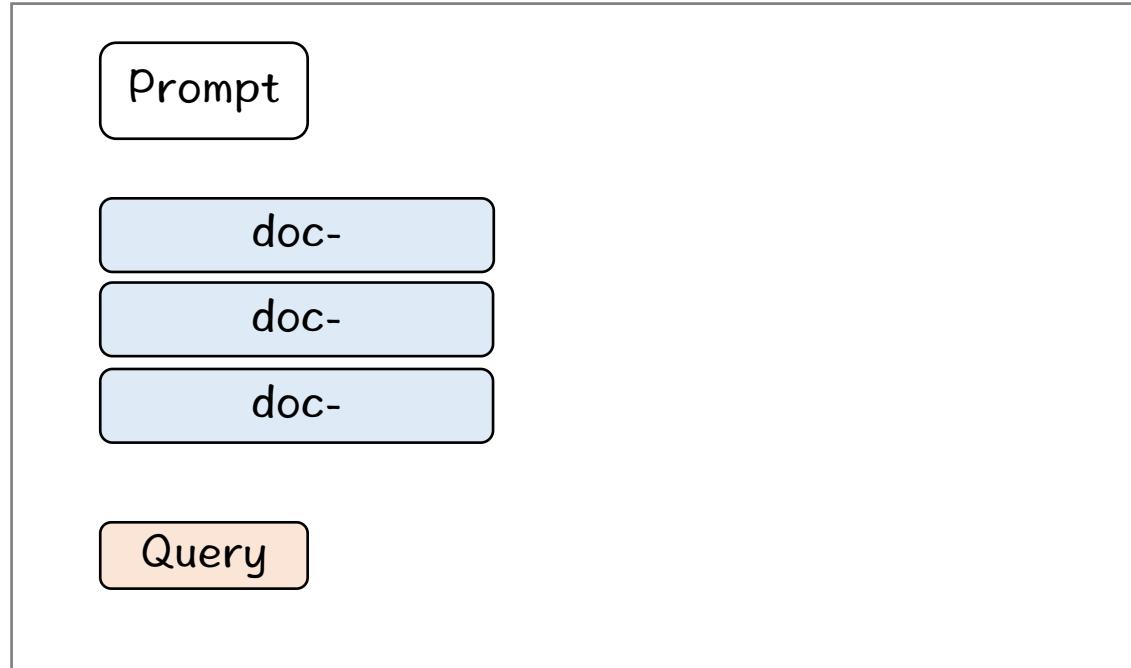
Original question: {question}

Contextual Compression DB



LLM

Hypothetical Document Embeddings (HyDE) DB



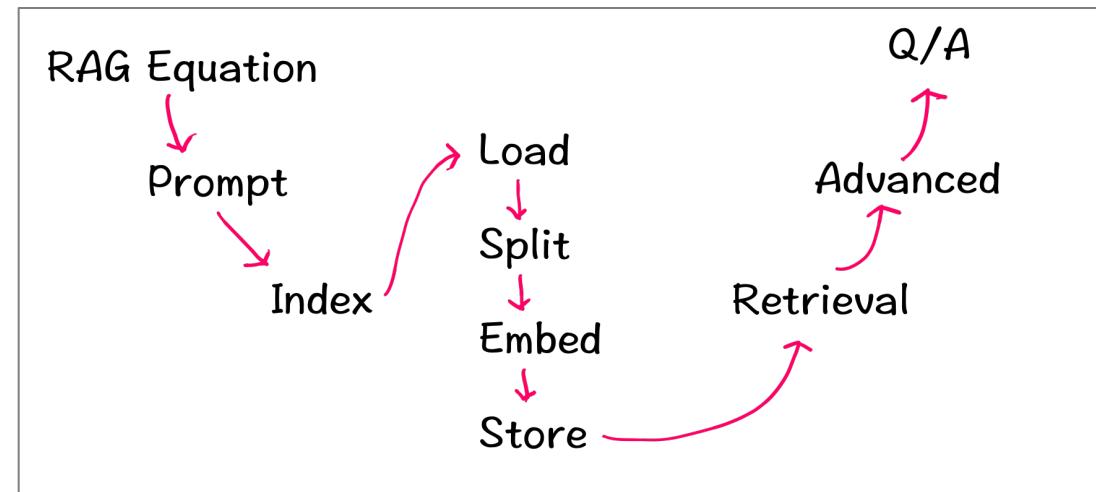
LLM

Q/A

Beginner's Guide to RAG - AI by Hand 🖊



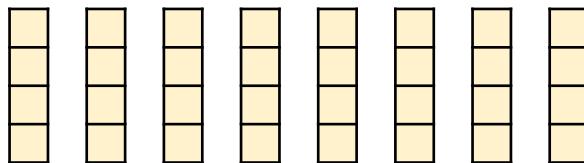
University of Colorado
Boulder



Download: <https://by-hand.ai/rag>

How does an LLM sample a sentence?

Input Embeddings



LLM

Probability Distributions

Vocab

I	.01
you	.01
they	.01
are	.01
am	.01
how	.50
why	.10
where	.10
who	.15
what	.10

.01
.01
.01
.40
.40
.05
.05
.05
.01
.01

.03
.50
.40
.01
.01
.01
.01
.01
.01

Random Numbers

.
.34
.52
.92
.65

How does an LLM sample a sentence?

