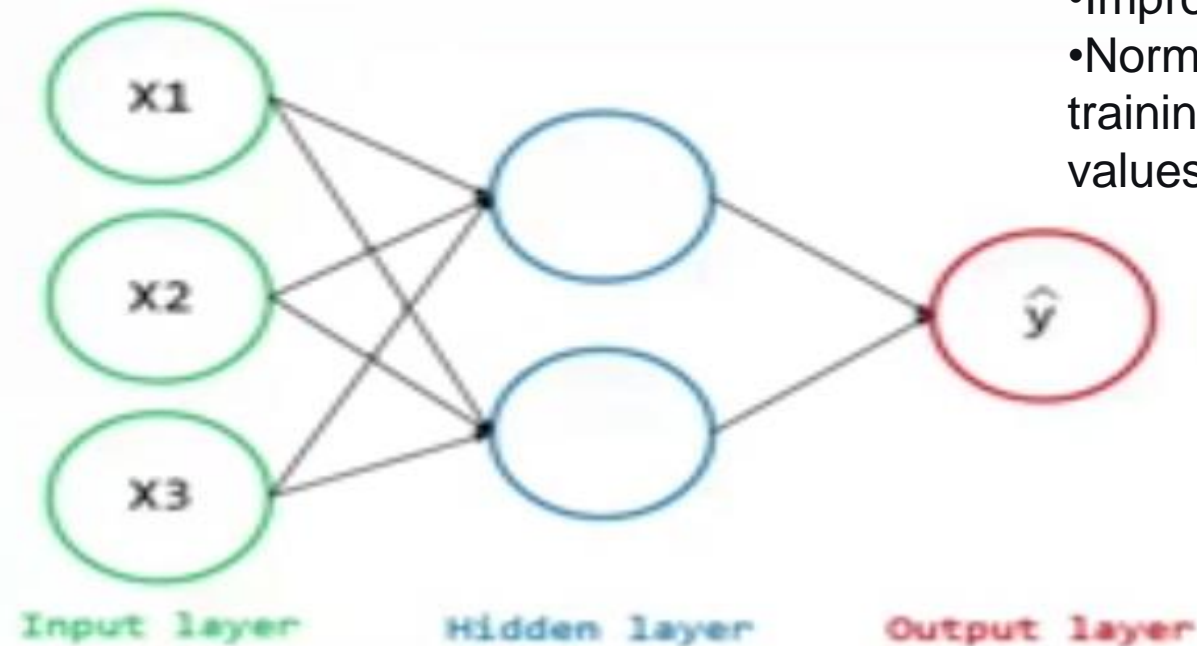


# Normalization-Transformer

Dr. Muhammad Safyan

- Normalization in deep learning refers to the process of transforming data or model outputs to have specific statistical properties, typically a mean of zero and a variance of one.

What do we normalize?



### Benefits of Normalization in Deep Learning

- Improved Training Stability:
- Normalization helps to stabilize and accelerate the training process by reducing the likelihood of extreme values that can cause gradients to explode or vanish.

# Benefits of Normalization in Deep Learning

- Improved Training Stability:
  - Normalization helps to stabilize and accelerate the training process by reducing the likelihood of extreme values that can cause gradients to explode or vanish.
- Faster Convergence:
  - By normalizing inputs or activations, models can converge more quickly because the gradients have more consistent magnitudes. This allows for more stable updates during back propagation.
- Mitigating Internal Covariate Shift:
  - Internal covariate shift refers to the change in the distribution of layer inputs during training. Normalization techniques, like batch normalization, help to reduce this shift, making the training process more robust.



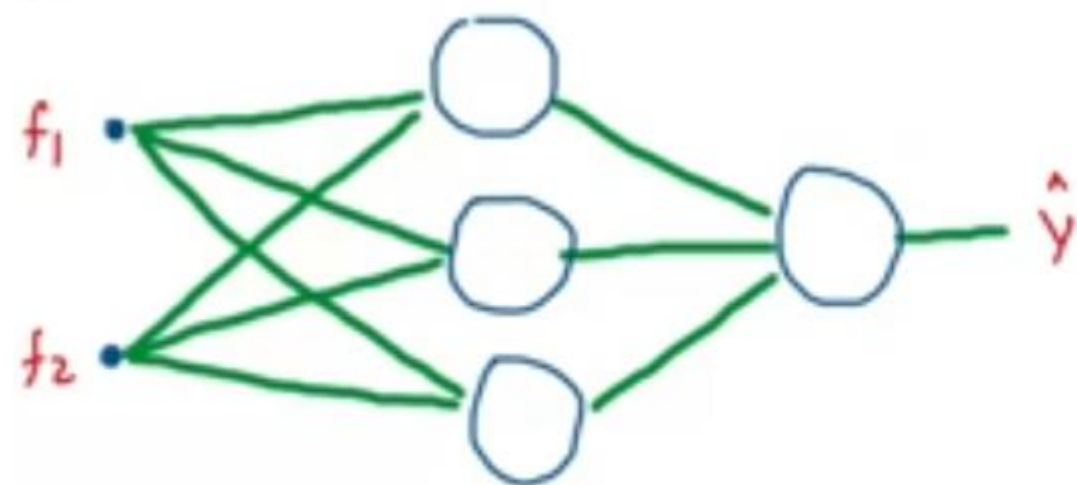
covariate shift



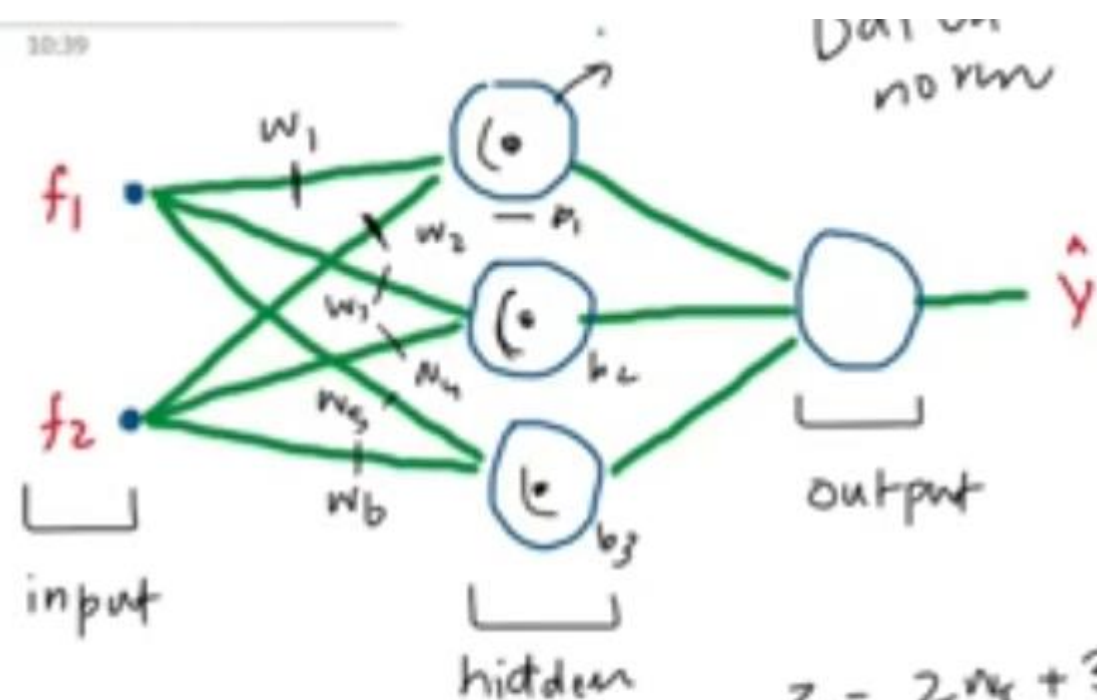
red roses

pred  
↓  
yellow white

10:39



$f_1$	$f_2$
2	3
1	1
5	4
6	1
7	1



$$z_3 = 2w_5 + 3w_6 + b_3 = 4$$

$$z_1 = 2w_1 + 3w_2 + b_1 = 7$$

$$z_2 = 2w_3 + 3w_4 + b_2 = 5$$

$f_1$	$f_2$	$z_1$	$z_2$	$z_3$
2	3	7	5	4
1	1	2	3	4
5	4	1	2	3
6	1	7	5	6
7	1	3	3	4

Below the table, arrows point from the  $z_1, z_2, z_3$  columns to labels  $h_1, \sigma_1, h_2, \sigma_2, h_3$  respectively.

↓  $y(2)$        $batches = 5$

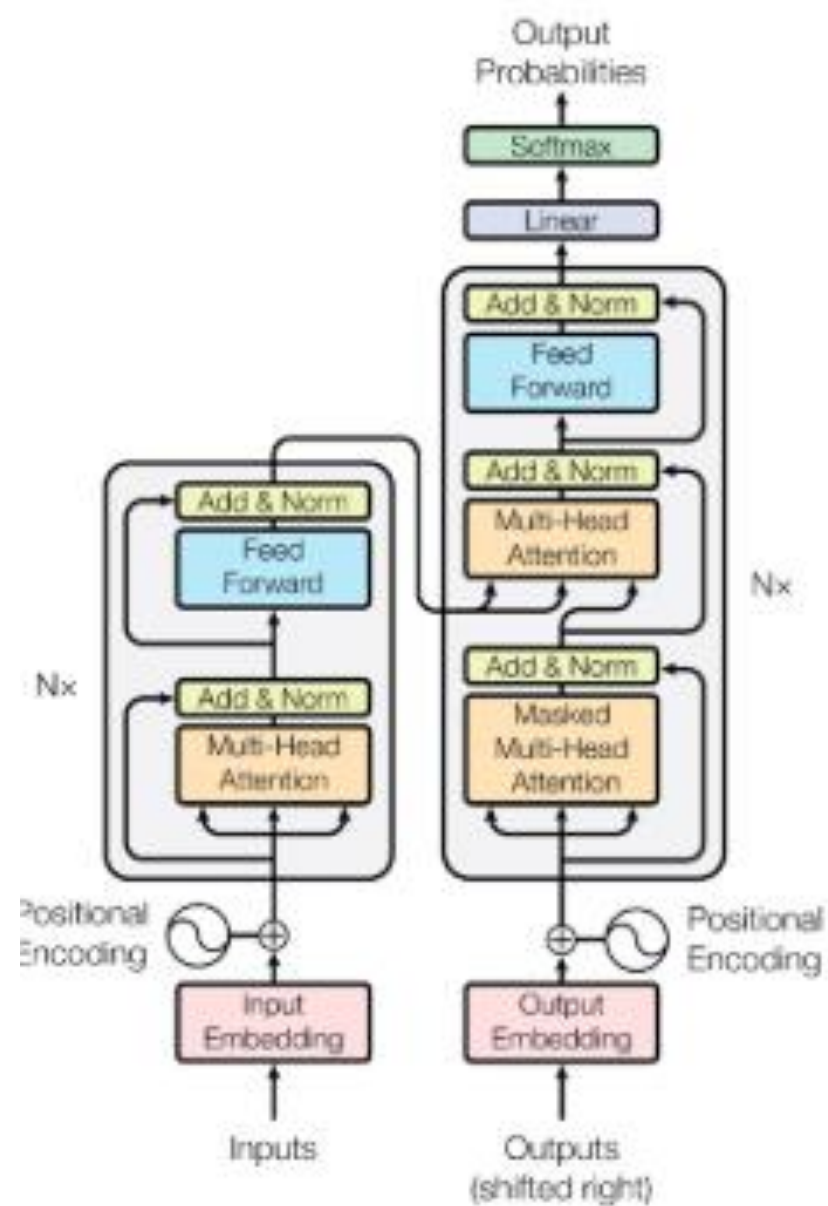
$f_1$	$f_2$	$z_1$	$z_2$	$z_3$
2	3	7	5	4
1	1	2	3	4
5	4	1	2	3
6	1	7	5	6
7	1	3	3	4

$\mu_1$     $\sigma_1$     $\mu_2$     $\sigma_2$     $\mu_3$     $\sigma_3$

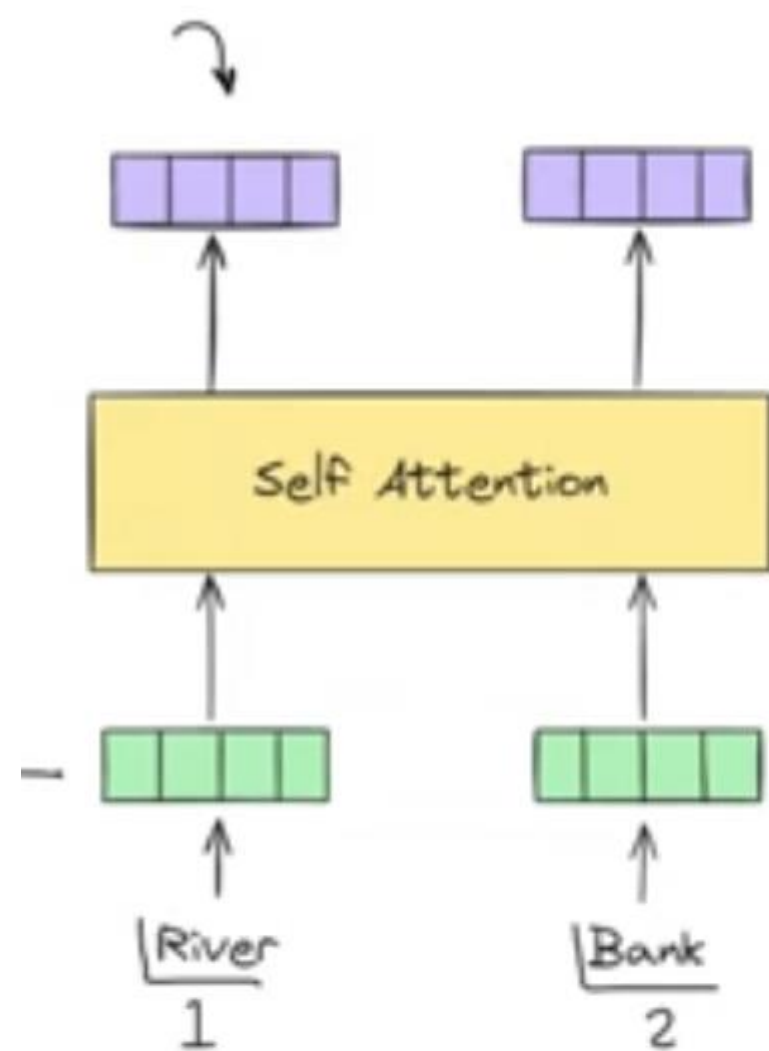
$v_3 + 3w_4 + b_4 = 5$

$$\frac{7 - \mu_1}{\sigma_1} = \frac{0.36}{(1)} \frac{y_1}{(1)} + \frac{\beta_1}{(0)} = 0.36$$

$$\frac{2 - \mu_1}{\sigma_1} = 0.71 y_1 + \beta_1 = 0.71$$







Review	Sentiment
Hi Babar	1
How are you today	0
I am good	0
You?	1

Embedding dimension - 3

Batch Size - 2

0.2	0.45	0.71
-----	------	------

IP

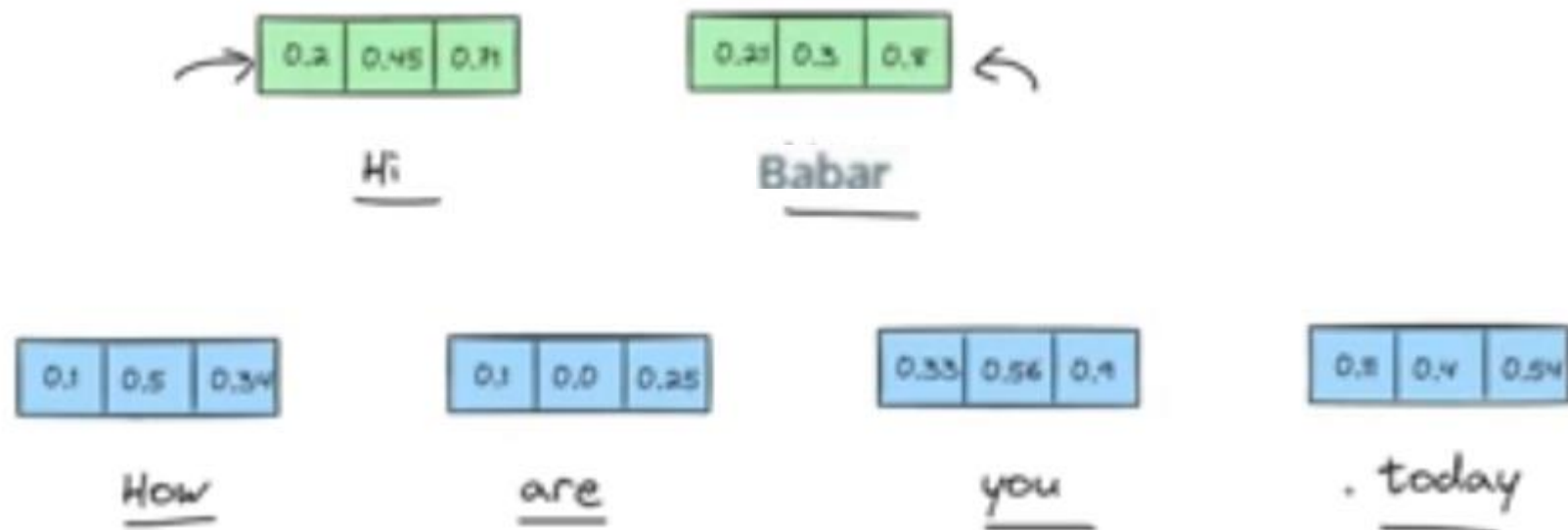
0.20	0.3	0.8
------	-----	-----

Output

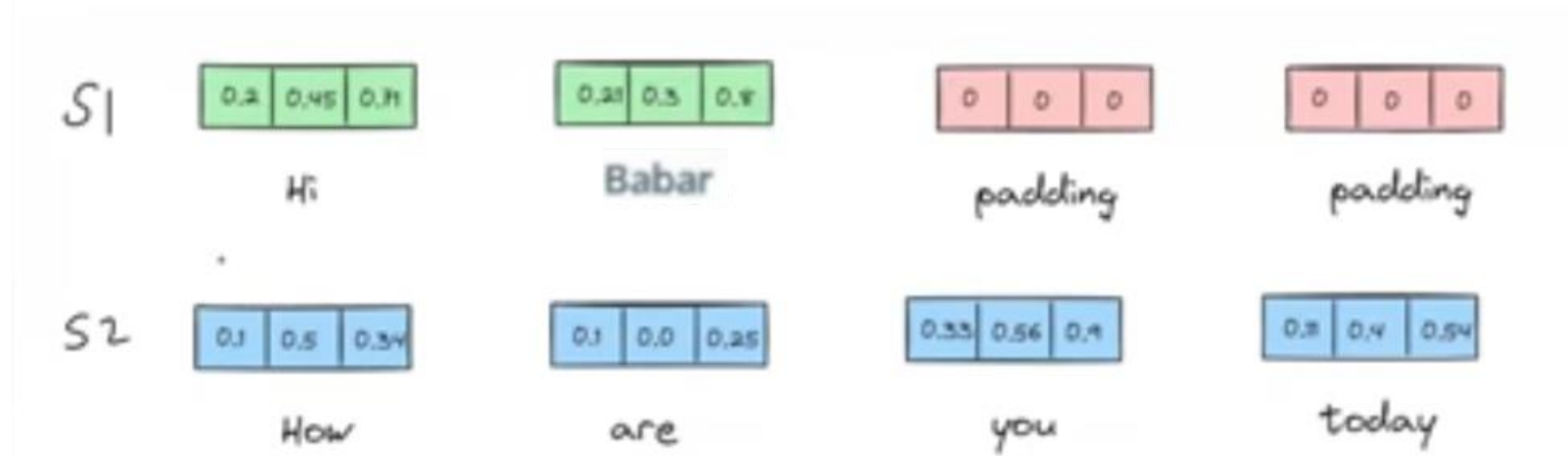
Review	Sentiment
Hi Babar	1
How are you today	0
I am good	0
You?	1

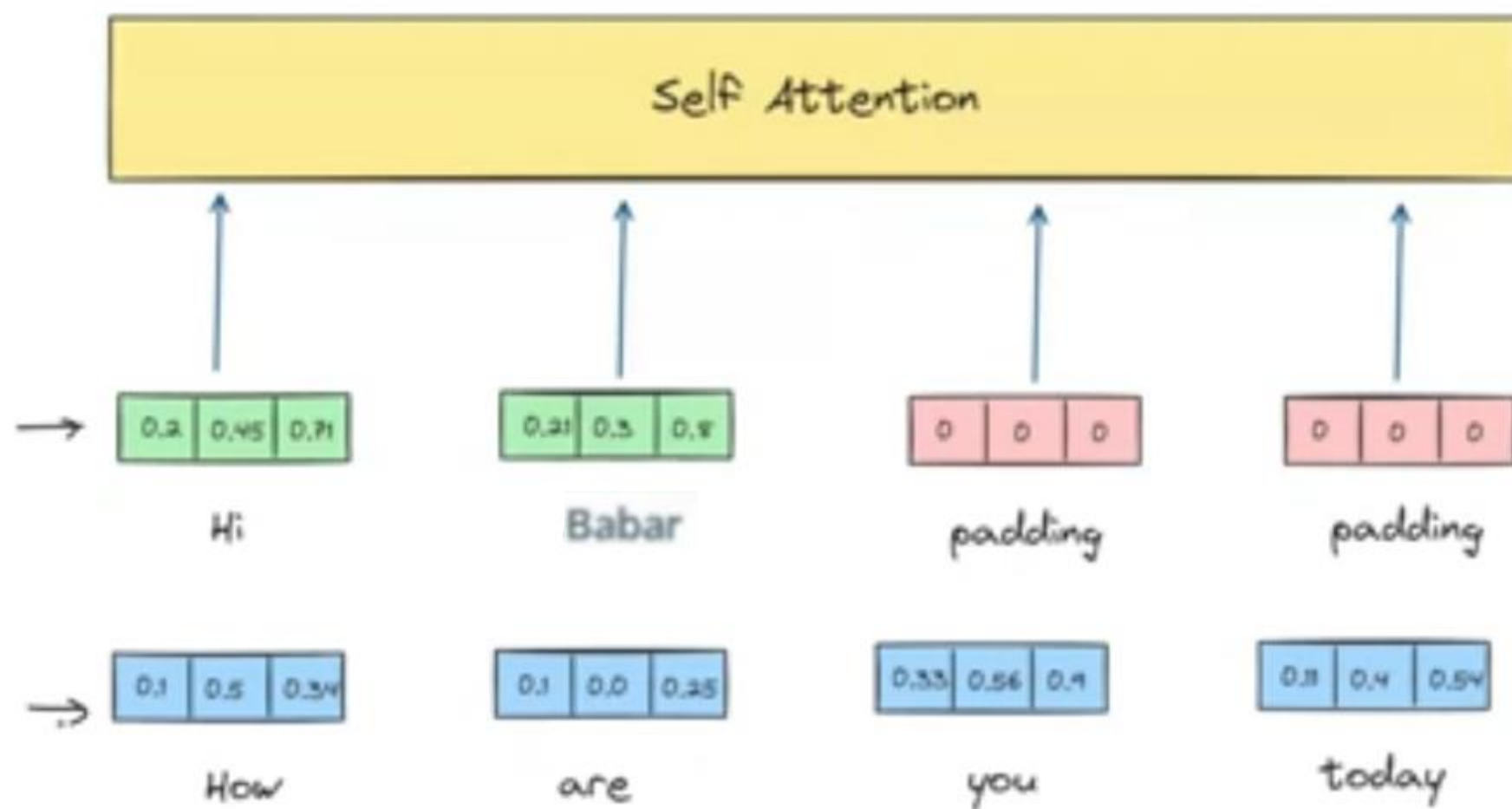
Embedding dimension - 3

Batch Size - 2

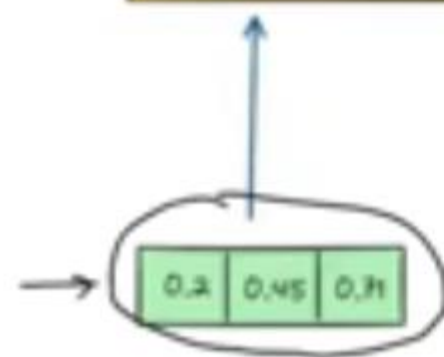


# padding

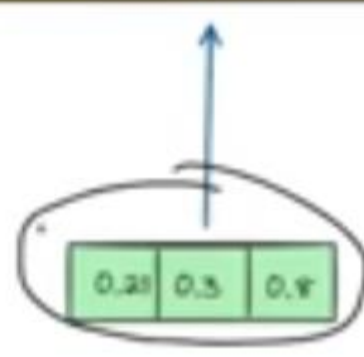




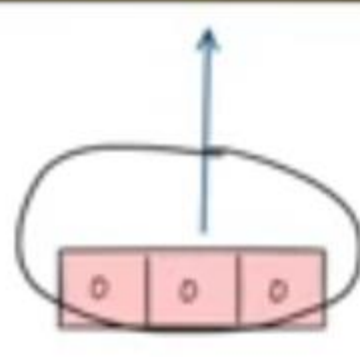
→ Self Attention



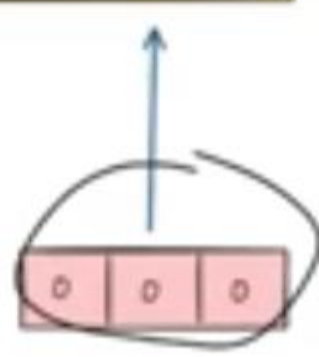
Hi



Babar

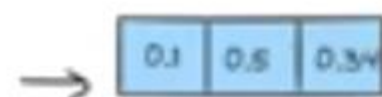


padding



padding

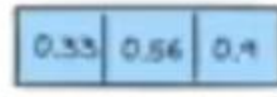
matrix



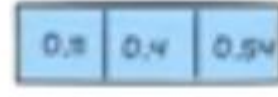
How



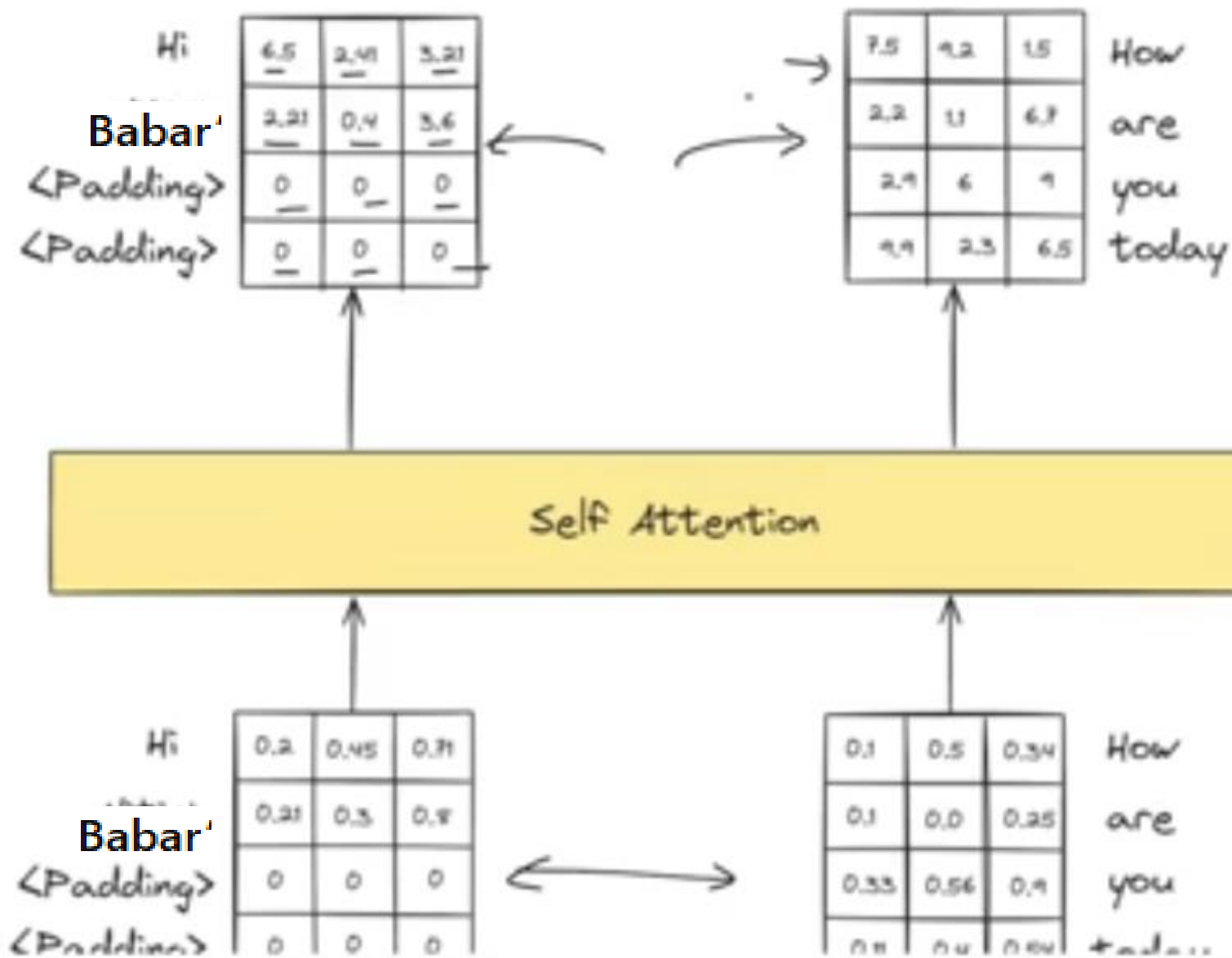
are



you



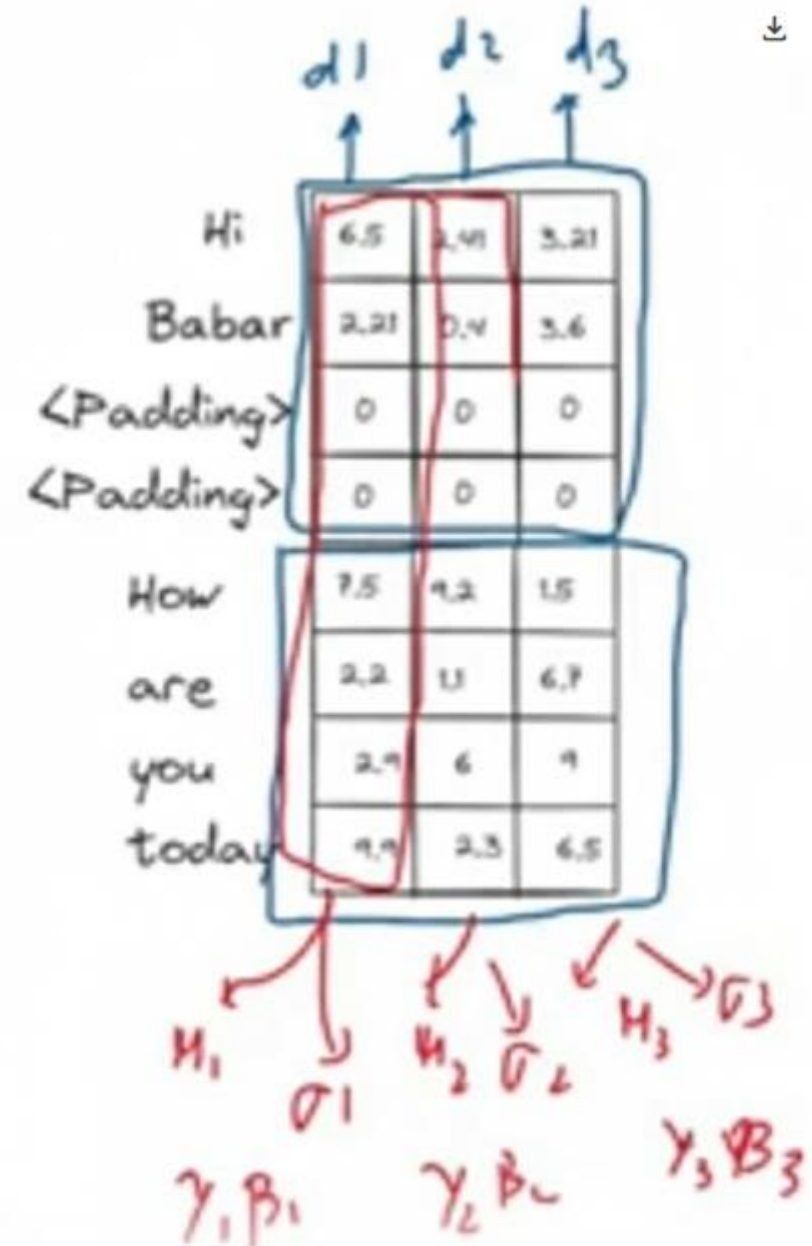
today



# Vertically stackup

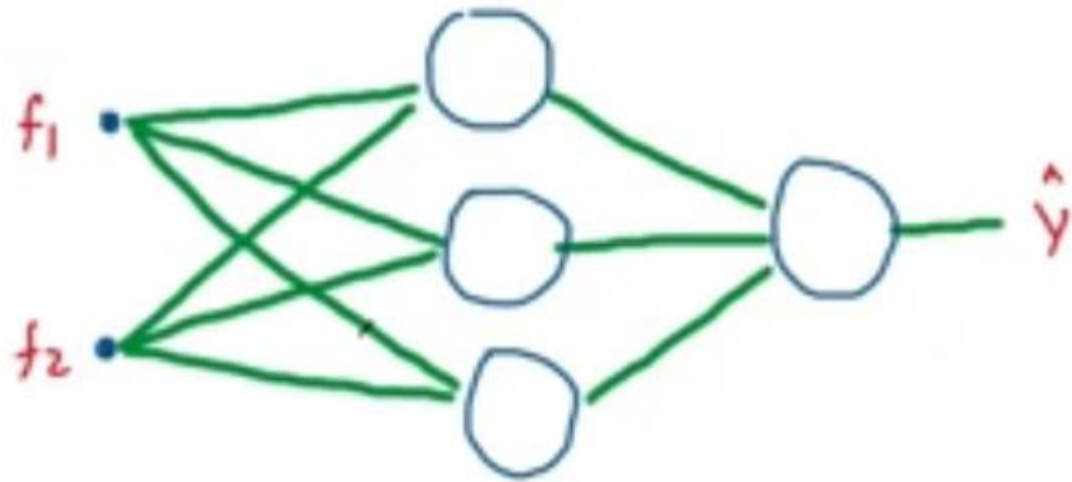
Hi	6.5	2.41	3.21
Babar	2.21	0.4	3.6
<Padding>	0	0	0
<Padding>	0	0	0
How	7.5	4.2	1.5
are	2.2	1.1	6.7
you	2.9	6	9
today	9.9	2.3	6.5

- Imagine
- 100 sentences batch
- Average sentence length 20
- Longest sentence 60.





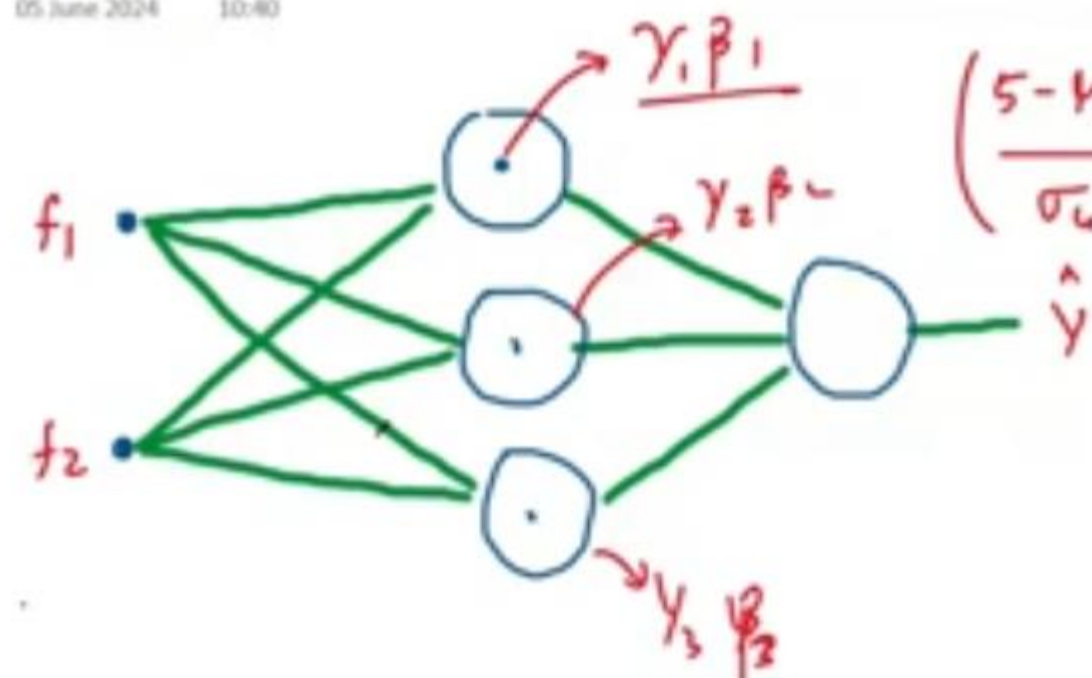
# Layer Normalization



$f_1$	$f_2$	$z_1$	$z_2$	$z_3$
2	3	7	5	4
1	1	2	3	4
5	4	1	2	3
6	1	7	5	6
7	1	3	3	4

# Layer Norm

05 June 2024 10:40



$$\left( \frac{2 - \mu_2}{\sigma_2} \right) \gamma_1 + \beta_1 \quad \left( \frac{4 - \mu_1}{\sigma_1} \right) \gamma_3 + \beta_3$$

$$\left( \frac{5 - \mu_4}{\sigma_4} \right) \gamma_2 + \beta_2$$

$$\left( \frac{5 - \mu_1}{\sigma_1} \right) \gamma_2 + \beta_2$$

$$\frac{7 - \mu_1}{\sigma_1} = \boxed{0.3 \gamma_1 + \beta_1}$$

across batch across feature

$f_1$	$f_2$	$z_1$	$z_2$	$z_3$
2	3	7	5	4
1	1	2	3	4
5	4	1	2	3
6	1	7	5	6
7	1	3	3	4

Annotations for the table:

- $\mu_1$  points to the first column ( $f_1$ ).
- $\sigma_1$  points to the first row ( $z_1$ ).
- $\mu_2$  points to the second column ( $f_2$ ).
- $\sigma_2$  points to the second row ( $z_2$ ).
- $\mu_5$  points to the fifth column ( $z_3$ ).
- $\sigma_5$  points to the fifth row ( $z_3$ ).

