

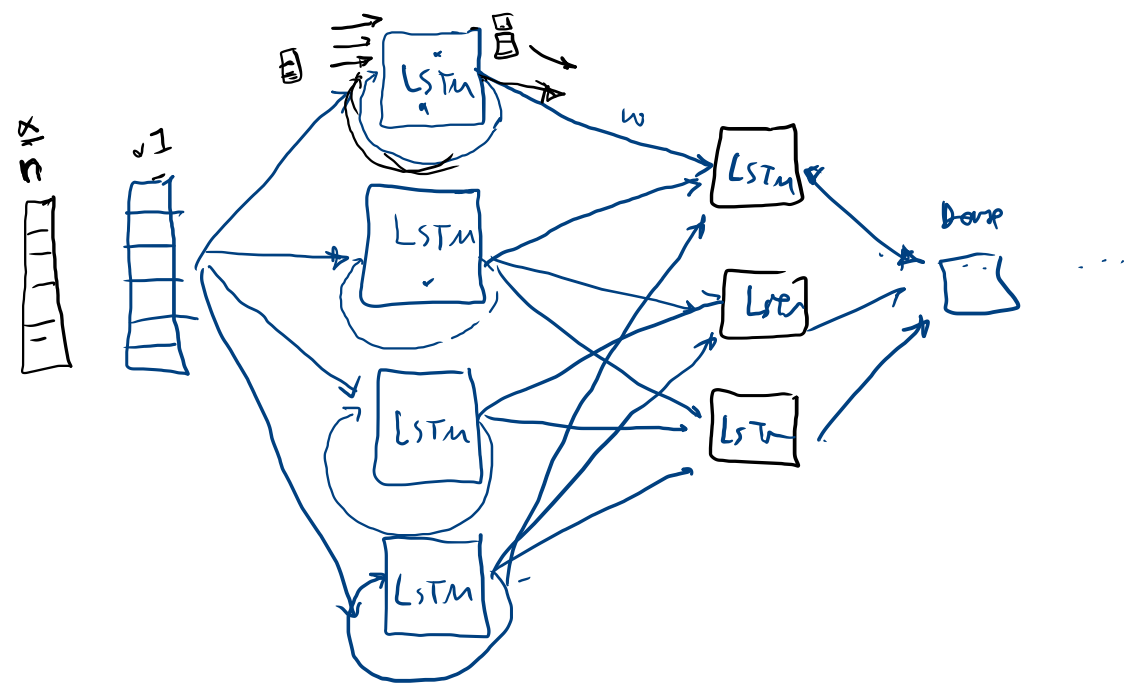
ViT (Vision Transformer)

Transformer Encoder / decoder

---

Vision Transformer

layer 1: LSTM (units: 5) / layer 2: LSTM (units: 3)



- detection ✓ → yolo v8 . → mAP : 73 → 90<sup>NT</sup>
- segmentation ✓ → deeplob v3

ViT → - classification → VGG



2018 ? → استخراج!



image → performance metrics → improve!

!  $\vec{u}, \vec{v}$  - Embedding ViT  $\vec{u}, \vec{v}$   
 $\vec{u}, \vec{v}$  - <sup>Embed</sup> LSTM CNN  $\vec{u}, \vec{v}$   
NO

## Data Hungry

۱.  $ViT$  بیشتر داده‌ها! - نیازمند دیتای بزرگتر!

ImageNet  
↓  
عدم امتداد

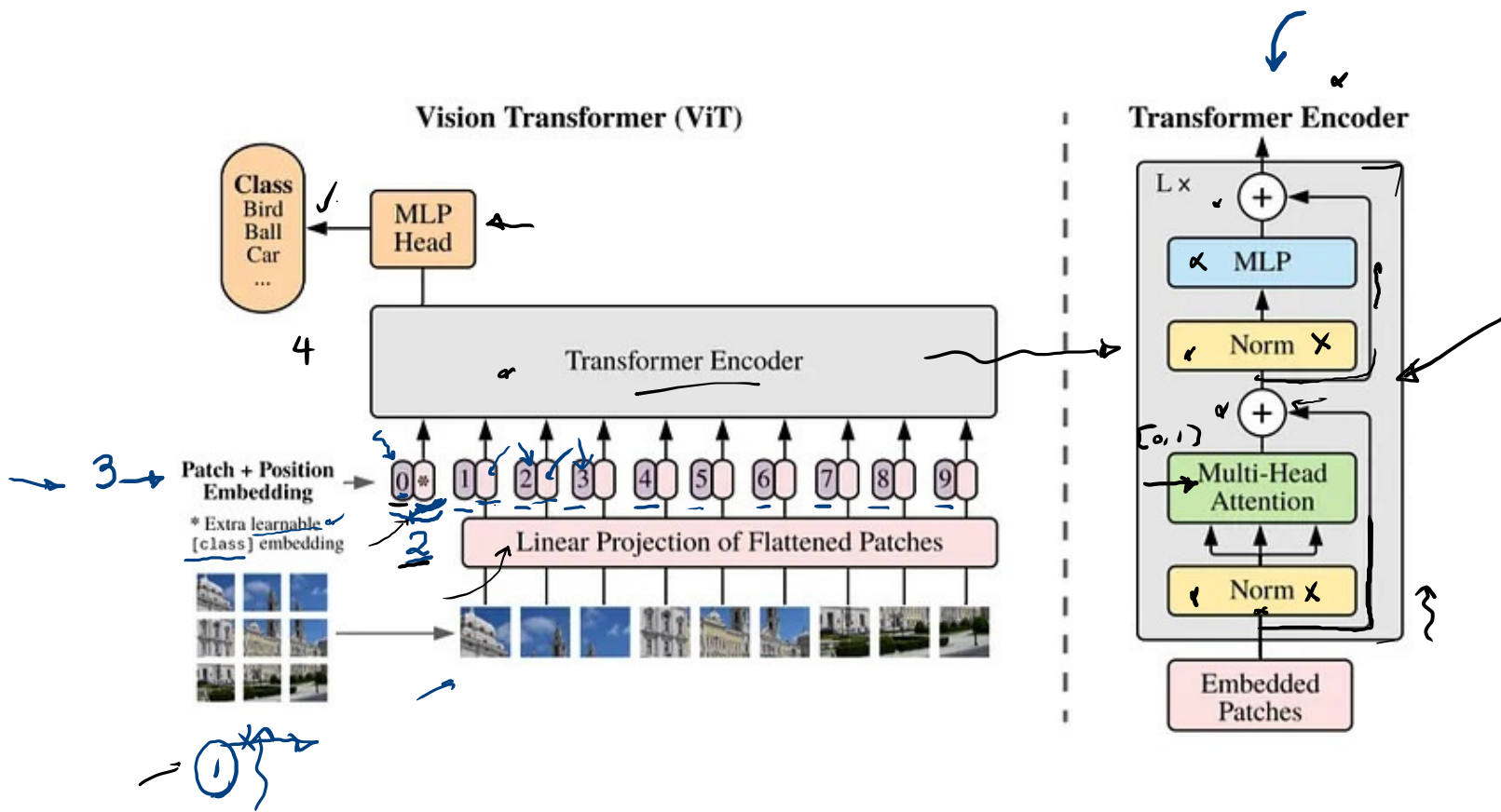
۲. Computation Cost.



$ViT \gg \gg CNN$

۳. هنوز دیگرها همده نتوانسته‌اند / این اکثریت هر  $CNN$ ،

لاستیک را به



1. preprocessing!

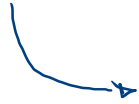


Image to patches!

1	2	3
4	5	6
7	8	9



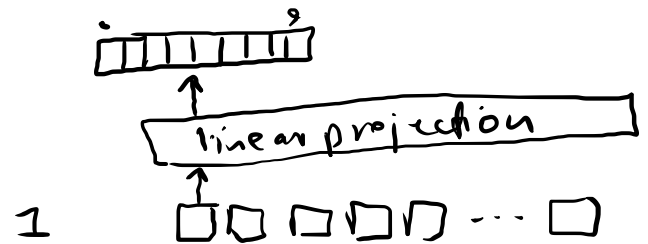
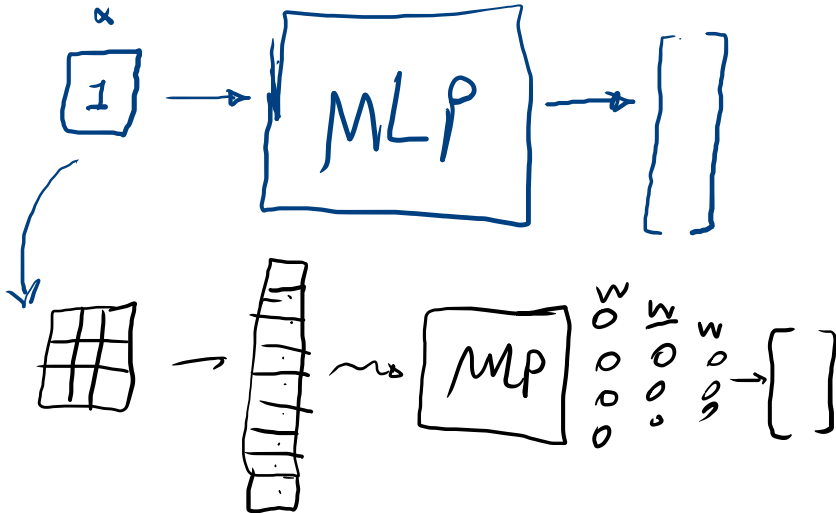
1 2 3 4 5 6 7 8 9

↓  
patch

Conv2D (kernel size = stride)

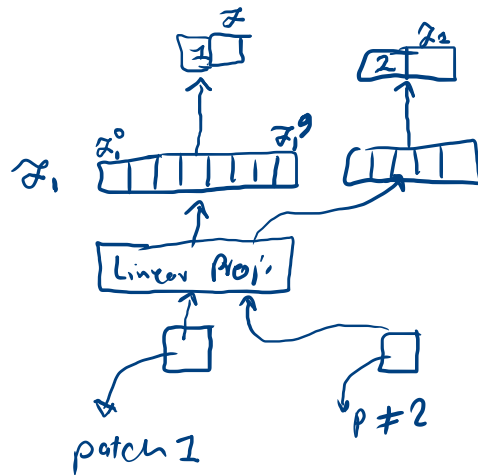
## 2. Linear projection!

نیم پاتچ  
کے لیے (1) سے!

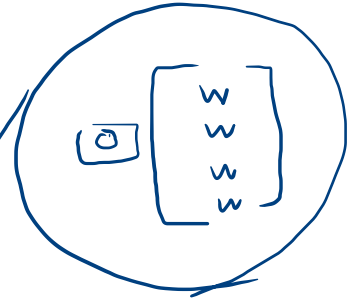


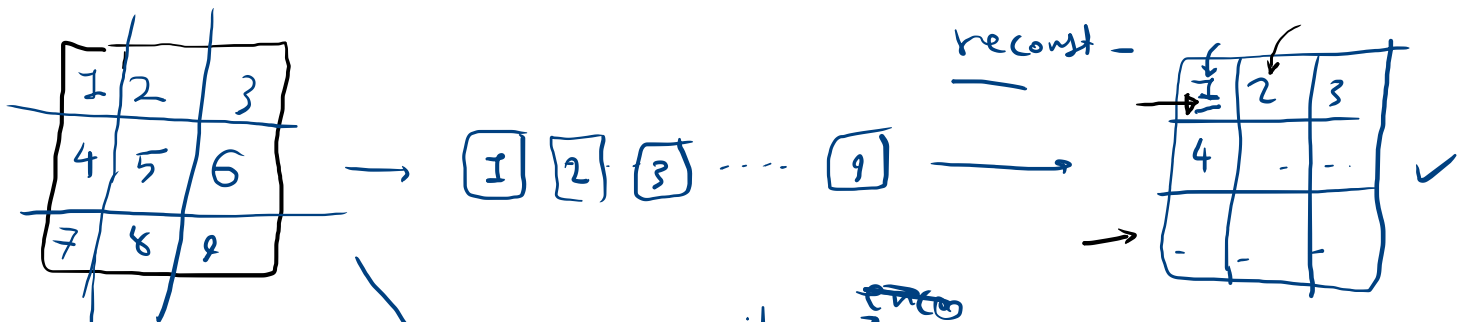


### 3. position Embedding!



why?





without positional ~~info~~ Embed.

4	3	8
9	7	1
2	6	5

Conv ~ { scaling, ~~in~~ rotation → not invariant.

Conv, pooling!

pool → scaling, invariant.

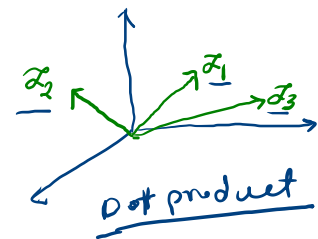
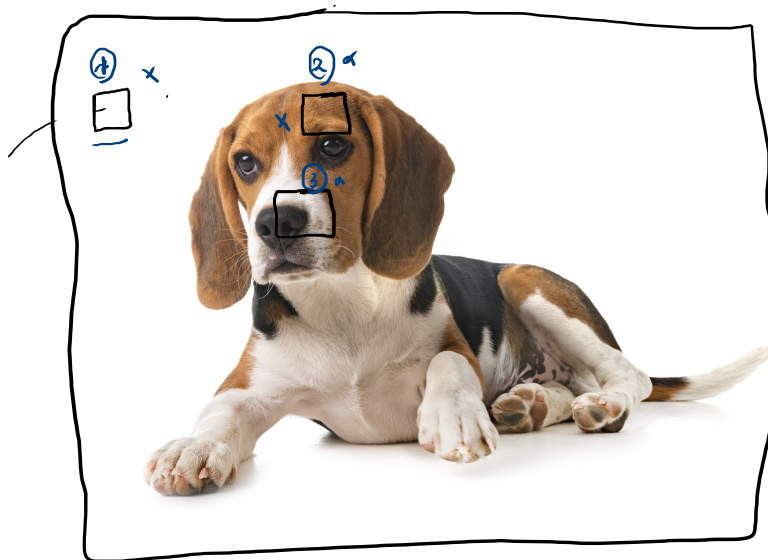
#### 4. Transformer encoder!

Attention!

$\mathcal{Z}_1$  - vector  $\underline{3}$  (9, 8, 7)

$\mathcal{Z}_2$  - vector

$\mathcal{Z}_3$  - vector



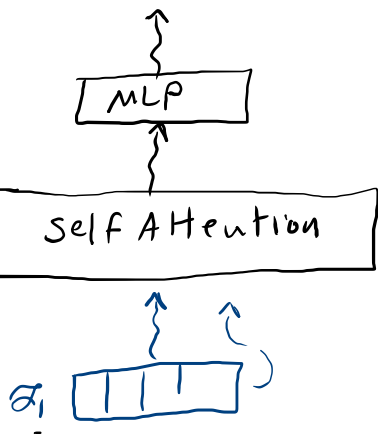
$$\mathcal{Z}_2 \mathcal{Z}_1 \ll \mathcal{Z}_3 \mathcal{Z}_2$$

	$\mathcal{Z}_1$	$\mathcal{Z}_2$	$\mathcal{Z}_3$
$\mathcal{Z}_1$	<u>1</u>	.3	.2
$\mathcal{Z}_2$	x	x	x
$\mathcal{Z}_3$	x	x	x

→ ✓

Attention!

①



Query: Feature of interest!

[ ]

$$\text{Attention}(\underline{Q}, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Normalization Factor

Value: original features to be scaled by probabilities!

Key

$$Q = [\underline{x_1}, x_2, \dots, x_9] \underline{w^Q}$$

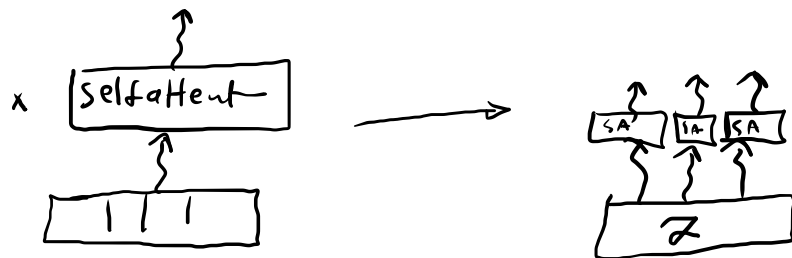
$$k: \text{key} : [\underline{x_1}, x_2, \dots, x_9] \underline{w^k}$$

$$v: [x_1, \dots, x_9] \underline{w^{\cancel{k}}}$$

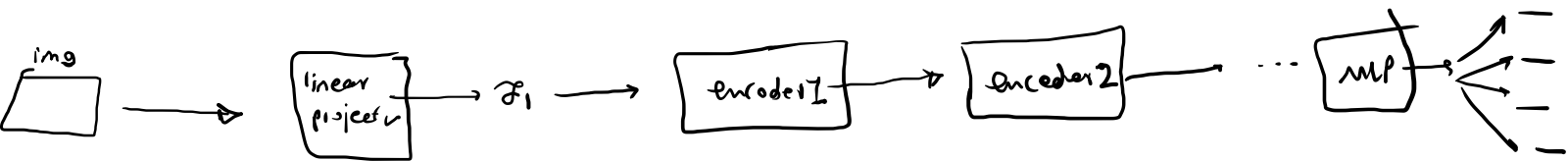
5. Improvement in Encoder.

↓ a. Normalization Layer  
and skip connection!

b. multi head self attention! (MHSA)







-  
End