

+

+

SA  
+  
+

موضع

# SAM (Segment Anything Model)

X

Instance Segmentation

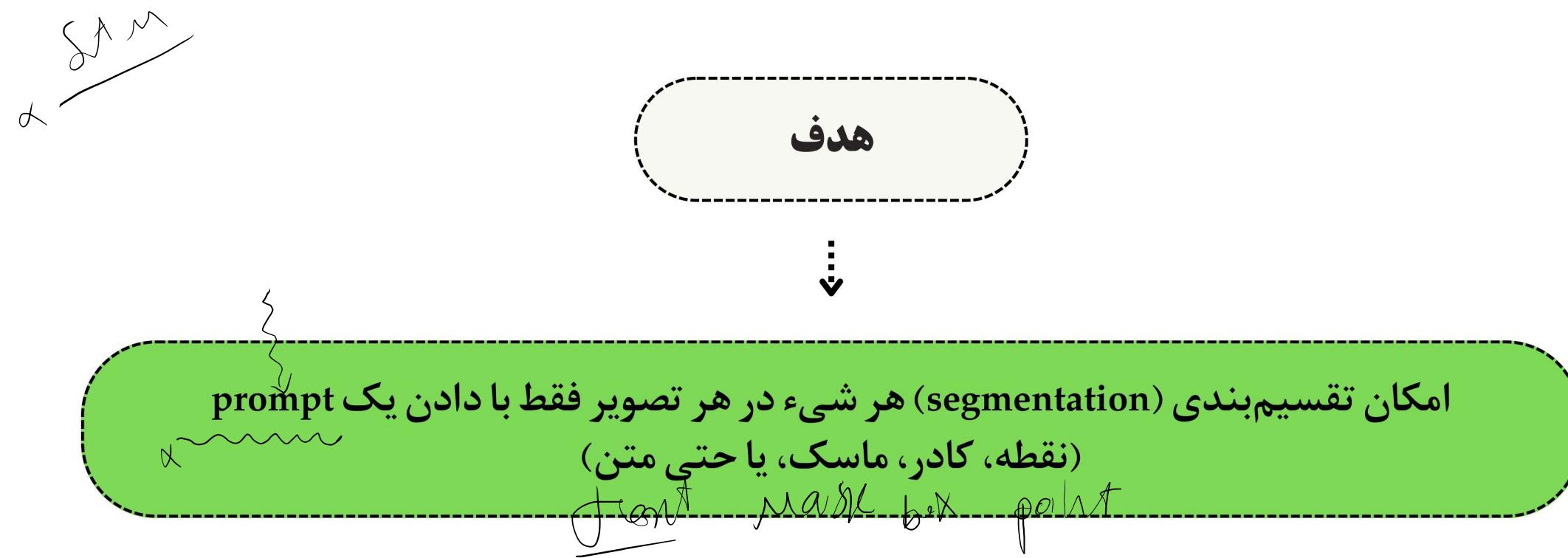
X

↙ [ Segment Anything ]

Alexander Kirillov<sup>1,2,4</sup> Eric Mintun<sup>2</sup> Nikhila Ravi<sup>1,2</sup> Hanzi Mao<sup>2</sup> Chloe Rolland<sup>3</sup> Laura Gustafson<sup>3</sup>  
Tete Xiao<sup>3</sup> Spencer Whitehead Alexander C. Berg Wan-Yen Lo Piotr Dollár<sup>4</sup> Ross Girshick<sup>4</sup>  
<sup>1</sup>project lead    <sup>2</sup>joint first author    <sup>3</sup>equal contribution    <sup>4</sup>directional lead

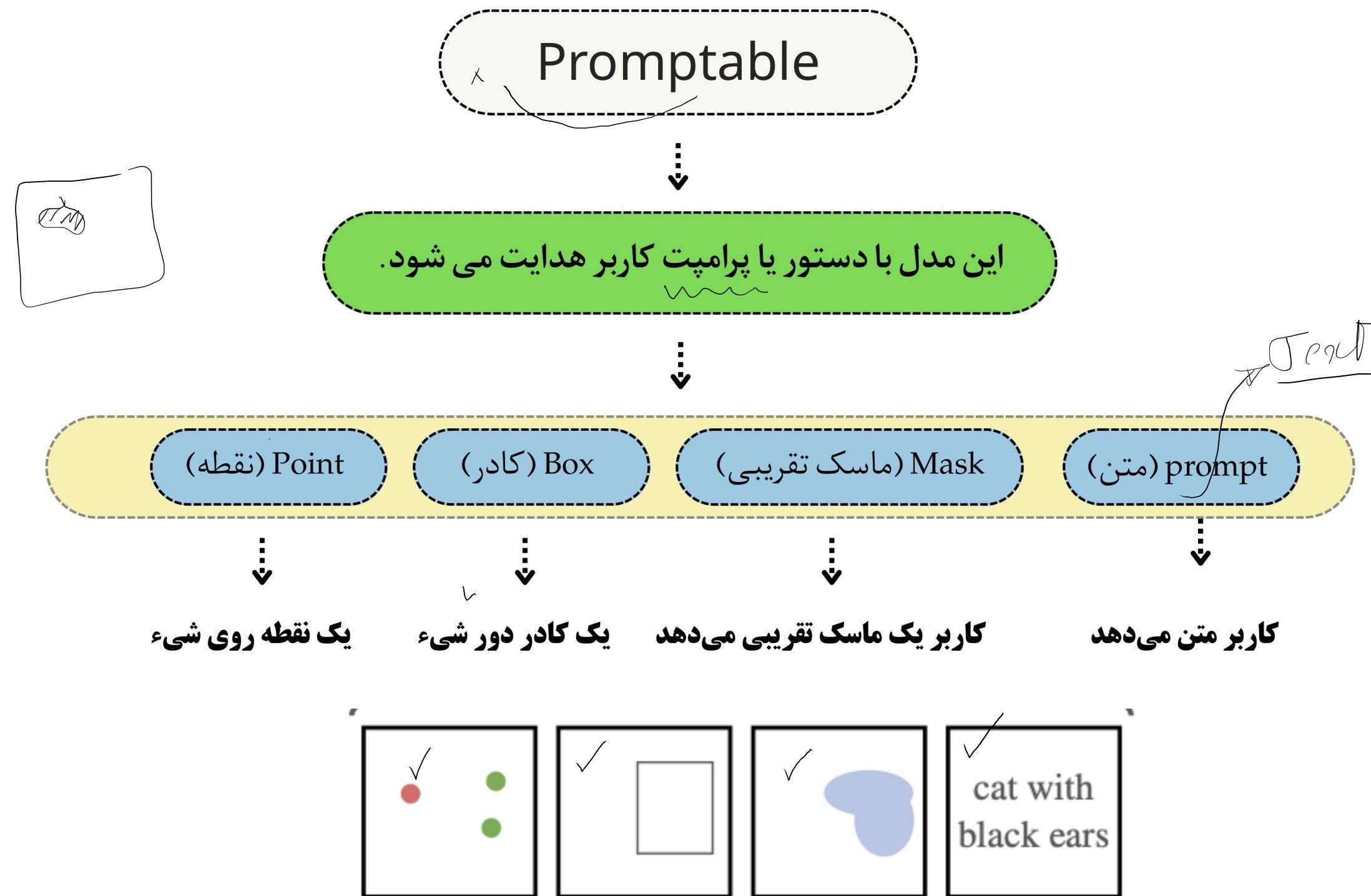
↔

Meta AI Research, FAIR



# ویژگی های SAM1

Instance  
segmentation

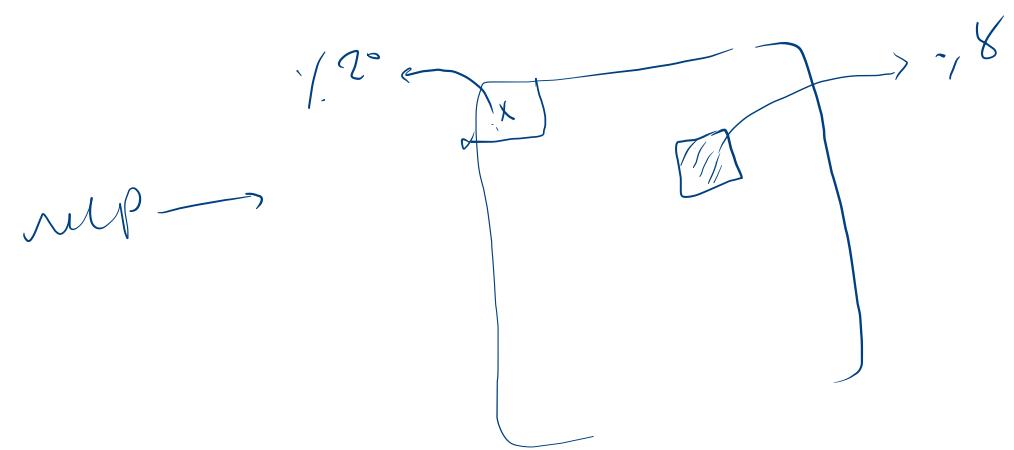


ویژگی های SAM1

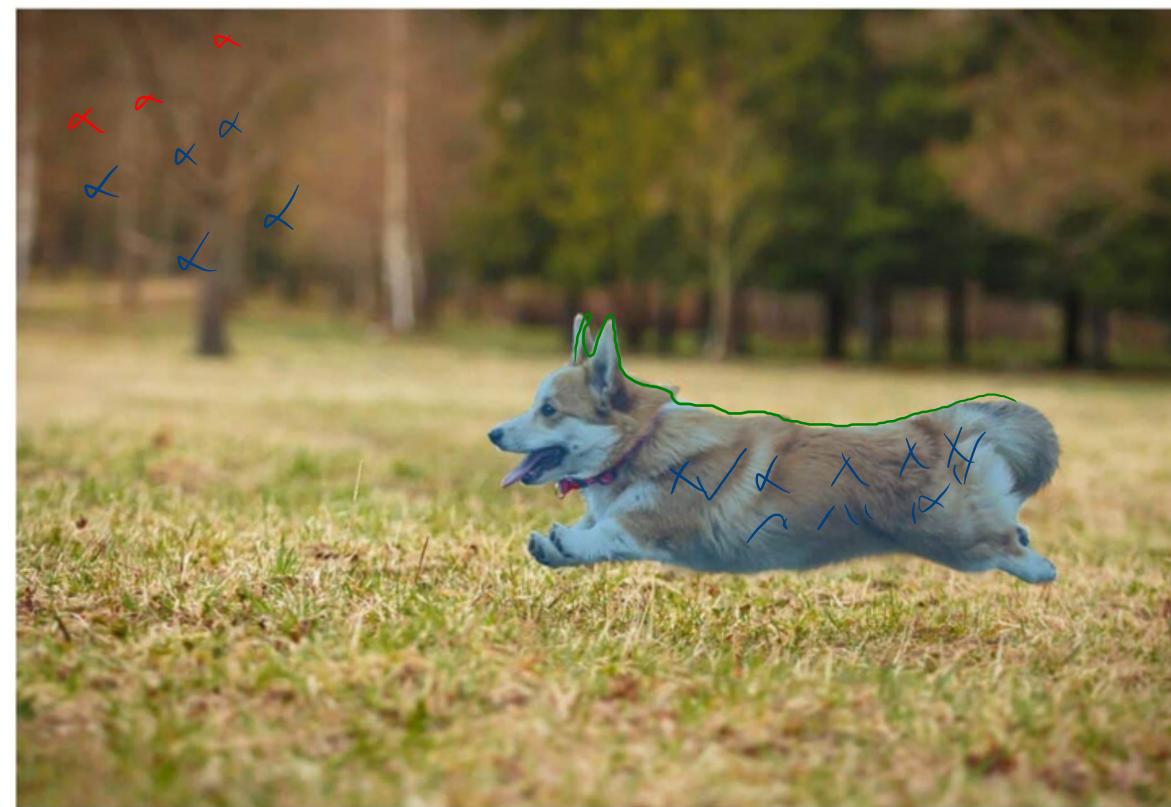
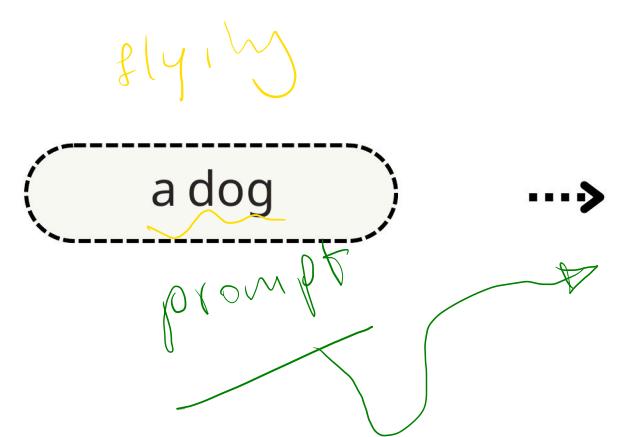
Coarse  
Segmentation



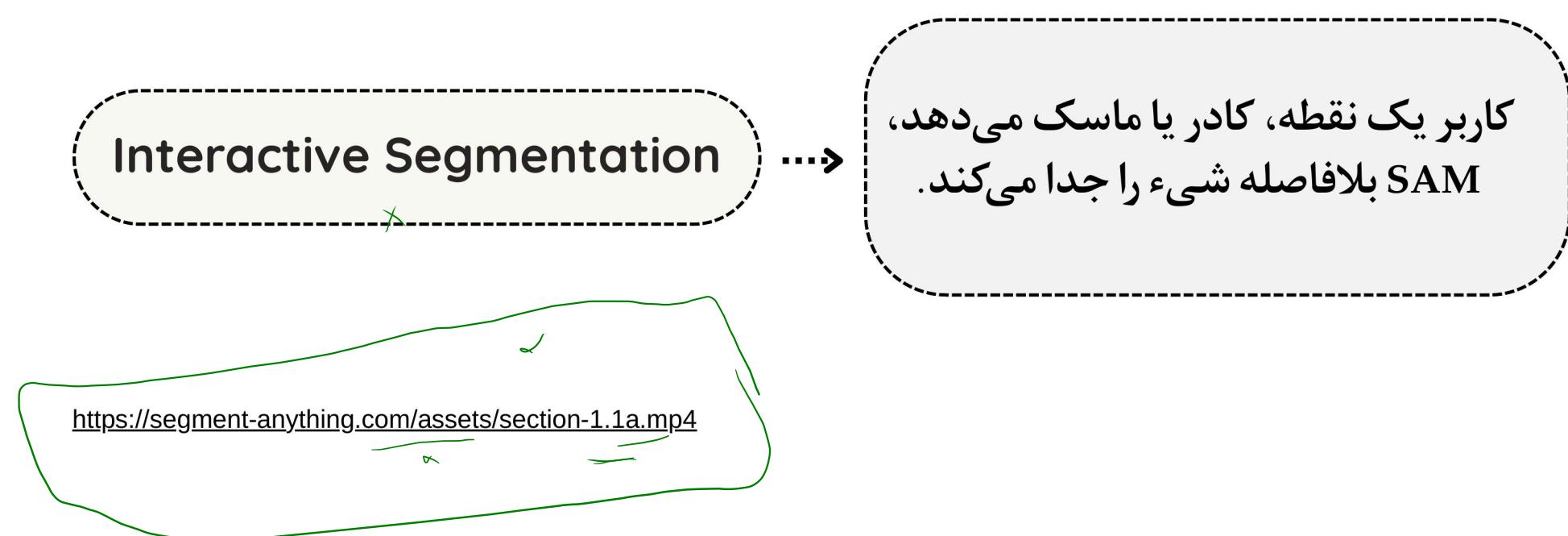
نموده خروجی



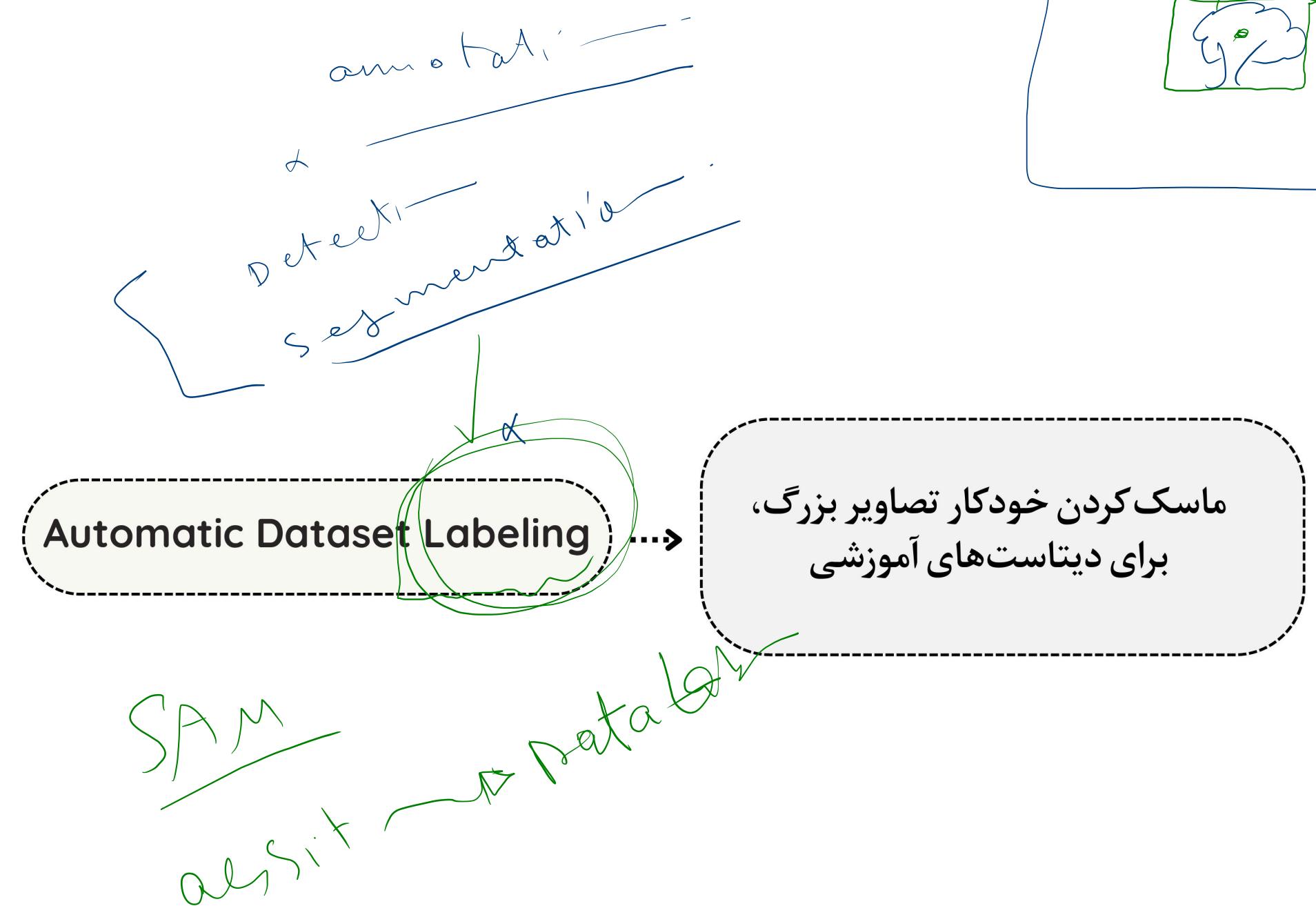
ویژگی های SAM1



# توانمندی های SAM1



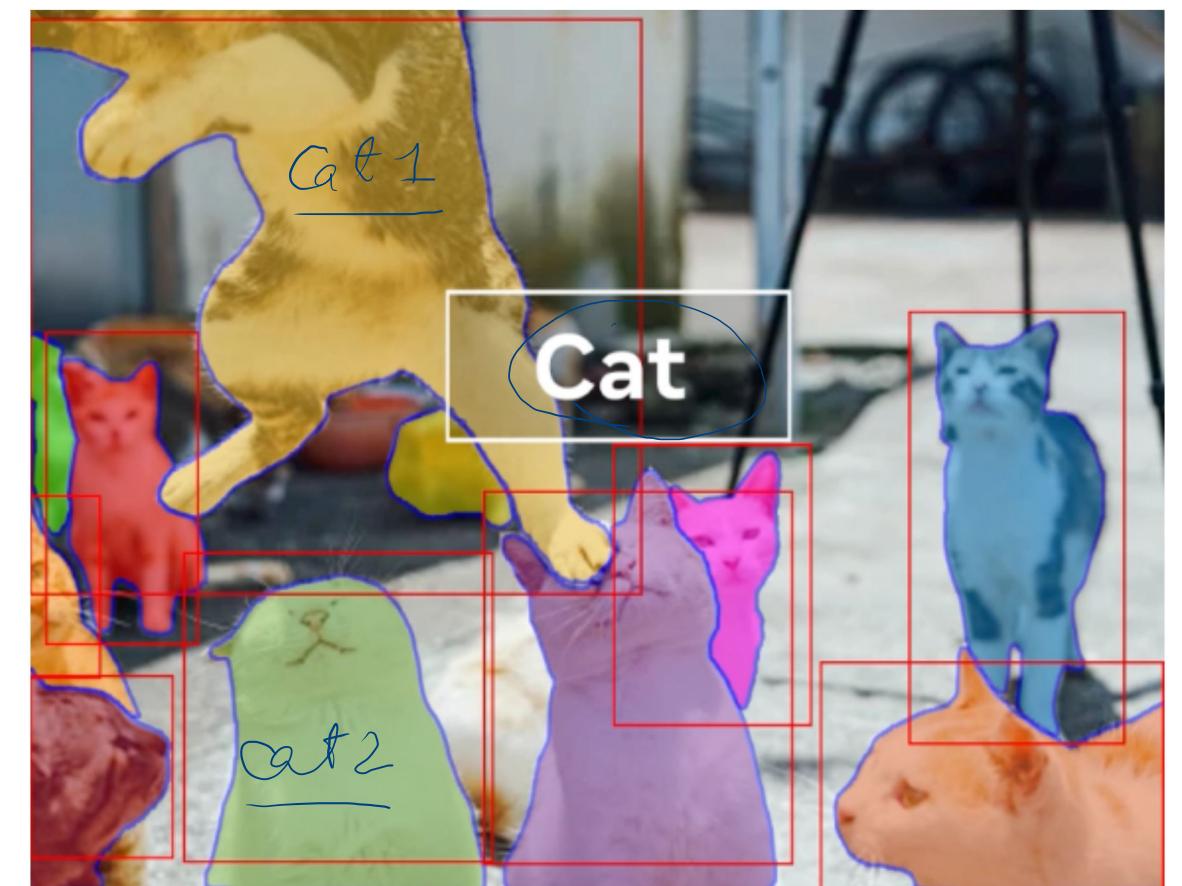
# توانمندی های SAM1



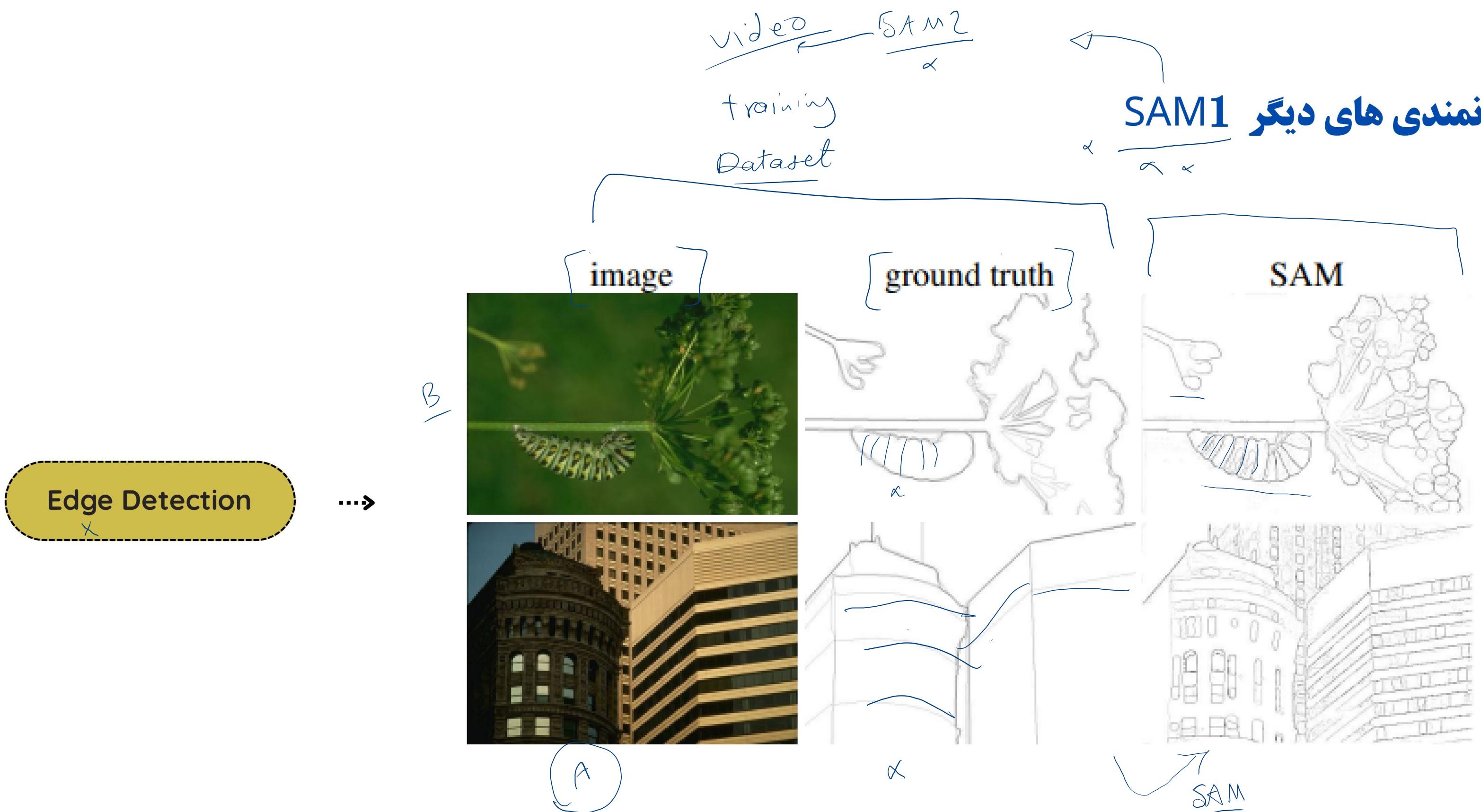
# توانمندی های SAM1

Instance Segmentation

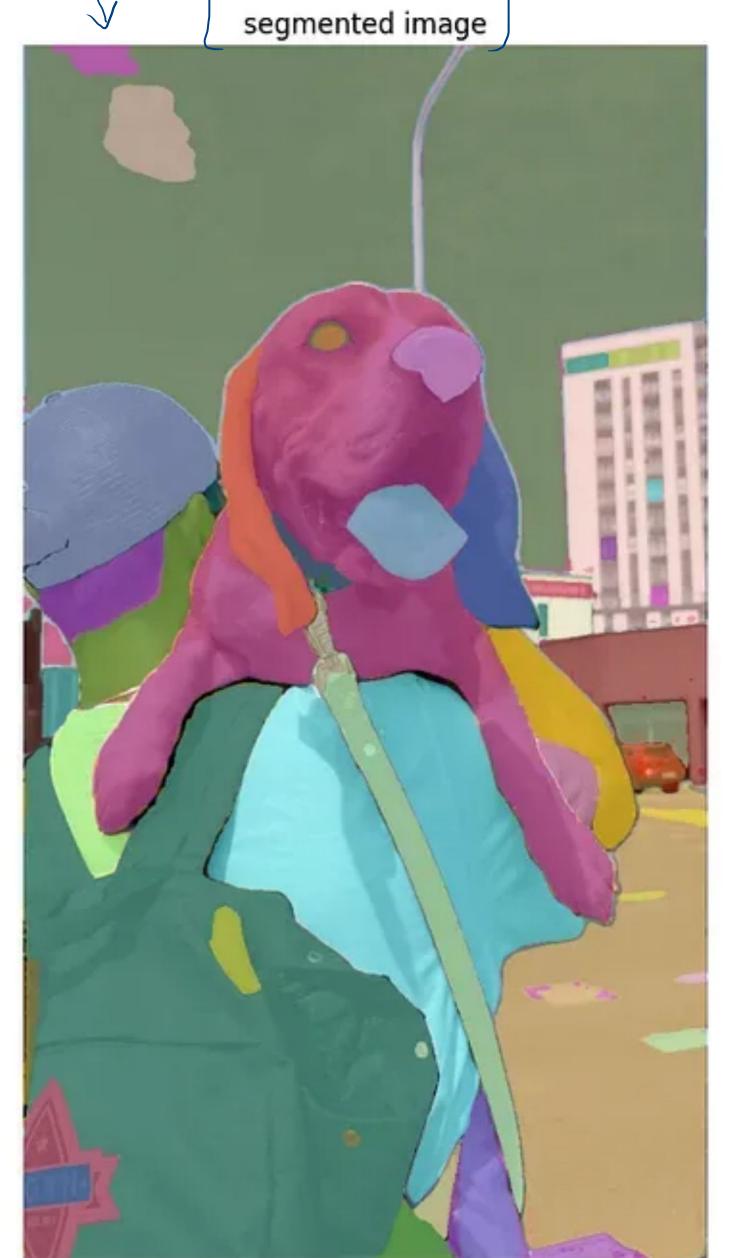
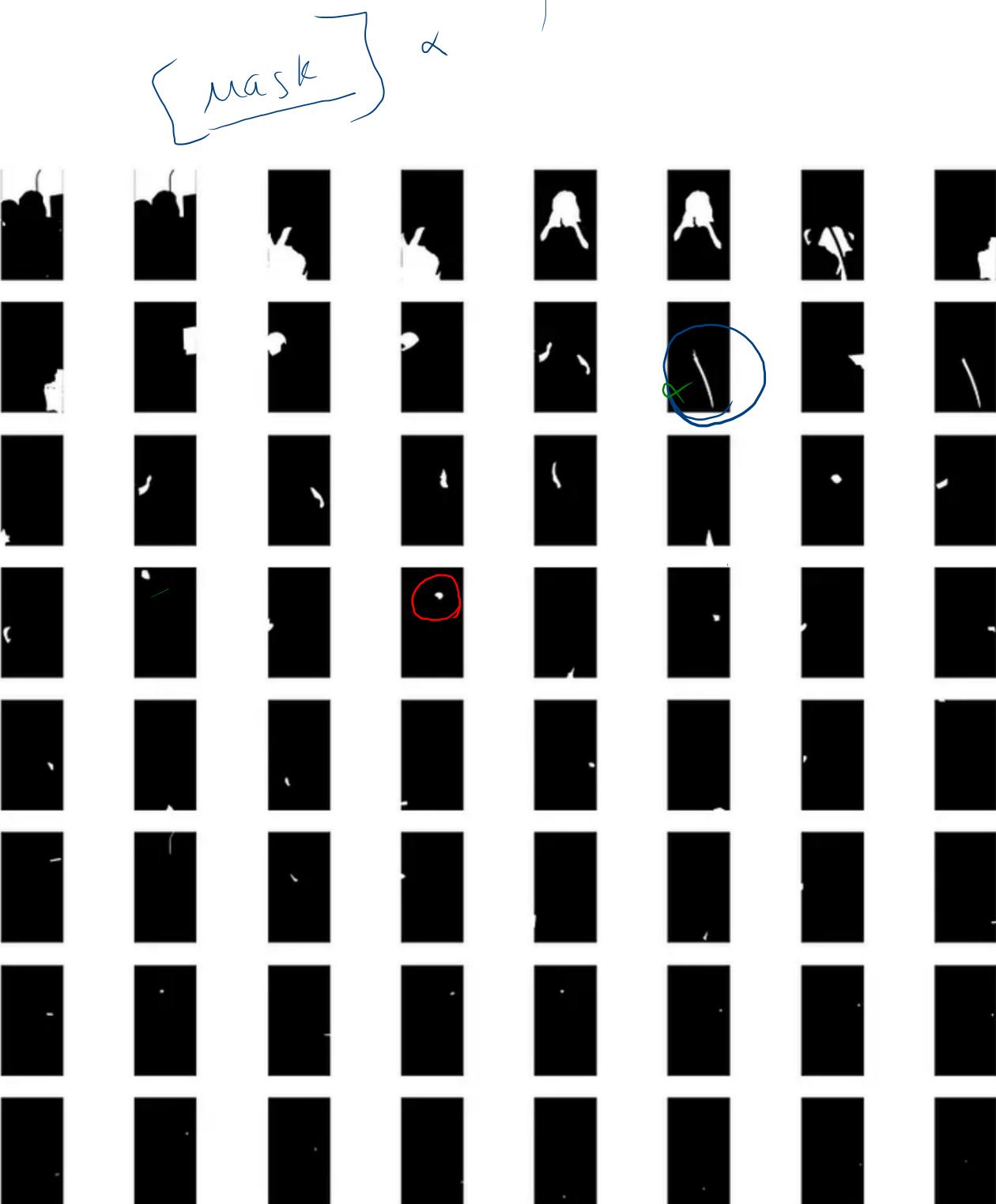
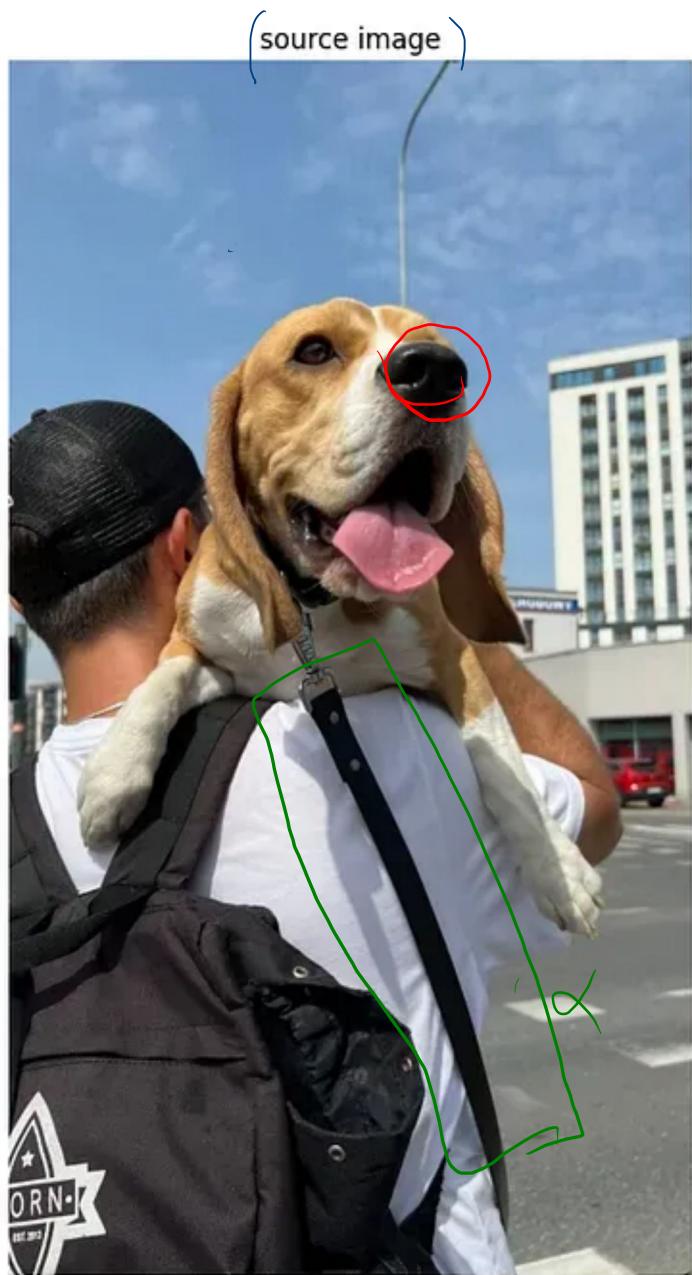
جداسازی اشیاء جداگانه در تصویر  
(مثل جدا کردن هر فرد در جمعیت).



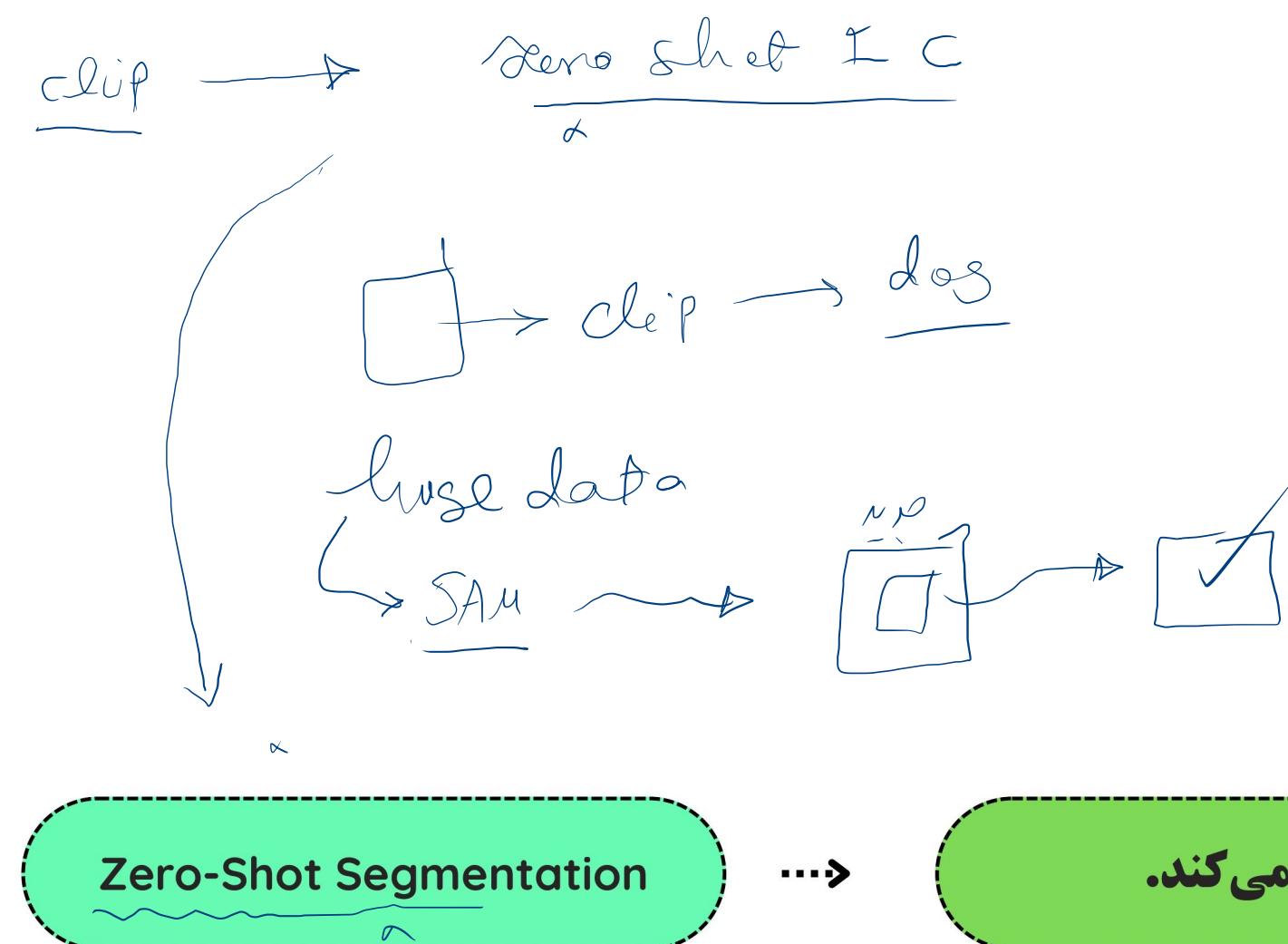
# توانمندی های دیگر



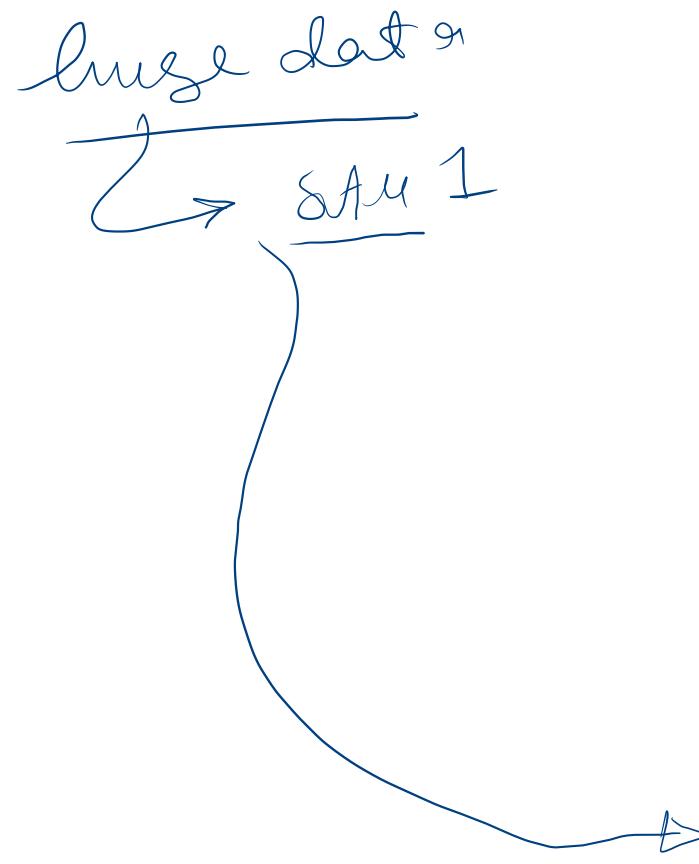
## ماسک در سگمنتیشن



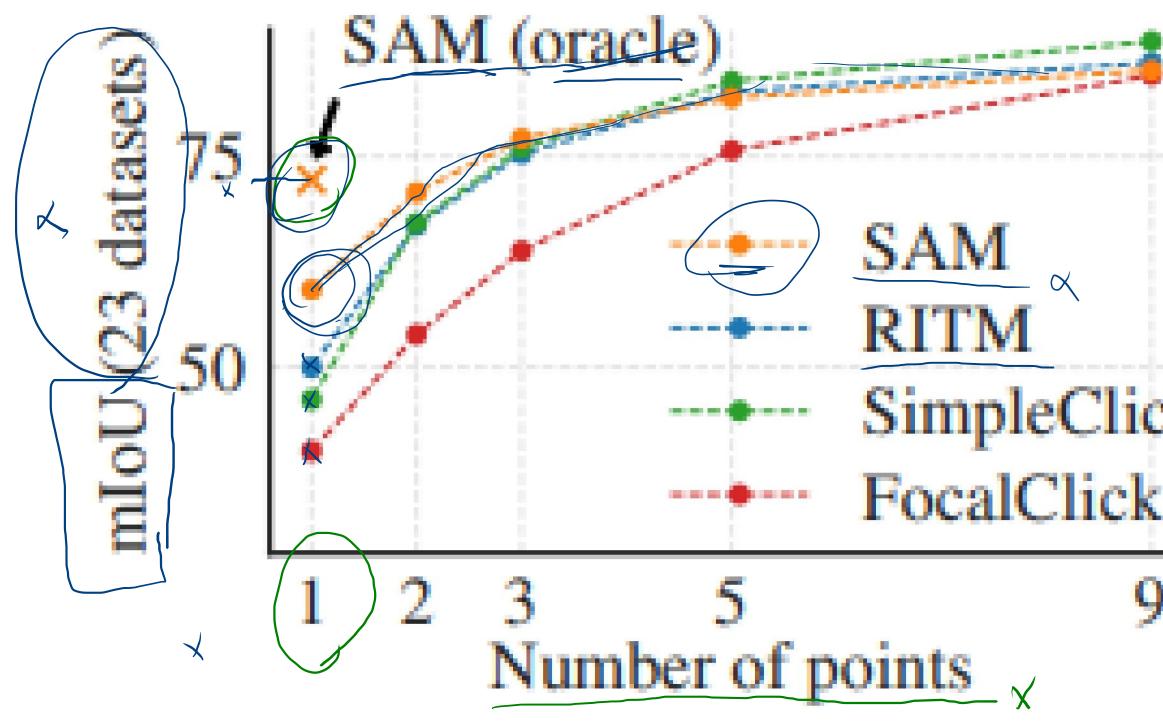
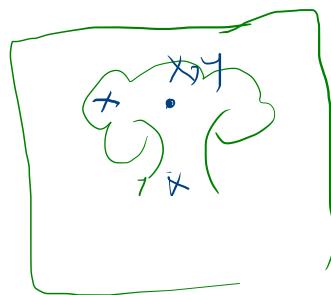
## توانمندی های دیگر SAM1



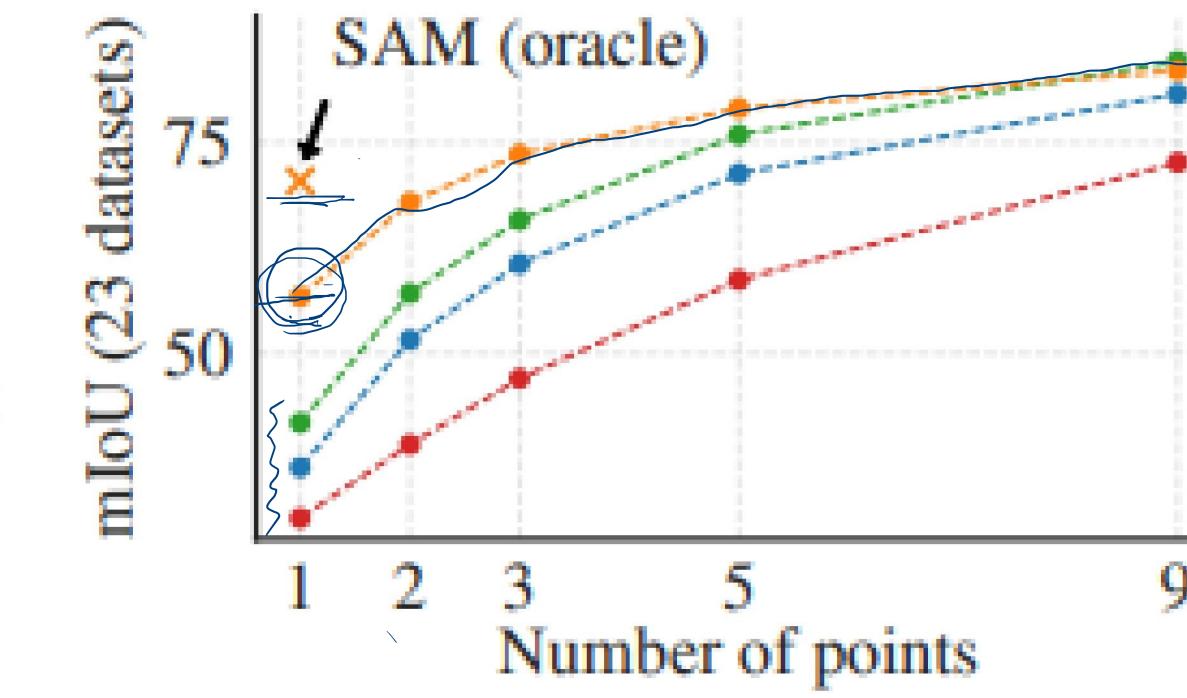
## توانمندی های دیگر SAM1



## مزیت های مدل SAM

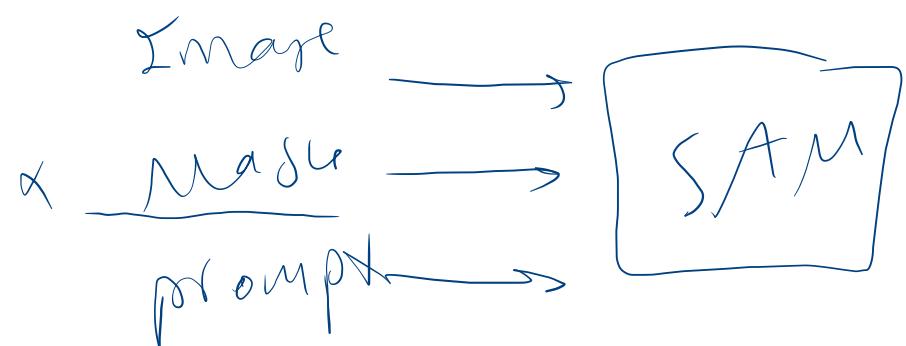
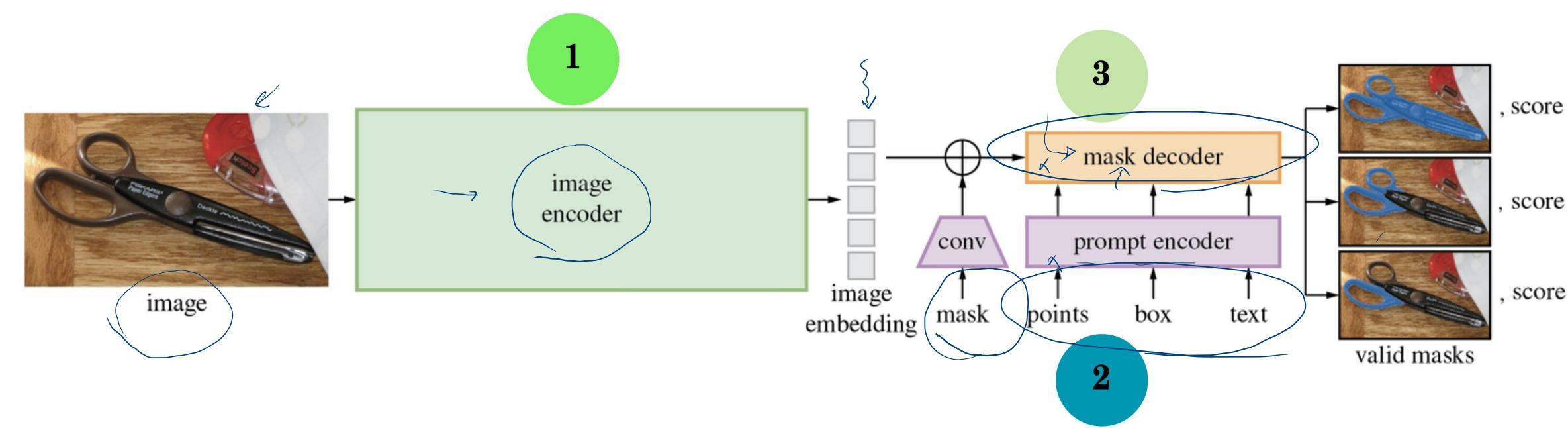


(c) Center points (default)



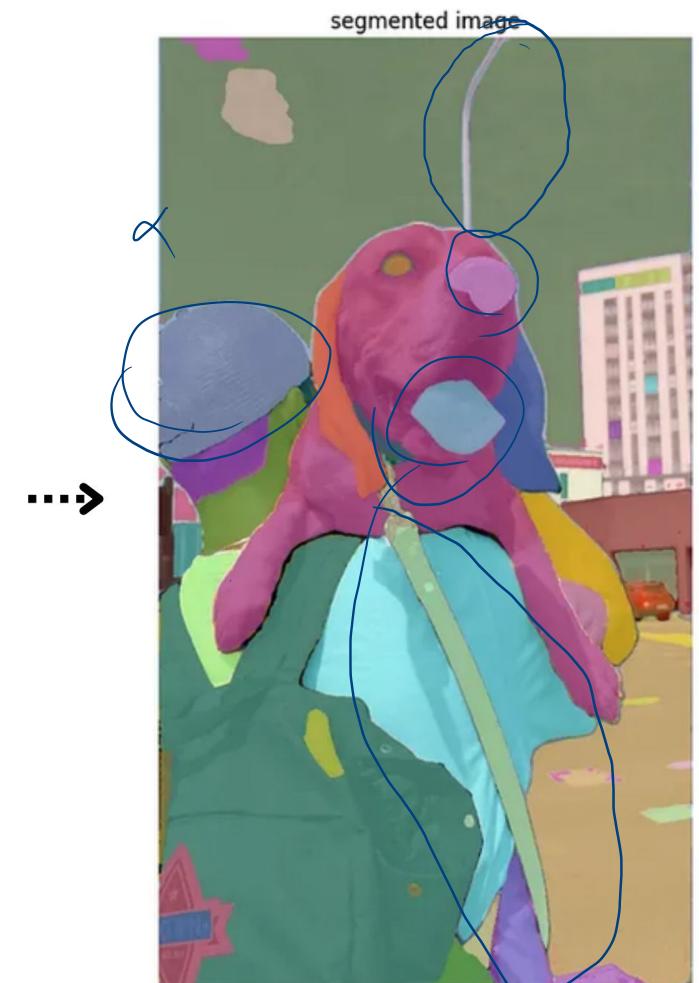
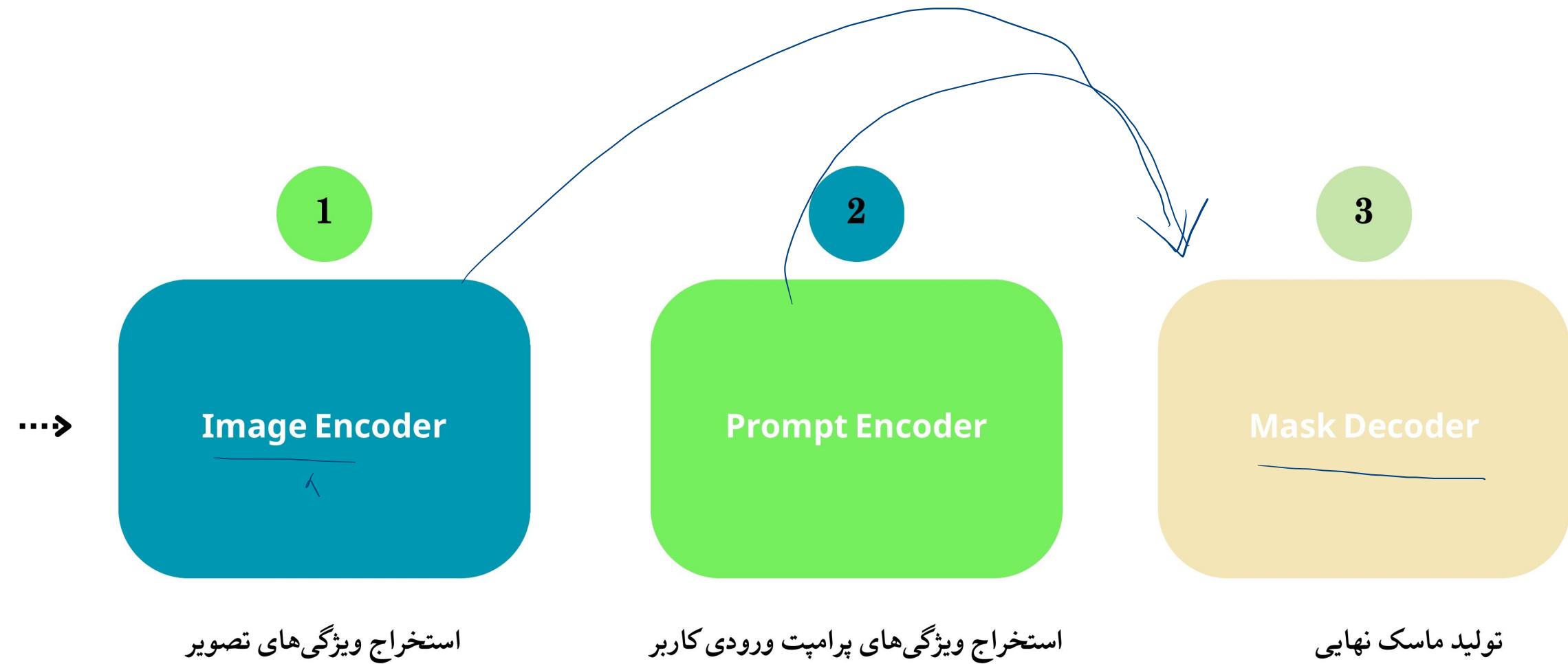
(d) Random points

# ساختار مدل SAM

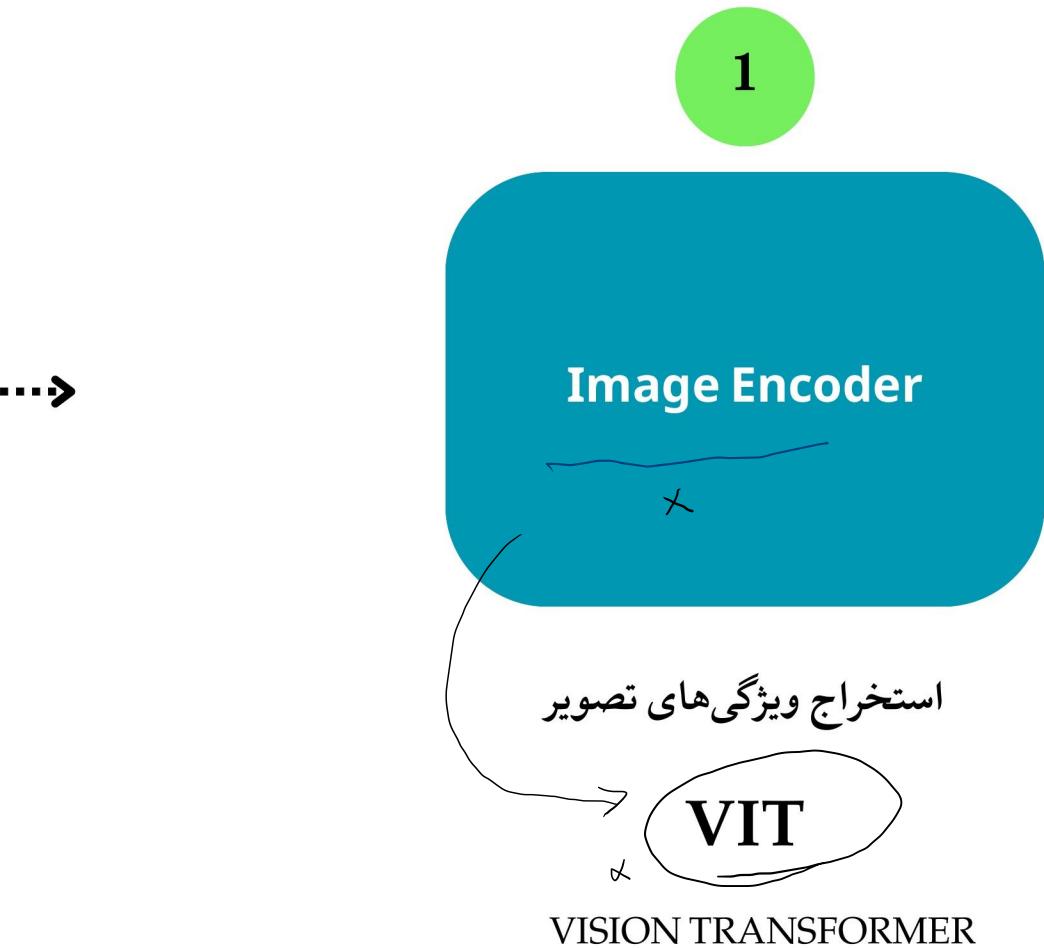
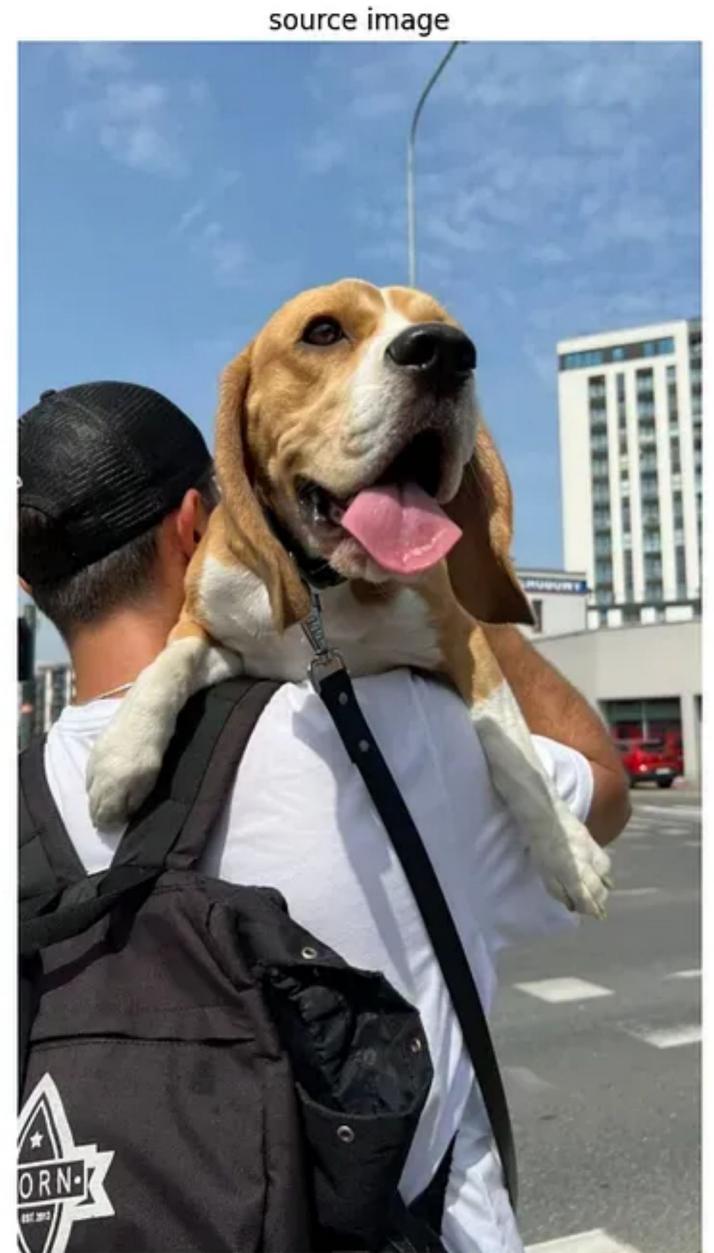


# ساختار مدل SAM

prompt  
magic



# ساختار مدل SAM

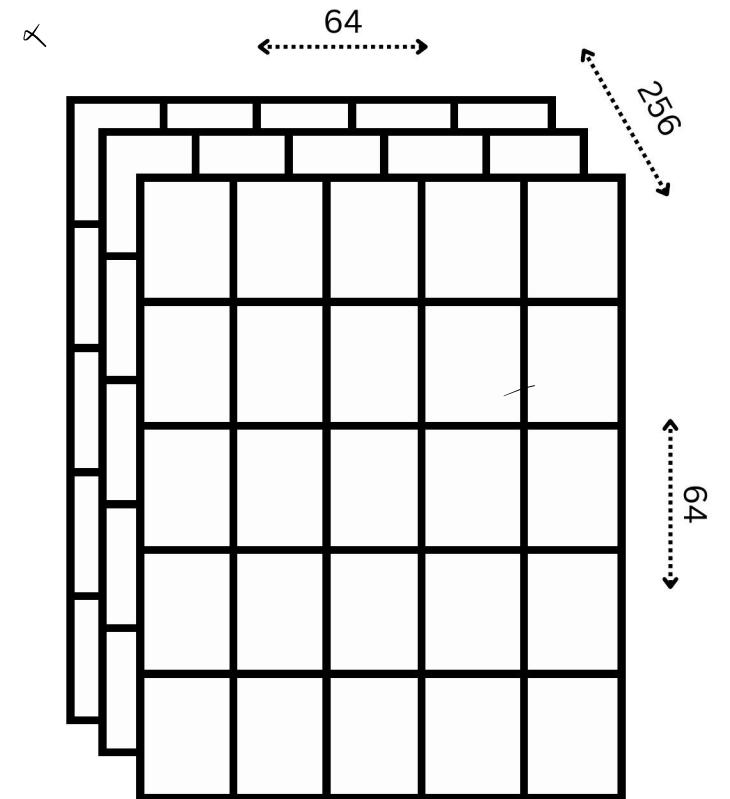


تصویر RGB (مثلاً  $512 \times 512$ )

ورودی

مدل

خروجی



( $256 \times 64 \times 64$  تصویری EMBEDDING مثل)

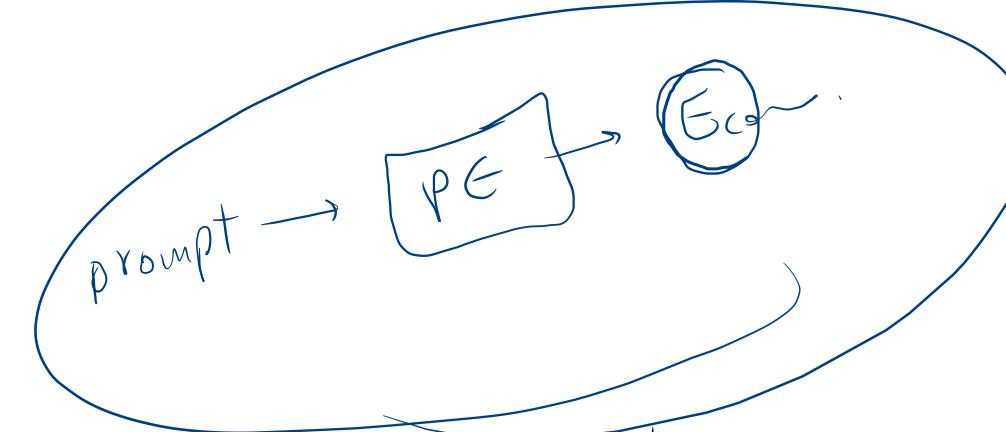
# Prompt Encoder



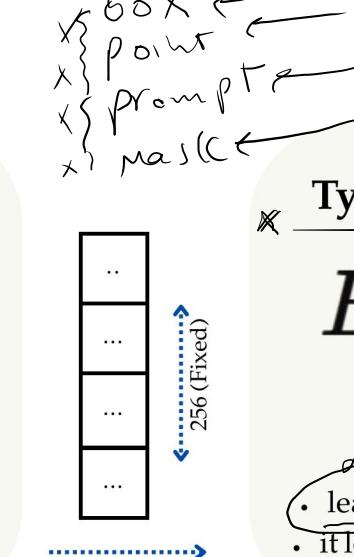
$x_1, y_1, x_2, y_2$

$$(x'_1, y'_1, x'_2, y'_2)$$

$$(x_1, y_1, x_2, y_2)$$



چگونه ورودی Bounding Box انکد می شود؟



1 Normalization([0,1])

- $x_{1,norm} = \frac{x_1}{W}, y_{1,norm} = \frac{y_1}{H}$
- $x_{2,norm} = \frac{x_2}{W}, y_{2,norm} = \frac{y_2}{H}$

\*Resolution Invariant

2 Positional Encoding

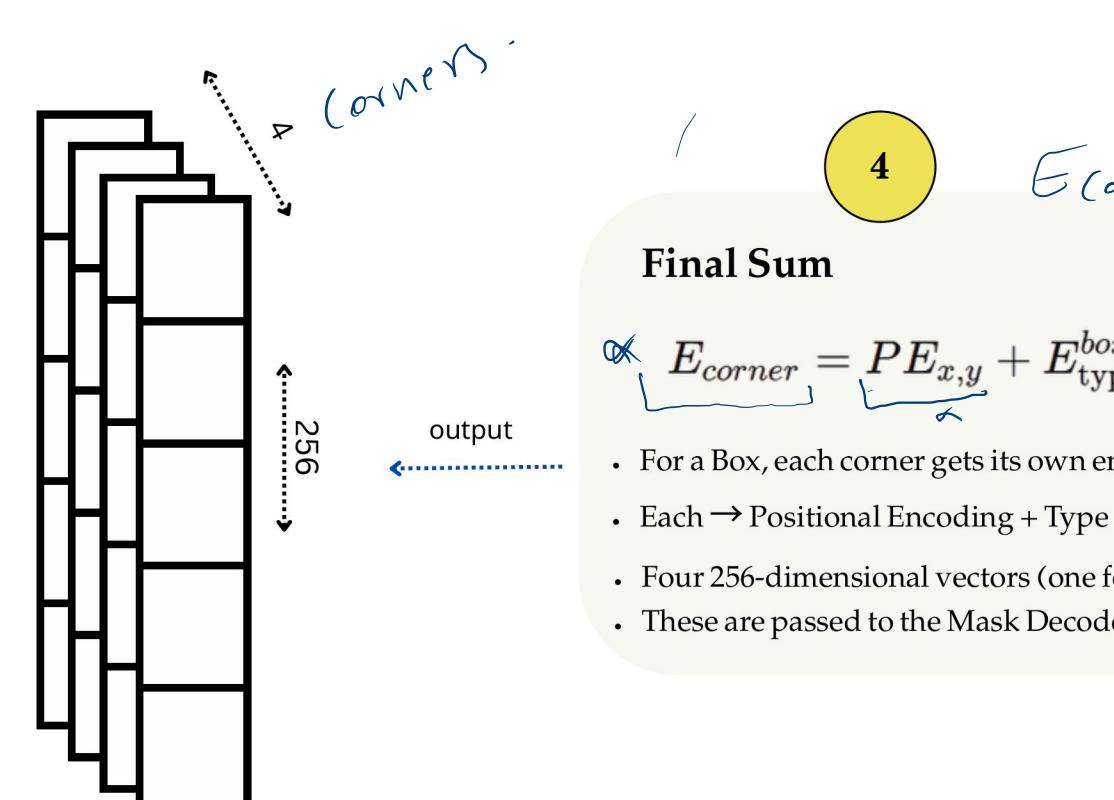
$PE_{x,y} = [\sin(\omega_1 x), \cos(\omega_1 x), \dots, \sin(\omega_d y), \cos(\omega_d y)]$

\*2D Sine-Cosine Positional Encoding  
 $\omega$  = fixed frequency bands.  
 It lets the model "understand" the location

3 Type Embedding

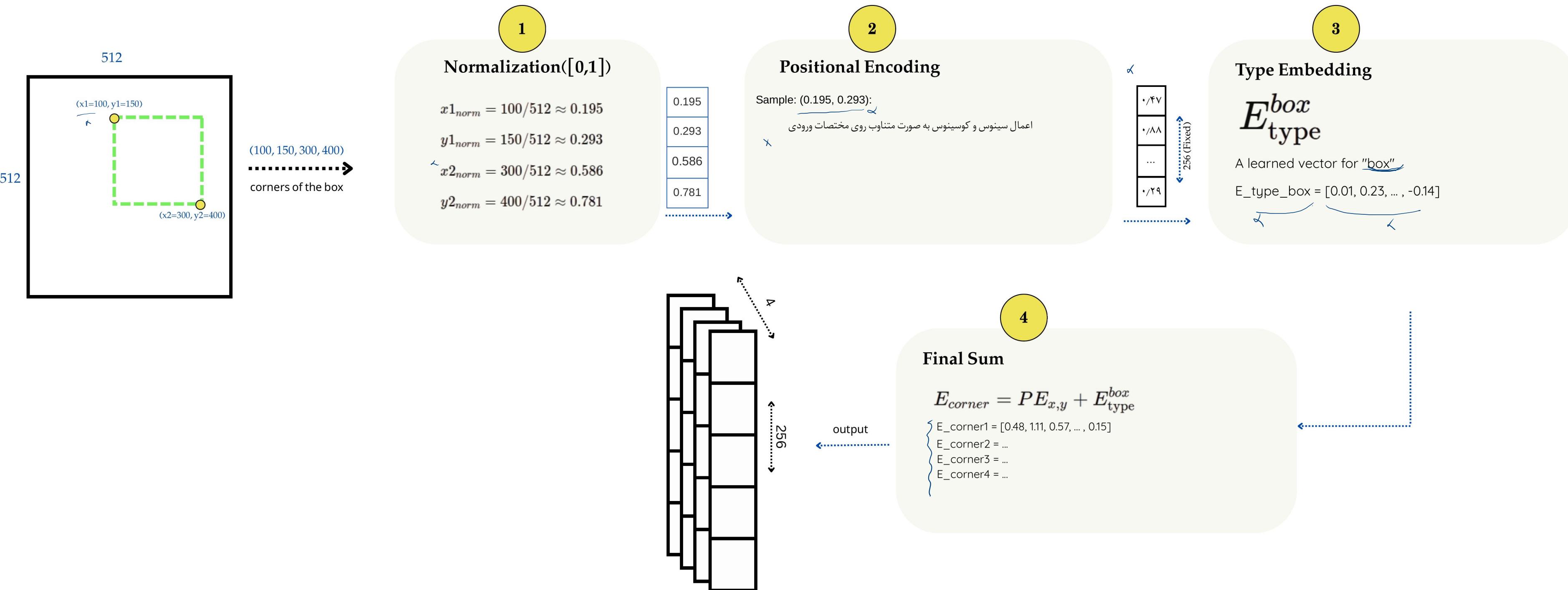
$E_{type}^{box}$

- learned Type Embedding
- it lets the model know "this is a box", not a point or mask



# Prompt Encoder

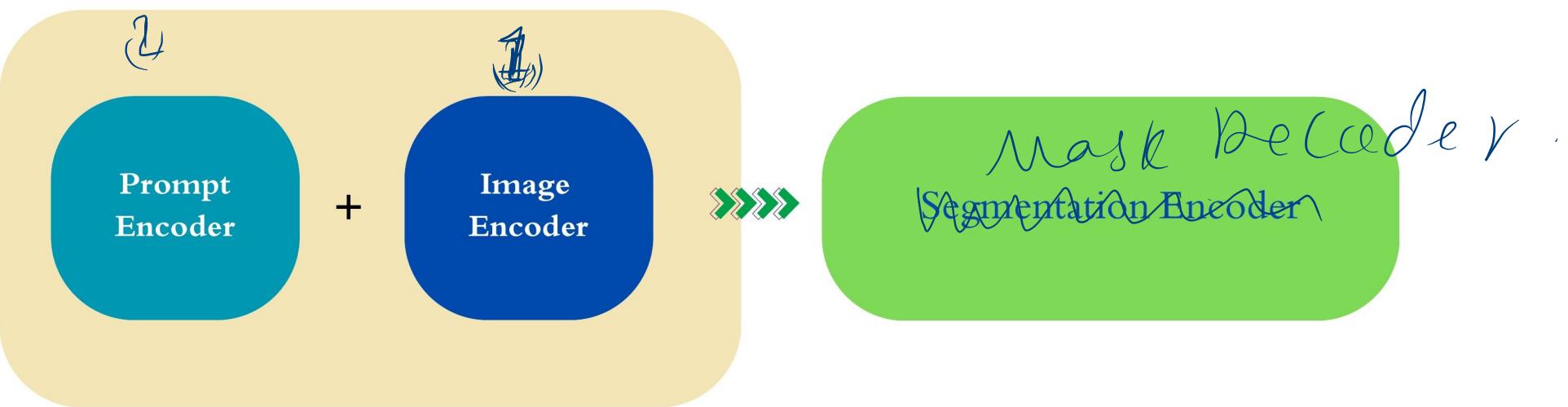
چگونه ورودی Bounding Box انکد می شود؟ (مثال عددی)



# Mask Decoder

3

## Mask Decoder



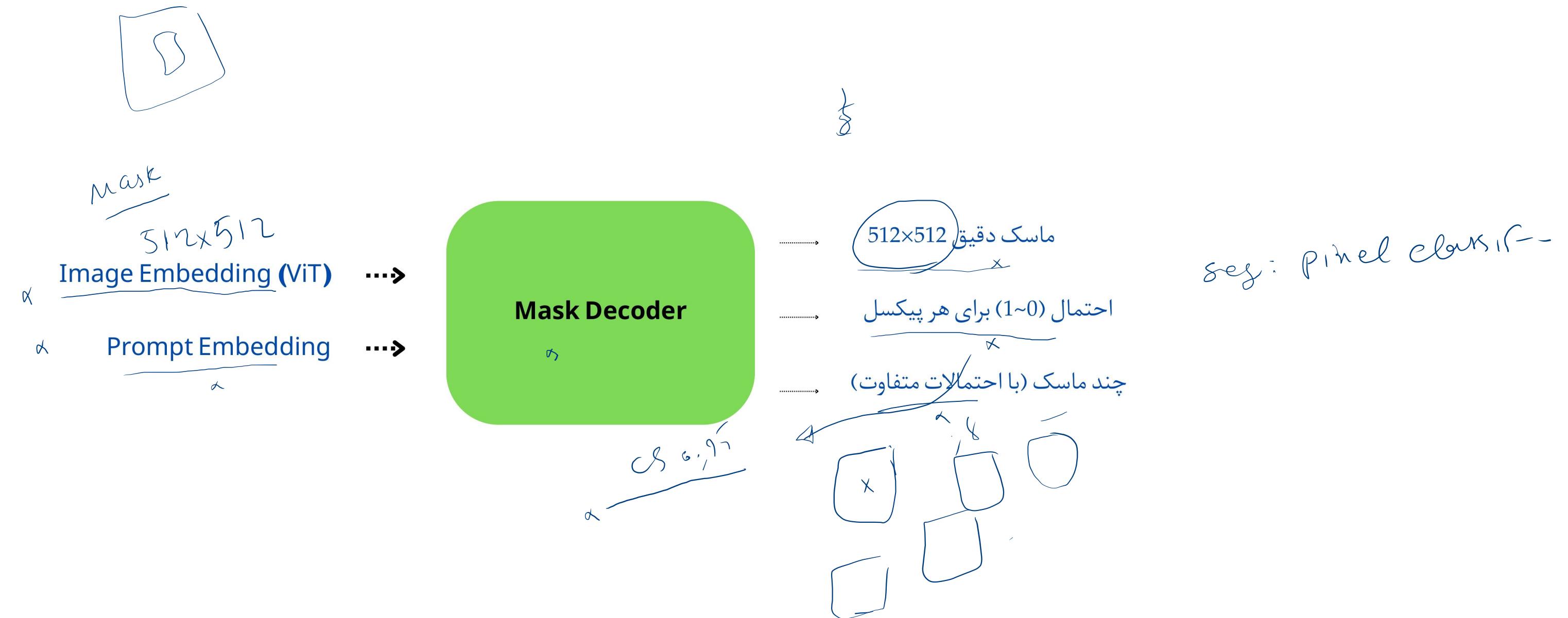
قلب مدل SAM است که  
همهی اطلاعات رو ترکیب  
می کند.

با توجه به دستور کاربر، مرز  
دقیق شیء رو جدا می کند.

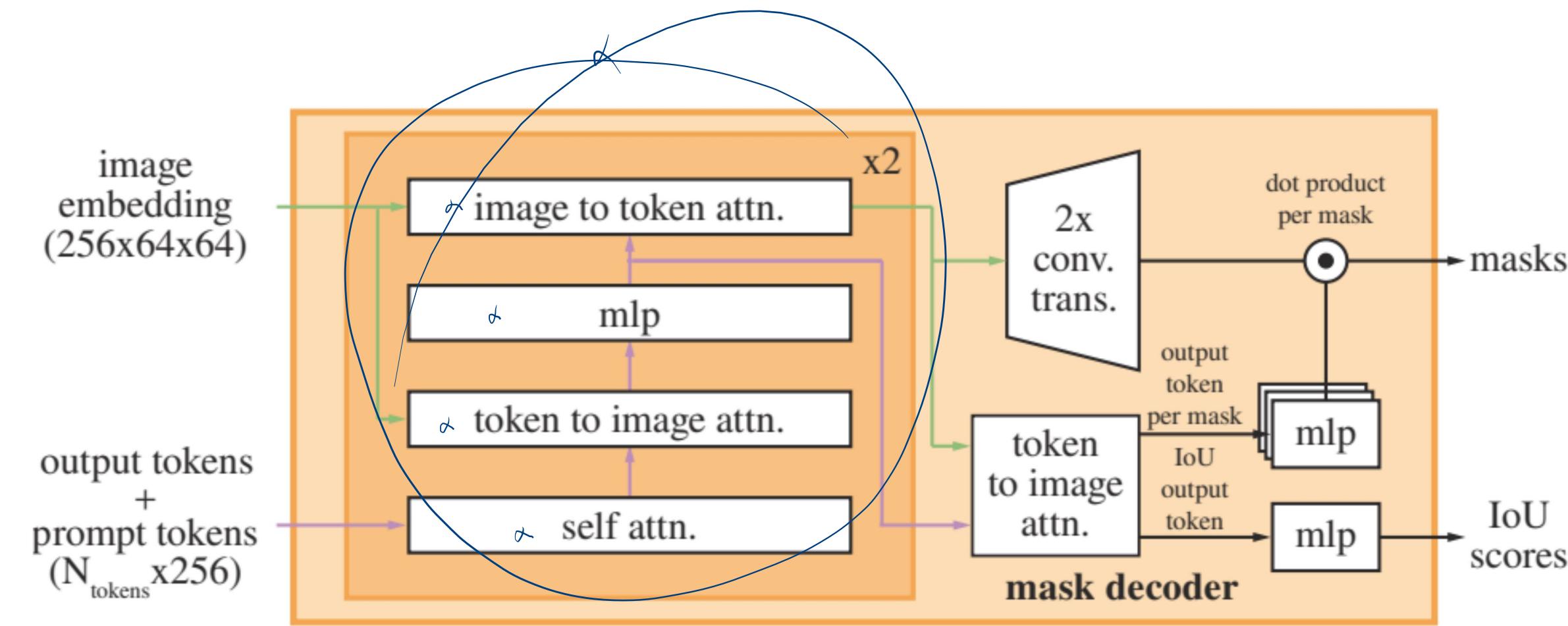
## Mask Encoder



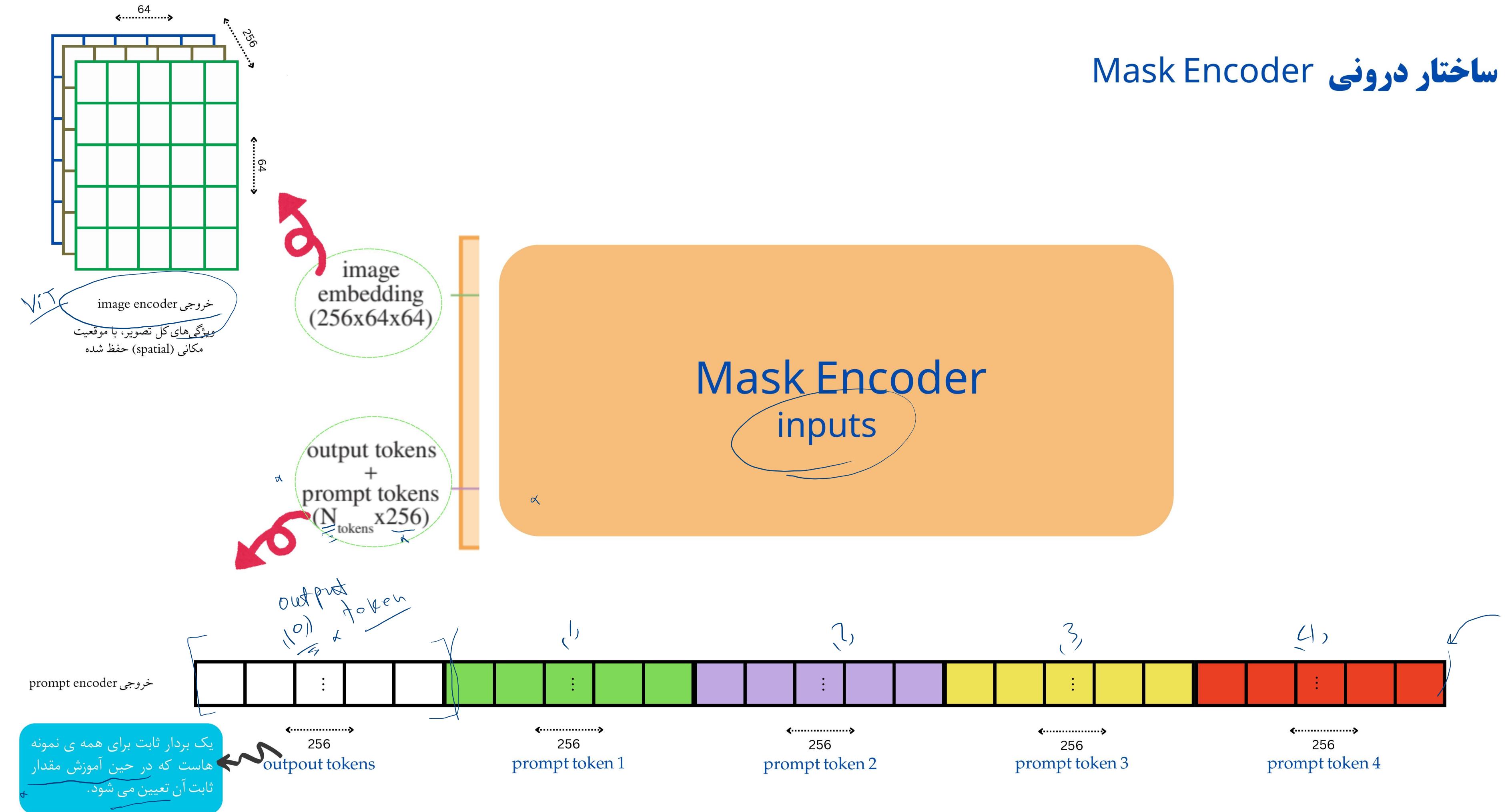
# Mask Encoder



## ساختار درونی Mask Encoder

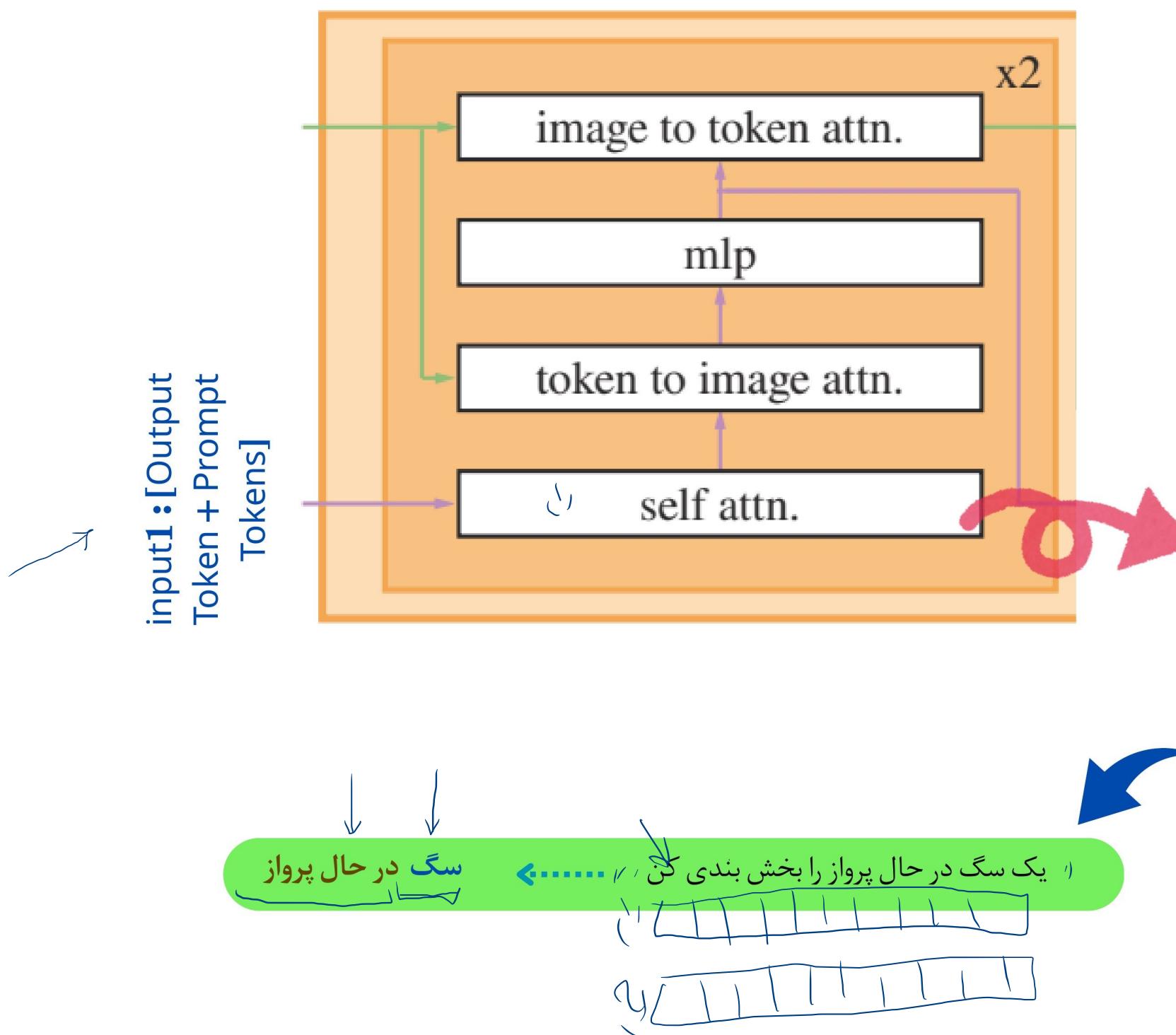


## ساختار درونی Mask Encoder



# ساختار درونی Mask Encoder

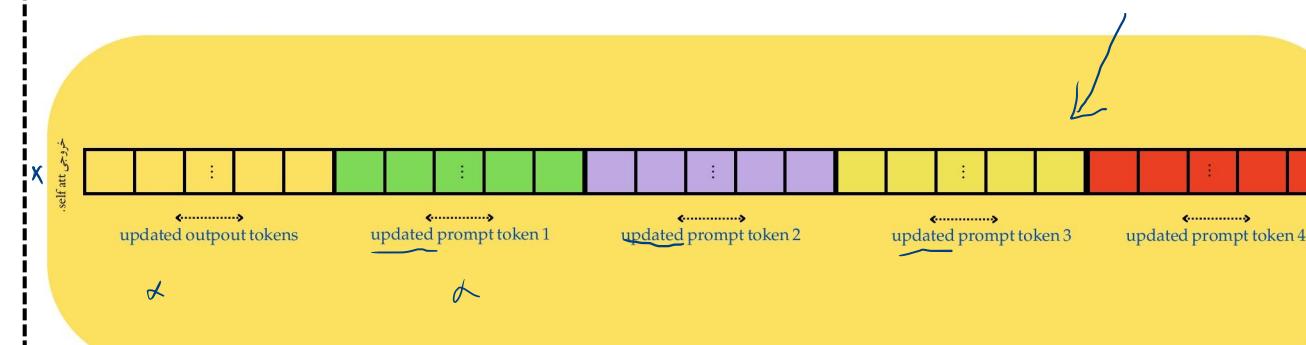
Self Attention



• به اشتراک گذاری اطلاعات بین Output Token و با Prompt tokens و مشخص کردن نقاط مهم پرامپت ورودی و توجه بیشتر به آن

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

- Query (Q) = [Output Token + Prompt Tokens]
- Key (K) = [Output Token + Prompt Tokens]
- Value (V) = [Output Token + Prompt Tokens]

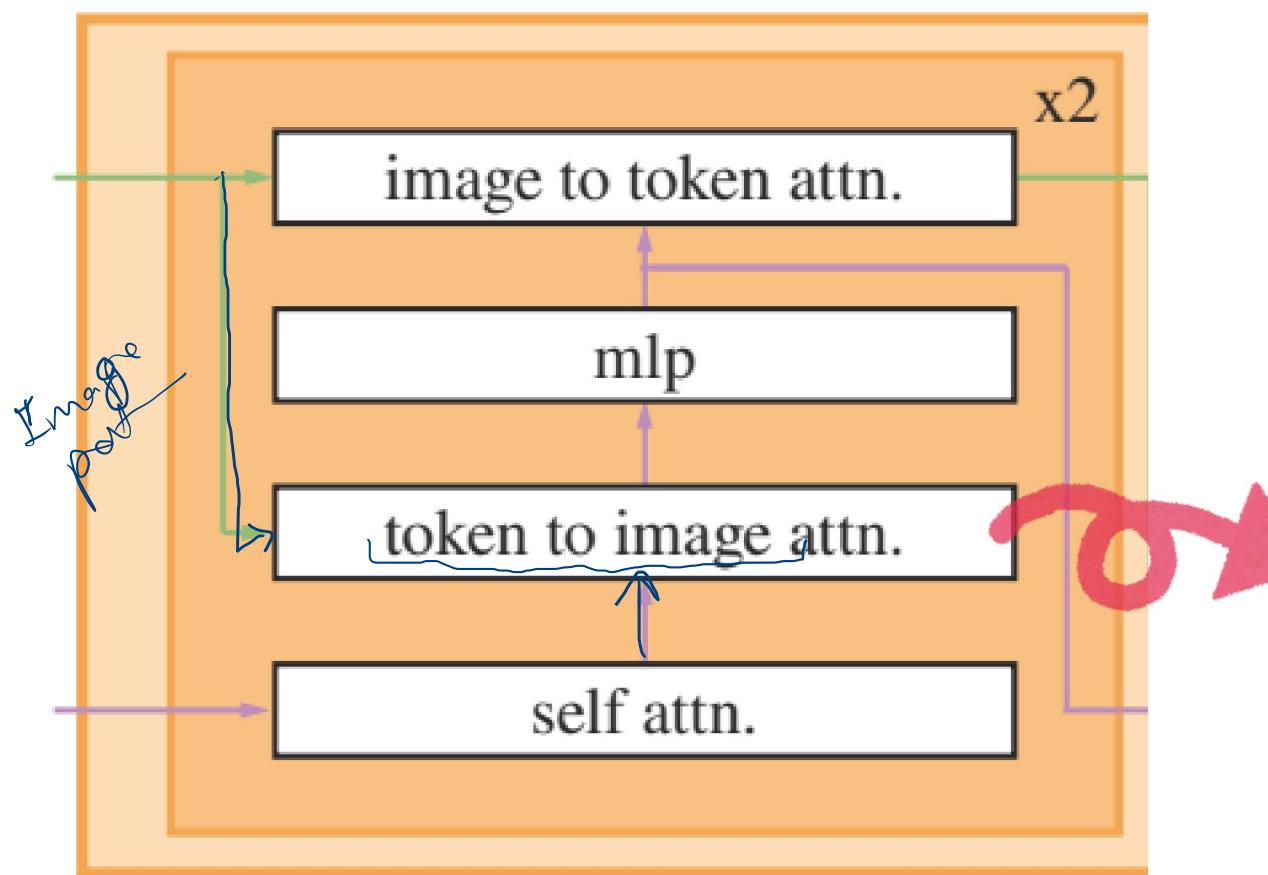


درباره self attention

خروجی بخش self attention

# ساختار درونی Mask Encoder

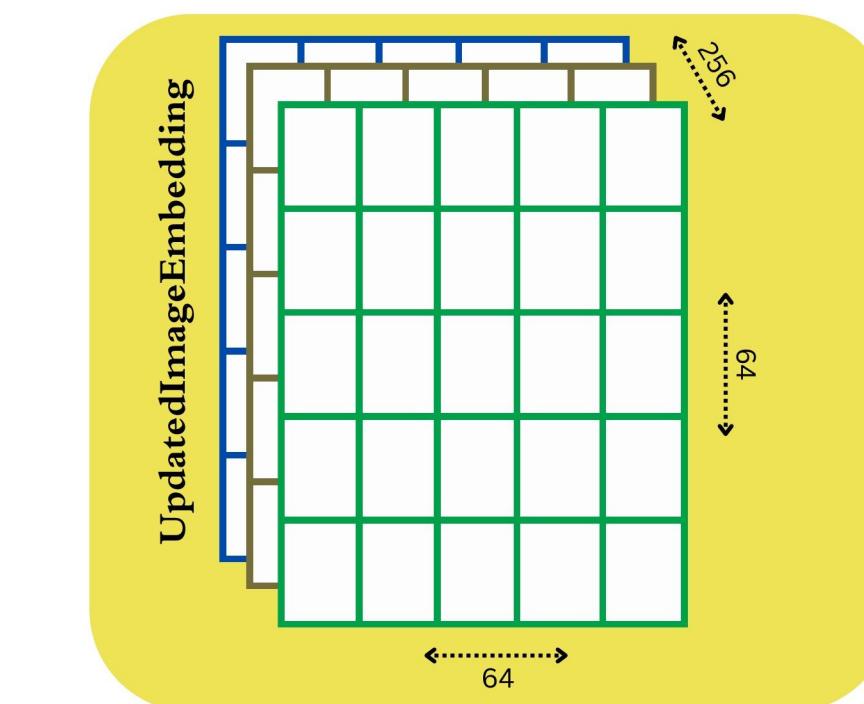
Token to Image Att.



درباره token to image att.  
چه چیزها و نکاتی مهمند. به نوعی در این بخش پرامپت نقش محدود کننده و شرط (Condition) را برای تصویر بازی میکند.

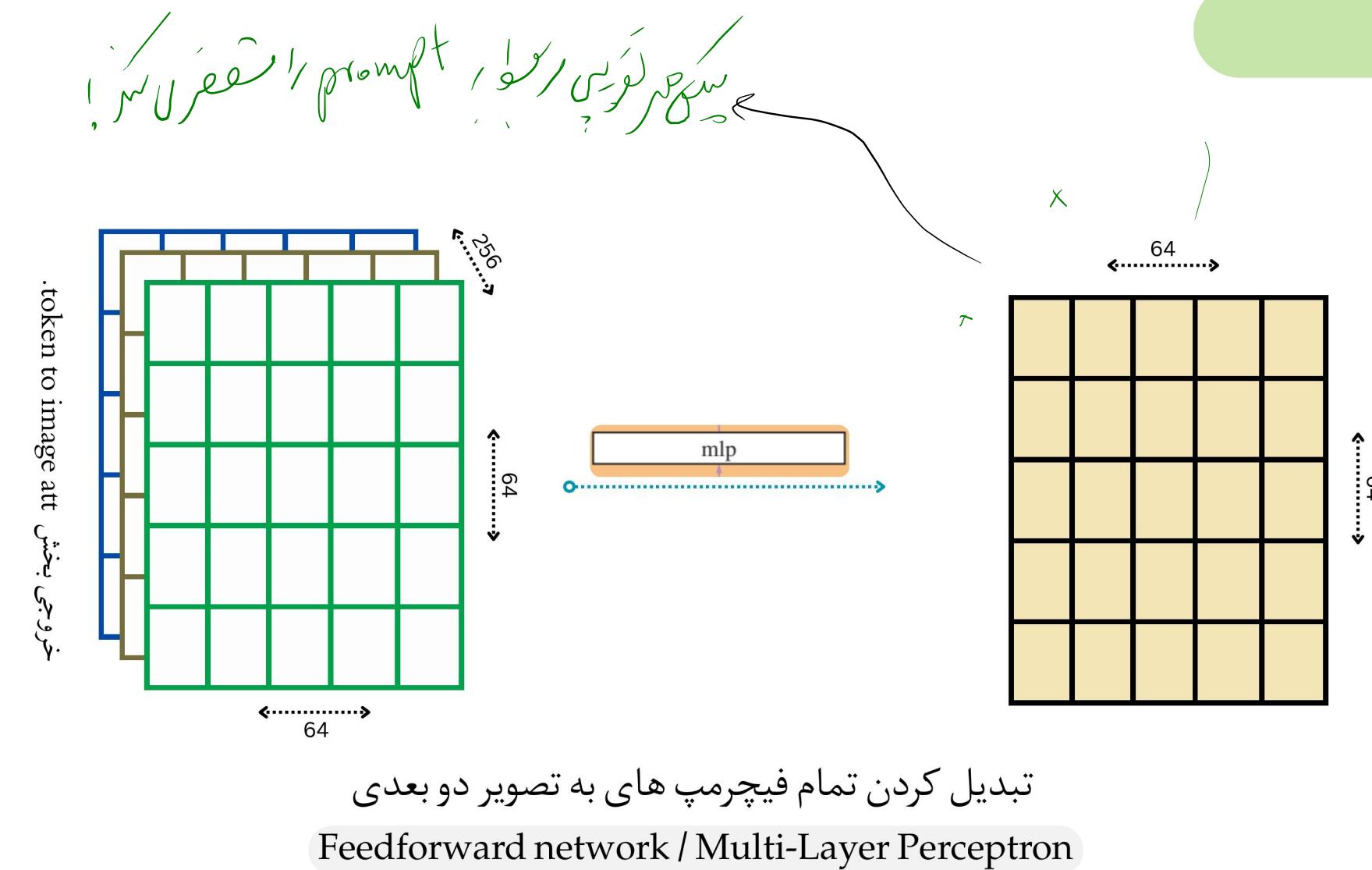
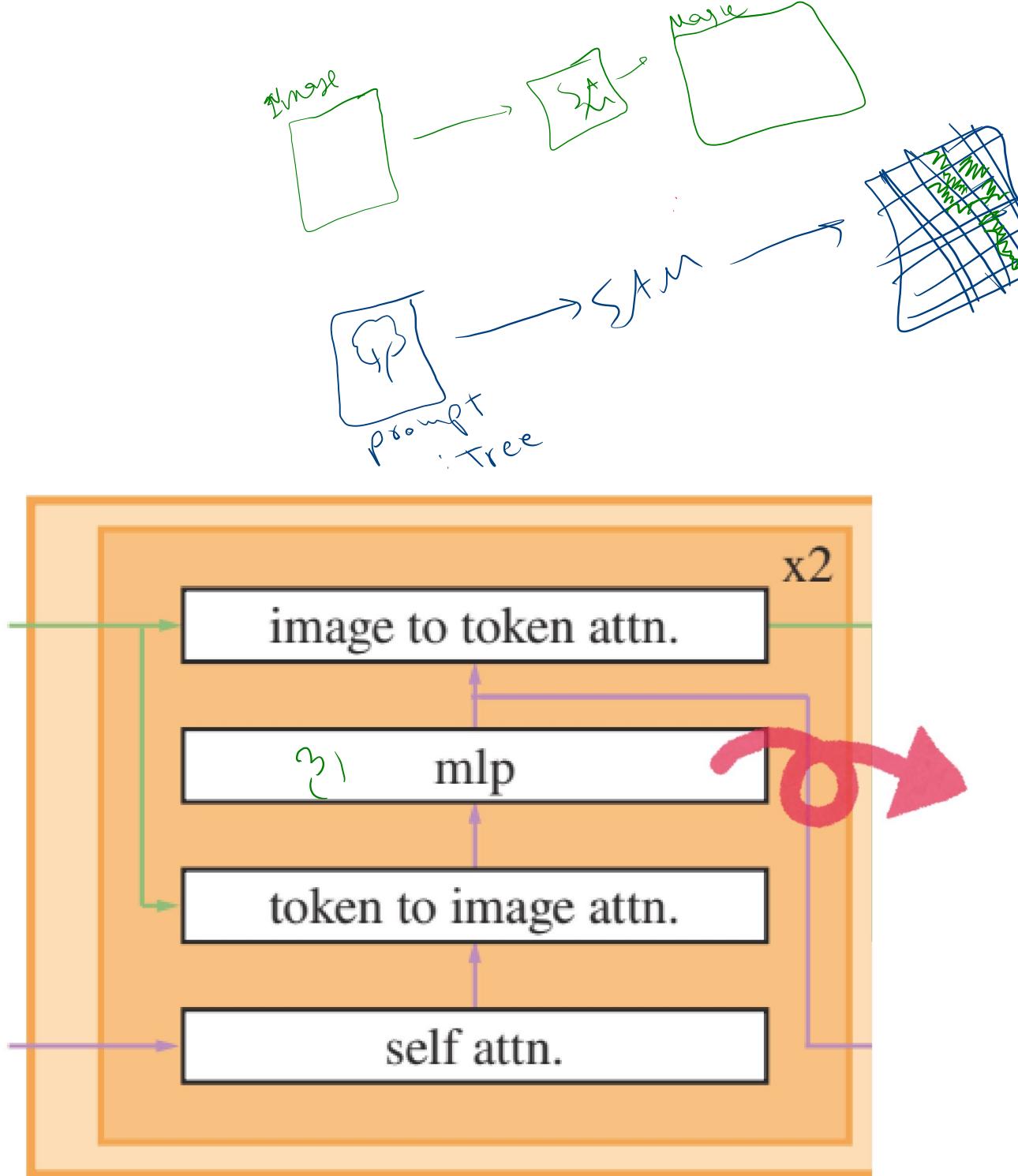
$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V$$

- Query (Q) = Image Embedding
  - Key (K) = Prompt+Output Tokens
  - Value (V) = Prompt+Output Tokens
- self att. *is* self att. *ge?*



خرجی بخش token to image att

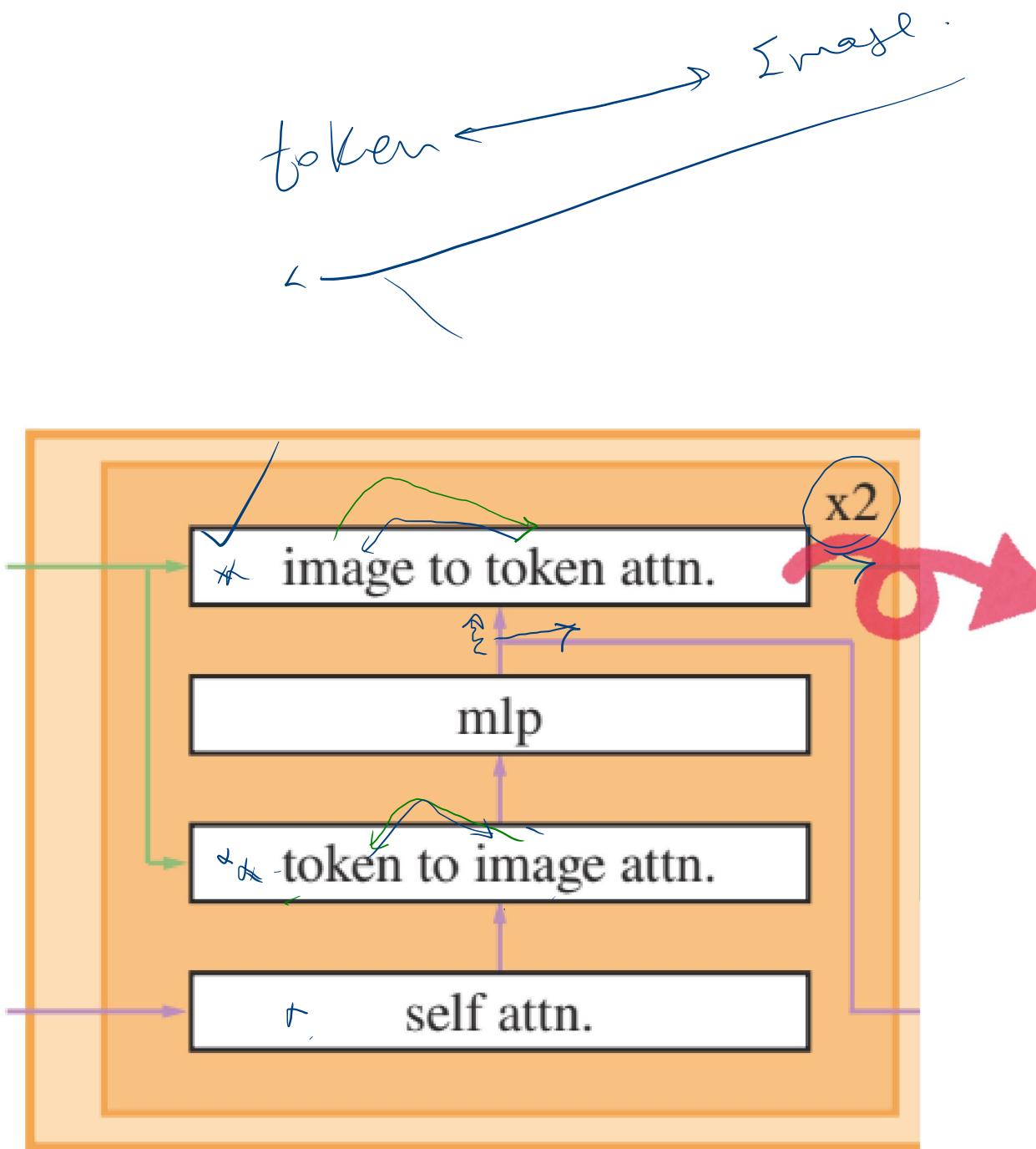
# ساختار درونی Mask Encoder



این بخش مشخص می کند که این Prompt مربوط به کدام قسمت تصویر است.

# ساختار درونی Mask Encoder

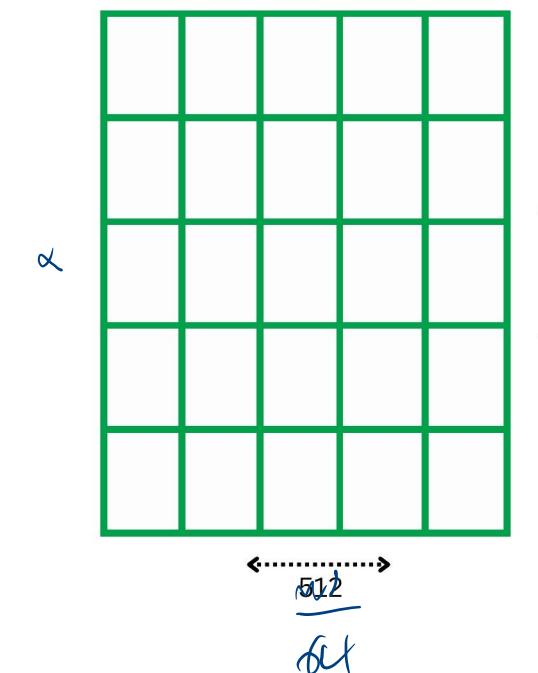
Image to token Att.



- در این بخش مدل با استفاده از Cross-Attention یاد می‌گیرد که کدام بخش تصویر به Prompt ورودی مربوط است.

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

- Query (Q) = output of mlp (updated prompt token) - 256 D Vector
- Key (K) = Image Token (ViT patches)
- Value (V) = Image Token (ViT patches)

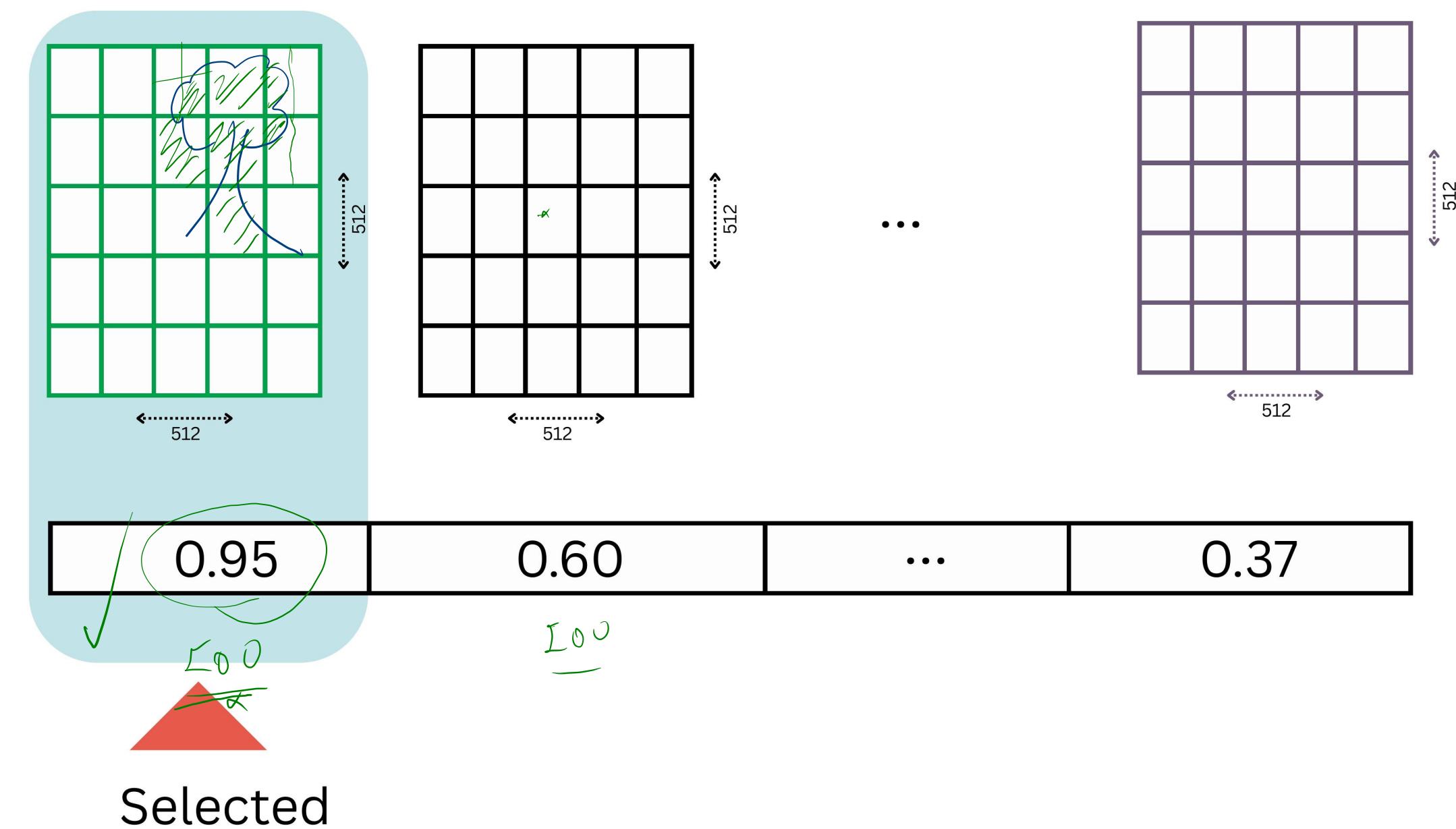
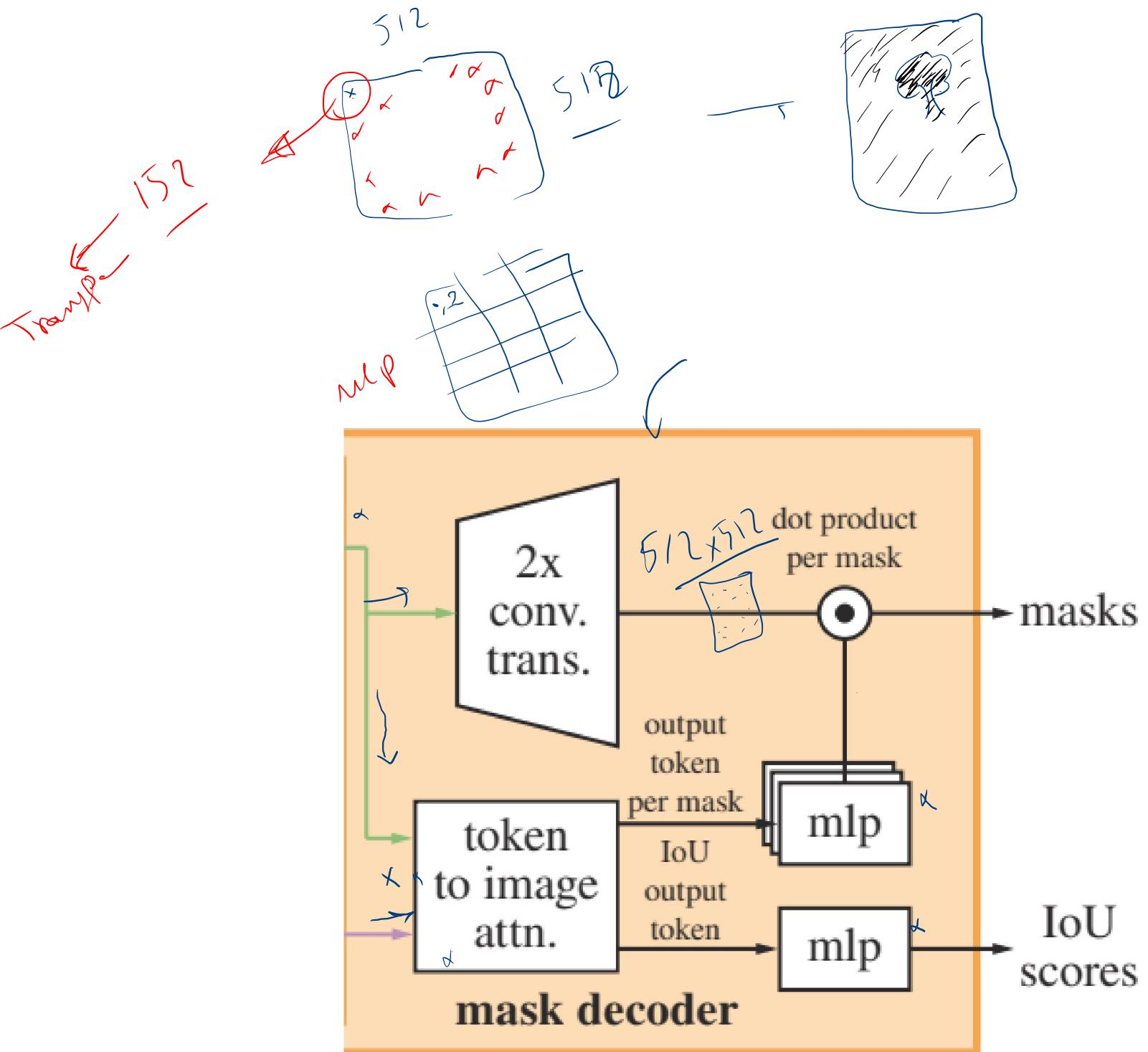


نحوی بحثی .image to token att

درباره .image to token att

# ساختار درونی Mask Encoder

Image to token Att.



The end

x Seiji