

موضوع

YOLO v12

قهرمان تشخیص اشیا :)

YOLOv12: Attention-Centric Real-Time Object Detectors



Yunjie Tian
University at Buffalo
yunjieti@buffalo.edu

Qixiang Ye
University of Chinese Academy of Sciences
qxye@ucas.ac.cn

David Doermann
University at Buffalo
doermann@buffalo.edu

موضع

YOLO v12

قهرمان تشخيص اشیا :)

YOLOv12: Attention-Centric Real-Time Object Detectors

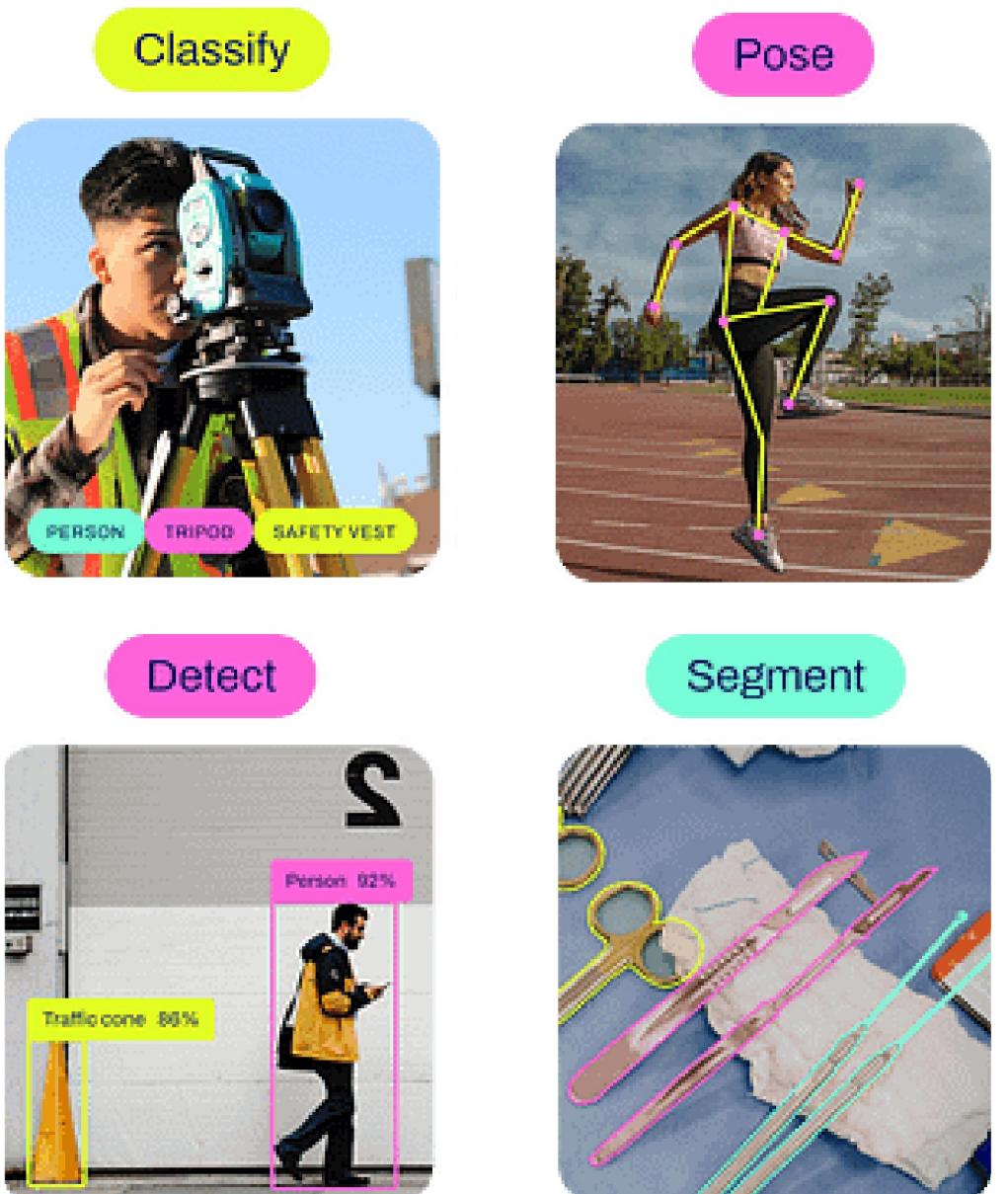


Yunjie Tian
University at Buffalo
yunjieti@buffalo.edu

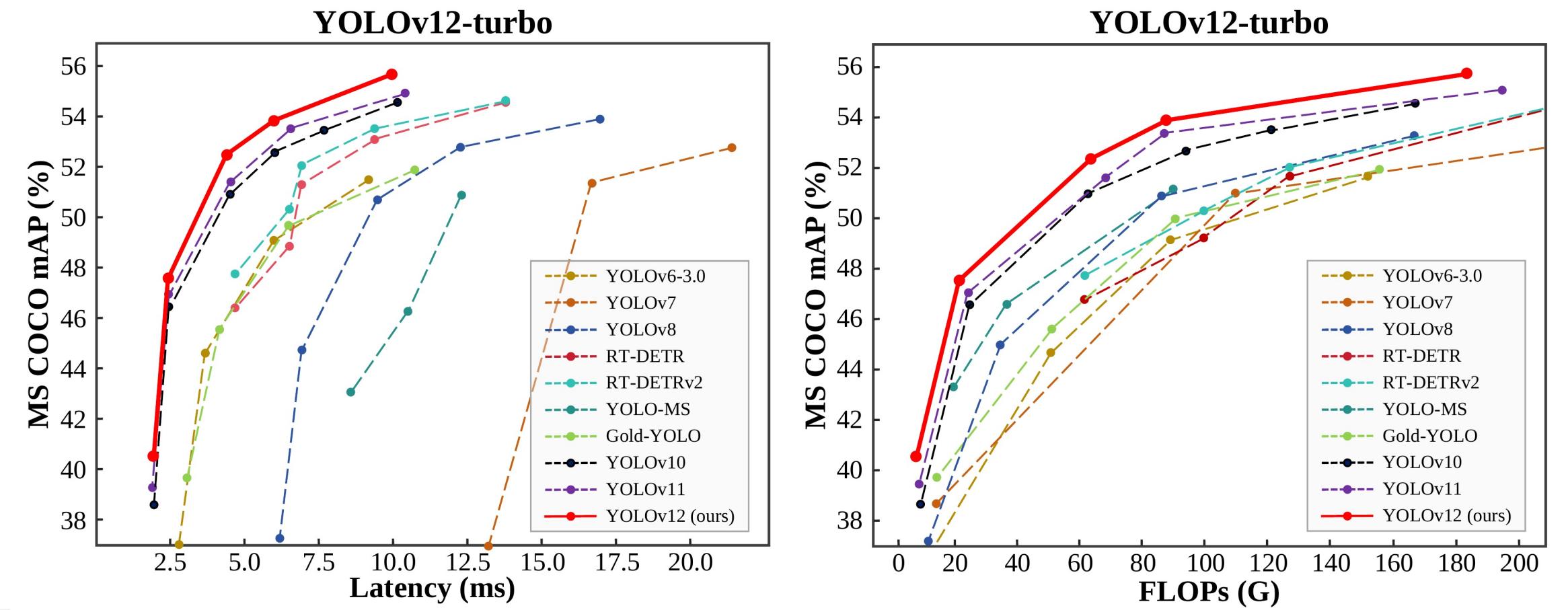
Qixiang Ye
University of Chinese Academy of Sciences
qxye@ucas.ac.cn

David Doermann
University at Buffalo
doermann@buffalo.edu

کاربردها و موارد استفاده YOLOV12



- شناسایی اشیاء (Object Detection)
- بخش بندی نمونه ها (instance segmentation)
- دسته بندی تصاویر (Image Classification)
- تخمین وضعیت بدن (Pose Estimation)
- تشخیص اشیا با باکس های مرزی مورب (OBB)



sunsmarterje/yolov12

YOLOv12: Attention-Centric Real-Time Object Detectors

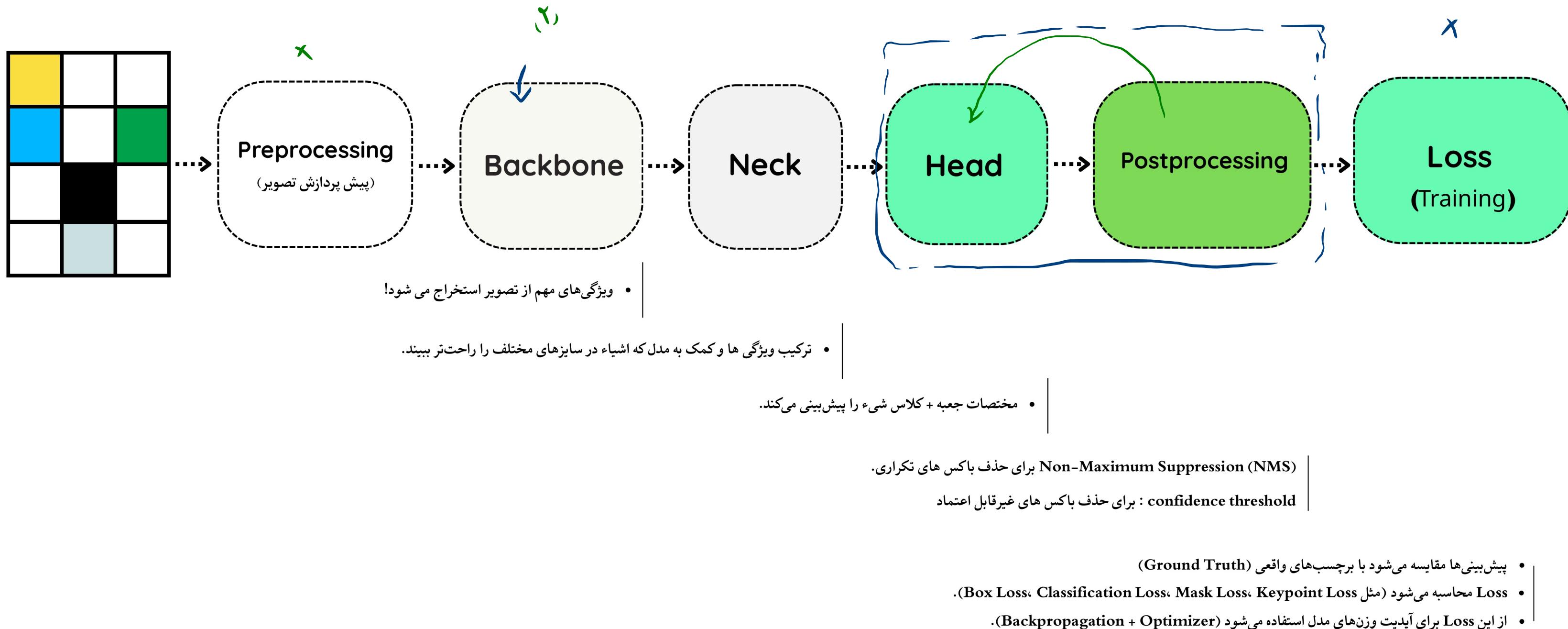
Contributors: 8 Issues: 69 Stars: 2k Forks: 201

sunsmarterje/yolov12: YOLOv12: Attention-Centric Real-Time Object Detectors

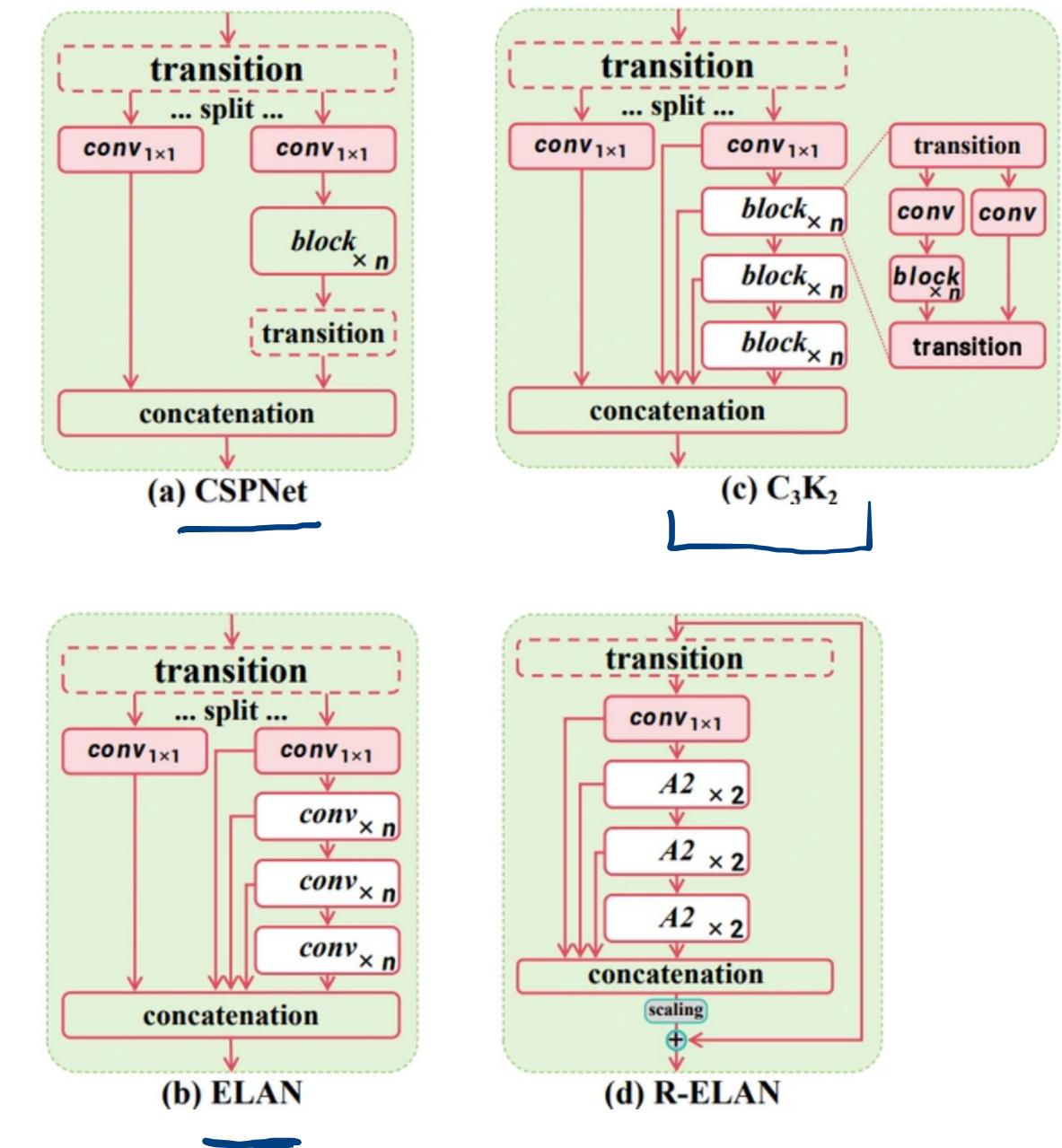
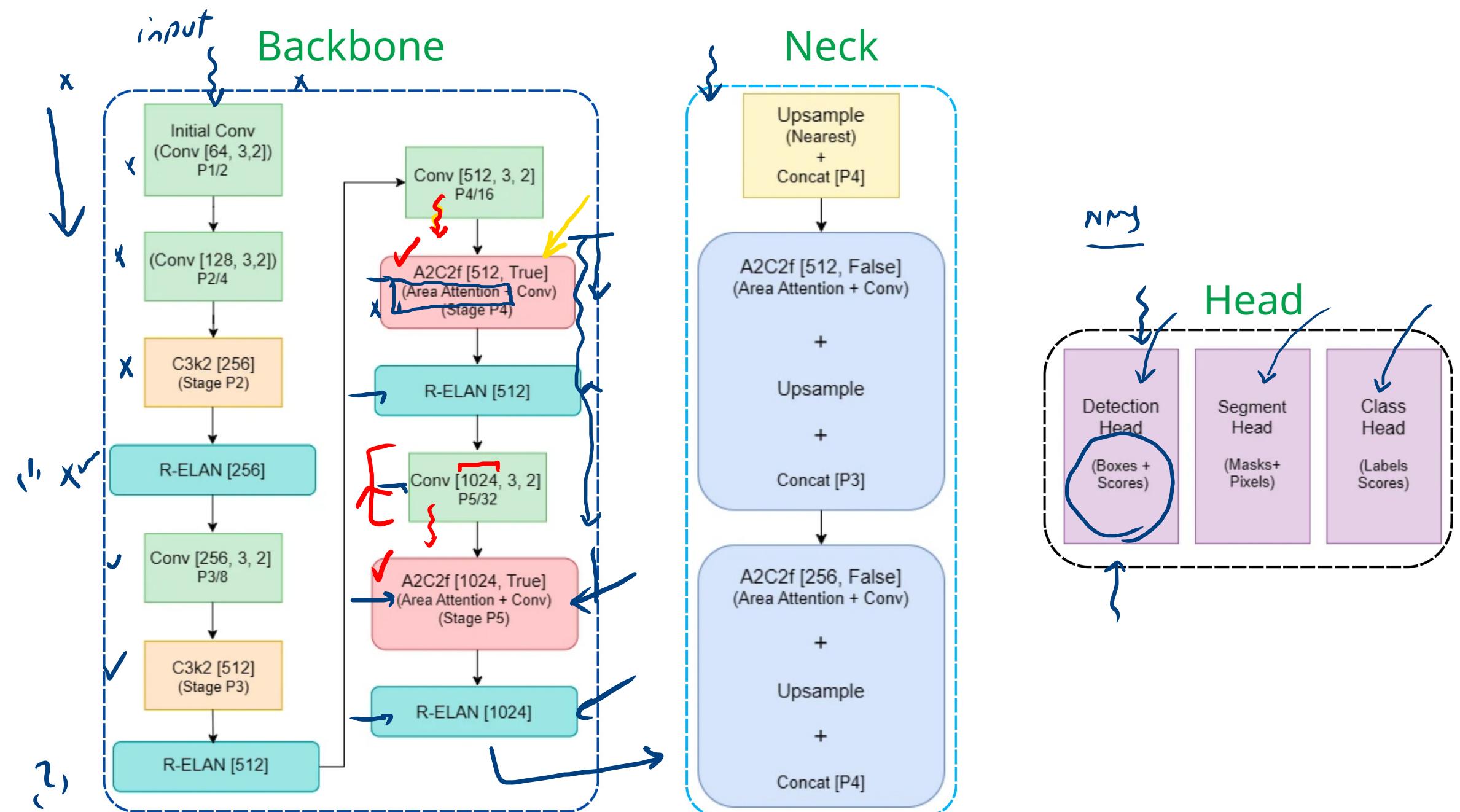
YOLOv12: Attention-Centric Real-Time Object Detectors - sunsmarterje/yolov12

[GitHub](#)

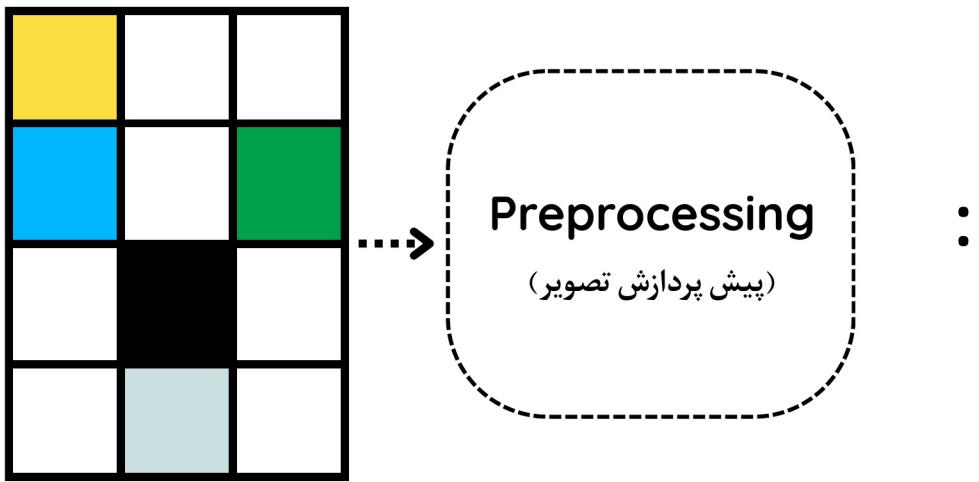
معماری کلی YOLOv12



معماری جزئی YOLOv12



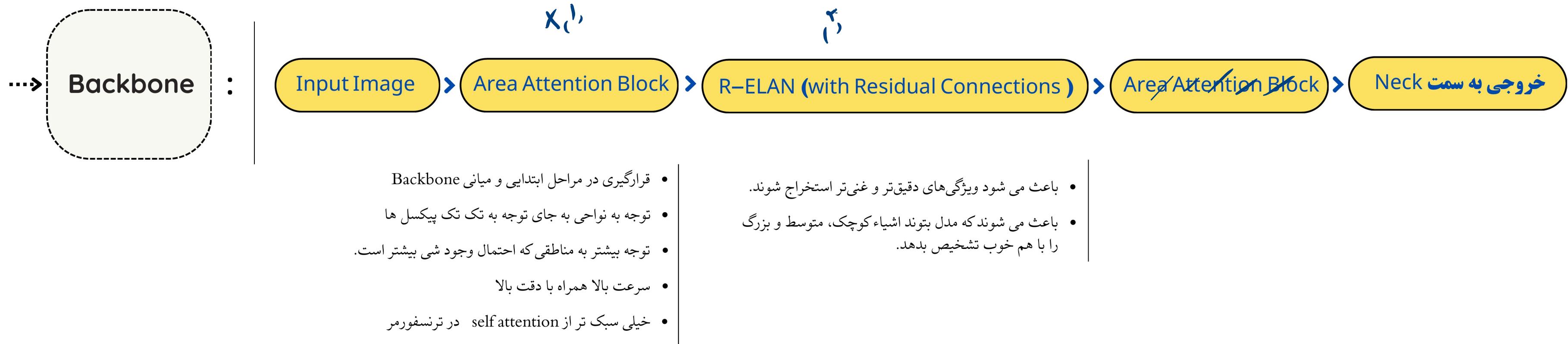
مرحله ۱: پیش پردازش



- **Resize : 640x640, ...**
- **Normalize : (Image/255)**
- **Data Augmentation Pipeline**
Flipping, Cropping, Scaling, Mosaic Augmentation , ...

مرحله ۲ : Backbone

در این مرحله ویژگی‌های ساده مثل لبه‌ها، تا ویژگی‌های پیچیده مثل شکل کلی اشیاء



مرحله ۱ از بخش Area Attention Block : قسمت Backbone



توجه بیشتر به مناطقی که احتمال وجود شی بیشتر است.

Area Attention Block

تمرکز روی نواحی به جای توجه تمرکز روی پیکسل ها
(برخلاف self attention)



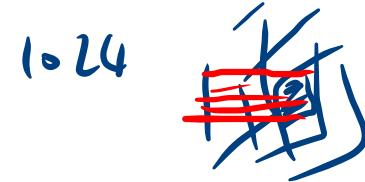
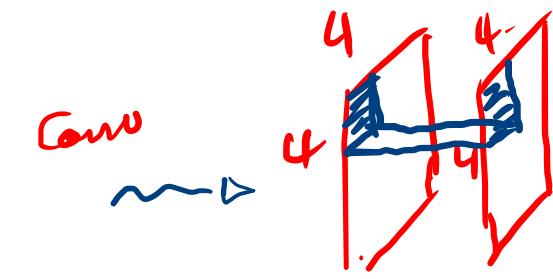
سرعت بیشتر نسبت به
کاهش هزینه محاسباتی

Area Attention

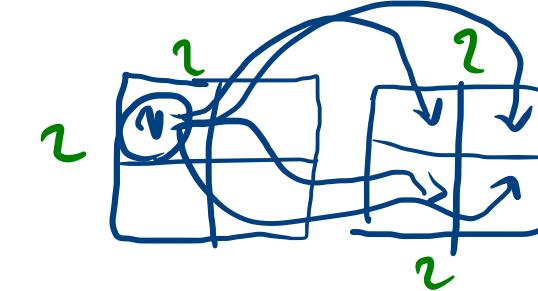
Yang Li¹ Lukasz Kaiser¹ Samy Bengio¹ Si Si¹

¹Google Research, Mountain View, CA, USA. Correspondence
to: Yang Li <liyang@google.com>

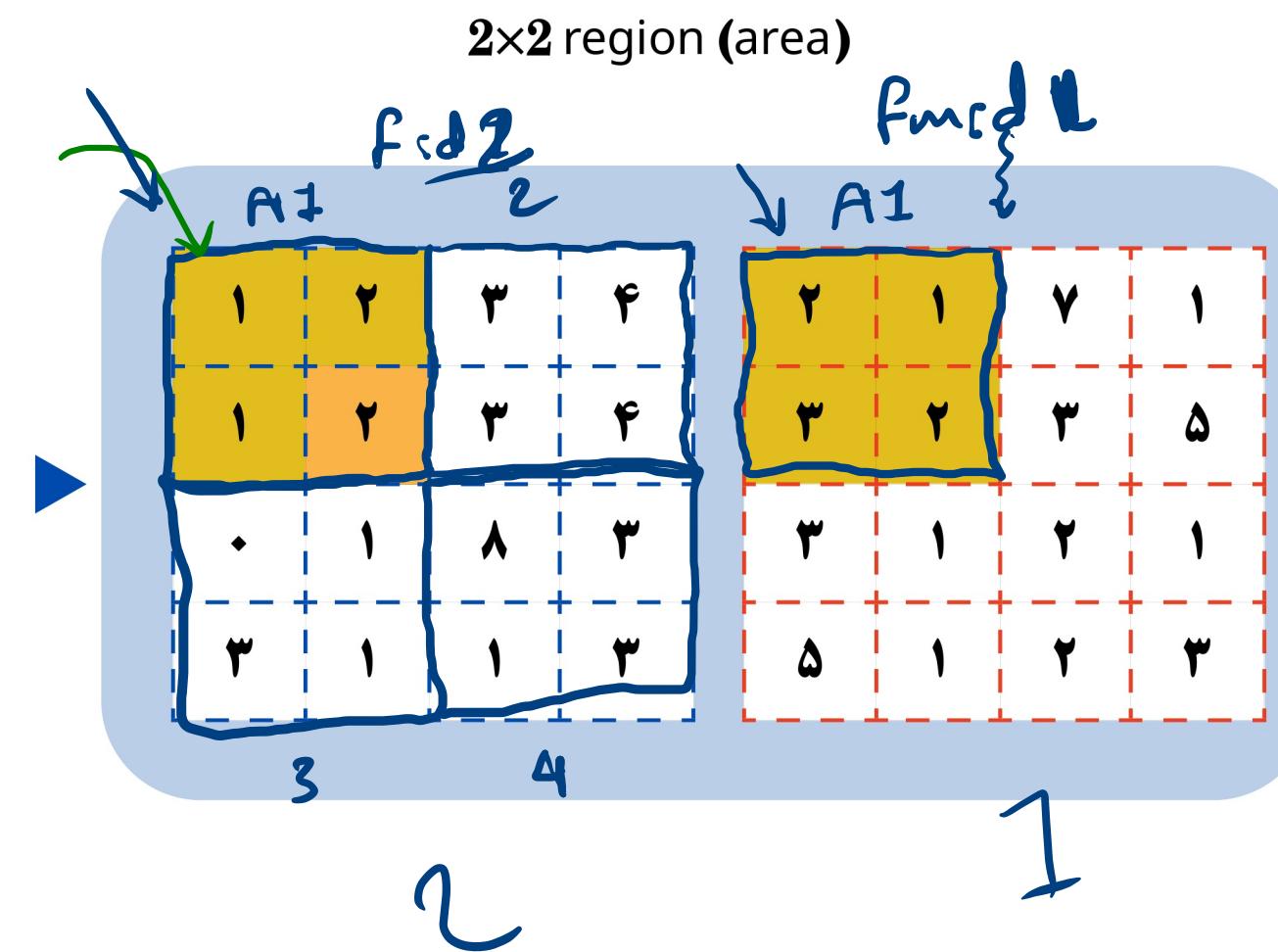
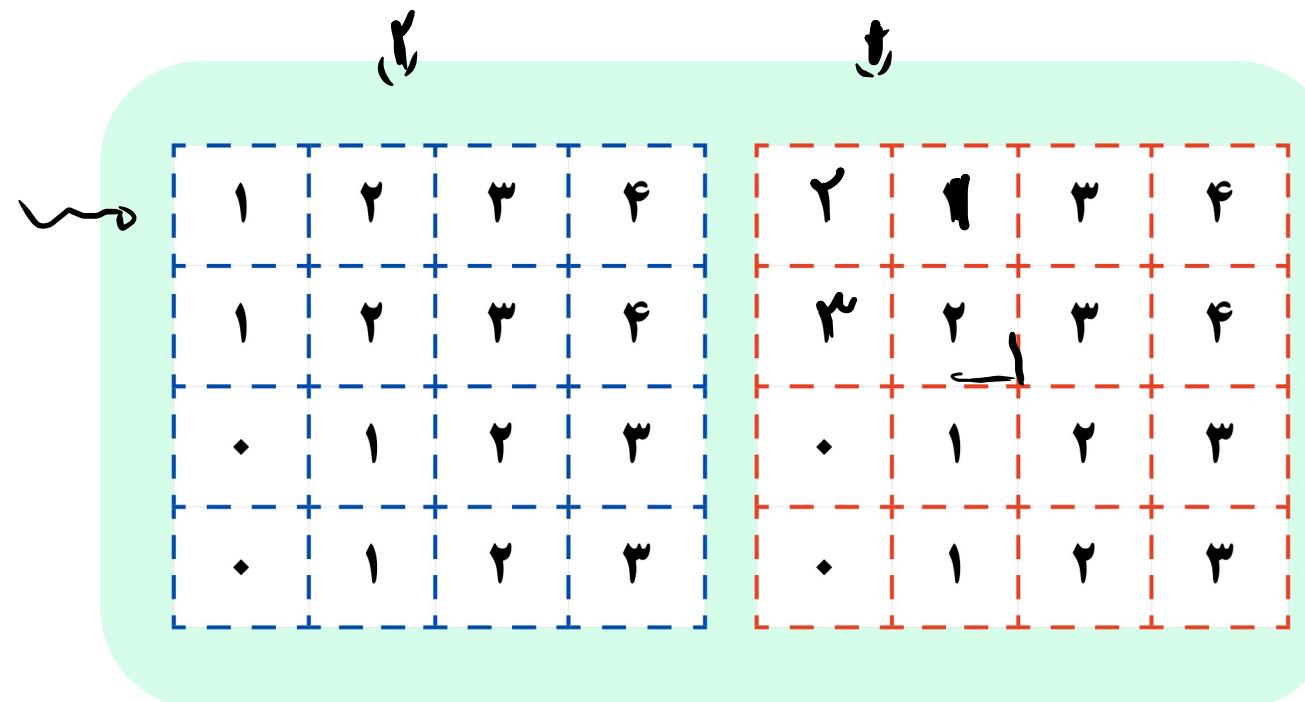
Area Attention Block - 1

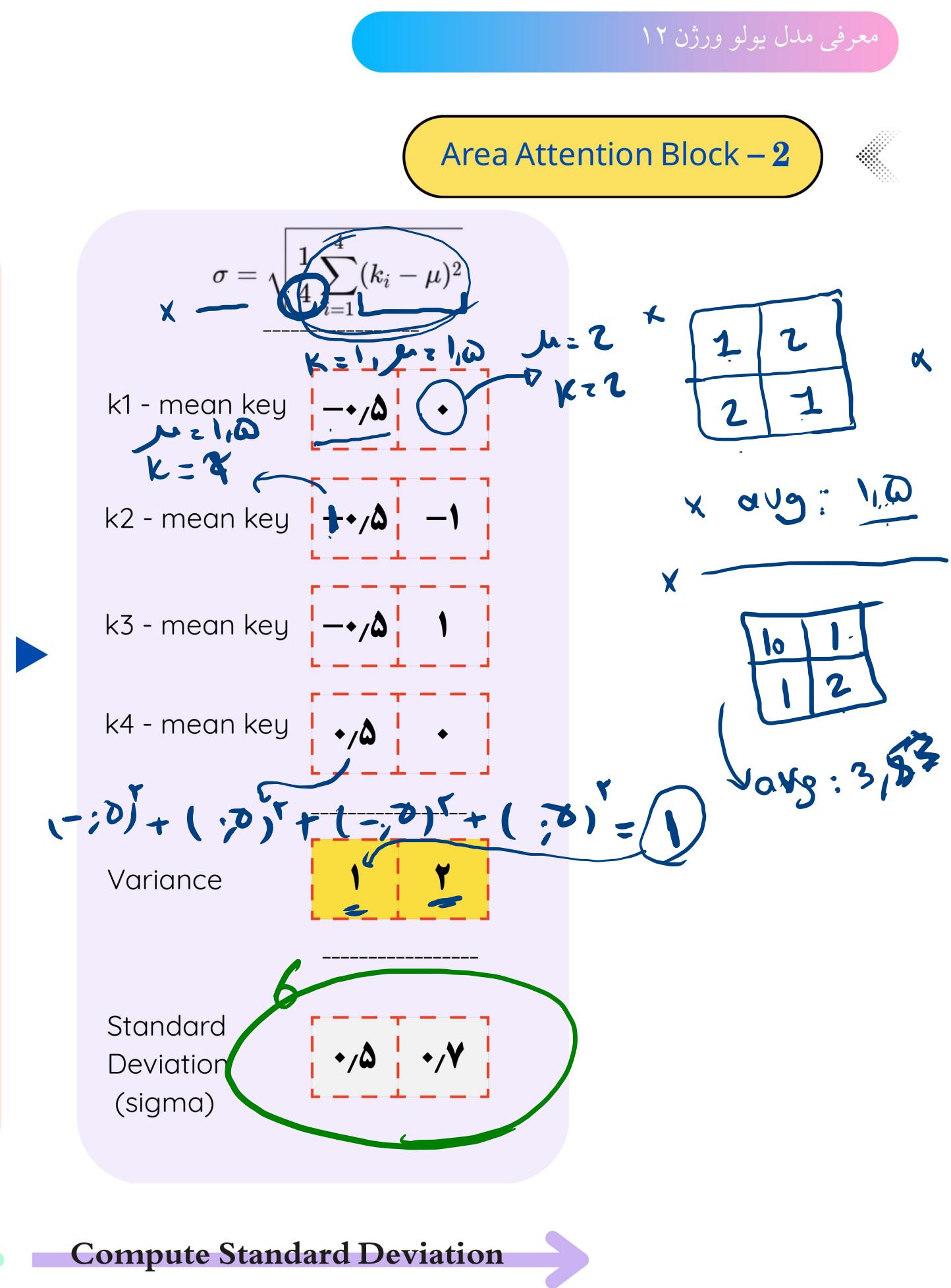
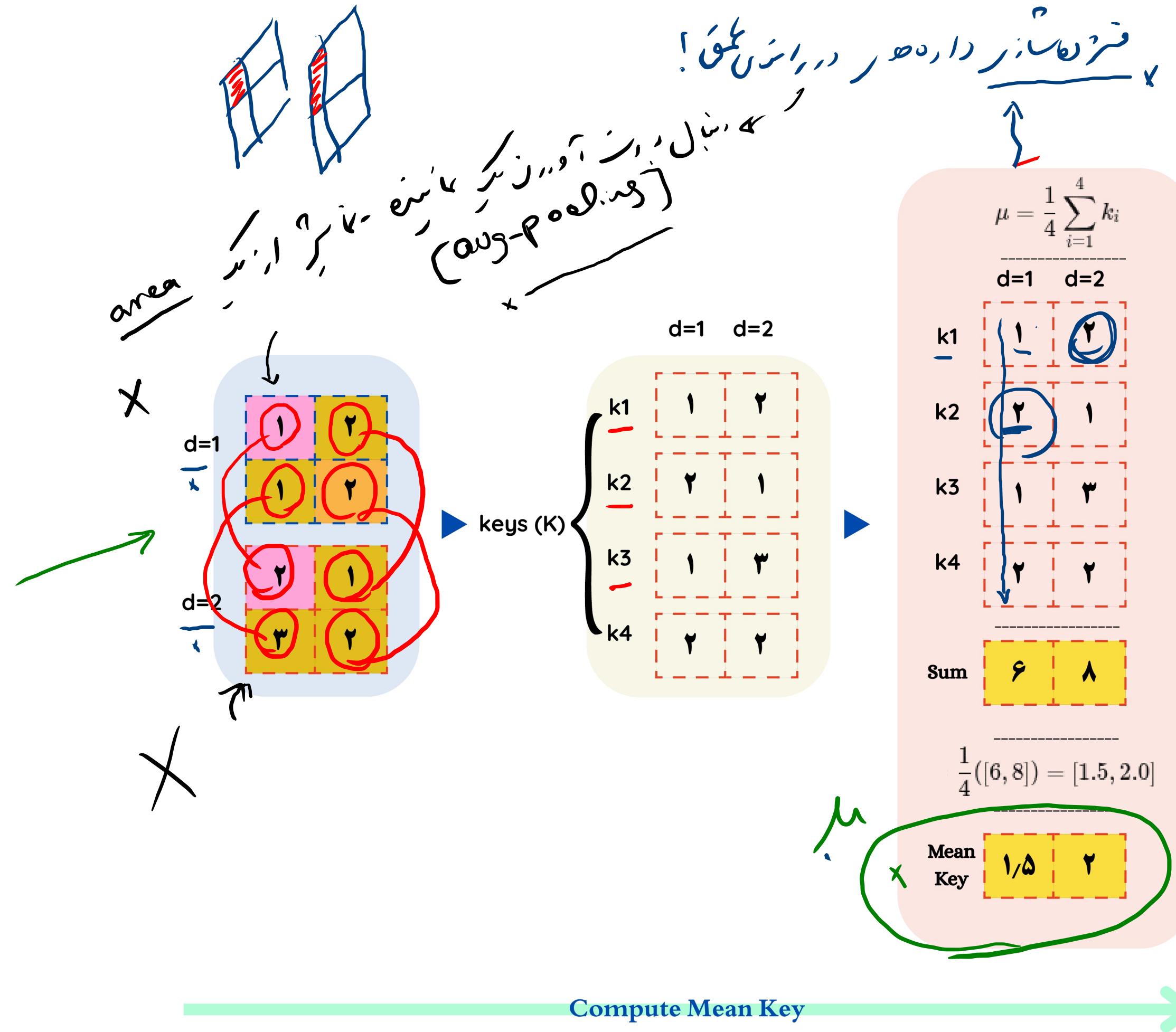


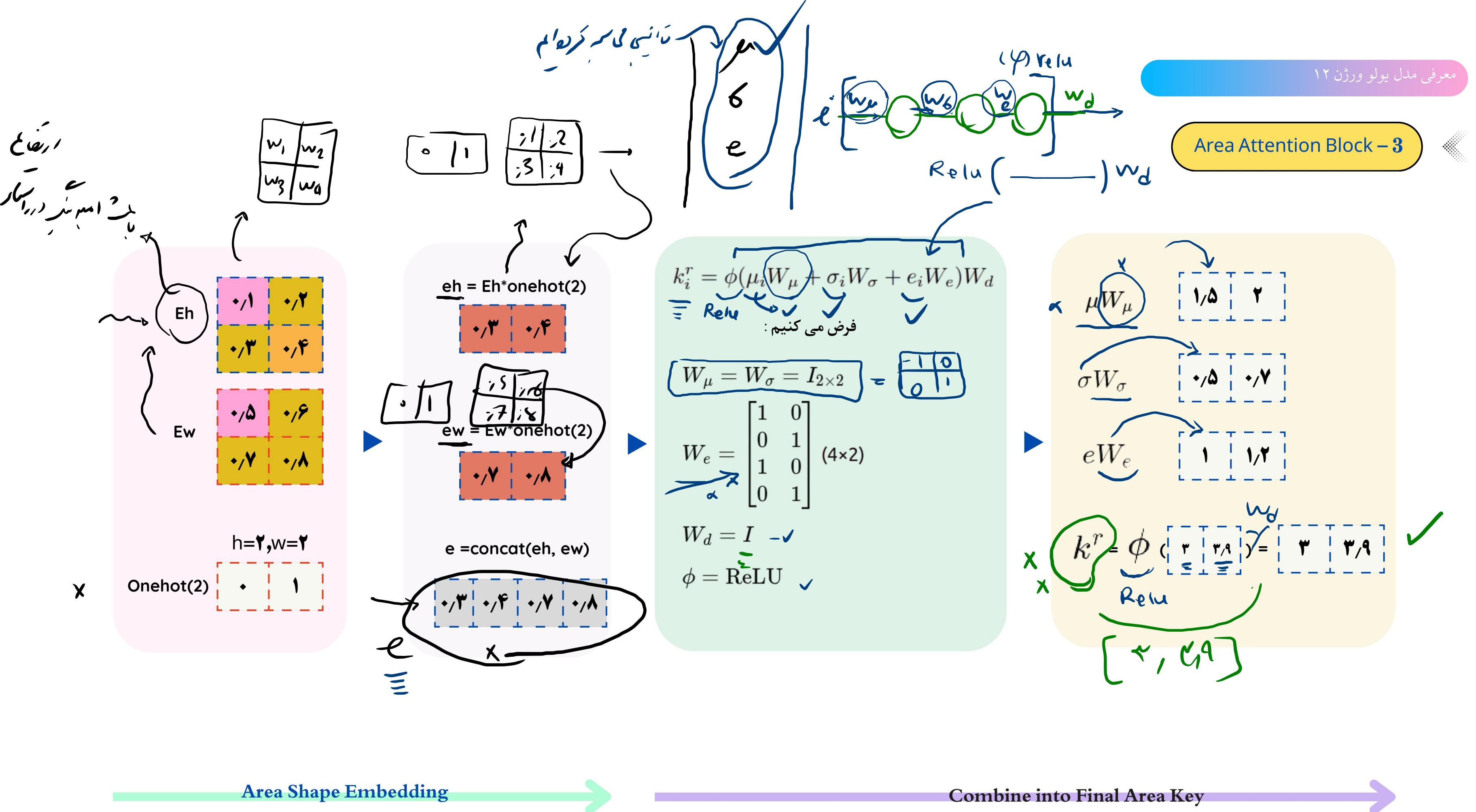
تصویر (فیچر مپ) ورودی
(shape $4 \times 4 \times 2$)



$$\begin{array}{l} h_{\text{max}} = 2 \\ w_{\text{max}} = 2 \end{array} \rightarrow \begin{array}{l} 1 \leq h \leq 2 \\ 1 \leq w \leq 2 \end{array}$$







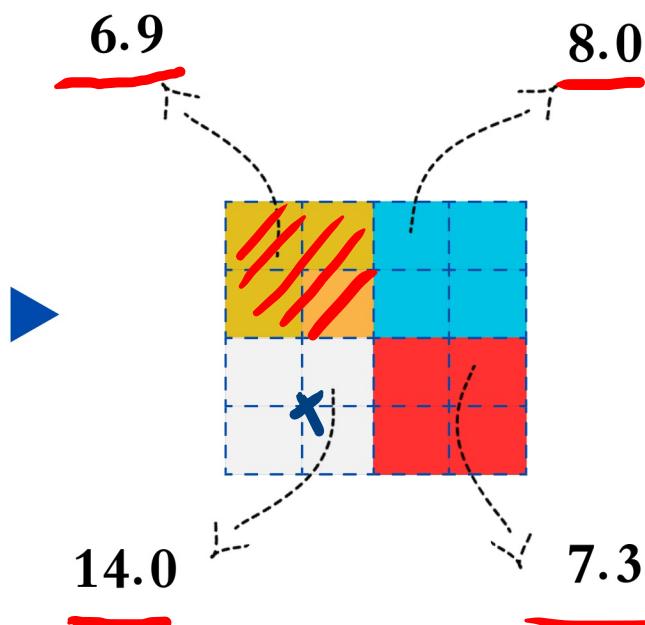
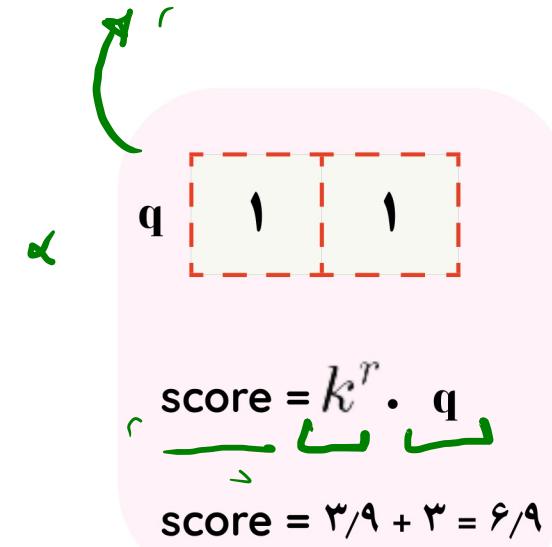
Area Attention Block - 4

$$\frac{e^i}{\sum e^i} = \frac{e^{6.9}}{e^{6.9} + e^{8.0} + e^{14.0} + e^{7.3}}$$

w
parameter.

trainable

وزن



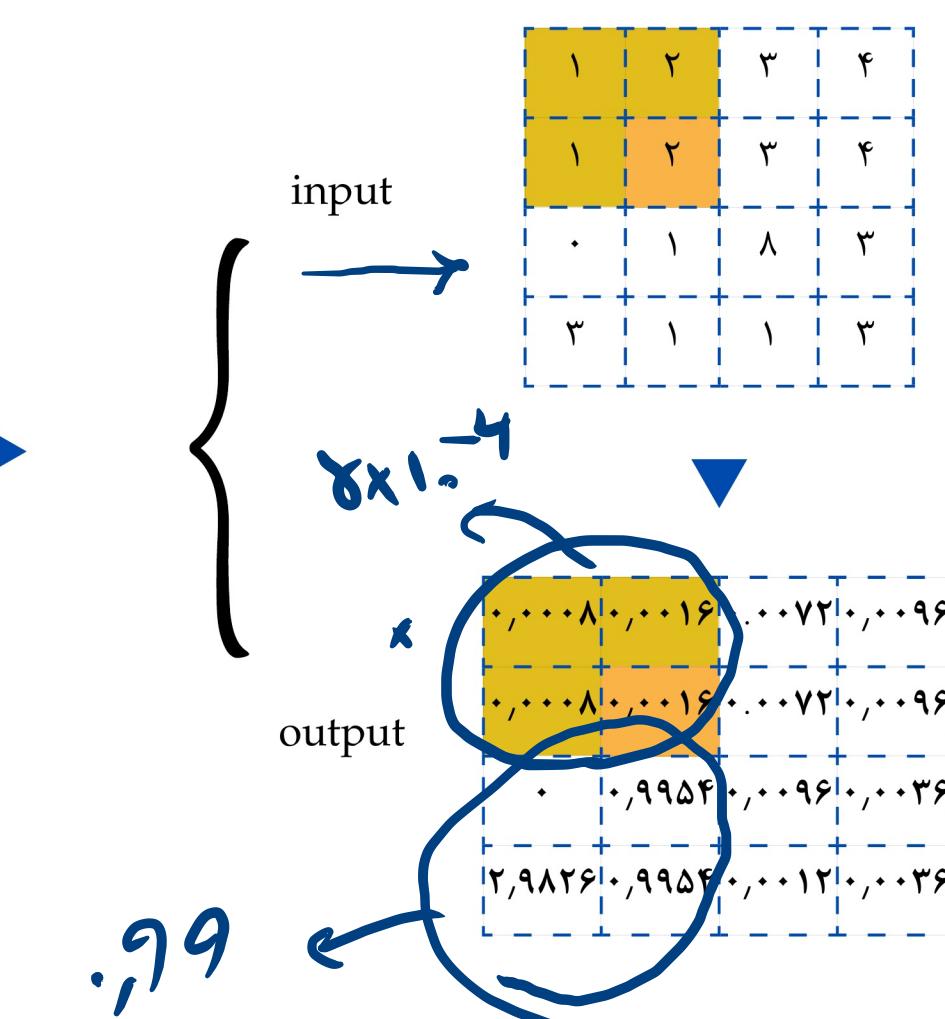
نمره توجه به یک ناحیه (area)
area

$$3,9 + 3 = 6,9$$

Scores = [6.9, 8.0, 14.0, 7.3]
attention weights = Softmax(Scores)

attention weights = [0.00082, 0.00247, 0.99549, 0.00123]

••••••••	*	top left	x
••••••••	*	top right	x
••••••••	*	bottom left	x
••••••••	*	bottom right	x

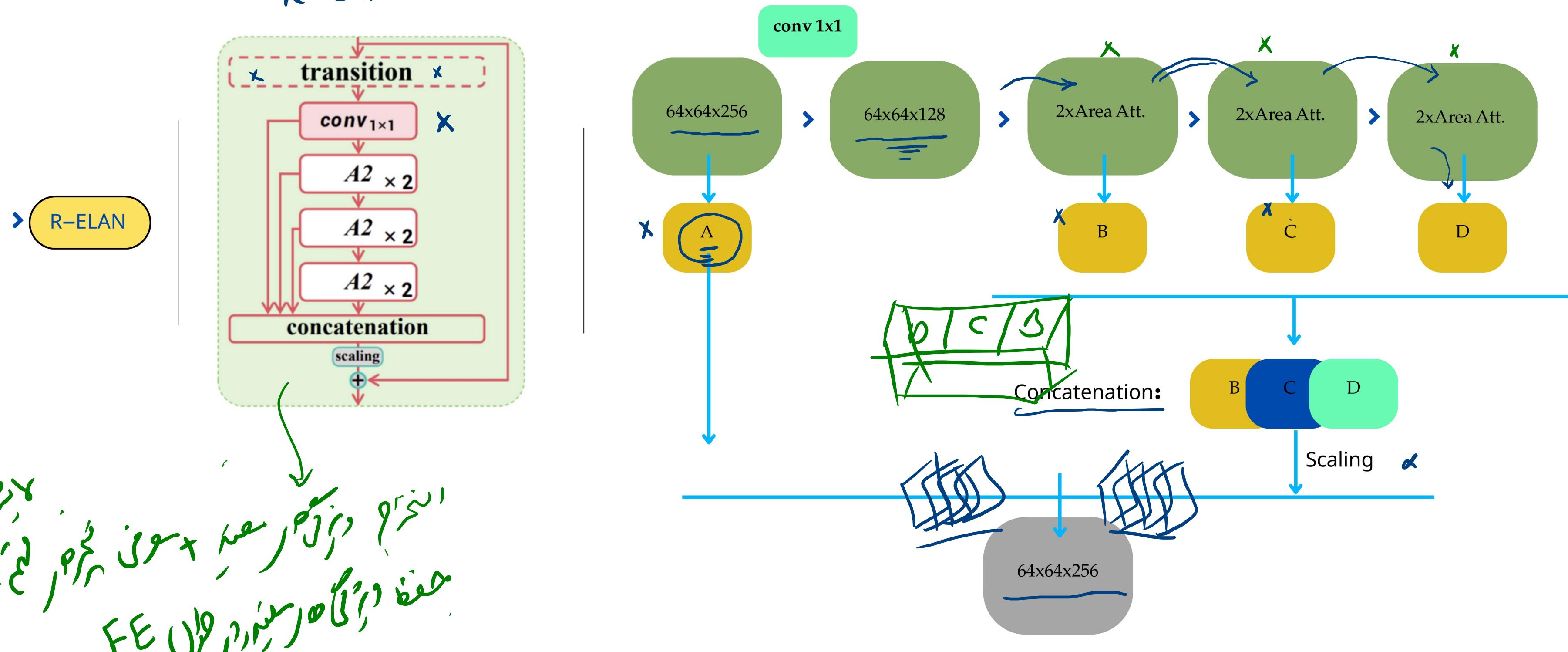


Query & Dot Product

Final Attention Output

Residual-Enhanced Efficient Layer Aggregation Network

ignores low-level features
Feature
R-ELAN



Image

2	3
8	9



1
2

1

= ?



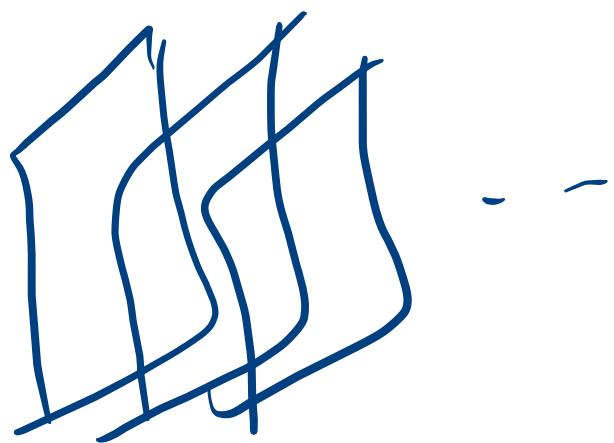
4	6
16	18

α



$\text{Conv}(128, (1,1))$

128





خروجی یولو ۱۲ با AREA ATTENTION



The End