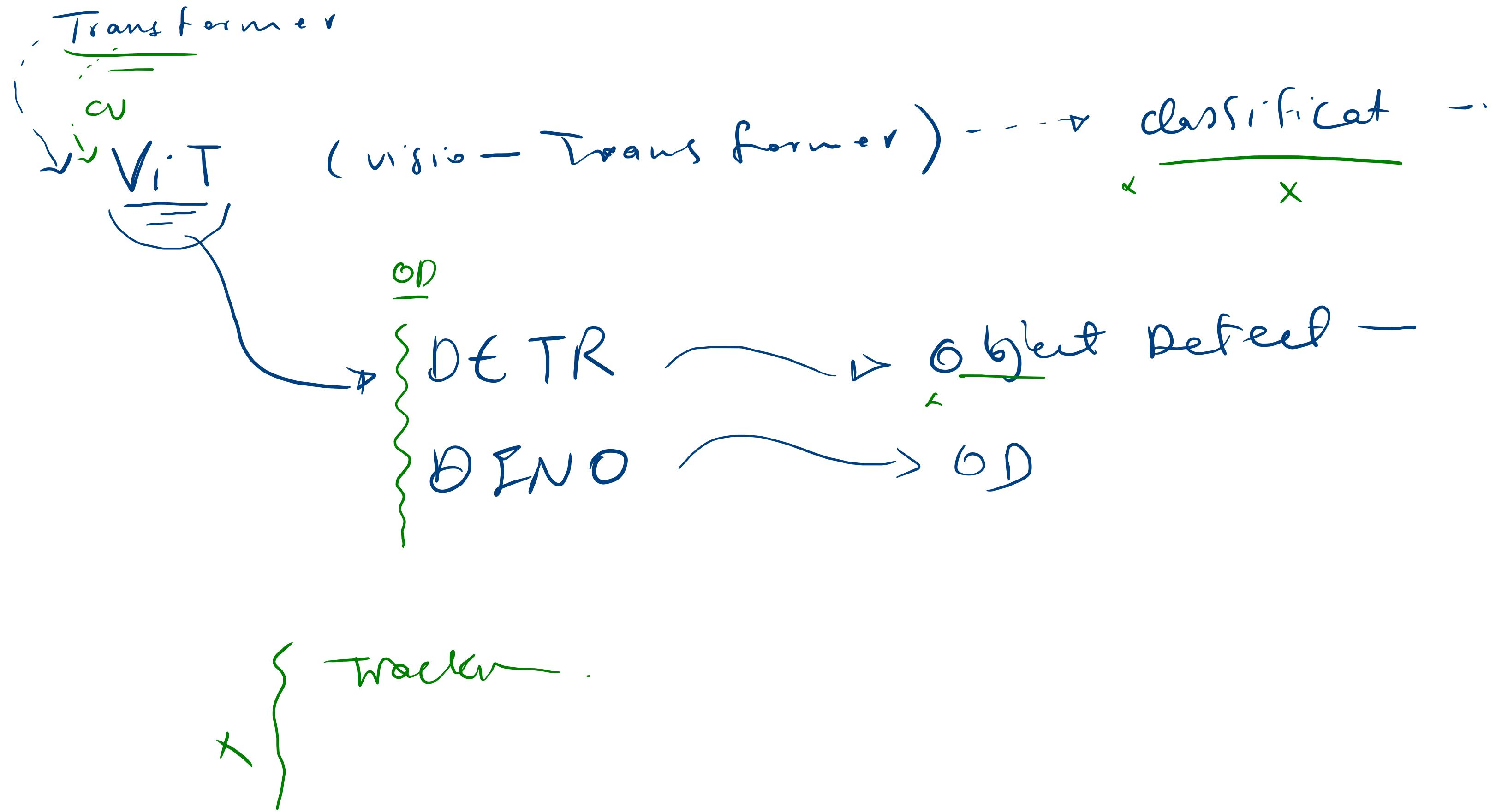


Train —

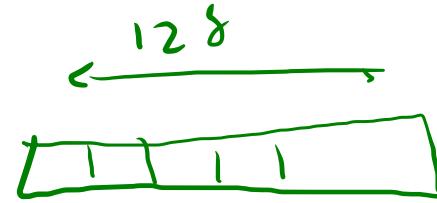
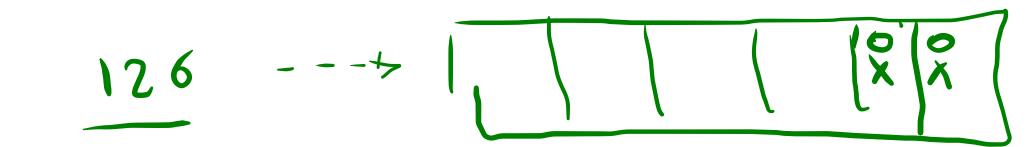
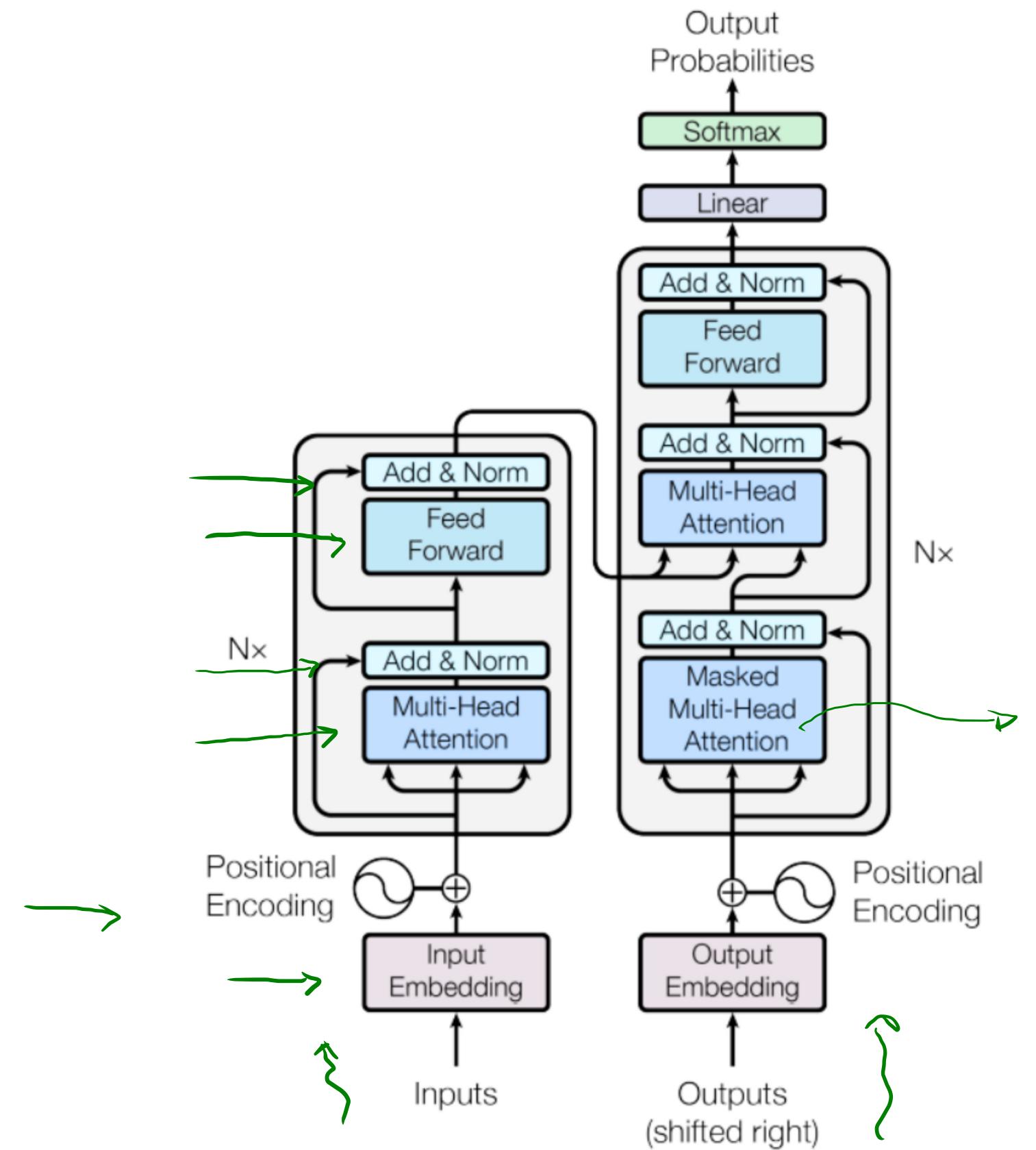
Detr

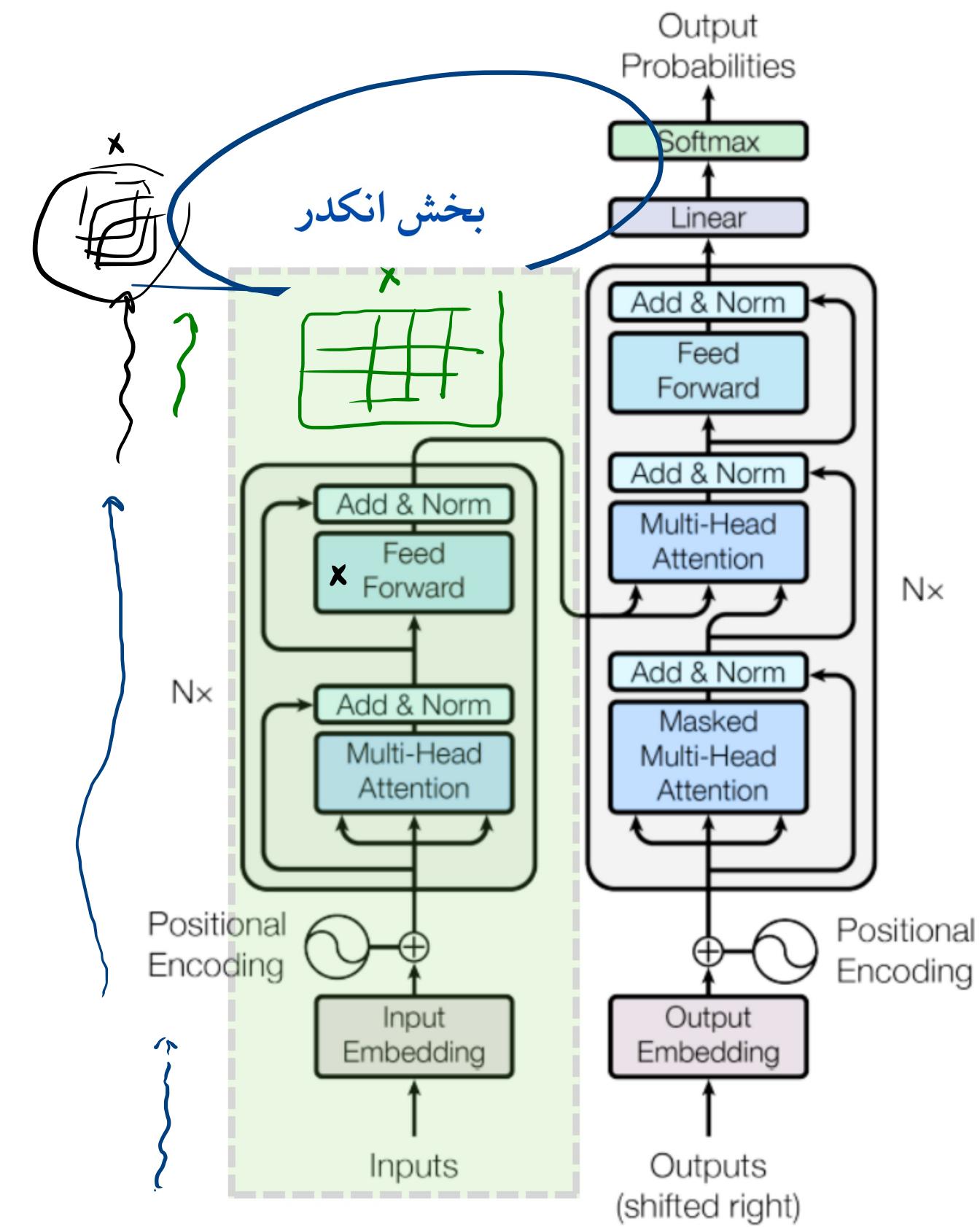
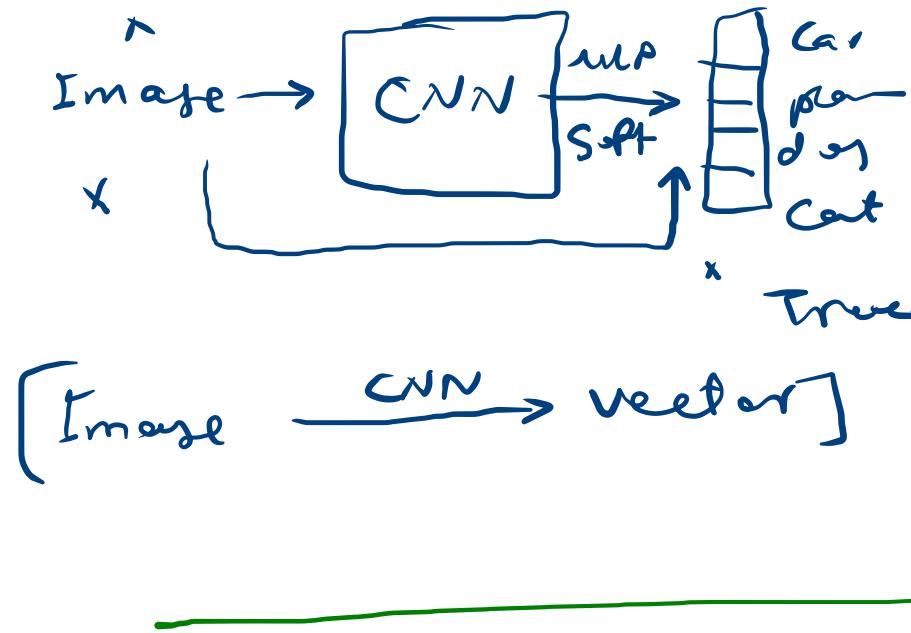
Origin



شبکه ویژن ترنسفورمر

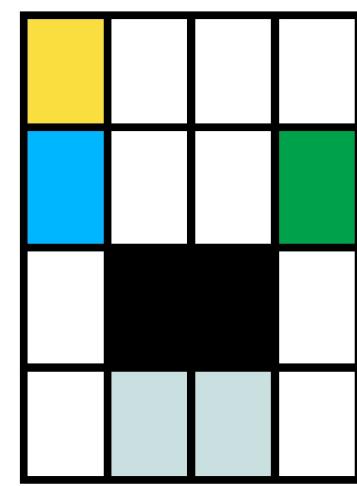
Vision Transformer
(ViT)





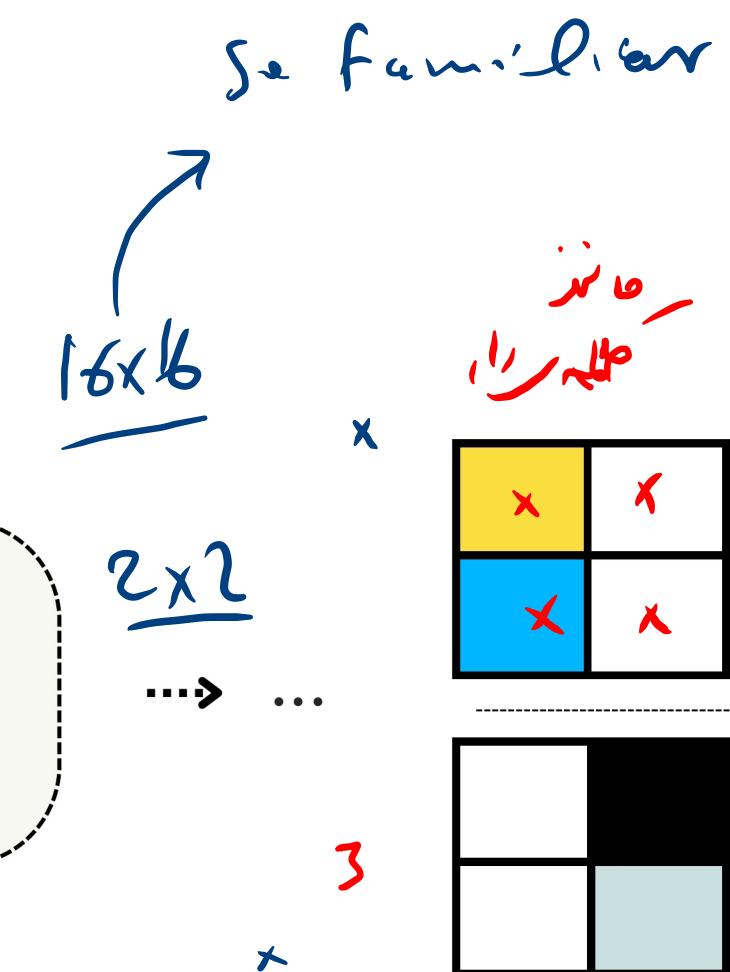
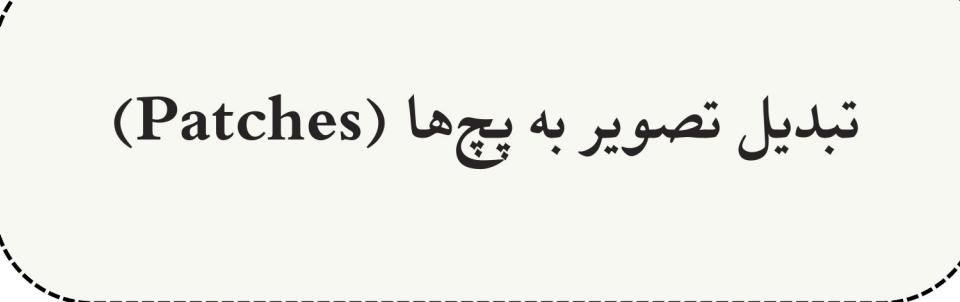
شبکه ViT از بخش انکدر ترانسفورمر استفاده می کند.

شبکه ViT از بخش انکدر ترانسفورمر استفاده می کند.



4

تبديل تصوير به پچها (Patches)



$I \rightarrow [1, 3, 8, 9]$

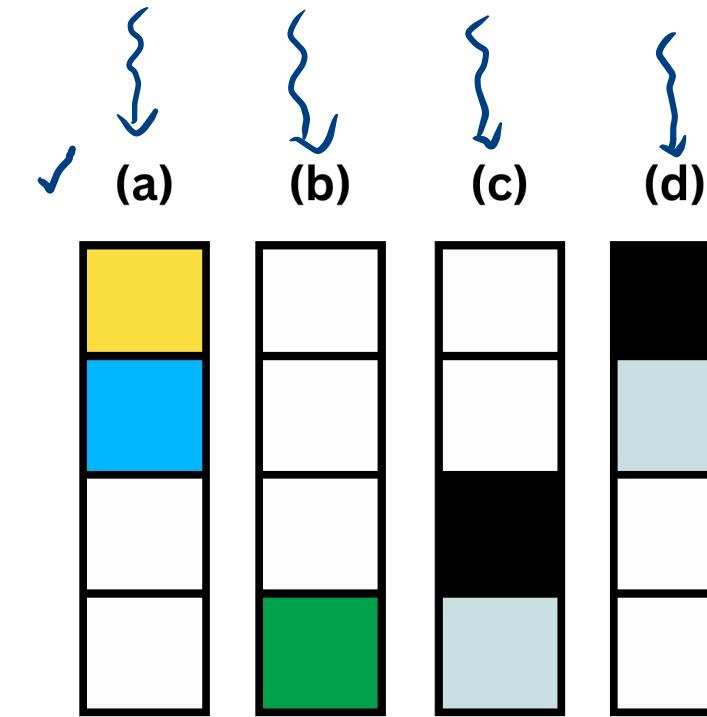
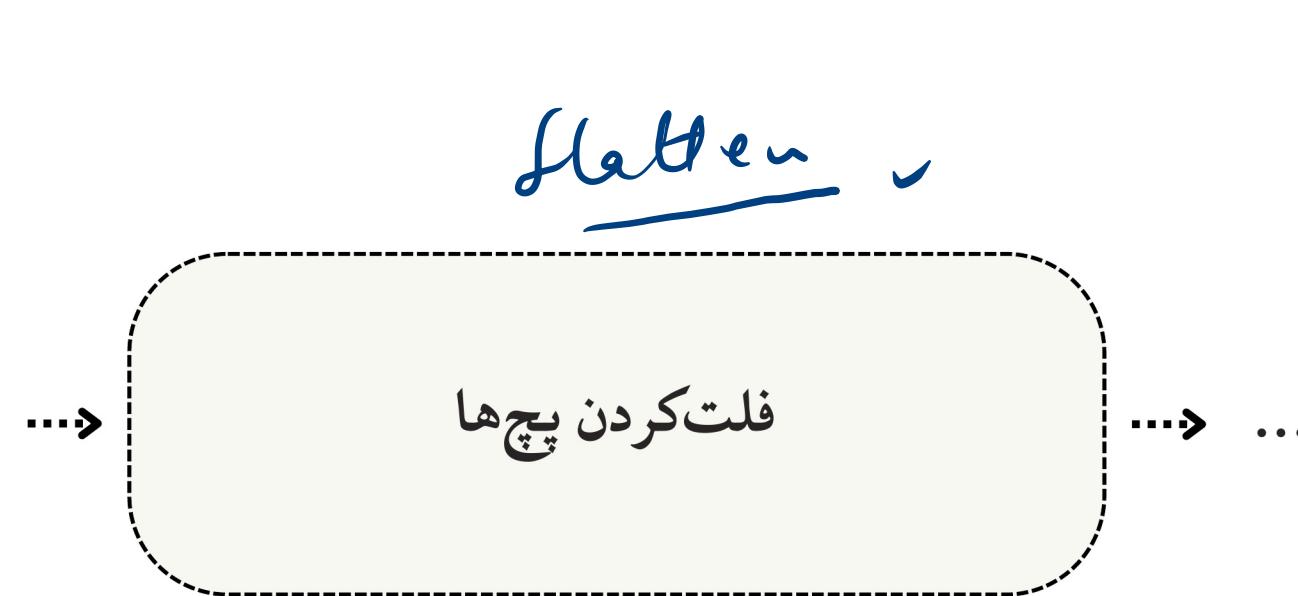
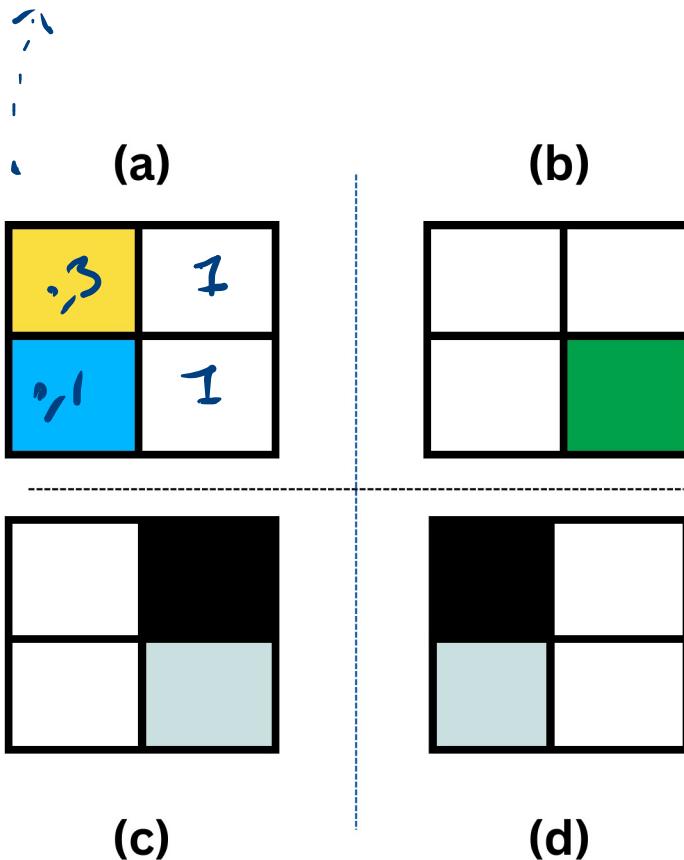
شبکه ویژن ترانسفورمر

patch $[0][3, 1, 1, 1]$

شبکه ویژن ترانسفورمر



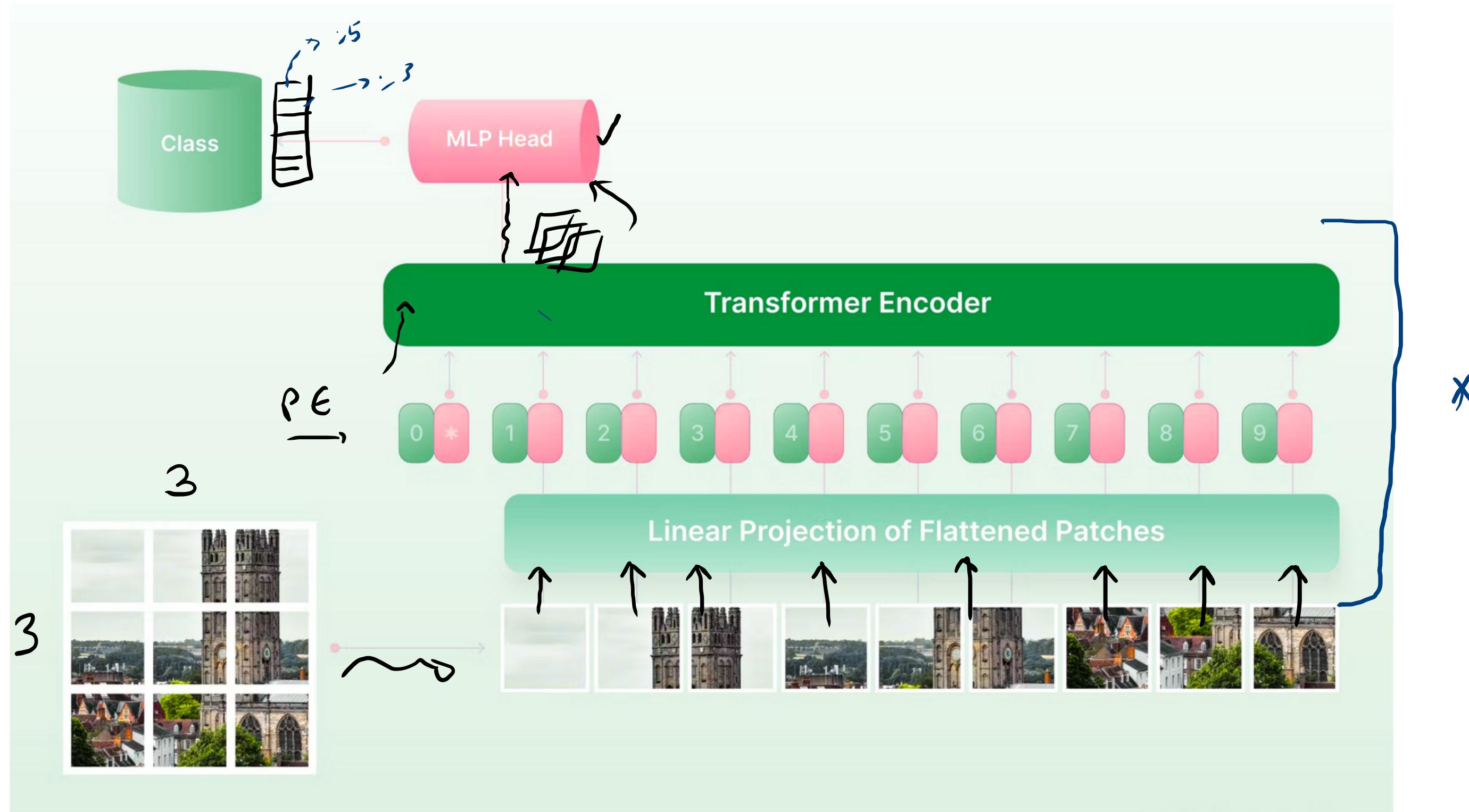
بخش دوم: فلت کردن پچ ها



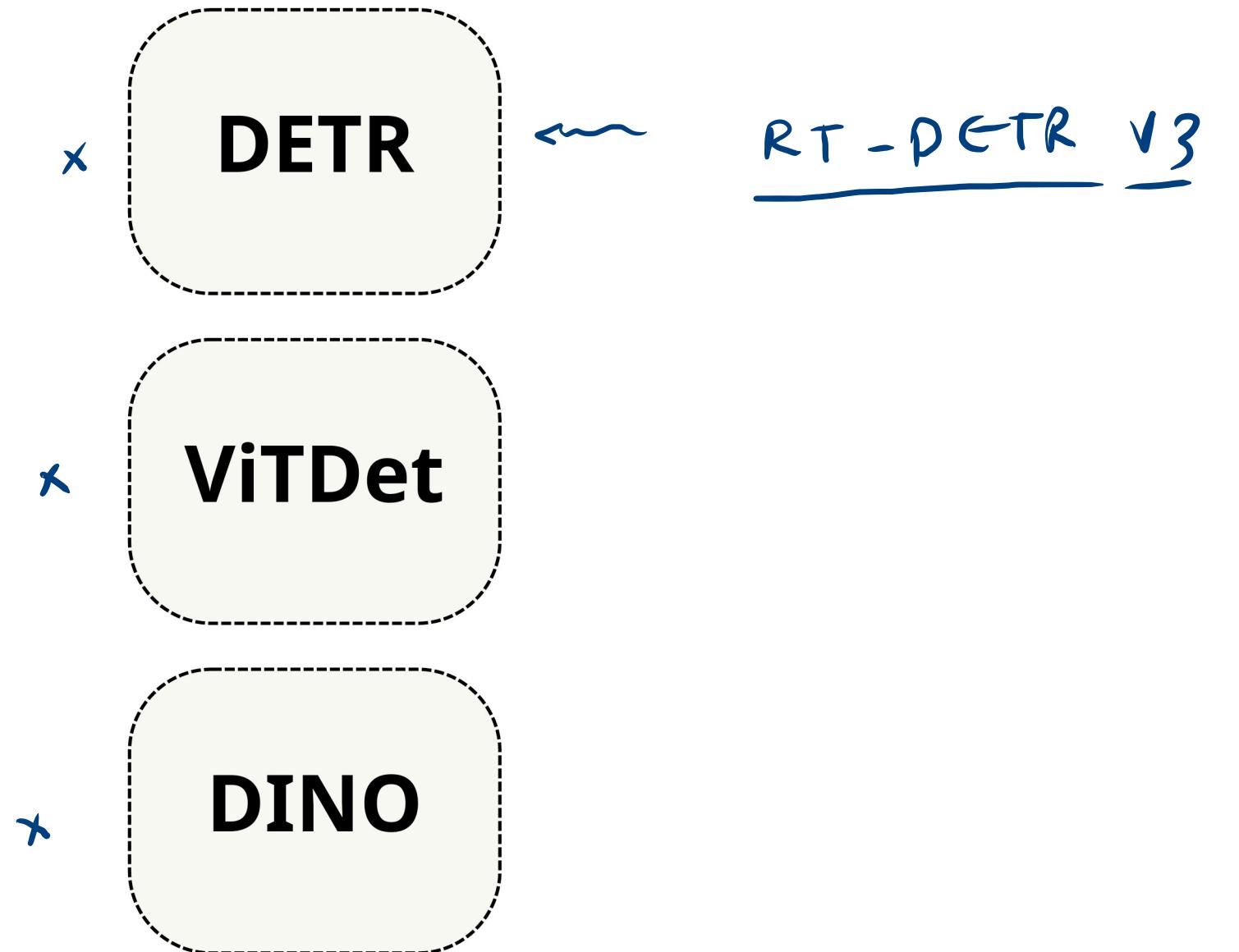
Resnet
— α

شبکه ویژن ترانسفورمر (برای کلاسیفای تصویر)

شبکه ViT از بخش انکدر ترانسفورمر استفاده می‌کند.



شبکه ویژن ترانسفورمر برای تشخیص اشیا
OD
شبکه ViT از بخش ازنکدکر ترانسفورمر استفاده می‌کند.



CLIP

Train

Image

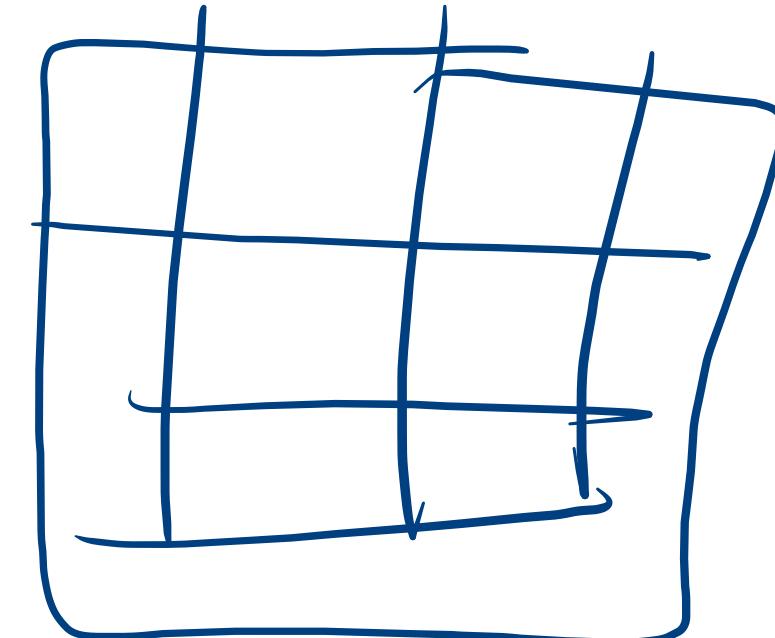
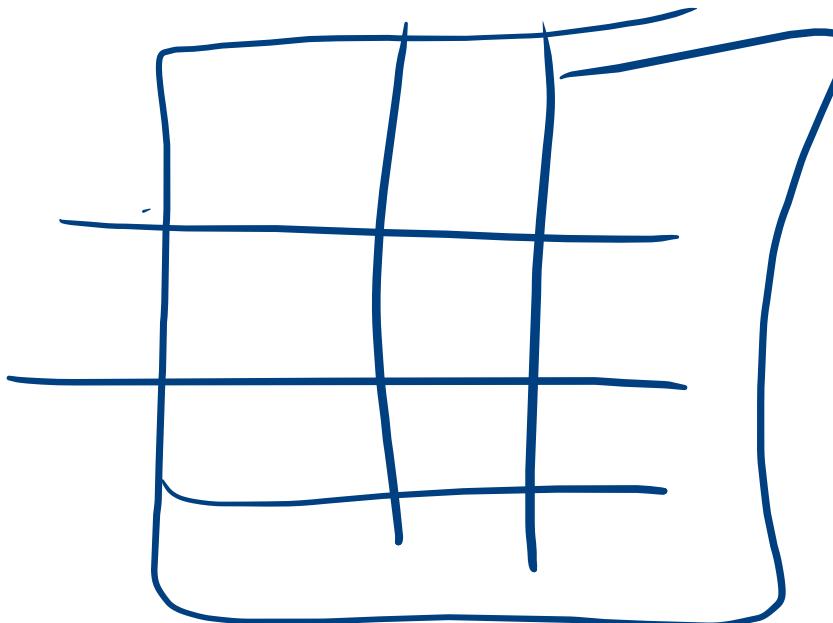
Image encoder

Tent

Tent encoder

Image encoder

Tent encoder



~10

loss

Image vector

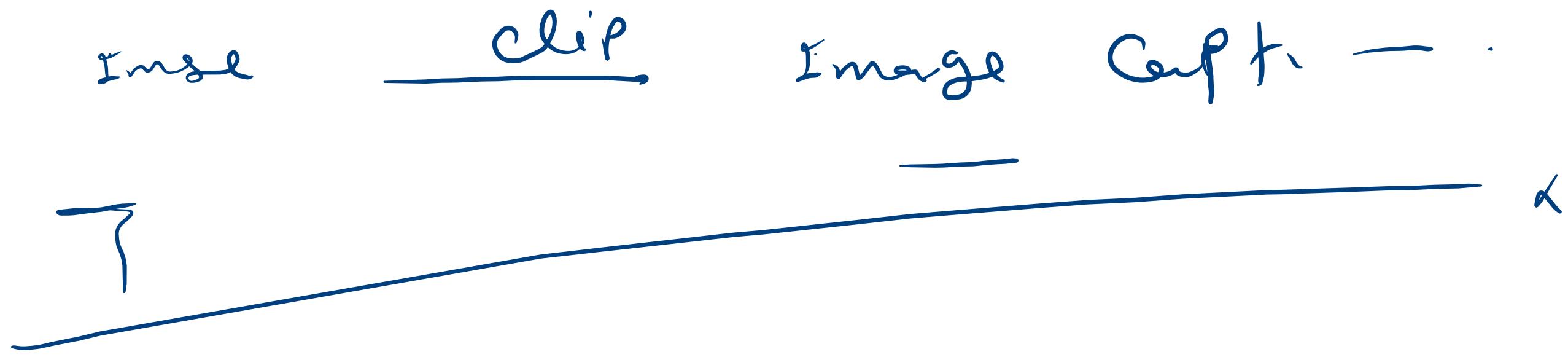
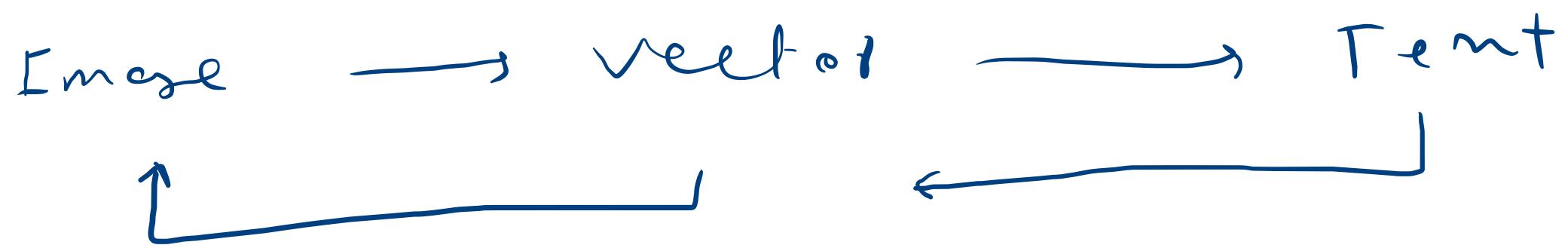
1

Tent vector

2

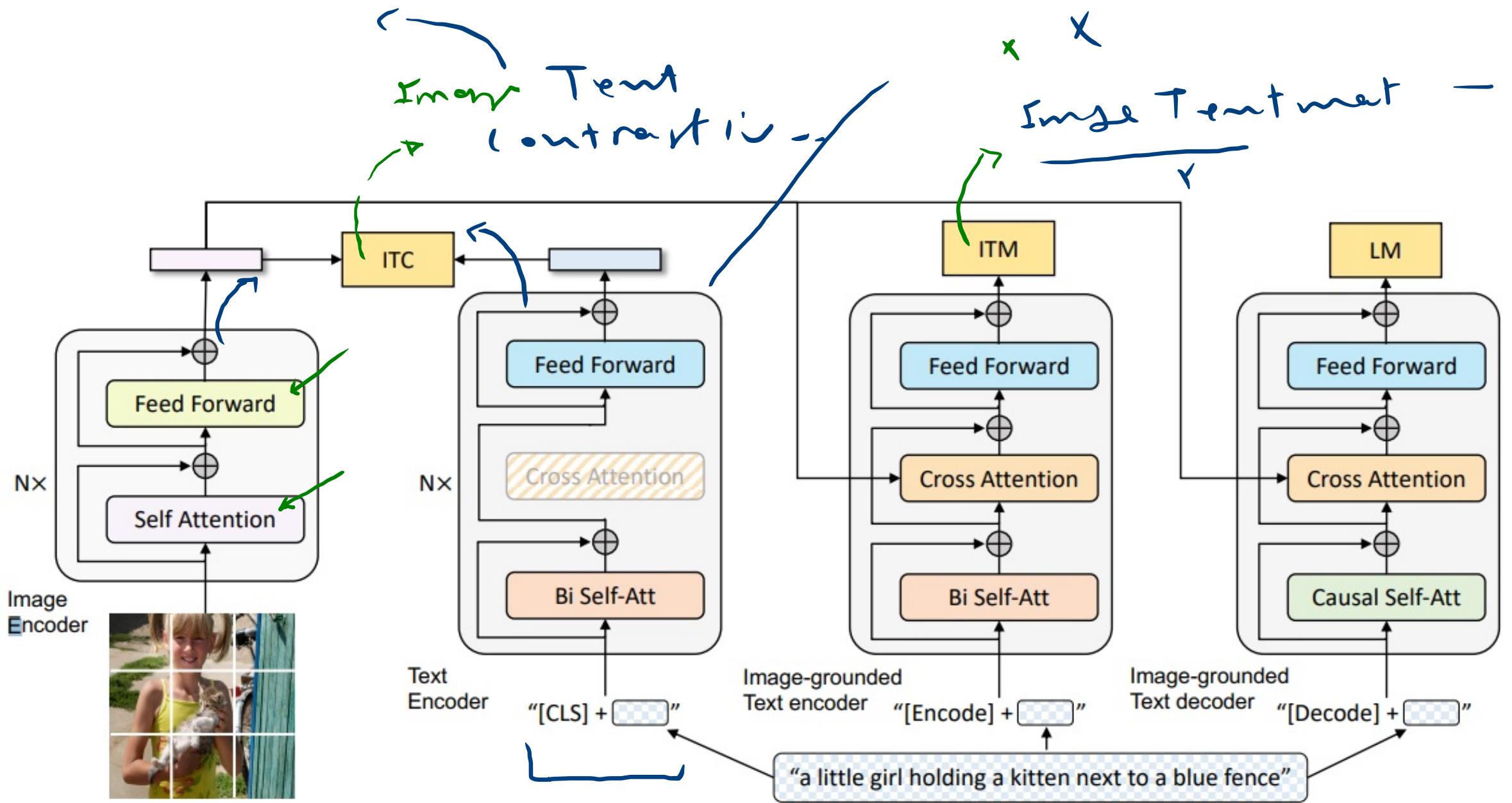
Similarity ✓

log



Image

enlarge train



BLEP > clip