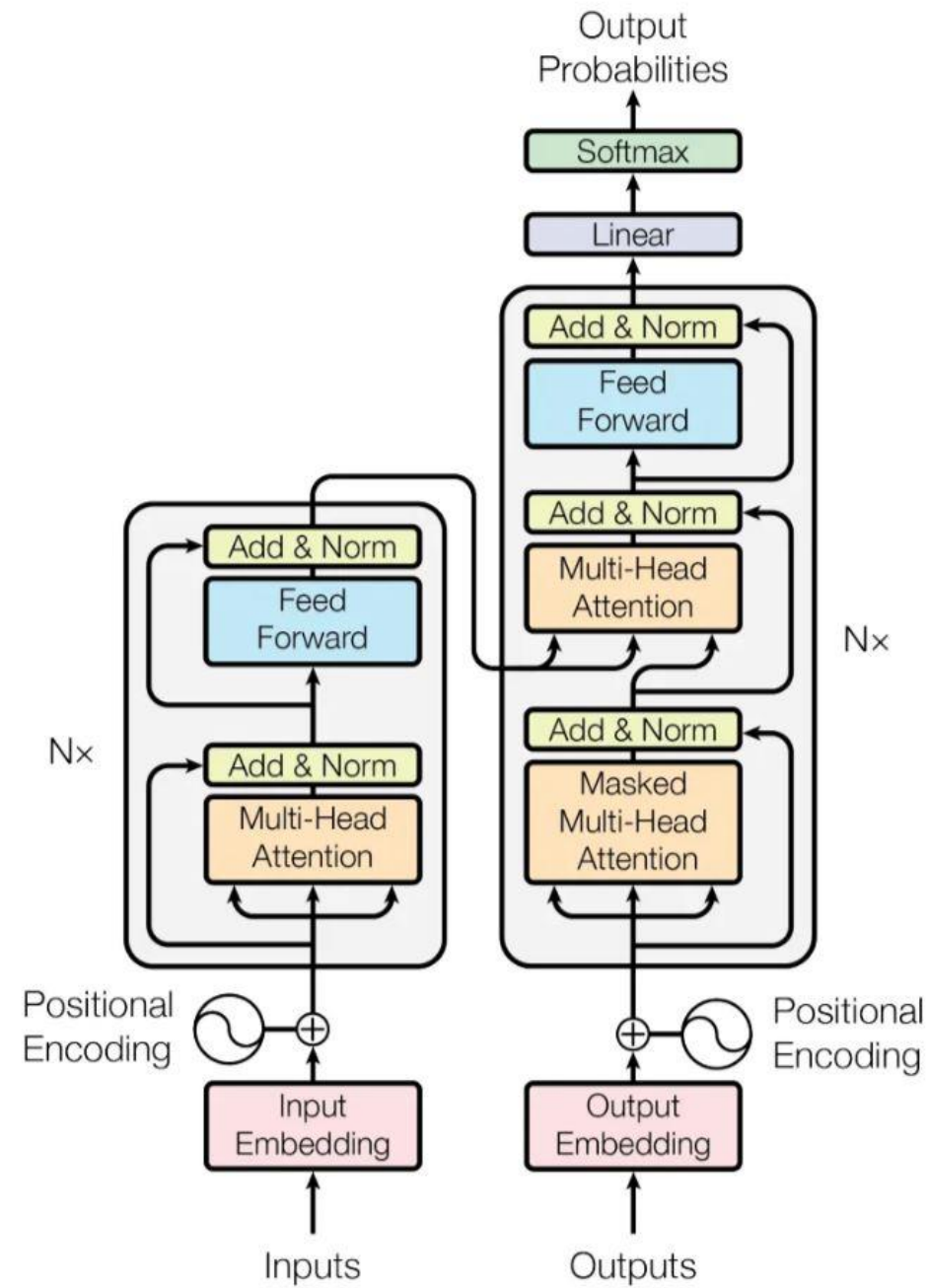


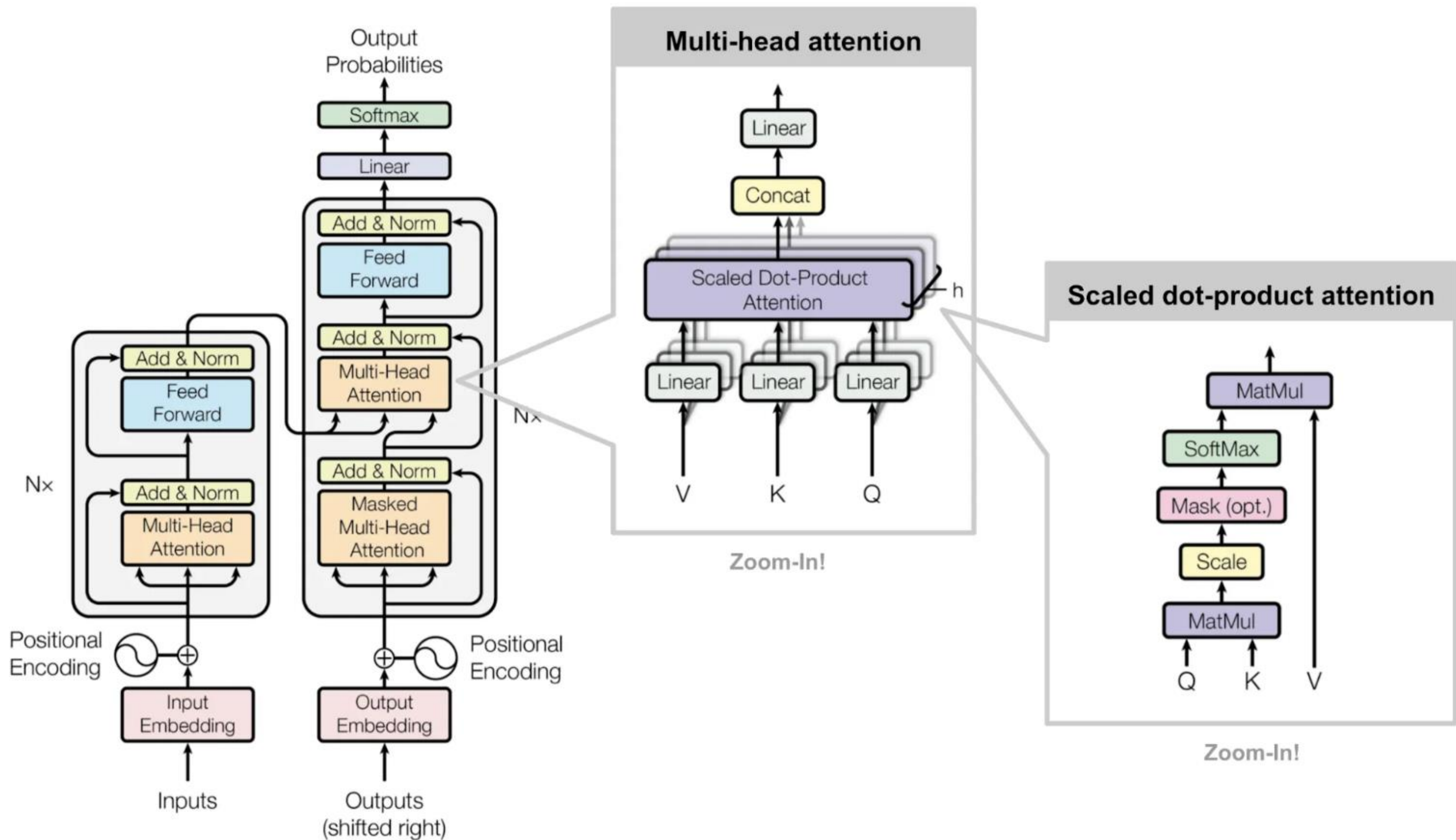
BERT

Encoder

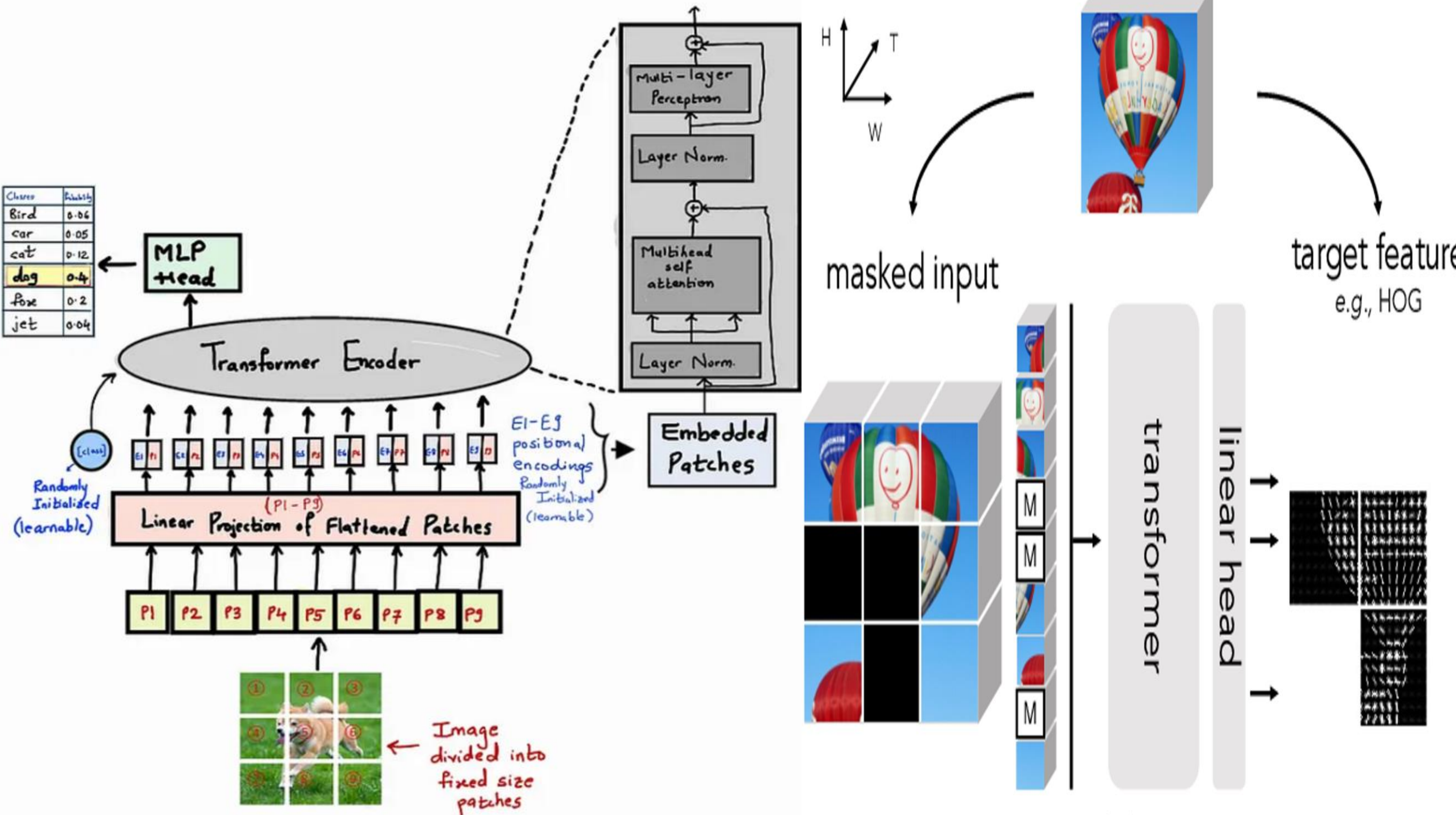


GPT

Decoder



Classes	Probability
Bird	0.06
car	0.05
cat	0.12
dog	0.4
fox	0.2
jet	0.04



Cropped Image

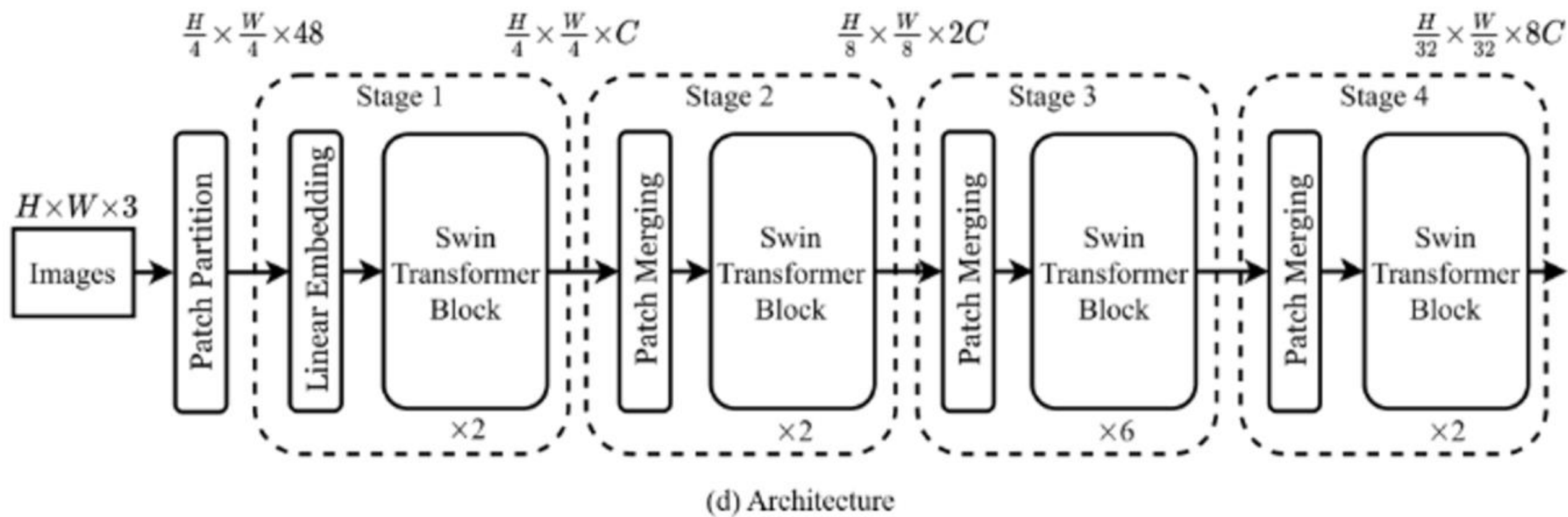
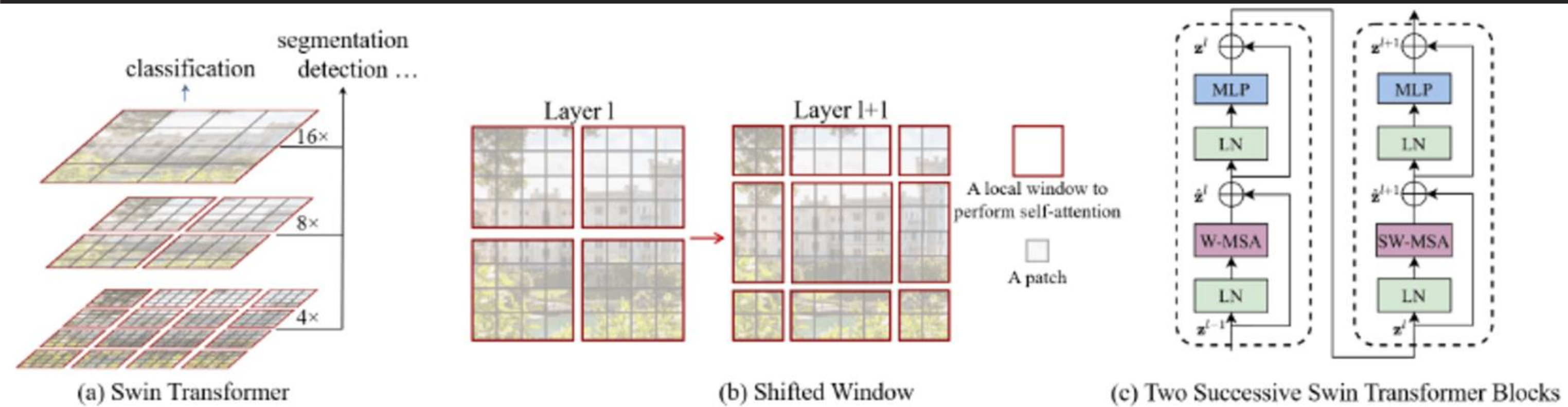


Image Patches



Flattened Image Patches

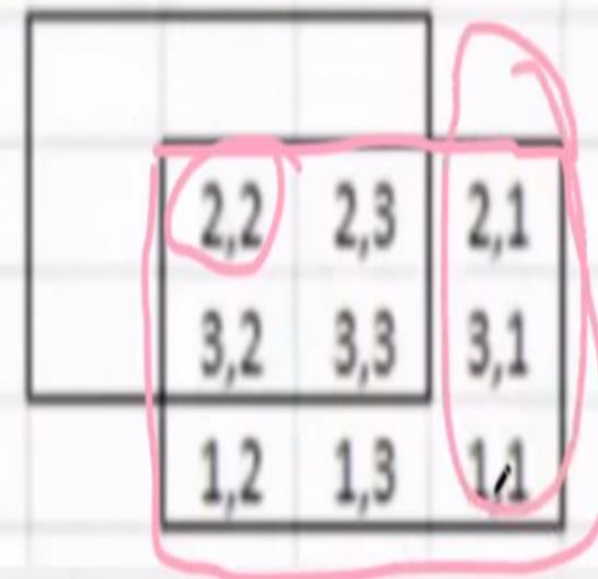
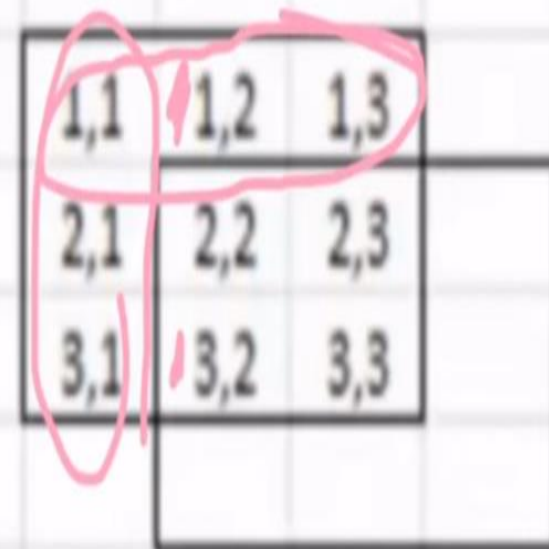
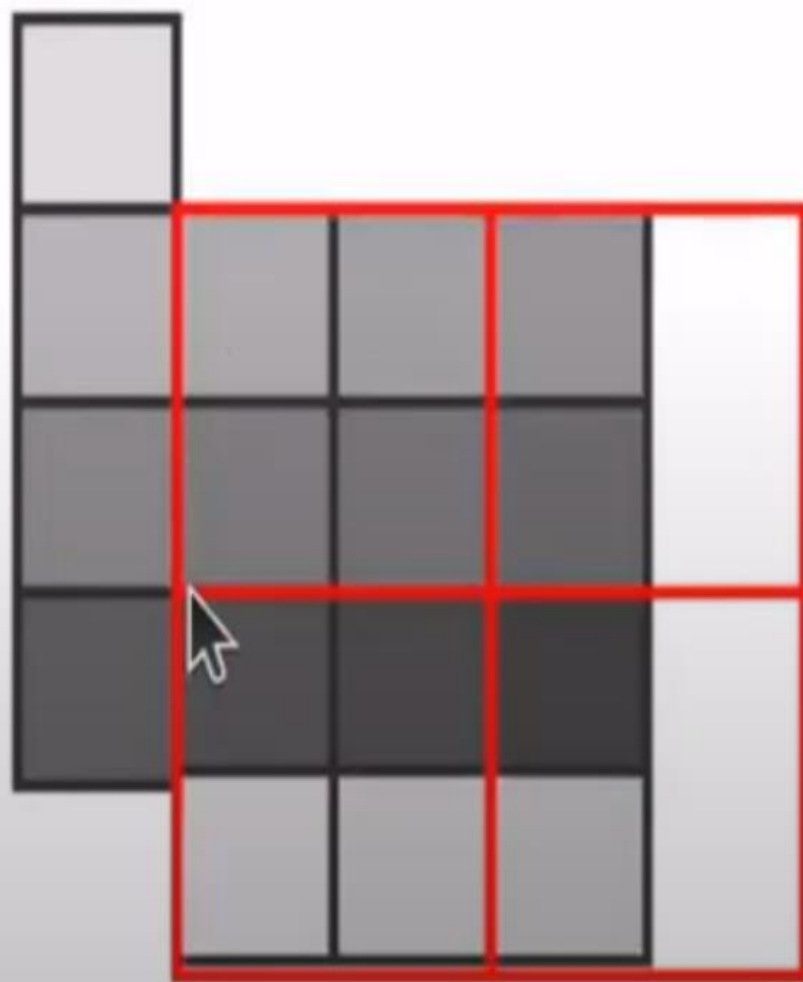




Shifted Window MSA

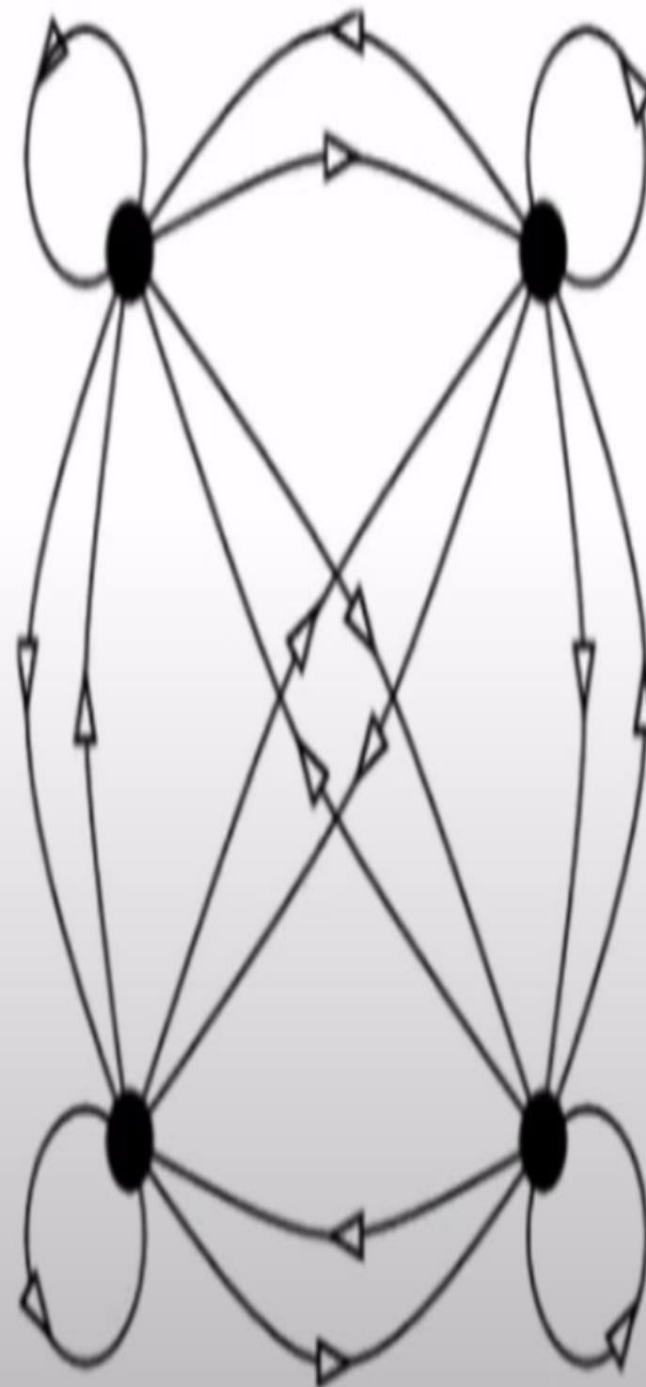
Step 2: For efficient batch computation, move patches into empty slots to create a complete window.

This is known as 'cyclic shift' in the paper.



Upper-Lower		2,2	2,3	2,1	3,2	3,3	3,1	1,2	1,3	1,1
(Last Row)	2,2	0	0	0	0	0	0	mask		
	2,3	0	0	0	0	0	0			
	2,1	0	0	0	0	0	0			
	3,2	0	0	0	0	0	0			
	3,3	0	0	0	0	0	0			
	3,1	0	0	0	0	0	0			
	1,2							0	0	0
	1,3							0	0	0
	1,1							0	0	0

-1, -1	-1, 0	-1, 1	-1, 2
0, -1	0, 0	0, 1	0, 2
1, -1	1, 0	1, 1	1, 2
2, -1	2, 0	2, 1	2, 2



Patch Merging

Assuming that $n=2$, and each group consists of 2×2 neighboring patches

Step 1: Split input image into groups of 2×2

Step 2: In each group, stack the patches depth-wise

Step 3: Combine the stacked groups

