



Swin Transformer

Артем Исмагилов

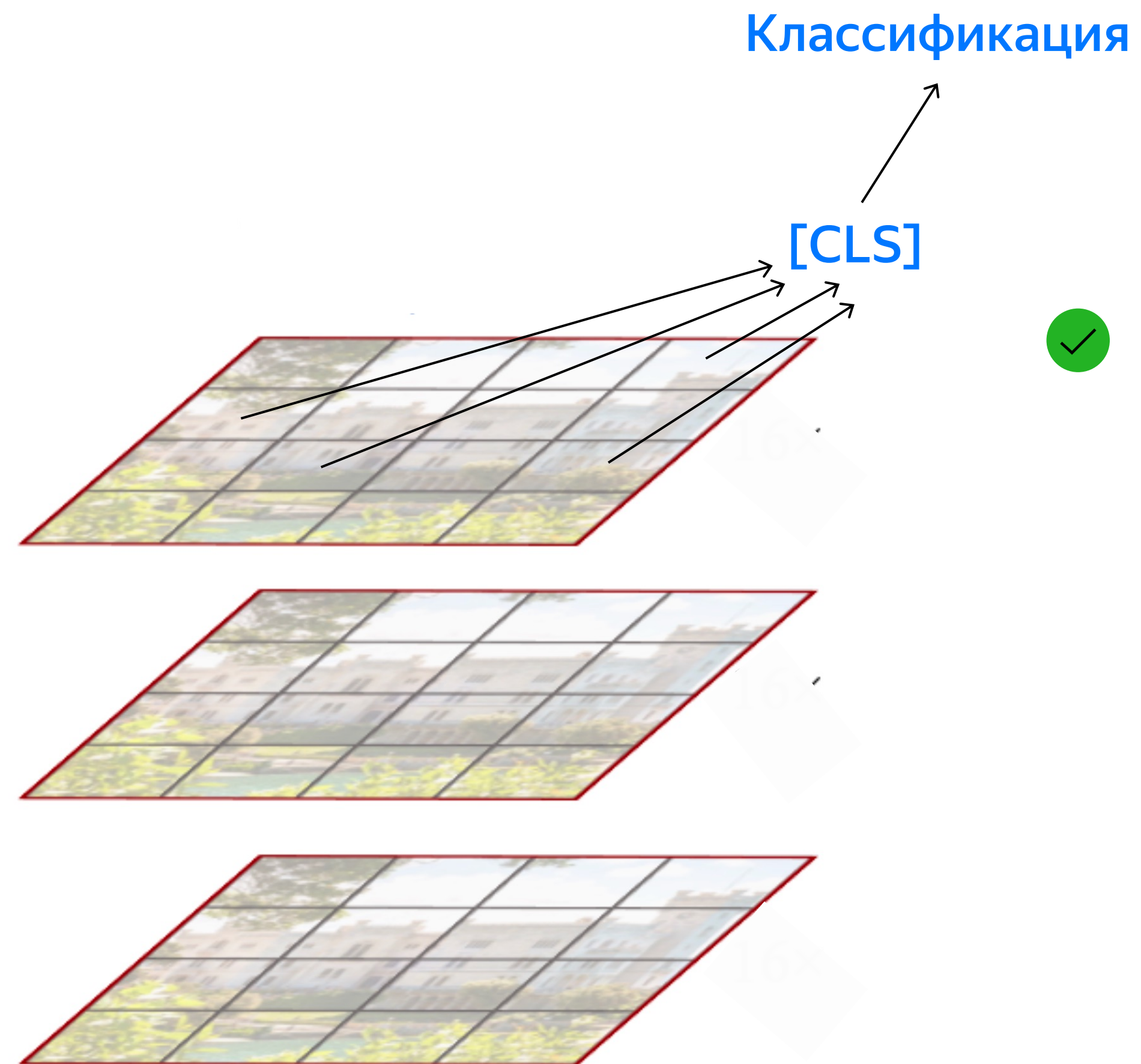
Intro

Swin Transformer

Как применять для разных задач

Zero Shot и One Shot детекция

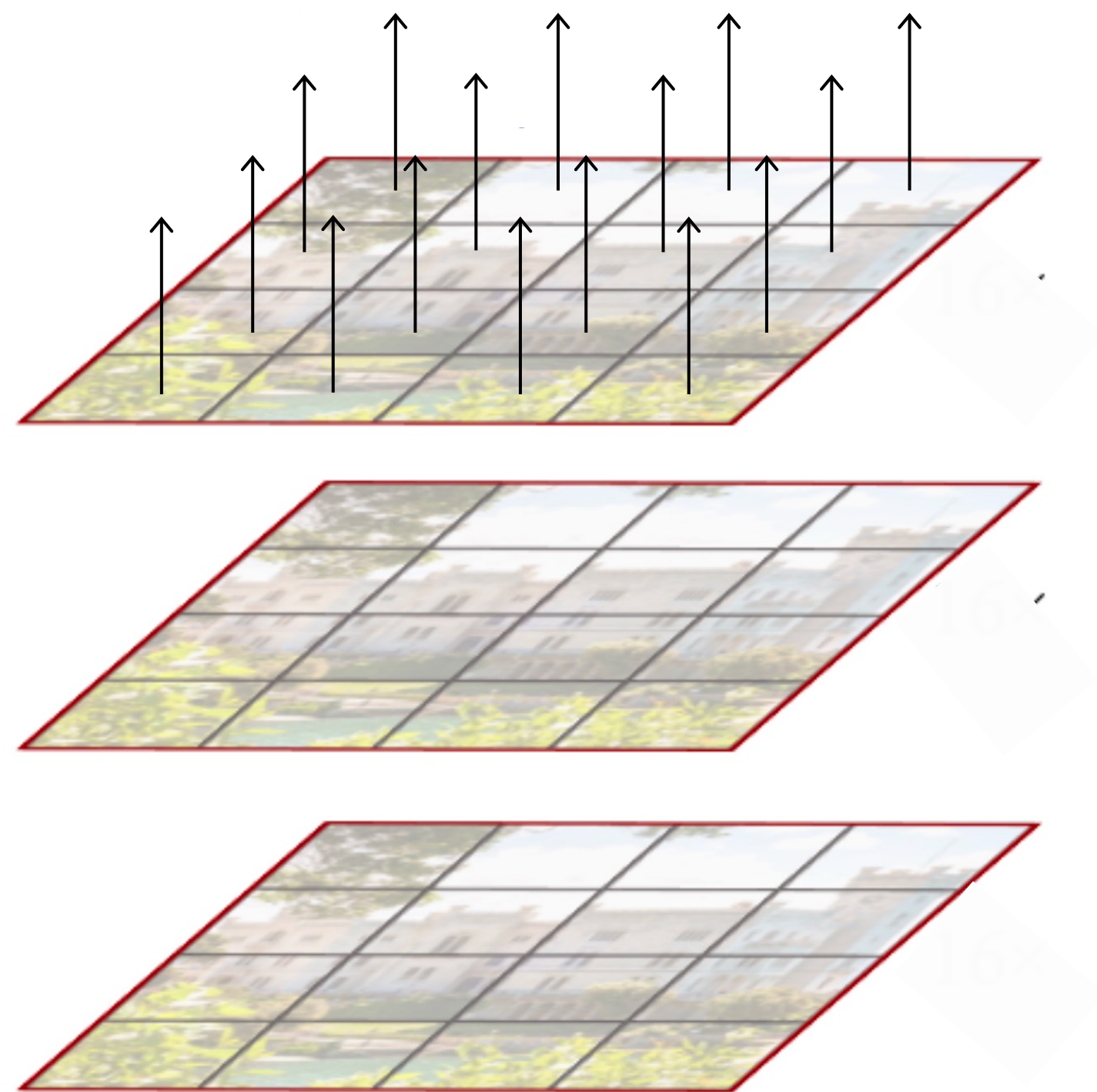
ViT



✓ Хорошо извлекает глобальные признаки

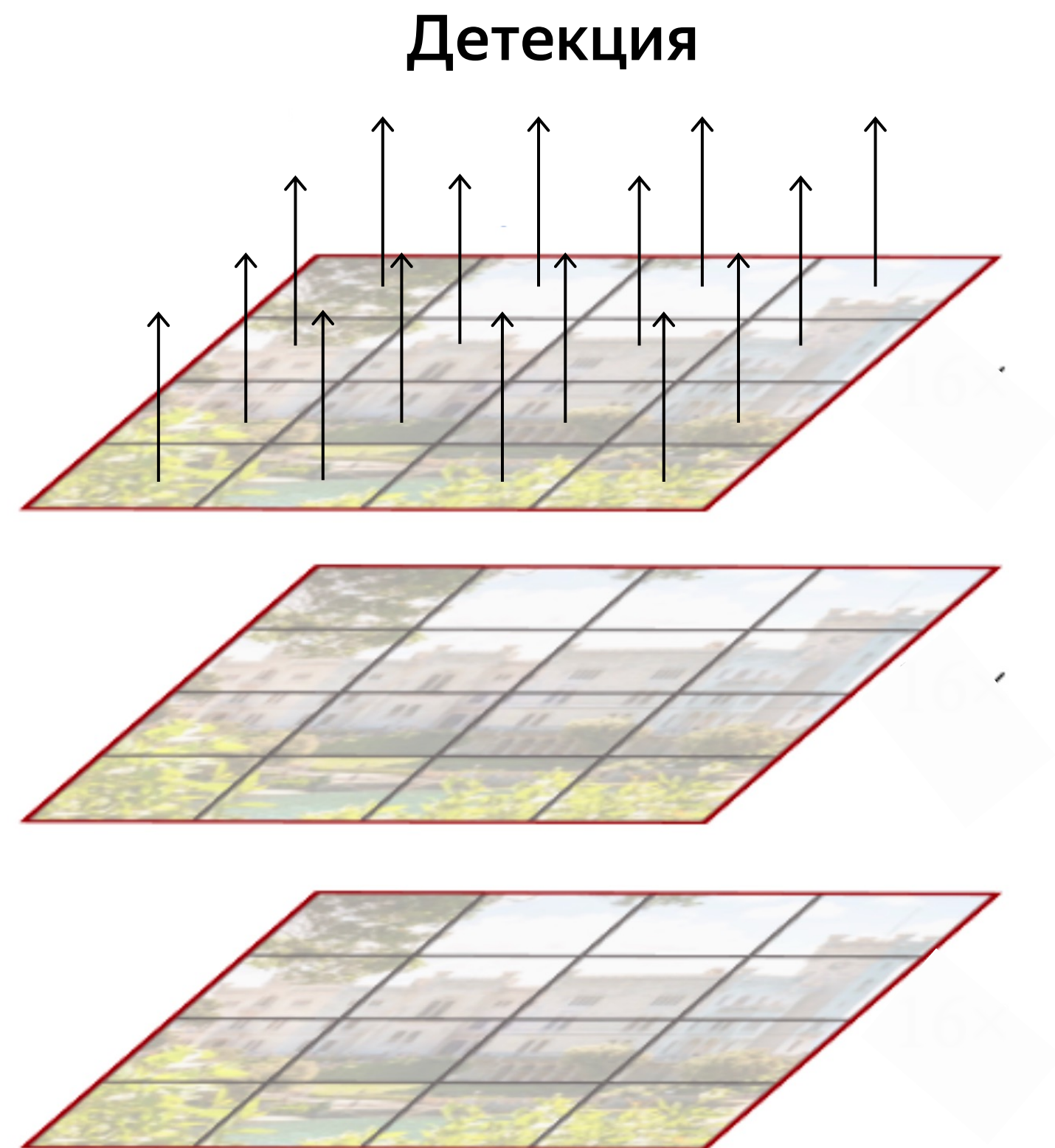
ViT

Детекция



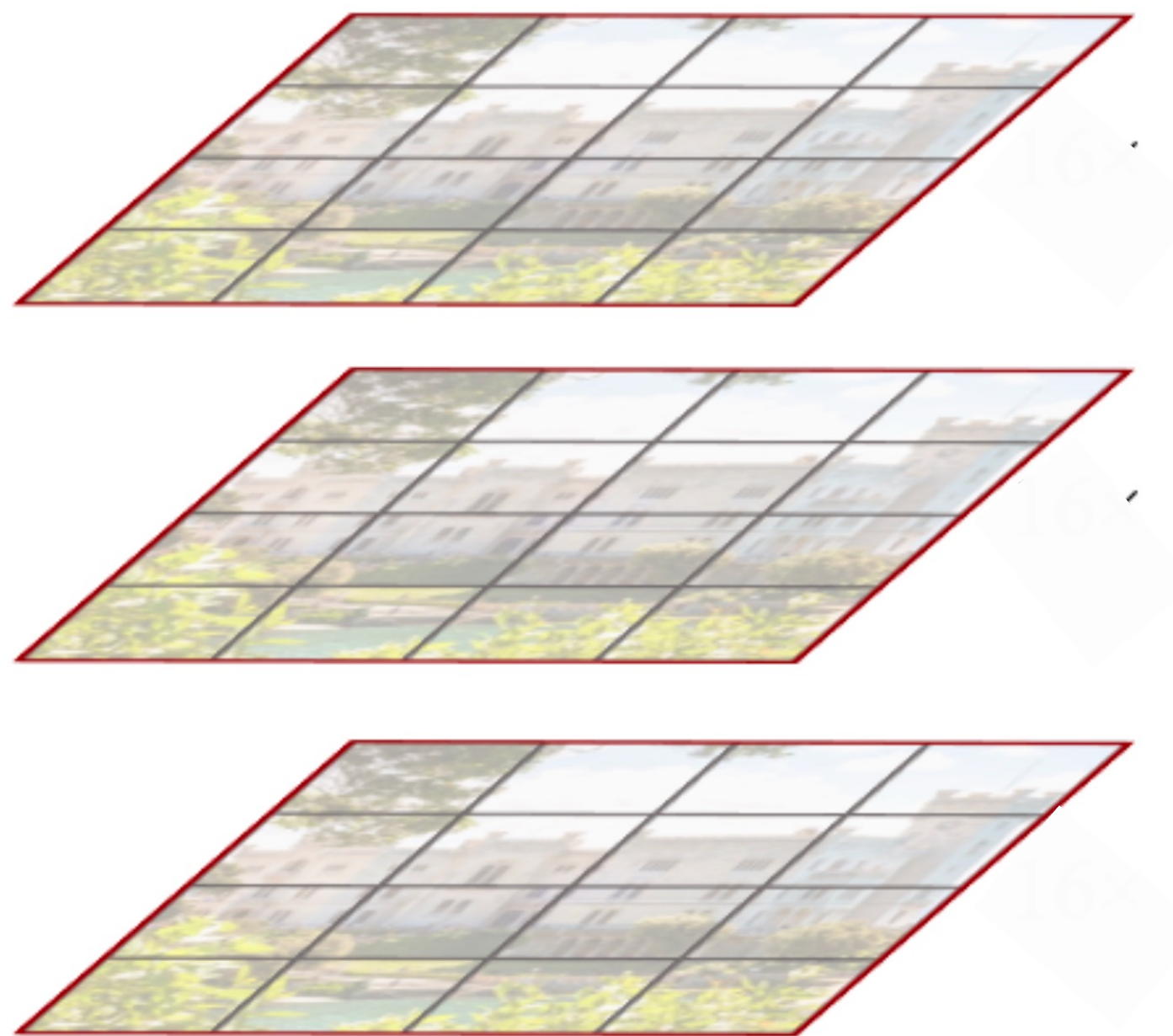
- ✓ Хорошо извлекает глобальные признаки
- ✓ Хорошо извлекает локальные признаки

ViT



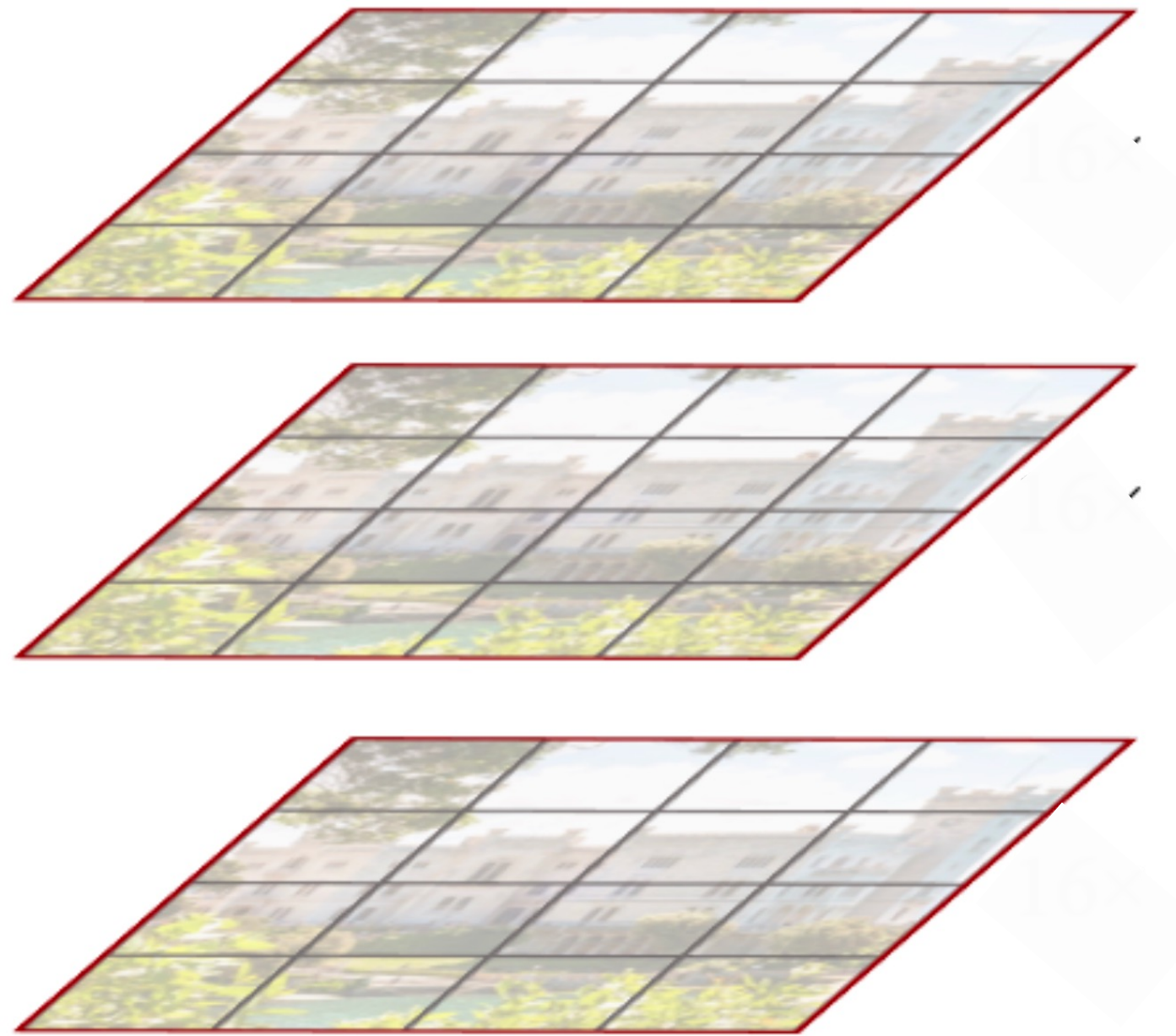
- ✓ Хорошо извлекает глобальные признаки
- ✓ Хорошо извлекает локальные признаки
- ✓ Локальные признаки содержат весь контекст

ViT



- ✓ Хорошо извлекает глобальные признаки
- ✓ Хорошо извлекает локальные признаки
- ✓ Локальные признаки содержат весь контекст
- ✗ В Self-Attention очень много патчей

ViT



- ✓ Хорошо извлекает глобальные признаки
- ✓ Хорошо извлекает локальные признаки
- ✓ Локальные признаки содержат весь контекст
- ✗ В Self-Attention очень много патчей
- ✗ Увеличить разрешение карты признаков очень дорого

Intro

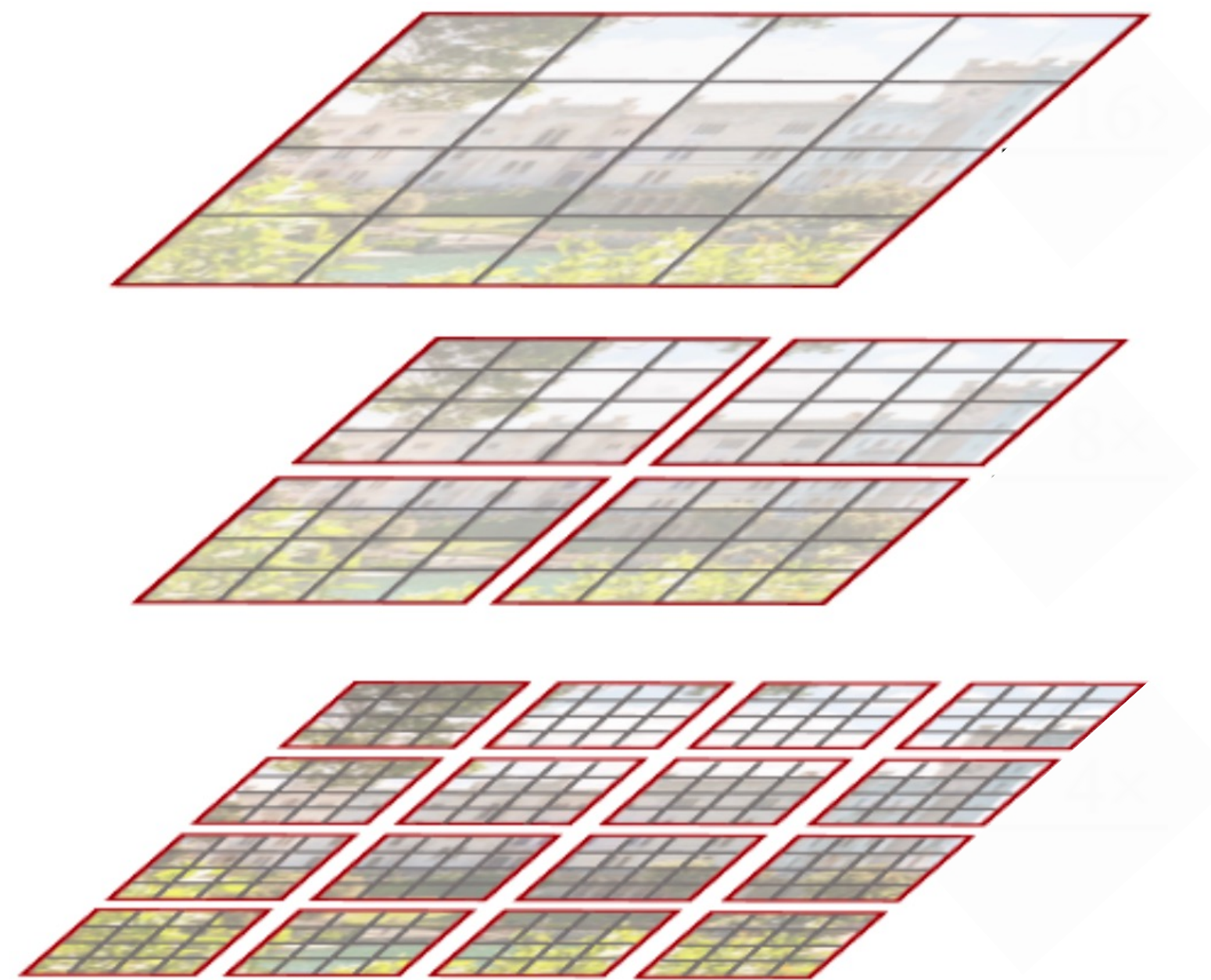
Swin Transformer

Как применять для разных задач

Zero Shot и One Shot детекция

Swin Transformer

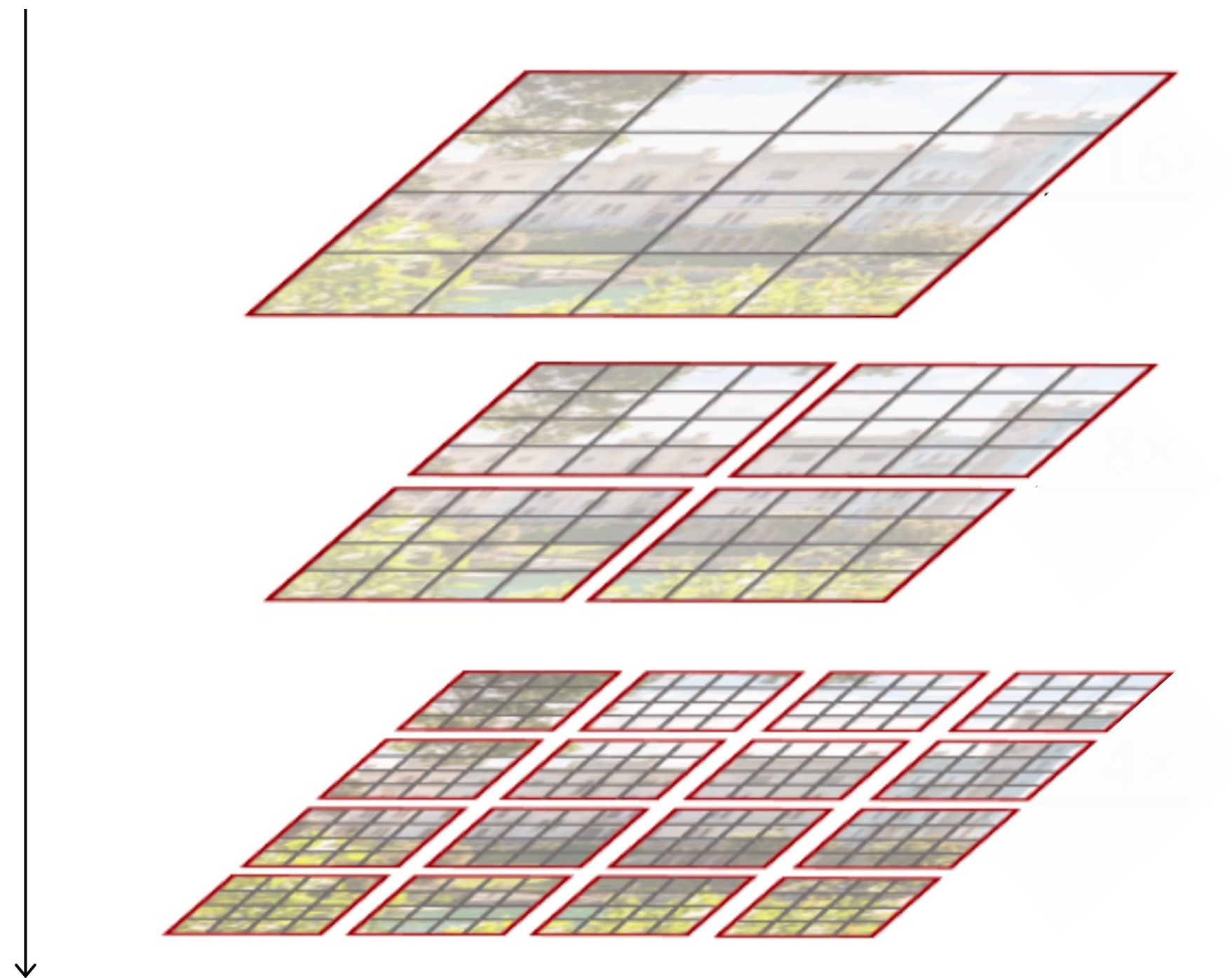
Shifted Windows



Swin Transformer

Shifted Windows

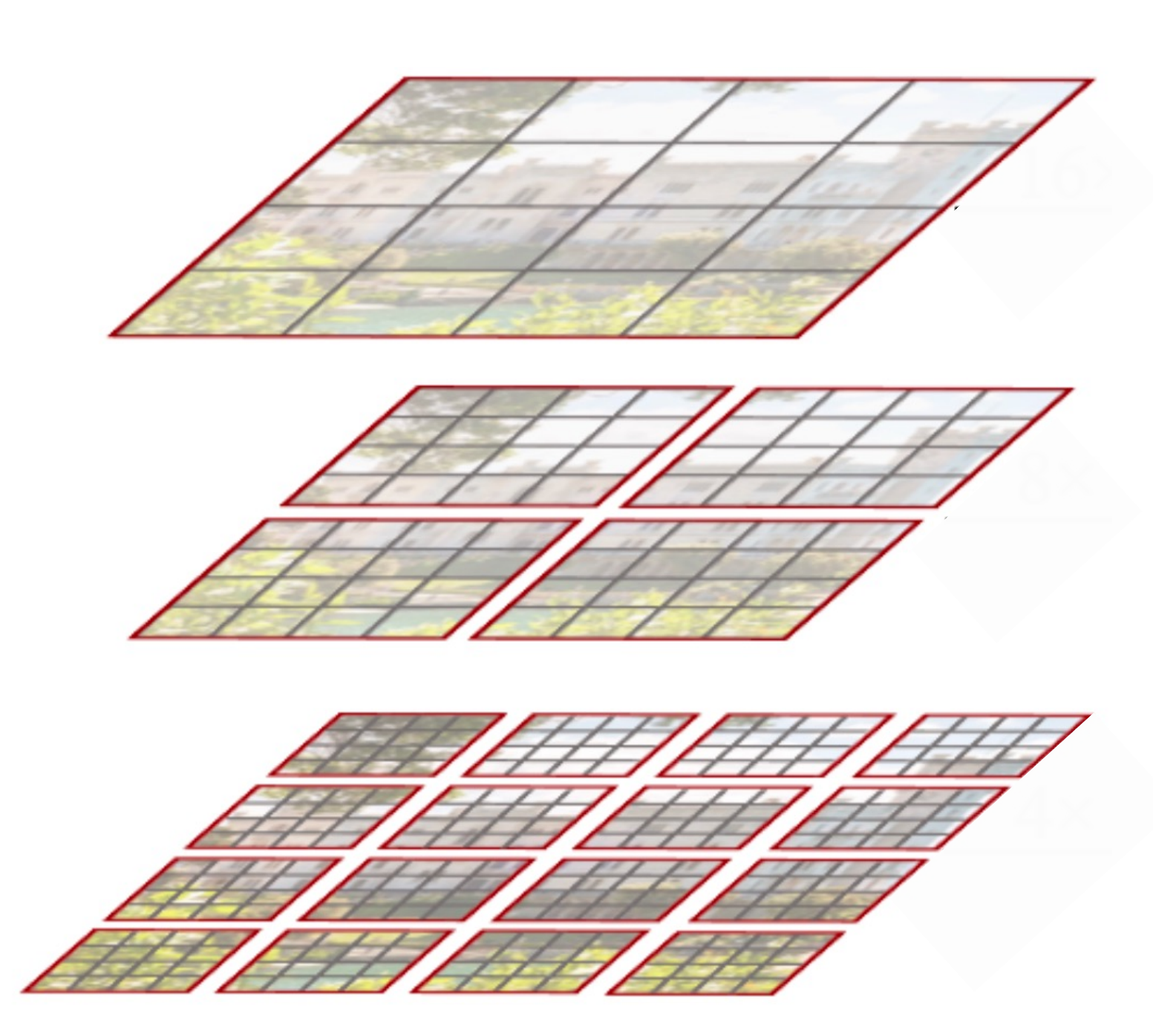
Больше разрешение



Swin Transformer

Shifted Windows

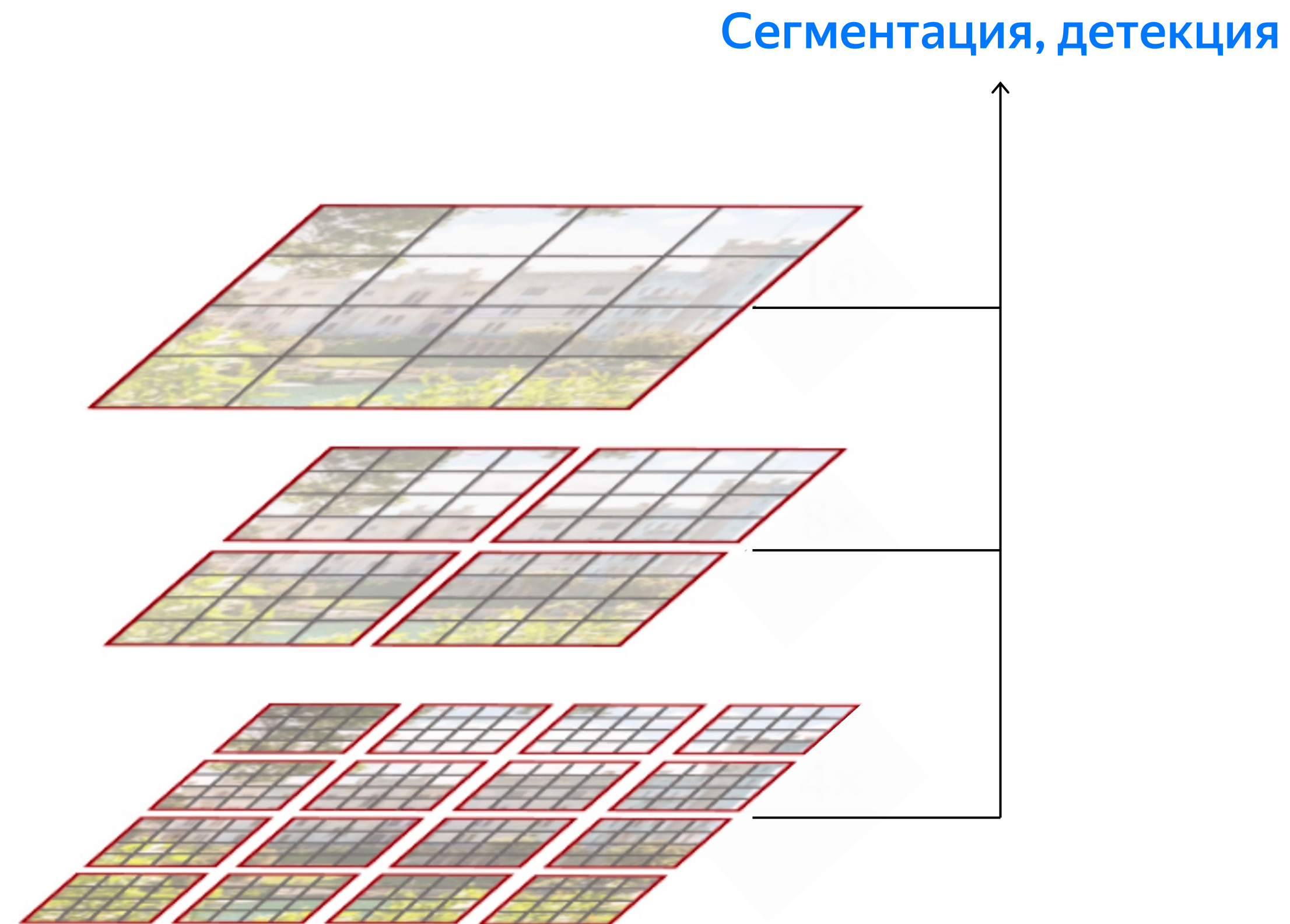
Больше разрешение



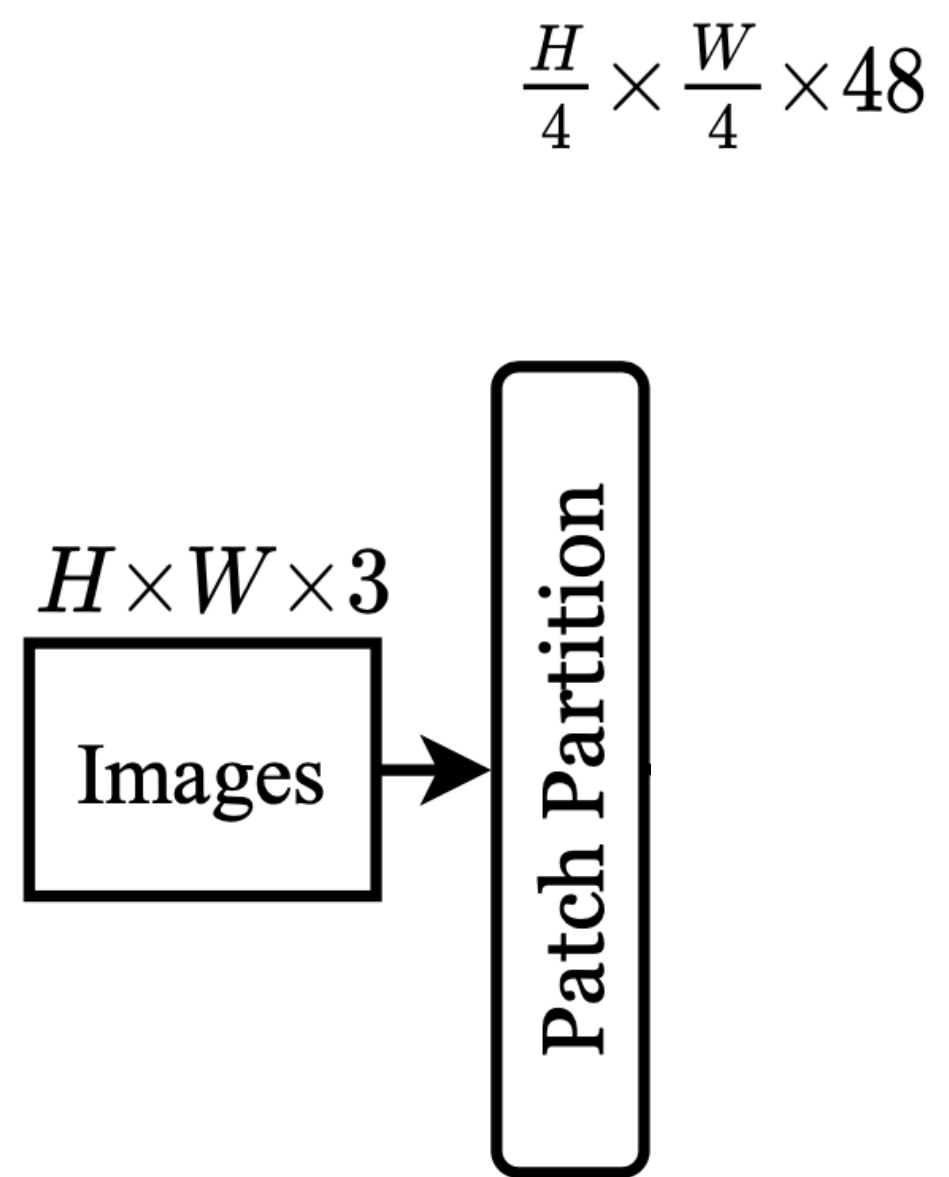
Больше контекст

Swin Transformer

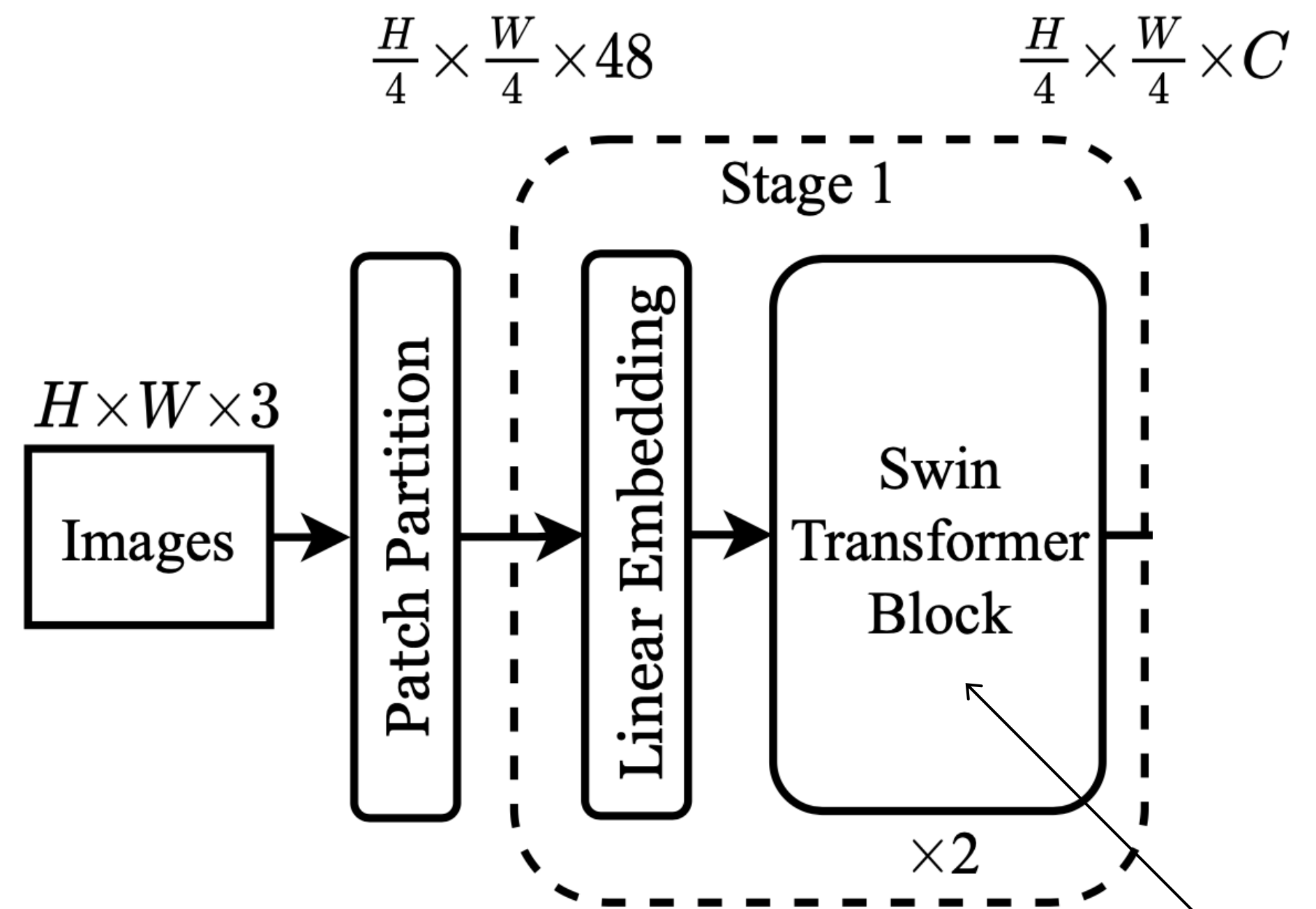
Shifted Windows



Swin Transformer

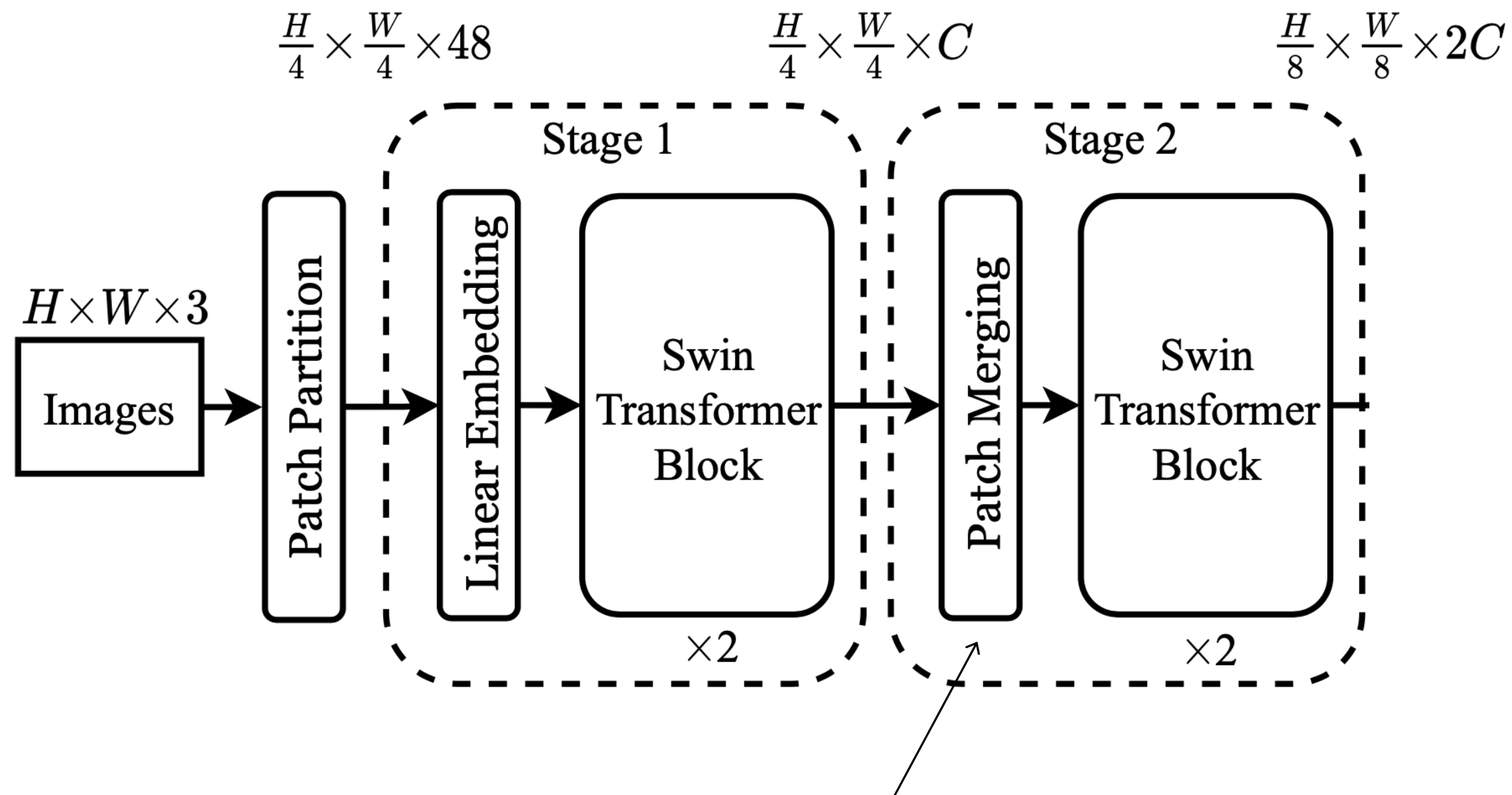


Swin Transformer



Обычный блок трансформера, но attention только внутри окон

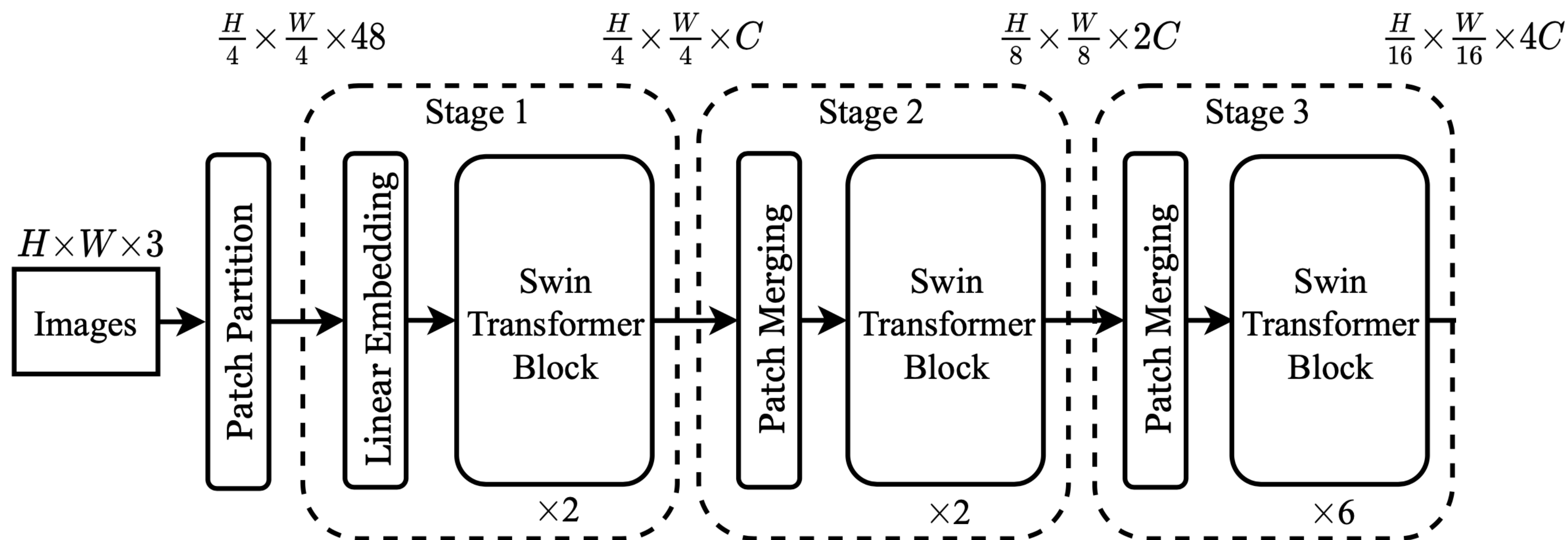
Swin Transformer



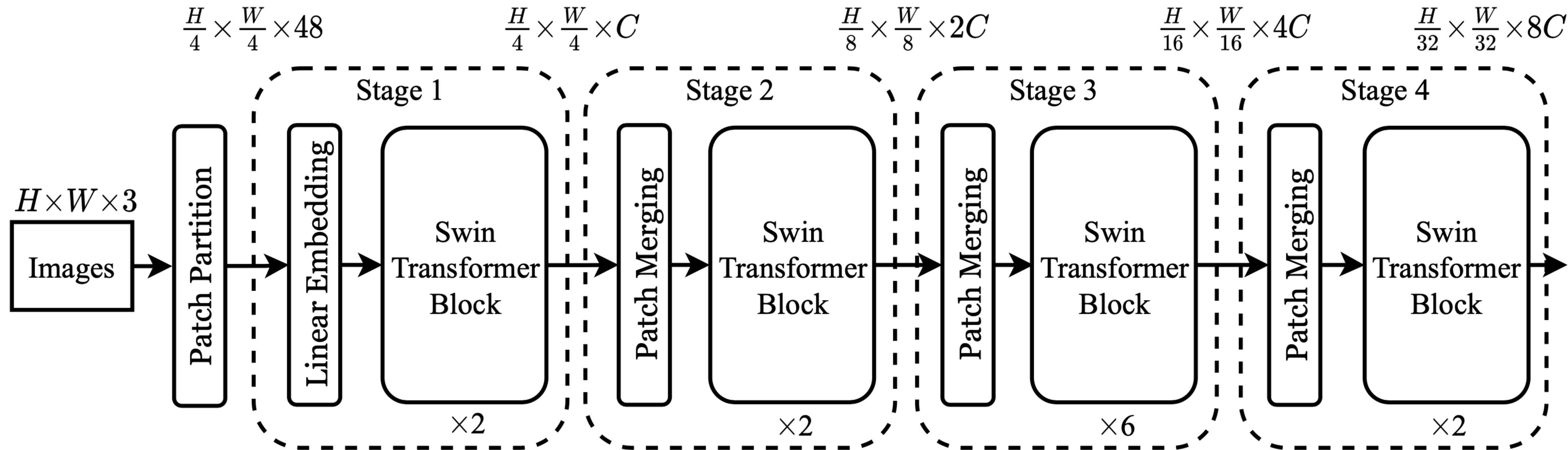
Конкатенируем эмбеддинги патчей в области 2×2
и уменьшаем размерность в 2 раза

Сторона окна Attention увеличивается в два раза

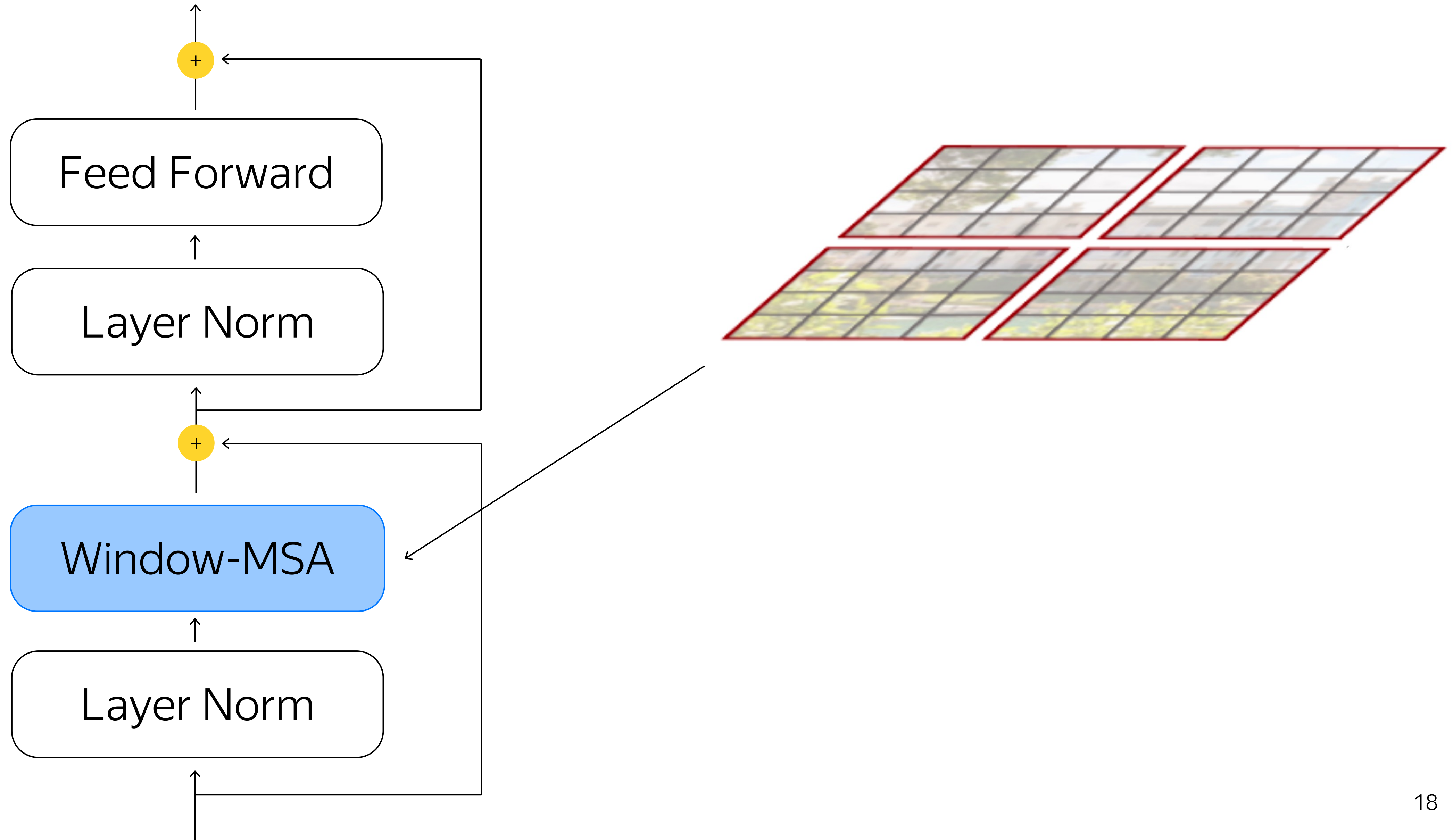
Swin Transformer



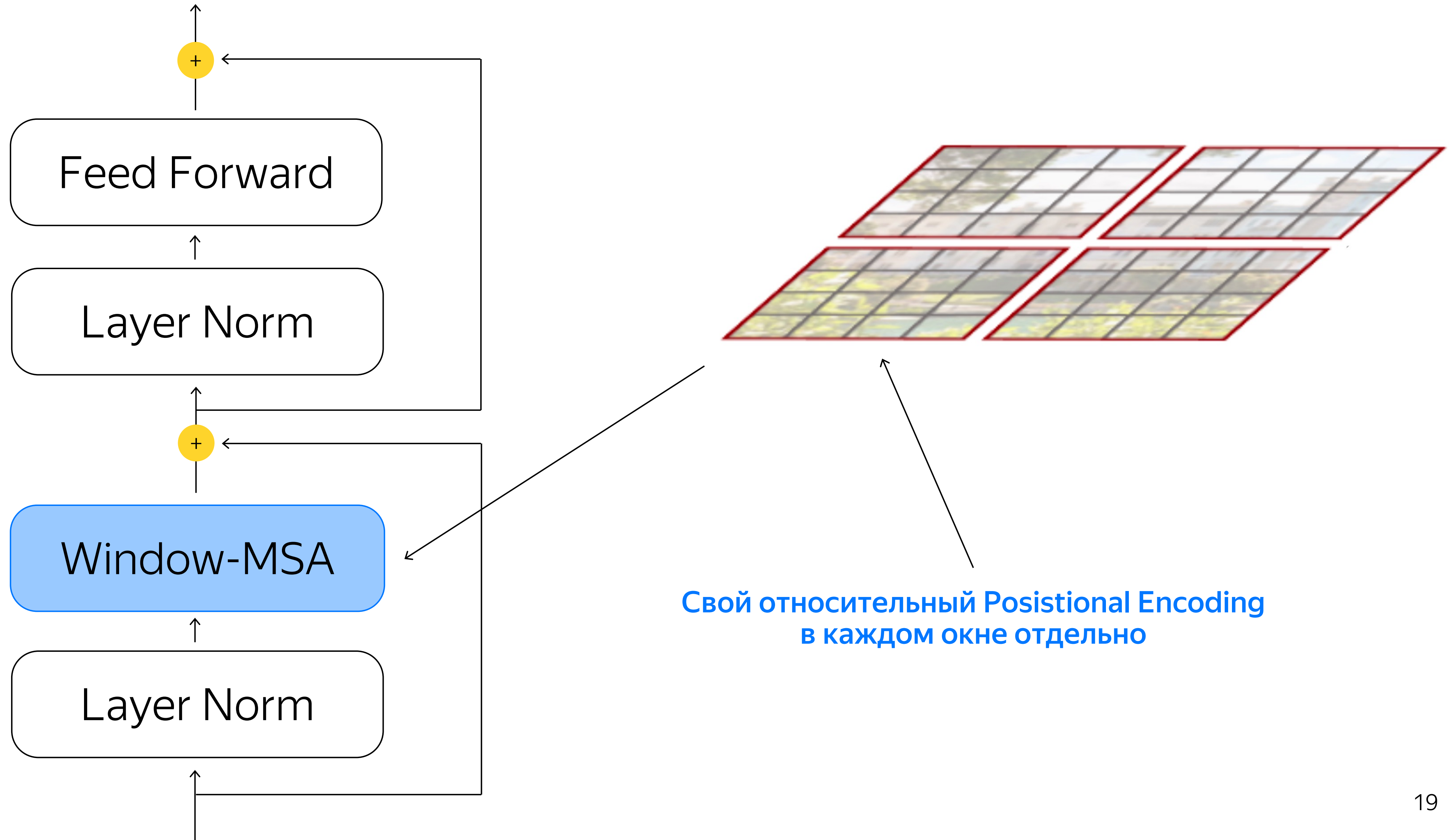
Swin Transformer



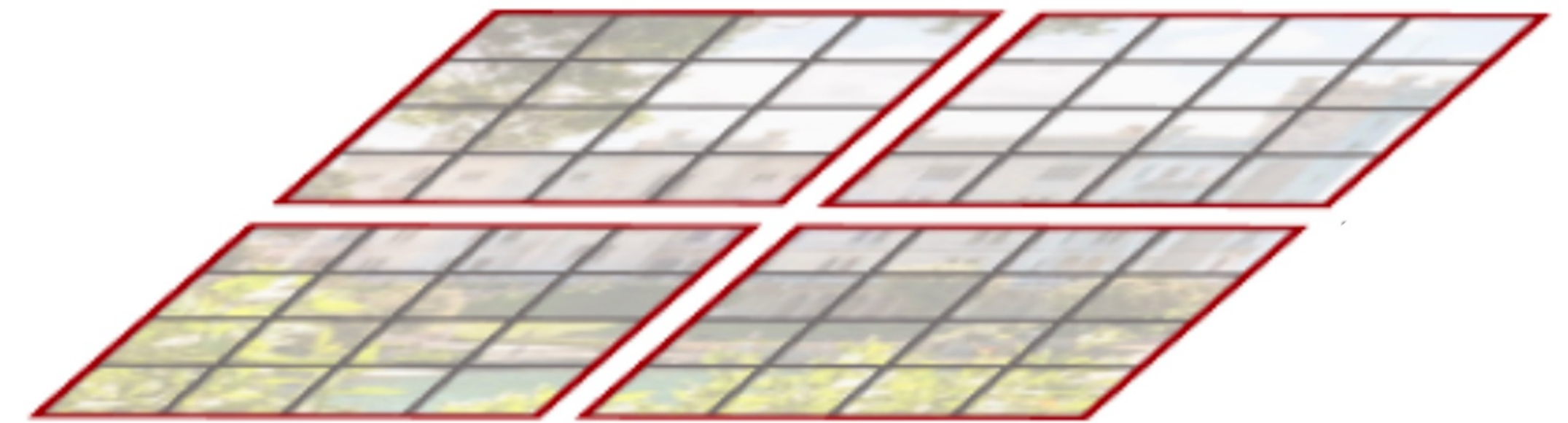
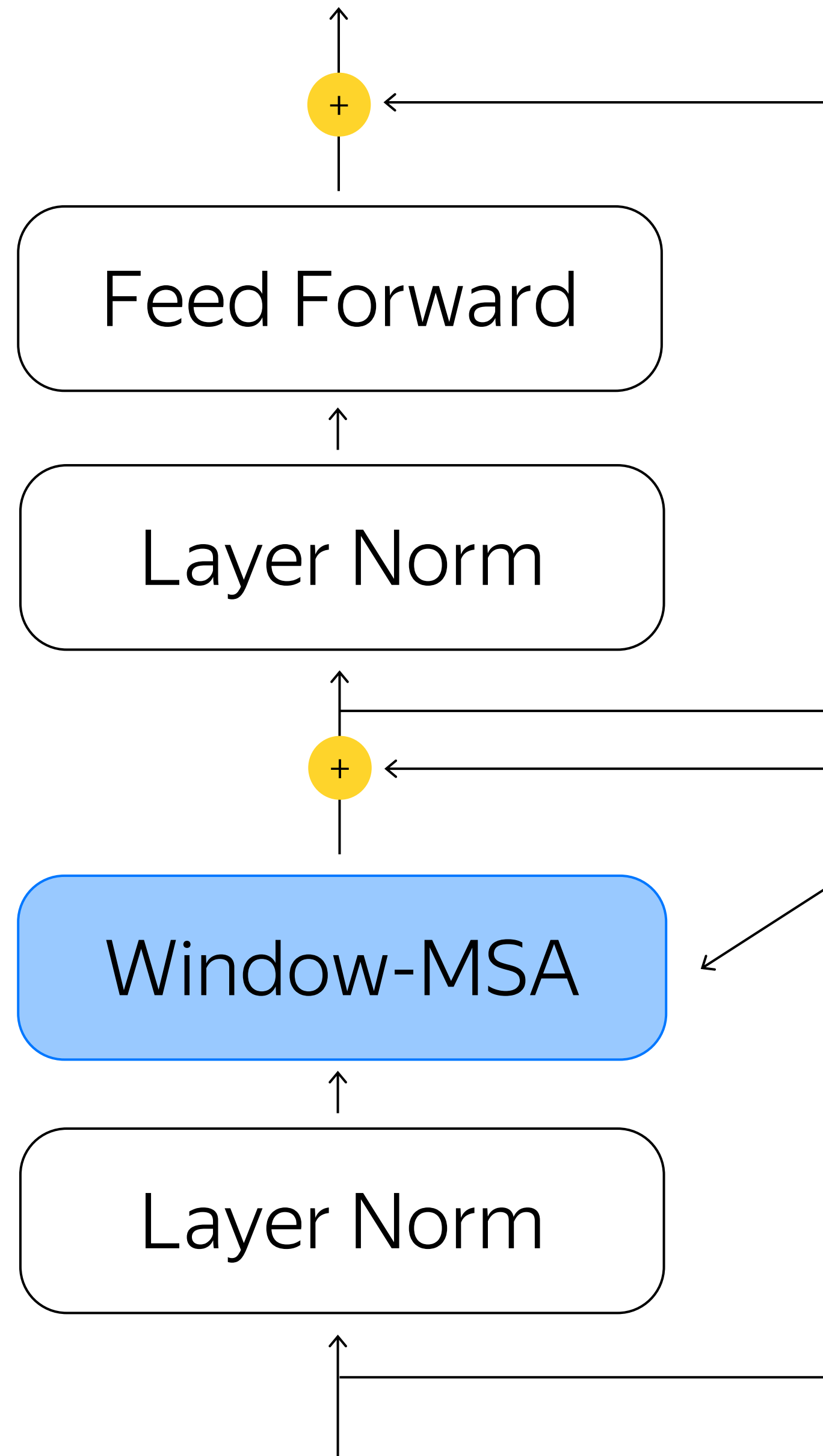
Swin Transformer Block



Swin Transformer Block



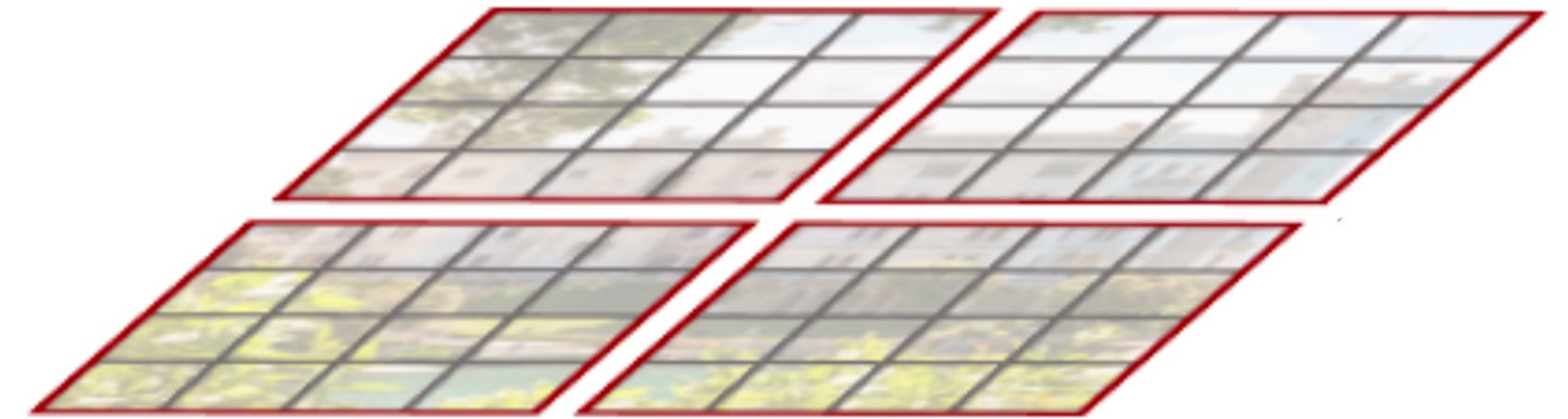
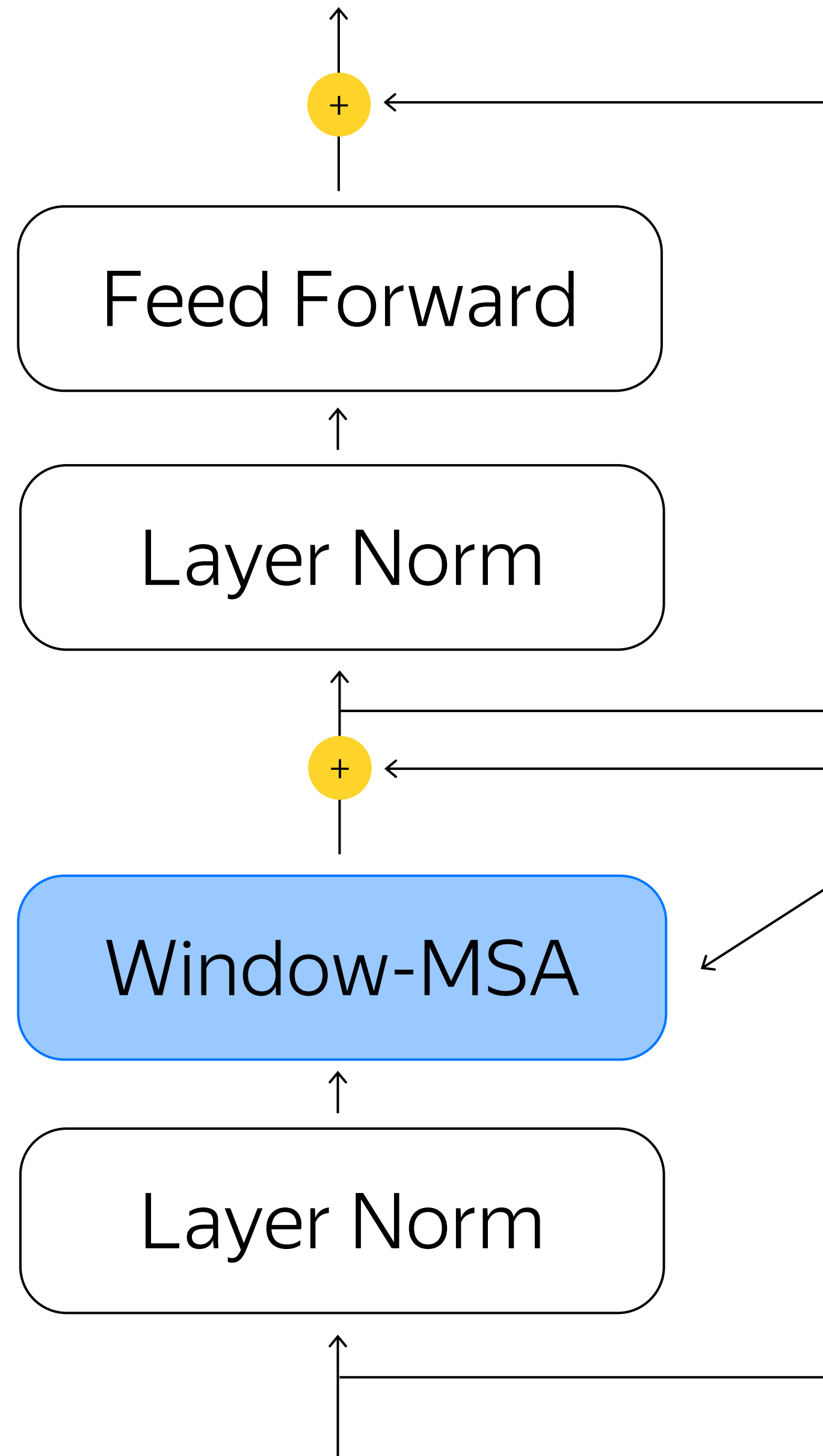
Swin Transformer Block



Обычный Self-Attention

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C,$$

Swin Transformer Block



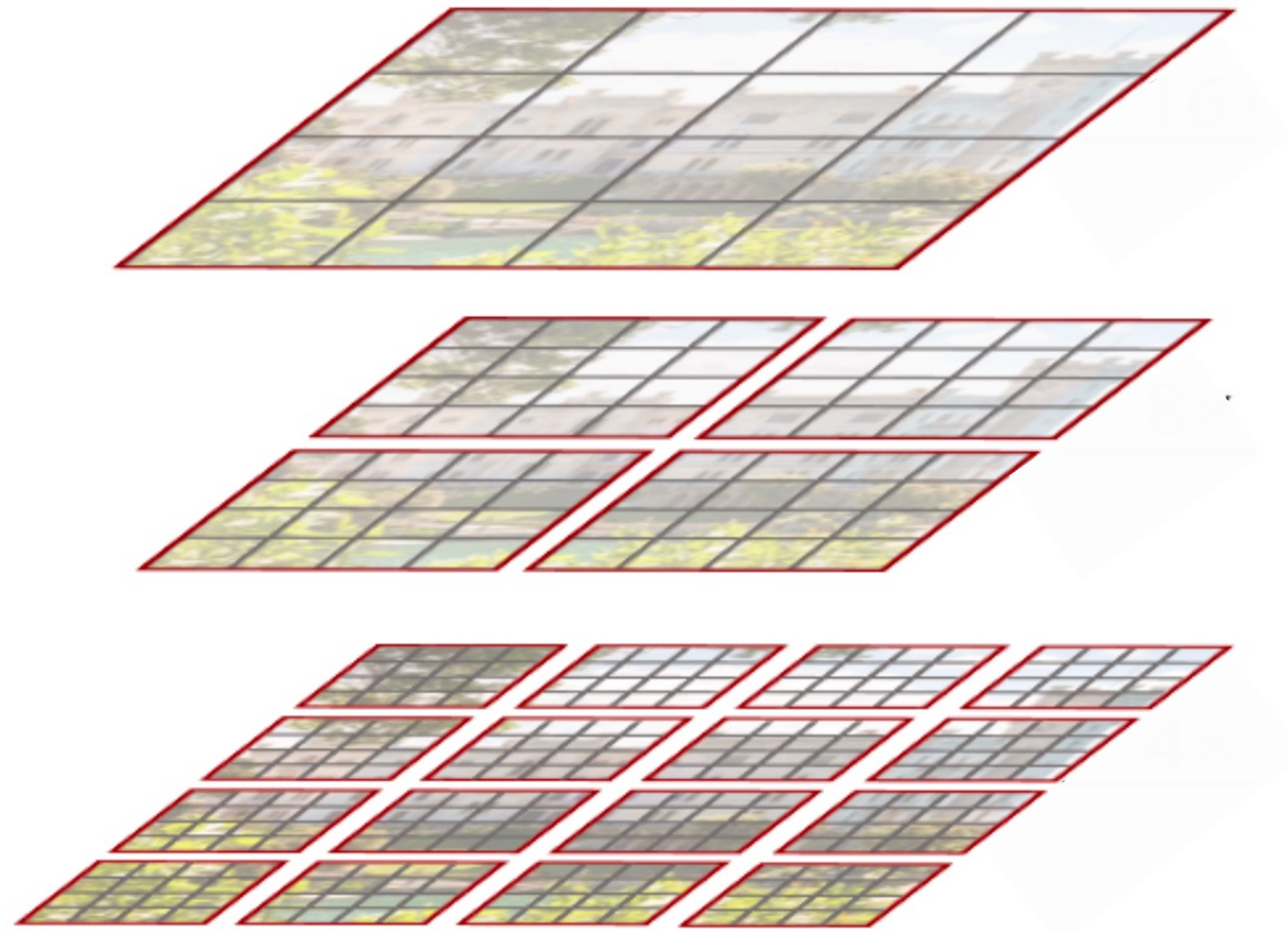
Обычный Self-Attention

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C,$$

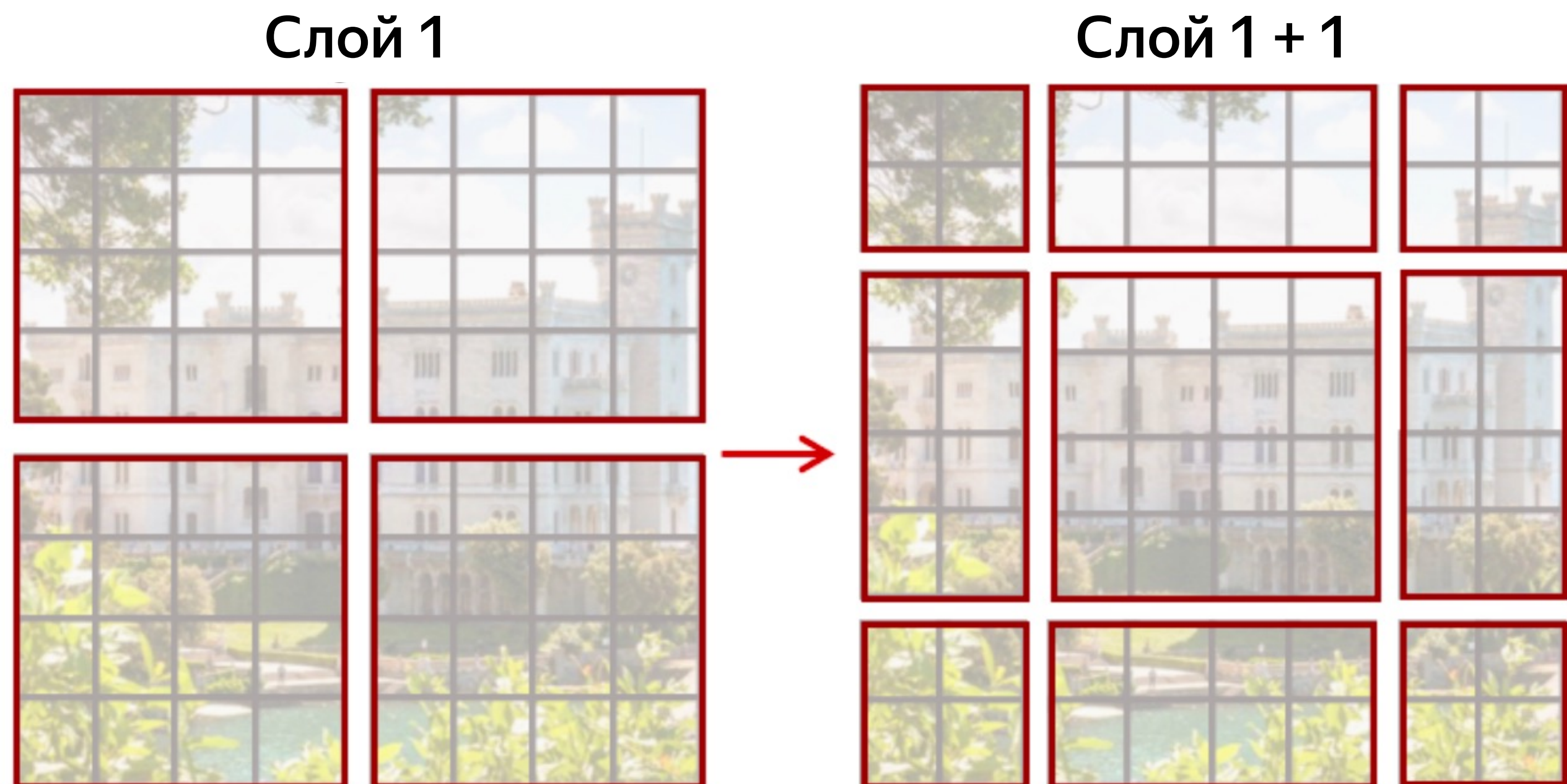
Window Self-Attention с окном M

$$\Omega(\text{W-MSA}) = 4hwC^2 + 2M^2hwC,$$

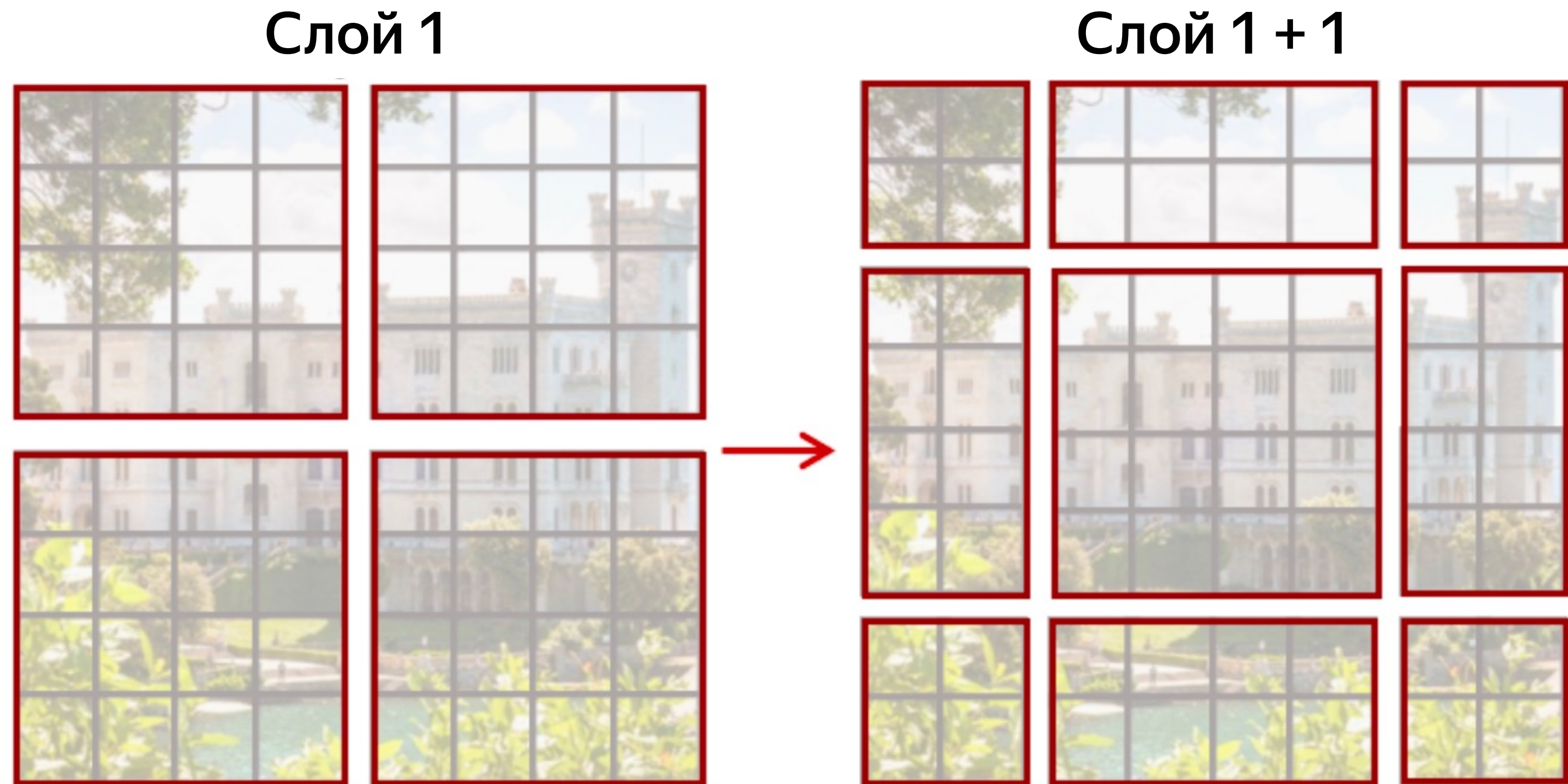
Причем тут shifted windows?



Причем тут shifted windows?



Причем тут shifted windows?



К Attention в этом случае добавляется маска

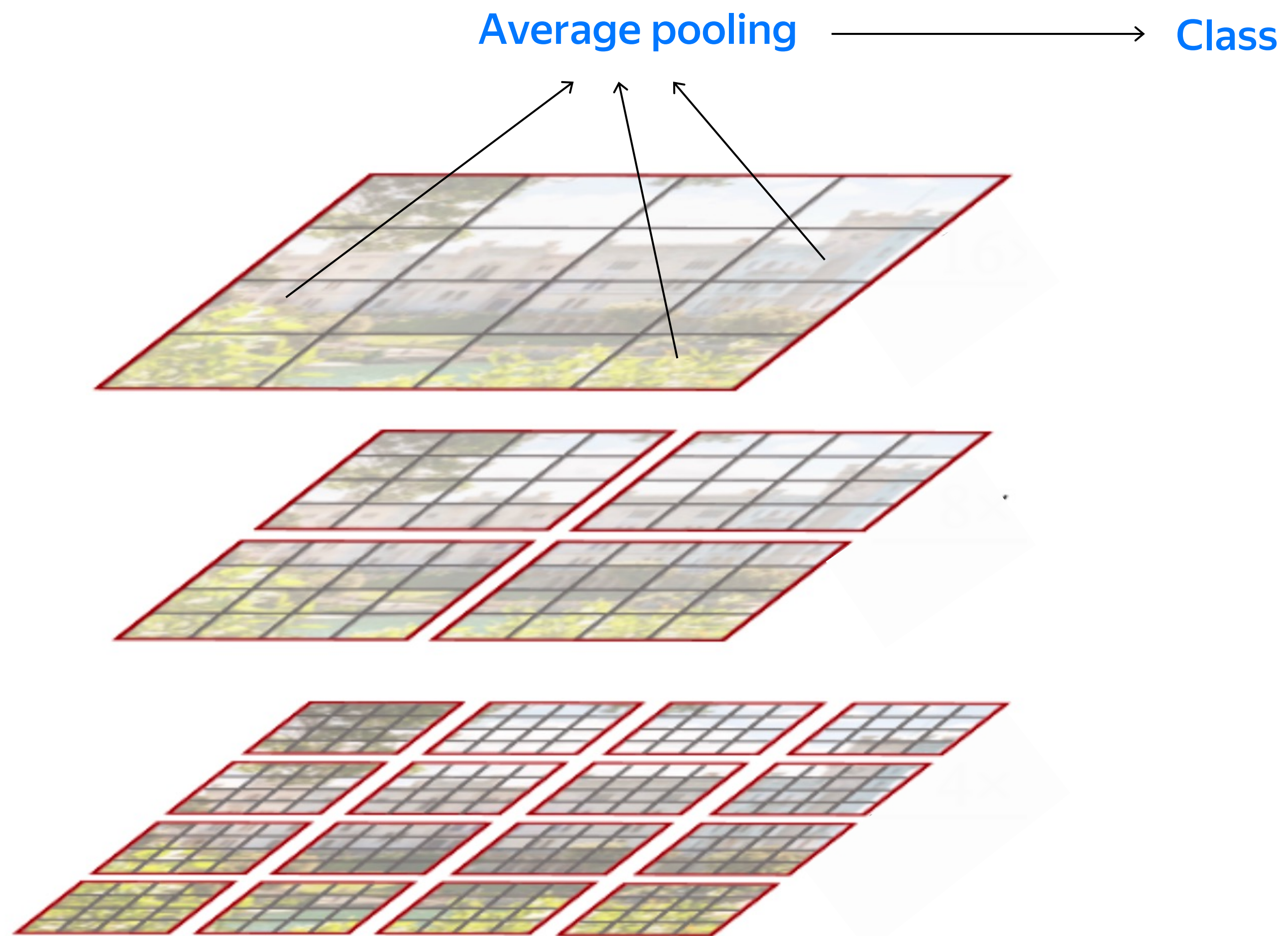
Intro

Swin Transformer

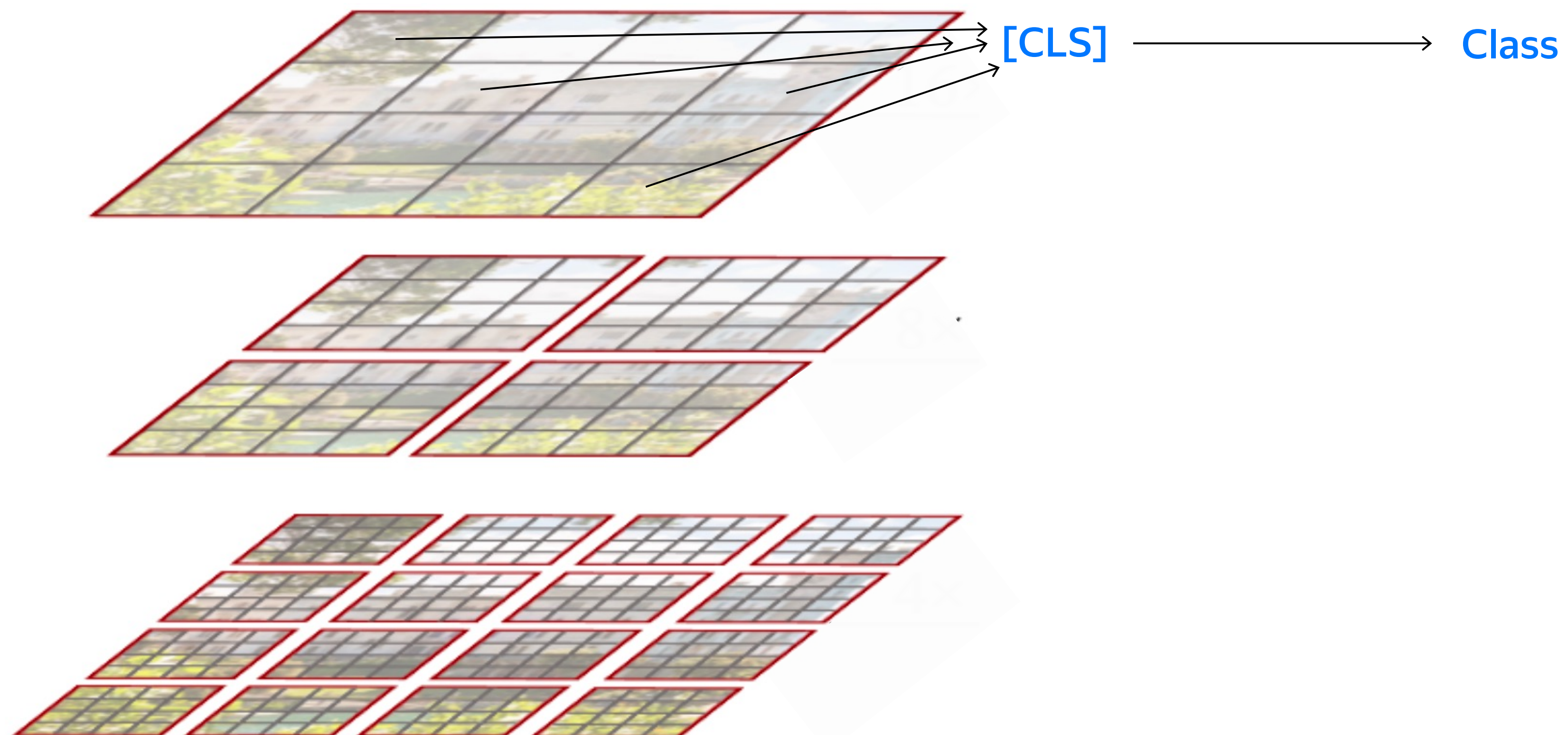
Как применять для разных задач

Zero Shot и One Shot детекция

Классификация для бедных



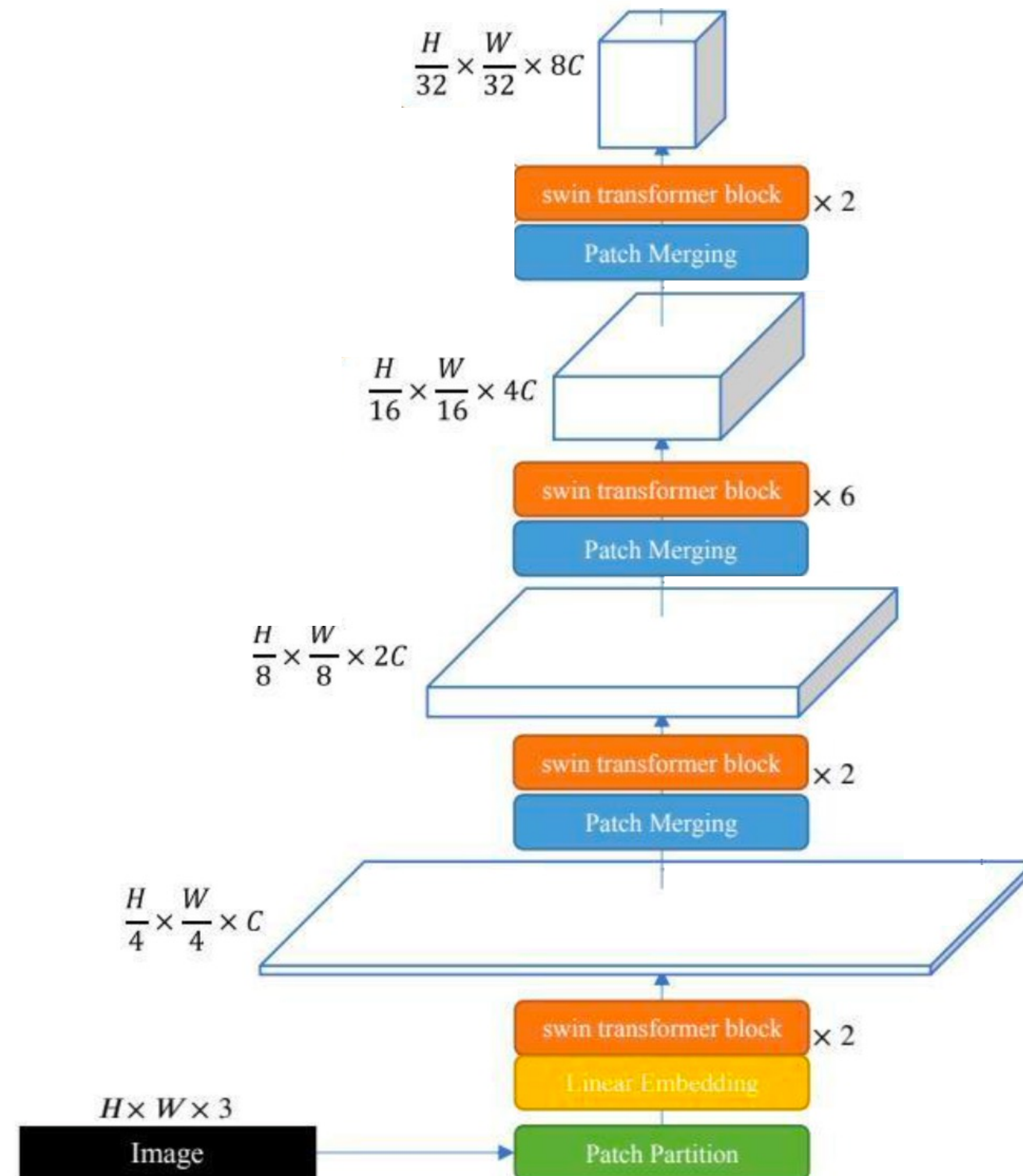
Классификация



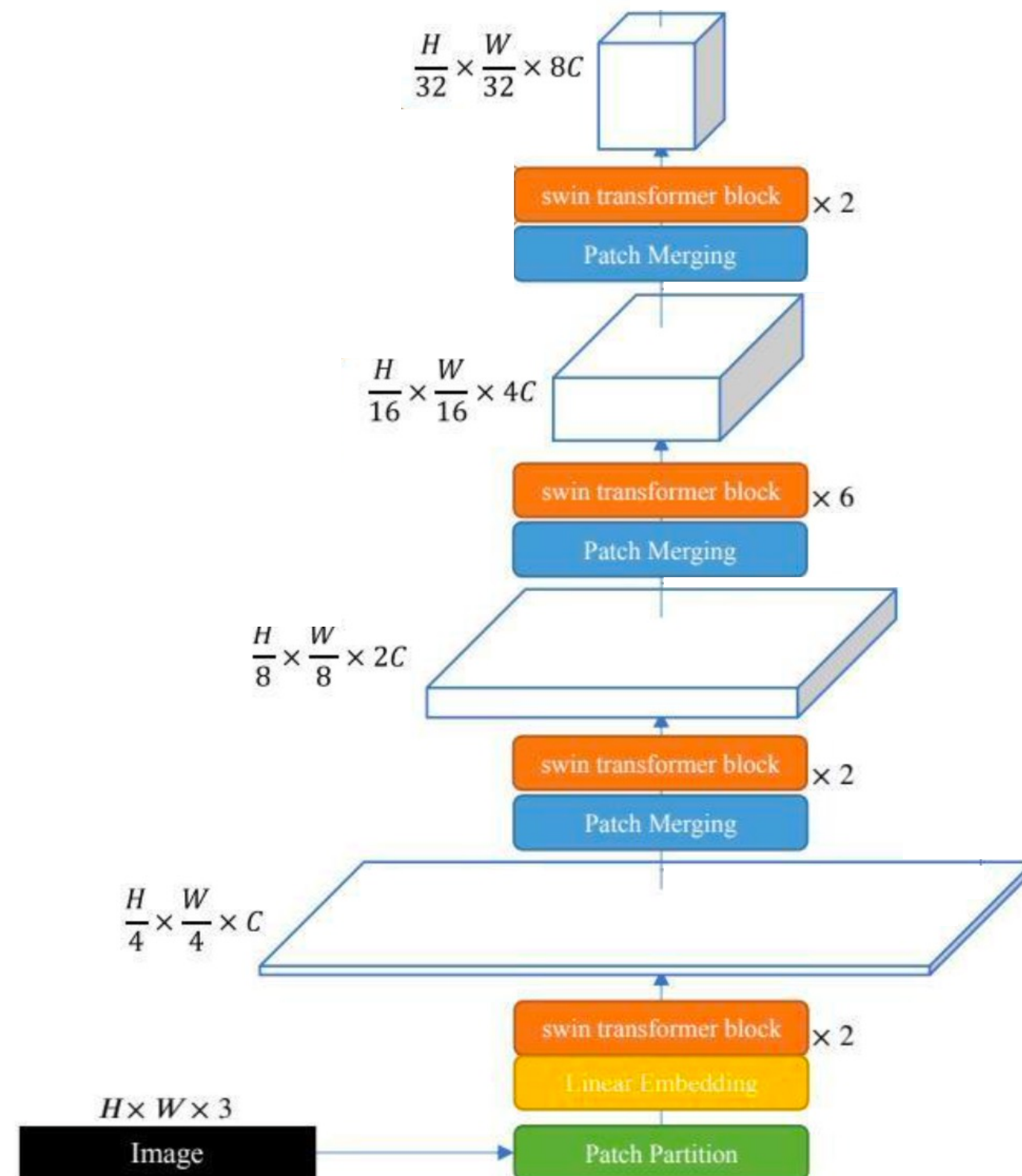
Классификация

Regular ImageNet-1K trained models					
method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
RegNetY-4G	224 ²	21M	4.0G	1156.7	80.0
RegNetY-8G	224 ²	39M	8.0G	591.6	81.7
RegNetY-16G	224 ²	84M	16.0G	334.7	82.9
EffNet-B3	300 ²	12M	1.8G	732.1	81.6
EffNet-B4	380 ²	19M	4.2G	349.4	82.9
EffNet-B5	456 ²	30M	9.9G	169.1	83.6
EffNet-B6	528 ²	43M	19.0G	96.9	84.0
EffNet-B7	600 ²	66M	37.0G	55.1	84.3
ViT-B/16	384 ²	86M	55.4G	85.9	77.9
ViT-L/16	384 ²	307M	190.7G	27.3	76.5
DeiT-S	224 ²	22M	4.6G	940.4	79.8
DeiT-B	224 ²	86M	17.5G	292.3	81.8
DeiT-B	384 ²	86M	55.4G	85.9	83.1
Swin-T	224 ²	29M	4.5G	755.2	81.3
Swin-S	224 ²	50M	8.7G	436.9	83.0
Swin-B	224 ²	88M	15.4G	278.1	83.5
Swin-B	384 ²	88M	47.0G	84.7	84.5

Сегментация

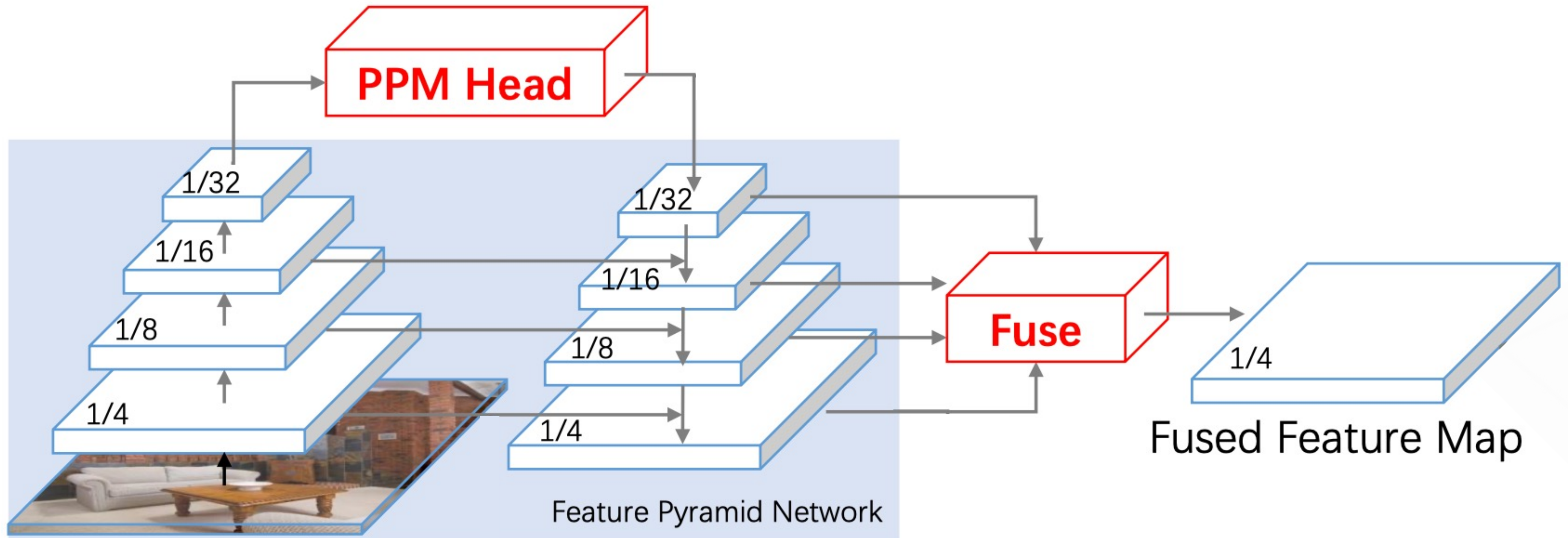


Сегментация

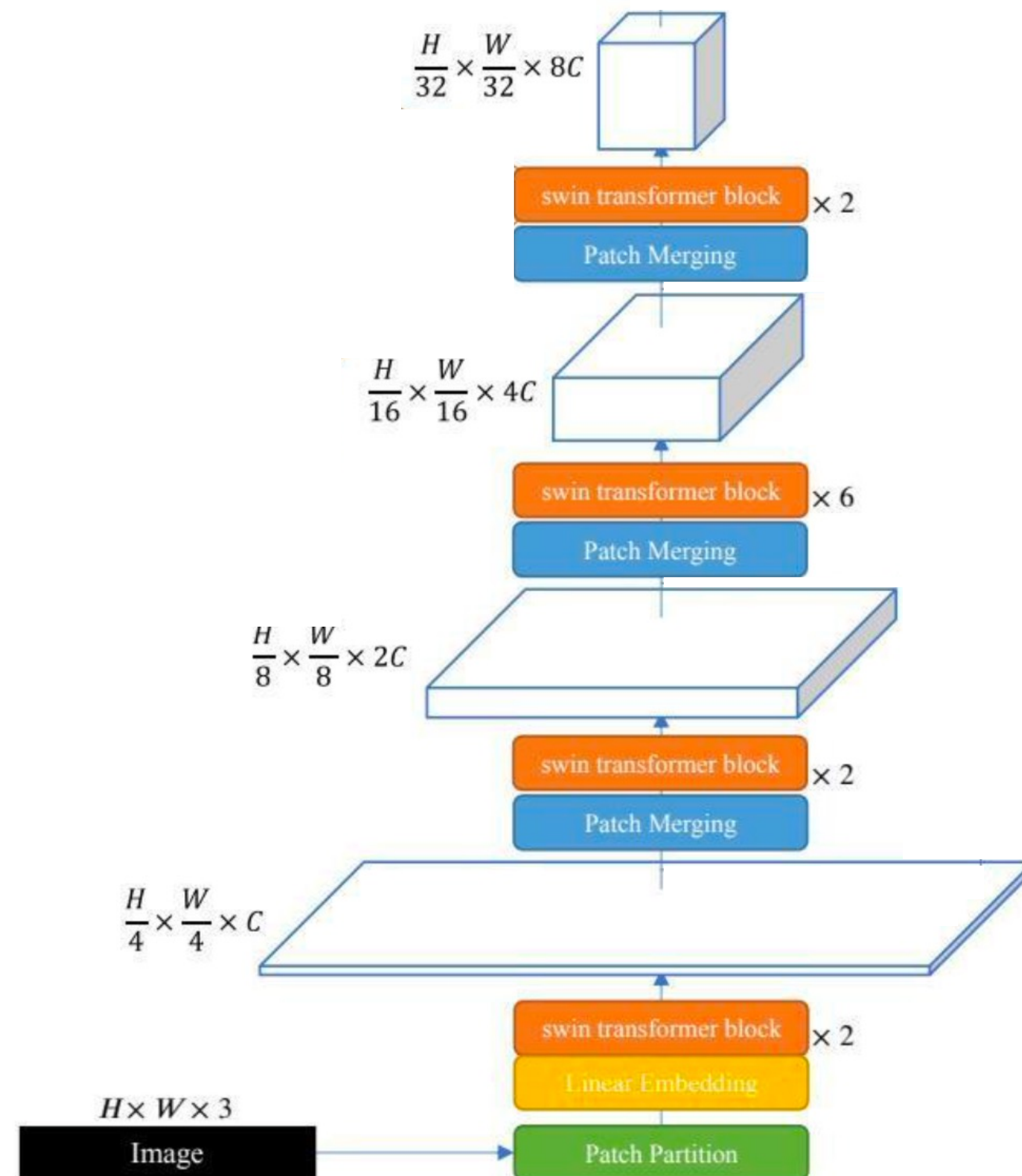


Очень напоминает сверточную сеть

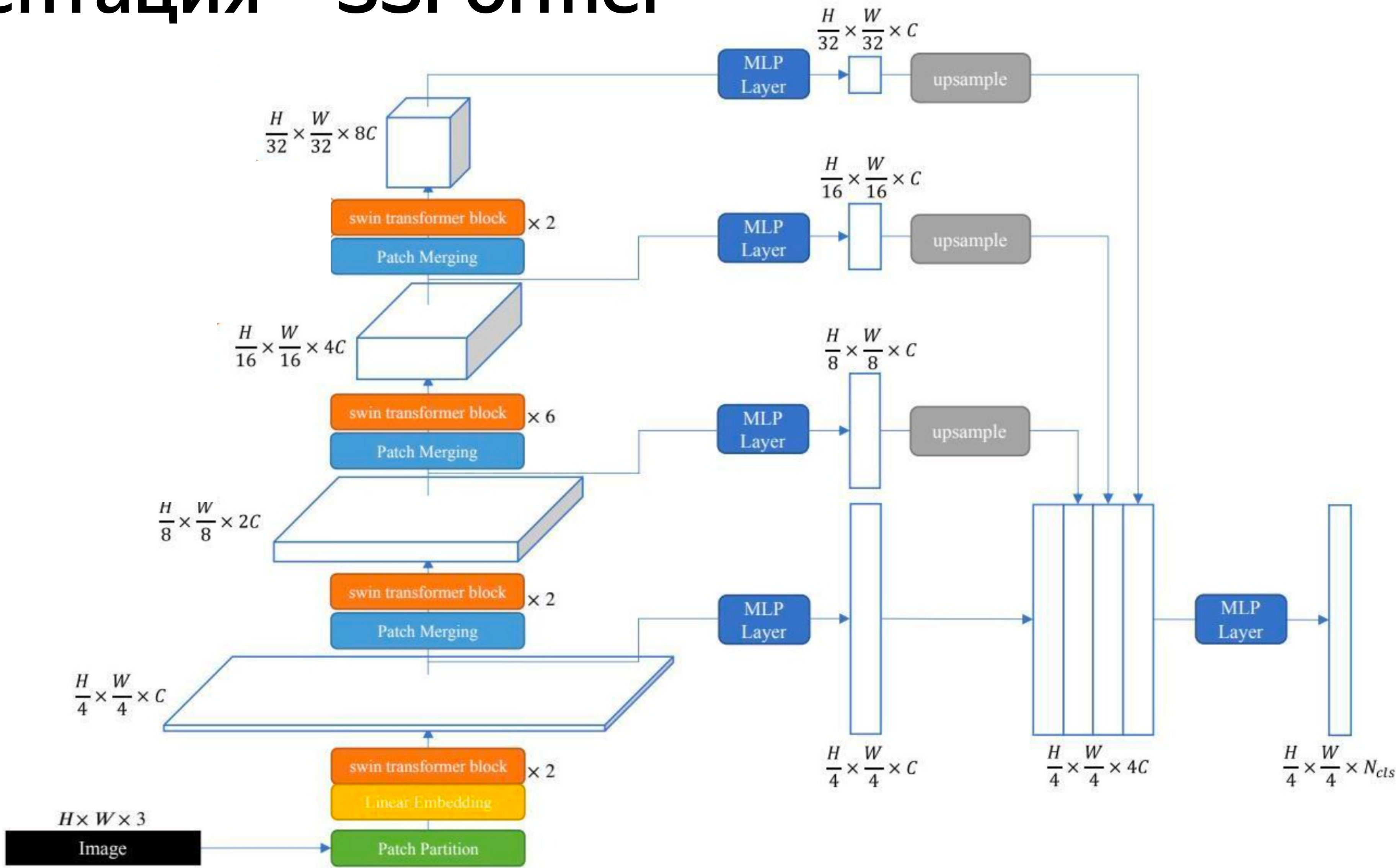
Сегментация – UperNet



Сегментация – SSFormer



Сегментация – SSFormer



Сегментация – ADE 20K

Method	Iteration	Params	Flops*	mIoU
FCN	160000	68.6M	275.7G	39.91
DeepLabV3+	160000	62.7M	255.1G	45.47
DMnet	160000	72.3M	273.6G	45.42
PSPNet	160000	68.1M	256.4G	44.39
PSANet	160000	73.1M	272.5G	43.74
SETR-PUP	160000	317.3M	362.1G	48.24
Swin-T	160000	121.3M	297.2G	50.31
SSformer	160000	87.5M	91.01G	47.71

*The Flops are calculated with the input images of 512×512 resolution.

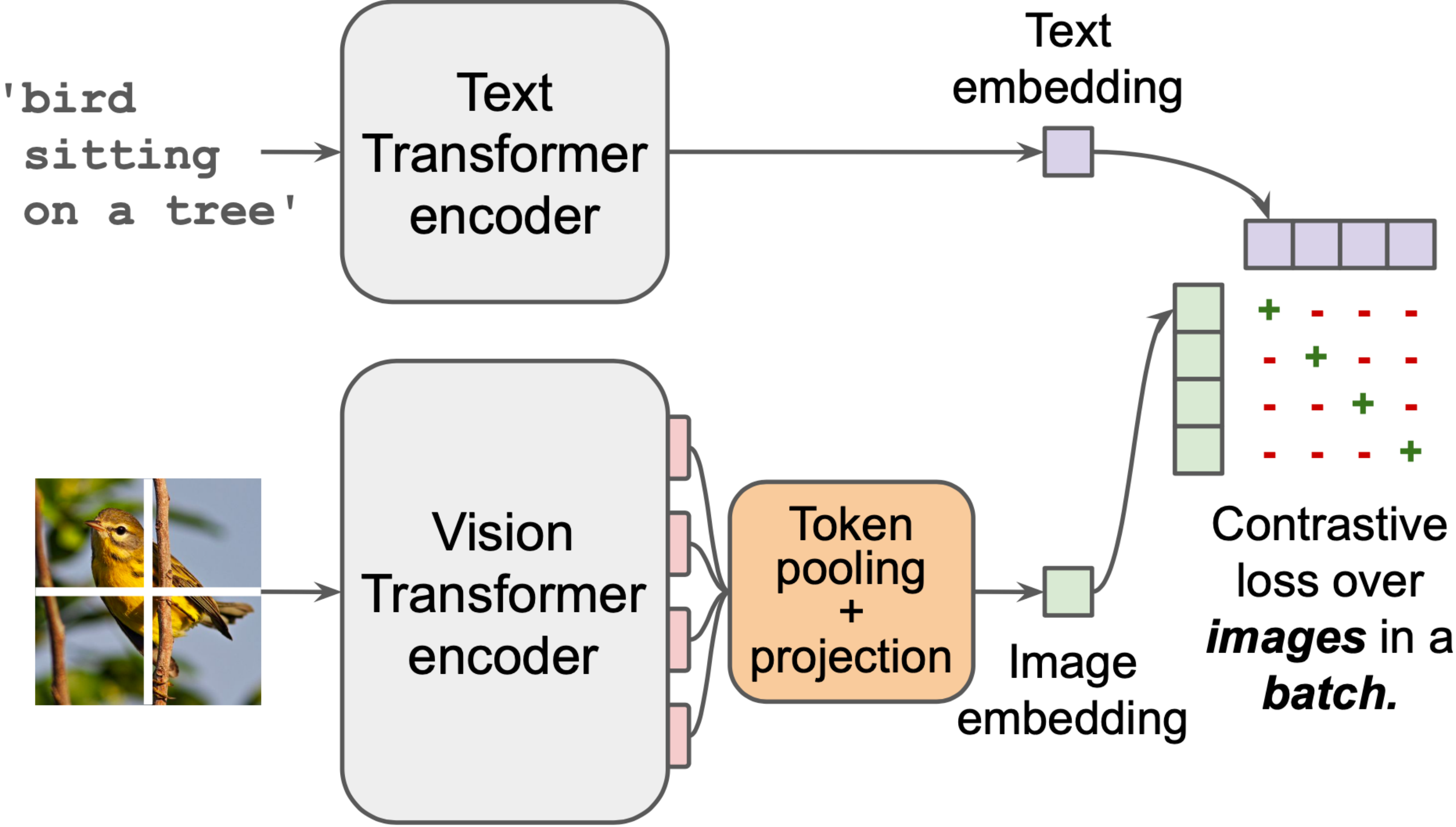
Intro

Swin Transformer

Как применять для разных задач

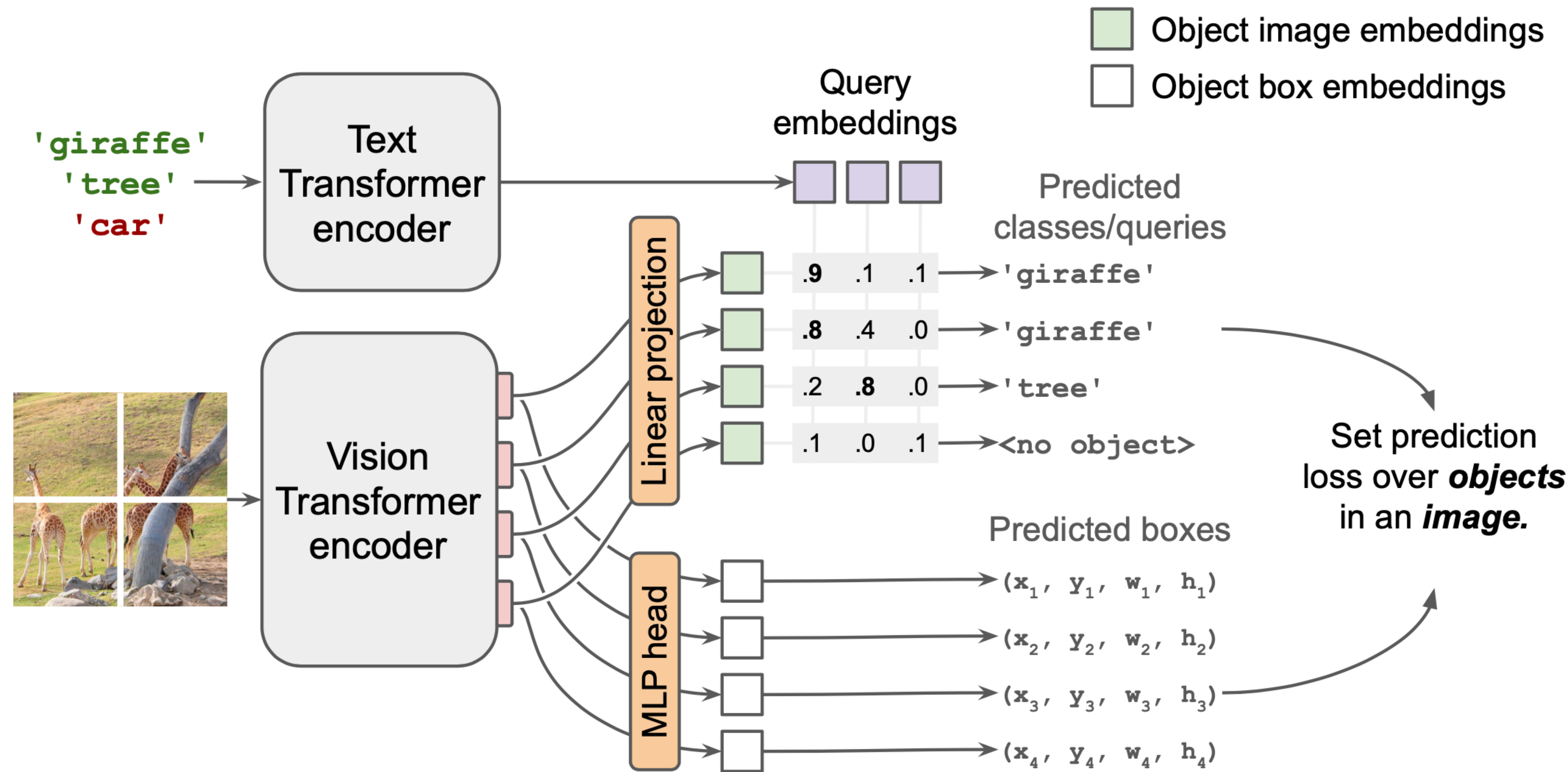
Zero Shot и One Shot детекция

OWL Detector



Учим модель
сопоставлять текст и
изображение
Аналогично CLIP

OWL Detector



OWL Detector

Можно взять encoder для изображений

