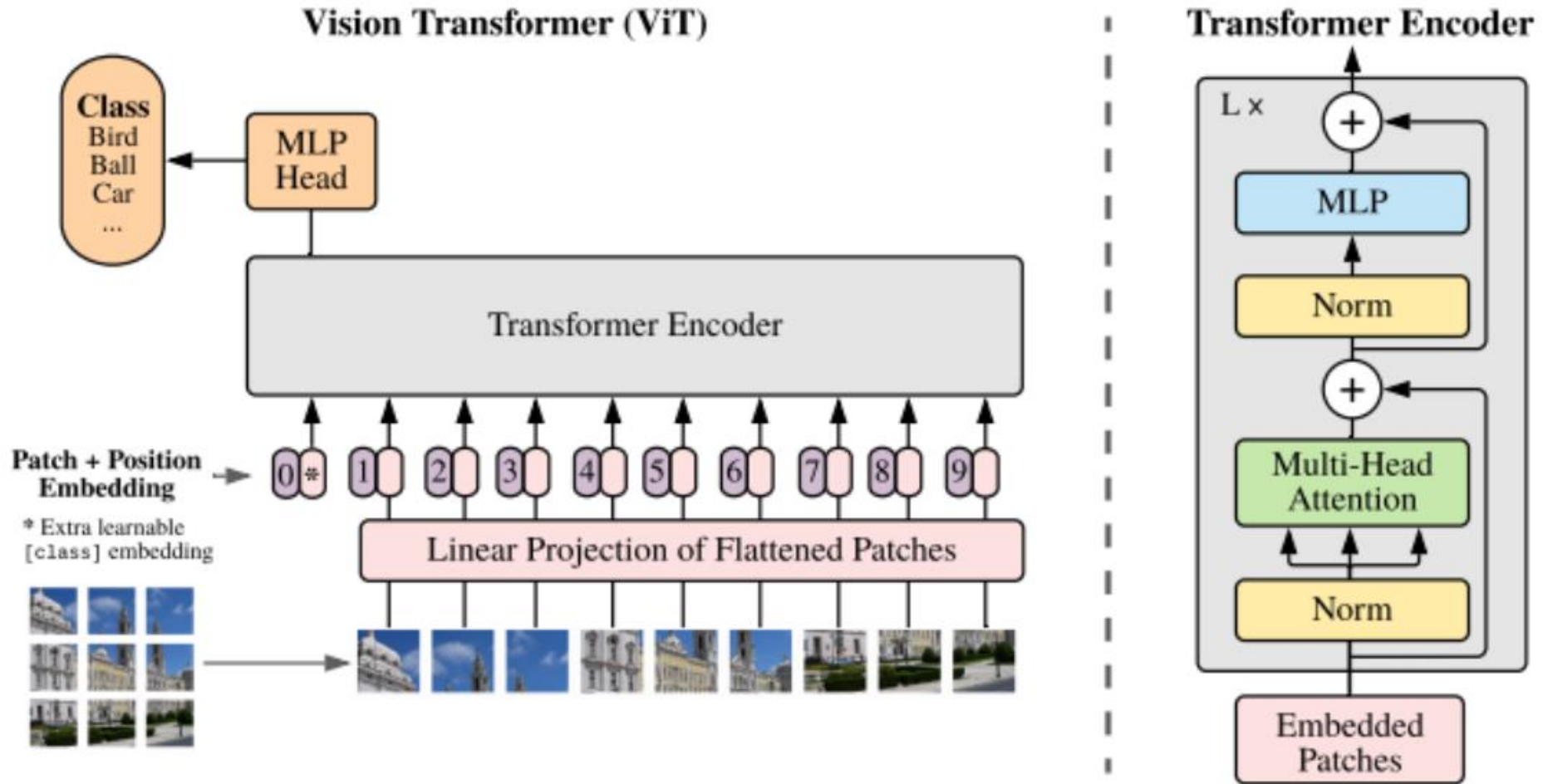


# Vision Transformer



- Motivation
- Methodology
- Result

# Motivation

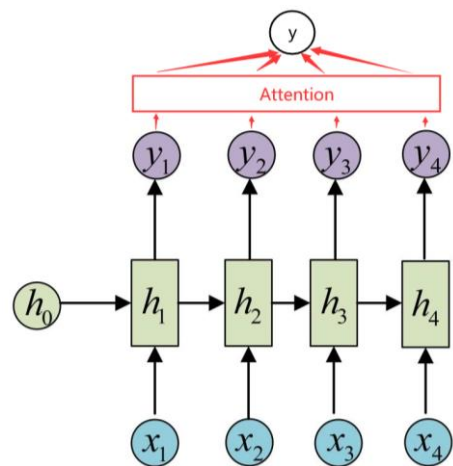
We have shown that the standard **BERT** recipe **is effective** on a wide range of model sizes

Time Sequential relation: Text

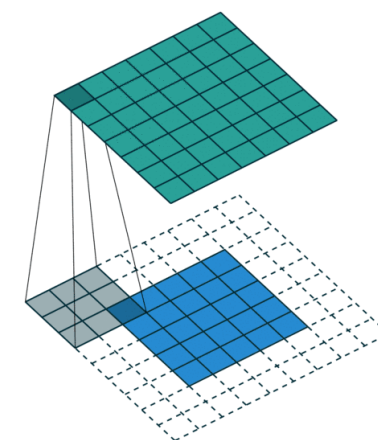


Spatial relation: Image

We have shown that the standard BERT recipe is **effective** on a wide range of model sizes



Attention, Transformer

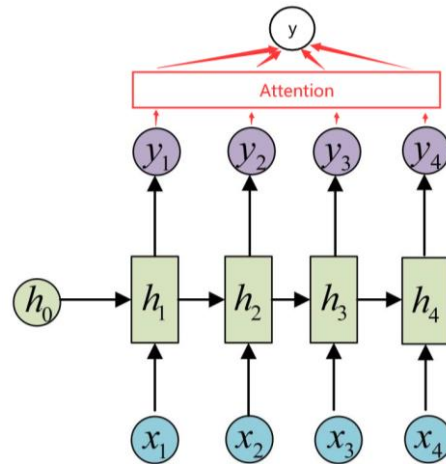


Convolution

Since Transformer yields considerable result in Time Sequential data,  
Can Transformer be applied to images?



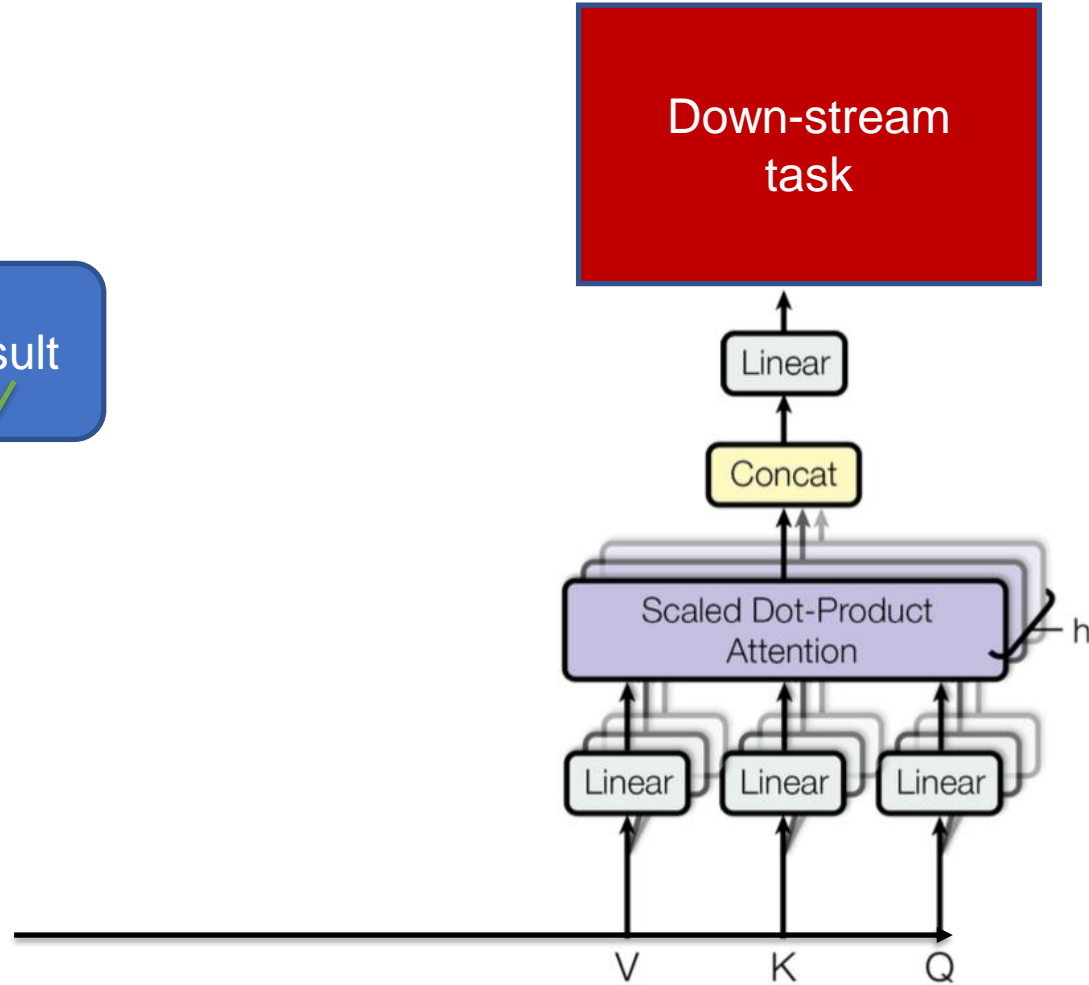
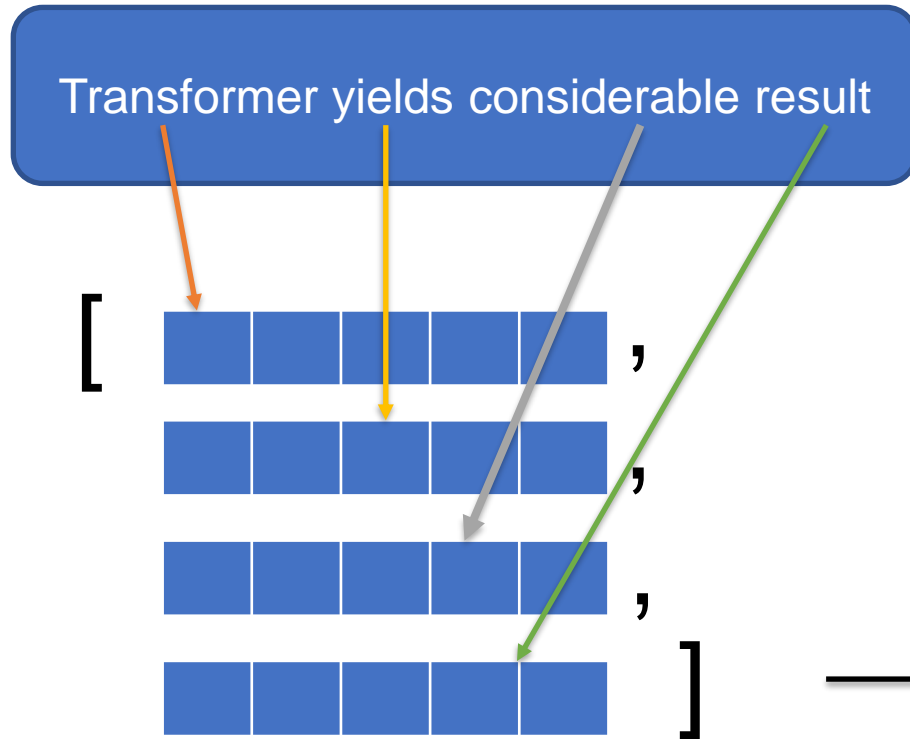
+

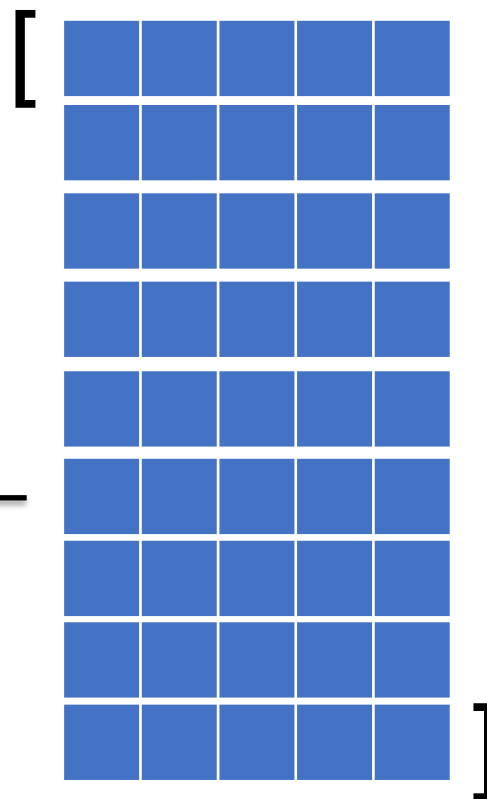


= ?

# Methodology

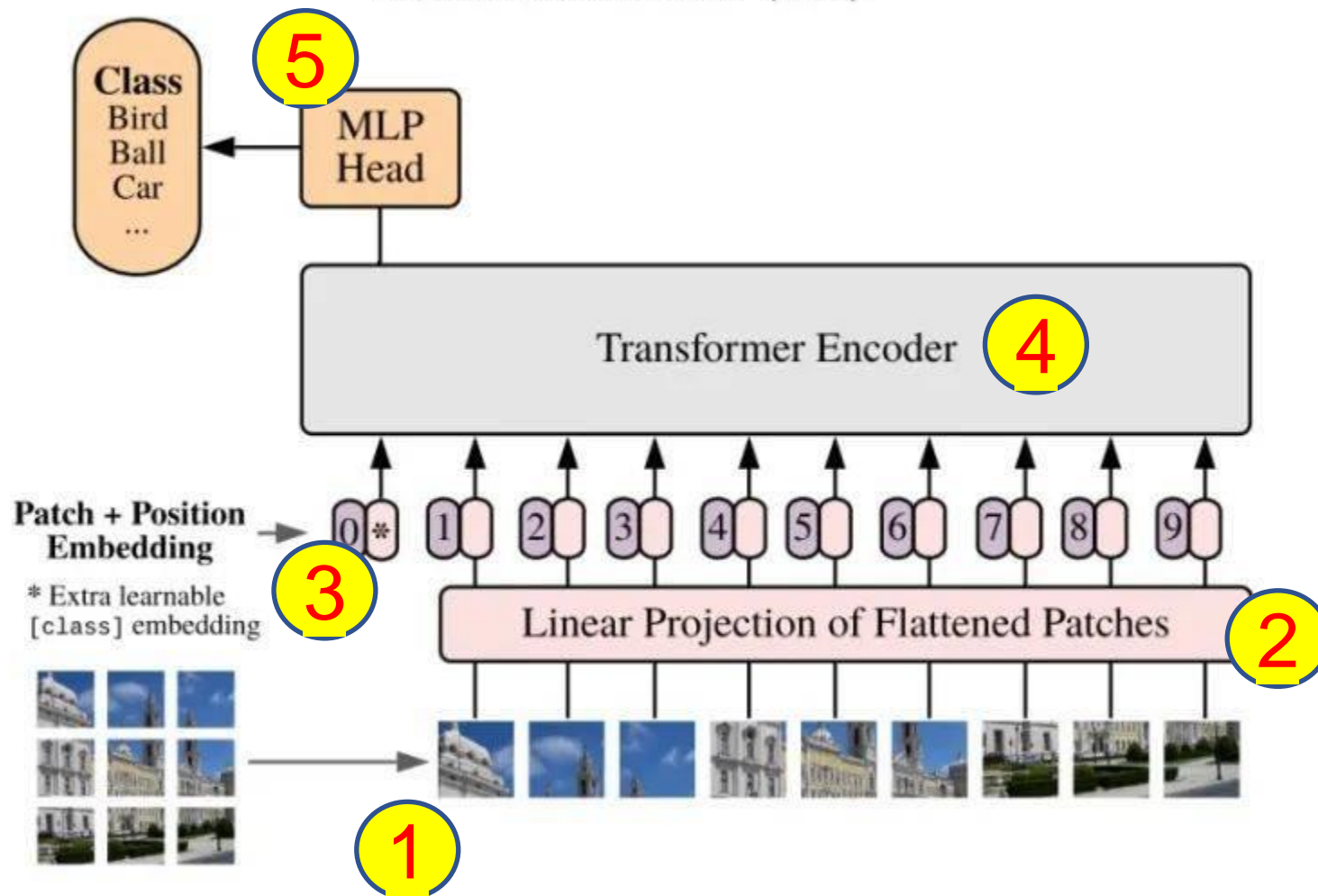
## Review: Transformer



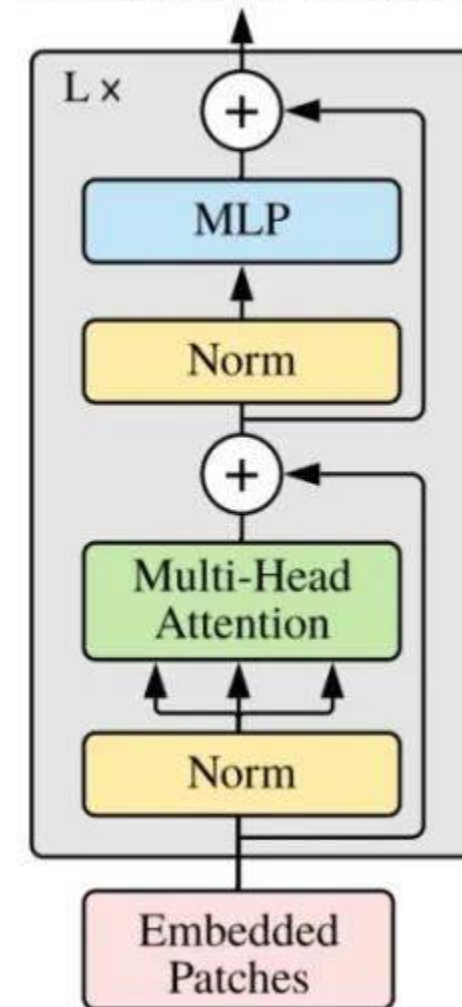


Transformer

## Vision Transformer (ViT)



## Transformer Encoder










# Result

## CIFAR100 Top1 Accuracy

1	<b>EffNet-L2</b> (SAM)	96.08	✓	<a href="#">Sharpness-Aware Minimization for Efficiently Improving Generalization</a>
2	<b>ViT-H/14</b>	94.55±0.04	✓	<a href="#">An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale</a>
3	<b>ViT-B-16</b> (ImageNet-21K-P pretrain)	94.2	✓	<a href="#">ImageNet-21K Pretraining for the Masses</a>
4	<b>CvT-W24</b>	94.09	✓	<a href="#">CvT: Introducing Convolutions to Vision Transformers</a>
5	<b>ViT-L/16</b>	93.90±0.05	✓	<a href="#">An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale</a>

# ImageNet Top1 Accuracy

1	CoAtNet-7	90.88%	2440M	✓	CoAtNet: Marrying Convolution and Attention for All Data Sizes			2021	<div>CNN</div> <div>Conv+Transformer</div> <div>JFT-3B</div>	
2	ViT-G/14	90.45%	1843M	✓	Scaling Vision Transformers			2021	<div>Transformer</div> <div>JFT-3B</div>	
3	CoAtNet-6	90.45%	1470M	✓	CoAtNet: Marrying Convolution and Attention for All Data Sizes			2021	<div>Conv+Transformer</div> <div>JFT-3B</div>	
4	ViT-MoE-15B (Every-2)	90.35%	14700M	✓	Scaling Vision with Sparse Mixture of Experts			2021	<div>Transformer</div> <div>JFT-3B</div>	
5	Meta Pseudo Labels (EfficientNet-L2)	90.2%	98.8%	480M	✓	Meta Pseudo Labels			2021	<div>EfficientNet</div> <div>JFT-300M</div>