

# An Image Is Worth 393 Areas:

Training image Transformers with Area-Attention

Osher Tidhar

Yoav Kurtz



Transformer

# Agenda

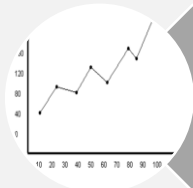


CAT

Motivation



Our Method



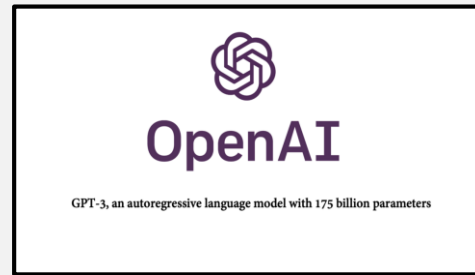
Results



Further Steps

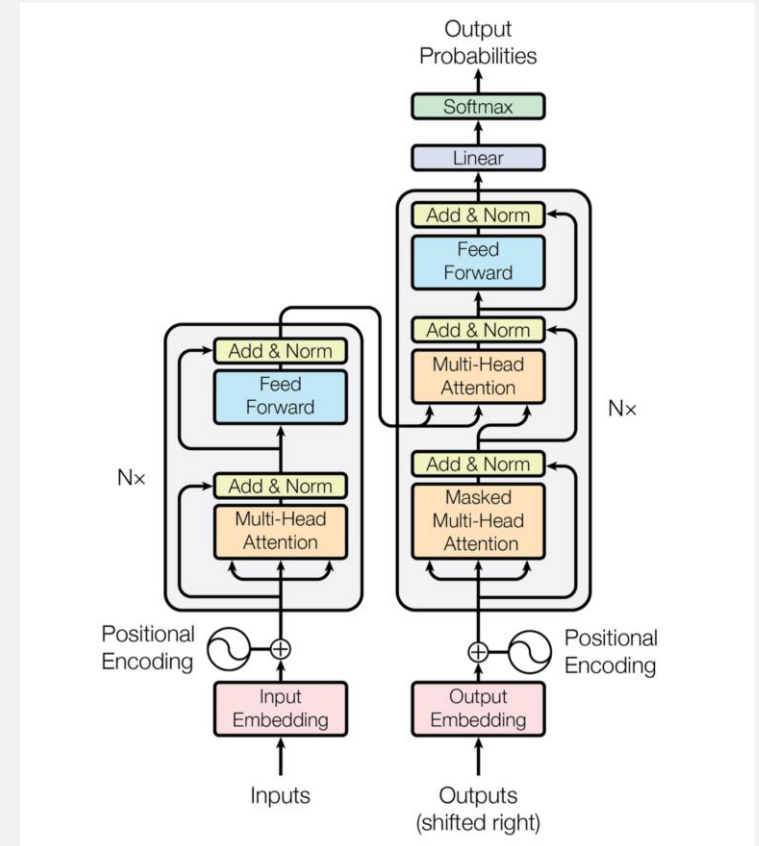
# Transformers

- Model of choice for NLP problems.



- Recently migrated to computer vision.

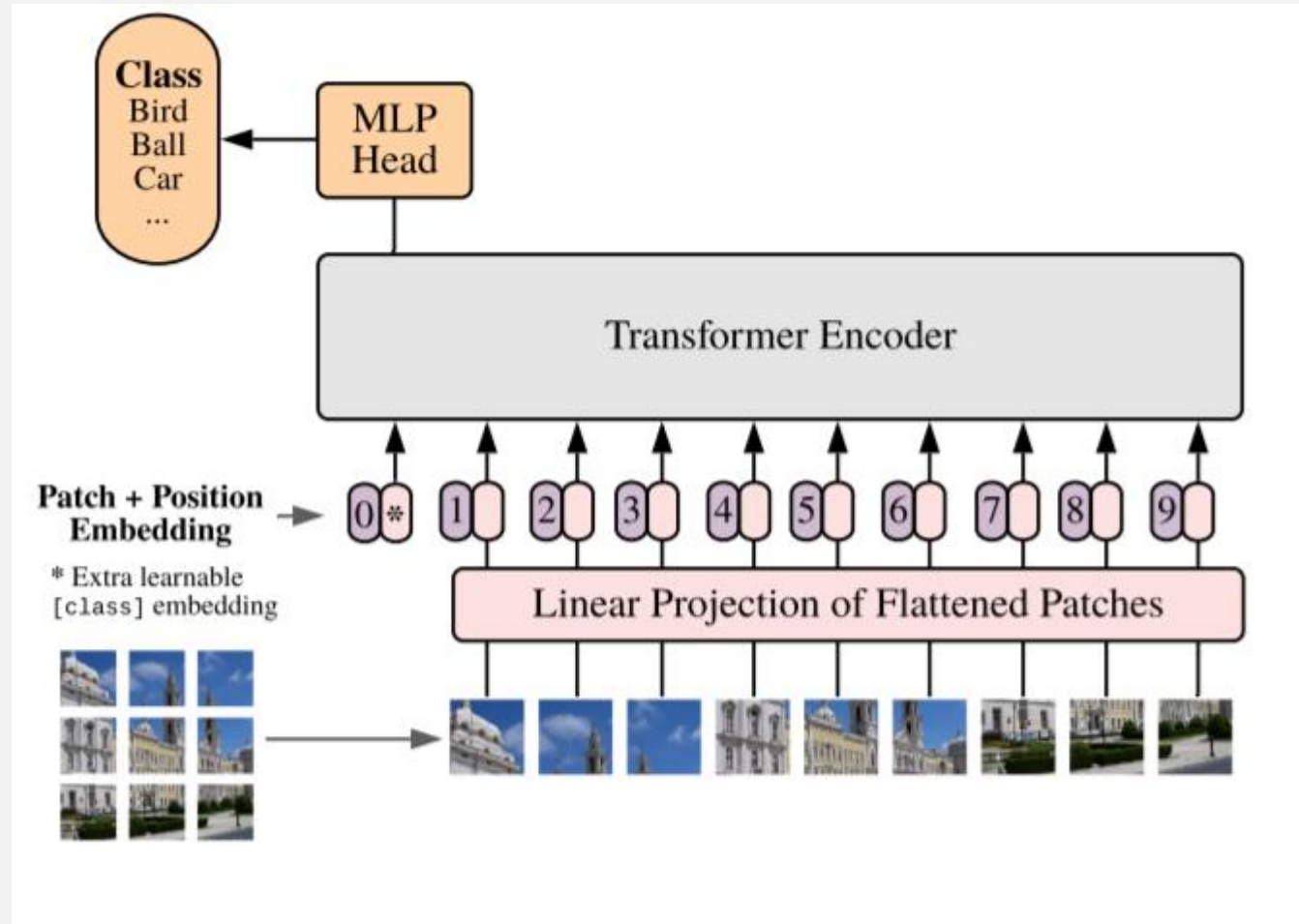
- Motivation
- Our Method
- Results
- Further Steps



Transformer Overview

# Vision Transformer<sup>1</sup>

- Motivation
- Our Method
- Results
- Further Steps

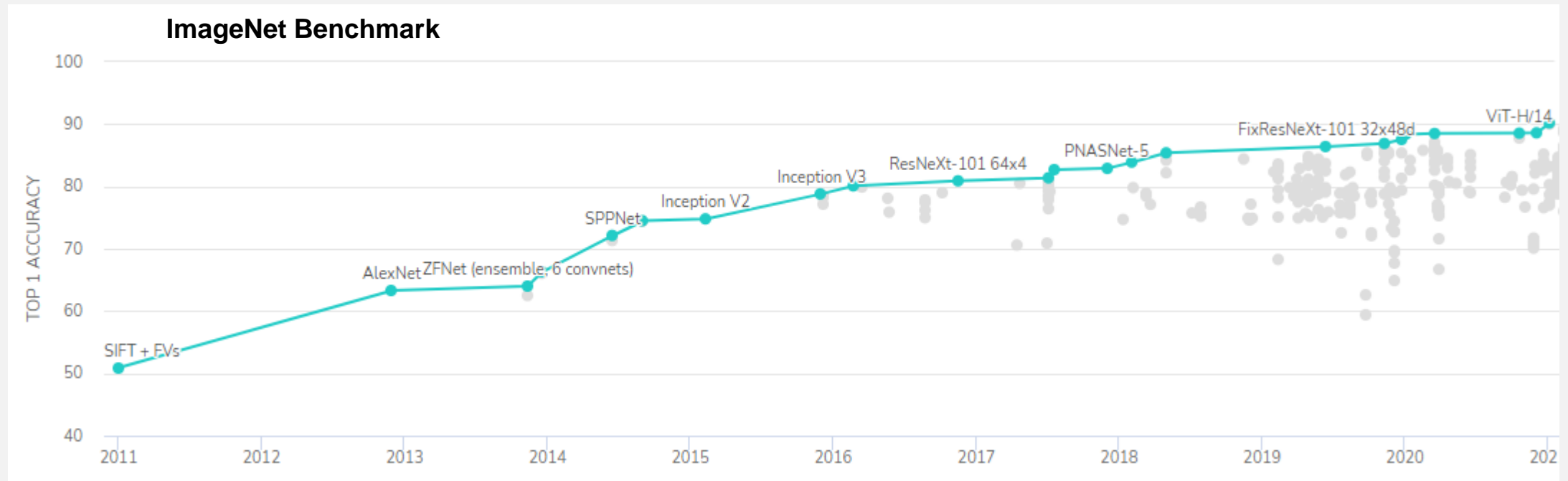


<sup>1</sup>[An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

# Transformers

- Motivation
- Our Method
- Results
- Further Steps

- Model of choice for NLP problems.
- Recently migrated to vision, showing **competitive**<sup>1,2</sup> results.



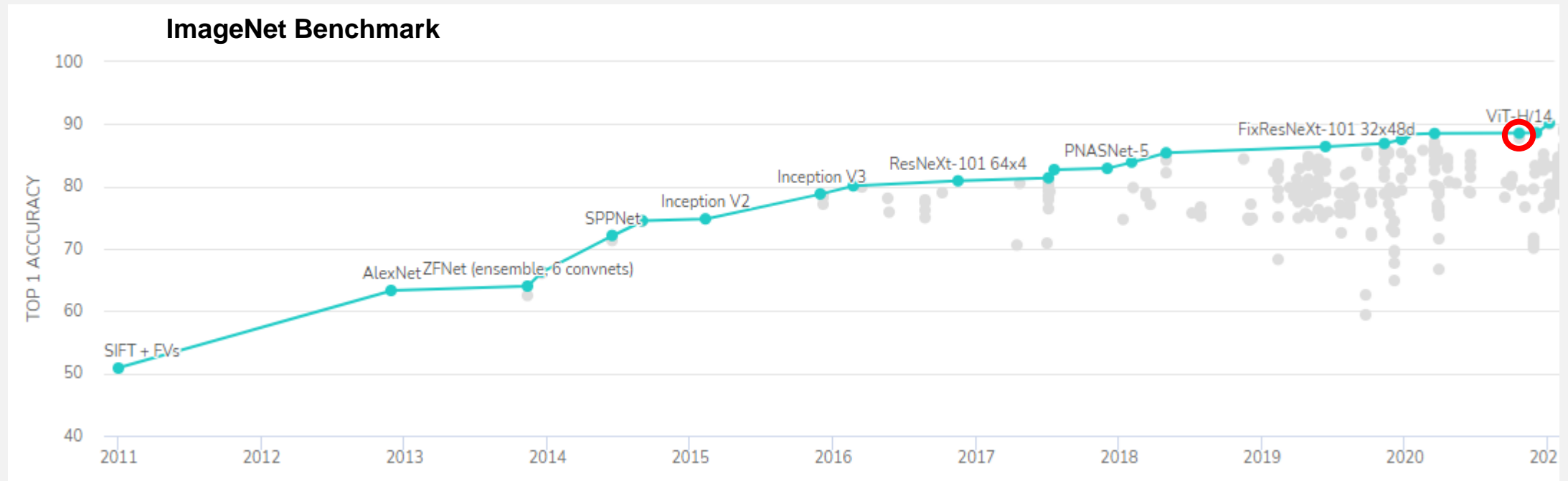
<sup>1</sup>[An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

<sup>2</sup>[Training data-efficient image transformers & distillation through attention](#)

# Transformers

- Motivation
- Our Method
- Results
- Further Steps

- Model of choice for NLP problems.
- Recently migrated to vision, showing **competitive**<sup>1,2</sup> results.



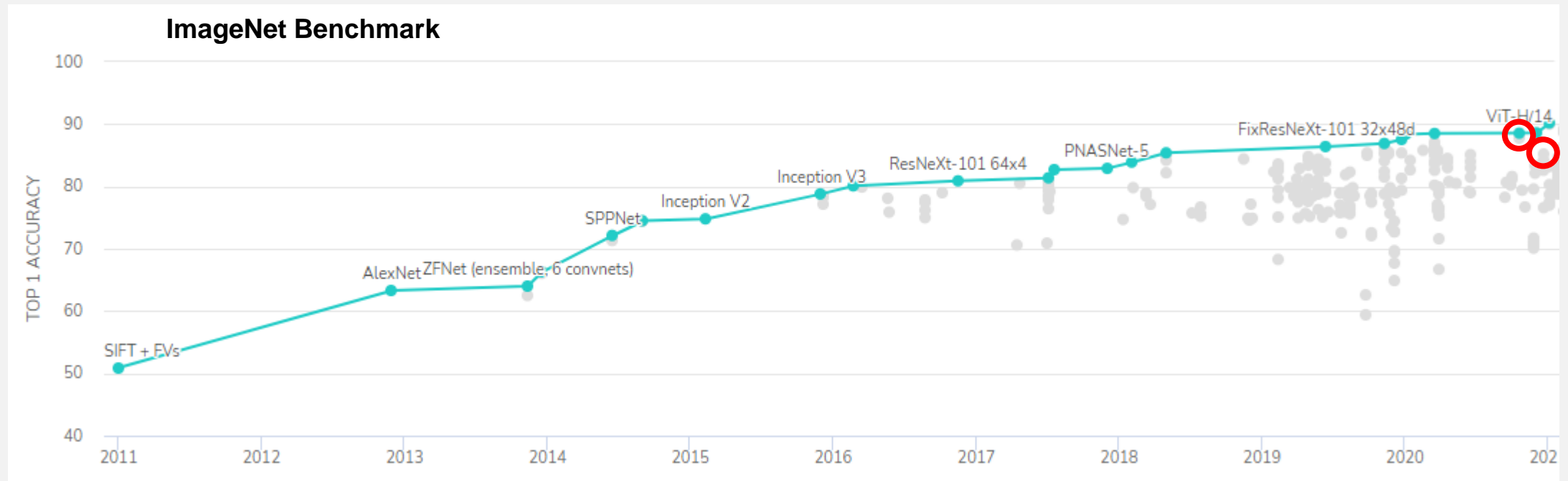
<sup>1</sup>[An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

<sup>2</sup>[Training data-efficient image transformers & distillation through attention](#)

# Transformers

- Motivation
- Our Method
- Results
- Further Steps

- Model of choice for NLP problems.
- Recently migrated to vision, showing **competitive**<sup>1,2</sup> results.



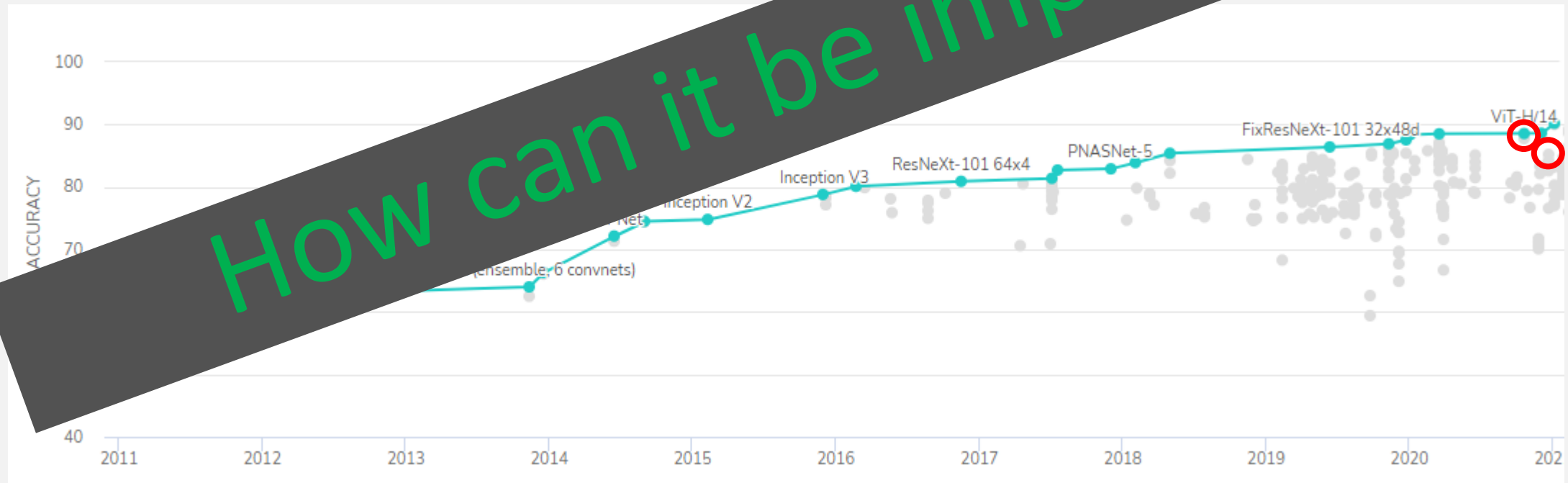
<sup>1</sup>[An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

<sup>2</sup>[Training data-efficient image transformers & distillation through attention](#)

# Transformers

- Motivation
- Our Method
- Results
- Future

- Model of choice for NLP problems.
- Recently migrated to vision, showing competitive performance



<sup>1</sup>[An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

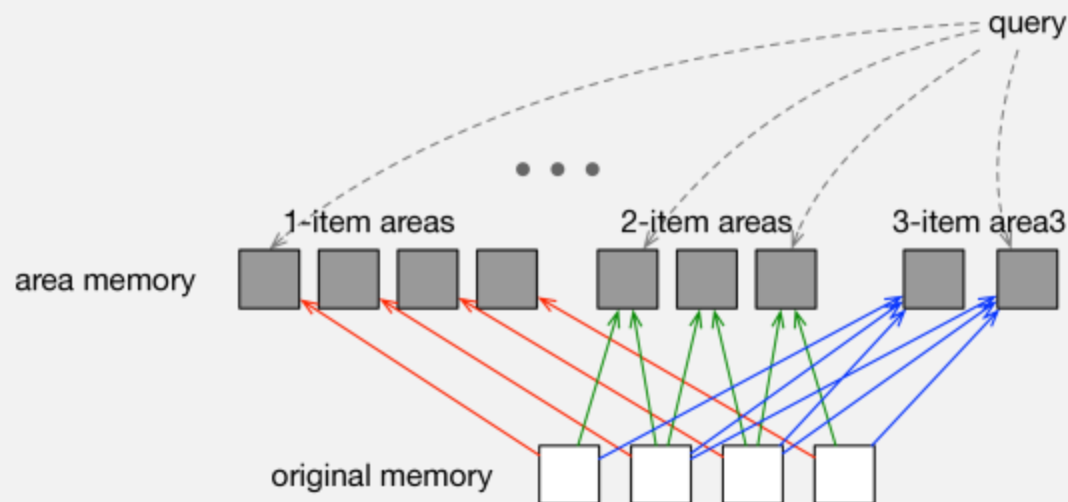
<sup>2</sup>[Training data-efficient image transformers & distillation through attention](#)



# Area-Attention<sup>1</sup>

- Motivation
- Our Method
- Results
- Further Steps

- Attending group of items in the memory that are structurally adjacent.
- Model can attend to combinations of items.



# Area-Attention<sup>1</sup>

- Motivation
- **Our Method**
- Results
- Further Steps

$Z_0$

$Z_1$

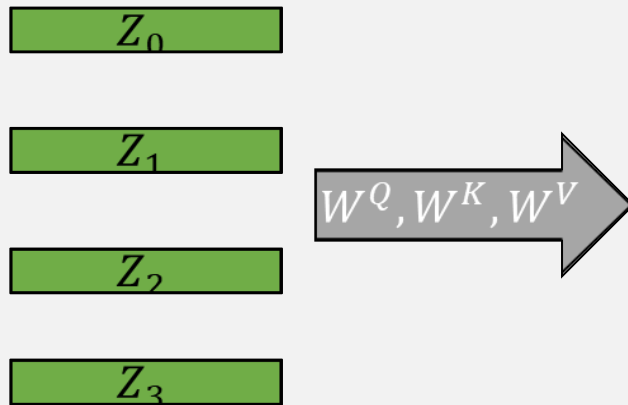
$Z_2$

$Z_3$

<sup>1</sup>[Area Attention](#)

# Area-Attention<sup>1</sup>

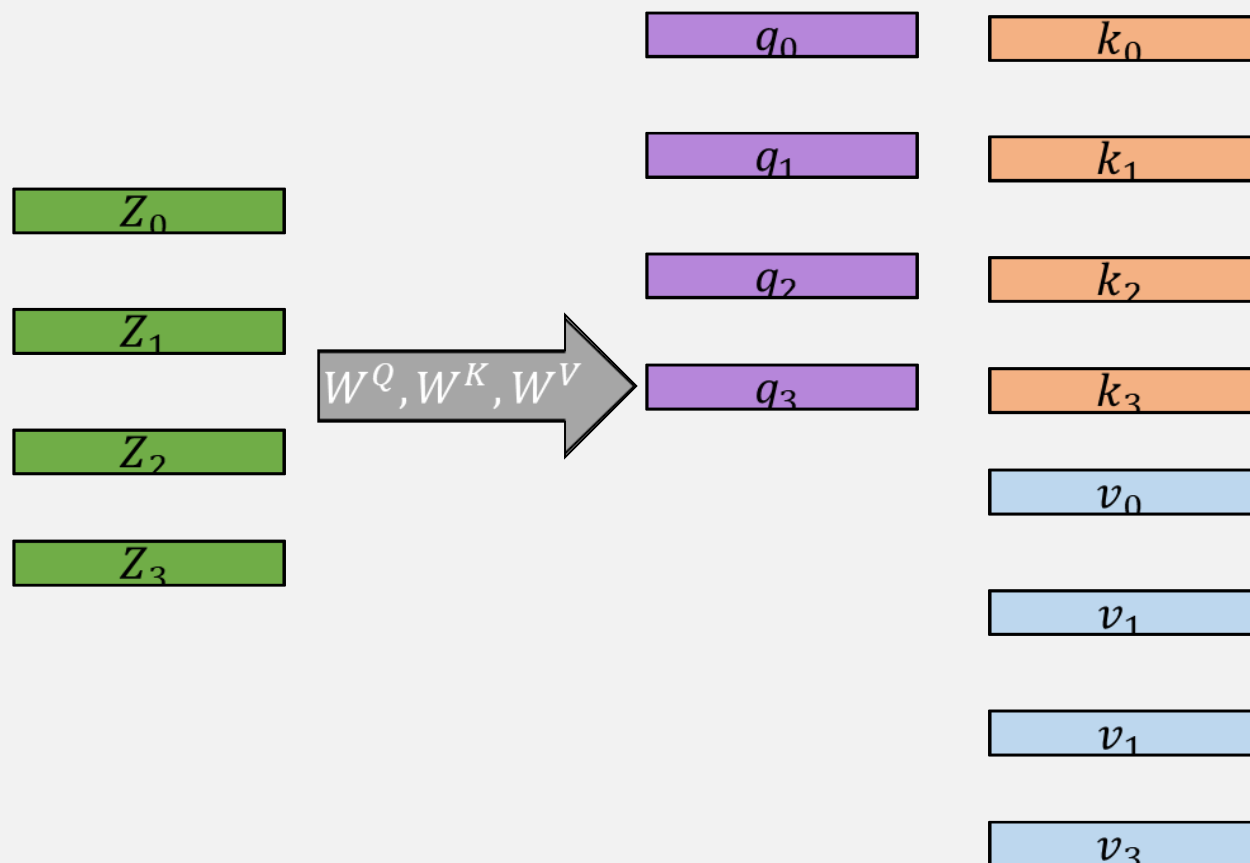
- Motivation
- **Our Method**
- Results
- Further Steps



<sup>1</sup>[Area Attention](#)

# Area-Attention<sup>1</sup>

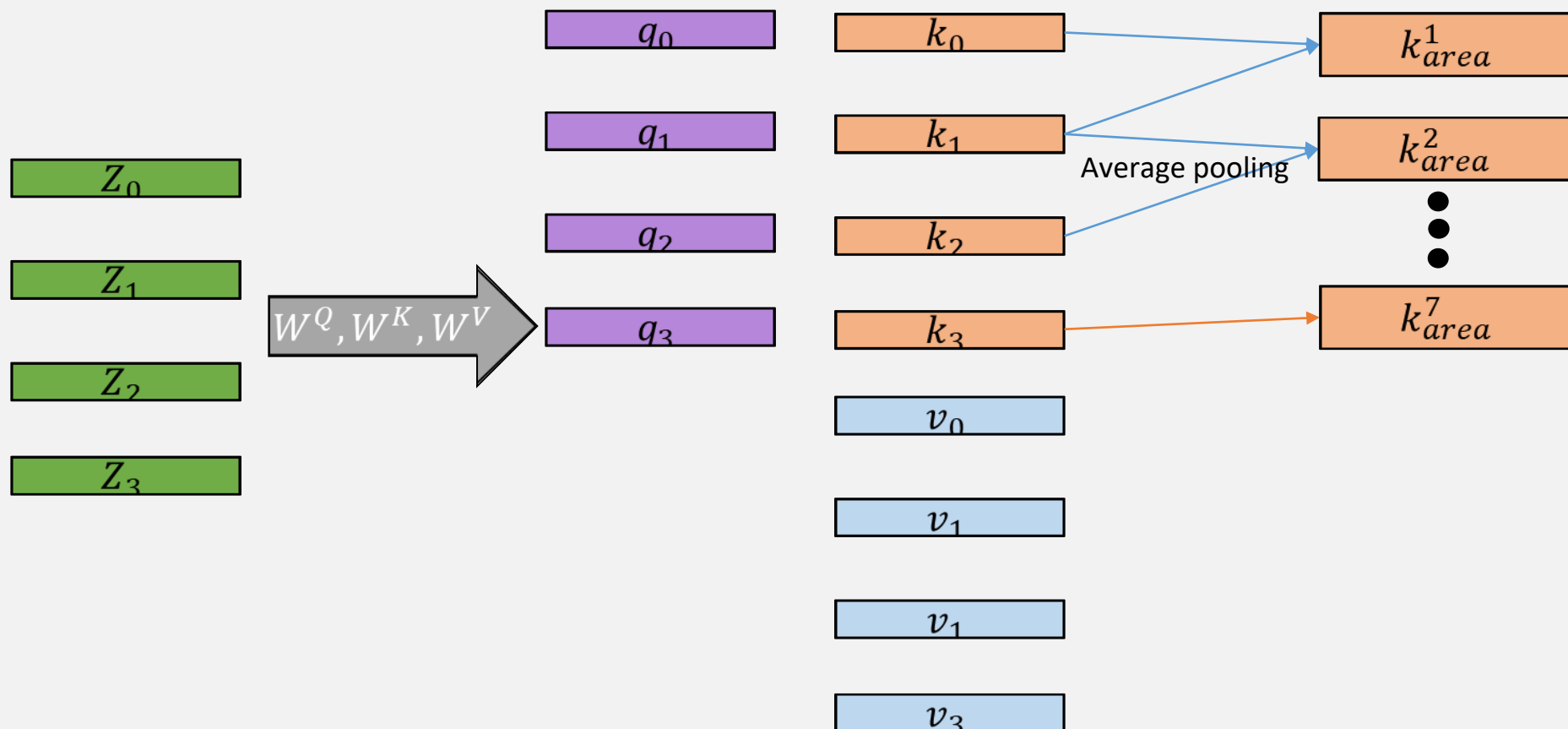
- Motivation
- Our Method
- Results
- Further Steps



<sup>1</sup>[Area Attention](#)

# Area-Attention<sup>1</sup>

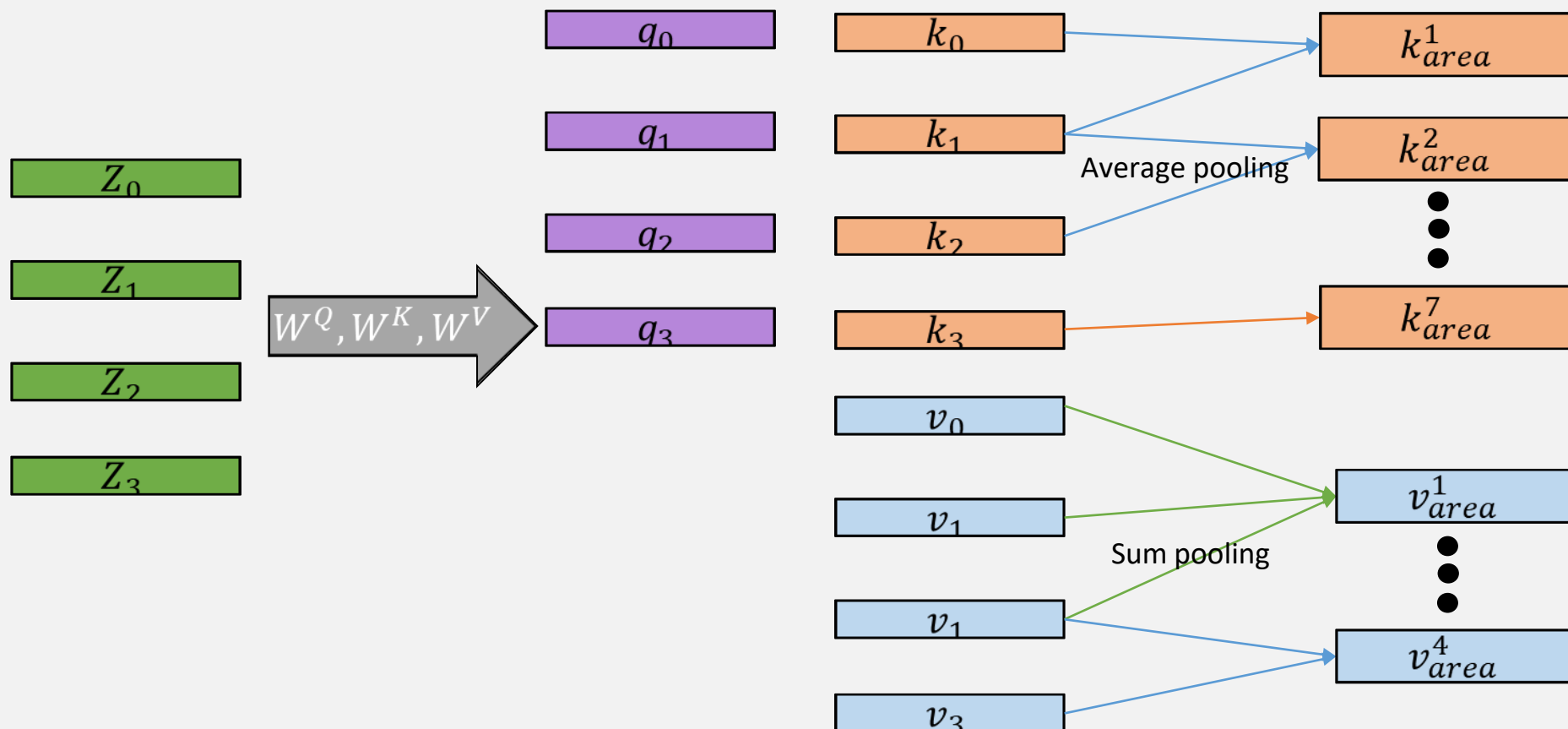
- Motivation
- Our Method
- Results
- Further Steps



<sup>1</sup>Area Attention

# Area-Attention<sup>1</sup>

- Motivation
- Our Method
- Results
- Further Steps

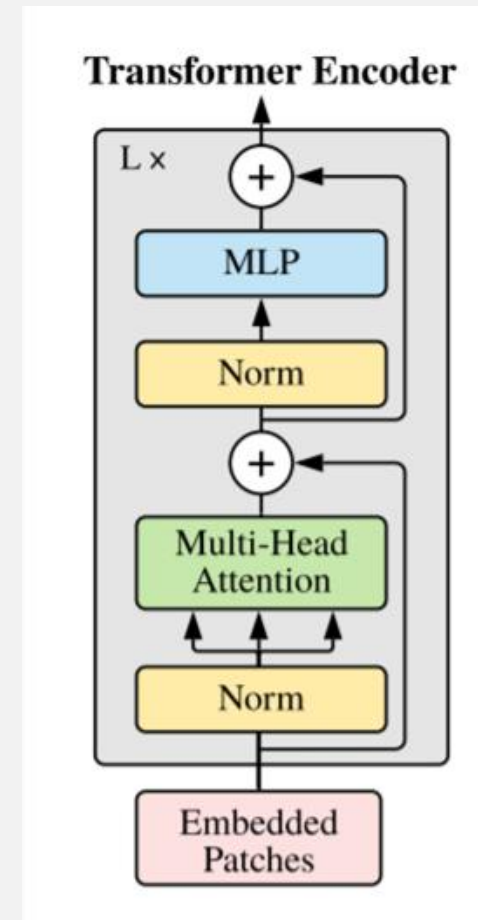


<sup>1</sup>Area Attention

# Vision Transformer + Area Attention

- Motivation
- Our Method
- Results
- Further Steps

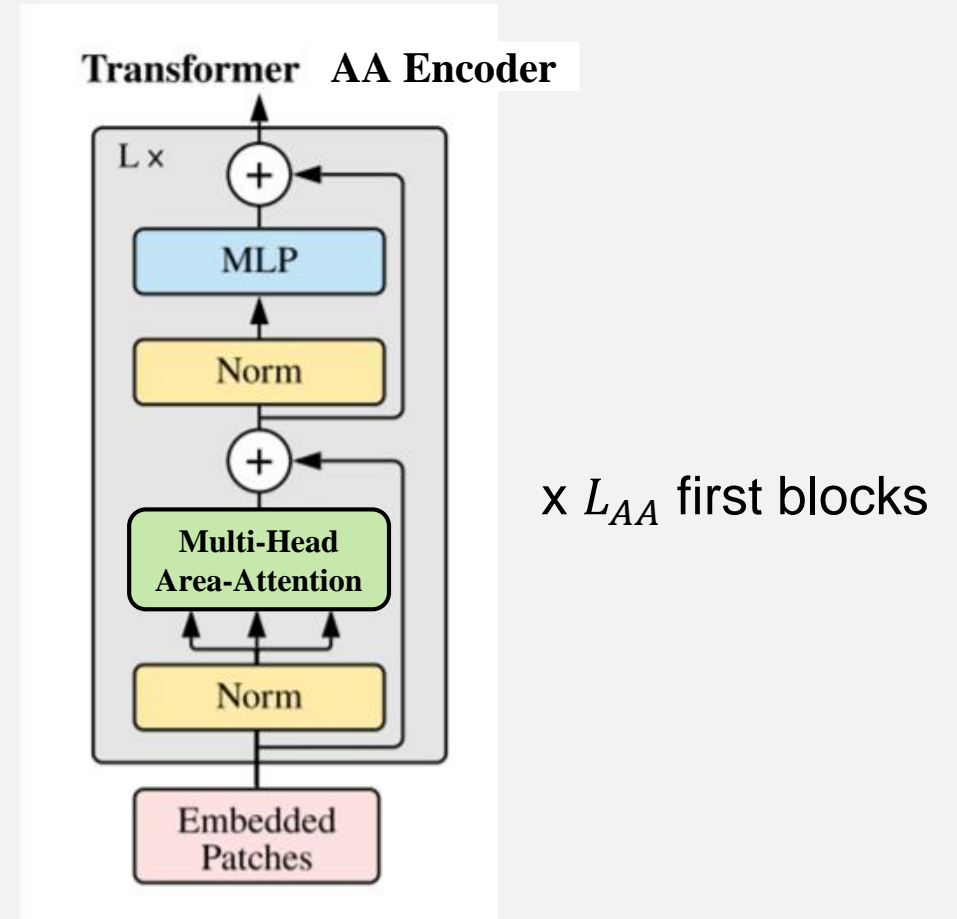
- Multi-head self-attention replaced with multi-head area-attention.
- Different AA configurations are tested.



# Vision Transformer + Area Attention

- Motivation
- Our Method
- Results
- Further Steps

- Multi-head self-attention replaced with multi-head area-attention.
- Different AA configurations are tested.





# Number of areas in ViT + AA

- Motivation
- Our Method
- Results
- Further Steps

- In ViT, each image is represented by patches of 16x16 pixels.
- In our model, what is the total number of areas that can be generated?

→ For the following configurations:

(H, W)	(P, P)	max area size
224x224	16x16	2

we got a sequence of length 197:

14x14 patch images + 1 token class.

which corresponds to 393 areas:

197 areas built of 1 element + 196 areas built of a combination of 2 adjacent elements.

# Choosing a dataset for our experiments

- Motivation
- Our Method
- Results
- Further Steps

- Dataset - CIFAR-10

	Train size	Test size	#classes
CIFAR-10	50,000	10,000	10

# Choosing models for our experiments

- Motivation
- Our Method
- Results
- Further Steps

- Original architecture of the Vision-Transformer model:

	Embedding	#heads	#layers	#params	Training resolution
ViT Base	768	12	12	86M	224

- Models we used: Vision-Transformers small<sup>1</sup> and tiny<sup>2</sup>

	Embedding	#heads	#layers	#params	Training resolution
ViT Small	384	6	12	22M	224
ViT Small + AA	384	6	12	22M	224
ViT Tiny	192	3	12	5M	224
ViT Tiny + AA	192	3	12	5M	224

# Choosing models for our experiments

- Motivation
- Our Method
- Results
- Further Steps

- Original architecture of the Vision-Transformer model:

	Embedding	#heads	#layers	#params	Training resolution
ViT Base	768	12	12	86M	224

- Models we used: Vision-Transformers small<sup>1</sup> and tiny<sup>2</sup>

	Embedding	#heads	#layers	#params	Training resolution
ViT Small	384	6	12	22M	224
ViT Small + AA	384	6	12	22M	224
ViT Tiny	192	3	12	5M	224
ViT Tiny + AA	192	3	12	5M	224

Same number of  
Parameters

# Accuracy achieved with ViT + AA

- Motivation
- Our Method
- Results
- Further Steps

- Accuracy of our pretrained weights on CIFAR-10 Testset:

- ViT-Small model:

	top-1 acc	top-5 acc	loss
<b>ViT-small + AA with max_size=2</b>	<b>91.44</b>	<b>99.6</b>	<b>0.408</b>
<b>ViT-small + AA with max_size=3</b>	<b>87.77</b>	<b>98.69</b>	<b>0.526</b>
<b>ViT-small + AA with max_size=4</b>	<b>84.61</b>	<b>98.17</b>	<b>0.602</b>
Only ViT-small	90.02	99.6	0.432

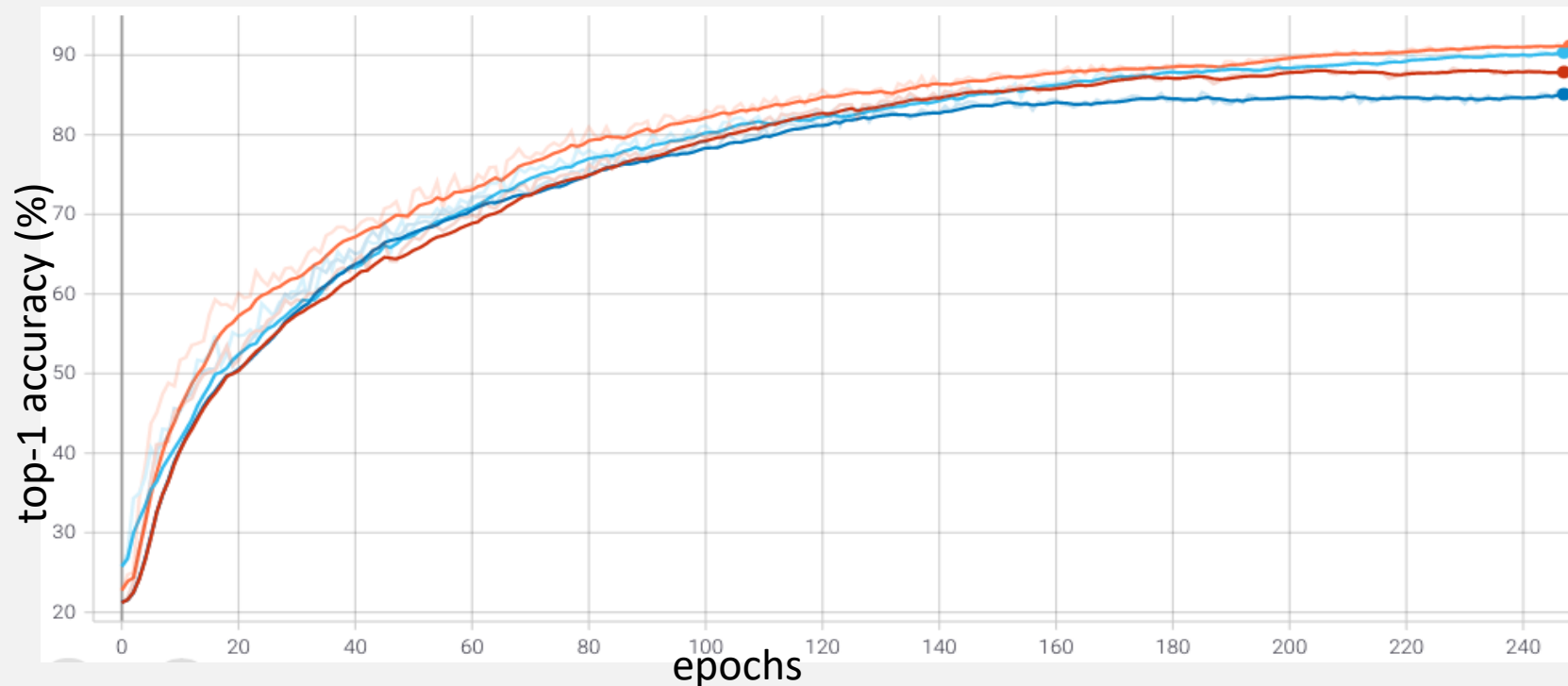
- ViT-Tiny model:

	top-1 acc	top-5 acc	loss
<b>ViT-tiny + AA with max_size=2</b>	<b>86.14</b>	<b>99.52</b>	<b>0.557</b>
Only ViT-tiny	85.49	99.49	0.576

# Accuracy achieved with ViT + AA

- Motivation
- Our Method
- Results
- Further Steps

- Only AA2 configuration achieves better accuracy than only ViT-S.
- Accuracy improves as we decrease the maximum size of an area (for the first 2 blocks).

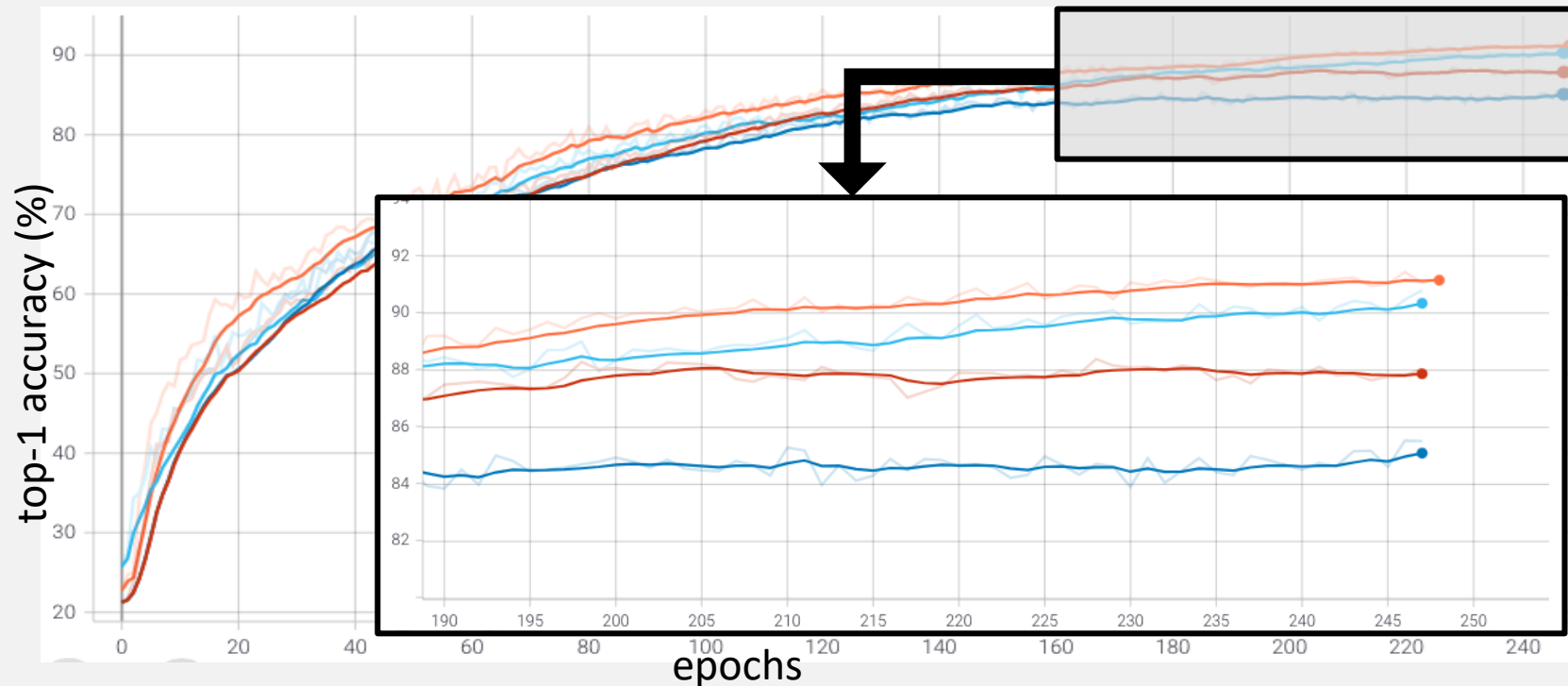


- ViT-S + AA with max size = 2
- Only ViT-S
- ViT-S + AA with max size = 3
- ViT-S + AA with max size = 4

# Accuracy achieved with ViT + AA

- Motivation
- Our Method
- Results
- Further Steps

- Only AA2 configuration achieves better accuracy than only ViT-S.
- Accuracy improves as we decrease the maximum size of an area (for the first 2 blocks).

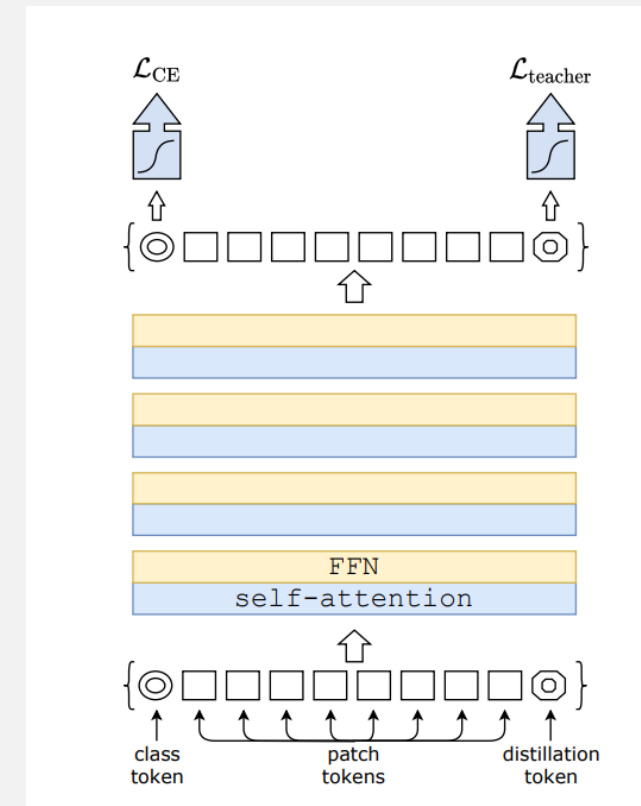


- ViT-S + AA with max size = 2
- Only ViT-S
- ViT-S + AA with max size = 3
- ViT-S + AA with max size = 4

# Further steps

- Motivation
- Our Method
- Results
- Further Steps

- Training AA-ViT using network distillation.
  - Plays the same role as the class token, except that it aims at reproducing the label estimated by the teacher.
  - Both tokens interact in the transformer through attention
  - Achieved results that were competitive with the results of convnets for Imagenet.





Questions?