

Transformers in Computer Vision

Sergey Zagoruyko

October 7, 2024

Course contents:

1. Transformers in Computer Vision
2. Transformers in Object Detection
3. Transformers in Autonomous Driving

What is convolution?

Convolution Example

Input:

1	4	-1	0	2	-2	1	3	3	1
---	---	----	---	---	----	---	---	---	---

Filter:

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

Output:

--	--	--	--	--	--	--

Source: François Fleuret's deep learning course

Convolution Example

Input:

1	4	-1	0	2	-2	1	3	3	1
---	---	----	---	---	----	---	---	---	---

Filter:

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

Output:

9						
---	--	--	--	--	--	--

Source: François Fleuret's deep learning course

Convolution Example

Input:

1	4	-1	0	2	-2	1	3	3	1
---	---	----	---	---	----	---	---	---	---

Filter:

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

Output:

9	0					
---	---	--	--	--	--	--

Source: François Fleuret's deep learning course

Convolution Example

Input:

1	4	-1	0	2	-2	1	3	3	1
---	---	----	---	---	----	---	---	---	---

Filter:

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

Output:

9	0	1				
---	---	---	--	--	--	--

Source: François Fleuret's deep learning course

Convolution Example

Input:

1	4	-1	0	2	-2	1	3	3	1
---	---	----	---	---	----	---	---	---	---

Filter:

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

Output:

9	0	1	3			
---	---	---	---	--	--	--

Source: François Fleuret's deep learning course

Convolution Example

Input:

1	4	-1	0	2	-2	1	3	3	1
---	---	----	---	---	----	---	---	---	---

Filter:

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

Output:

9	0	1	3	-5		
---	---	---	---	----	--	--

Source: François Fleuret's deep learning course

Convolution Example

Input:

1	4	-1	0	2	-2	1	3	3	1
---	---	----	---	---	----	---	---	---	---

Filter:

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

Output:

9	0	1	3	-5	-3	
---	---	---	---	----	----	--

Source: François Fleuret's deep learning course

Convolutions in 2D

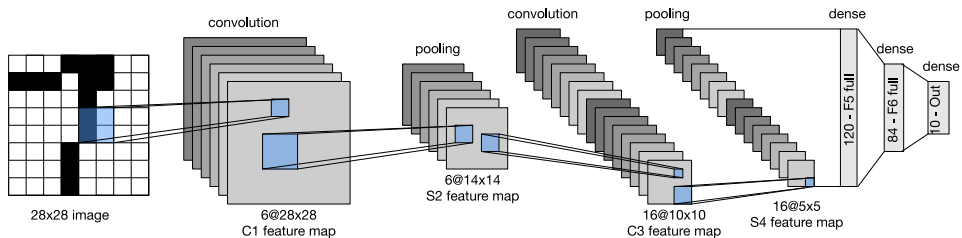
1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

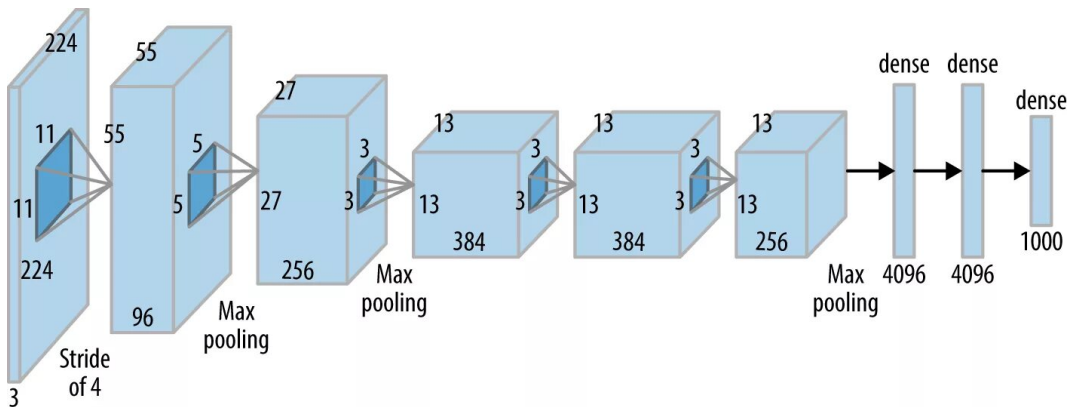
LeNet



LeNet-5 in 1998

<https://en.wikipedia.org/wiki/LeNet>

AlexNet



AlexNet in 2012

<https://machinelearningmastery.ru/the-w3h-of-alexnet-vggnet-resnet-and-inception-7baaaecccc96/>

Inductive Bias in Neural Networks

- ▶ **Definition:**

- ▶ Inductive bias refers to the set of assumptions a learning algorithm makes to generalize from the training data to unseen data.

Inductive Bias in Neural Networks

- ▶ **Definition:**

- ▶ Inductive bias refers to the set of assumptions a learning algorithm makes to generalize from the training data to unseen data.

- ▶ **Importance:**

- ▶ Helps in guiding the learning process and making predictions.
 - ▶ Determines the types of patterns a model can learn.

Inductive Bias in Neural Networks

- ▶ **Definition:**

- ▶ Inductive bias refers to the set of assumptions a learning algorithm makes to generalize from the training data to unseen data.

- ▶ **Importance:**

- ▶ Helps in guiding the learning process and making predictions.
- ▶ Determines the types of patterns a model can learn.

- ▶ **Examples in Neural Networks:**

- ▶ *Convolutional Neural Networks (CNNs)*: Assumption of spatial hierarchies in images.
- ▶ *Recurrent Neural Networks (RNNs)*: Assumption of sequential dependencies in time-series data.

Inductive Bias in Neural Networks

► Definition:

- Inductive bias refers to the set of assumptions a learning algorithm makes to generalize from the training data to unseen data.

► Importance:

- Helps in guiding the learning process and making predictions.
- Determines the types of patterns a model can learn.

► Examples in Neural Networks:

- *Convolutional Neural Networks (CNNs)*: Assumption of spatial hierarchies in images.
- *Recurrent Neural Networks (RNNs)*: Assumption of sequential dependencies in time-series data.

► Trade-offs:

- Stronger biases can lead to faster learning but may reduce flexibility.
- Weaker biases increase flexibility but may require more data.

Two pillars of deep learning:

Two pillars of deep learning:

- ▶ Large amounts of rich diverse data

Two pillars of deep learning:

- ▶ Large amounts of rich diverse data
- ▶ Large amounts of compute

AlexNet

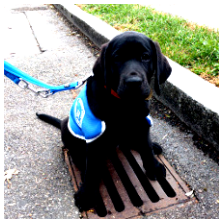
Visualizing parameters



Francois Fleuret's Deep Learning Course

AlexNet

Visualizing activations



Francois Fleuret's Deep Learning Course

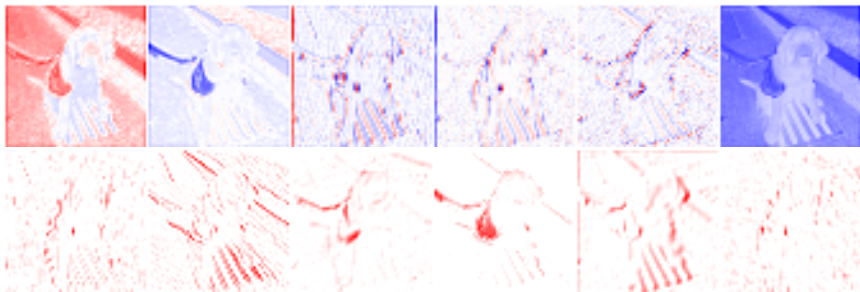
AlexNet

Visualizing activations



AlexNet

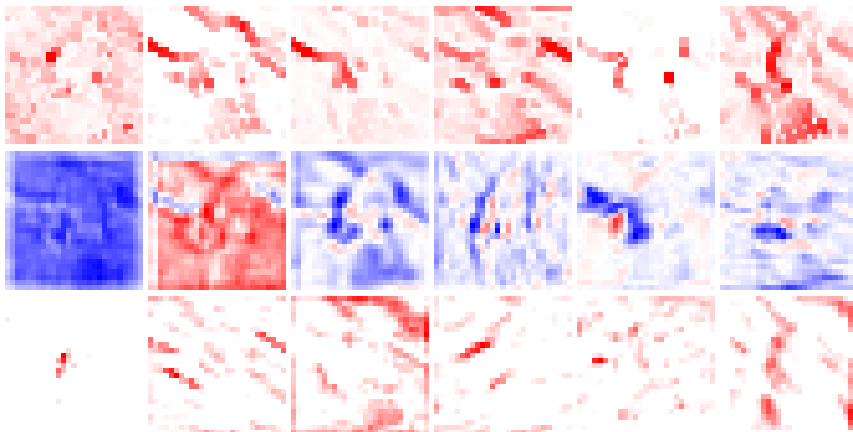
Visualizing activations



Francois Fleuret's Deep Learning Course

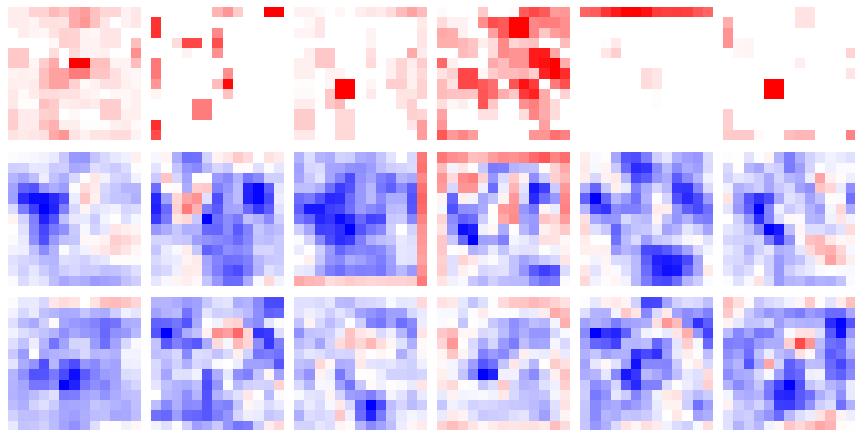
AlexNet

Visualizing activations



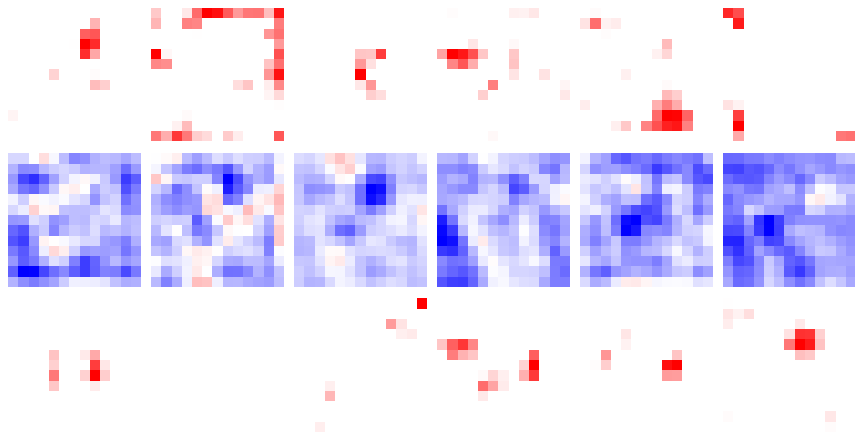
AlexNet

Visualizing activations

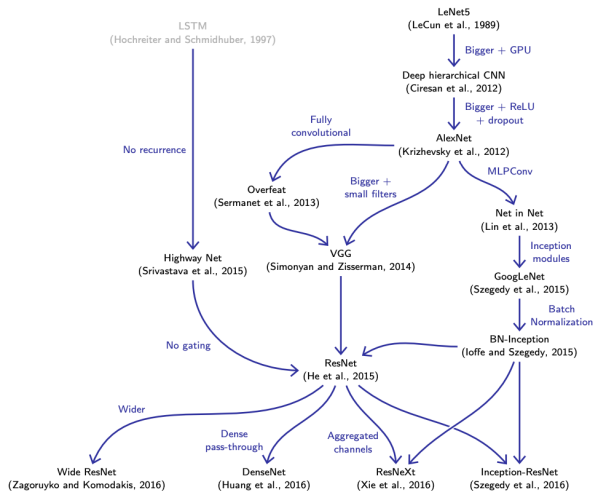


AlexNet

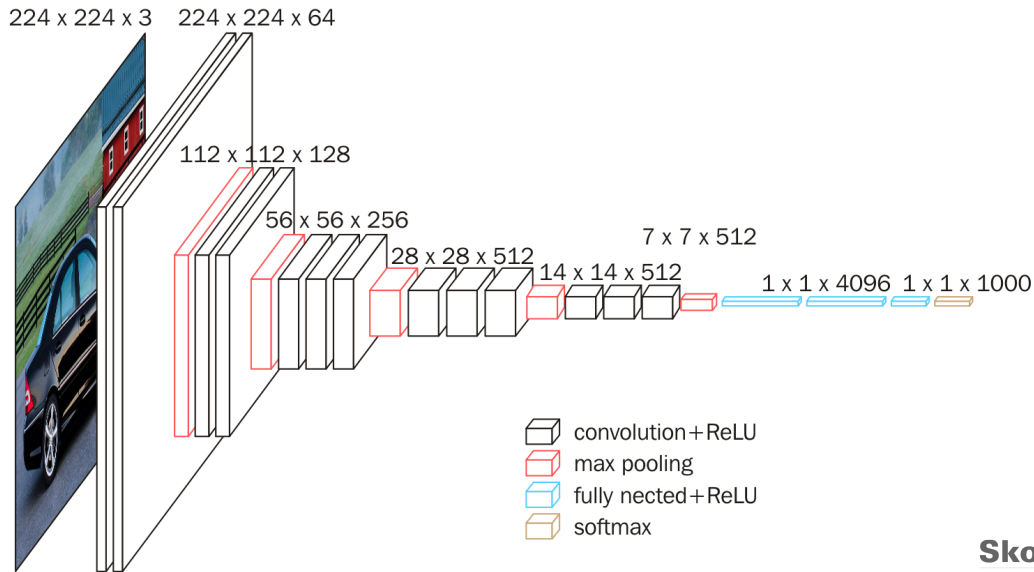
Visualizing activations



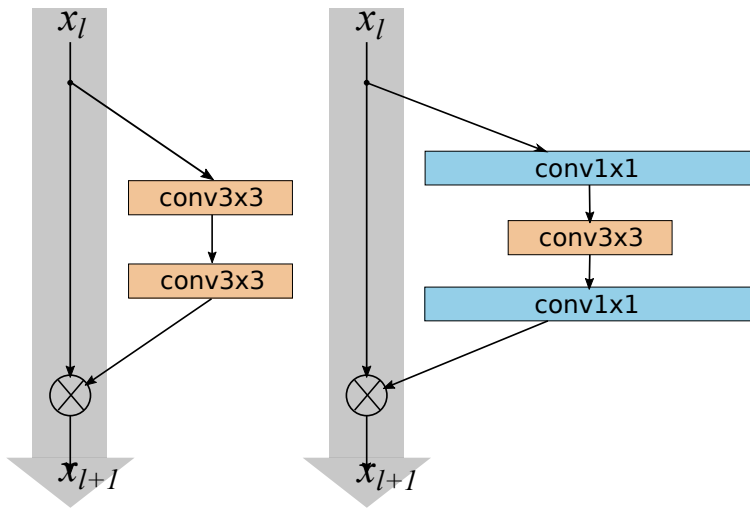
AlexNet



VGG



ResNet



(Zagoruyko and Komodakis 2016)

What is a Transformer?

Self-attention

Core operation in the Transformer:

$$Q = W_q X$$

$$K = W_k X$$

$$V = W_v X$$

$$Z = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) V$$

Self-attention

Core operation in the Transformer:

$$Q = W_q X$$

$$K = W_k X$$

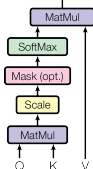
$$V = W_v X$$

$$Z = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) V$$

- ▶ Quadratic cost
- ▶ Input order equivariant

Multi-head attention

Scaled Dot-Product Attention



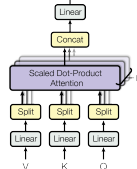
$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}}) V$$

$$MultiHead(Q, K, V) = Concat(H_1, ..., H_h) W^O$$

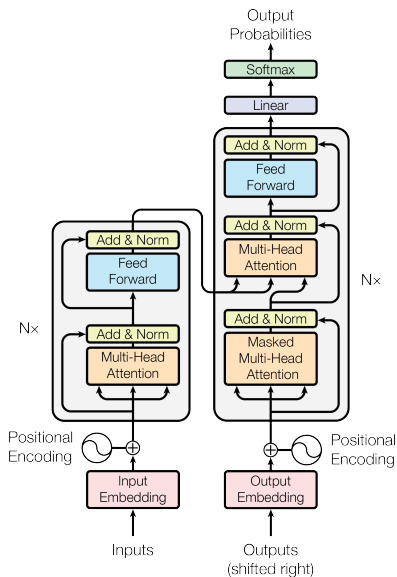
$$H_i = Attention(QW_i^Q, KW_i^K, VW_i^V), i = 1, ..., h$$

$$W_i^Q \in \mathbb{R}^{d_{model} \times d_k}, W_i^K \in \mathbb{R}^{d_{model} \times d_k}, W_i^V \in \mathbb{R}^{d_{model} \times d_v}, W_i^O \in \mathbb{R}^{hd_v \times d_{model}}$$

Multi-head Attention

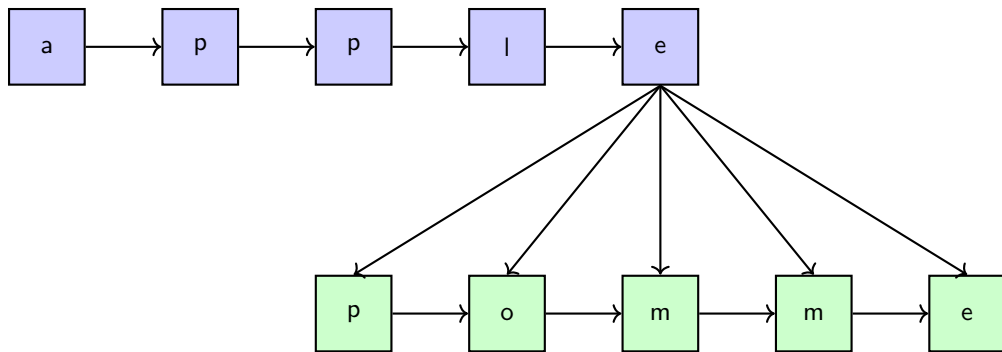


Original Transformer (Vaswani et al. 2017)



Seq2Seq Translation: Apple to Pomme

Encoder



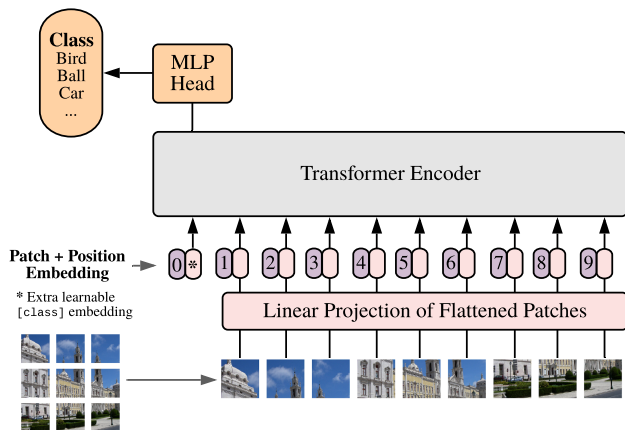
ViT: Vision Transformers

Vision Transformer, ViT (Dosovitskiy et al. 2020) -

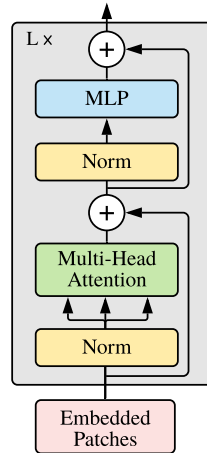
Inspired by the Transformer scaling successes in NLP, we experiment with applying a standard Transformer directly to images, with the fewest possible modifications. To do so, we split an image into patches and provide the sequence of linear embeddings of these patches as an input to a Transformer. Image patches are treated the same way as tokens (words) in an NLP application. We train the model on image classification in supervised fashion. (Dosovitskiy et al. 2020)

ViT: Vision Transformers

Vision Transformer (ViT)



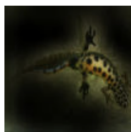
Transformer Encoder



(Dosovitskiy et al. 2020)

ViT attention

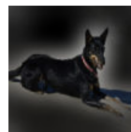
1



2



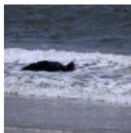
3



9



10

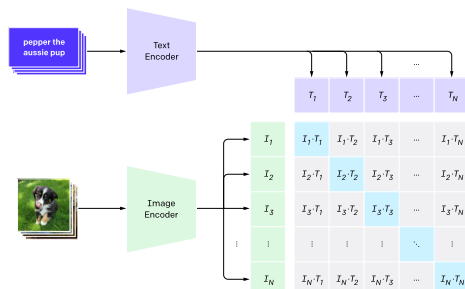


11



What is CLIP?

1. Contrastive pre-training



- ▶ **CLIP** stands for **C**ontrastive **L**anguage–**I**mage **P**retraining.
- ▶ Developed by OpenAI, it combines text and image understanding.
- ▶ Utilizes a large dataset of text-image pairs for training.

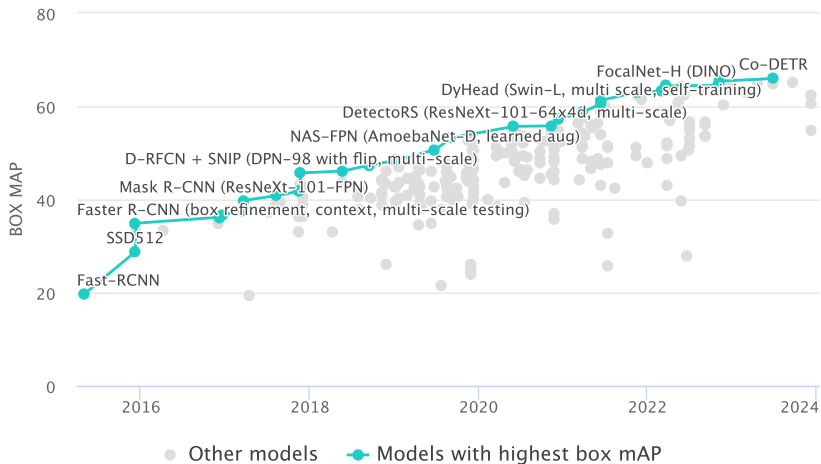
<https://openai.com/index/clip>

Object Detection

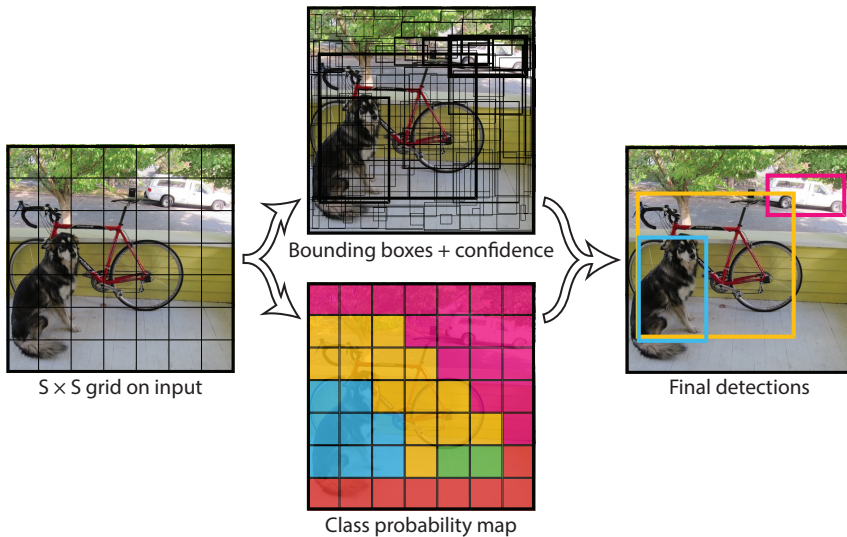
Object Detection



Object Detection



YOLO



Classical approach to detection

- Popular approach: detection $:=$ classification of boxes

Classical approach to detection

- ▶ Popular approach: detection $:=$ classification of boxes
- ▶ Requires selecting a subset of candidate boxes

Classical approach to detection

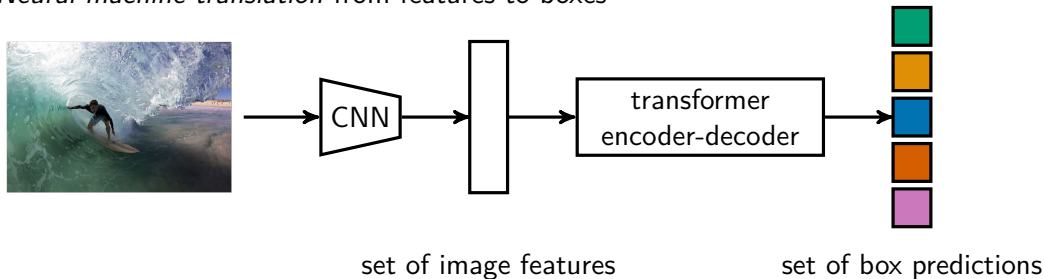
- ▶ Popular approach: detection $:=$ classification of boxes
- ▶ Requires selecting a subset of candidate boxes
- ▶ Regression step to refine the predictions

Classical approach to detection

- ▶ Popular approach: detection $:=$ classification of boxes
- ▶ Requires selecting a subset of candidate boxes
- ▶ Regression step to refine the predictions
- ▶ Typically non-differentiable

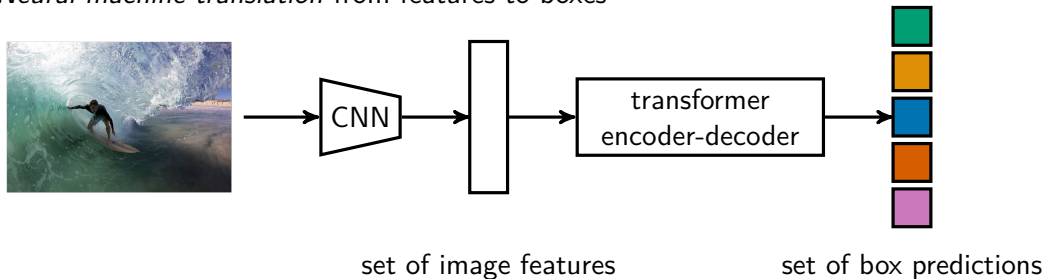
DETR: Rethinking object detection

- *Neural machine translation* from features to boxes



DETR: Rethinking object detection

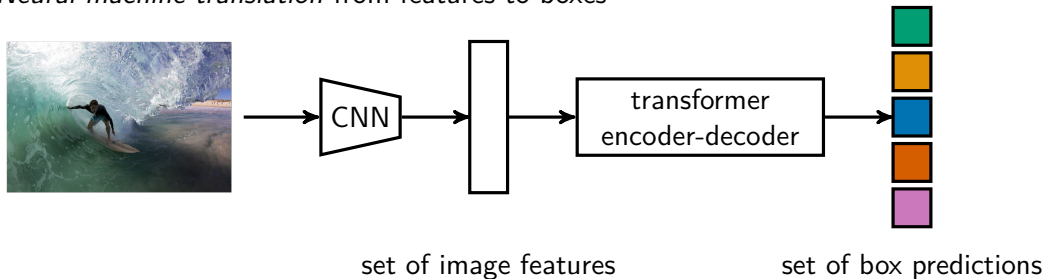
- *Neural machine translation* from features to boxes



- End-to-end parallel set prediction

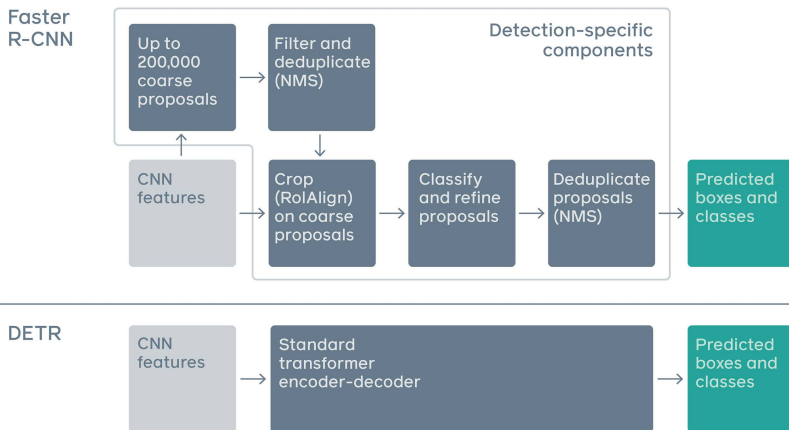
DETR: Rethinking object detection

- *Neural machine translation* from features to boxes

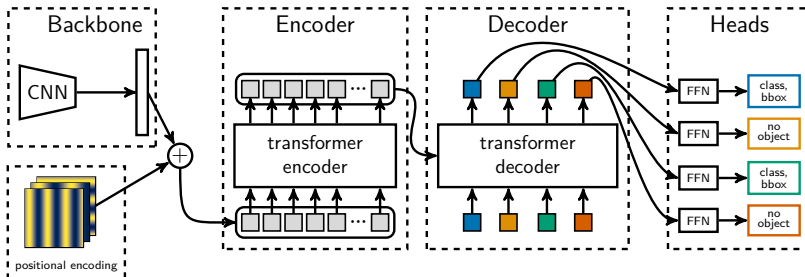


- End-to-end parallel set prediction
- Global scene reasoning

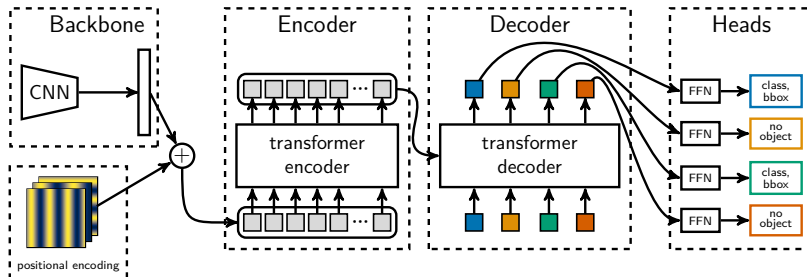
Streamlined detection pipeline



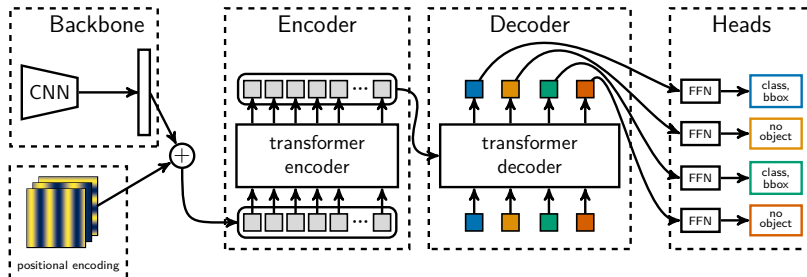
- We use standard ResNet from torchvision



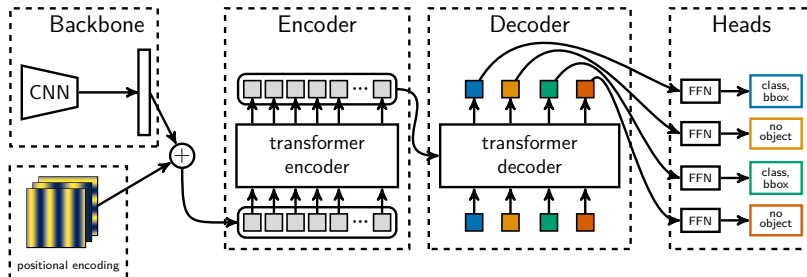
- We use standard ResNet from torchvision
- Pretraining on Imagenet is key (labels or SSL)



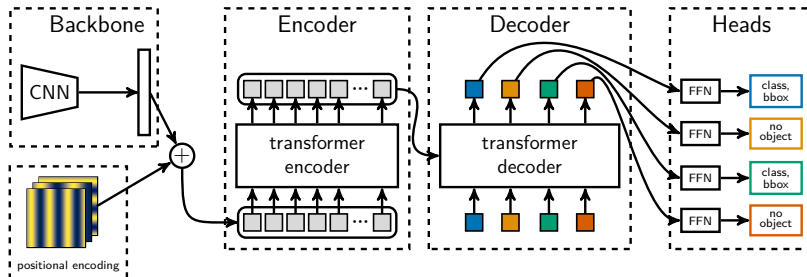
- We use 2D sine/cosine embeddings



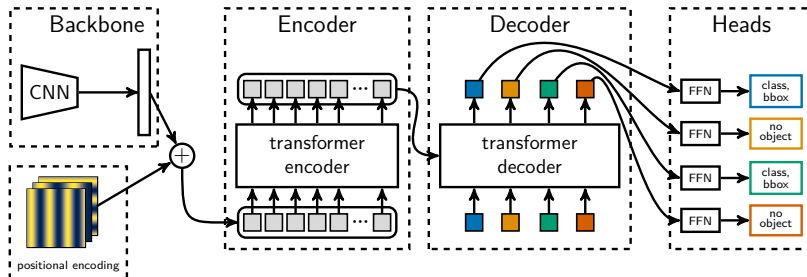
- We use 2D sine/cosine embeddings
- Embeddings are added to features



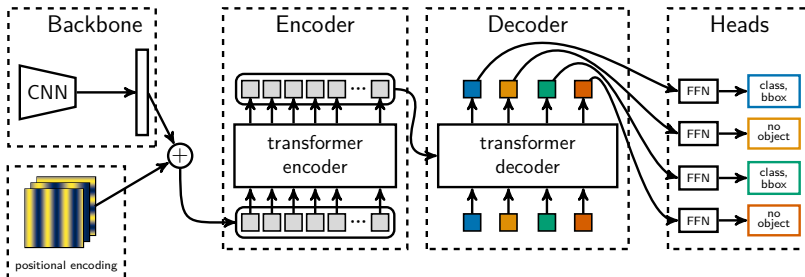
- We use 6 layers of transformer encoder



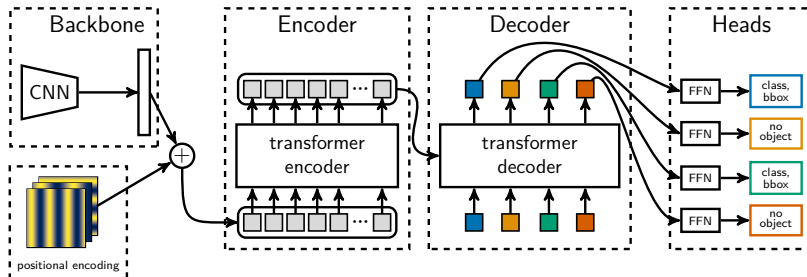
- We use 6 layers of transformer encoder
- Global reasoning through attention



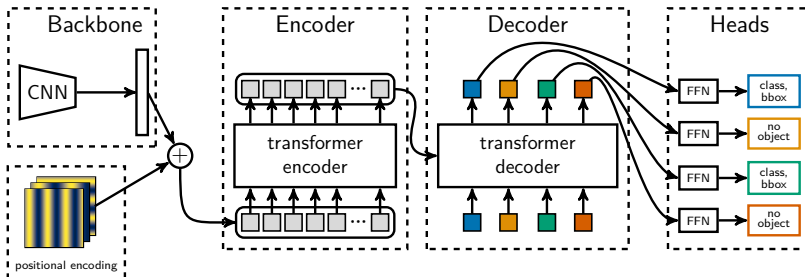
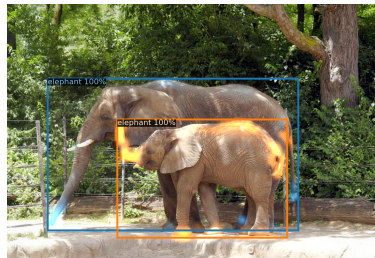
- We use 6 layers of transformer encoder
- Global reasoning through attention
- Starts separating instances



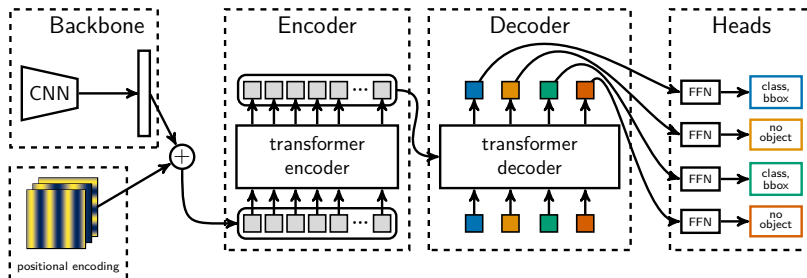
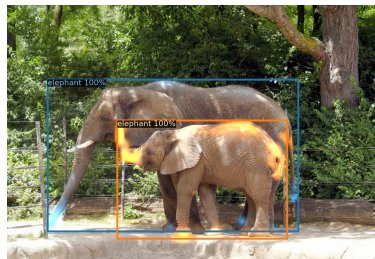
- We use 6 layers of transformer decoder;



- We use 6 layers of transformer decoder;
- Attention focuses on extremities



- We use 6 layers of transformer decoder;
- Attention focuses on extremities
- Predictions are refined at each layer **in parallel**



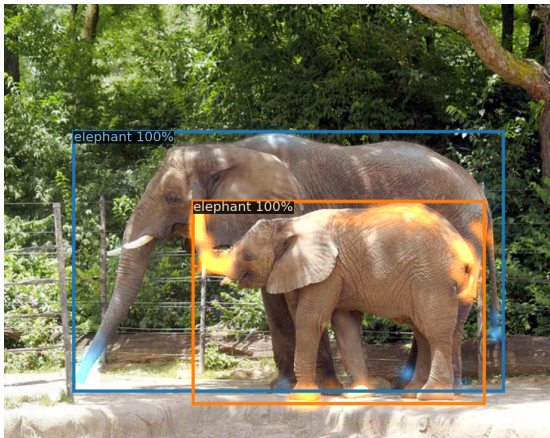
Decoder attention weights

6 decoder layers of:

- ▶ self-attention
- ▶ enc-dec attention
- ▶ FFN
- ▶ LayerNorm

All outputs are decoded *in parallel*

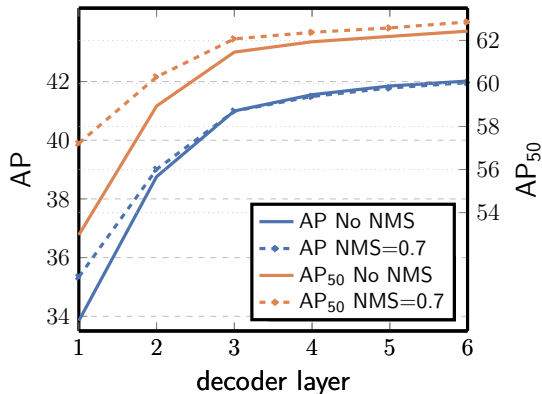
Attention focuses on extremities



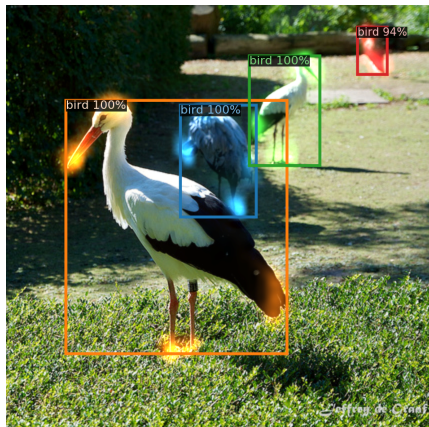
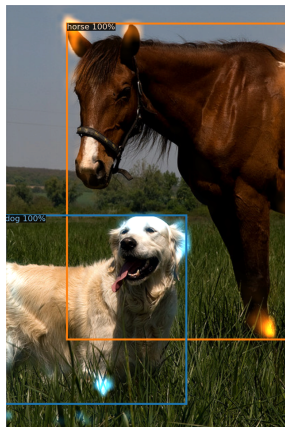
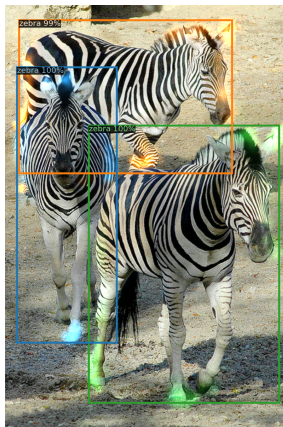
Decoder: no NMS needed

- ▶ AP/AP₅₀ goes up in lower layers (no communication)
- ▶ AP goes down in the last layers
- ▶ AP₅₀ goes up slightly

There is no need for NMS in DETR

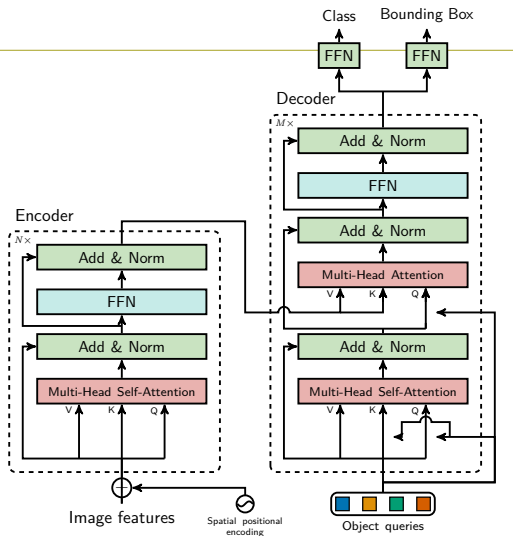


Decoder attention weights



Transformer architecture variations

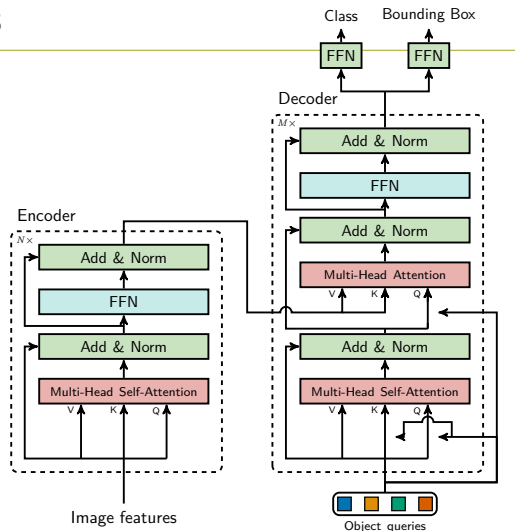
- Base transformer
39.2 AP^a



^aAshish Vaswani et al. (2017). "Attention is All you Need". In: *NeurIPS*.

Transformer architecture variations

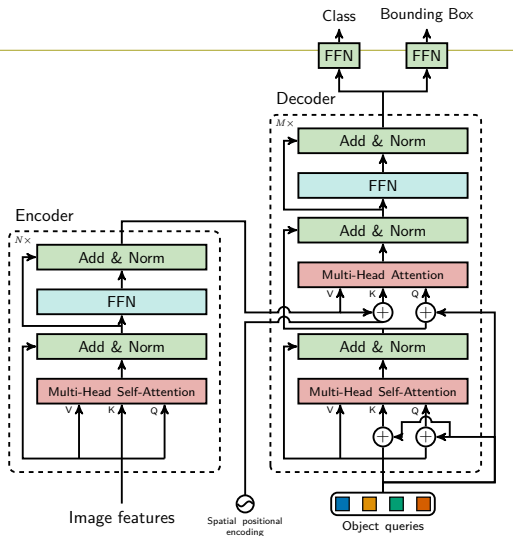
- Base transformer
39.2 AP^a
- No input positional encoding
32.8 AP



^aAshish Vaswani et al. (2017). "Attention is All you Need". In: *NeurIPS*.

Transformer architecture variations

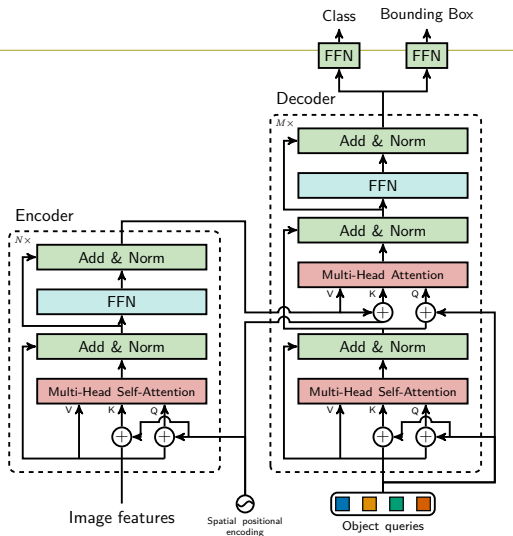
- Base transformer
39.2 AP^a
- No input positional encoding
32.8 AP
- Encodings in decoder attentions only
39.3 AP



^aAshish Vaswani et al. (2017). "Attention is All you Need". In: *NeurIPS*.

Transformer architecture variations

- ▶ Base transformer
39.2 AP^a
- ▶ No input positional encoding
32.8 AP
- ▶ Encodings in decoder attentions only
39.3 AP
- ▶ Encodings in all attentions
40.6 AP

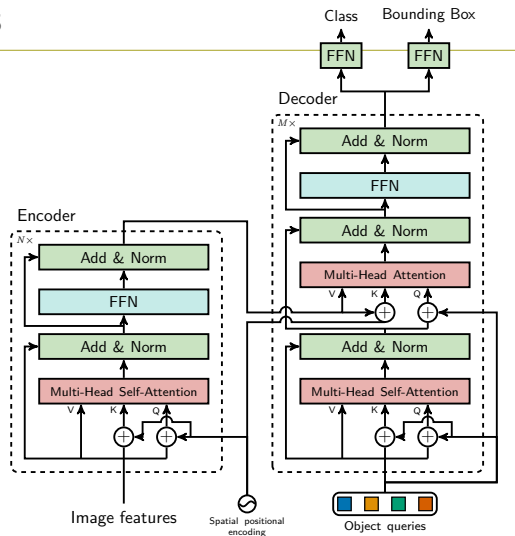


^aAshish Vaswani et al. (2017). "Attention is All you Need". In: *NeurIPS*.

Transformer architecture variations

- Base transformer
39.2 AP^a
- No input positional encoding
32.8 AP
- Encodings in decoder attentions only
39.3 AP
- Encodings in all attentions
40.6 AP

All transformer parts are contributing!



^aAshish Vaswani et al. (2017). "Attention is All you Need". In: *NeurIPS*.