

# Discussion of transformer models

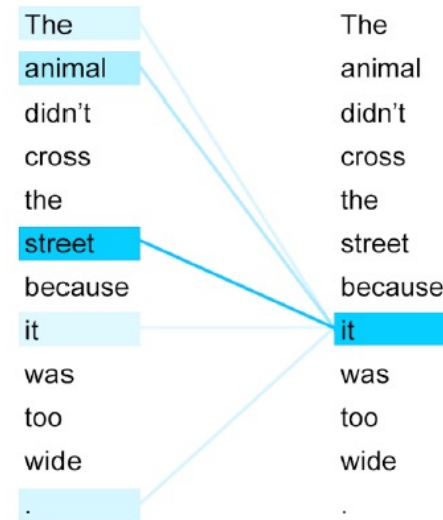
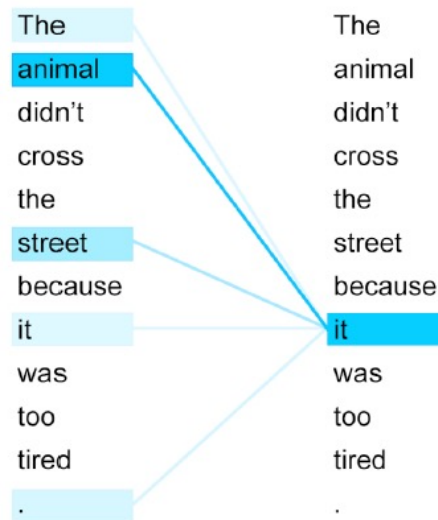
**Mitko Veta**  
M.Veta@tue.nl

## Attention in language:

“This GitHub page contains all the general information about the course and the study materials.”

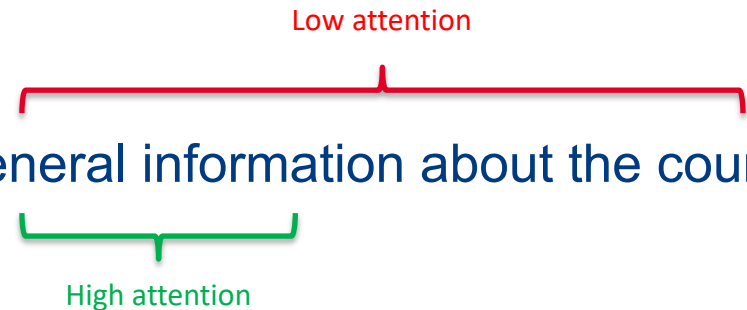
Low attention

High attention

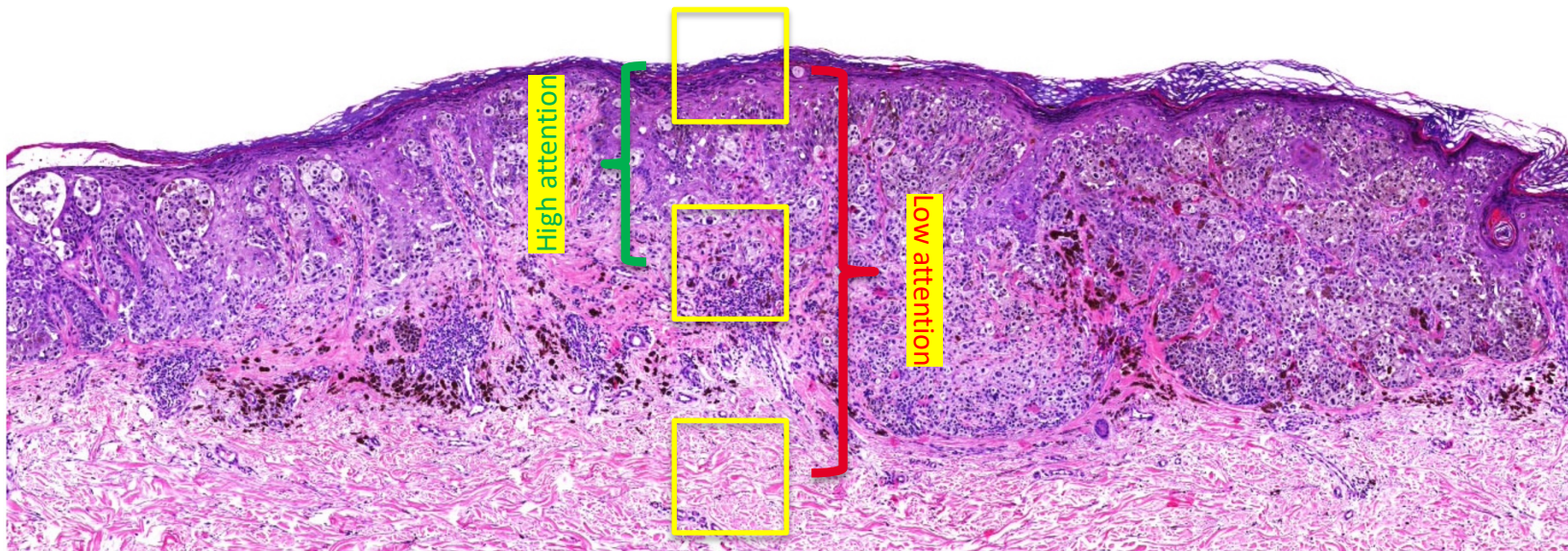


## Attention in language:

“This GitHub page contains all the general information about the course and the study materials.”



## Attention in vision (invasive melanoma histology example):



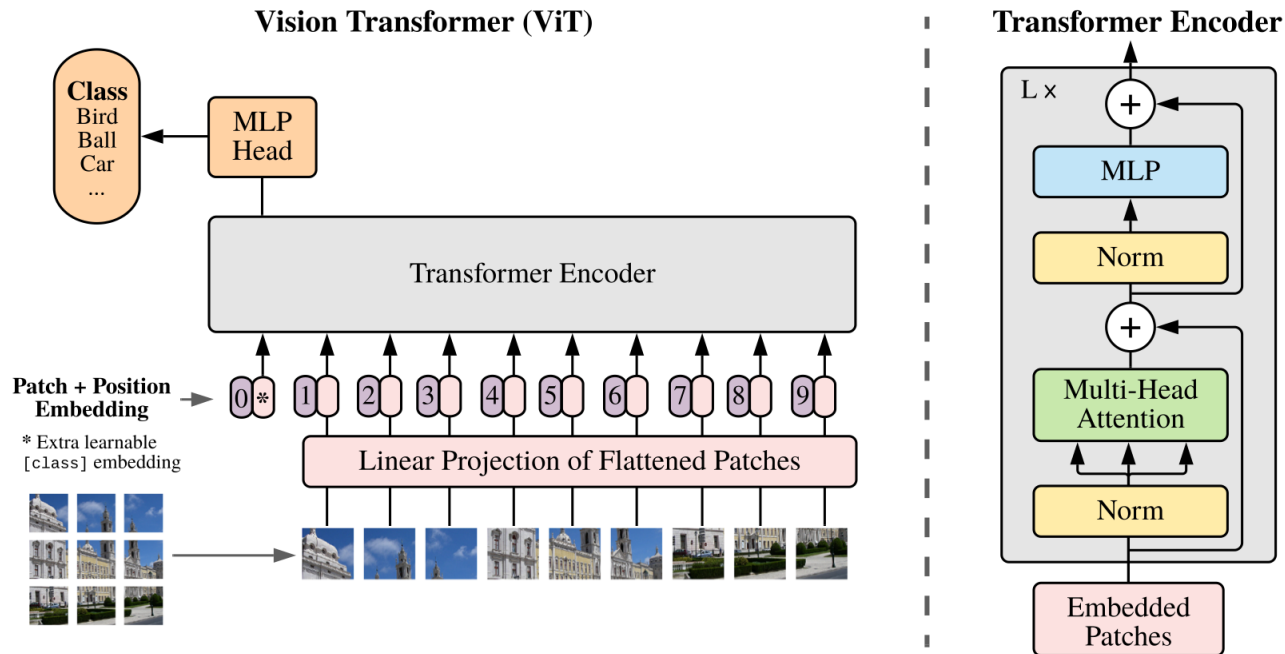
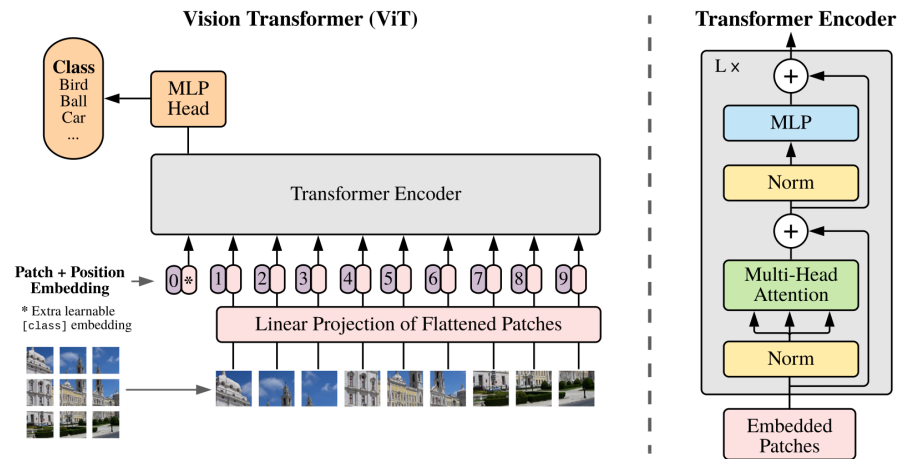


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).



```
def EncoderBlock(X):
    Z = LayerNorm(MultiHeadAttn(Q=X, K=X, V=X) + X)
    E = LayerNorm(FeedForward(Z) + Z)
    return E
```

```
def Encoder(X, N):
    E = POS(Embed(X))
    for n in range(N):
        E = EncoderBlock(E)
    return E
```

## Topic 1: Data requirements

**Q:** Why do transformers generalize worse when trained on insufficient amounts of data compared to CNN?

**Q:** In the vision transformers paper, it was mentioned that on smaller dataset, models like ResNet perform better than ViT. Given that many medical imaging datasets are orders of magnitude smaller than the "smaller datasets" that were used in that study, how can we best use ViTs in the biomedical field? Are techniques like transfer learning from larger non-medical datasets be able to be used as pretraining?

**Q:** Considering the larger data requirements of ViTs compared to its CNN counterpart, is it possible that ViTs are adopted for medical applications?

## Topic 1.5: Inductive bias

**Q:** Why is the inductive bias lower for transformers compared to convolutional neural networks?

**Q:** Transformers have less **biases** than CNNs but still manage to roughly learn the same things (look at the internal representations of the 16x16 paper) and are able to outcompete them in performance relative to training time. Which **biases** do transformers still have and which should we hope to remove to improve performance even more?

Inductive bias ~ assumptions about the data built into your model.

In CNNs:

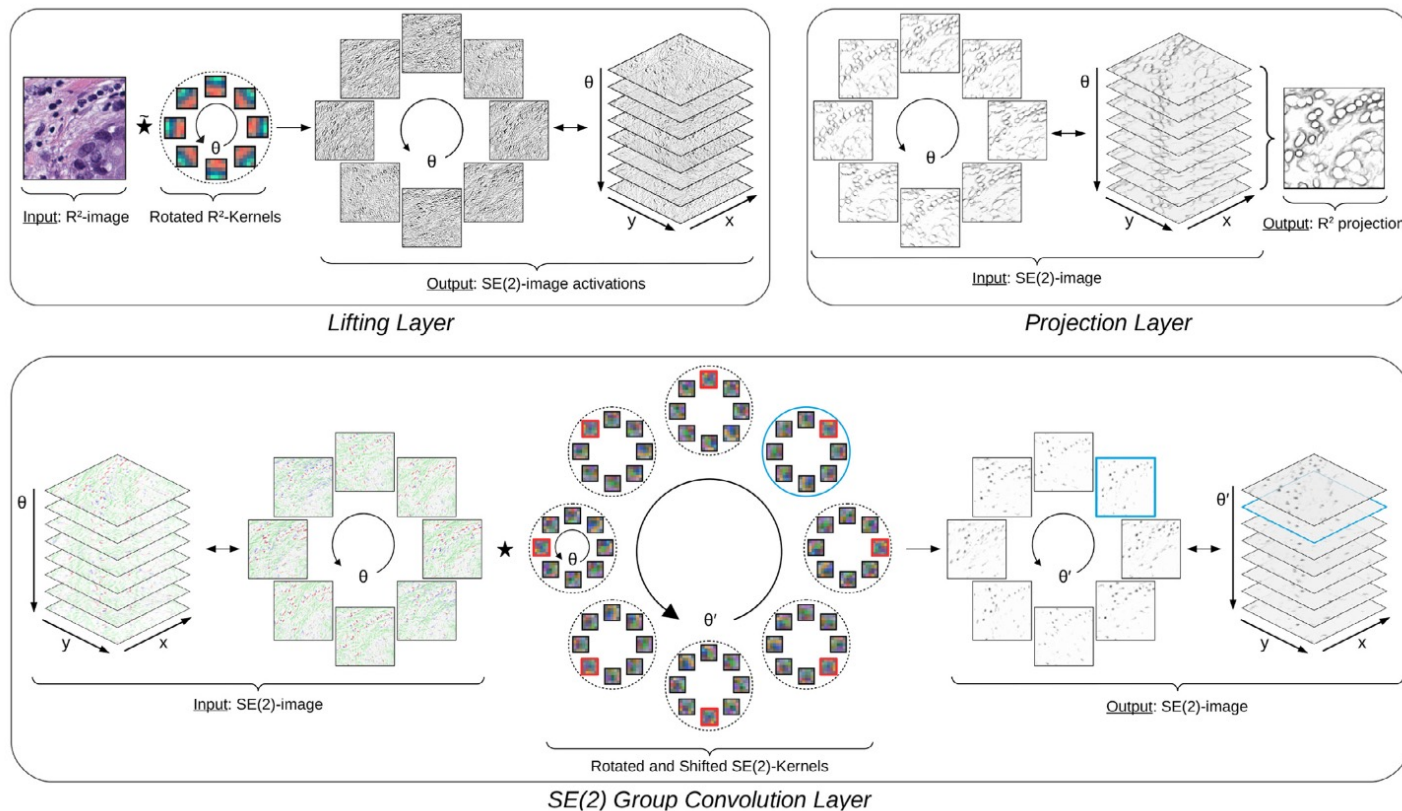
- Locality (small kernels)
- Equivariance (weight sharing)
- Invariance (pooling)

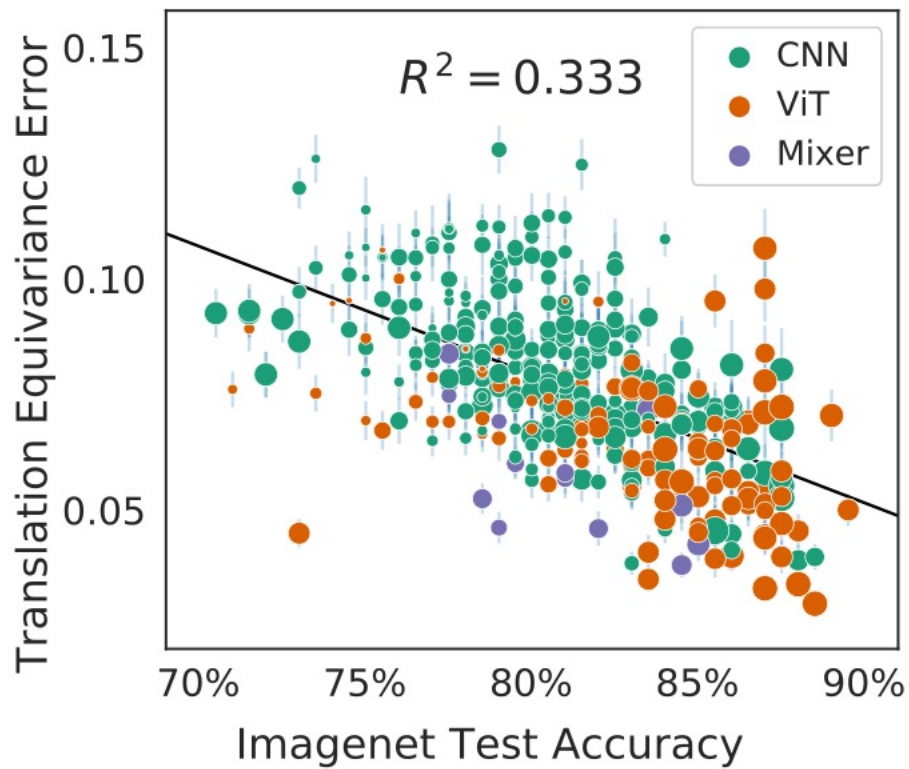
In transformers:

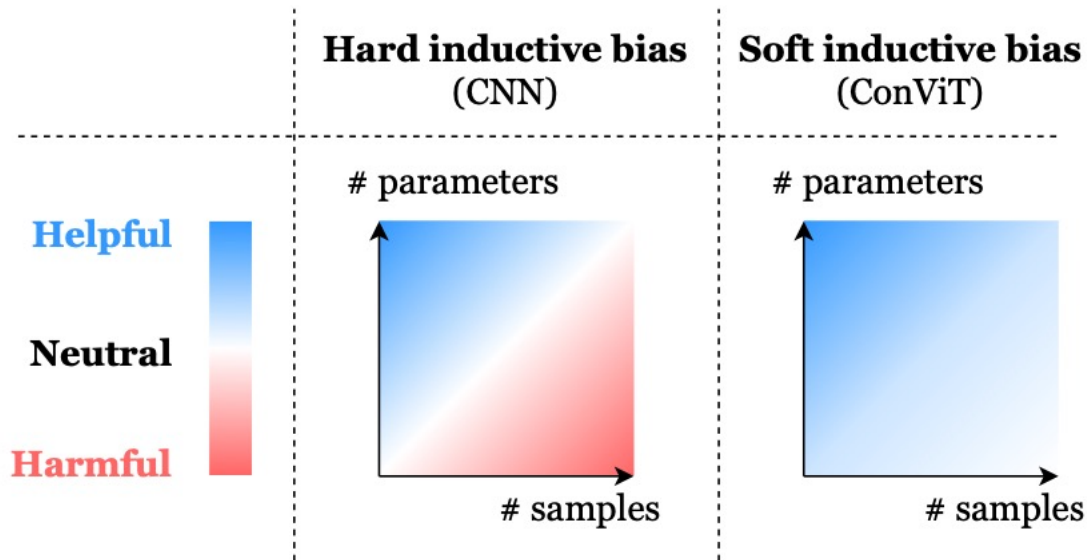
- Locality (cutting into patches)
- Equivariance (linear projection)
- (Both are "weaker" inductive biases)



Example from our own work: rotational equivariance and invariance (in addition to translational)  $\rightarrow$  higher inductive bias:







**Figure 1. Soft inductive biases can help models learn without being restrictive.** Hard inductive biases, such as the architectural constraints of CNNs, can greatly improve the sample-efficiency of learning, but can become constraining when the size of the dataset is not an issue. The soft inductive biases introduced by the ConViT avoid this limitation by vanishing away when not required.

## Topic 2: Handling biomedical data

**Q:** How can vision transformers be adapted to handle other types of medical data, such as combining medical images with health records?

**Q:** (Apart from text and images,) to which topic could transformers be applied additionally? (You need lots of data as we can see from the second paper, so that is a constraint to keep in mind)

**Q:** What are the advantages and disadvantages of vision transformer models compared to state-of-the-art convolutional neural network models for (medical) image analysis?

## Disadvantages:

Need more data to train

Computationally demanding

## Advantages:

Less assumptions about the data

Include global information in the early layers

Transformers are easily adaptable to multi-modal data - Why?

## Topic 3: Self-supervised learning

**Q:** How does the self-Distillation with NO-labels (DINO) method for pretraining vision transformers work?

**Q:** In the case of histopathology, visual transformers enable us to perform well both in classification tasks on easily discernible visual concepts as well as capturing contextual and hierarchical information in the tumor microenvironment. However, pretraining such models is both limited due to the high data dimensionality and low amount of data points as well as being computationally intractable on commercial workstations. How can we deal with the tradeoff between performance and computational complexity and is it worth the effort over picking a simpler computer vision algorithm or neural network architecture?

DINO is a method for self-supervised pre-training.

Self-supervision ~ supervised training with "intrinsic" labels on a pretext task. Goal is to learn good image representations that can be transferred to a downstream task.

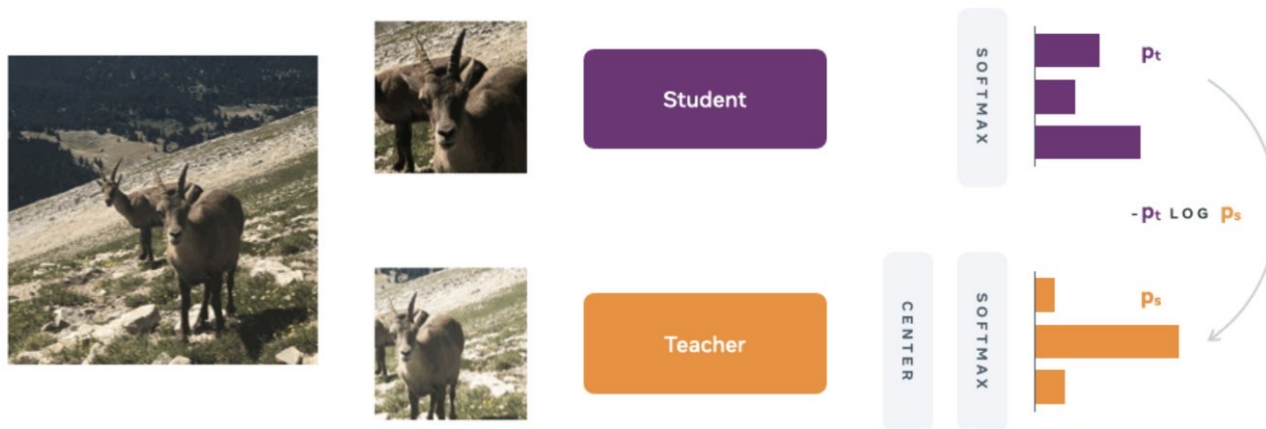
Advantages:

- Can use large, unlabelled datasets for pretraining
- "Richer" training signal compared to narrow supervised tasks

## Example pretext tasks:

- Predict one part of the image from another
- Predict if two sub-images originate from the same or from different images
- Force two different transformations of the same image to have the same representation

## DINO:





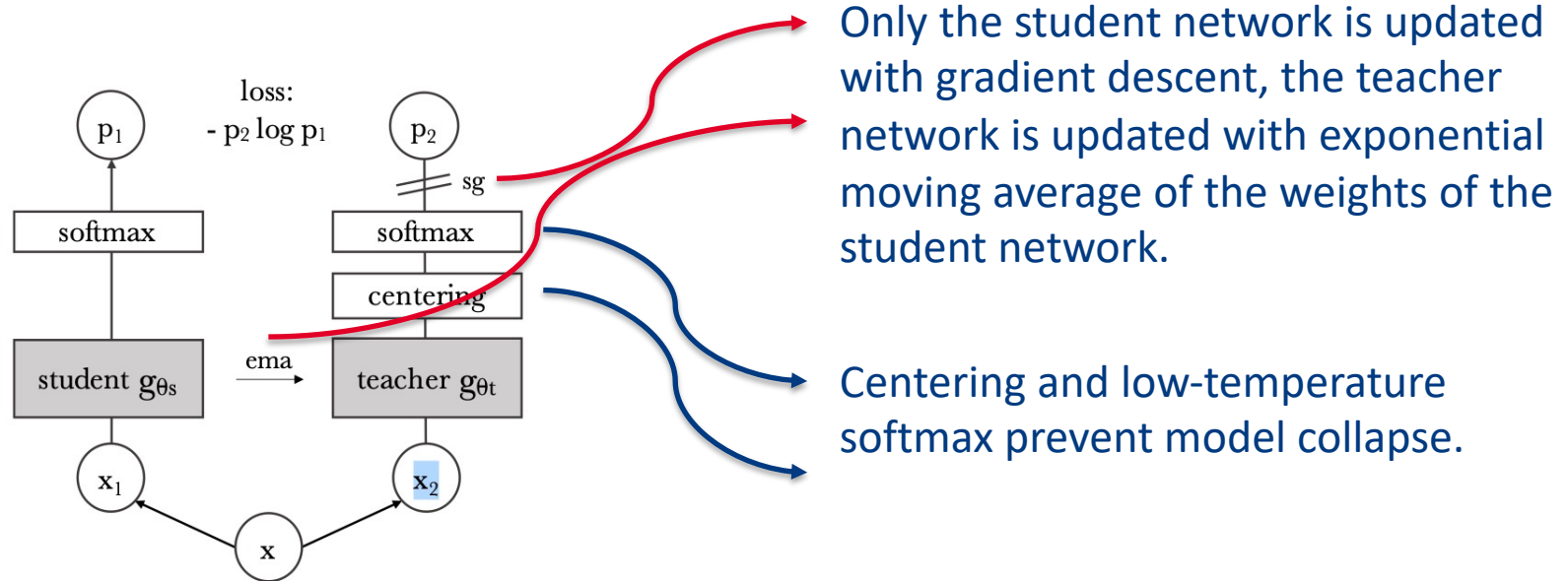


Figure 2: **Self-distillation with no labels.** We illustrate DINO in the case of one single pair of views ( $x_1$ ,  $x_2$ ) for simplicity. The model passes two different random transformations of an input image to the student and teacher networks. Both networks have the same architecture but different parameters. The output of the teacher network is centered with a mean computed over the batch. Each networks outputs a  $K$  dimensional feature that is normalized with a temperature softmax over the feature dimension. Their similarity is then measured with a cross-entropy loss. We apply a stop-gradient (sg) operator on the teacher to propagate gradients only through the student. The teacher parameters are updated with an exponential moving average (ema) of the student parameters.

$x_1$ : global and local views of the image

$x_2$ : local views of the image

## Emerging properties:

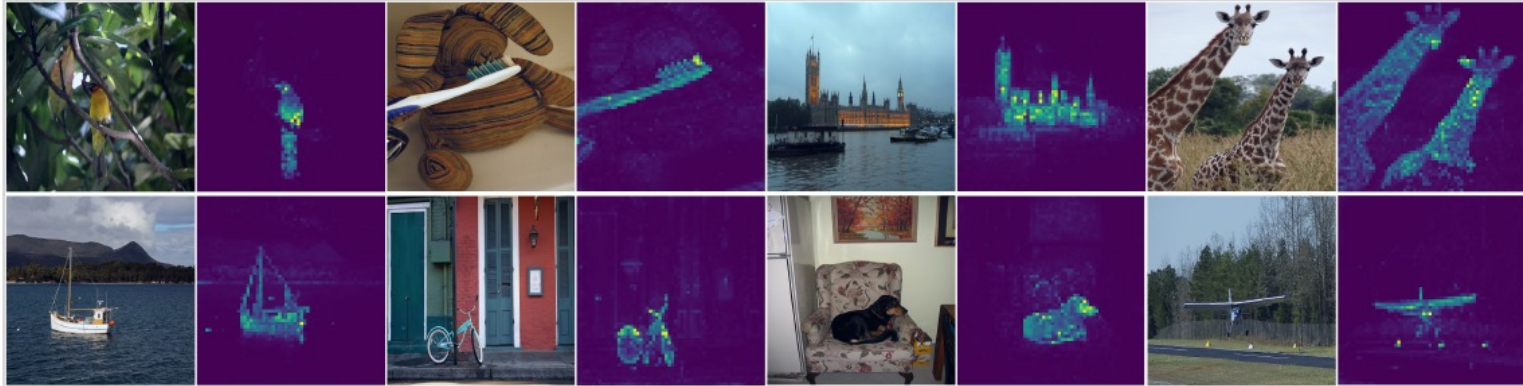


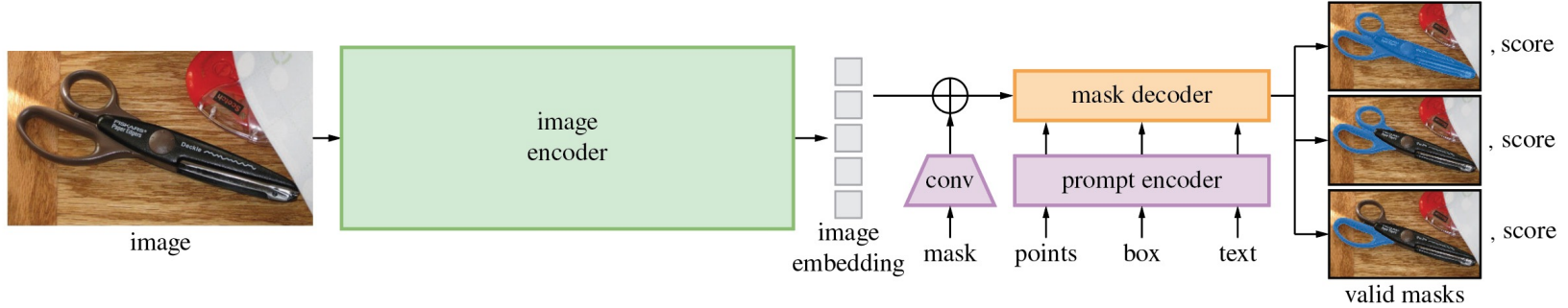
Table 10:  $k$ -NN and linear evaluation for ViT-S/16 and ResNet-50 pre-trained with DINO. We use ImageNet-1k [60] (“Inet”), Places205 [84], PASCAL VOC [24] and Oxford-102 flowers (“FLOWERS”) [46]. ViT trained with DINO provides features that are particularly  $k$ -NN friendly.

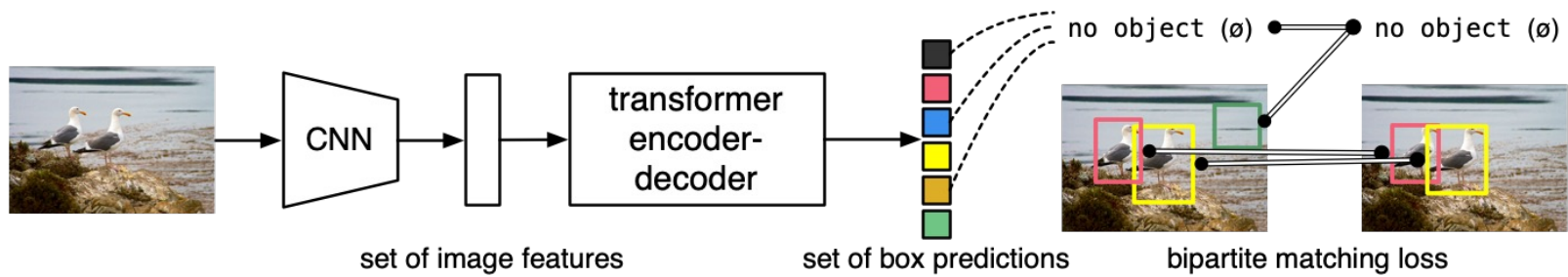
|                  | Logistic |       |          | $k$ -NN |       |          |
|------------------|----------|-------|----------|---------|-------|----------|
|                  | RN50     | ViT-S | $\Delta$ | RN50    | ViT-S | $\Delta$ |
| Inet 100%        | 72.1     | 75.7  | 3.6      | 67.5    | 74.5  | 7.0      |
| Inet 10%         | 67.8     | 72.2  | 4.4      | 59.3    | 69.1  | 9.8      |
| Inet 1%          | 55.1     | 64.5  | 9.4      | 47.2    | 61.3  | 14.1     |
| Pl. 10%          | 53.4     | 52.1  | -1.3     | 46.9    | 48.6  | 1.7      |
| Pl. 1%           | 46.5     | 46.3  | -0.2     | 39.2    | 41.3  | 2.1      |
| VOC07            | 88.9     | 89.2  | 0.3      | 84.9    | 88.0  | 3.1      |
| FLOWERS          | 95.6     | 96.4  | 0.8      | 87.9    | 89.1  | 1.2      |
| Average $\Delta$ |          |       | 2.4      |         |       | 5.6      |

## **Topic 4: Transformers for image segmentation and object detection**

**Q:** One of the papers mentioned that at the time of publishing transformers had not yet been used for image detection or segmentation purposes. Is that still the case and what changes to the model architecture would be required for these applications to be possible?

**Q:** How could transformers be used for detection/ segmentation in images? (for example to outline an object and give as output what the object is)





## Other, technical questions

**Q:** In transformer models, which of the following vectors have to be of equal length: key, query, value?

**Q:** Does the order of the input sequence elements matter?