

# Few-shot object detection using self-supervised learning

STUDENT TEAM: Daniel Reisenbüchler, Daniel Sens, Onat Sahin, Rahul PS

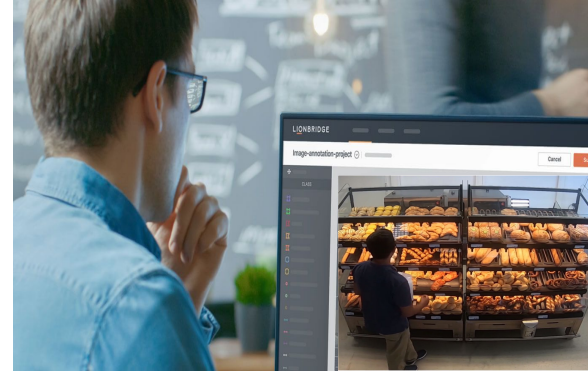
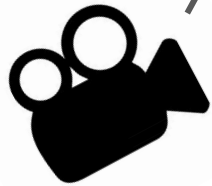
SCIENTIFIC LEAD: Mathias Sundholm  
Hamdi Belhassen

PROJECT LEAD: Dr. Ricardo Acevedo Cabra

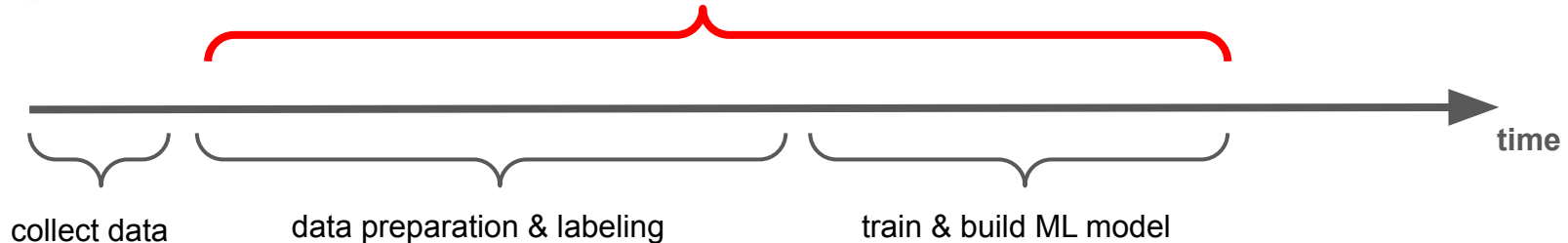
TUM CO-MENTOR: Michael Rauchensteiner



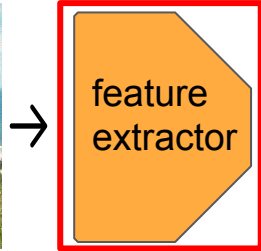
# High-level project description - Recap



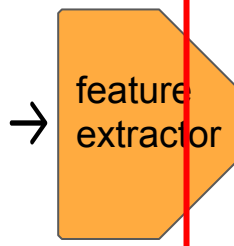
**goal: reduce amount of time**



# Recap object detection - 2 stage training process

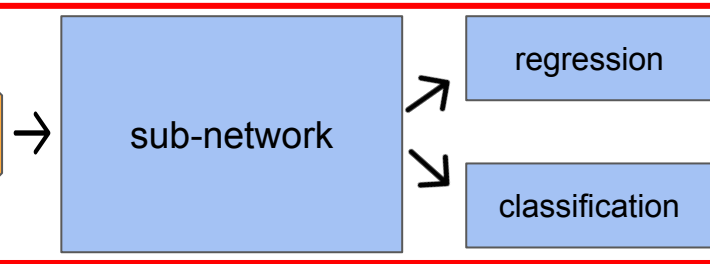


pretrain on ImageNet

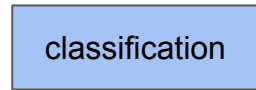
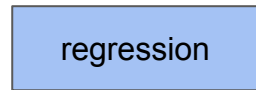


freeze

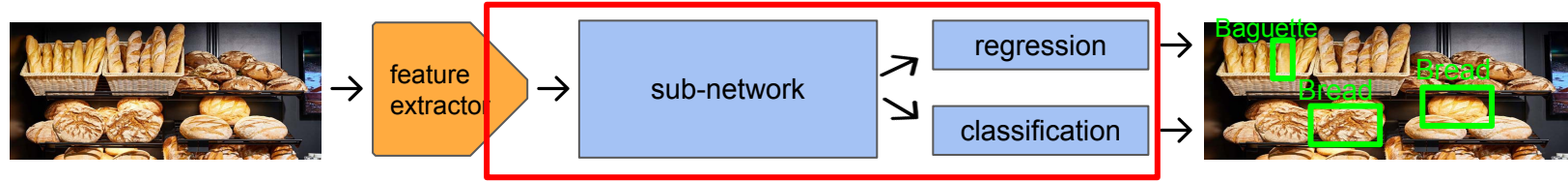
finetune



train on object detection dataset (e.g. MS COCO)

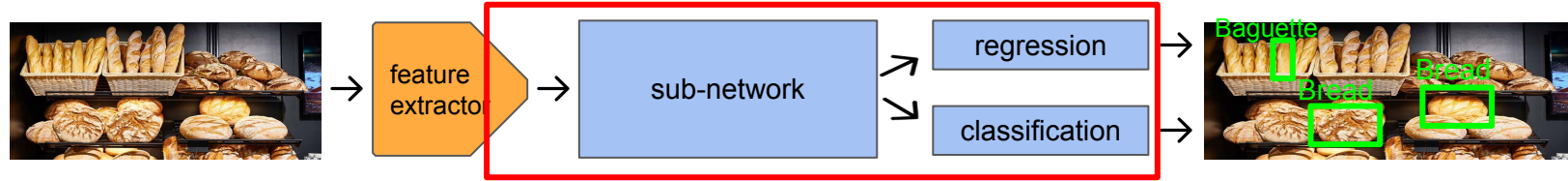


# Problem description - task environments

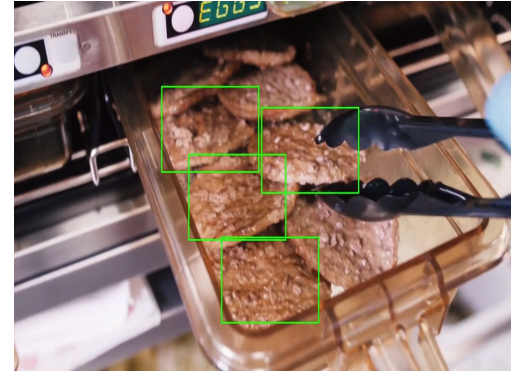
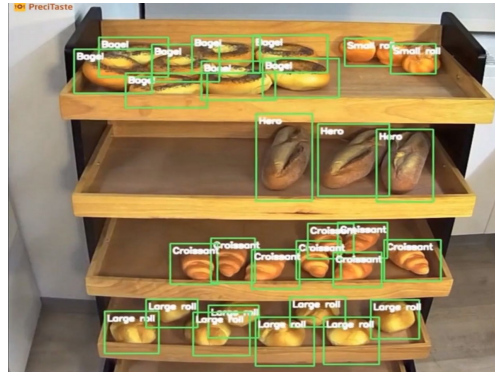


even fine-tuning needs hundreds of labeled examples to yield high accuracies

# Problem description - task environments



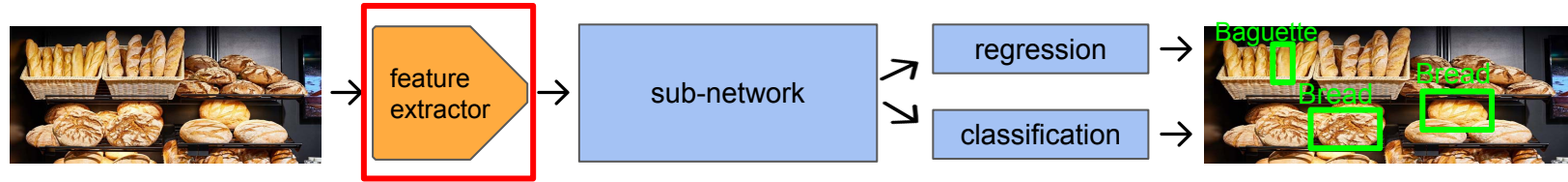
↻ label specific customer data and finetune the model again and again  
➔ **PROBLEM: domain shift for every new customer.**



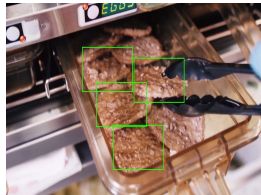
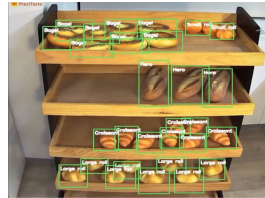
different customers - different task environments



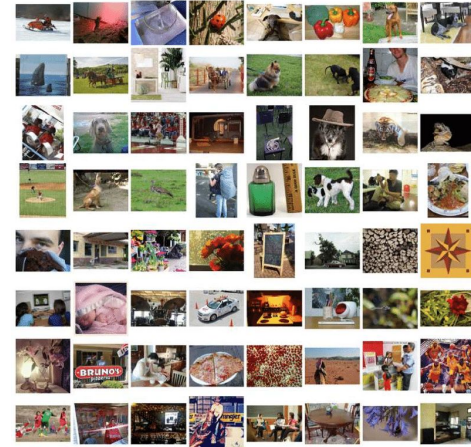
# Problem description - task environments



**Feature extractor trained on a general dataset like ImageNet may not be suitable**

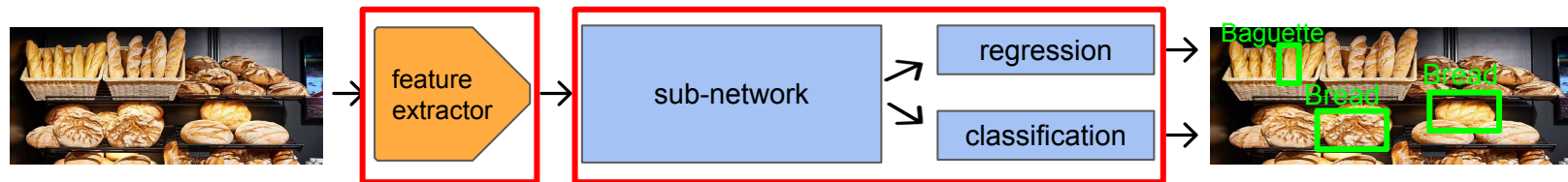


Images from PreciBake



Images from ImageNet

# Problem description - proposed solutions



**PROBLEM 1: A good embedding that works well with task environments is necessary**

**SOLUTION:** Use self-supervised learning to pre-train a backbone with unlabeled images similar to task environment.  
(For example: Food images instead of ImageNet)

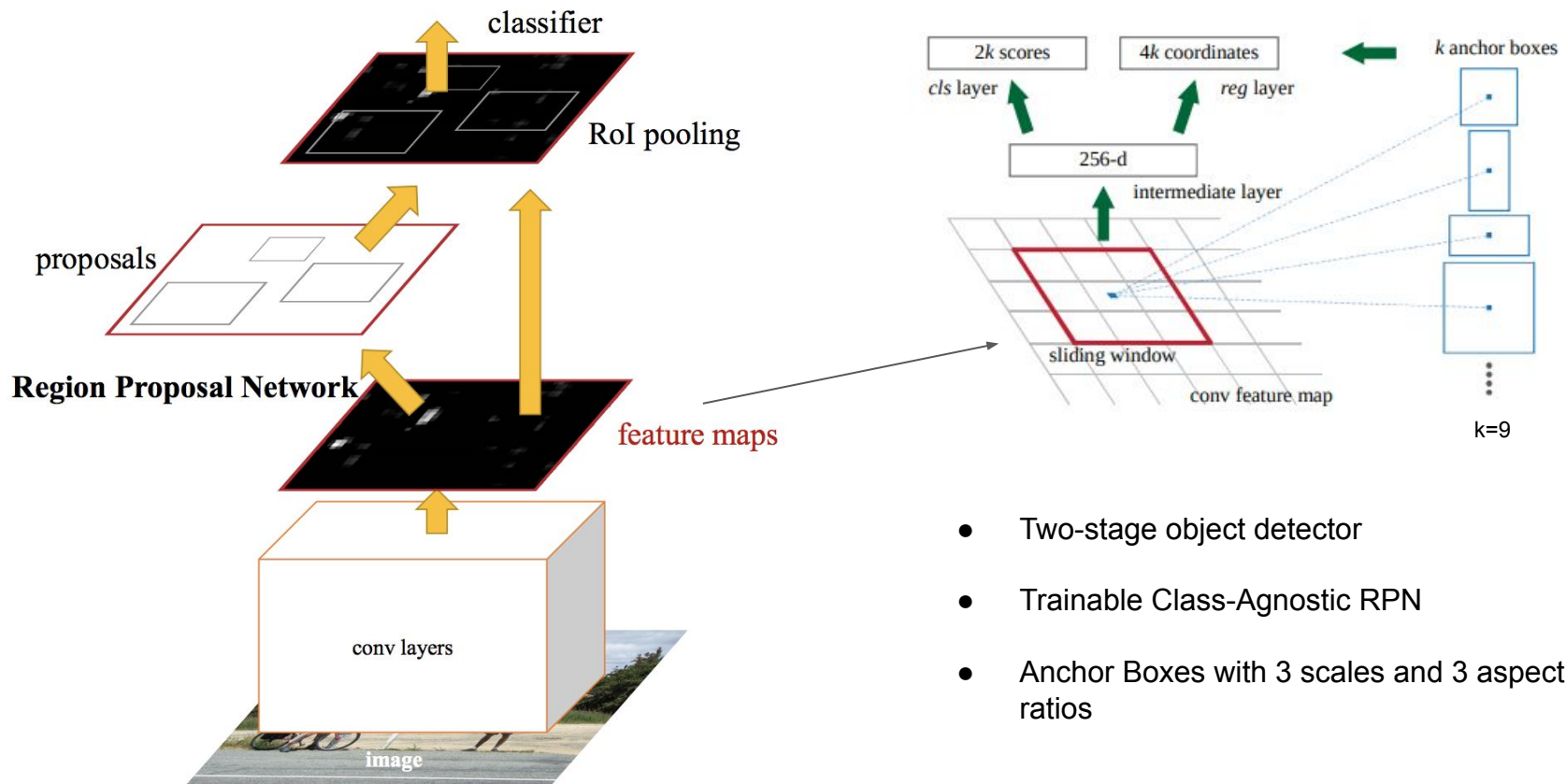
Use Faster R-CNN to experiment on backbones

**PROBLEM 2: Fine-tuning for new task environments should be effortless**

**SOLUTION:** Use few-shot object detectors

Use the backbones obtained with Attention-RPN with Multi Relation Detector

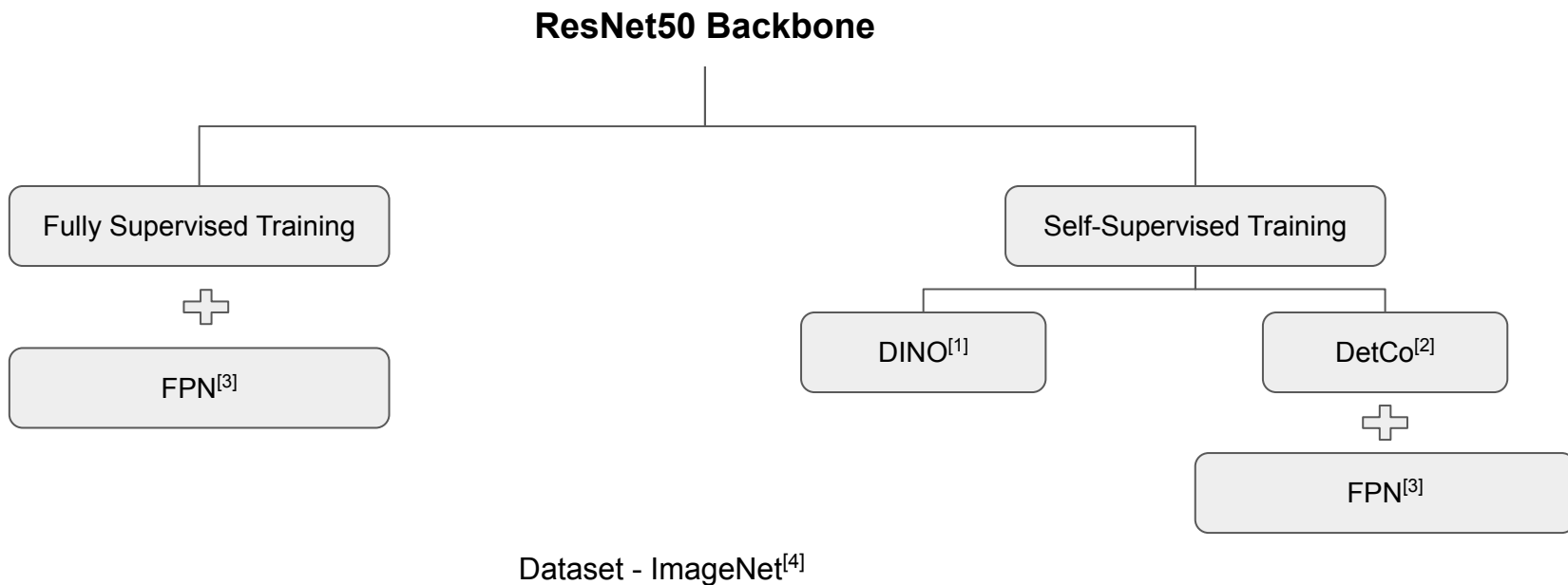
# Faster R-CNN for Object Detection<sup>[1]</sup>





# Experiments on Self-Supervised Backbones for Faster RCNN (with ResNet50)

## Overview of Experiments



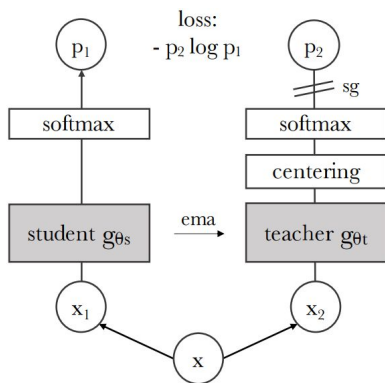
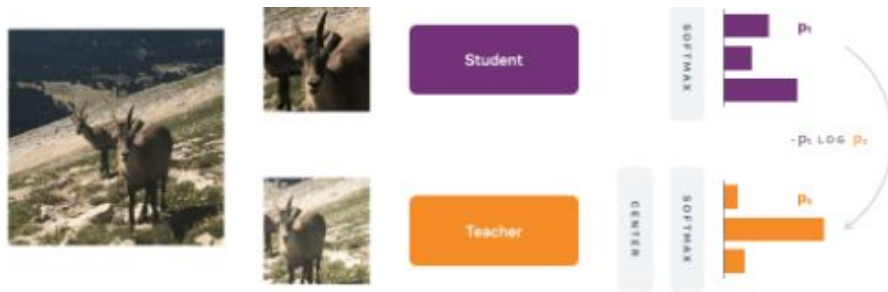
[1] Mathilde Caron, et al. "Emerging Properties in Self-Supervised Vision Transformers." (2021).

[2] Enze Xie, et al. "DetCo: Unsupervised Contrastive Learning for Object Detection." (2021).

[3] Tsung-Yi Lin et al. Feature Pyramid Networks for Object Detection. (2017).

[4] Olga Russakovsky et al. "ImageNet Large Scale Visual Recognition Challenge". In: International Journal of Computer Vision (IJCV) 115.3 (2015),

# DINO: Self-distillation with no labels<sup>[1]</sup>



- Student and Teacher network have the same architecture
- Student network weights updated by SGD
- Teacher network updated by ema of the student weights

$$\theta_t \leftarrow \lambda \theta_t + (1 - \lambda) \theta_s,$$

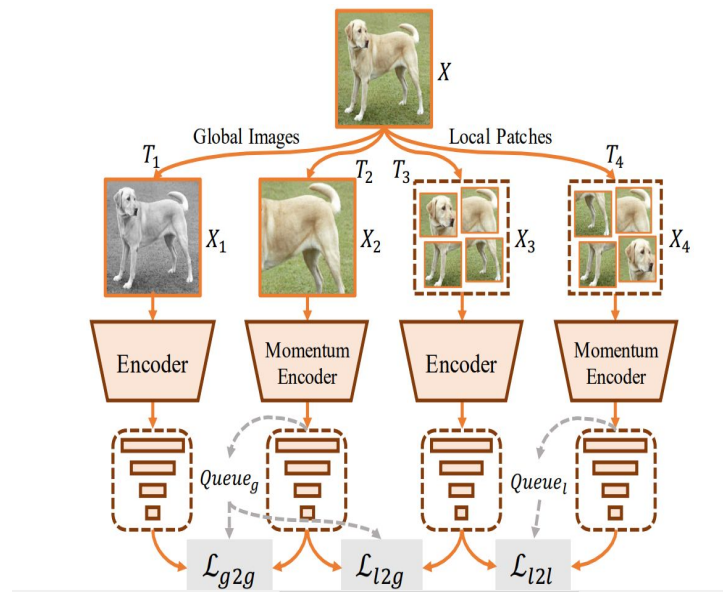
$$P_s(x)^{(i)} = \frac{\exp(g_{\theta_s}(x)^{(i)} / \tau_s)}{\sum_{k=1}^K \exp(g_{\theta_s}(x)^{(k)} / \tau_s)}$$

Centering to prevent collapse

$$g_t(x) \leftarrow g_t(x) + c$$

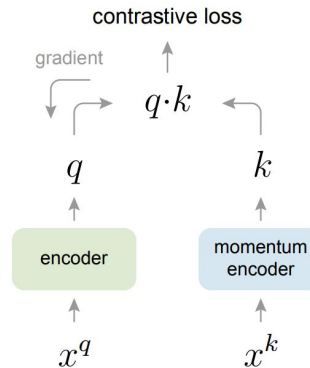
$$c \leftarrow mc + (1 - m) \frac{1}{B} \sum_{i=1}^B g_{\theta_t}(x_i)$$

# DetCo: Unsupervised Contrastive Learning for Object Detection<sup>[2]</sup>



$$\mathcal{L}(\mathbf{I}_q, \mathbf{I}_k, \mathbf{P}_q, \mathbf{P}_k) = \sum_{i=1}^4 w_i \cdot (\mathcal{L}_{g \leftrightarrow g}^i + \mathcal{L}_{l \leftrightarrow l}^i + \mathcal{L}_{g \leftrightarrow l}^i)$$

- It is important to have strong discriminative ability at each stage of the network
- Local patch features are as important as global representations of an image



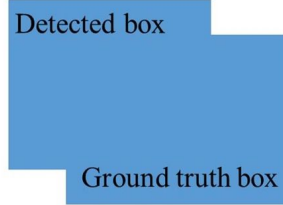
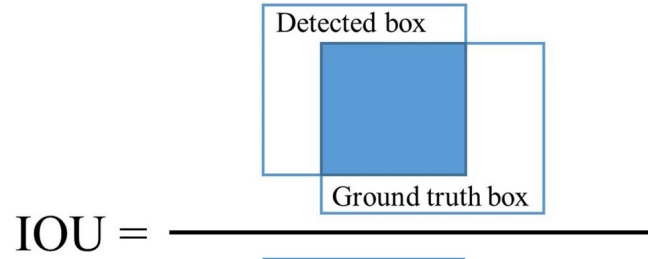
Momentum Contrast (MoCo)<sup>[3]</sup>

$$\mathcal{L}_{g \leftrightarrow g}(\mathbf{I}_q, \mathbf{I}_k) = -\log \frac{\exp(q^g \cdot k_+^g / \tau)}{\sum_{i=0}^K \exp(q^g \cdot k_i^g / \tau)}$$

[2] Enze Xie, et al. "DetCo: Unsupervised Contrastive Learning for Object Detection." (2021).

[3] Kaiming He et al. Momentum Contrast for Unsupervised Visual Representation Learning. (2020)

# Metrics for Object Detection



$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

AP

```
for classes in [1 :1 :K]
    for T in [0.50 :0.05:0.95]
        Calculate AP
        mAP+=AP
    mAP = mAP/10
mAP = mAP/K
```

AP50 | AP75

```
for classes in [1 :1 :K]
    for T ∈ {0.50} || {0.75}
        Calculate AP
    mAP = AP/K
```

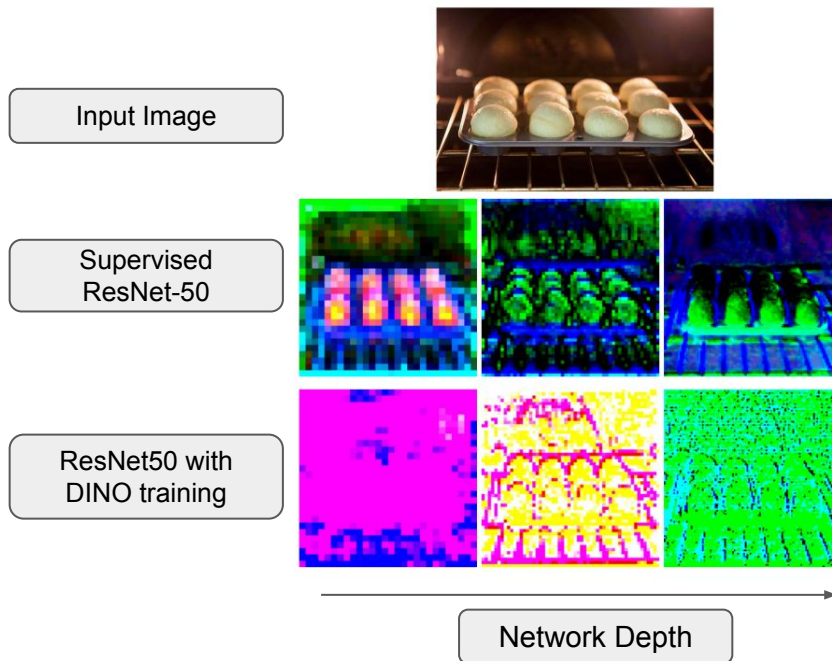
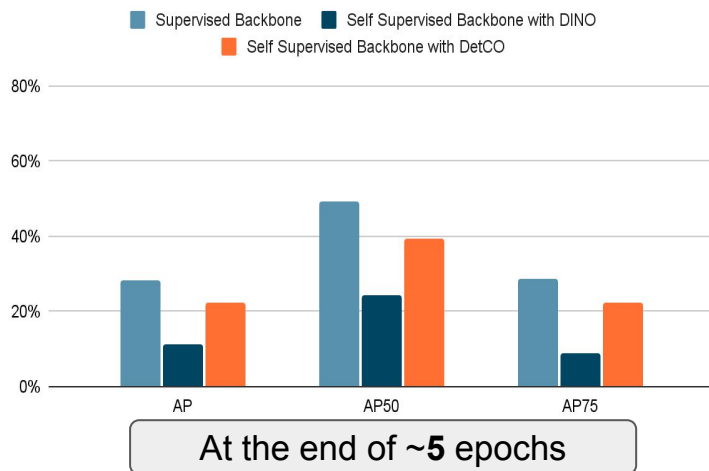
IOU > Threshold **T** → **True Positive**

# Solving Problem 1: Self-Supervised Learning

## Insight 1

It is important for backbones to preserve spatial context for object detection. DINO (which uses contrastive losses between different image crops) performs poorer than DetCO, which uses multi-stage contrastive losses in global as well as local patch scales.

AP Results for different ResNet50 BackBones on COCO dataset



**Pre-training for object classification may not translate to object detection.**

# Solving Problem 1: Self-Supervised Learning

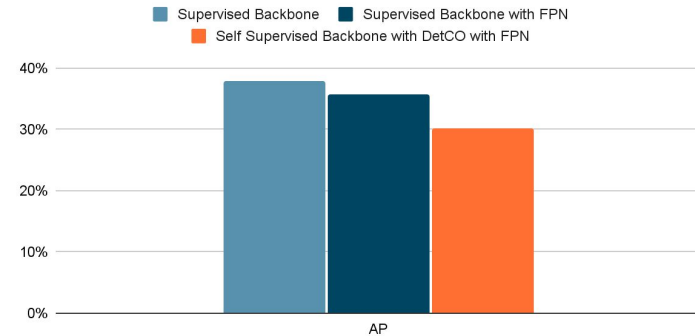
## Insight 2

Using FPN with ResNet speeds up training significantly, while giving similar performance for the same amount of epochs.



Feature Pyramid Network

AP Results for different ResNet50 BackBones with FPN on COCO dataset

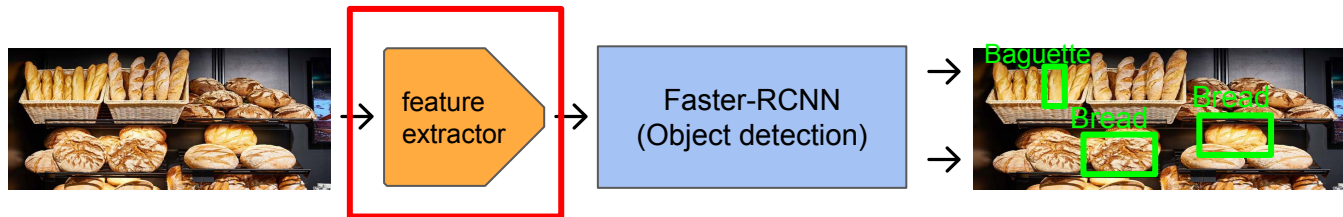


At the end of 12 epochs

Backbone	Training time per epoch
Faster RCNN without FPN	11 hours
Faster RCNN FPN	~ 4 hours



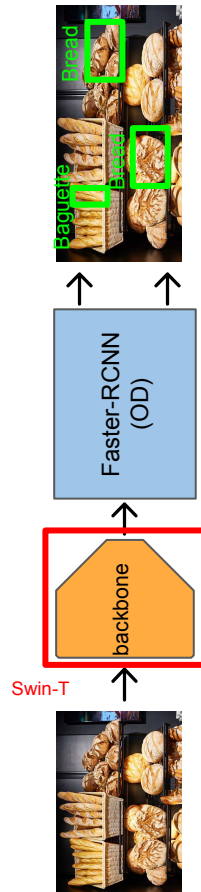
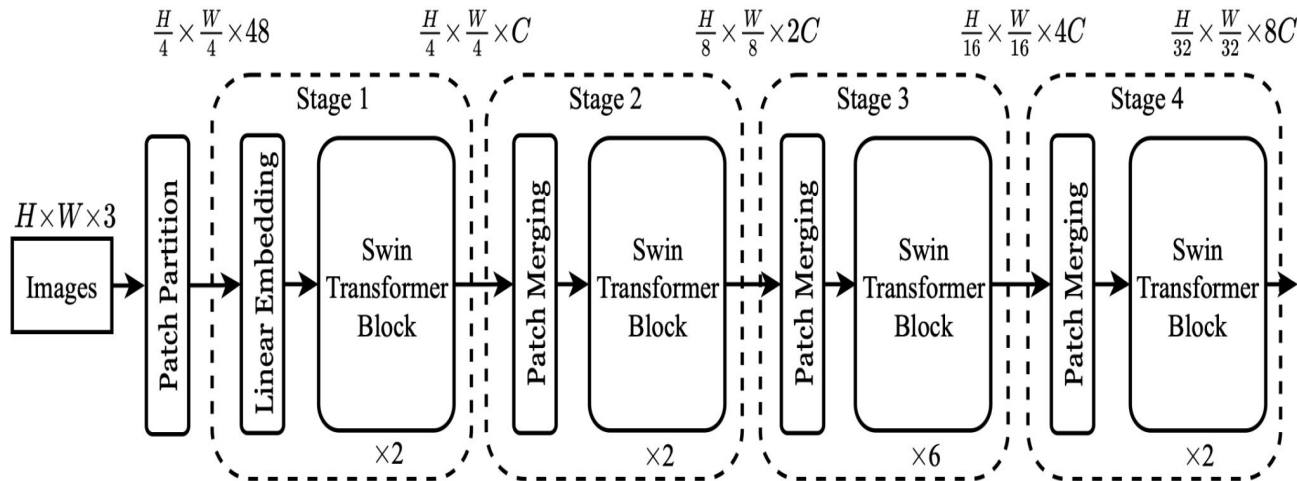
# Overview



**feature extractor initialized with self-supervised pretrained weights**

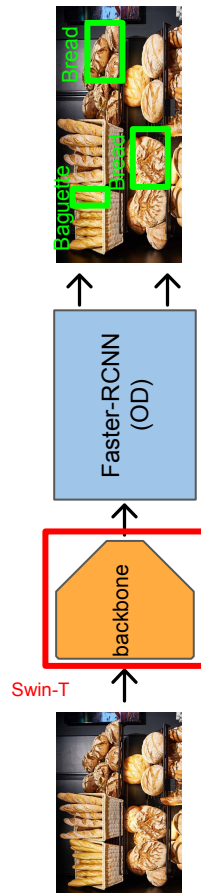
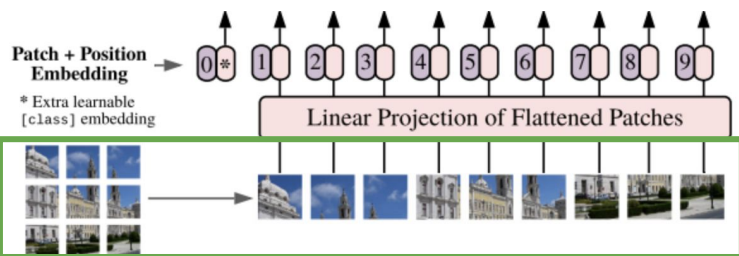
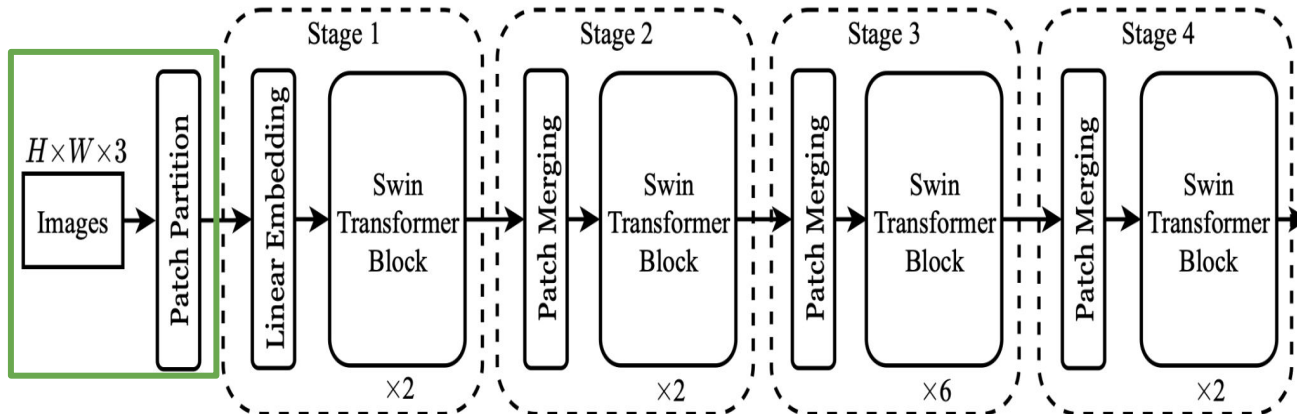


# Swin Transformer

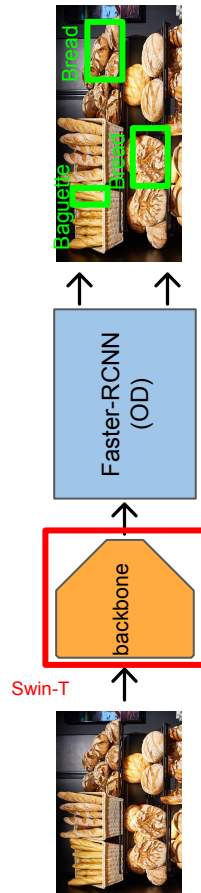
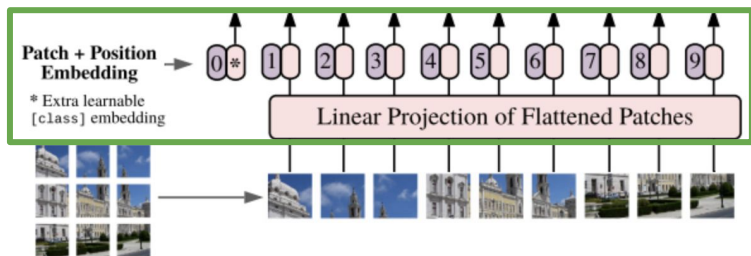
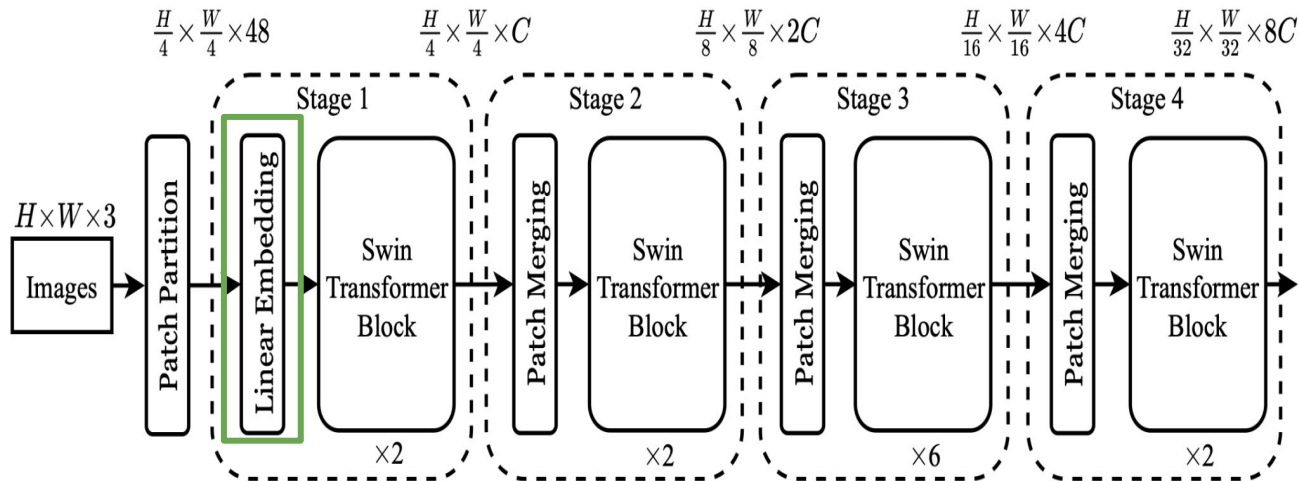


# Swin Transformer

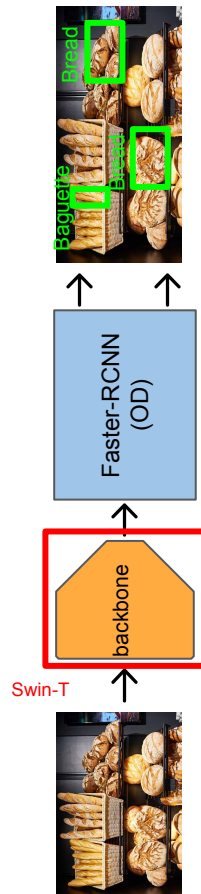
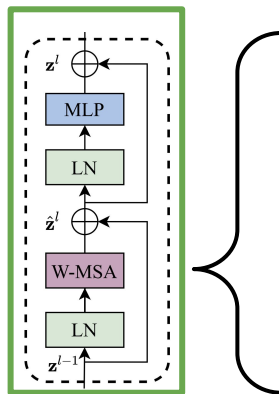
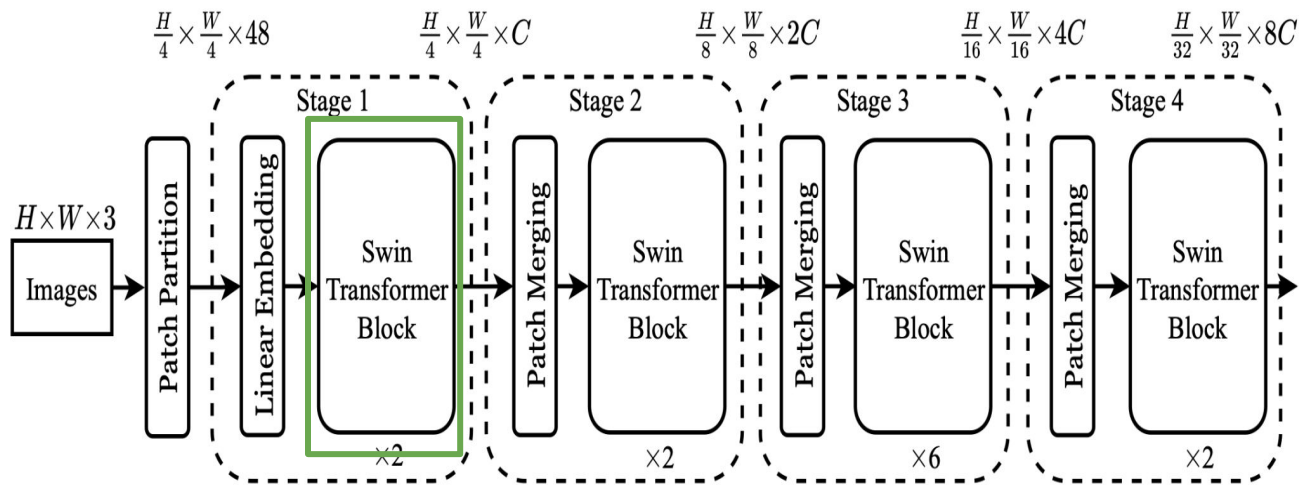
$$\frac{H}{4} \times \frac{W}{4} \times 48 \quad \frac{H}{4} \times \frac{W}{4} \times C \quad \frac{H}{8} \times \frac{W}{8} \times 2C \quad \frac{H}{16} \times \frac{W}{16} \times 4C \quad \frac{H}{32} \times \frac{W}{32} \times 8C$$



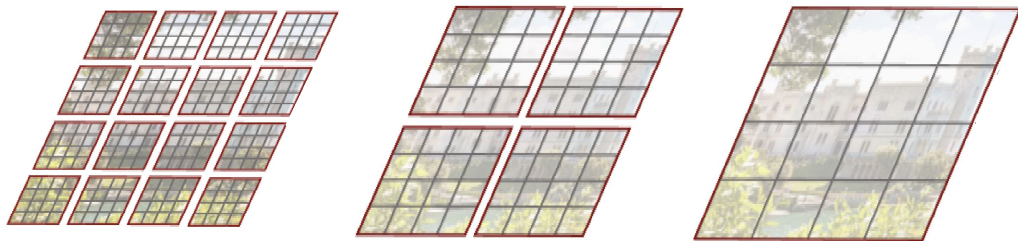
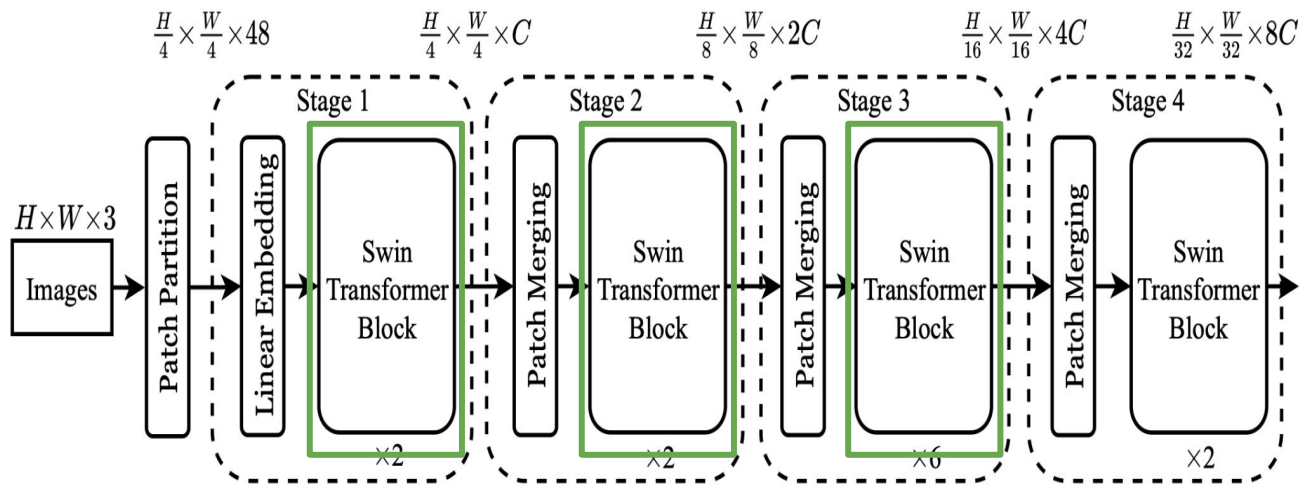
# Swin Transformer



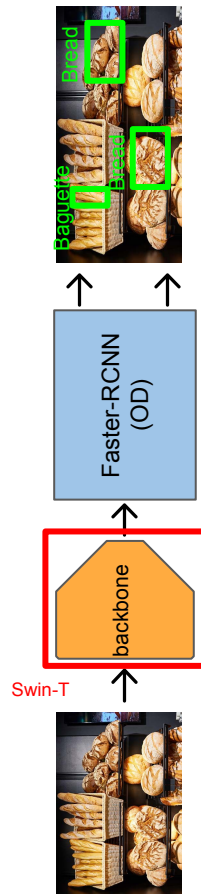
# Swin Transformer



# Swin Transformer

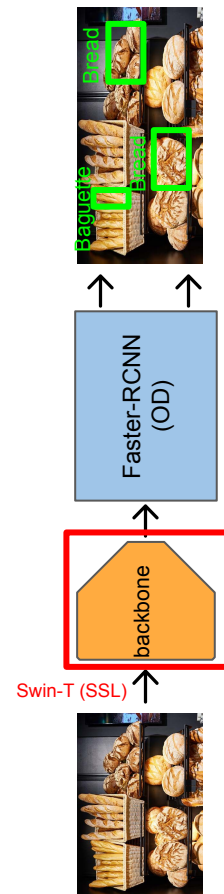
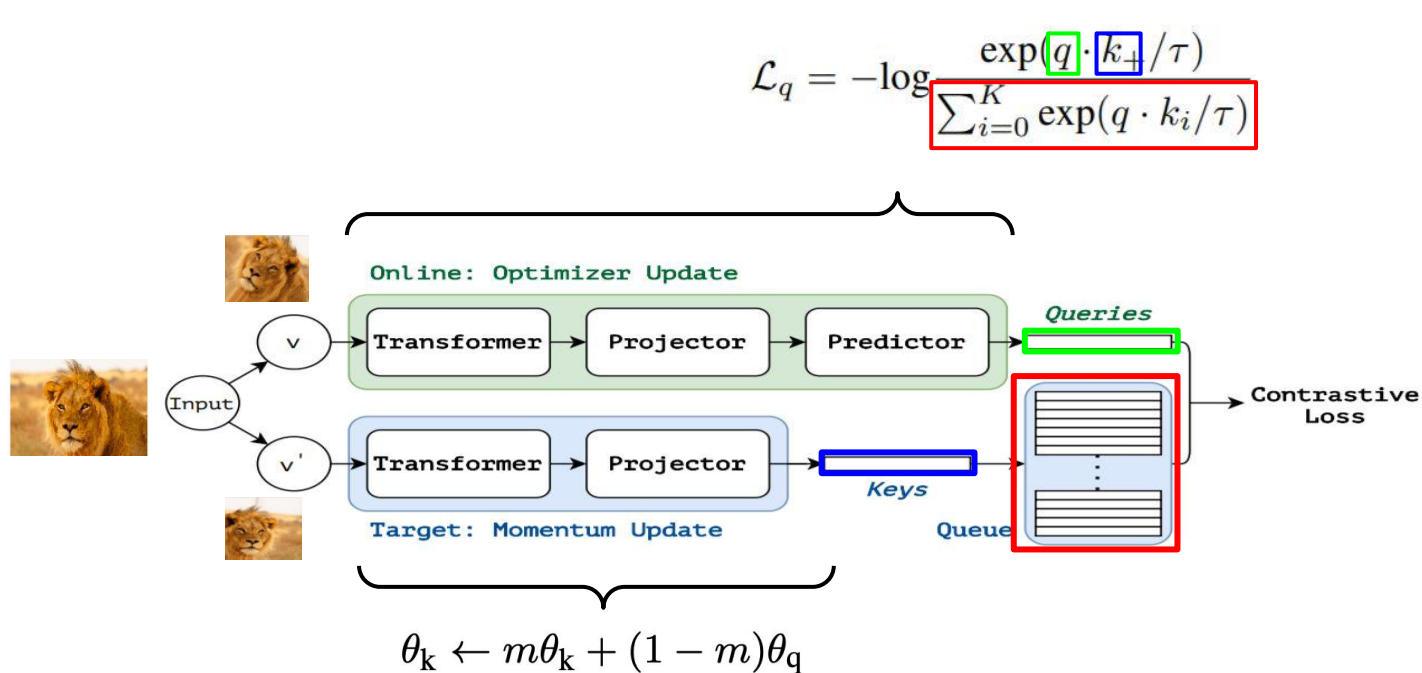


Attention in subsequent stages

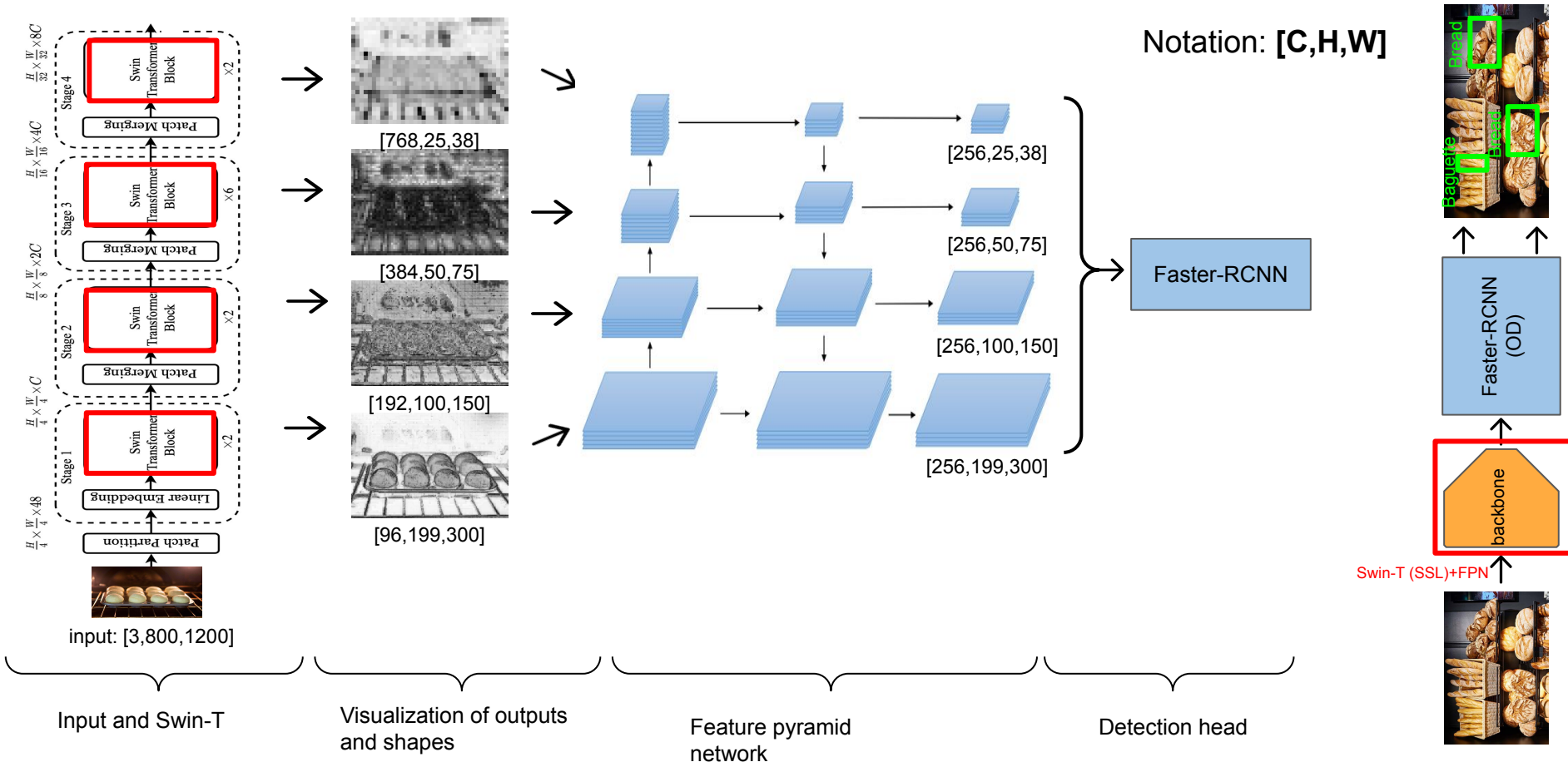




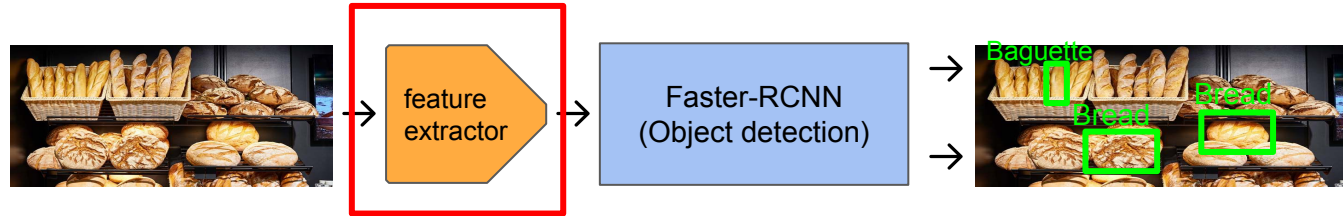
# MoBY - SSL for Swin-T



# Swin Transformer + FPN and Visualization



# Results: SSL in object detection

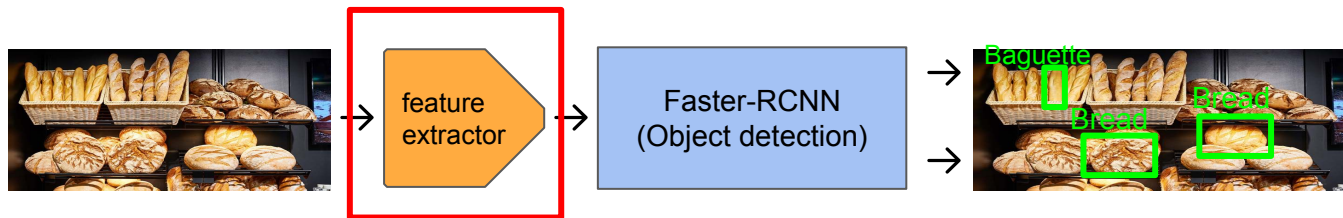


**feature extractor initialized with self-supervised pretrained weights**

Backbone	Initialization	AP
ResNet50 + FPN	supervised	37.93
ResNet50 + FPN	self-supervised (DetCo)	30.18
Swin-T + FPN	supervised	38.89
Swin-T + FPN	self-supervised (MoBY)	38.91

# Overview

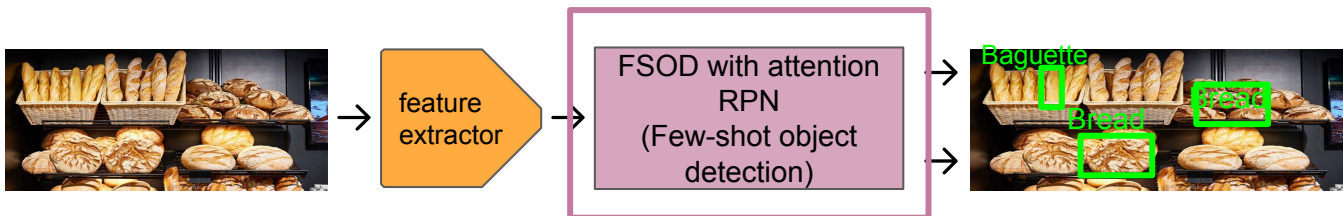
## Step 1: SSL



feature extractor initialized with self-supervised pretrained weights



## Step 2: FSOD

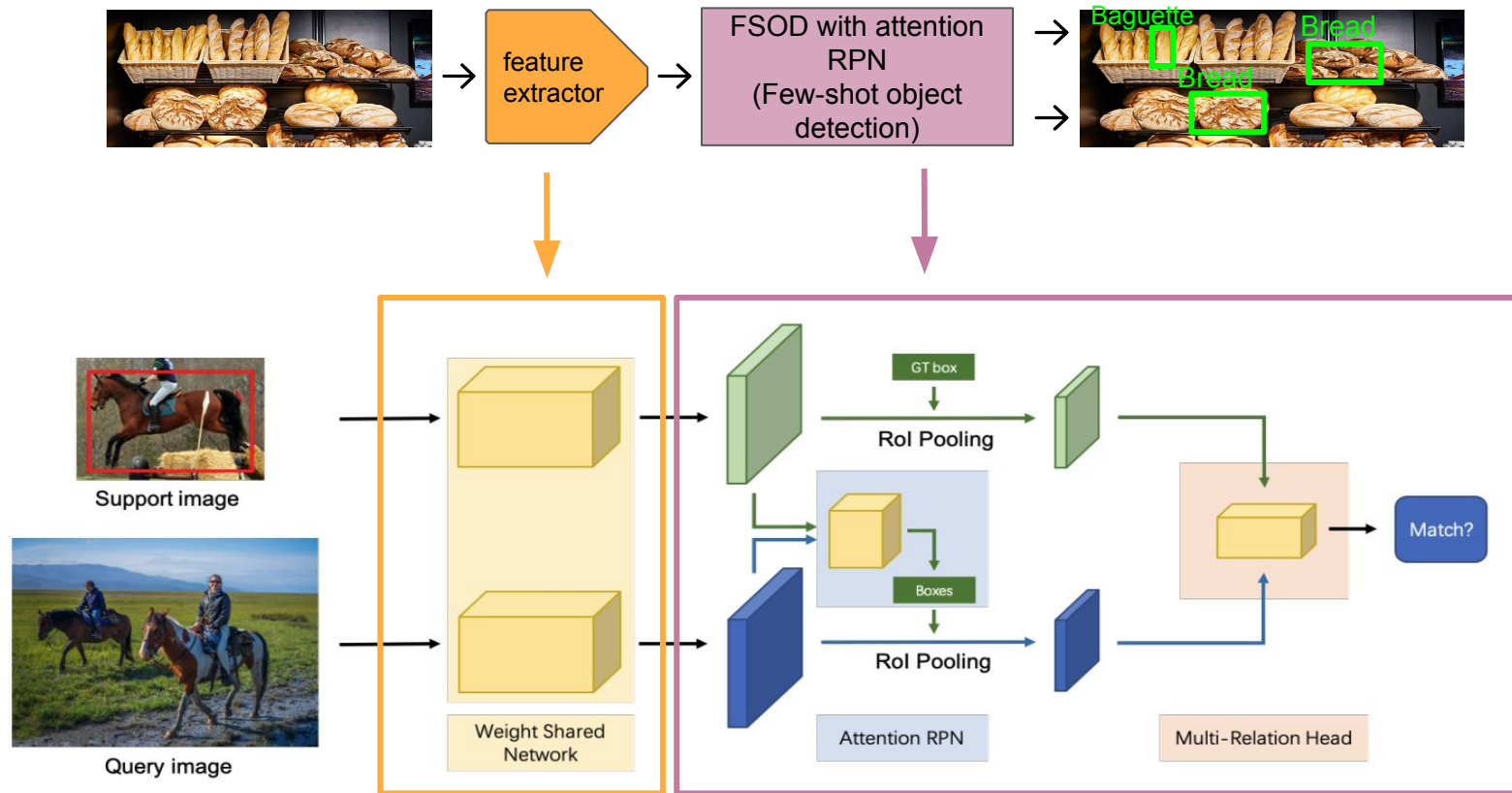


feature extractor initialized self-supervised + change Faster-RCNN to FSOD setup  
(fixed)

# Few-shot object detection

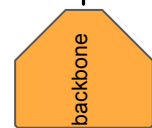
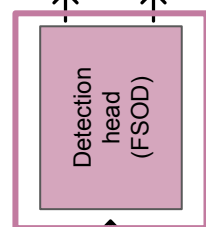
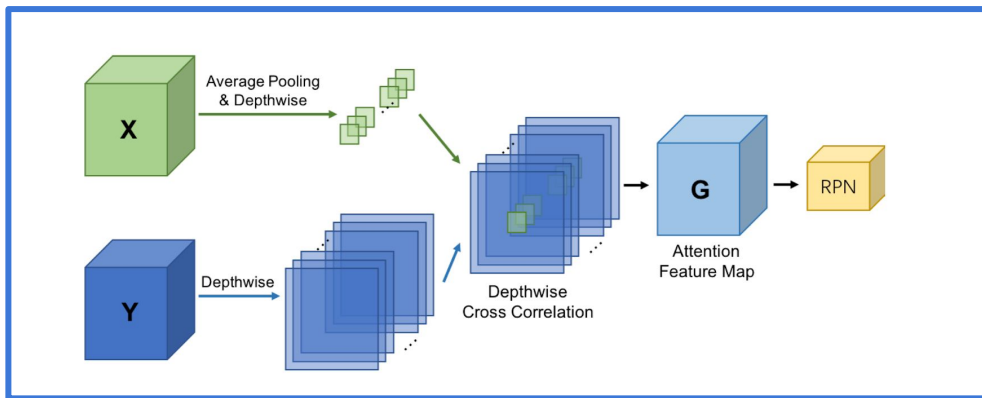
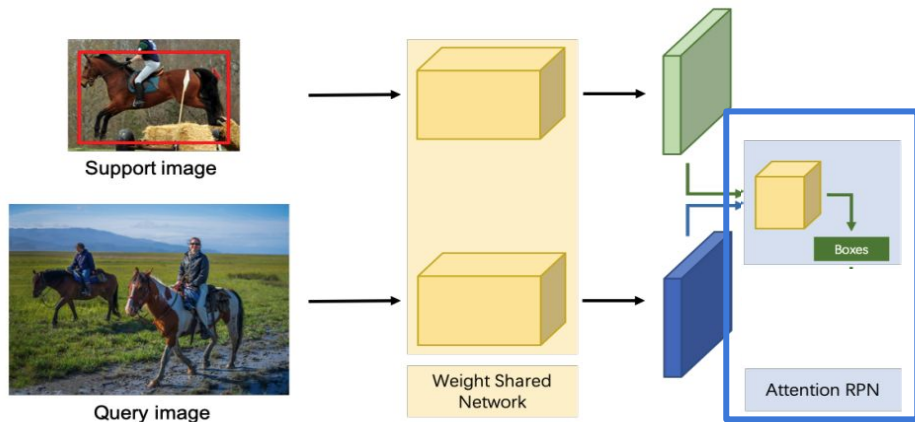
- Given two sets of categories  $\mathcal{C}_{base}$  and  $\mathcal{C}_{novel}$ , where  $\mathcal{C}_{base} \cap \mathcal{C}_{novel} = \emptyset$
- Detect objects of  $\mathcal{C}_{base} \cup \mathcal{C}_{novel}$  by learning from datasets  $\mathcal{D}_{base}$  and  $\mathcal{D}_{novel}$
- $\mathcal{D}_{base}$  contains abundant annotated instances of  $\mathcal{C}_{base}$  (e.g. 500k instances)
- $\mathcal{D}_{novel}$  contains  $K$  annotated instances of  $|\mathcal{C}_{novel}| = N$  categories, where typically  $K \in [1, 30] \cap \mathbb{N}$
- Used model trains only on  $\mathcal{D}_{base}$  and at test time utilizes  $\mathcal{D}_{support} = \mathcal{D}_{novel}$

# FSOD - Detection Head

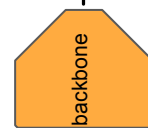
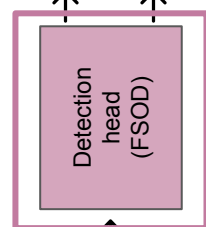
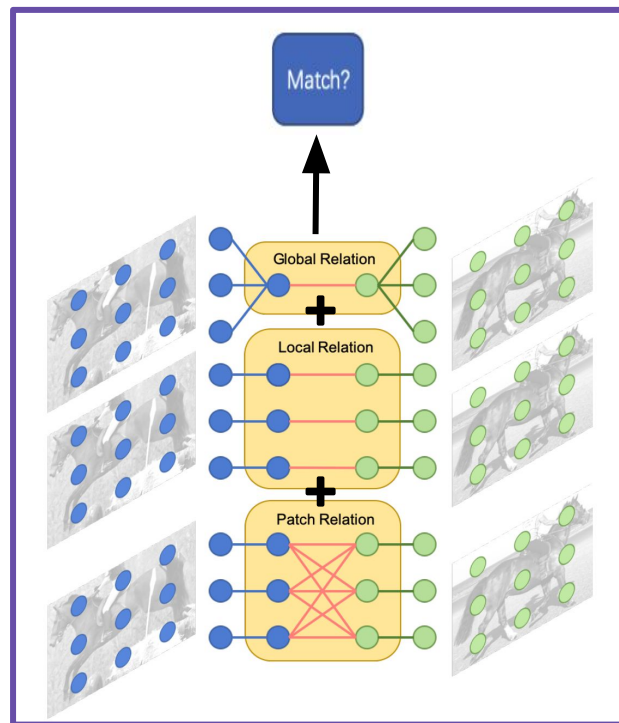
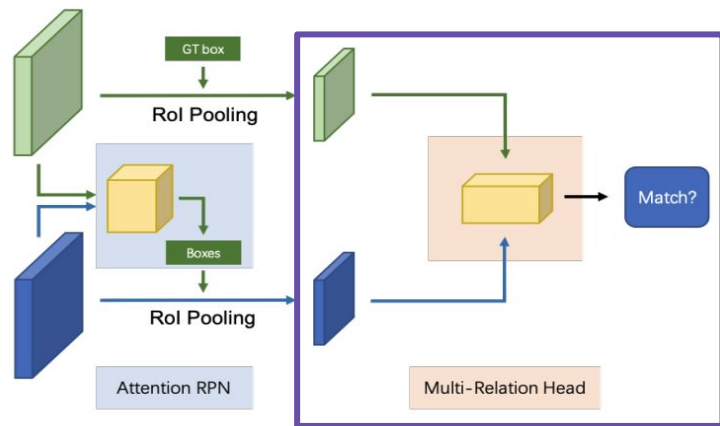




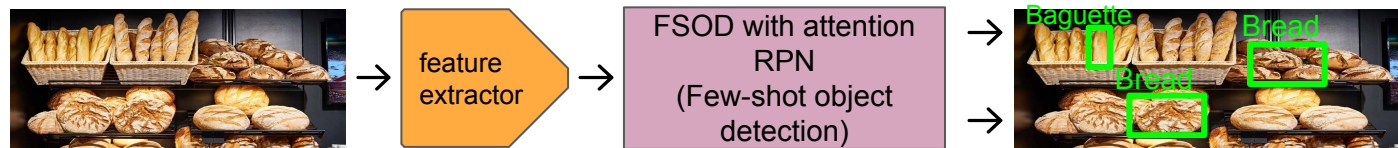
# Attention RPN



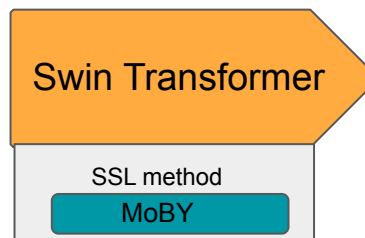
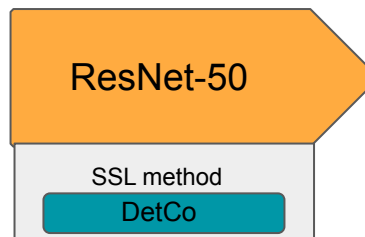
# Rol Pooling



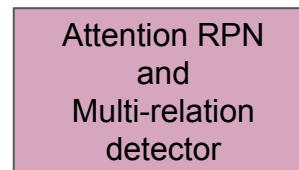
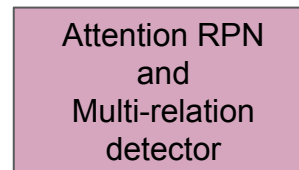
# FSOD - Experiments



## Feature Extractor



## Detection head



- Froze the first 2 stages of each backbone
  - => Low level features do not change during training
  - => Causes poor generalisation to novel Classes
- SwinT with batch size 2 and ResNet50 with batch size 4
  - => gradient has very high variance
  - => need to reduce LR to compensate for it, but then we learn very slow
- Attention RPN is hardcoded w.r.t. Its expected input

- We showed that general object detection can benefit from SSL
  - => Usability of self-supervised method depends on the target task
  - => Could be further improved if SSL is applied to target data
- We tried to assess if FSOD could also benefit from SSL
  - => need more computational resources for efficient research
  - => research in this topic still in early stages
  - => could be used as a label support tool

Thank you!  
Questions?