

Proposal-Contrastive Pretraining for Object Detection from Fewer Data

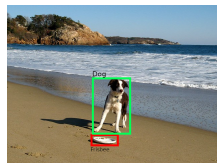
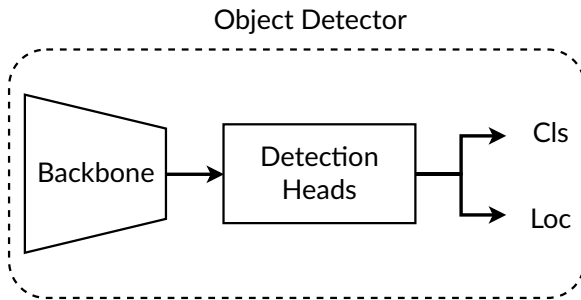
Quentin Bouniot^{1,2} Romaric Audigier¹ Angelique Loesch¹ Amaury Habrard^{2,3}

¹CEA-List

²Université Jean Monnet

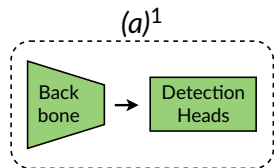
³Institut Universitaire de France





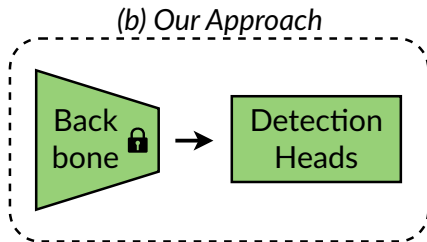
- ▶ Detectors composed of **backbone model** and **detection-specific heads**.
- ▶ Predict **class (Cls)** and **location (Loc)** for each objects in an image.

Overall Pretraining



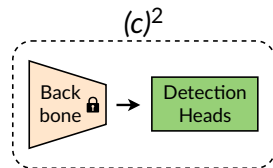
✓ Consistency

✗ Costly



✓ Consistency

✓ Less costly



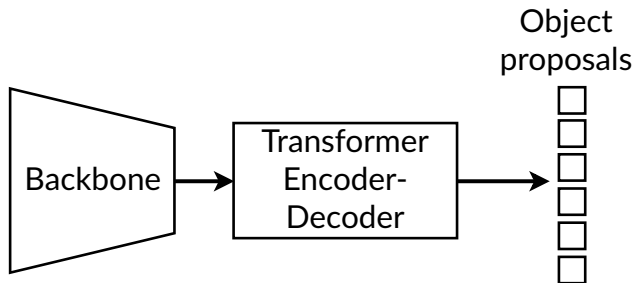
✗ Discrepancy

✓ Less costly

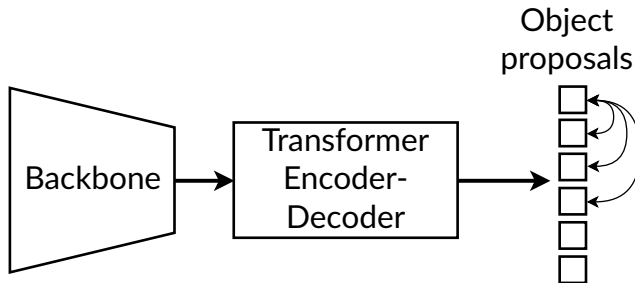
¹Fangyun Wei et al. "Aligning pretraining for detection via object-level contrastive learning". In: *NeurIPS*. 2021

²Zhigang Dai et al. "Up-DETR: Unsupervised pre-training for object detection with transformers". In: *CVPR*. 2021; Amir Bar et al. "Detreg: Unsupervised pretraining with region priors for object detection". In: *CVPR*. 2022

- 1 Context
- 2 Proposal Selection Contrast (ProSeCo)
 - Idea
 - Proposal-Contrastive Learning
 - Avoiding Collapse
- 3 Experimental Results
 - Comparison with state of the art
 - Ablation Studies
- 4 Conclusion



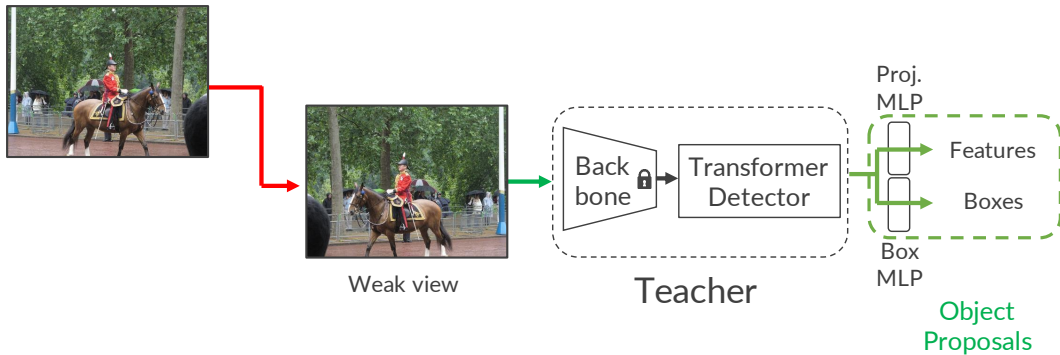
- Transformer-based detectors generates N proposals $\gg k$ objects in images.



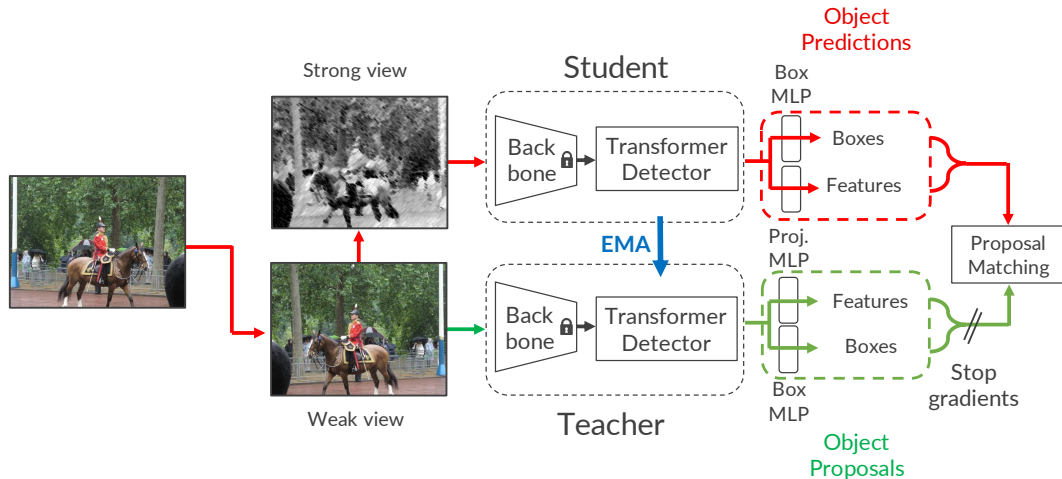
- Transformer-based detectors generates N proposals $\gg k$ objects in images.

Contribution: Contrastive learning **between** proposals.

Proposal-Contrastive Learning



Proposal-Contrastive Learning



- **Object Proposals** from **Teacher** are matched with **Predictions** from **Student**.

Unsupervised Proposal Matching

$$\hat{\sigma}_i^{\text{prop}} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_{j=1}^N \mathcal{L}_{\text{prop_match}}(\underbrace{\mathbf{y}_{(i,j)}}_{\text{Object Proposals}}, \underbrace{\hat{\mathbf{y}}_{(i,\sigma(j))}}_{\text{Object Predictions}})$$

Diagram illustrating the Unsupervised Proposal Matching equation. The equation is $\hat{\sigma}_i^{\text{prop}} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_{j=1}^N \mathcal{L}_{\text{prop_match}}(\mathbf{y}_{(i,j)}, \hat{\mathbf{y}}_{(i,\sigma(j))})$. Annotations include: a green arrow pointing from "Object Proposals" to $\mathbf{y}_{(i,j)}$; a blue arrow pointing from "Permutations of N elements" to $\sigma \in \mathfrak{S}_N$; and a red arrow pointing from "Object Predictions" to $\hat{\mathbf{y}}_{(i,\sigma(j))}$.

- **Proposal** j found by the **teacher** associated to **prediction** $\hat{\sigma}_i^{\text{prop}}(j)$ of the **student**.

Unsupervised Proposal Matching

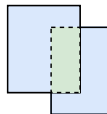
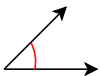
$$\hat{\sigma}_i^{\text{prop}} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_{j=1}^N \mathcal{L}_{\text{prop_match}}(\underbrace{\mathbf{y}_{(i,j)}}_{\text{Object Proposals}}, \underbrace{\hat{\mathbf{y}}_{(i,\sigma(j))}}_{\text{Object Predictions}})$$

Annotations:
 - \mathfrak{S}_N : Permutations of N elements
 - $\mathbf{y}_{(i,j)}$: Object Proposals
 - $\hat{\mathbf{y}}_{(i,\sigma(j))}$: Object Predictions

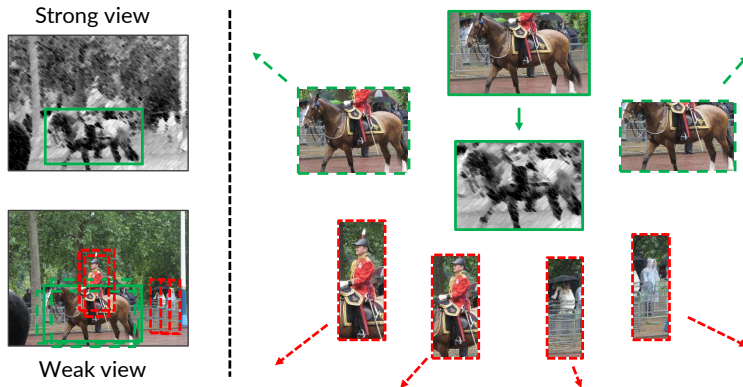
- **Proposal** j found by the **teacher** associated to **prediction** $\hat{\sigma}_i^{\text{prop}}(j)$ of the **student**.

Matching Cost $\mathcal{L}_{\text{prop_match}}$ **depends on:**

- features similarity
- L_1 loss of box coordinates
- generalized IoU loss

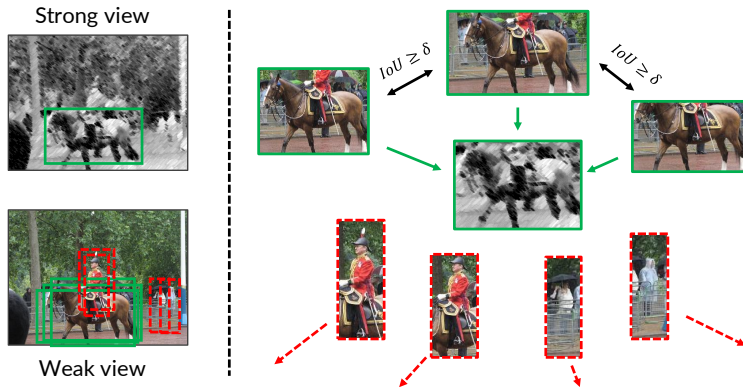


Naive way



× Close proposals considered as **negative** examples.

Localization-aware Contrastive loss



✓ **Overlapping proposals are considered as positive examples.**

Soft Contrastive Estimation (SCE) loss function³

Relations between proposals

Temperature

$$p'_{(in,jm)} = \frac{\mathbb{1}_{i \neq n} \mathbb{1}_{j \neq m} \exp(\mathbf{z}_{(i,j)} \cdot \mathbf{z}_{(n,m)} / \tau_t)}{\sum_{k=1}^{N_b} \sum_{l=1}^N \mathbb{1}_{i \neq k} \mathbb{1}_{j \neq l} \exp(\mathbf{z}_{(i,j)} \cdot \mathbf{z}_{(k,l)} / \tau_t)}$$

Features of Object Proposals

³Julien Denize et al. "Similarity contrastive estimation for self-supervised soft contrastive learning". In: WACV. 2023.

Soft Contrastive Estimation (SCE) loss function³

Relations between proposals

Temperature

Features of Object Proposals

Features of Object Predictions

Contrastive aspect between predictions and proposals

$$p'_{(in,jm)} = \frac{\mathbb{1}_{i \neq n} \mathbb{1}_{j \neq m} \exp(\mathbf{z}_{(i,j)} \cdot \mathbf{z}_{(n,m)} / \tau_t)}{\sum_{k=1}^{N_b} \sum_{l=1}^N \mathbb{1}_{i \neq k} \mathbb{1}_{j \neq l} \exp(\mathbf{z}_{(i,j)} \cdot \mathbf{z}_{(k,l)} / \tau_t)}$$
$$p''_{(in,jm)} = \frac{\exp(\mathbf{z}_{(i,j)} \cdot \hat{\mathbf{z}}_{(n,m)} / \tau)}{\sum_{k=1}^{N_b} \sum_{l=1}^N \exp(\mathbf{z}_{(i,j)} \cdot \hat{\mathbf{z}}_{(k,l)} / \tau)}$$

³Julien Denize et al. "Similarity contrastive estimation for self-supervised soft contrastive learning". In: WACV. 2023.

Localization-aware similarity distribution

$$w_{(in,jm)}^{\text{Loc}} = \lambda_{\text{SCE}} \cdot \mathbb{1}_{i=n} \mathbb{1}_{IoU_i(j,m) \geq \delta} + (1 - \lambda_{\text{SCE}}) \cdot p'_{(in,jm)}$$

IoU between proposals in same image above threshold

Localization-aware similarity distribution

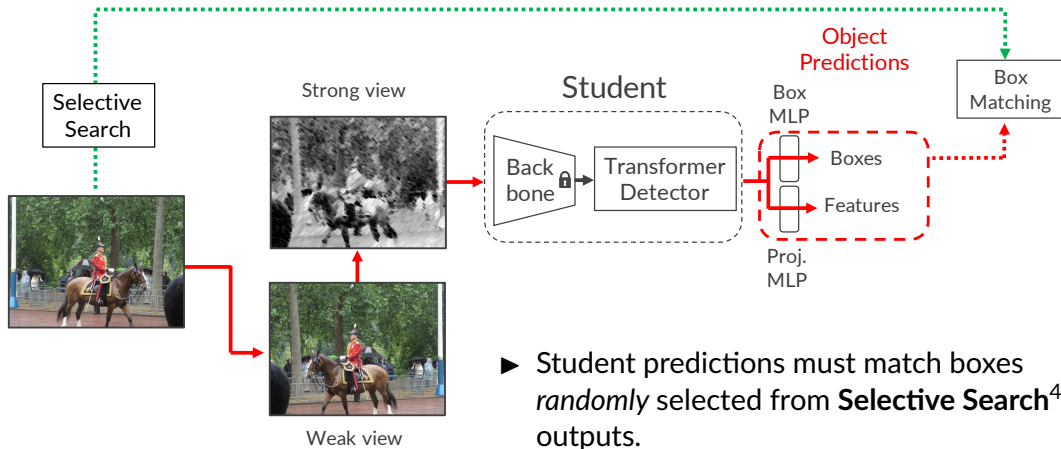
$$w_{(in,jm)}^{\text{Loc}} = \lambda_{\text{SCE}} \cdot \mathbb{1}_{i=n} \mathbb{1}_{IoU_i(j,m) \geq \delta} + (1 - \lambda_{\text{SCE}}) \cdot p'_{(in,jm)}$$

IoU between proposals in same image above threshold

Localized SCE (LocSCE) function

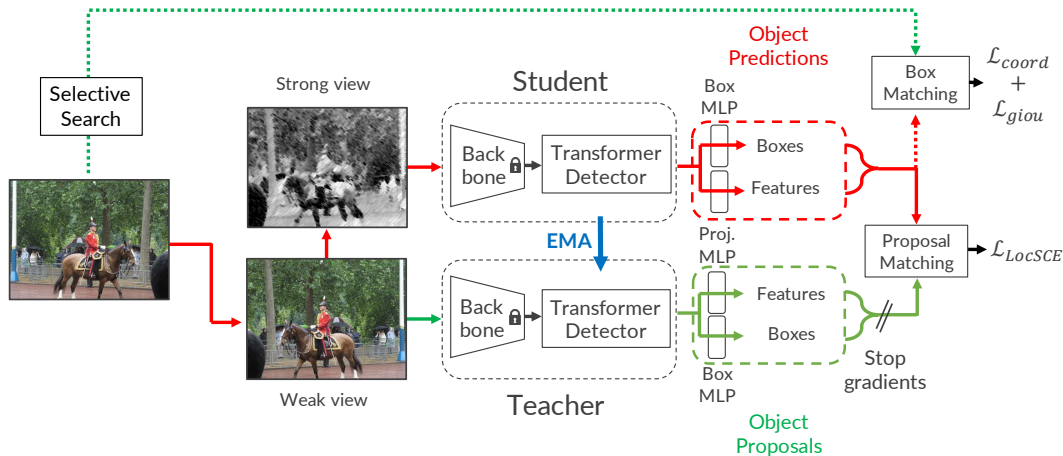
$$\mathcal{L}_{\text{LocSCE}}(\mathbf{y}, \hat{\mathbf{y}}, \hat{\sigma}^{\text{prop}}) = - \frac{1}{N_b N} \sum_{i=1}^{N_b} \sum_{n=1}^{N_b} \sum_{j=1}^N \sum_{m=1}^N w_{(in,jm)}^{\text{Loc}} \log(p''_{(in,j\hat{\sigma}_n^{\text{prop}}(m))})$$

Effective batch size



⁴Jasper RR Uijlings et al. "Selective search for object recognition". In: *IJCV*. 2013.

Proposal Selection Contrast (ProSeCo)



- Full pretraining procedure with both **contrastive** and **localization** learning.

Pretraining	Detector	Mini-COCO		
		1% (1.2k)	5% (5.9k)	10% (11.8k)
Supervised	Def. DETR	13.0	23.6	28.6
SwAV ⁵	Def. DETR	13.3	24.5	29.5
SCRL ⁶	Def. DETR	16.4	26.2	30.6
DETR ⁷	Def. DETR	15.9	26.1	30.9
Supervised	Mask R-CNN	–	19.4	24.7
SoCo* ⁸	Mask R-CNN	–	26.8	31.1
<i>ProSeCo (Ours)</i>	Def. DETR	18.0	28.8	32.8

⁵Mathilde Caron et al. "Unsupervised learning of visual features by contrasting cluster assignments". In: *NeurIPS*. 2020.

⁶Byungseok Roh et al. "Spatially consistent representation learning". In: *CVPR*. 2021.

⁷Amir Bar et al. "Detreg: Unsupervised pretraining with region priors for object detection". In: *CVPR*. 2022.

⁸Fangyun Wei et al. "Aligning pretraining for detection via object-level contrastive learning". In: *NeurIPS*. 2021.

Method	FSOD-test	FSOD-train	PASCAL VOC	Mini-VOC	
	100% (11k)	100% (42k)	100% (16k)	5% (0.8k)	10% (1.6k)
Supervised	39.3	42.6	59.5	33.9	40.8
DETR ⁹	43.2	43.3	63.5	43.1	48.2
<i>ProSeCo (Ours)</i>	46.6	47.2	65.1	46.1	51.3

- ✓ **ProSeCo** improves over SOTA on all datasets considered, with **various amount** of labeled data.

⁹ Amir Bar et al. "Detreg: Unsupervised pretraining with region priors for object detection". In: CVPR. 2022.

Pretraining	Dataset	mAP
ProSeCo w/ SwAV	COCO	27.4
ProSeCo w/ SwAV	IN	27.8
DETRReg w/ SCRL	IN	28.0
ProSeCo w/ SCRL	IN	28.8

Loss	δ	mAP
SCE	1.0	26.1
<i>LocSCE (Ours)</i>	0.2	27.0
<i>LocSCE (Ours)</i>	0.7	27.1
<i>LocSCE (Ours)</i>	0.5	27.8

- **Dataset diversity** more important than closeness to downstream task
- ✓ **Consistency** in the features improves performance
- ✓ **Location of proposals** helps for introducing **easy positives** for contrastive learning

We propose ProSeCo, a Proposal-Contrastive Pretraining strategy for Object Detection with Transformers.







- ✓ Leverage high number of Object Proposals for **Proposal-Contrastive Learning**.
- ✓ Our **ProSeCo improves performance** when training with limited labeled data.
- ✓ **Consistency** with object-level features is important for Object Detection.
- ✓ **Location information** helps for Proposal-Contrastive learning.



Thank You !

Do not hesitate to contact us for question !

Bouniot et al., "Proposal-Contrastive Pretraining for Object Detection from Fewer Data"



-  Fangyun Wei et al. “Aligning pretraining for detection via object-level contrastive learning”. In: *NeurIPS*. 2021.
-  Zhigang Dai et al. “Up-DETR: Unsupervised pre-training for object detection with transformers”. In: *CVPR*. 2021.
-  Amir Bar et al. “Detreg: Unsupervised pretraining with region priors for object detection”. In: *CVPR*. 2022.
-  Julien Denize et al. “Similarity contrastive estimation for self-supervised soft contrastive learning”. In: *WACV*. 2023.
-  Jasper RR Uijlings et al. “Selective search for object recognition”. In: *IJCV*. 2013.
-  Mathilde Caron et al. “Unsupervised learning of visual features by contrasting cluster assignments”. In: *NeurIPS*. 2020.

-  Byungseok Roh et al. “Spatially consistent representation learning”. In: *CVPR*. 2021.
-  Quentin Bouniot et al. “Proposal-Contrastive Pretraining for Object Detection from Fewer Data”. In: *ICLR*. 2023.