

JUNE 18-22, 2023

CVPR

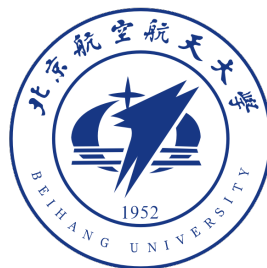


Siamese DETR

Zeren Chen^{1,2}, Gengshi Huang², Wei Li³, Jianing Teng², KunWang²,
Jing Shao², Chen Change Loy³, Lu Sheng¹

¹Beihang University, ²SenseTime Research, ³Nanyang Technological University

Tag: WED-PM-321



NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE

S-LAB
FOR ADVANCED
INTELLIGENCE



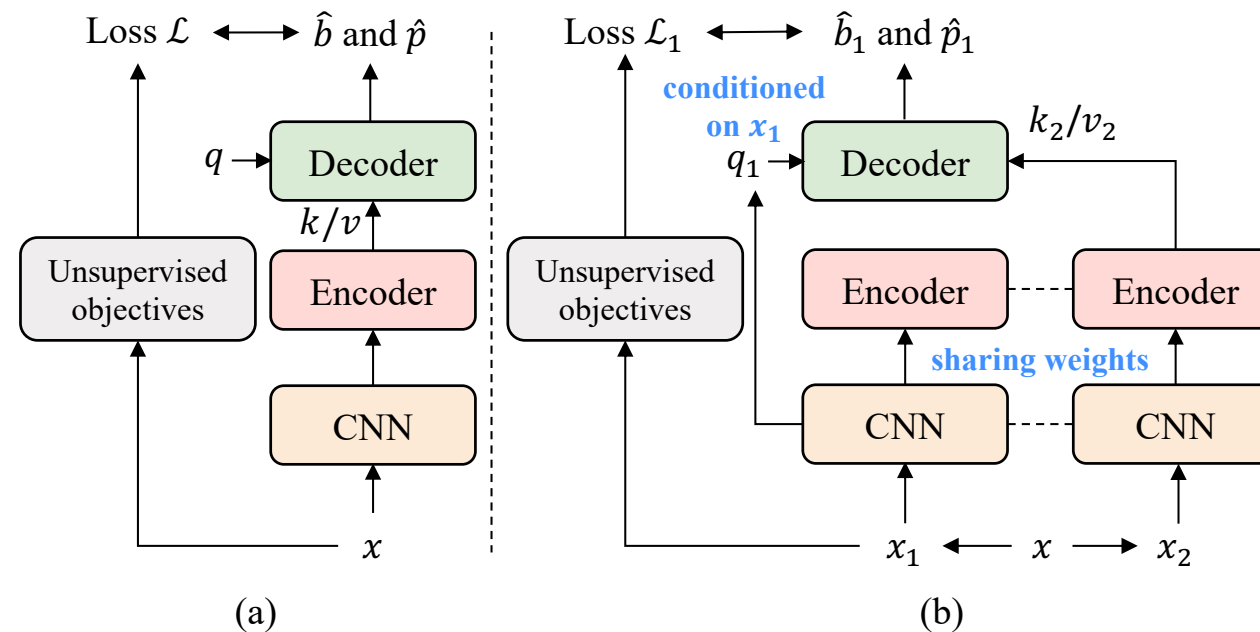
Code is available at <https://github.com/Zx55/SiameseDETR>

Overview

- Combine Siamese networks with cross-attention mechanism in DETR.
- Two newly-designed self-supervised pretext tasks.
 - **Multi-View Region Detection**
 - **Multi-View Semantic Discrimination**
- Siamese DETR outperforms its counterpart with multiple DETR variants on the COCO and PASCAL VOC benchmark.

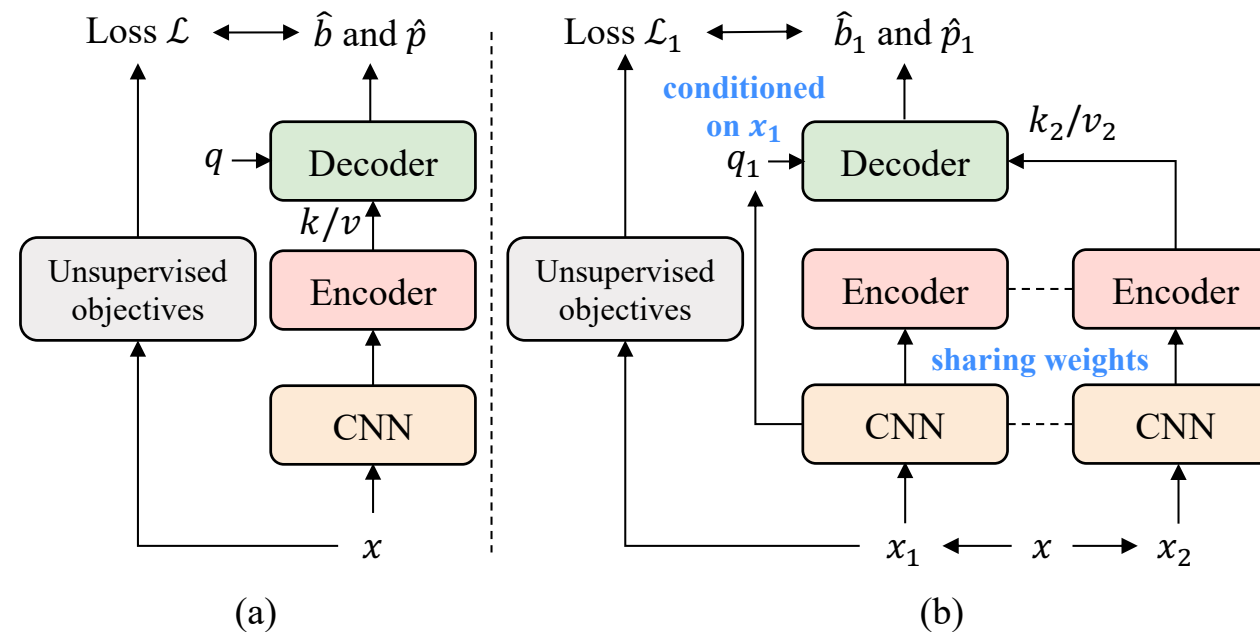
Introduction & Motivation

- Objective: Design a self-supervised learning approach for DETR pretraining to alleviate the massive appetite for labeled data in training DETR.
- Existing self-supervised learning approaches cannot be extended to DETR effectively (e.g., SimCLR, MoCo).



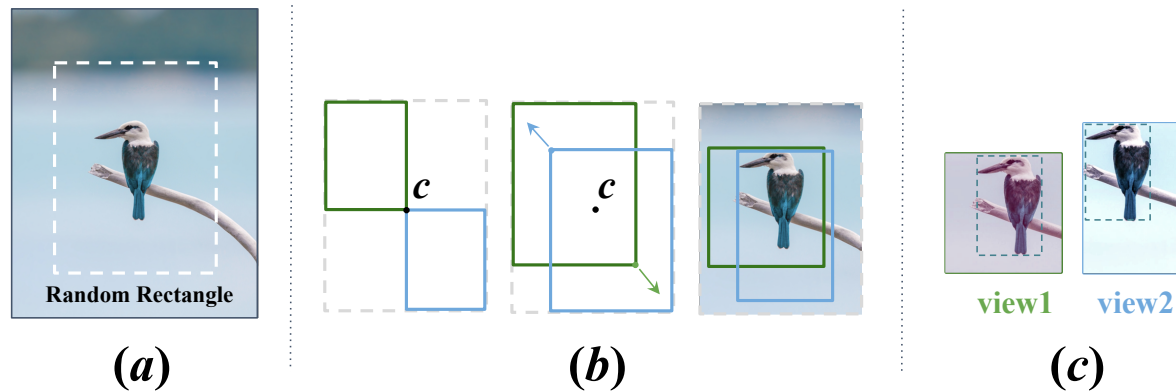
Introduction & Motivation

- Several recent attempts (e.g., UP-DETR, DETReg) follows a single-view paradigm (see a), ignoring the ability of learning view-invariant representation.
- Siamese DETR learn view-invariant and detection-oriented representation through two pretext tasks (see b).



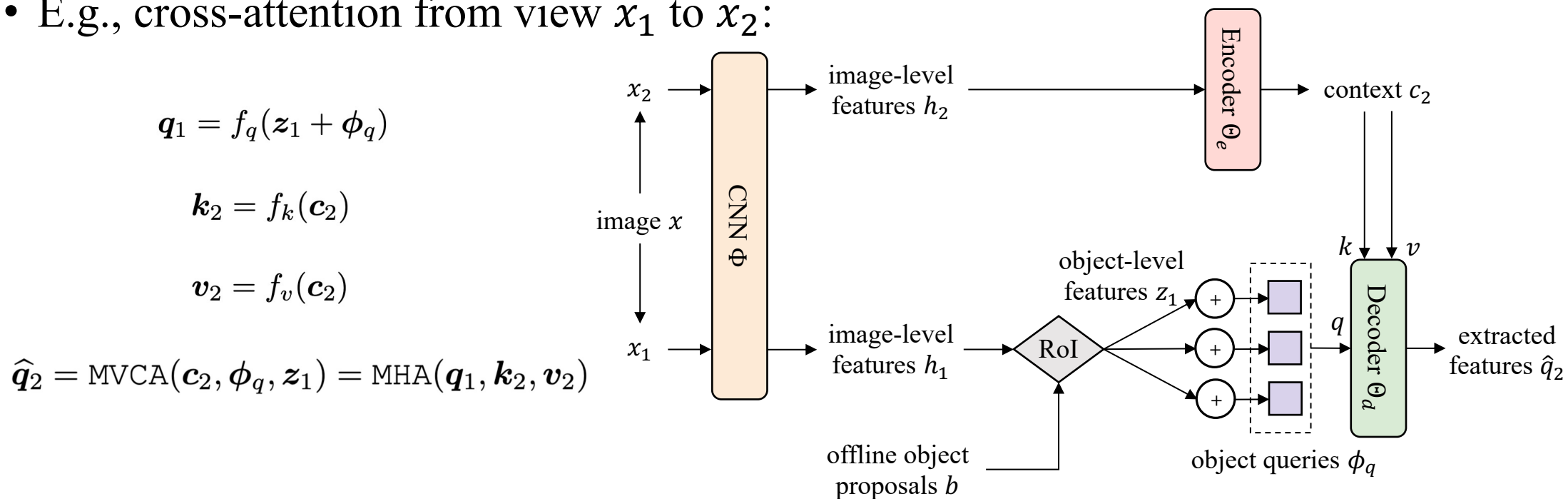
View Construction

- Generate two views based on an IoU-constrained policy.
 - Generate two rectangles and keep their IoU larger than a threshold.
 - Crop two rectangles and apply augmentations (following SimSiam).
 - Generate offline object proposals with Edgeboxes in the overlapping area.



Cross-View Learning

- We propose Multi-View Cross-Attention for multi-view representation learning.
- E.g., cross-attention from view x_1 to x_2 :



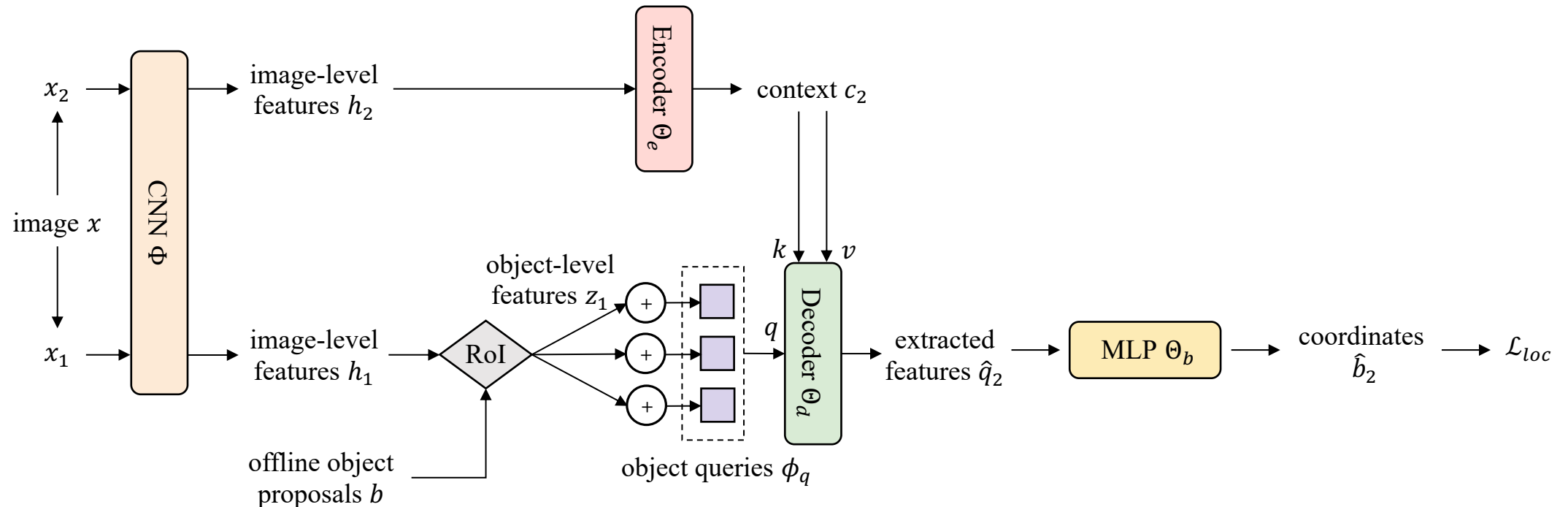
Here, $\hat{\mathbf{q}}_2$ are supposed to be semantically consistent with the corresponding region features \mathbf{z}_1 on view x_1 .

Learning to Locate

- Locate the region in view x_2 that is relative to the region feature z_1 .

$$\hat{\mathbf{b}}_2 = f_{box}(\hat{\mathbf{q}}_2) \in \mathbb{R}^{N \times 4}$$

$$\mathcal{L}_{loc} = \ell_{box}(\hat{\mathbf{b}}_2, \mathbf{b}_2) + \ell_{box}(\hat{\mathbf{b}}_1, \mathbf{b}_1)$$



Learning to Discriminate

- Global discrimination

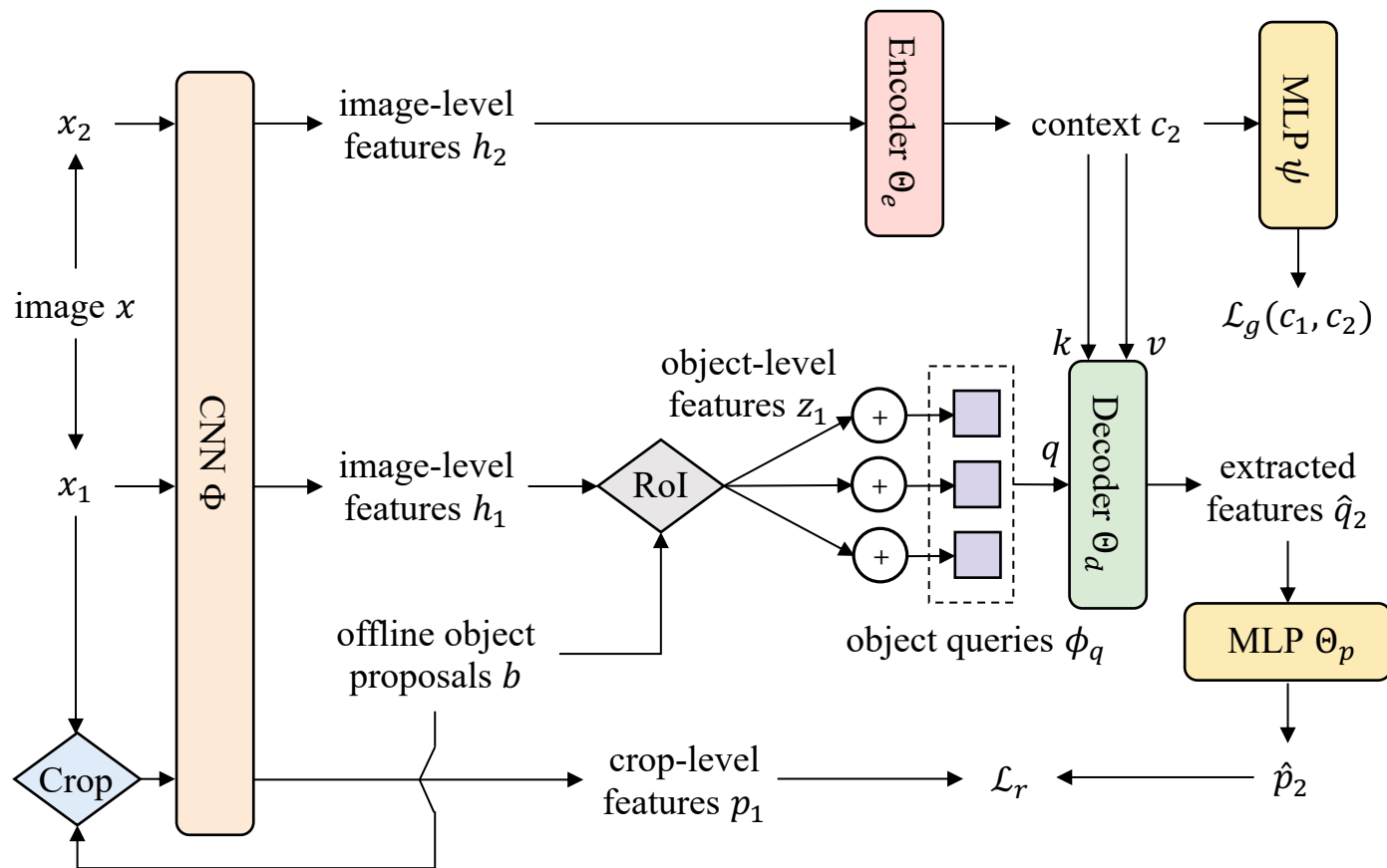
$$\mathcal{L}_g = \mathcal{C}[\text{MLP}(\mathbf{c}_1), \text{detach}(\mathbf{c}_2)] + \mathcal{C}[\text{MLP}(\mathbf{c}_2), \text{detach}(\mathbf{c}_1)]$$

- Regional discrimination

$$\mathbf{p}_1 = \text{Backbone}(\text{Crop}(\mathbf{x}_1, \mathbf{b}_1))$$

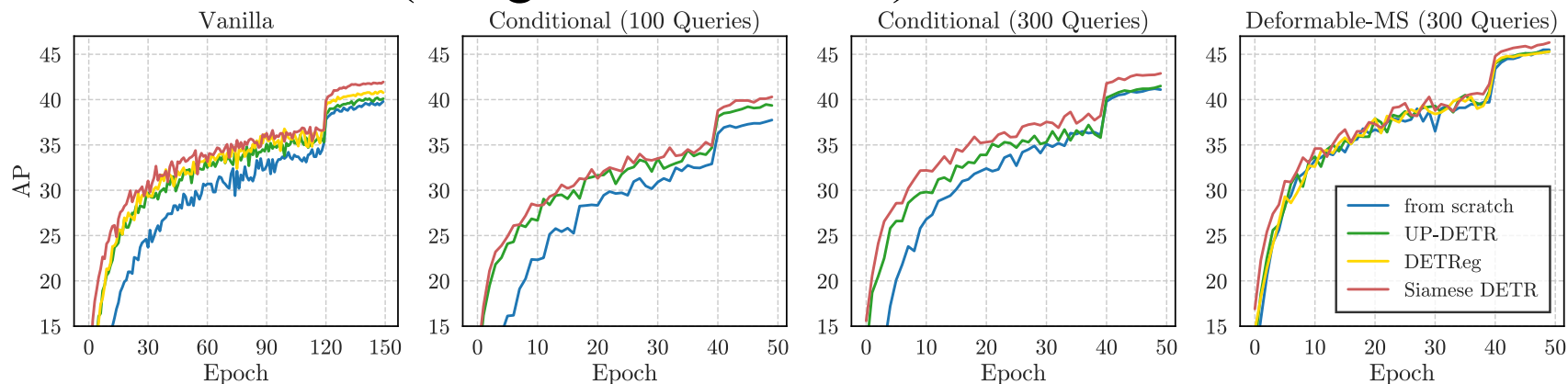
$$\hat{\mathbf{p}}_2 = f_{sem}(\hat{\mathbf{q}}_2) \in \mathbb{R}^{N \times C'}$$

$$\mathcal{L}_r = \mathcal{D}(\hat{\mathbf{p}}_2, \mathbf{p}_1) + \mathcal{D}(\hat{\mathbf{p}}_1, \mathbf{p}_2)$$

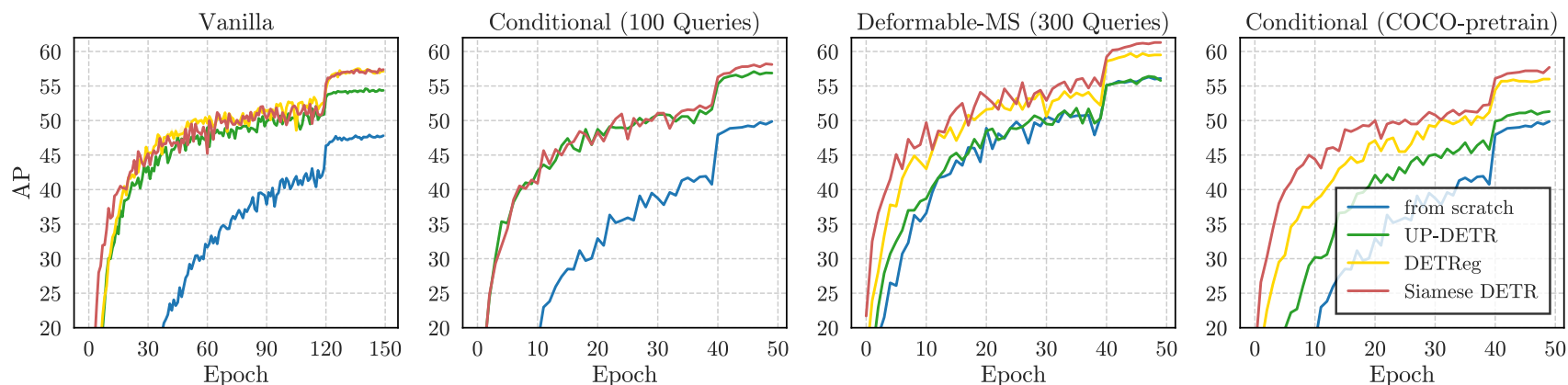


Experiments

- MS-COCO benchmark (ImageNet \rightarrow COCO)



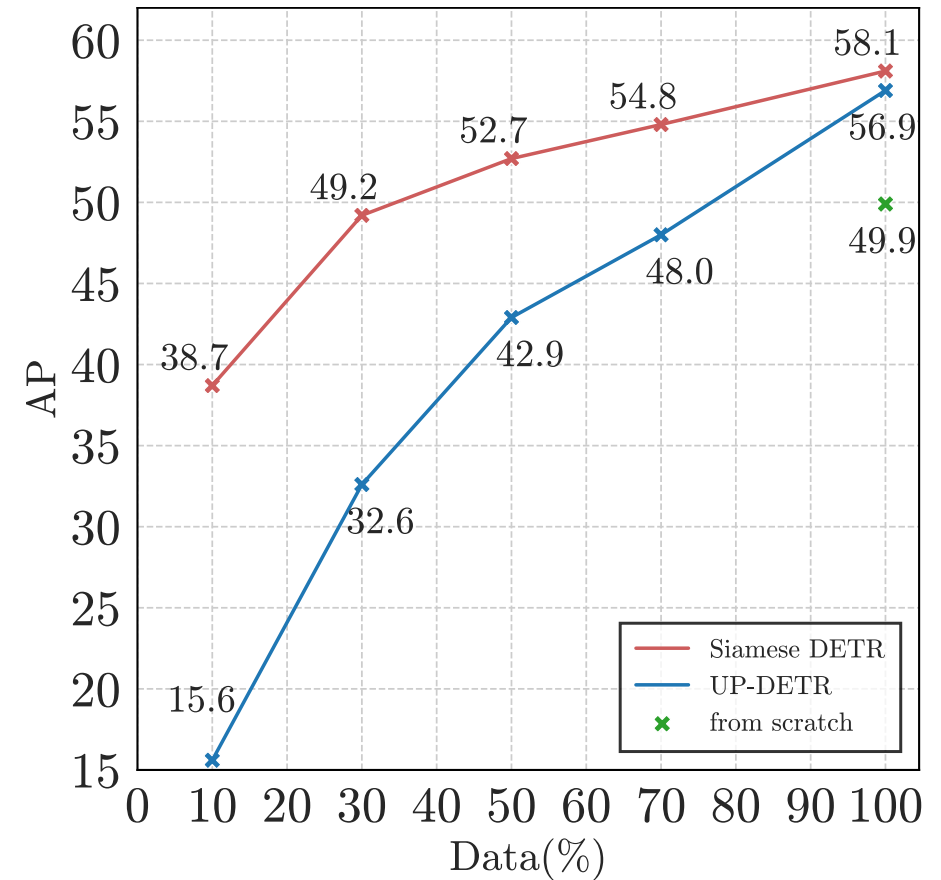
- PASCAL VOC benchmark (ImageNet \rightarrow VOC / COCO \rightarrow VOC)



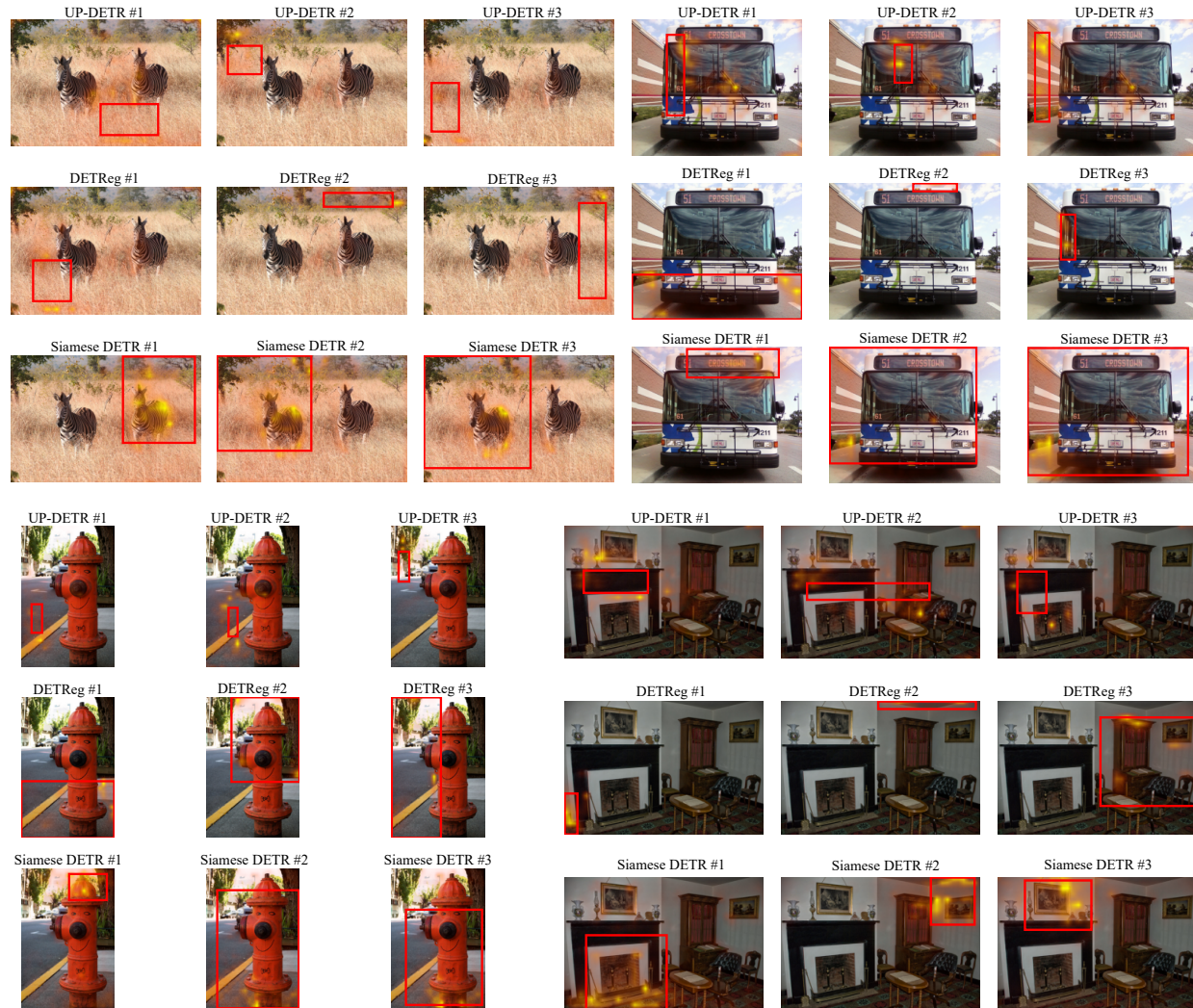
Experiments

Method	Dataset	Proposals	AP
UP-DETR	ImageNet→COCO	Random	39.4
DETRReg	ImageNet→COCO	Random	40.3
ours	ImageNet→COCO	Random	40.4
UP-DETR	ImageNet→COCO	Edgeboxes	39.3
DETRReg	ImageNet→COCO	Edgeboxes	40.3
ours	ImageNet→COCO	Edgeboxes	40.5
UP-DETR	COCO→VOC	Random	51.3
DETRReg	COCO→VOC	Random	51.9
ours	COCO→VOC	Random	54.9
DETRReg	COCO→VOC	SelectiveSearch	55.9
ours	COCO→VOC	SelectiveSearch	56.2
UP-DETR	COCO→VOC	Edgeboxes	57.0
DETRReg	COCO→VOC	Edgeboxes	56.3
ours	COCO→VOC	Edgeboxes	57.7

DETR	VOC		VOC 10%	
	DAB-DETR	DN-DETR	DAB-DETR	DN-DETR
<i>from scratch</i>	57.9	58.9	32.2	32.9
Siamese DETR	62.2 (+4.3)	63.4 (+4.5)	41.8 (+9.6)	43.6 (+10.7)



Visualization



Thanks for your watching!