

Siamese DETR

Zeren Chen^{1,2*}, Gengshi Huang^{2*}, Wei Li³, Jianing Teng², Kun Wang²,
Jing Shao², Chen Change Loy³, Lu Sheng^{1†}
¹ School of Software, Beihang University, ² SenseTime Research,
³ S-Lab, Nanyang Technological University.

{czr1604, lsheng}@buaa.edu.cn, {wei.l, ccloy}@ntu.edu.sg, huanggengshi@gmail.com
wangkun@sensetime.com, {tengjianing, shaojing}@senseauto.com.

Abstract

Recent self-supervised methods are mainly designed for representation learning with the base model, e.g., ResNets or ViTs. They cannot be easily transferred to DETR, with task-specific Transformer modules. In this work, we present **Siamese DETR**, a **Siamese** self-supervised pretraining approach for the Transformer architecture in **DETR**. We consider learning view-invariant and detection-oriented representations simultaneously through two complementary tasks, i.e., localization and discrimination, in a novel multi-view learning framework. Two self-supervised pretext tasks are designed: (i) **Multi-View Region Detection** aims at learning to localize regions-of-interest between augmented views of the input, and (ii) **Multi-View Semantic Discrimination** attempts to improve object-level discrimination for each region. The proposed **Siamese DETR** achieves state-of-the-art transfer performance on COCO and PASCAL VOC detection using different DETR variants in all setups. Code is available at <https://github.com/Zx55/SiameseDETR>.

1. Introduction

Object detection with Transformers (DETR) [3] combines convolutional neural networks (CNNs) and Transformer-based encoder-decoders, viewing object detection as an end-to-end set prediction problem. Despite its impressive performance, DETR and its variants still rely on large-scale, high-quality training data. It generally requires huge cost and effort to collect such massive well-annotated datasets, which can be prohibited in some privacy-sensitive applications such as medical imaging and video surveillance.

Recent progress in multi-view self-supervised represen-

*Equal contribution.

†Corresponding author.

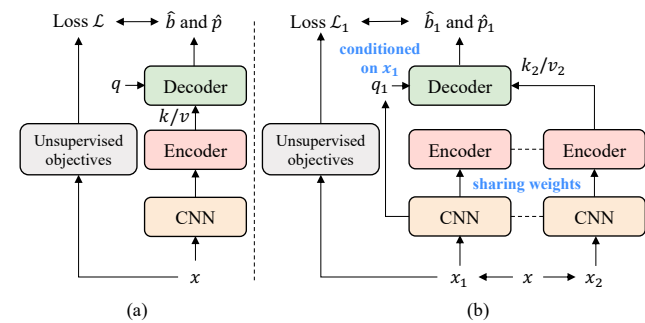


Figure 1. Comparison between single-view and multi-view detection pretraining for DETR. (a) The single-view framework, e.g., UP-DETR [9] and DETReg [1], perform self-supervised representation learning using unsupervised objectives generated on the single view, e.g., random patches (UP-DETR) or pseudo labels (DETRReg), leading to a small information gain during pretraining. (b) The proposed multi-view Siamese DETR for DETR pretraining. Here, \hat{b} and \hat{p} denote box and semantic predictions. q , k and v denote query, key and value in DETR, respectively.

tation learning [4, 6–8, 14, 15, 21] can potentially alleviate the appetite for labeled data in training DETR for object detection. However, these self-supervised learning approaches mainly focus on learning generalizable representations with base models, such as ResNets [17] and ViTs [11]. It is unclear how these approaches can be effectively extended to DETR with task-specific Transformers modules that are tailored for end-to-end object detection.

Designing self-supervised pretext tasks for pretraining the Transformers in DETR is a challenging and practical problem, demanding representations that could benefit object detection, beyond just learning generic representation. Several attempts have been made to address this issue. For example, UP-DETR [9] introduces an unsupervised pretext task based on random query patch detection, predicting bounding boxes of randomly-cropped query patches in the

given image. Recent DETReg [1] employs a pre-trained SwAV [5] and offline Selective Search proposals [28] to provide pseudo labels for DETR pertaining. In general, both UP-DETR and DETReg follow a single-view pretraining paradigm (see Figure 1 (a)), without exploring the ability of learning view-invariant representations demonstrated in existing multi-view self-supervised approaches.

In this work, we are interested in investigating the effectiveness of multi-view self-supervised learning for DETR pre-training. Different from conventional multi-view framework [5, 6, 15], we combine the Siamese network with the cross-attention mechanism in DETR, presenting a Siamese self-supervised pretraining approach, named Siamese DETR, with two proposed self-supervised pretext tasks dedicated to view-invariant detection pretraining. Specifically, given each unlabeled image, we follow [1, 31] to obtain the offline object proposals and generate two augmented views guided by Intersection over Union (IoU) thresholds. As illustrated in Figure 1 (b), by directly locating the query regions between augmented views and maximizing the discriminative information at both global and regional levels, Siamese DETR can learn view-invariant representations with localization and discrimination that are aligned with downstream object detection tasks during pre-training. Our contributions can be summarized as below:

- We propose a novel Siamese self-supervised approach for the Transformers in DETR, which jointly learns view-invariant representations with discrimination and localization. In particular, we contribute two new designs of self-supervised pretext tasks specialized for multi-view detection pretraining.
- Without bells and whistles, Siamese DETR outperforms UP-DETR [9] and DETReg [1] with multiple DETR variants, such as Conditional [26] and Deformable [38] DETR, on the COCO and PASCAL VOC benchmarks, demonstrating the effectiveness and versatility of our designs.

2. Related Work

Object Detection with Transformers. DETR [3] integrates CNNs with Transformers [29], effectively eliminating the need for hand-crafted components, such as rule-based training target assignment, anchor generation, and non-maximum suppression. Several recent DETR variants [13, 20, 26, 35, 38] have been proposed to improve the attention mechanism and bipartite matching in Transformers. For example, [38] only attend to a small set of key sampling points around a reference for faster convergence. [26] learn a conditional spatial query from the decoder embedding for decoder multi-head cross-attention. [20] introduce a denoising pipeline to reduce the difficulty of bipartite matching.

In contrast, we explore another paradigm to improve the representations of Transformers for DETR via self-supervised pretraining. [9] design a pretext task based on random query patch detection. [1] train DETR using pseudo labels generated by pretrained SwAV [5] and offline proposals [28]. While similarly following the existing *pre-train and fine-tune* paradigm, our work significantly differs from UP-DETR and DETReg. To our knowledge, we make the first attempt to combine the Siamese network with a cross-attention mechanism in DETR, enabling the model to learn view-invariant localizing ability during pretraining.

Self-supervised Pretraining. One of the main approaches for self-supervised learning is to compare different augmented views of the same data instances in the representation space. Some notable studies include that by [6], which presents a simple framework by removing the requirements of specialized architectures or memory banks. [15] employ a momentum encoder with a dynamic dictionary look-up and retrieve more negative samples by using large dictionary sizes. [8] explore simple Siamese networks to learn meaningful representations with positive pairs only. [4] enforce consistency between cluster assignments produced for different augmentations of the same image. In addition, several attempts [18, 30–32, 36] have been made to learn detection-oriented representations directly using intrinsic cues, such as mask predictions [18, 36], offline region proposals [31, 33], and joint global-local partitions [30, 32]. While different in the specific learning strategies, all these works focus on learning discriminative representations for base models, which are insufficient for transfer learning in DETR. In addition, the pretext tasks in these methods cannot be directly applied to the Transformers in DETR.

Siamese Networks. Siamese networks [2] are weight-sharing neural networks and usually take two inputs for comparison, which are widely adopted in many applications, such as face verification [27], person re-identification [37], one-shot learning [19] and semi-supervised learning [24, 34]. Recent advances in self-supervised learning [6–8, 14] are also built upon Siamese networks, motivating us to explore the Siamese architecture for pretraining the Transformers in DETR.

3. Siamese DETR

An overview of our Siamese DETR architecture is presented in Figure 2, which illustrates the main pipeline of our multi-view detection pretraining. We first revisit the DETR in Section 3.1. We then describe the view construction algorithm in Section 3.2 and the multi-view pretraining paradigm of Siamese DETR in Section 3.3, powered by two specially designed self-supervised pretext tasks for learning to detect objects.

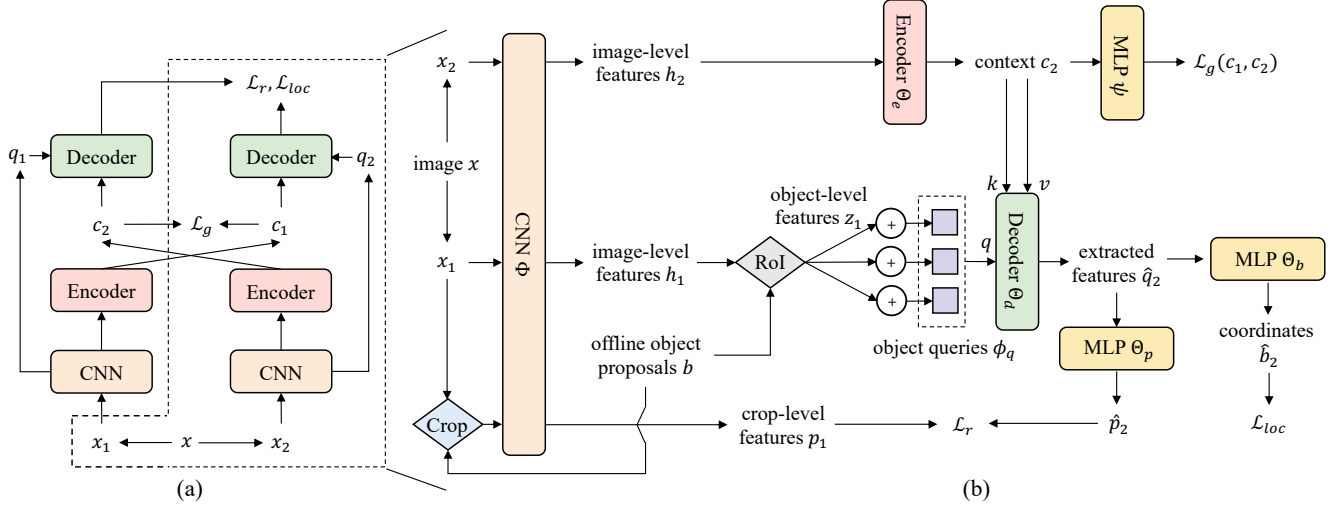


Figure 2. **(a)** Overall architecture of proposed Siamese DETR. **(b)** The forward process of one view in our symmetrical pipeline. We perform region detection (\mathcal{L}_{loc}) and semantic discrimination (\mathcal{L}_g and \mathcal{L}_r) in a multi-view fashion. Given the conditional input of region features in one view, we aim at locating and discriminating their corresponding regions in another view.

3.1. Revisiting DETR

A typical DETR model consists of two modules: (i) a backbone model, *i.e.*, CNNs, for feature extraction, (ii) Transformers with encoder-decoders architecture for set prediction, built by stacking multi-head attentions [29]. The backbone model extracts the image-level features $\mathbf{h} = \text{Backbone}(\mathbf{x})$ for a given image $\mathbf{x} \in \mathbb{R}^{3 \times H_0 \times W_0}$. Then, the Transformer encoder takes \mathbf{h} as inputs, encoding the image features as global context $\mathbf{c} \in \mathbb{R}^{C \times H_1 \times W_1}$:

$$\mathbf{c} = \text{Encoder}(\mathbf{h}). \quad (1)$$

Note that we omit the positional embedding ϕ_p in the description for clarity. The cross-attention mechanism is a general form of multi-head attention (MHA) in the Transformer decoder, which calculates the weighted sum $\hat{\mathbf{q}} \in \mathbb{R}^{N \times C}$ between the flattened global context $\mathbf{c} \in \mathbb{R}^{H_1 W_1 \times C}$ and N object queries $\phi_q \in \mathbb{R}^{N \times C}$ for further box and categorical prediction.

$$\begin{aligned} \hat{\mathbf{q}} &= \text{CrossAtten}(\mathbf{c}, \phi_q) = \text{MHA}(\mathbf{q}, \mathbf{k}, \mathbf{v}) \\ &= \sum \text{Softmax}\left(\frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d_k}}\right)\mathbf{v}. \end{aligned} \quad (2)$$

The attention weights and summations are obtained by the query \mathbf{q} , key \mathbf{k} , and value \mathbf{v} , which are the linear mapping of the context \mathbf{c} and the object queries ϕ_q :

$$\mathbf{q} = f_q(\phi_q); \quad \mathbf{k} = f_k(\mathbf{c}); \quad \mathbf{v} = f_v(\mathbf{c}). \quad (3)$$

Here f_q, f_k, f_v denote the projection for query, key and value in the cross-attention module.

In this work, we aim to pre-train the Transformers of DETR in a self-supervised way, extending the boundary of

existing self-supervised pretraining. Motivated by recent advances [6–8, 14], we propose Siamese DETR, a Siamese multi-view self-supervised framework designed for Transformers in DETR, in which the model parameters and the learnable object queries in two DETRs are all shared. Following UP-DETR [9], Siamese DETR aims at learning representations with Transformers in self-supervised pretraining while keeping the backbone model frozen.

3.2. View Construction

We start with generating two views $\{\mathbf{x}_1, \mathbf{x}_2\}$ for each unlabeled image \mathbf{x} , allowing the model to learn view-invariant object-level representations in self-supervised detection pretraining. As illustrated in Figure 3, we introduce an IoU-constrained policy to balance the shared information between two views. First, we generate a random rectangle within the image, which covers most content (50% to 100%). Then the center point of the rectangle is used as the anchor to create two sub-rectangles. By randomly expanding the sub-rectangles along the diagonal, we obtain two rectangles with the IoU larger than a threshold $\tau = 0.5$. Two rectangles are cropped from the image as the final two views $\{\mathbf{x}_1, \mathbf{x}_2\}$. We further apply randomly and independently sampled transformations on two views $\{\mathbf{x}_1, \mathbf{x}_2\}$ following the augmentation pipeline in [8]. We also apply a box jitter processing following [31] to encourage variance of scales and locations of object proposals across views.

The two augmented views $\{\mathbf{x}_1, \mathbf{x}_2\}$ are visually distinct but share adequate semantic content. Following [1, 31], we generate offline object proposals \mathbf{b} using unsupervised EdgeBoxes [39] in the overlapping area between the two views and randomly select $n = 10$ corresponding object

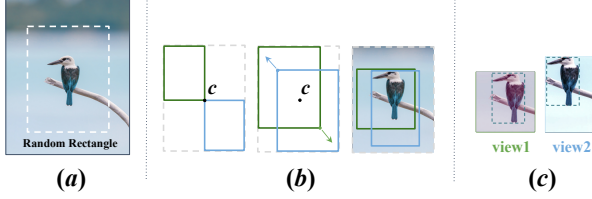


Figure 3. View construction using an IoU-constrained policy. (a) Generate a random rectangle within the image. (b) Generate two sub-rectangles on both sides of the center point c . Expand two rectangles along the diagonal while keeping the IoU larger than a threshold τ . Crop two sub-rectangles and apply augmentations. (c) Generate offline object proposals in the overlapping area.

proposals in two views $\{b_1, b_2\}$ from b . These object proposals can provide proper objectness priors to learn object-level representations during pretraining.

3.3. Multi-view Detection Pretraining

Given a region extracted from one view, we train the model for answering two questions: (1) *where is the corresponding region in another view?* (2) *what is the region, i.e., does it semantically similar?* In the following, we describe two pretext tasks designed in Siamese DETR: (1) learning to locate by **Multi-View Region Detection** and (2) learning to discriminate by **Multi-View Semantic Discrimination**.

Preparation. We take the augmented views $\{x_1, x_2\}$ as inputs to the backbone model to obtain the image-level features: $\{h_1, h_2\} = \text{Backbone}(\{x_1, x_2\})$. Then, the object-level region features $\{z_1, z_2\}$ are extracted from the image-level features $\{h_1, h_2\}$ based on each object proposal $\{b_1, b_2\}$ using RoIAlign [16]:

$$\begin{aligned} z_1 &= \text{RoIAlign}(h_1, b_1); \\ z_2 &= \text{RoIAlign}(h_2, b_2). \end{aligned} \quad (4)$$

We also obtain the corresponding crop-level region features $\{p_1, p_2\} \in \mathbb{R}^{C \times H_3 \times W_3}$ by first cropping from the augmented views and then extraction:

$$\begin{aligned} p_1 &= \text{Backbone}(\text{Crop}(x_1, b_1)); \\ p_2 &= \text{Backbone}(\text{Crop}(x_2, b_2)). \end{aligned} \quad (5)$$

By default, both $\{z_1, z_2\}$ and $\{p_1, p_2\}$ are processed with global average pooling before further usage. Only the forward pass is involved in preparations, as we focus on pretraining the Transformers.

Multi-View Cross-Attention. We propose a Multi-View Cross-Attention (MVCA) mechanism that extends the cross-attention module in DETR for multi-view representation learning. With the introduced notions for two views $\{(c_1, z_1, p_1, b_1), (c_2, z_2, p_2, b_2)\}$, we formulate the cross-attention from view x_1 to view x_2 as follows:

$$\hat{q}_2 = \text{MVCA}(c_2, \phi_q, z_1) = \text{MHA}(q_1, k_2, v_2), \quad (6)$$

where the query q_1 , key k_2 and value v_2 are given by:

$$q_1 = f_q(z_1 + \phi_q); \quad k_2 = f_k(c_2); \quad v_2 = f_v(c_2). \quad (7)$$

We add the region features z_1 from view x_1 to the object queries ϕ_q so that with conditional input of region features z_1 , the object queries ϕ_q can extract the relevant features \hat{q}_2 from the global context c_2 of view x_2 . Here, \hat{q}_2 is supposed to be aggregated features on view x_2 that are semantically consistent with the corresponding region features z_1 on view x_1 .

Learning to locate: The MVCA mechanism allows us to conduct **Multi-View Region Detection** directly. Specifically, with the input of each region feature z_1 from view x_1 and its extracted feature \hat{q}_2 , our goal is to locate the region in view x_2 that is relative to the region feature z_1 . We apply a prediction head f_{box} for box prediction:

$$\hat{b}_2 = f_{box}(\hat{q}_2) \in \mathbb{R}^{N \times 4}. \quad (8)$$

After performing bipartite matching [3], we calculate the multi-view symmetrical localization loss as:

$$\mathcal{L}_{loc} = \ell_{box}(\hat{b}_2, b_2) + \ell_{box}(\hat{b}_1, b_1), \quad (9)$$

where ℓ_{box} is a combination of generalized IoU loss and ℓ_1 loss the same as [3].

Learning to discriminate: Due to the unavailability of semantic label information, we propose **Multi-View Semantic Discrimination** to learn to discriminate at both global and regional levels. First, we apply a prediction head f_{sem} for further discriminative learning:

$$\hat{p}_2 = f_{sem}(\hat{q}_2) \in \mathbb{R}^{N \times C'}. \quad (10)$$

Considering that the context c of each view contains global contextual information, we maximize the similarity of the encoded context between two augmented views. Following [8], we apply a three-layer MLP (FC-BN-ReLU) and compute the global discrimination loss \mathcal{L}_g symmetrically as:

$$\begin{aligned} \mathcal{L}_g &= \mathcal{C}[\text{MLP}(c_1), \text{detach}(c_2)] + \\ &\quad \mathcal{C}[\text{MLP}(c_2), \text{detach}(c_1)], \end{aligned} \quad (11)$$

where \mathcal{C} is the negative cosine similarity. In addition, due to the semantic consistency between the input region features z_1 and the extracted features \hat{q}_2 , we consider maximizing the semantic consistency for each region. Here, despite representing the same instances, the crop-level region features p_1 can provide more discriminative information than object-level region features z_1 . It is because the object-level region features z_1 are extracted from the image-level features (See Equation 4 and 5) and contain an aggregate of surrounding contexts with less discriminative information

on themselves, especially for small regions. It motivates us to replace object-level region features z_1 with crop-level region features p_1 as the learning objectives. Also, to avoid potential training collapse, we reconstruct the crop-level region features p_1 . Finally, we formulate the semantic consistency objective to improve the region discrimination as:

$$\mathcal{L}_r = \mathcal{D}(\hat{p}_2, p_1) + \mathcal{D}(\hat{p}_1, p_2), \quad (12)$$

where \mathcal{D} is the normalized ℓ_2 distance.

Loss Function. Formally, the overall loss function for Siamese DETR is formulated as:

$$\mathcal{L} = \lambda_0 \mathcal{L}_r + \lambda_1 \mathcal{L}_g + \lambda_2 \mathcal{L}_{loc}, \quad (13)$$

where $\lambda_{0/1/2}$ are the loss weighting hyper-parameters.

3.4. Discussion

Compared with previous methods [1, 9] and conventional multi-view SSL methods like MoCo [15], we introduce siamese network in the pre-training pipeline, exploring the view-invariant localization ability for DETR pretraining. Through combining the characteristics of cross-attention with siamese network, Siamese DETR can aggregate the latent information from key view given the cue of the query view, then localizing and discriminating the corresponding regions. It is completely different from the common multi-view techniques, *i.e.*, simply contrasting the output of the DETR. The proposed Multi-View Cross-Attention enables model to learn to view-invariant localization ability, providing better priors in downstream tasks compared with UP-DETR, which spends most of pre-training time detecting the background classes. Meanwhile, DETReg [1], which can be viewed as detecting the regions generated by proposals and pseudo labels, does not consider view-invariant detection pretraining as well.

4. Experiments

4.1. Implementation Details

Architecture. Siamese DETR consists of a frozen ResNet-50 backbone, pretrained by SwAV [5] on ImageNet, and a Transformer with encoder-decoder architecture. Both the Transformer encoder and decoder are stacked with 6 layers of 256 dimensions and 8 attention heads. To verify the performance and generalization of our design, we compare Siamese DETR with a baseline model without pretraining (denoted as *from scratch*), UP-DETR [9], and DETReg [1]. We use three DETR variants, *i.e.*, original DETR [3] (denoted as *Vanilla DETR*), Conditional DETR [26], and Deformable DETR (Single-Scale and Multi-Scale, denoted as Deform-SS and Deform-MS, respectively) [38]. For a fair comparison, the number of object queries is 100 in Vanilla DETR and Conditional DETR. We also use 300 queries

in Conditional DETR and Deformable DETR. Besides, we also provide the transfer results of more advanced DETR-like architecture (*e.g.*, DAB-DETR [23]) in the appendix. We implement Siamese DETR based on the MMSelfSup¹.

Evaluation Protocol. We follow the evaluation protocol in UP-DETR. Specifically, we first pretrain DETR variants on ImageNet (~ 1.28 million images) [10] or COCO train2017 ($\sim 118k$ images) [22] separately. Then, the ImageNet-pretrained models are finetuned on COCO train2017 or PASCAL VOC trainval107+12 ($\sim 16.5k$ images) [12] separately, while COCO-pretrained models are finetuned on PASCAL VOC trainval107+12. We report COCO-style metrics, including AP, AP₅₀, AP₇₅, AP_s, AP_m, AP_l, in both COCO val2017 and PASCAL VOC test2007 benchmarks.

Pretraining. We use an AdamW [25] optimizer with a total batch size of 256 on ImageNet and 64 on COCO, a learning rate of 1×10^{-4} and a weight decay of 1×10^{-4} . We adopt a full schedule of 60 epochs, and the learning rate decays at 40 epochs, denoted as the 40/60 schedule for brevity. Unless specified, UP-DETR uses random boxes, DETReg uses proposals generated by Selective Search [28], and Siamese DETR uses Edgeboxes in experiments.

Finetuning. For Vanilla DETR, we adopt 120/150 and 40/50 schedules in COCO and PASCAL VOC benchmarks. The initial learning rates of the Transformer and backbone are set to 1×10^{-4} and 5×10^{-5} . For the other two DETR variants, we report the result under the 40/50 schedule. The initial learning rates of the Transformer and backbone are set to $1 \times 10^{-4}/2 \times 10^{-4}$ and $5 \times 10^{-5}/2 \times 10^{-5}$ in Conditional/Deformable DETR, respectively. The batch size of all setups in finetuning is set to 32 for a fair comparison.

4.2. Main Results

COCO Object Detection. Table 1 shows the transfer results on COCO. Siamese DETR achieves the best performance using three different DETR variants on all setups. Especially for Deformable DETR of Multi-Scale, Siamese DETR boosts the model upon baseline more significantly than UP-DETR and DETReg, demonstrating the compatibility of our design with different DETR architectures.

Besides, when adopting a stricter IoU threshold in metrics, *e.g.*, from AP₅₀ to AP₇₅, Siamese DETR achieves a more considerable performance lead in most cases. It suggests the representations learned by Siamese DETR provide a stronger localization prior. We further illustrate AP metrics of Siamese DETR and UP-DETR using different IoU thresholds in Figure 4. Specifically, Siamese DETR performs well in localizing small objects and draws the gap against UP-DETR for medium and large objects when the IoU threshold is greater than 0.7.

¹<https://github.com/open-mmlab/mmselfsup>. Apache-2.0 License.

Table 1. Comparisons of Siamese DETR with supervised/UP-DETR/DETRReg in the COCO detection benchmark. The results of all models are achieved by officially-released repositories and pretrained models. Here “#epoch” denotes the number of epochs in downstream finetuning.

Method	Backbone	DETR	#query	#epoch	AP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l
<i>from scratch</i>	Sup. R50	Vanilla	100	150	39.5	60.3	41.4	17.5	43.0	59.1
<i>from scratch</i>	SwAV R50	Vanilla	100	150	39.7	60.3	41.7	18.5	43.8	57.5
UP-DETR	SwAV R50	Vanilla	100	150	40.5	60.8	42.6	19.0	44.4	60.0
DETRReg	SwAV R50	Vanilla	100	150	41.9	61.9	44.1	19.1	45.7	61.5
ours	SwAV R50	Vanilla	100	150	42.0	63.1	44.2	19.6	46.0	61.9
<i>from scratch</i>	SwAV R50	Conditional	100	50	37.7	59.6	39.2	17.1	41.7	56.3
UP-DETR	SwAV R50	Conditional	100	50	39.4	61.2	41.0	18.1	43.0	58.7
DETRReg	SwAV R50	Conditional	100	50	40.2	61.8	42.0	19.1	43.7	60.0
ours	SwAV R50	Conditional	100	50	40.5	61.6	42.6	19.5	44.2	60.1
<i>from scratch</i>	SwAV R50	Conditional	300	50	41.1	62.3	43.4	20.6	45.0	59.4
UP-DETR	SwAV R50	Conditional	300	50	41.5	63.2	43.6	21.3	45.4	60.2
ours	SwAV R50	Conditional	300	50	43.0	64.2	45.6	22.0	47.2	61.8
<i>from scratch</i>	SwAV R50	Deform-SS	300	50	40.3	60.9	42.9	20.1	44.8	57.2
UP-DETR	SwAV R50	Deform-SS	300	50	40.8	61.8	43.4	20.4	45.1	59.1
ours	SwAV R50	Deform-SS	300	50	42.1	62.8	44.7	22.3	46.6	59.9
<i>from scratch</i>	SwAV R50	Deform-MS	300	50	45.5	64.2	49.4	27.8	49.2	59.4
UP-DETR	SwAV R50	Deform-MS	300	50	45.3	64.5	49.6	26.0	49.2	59.9
DETRReg	SwAV R50	Deform-MS	300	50	45.5	64.1	49.9	26.9	49.5	59.6
ours	SwAV R50	Deform-MS	300	50	46.3	64.6	50.5	28.1	50.1	61.5

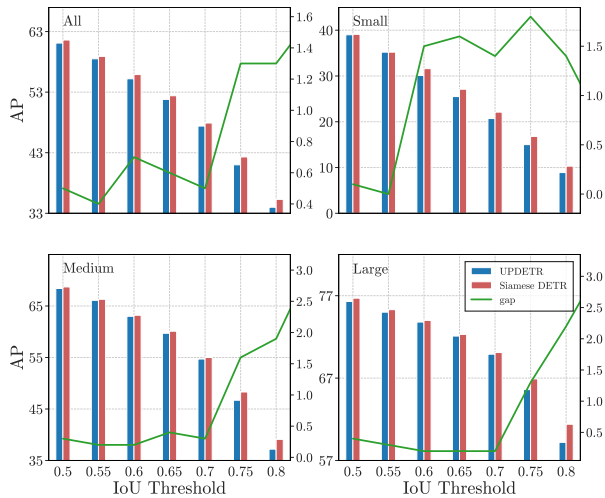


Figure 4. AP on COCO using different IoU Thresholds.

PASCAL VOC Detection. Table 2 shows the transfer results on PASCAL VOC. Similar to conclusions on COCO, Siamese DETR achieves the best performance among all approaches on PASCAL VOC. We also report the result of COCO-pretrained models. Siamese DETR is 6.4 AP better than UP-DETR and 1.8 AP better than DETRReg, which verifies the compatibility of Siamese DETR with different pretraining datasets, especially the scene-centric COCO.

4.3. Ablations

Effectiveness of Two Proposed Pretext Tasks. We train five Siamese DETR variants using Conditional DETR on ImageNet and finetune them on PASCAL VOC. Results are shown in Table 3. We treat UP-DETR as the baseline (56.9 AP on PASCAL VOC), which performs single-view patch detection with Transformer and reconstructs the decoder’s output with its input patches.

By extending single-view detection into a multi-view manner, (a) obtains a competitive result of 57.1 AP, suggesting that the view-invariant representations can perform better in downstream detection tasks. We further maximize the multi-view semantic consistency in terms of global and region discrimination, improving the transfer performance by 0.2 AP and 0.6 AP in (b) and (d), respectively. Global discrimination (b) brings smaller performance gains than region discrimination (d) in detection tasks, suggesting that it is impractical to directly apply existing instance discrimination pretext tasks [15, 19] on Transformers of DETR in detection-oriented pretraining tasks. We also notice that (d) using crop-level region features achieves better performance than (c) using object-level region features, which verifies more discriminative information in crop-level region features. Finally, (e) with both Multi-View Region Detection and Multi-View Semantic Discrimination yields the best result of 58.1 AP.

Object Proposals. Edgeboxes [39] in Siamese DETR provide objectness priors during pretraining. We attempt to

Table 2. Comparisons of Siamese DETR with supervised/UP-DETR/DETReg in the PASCAL VOC detection benchmark. The results of all models are achieved by officially-released repositories and pretrained models. Here “#epoch” denotes the number of epochs in downstream finetuning.

Method	Backbone	DETR	Pretrain Dataset	#query	#epoch	AP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l
<i>from scratch</i>	SwAV R50	Vanilla	-	100	50	28.5	47.5	29.4	1.3	7.4	40.3
UP-DETR	SwAV R50	Vanilla	ImageNet	100	50	50.0	73.5	53.4	6.5	29.5	64.1
DETReg	SwAV R50	Vanilla	ImageNet	100	50	53.8	76.5	57.3	8.3	35.4	67.5
ours	SwAV R50	Vanilla	ImageNet	100	50	54.4	77.4	57.6	7.7	35.0	68.4
<i>from scratch</i>	SwAV R50	Vanilla	-	100	150	47.8	73.8	50.9	5.4	27.6	61.4
UP-DETR	SwAV R50	Vanilla	ImageNet	100	150	54.4	78.1	58.6	10.5	35.8	67.5
DETReg	SwAV R50	Vanilla	ImageNet	100	150	57.0	79.7	61.6	11.5	39.5	70.2
ours	SwAV R50	Vanilla	ImageNet	100	150	57.4	80.3	62.2	11.6	39.3	71.0
<i>from scratch</i>	SwAV R50	Conditional	-	100	50	49.9	78.2	55.3	8.1	33.5	65.1
UP-DETR	SwAV R50	Conditional	ImageNet	100	50	56.9	81.5	61.6	11.3	39.2	69.8
DETReg	SwAV R50	Conditional	ImageNet	100	50	57.5	82.0	62.4	11.8	40.7	71.0
ours	SwAV R50	Conditional	ImageNet	100	50	58.1	81.6	62.8	12.2	40.6	71.5
<i>from scratch</i>	SwAV R50	Deform-SS	-	300	50	53.8	79.5	59.1	11.8	39.8	65.7
UP-DETR	SwAV R50	Deform-SS	ImageNet	300	50	54.0	79.3	58.8	11.4	38.6	66.1
ours	SwAV R50	Deform-SS	ImageNet	300	50	58.0	81.8	64.0	14.0	43.3	70.3
<i>from scratch</i>	SwAV R50	Deform-MS	-	300	50	56.1	80.7	61.9	17.4	42.7	66.4
UP-DETR	SwAV R50	Deform-MS	ImageNet	300	50	56.4	80.9	62.3	17.3	41.3	67.4
DETReg	SwAV R50	Deform-MS	ImageNet	300	50	59.7	82.0	66.4	18.2	46.4	70.4
ours	SwAV R50	Deform-MS	ImageNet	300	50	61.2	82.9	67.7	19.3	47.1	72.2
<i>from scratch</i>	SwAV R50	Conditional	-	100	50	49.9	78.2	55.3	8.1	33.5	65.1
UP-DETR	SwAV R50	Conditional	COCO	100	50	51.3	79.0	55.3	9.5	35.3	63.7
DETReg	SwAV R50	Conditional	COCO	100	50	55.9	80.0	61.6	11.0	39.3	68.5
ours	SwAV R50	Conditional	COCO	100	50	57.7	80.9	62.5	11.0	40.4	70.9

Table 3. Ablations on two proposed pretext tasks, *i.e.*, Multi-View Region Detection and Multi-View Semantic Discrimination. The notation “R-O” denotes region discrimination using object-level region features, “R-C” denotes region discrimination using crop-level region features, and “G” denotes global discrimination.

Method	Region Det.	Semantic Disc.	AP
UP-DETR	single-view	R-O	56.9
ours (a)	multi-view	-	57.1
ours (b)	multi-view	G	57.3
ours (c)	multi-view	R-O	57.3
ours (d)	multi-view	R-C	57.7
ours (e)	multi-view	R-C + G	58.1

replace it with boxes generated randomly or by Selective Search. All experiments are conducted using Conditional DETR. The results are shown in Table 4.

When pre-trained on object-centric datasets like ImageNet, which contains one single object in the center of the image, better objectness priors bring little improvement. In this case, Siamese DETR still outperforms its counterparts by about 0.1 AP using random proposals or Edgeboxes.

When pre-trained on scene-centric datasets like COCO, which contains multiple objects in the image, great improvements are found in all methods after applying better

Table 4. Comparisons of using different object proposals. The notation “A→B” denotes that the model is pretrained on dataset “A” and then finetuned on dataset “B”.

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

objectness priors from proposals. Using random proposals, Siamese DETR outperforms UP-DETR by 3.6 AP and DETReg by 3.0 AP in the COCO→VOC setup, which suggests Siamese DETR learns better detection-oriented representations without any objectness priors. When replaced with Selective Search and Edgeboxes, the performance gaps are

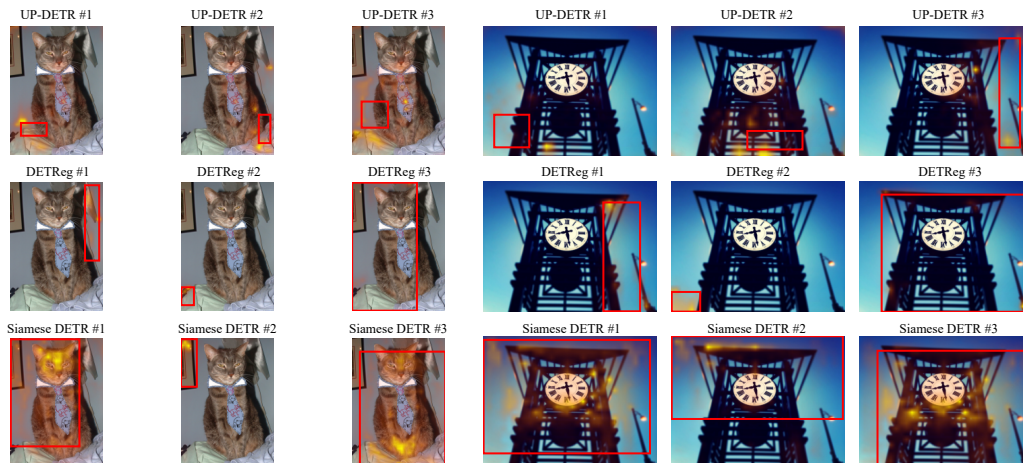


Figure 5. Visualization of box predictions and attention maps in downstream tasks. All these models (Vanilla DETRs) are initialized by Siamese DETR, UP-DETR, and DETReg without fine-tuning.

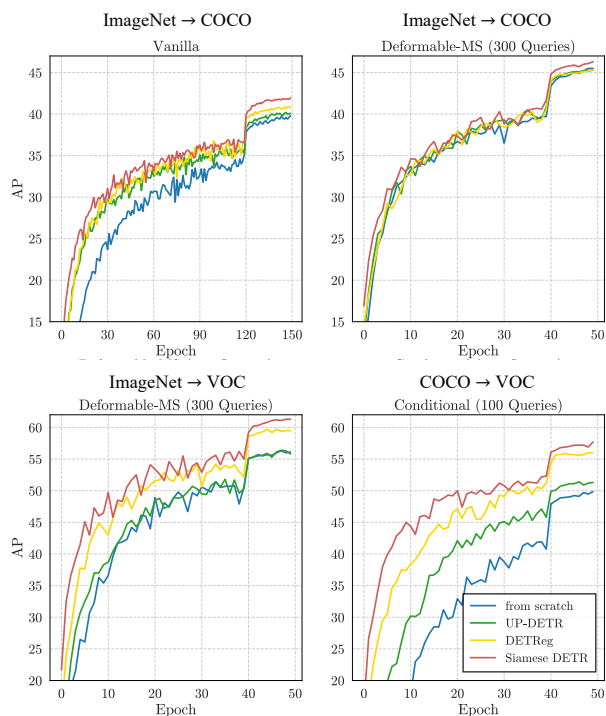


Figure 6. Convergence curves.

alleviated. In this case, our Siamese DETR with Edgeboxes achieves the best performance among all setups.

Convergence. The convergence curves of three DETR variants on downstream COCO and PASCAL VOC are illustrated in Figure 6. Compared with UP-DETR, DETReg, and the *from scratch* model, Siamese DETR converges faster and outperforms its counterparts by significant margins.

Visualization. We provide qualitative results for further understanding the advantage of Siamese DETR. In down-

stream tasks, we use three Vanilla DETRs, initialized by Siamese DETR, UP-DETR, and DETReg pretraining models. Figure 5 illustrates their box predictions and corresponding attention maps of decoder. Queries in Siamese DETR have stronger objectness priors, predicting more available box proposals overlapped with objects in the image. Benefitting from discriminative representations, cross-attention in Siamese DETR places more focus on the objects in the proposals. These qualitative results verify the transferability of Siamese DETR in downstream tasks.

5. Conclusion and Limitations

In this paper, we propose Siamese DETR, a novel self-supervised pretraining method for DETR. With two newly-designed pretext tasks, we directly locate the query regions in a cross-view manner and maximize multi-view semantic consistency, learning localization and discrimination representations transfer to downstream detection tasks. Siamese DETR achieves better performance with three DETR variants in COCO and PASCAL VOC benchmark against its counterpart. Despite the great potential for pretraining DETR, Siamese DETR has a limitation in that it still relies on a pre-trained CNN, *e.g.*, SwAV, without integrating CNN and Transformer into a unified pretraining paradigm. In our future work, a more efficient framework for the end-to-end DETR pretraining is desirable.

Acknowledgement. This study is supported by National Key Research and Development Program of China (2021YFB1714300). It is also supported under the RIE2020 Industry Alignment Fund - Industry Collaboration Projects (IAF-ICP) Funding Initiative, National Natural Science Foundation of China (62132001), Singapore MOE AcRF Tier 2 (MOE-T2EP20120-0001), as well as cash and in-kind contribution from the industry partner(s).

References

- [1] Amir Bar, Xin Wang, Vadim Kantorov, Colorado J Reed, Roei Herzig, Gal Chechik, Anna Rohrbach, Trevor Darrell, and Amir Globerson. Detreg: Unsupervised pretraining with region priors for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14605–14615, 2022. 1, 2, 3, 5
- [2] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. Signature verification using a “Siamese” time delay neural network. *NeurIPS*, 6:737–744, 1993. 2
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *eccv*, pages 213–229. Springer, 2020. 1, 2, 4, 5
- [4] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *NeurIPS*, 2020. 1, 2
- [5] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. 2020. 2, 5
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020. 1, 2, 3
- [7] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv*, 2020. 1, 2, 3
- [8] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, pages 15750–15758, 2021. 1, 2, 3, 4
- [9] Zhigang Dai, Bolun Cai, Yugeng Lin, and Junying Chen. Up-detr: Unsupervised pre-training for object detection with transformers. In *CVPR*, pages 1601–1610, 2021. 1, 2, 3, 5
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009. 5
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2020. 1
- [12] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2):303–338, 2010. 5
- [13] Peng Gao, Minghang Zheng, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fast convergence of detr with spatially modulated co-attention. In *ICCV*, pages 3621–3630, October 2021. 2
- [14] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *NeurIPS*, 2020. 1, 2, 3
- [15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020. 1, 2, 5, 6
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969, 2017. 4
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 1
- [18] Olivier J. Hénaff, Skanda Koppula, Jean-Baptiste Alayrac, Aaron van den Oord, Oriol Vinyals, and João Carreira. Efficient visual pretraining with contrastive detection. In *ICCV*, pages 10086–10096, October 2021. 2
- [19] Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. Siamese neural networks for one-shot image recognition. In *ICMLW*, volume 2. Lille, 2015. 2, 6
- [20] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13619–13627, 2022. 2
- [21] Junnan Li, Pan Zhou, Caiming Xiong, and Steven CH Hoi. Prototypical contrastive learning of unsupervised representations. *ICLR*, 2021. 1
- [22] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014. 5
- [23] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. Dab-detr: Dynamic anchor boxes are better queries for detr. *arXiv preprint arXiv:2201.12329*, 2022. 5
- [24] Yen-Cheng Liu, Chih-Yao Ma, Zijian He, Chia-Wen Kuo, Kan Chen, Peizhao Zhang, Bichen Wu, Zsolt Kira, and Peter Vajda. Unbiased teacher for semi-supervised object detection. *arXiv preprint arXiv:2102.09480*, 2021. 2
- [25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv*, 2017. 5
- [26] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *ICCV*, pages 3651–3660, 2021. 2, 5
- [27] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, pages 1701–1708, 2014. 2
- [28] Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective search for object recognition. *IJCV*, 104(2):154–171, 2013. 2, 5
- [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 2, 3
- [30] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning for self-supervised visual pre-training. In *CVPR*, pages 3024–3033, 2021. 2

- [31] Fangyun Wei, Yue Gao, Zhirong Wu, Han Hu, and Stephen Lin. Aligning pretraining for detection via object-level contrastive learning. *arXiv*, 2021. [2](#), [3](#)
- [32] Enze Xie, Jian Ding, Wenhai Wang, Xiaohang Zhan, Hang Xu, Peize Sun, Zhenguo Li, and Ping Luo. Detco: Unsupervised contrastive learning for object detection. In *ICCV*, pages 8392–8401, 2021. [2](#)
- [33] Jiahao Xie, Xiaohang Zhan, Ziwei Liu, Yew Soon Ong, and Chen Change Loy. Unsupervised object-level representation learning from scene images. *arXiv*, 2021. [2](#)
- [34] Mengde Xu, Zheng Zhang, Han Hu, Jianfeng Wang, Lijuan Wang, Fangyun Wei, Xiang Bai, and Zicheng Liu. End-to-end semi-supervised object detection with soft teacher. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3060–3069, 2021. [2](#)
- [35] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. [2](#)
- [36] Yucheng Zhao, Guangting Wang, Chong Luo, Wenjun Zeng, and Zheng-Jun Zha. Self-supervised visual representations learning by contrastive mask prediction. In *ICCV*, pages 10160–10169, 2021. [2](#)
- [37] Meng Zheng, Srikrishna Karanam, Ziyang Wu, and Richard J Radke. Re-identification with consistent attentive Siamese networks. In *CVPR*, pages 5735–5744, 2019. [2](#)
- [38] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *ICLR*, 2020. [2](#), [5](#)
- [39] C Lawrence Zitnick and Piotr Dollár. Edge boxes: Locating object proposals from edges. In *ECCV*, pages 391–405. Springer, 2014. [3](#), [6](#)