CROSS RESOLUTION ENCODING-DECODING FOR DETECTION TRANSFORMERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Detection Transformers (DETR) are renowned object detection pipelines, however computationally efficient multiscale detection using DETR is still challenging. In this paper, we propose a Cross-Resolution Encoding-Decoding (CRED) mechanism that allows DETR to achieve the accuracy of high-resolution detection while having the speed of low-resolution detection. CRED is based on two modules; Cross Resolution Attention Module (CRAM) and One Step Multiscale Attention (OSMA). CRAM is designed to transfer the knowledge of low-resolution encoder output to a high-resolution feature. While OSMA is designed to fuse multiscale features in a single step and produce a feature map of a desired resolution enriched with multiscale information. When used in prominent DETR methods, CRED delivers accuracy similar to the high-resolution DETR counterpart in roughly 50% fewer FLOPs. Specifically, state-of-the-art DN-DETR, when used with CRED (calling CRED-DETR), becomes 76% faster, with ~ 50% reduced FLOPs than its high-resolution counterpart with 202G FLOPs on MS-COCO benchmark. We plan to release pretrained CRED-DETRs for use by the community.

023 024 025

026 027

003 004

006

008

009 010

011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Detection Transformers (DETR) (Carion et al., 2020) are end-to-end object detection frameworks without post-processing, such as anchor boxes, box matching, and non-maximal suppression (NMS) that are required in ConvNet-based detectors (Ren et al., 2015; Liu et al., 2016). Since their first entry (Carion et al., 2020), DETRs have evolved, mainly regarding query design (Gao et al., 2021b; Zhu et al., 2021; Meng et al., 2021; Wang et al., 2021; Liu et al., 2022; Li et al., 2022) to improve the detection performance, while recent DETR pipelines achieve that via salient points (Liu et al., 2023) or unsupervised pretraining (Chen et al., 2023).

State-of-the-art DETRs exploit high-resolution features (known as Dilated Convolution or DC variant)
(Liu et al., 2022; Li et al., 2022) or dense multi-scale features (Zhang et al., 2023b) to push the
detection accuracy further. However, they suffer from high computation complexity w.r.t. their
accuracy gains. For example, DAB DETR (Liu et al., 2022) with high-resolution features improved
2.3AP, but it introduces a 114% rise in FLOPs.

The primary reason is the quadratic computational complexity (Dai et al., 2021) of the Transformer's attention mechanism (Zhang et al., 2023a) w.r.t. the spatial size, i.e., a Transformer is $O(H^2W^2)$ complex in processing a feature map of a spatial size of $H \times W$. Hence, doubling the resolution of the encoder input quadruples its computations while also affecting the decoder. Although deformable attention (Zhu et al., 2021) addresses this issue, it causes additional runtime overhead due to irregular memory accesses.

To address this issue, recent IMFA (Zhang et al., 2023a) proposes using a low-resolution feature map and Top-K sparsely sampled higher-resolution features (Figure 1). IMFA exhibits improvements in accuracy with lower FLOPs. However, sparse sampling incurs memory access costs due to irregular memory access, which becomes prevalent with more samples. Recent (Li et al., 2023) exploits attention among only interleaved tokens for reducing computations in multiscale attention in the encoder. (Zhao et al., 2024b) works on reduced resolution images (640×640) instead of high resolution settings (800×1200) on standard DETR.

053 Our aim aligns with improving DETR speed and accuracy by using multiscale features at highresolution settings. However, our key motivation is based on our finding that the encoder consists



Figure 1: Left: Single-Scale and/ DC DETR. Middle: IMFA DETR (Zhang et al., 2023a). Right: CRED DETR (Ours). Multiple arrows between two modules indicate layerwise refinement. Stage-1 features are generally not used due to large resolution and small receptive field.

066 of most computations relative to the decoder. Therefore, we propose a principal design change in 067 the DETR pipeline, i.e., feeding the encoder with low-resolution while feeding the decoder with a 068 high-resolution feature obtained from the backbone. By keeping the encoder input low-resolution, we save computations, while by keeping the decoder input high-resolution, we provide the decoder 069 access to fine-grained details. Since the high-resolution features from the backbone lack large spatial context (Carion et al., 2020), we develop a novel Cross Resolution Attention Module (CRAM). It 071 utilizes encoder output that has a global context and transfers this information into the high-resolution 072 feature map. When this feature map is fed to the decoder, the decoder has access to the fine-grained 073 details and the global context, thus improving the accuracy and runtime. 074

Then, by exploiting this capability of CRAM to transfer the information from low-resolution encoder 075 output to high-resolution feature, we propose to reduce the resolution of an encoder further to save 076 computations. This behavior is intended to develop faster DETRs offering speed-accuracy tradeoffs. 077 However, feeding the encoder naively with reduced resolution degrades its performance. Hence, we 078 devise a novel module called One Step Multiscale Attention (OSMA), which attends to multiscale 079 information in one step and can produce tokens or feature maps of the desired resolution enriched with multiscale information. When the encoder is fed the tokens produced by OSMA at aggressively 081 low resolution, accuracy is considerably improved while avoiding any runtime overhead relative to 082 the baseline which was fed with a low-resolution backbone output.

083 We name our overall approach as Cross-Resolution Encoding-Decoding (CRED), shown in Figure 1. 084 We demonstrate that our CRED-enhanced DETR can attain an AP (average precision) equivalent 085 to the original DETR's high-resolution counterpart in 50% fewer FLOPs at 76% improved runtime. For instance, DN-DETR (Li et al., 2022)+CRED reduces FLOPs from 202G to 103G ($\sim 50\%$) and 087 improves FPS from 13FPS to 23FPS ($\sim 76\%$) without losing accuracy. In addition, applying CRED in DETR variants (Meng et al., 2021; Liu et al., 2022; Zhang et al., 2023a) consistently improves their accuracy and runtime compared with their DC variants. To improve runtime further, we half 090 the encoder resolution via OSMA. Interestingly, we only observe -1AP in this configuration while the runtime is further improved by 84% compared with the vanilla DC variant. This signifies the 091 potential of Cross-Resolution Encoding-Decoding mechanism in DETRs. 092

093 094

063

064

065

2 PRELIMINARY

096 This section revisits the architectural design of vanilla DETR (Carion et al., 2020) and advanced DETRs (Liu et al., 2022; Li et al., 2022; Zhu et al., 2021). DETR comprises a backbone, a Transformer 098 encoder, and a Transformer decoder. In a backbone (Li et al., 2022), it is common to keep five stages, each operating at a resolution half of its previous stage. Thus, the final stage (stage-5) runs at a 099 *stride of 32* at the original resolution. Once an image $I \in \mathbb{R}^{3 \times H \times W}$ is fed into a backbone (Figure 1), 100 the backbone output $F_b \in \mathbb{R}^{C \times H_0 \times W_0}$ is fed to a Transformer encoder, producing encoded feature 101 embeddings or tokens $F_e \in \mathbb{R}^{d \times H_0 \times W_0}$. These embeddings are fed to the Transformer decoder to 102 produce a fixed number of queries (N_q) , each representing an object detectable in the image. The 103 queries are passed through two Feedforward neural Networks (FFN) to obtain the object class and its 104 bounding box. During training, bipartite or Hungarian matching (Carion et al., 2020) is performed 105 for one-to-one assignments of ground truth and predictions. 106

In DETR pipelines (Carion et al., 2020; Liu et al., 2022; Li et al., 2022), the embeddings F_e are directly fed to the decoder where the N_q queries interact with $H_0 \times W_0$ features in F_e via cross attention.



Figure 2: CRAM: Cross Resolution Attention Module.

124 The encoder complexity is $O(H_0^2 W_0^2)$ in its self-attention, whereas for the decoder, it is $O(N_q H_0 W_0)$ in 125 cross attention and $O(N_a^2)$ in its self-attention. In DETRs, N_a is relatively smaller, i.e., 300; however, 126 H_0W_0 is quite a large number depending on the image size, e.g., 1280×800 (Liu et al., 2022; Li 127 et al., 2022). Hence, most of the computations are concentrated in the encoder. When opting for 128 high-resolution detection, known as the DC variant (Liu et al., 2022; Li et al., 2022), the backbone 129 output stride is set lower than 32. For this purpose, stage-5 of the backbone is set to run at stride 16 (Meng et al., 2021; Liu et al., 2022; Li et al., 2022) w.r.t. the image, thereby doubling the resolution 130 131 of F_b . This leads to quadratic growth in the encoder computations due to the increased resolution (Dai et al., 2021). 132

In this paper, we rethink the encoder-decoder information flow while leveraging multiscale features in a computationally efficient manner. Summarily, we propose an approach that feeds the encoder with a stride \geq 32 while feeding the decoder with a stride \leq 32 to harness the best of both worlds, i.e., the accuracy of high-resolution detection and speed of low-resolution detection. To our knowledge, such a mechanism has not been demonstrated in DETRs.

3 Method

123

138 139

140

145

147

We propose to enhance the DETR design with the two modules. *Firstly*, we develop CRAM
 which enables Cross-Resolution Encoding-Decoding. *Secondly*, we develop OSMA module, which
 facilitates generating feature maps enriched with multiscale information in a computationally efficient
 manner and further enhances the performance of CRED.

- 146 3.1 CROSS RESOLUTION ATTENTION MODULE (CRAM)
- As mentioned in Sec. 2, high-resolution detection improves DETR accuracy, but high-resolution input to the encoder is compute-intensive due to its quadratic complexity (Dai et al., 2021). Hence, we feed the encoder with low resolution while feeding the decoder with high resolution, calling it *cross-resolution encoding-decoding paradigm*. However, the information flow between them no longer exists due to the different input sources to the encoder and decoder. Hence, we develop CRAM that acts as a bridge between the encoder and decoder in the proposed cross-resolution encoding decoding paradigm. The overall design of our approach (CRAM) is shown in Figure 2.

Consider two feature maps $X \in \mathbb{R}^{C \times h \times w}$ and $Y \in \mathbb{R}^{C \times H \times W}$. In our Cross-Resolution Encoding-Decoding approach, the low-resolution feature X is fed to the encoder layers, while the high-resolution feature Y is fed to the attention transfer module. To transfer the knowledge of the encoder embeddings ($X_e = \text{Encoder}(X), X_e \in \mathbb{R}^{C \times h \times w}$) to Y, we spatially upsample ($\hat{X}_e = \text{Upsample}(X_e)|_{H,W} \hat{X}_e \in \mathbb{R}^{C \times H \times W}$) the encoder output to match the resolution of Y, and a concatenation operation is performed.

Now, we use a linear projection layer, which combines the features information of the highresolution features and the upsampled low-resolution feature from the encoder. The output $(Z = \text{Linear}(\text{Cat}(Y, X_e)))$ is normalized via LayerNorm (Ba et al., 2016) and passed through SiLU



Figure 3: Local aggregation of the multiscale features in OSMA for $g_0 = 1$ which produces $Q \in \mathbb{R}^{N_g \times T \times C}$

173 activation (Hendrycks & Gimpel, 2016) ($\hat{Y}_e = \text{SiLU}(\text{LayerNorm}(Z)), \ \hat{Y}_e \in \mathbb{R}^{C \times H \times W}$). We use the 174 residual connections to facilitate smoother optimization. This process is repeated as many times as 175 the number of encoder layers to refine the high-resolution Y by transferring the knowledge embedded 176 in the encoder layers. The input to CRAM is initialized with the high-resolution feature from an 177 intermediate backbone stage (Sec. 2), whereas the final \hat{Y}_e is fed to the decoder for the predictions 178

179 **Computation Complexity** An encoder layer has $\mathcal{O}(h^2w^2)$ complexity to process its input $X_e \in$ $\mathbb{R}^{h \times w}$. In our case, we feed the encoder with $X_e \in \mathbb{R}^{H/r \times W/r}$, where $r \ge 2$, and feed the decoder with 180 $X_d \in \mathbb{R}^{H \times W}$. In aggregation, we have a total complexity of $\mathscr{O}(\frac{H^2W^2}{r^4} + HW)$, which is far lower than 182 the vanilla encoder complexity of $\mathcal{O}(H^2W^2)$. With this strategy, we achieve accuracy equivalent to when the encoder is fed with the high-resolution feature while having at least 50% fewer overall computations for r = 2. Specifically, w.r.t. vanilla encoder with high resolution, a FLOP saving of 185 50% is obtained. See Sec 4.2 for the computational complexity analysis.

187 188

189

190

191

192

181

183

184

171

172

Why Could Cross Resolution Attention Transfer Improve Performance? In vanilla DETR design, the encoder embeddings are produced from stage-5 (low-resolution) and have global receptive filed via self-attention (Dosovitskiy et al., 2020), which are fed to the decoder for the predictions. However, these embeddings do not have fine details of smaller objects. On the other hand, the high-resolution input Y (earlier stages, e.g., stage-4) to CRAM has a small receptive field but has details of smaller objects.

193 In CRAM, the concatenation operation followed by a linear layer infuses the local and global context 194 to produce fine-grained, high-resolution features, similar to what earlier semantic segmentation 195 approaches (Zhao et al., 2017) used for improving performance. In the same way, with this operation, 196 the high-resolution feature Y acquires a global receptive field when concatenated with the encoder 197 embedding or tokens X_{e} . After the layerwise refinement, it is fed to the decoder, which improves the accuracy and speed; even the encoder still functions at a smaller resolution.

199 200

201

3.2 ONE STEP MULTISCALE ATTENTION (OSMA)

202 In single-scale operation of DETRs (Liu et al., 2022; Li et al., 2022), the backbone output F_b is fed 203 to the encoder. However, F_b does not have direct access to the multiscale information. Whereas 204 multiscale DETRs feed the encoder with either sampled (Zhu et al., 2021) or dense multiscale features 205 (Zhang et al., 2023b), which is computationally heavy.

206 We propose OSMA, which produces F_o enriched with multiscale information, i.e., the best of 207 both single-scale and multiscale methods without progressive fusion, i.e., fusing two scales at a 208 time (yol)(Figure 1). In general, F_b has a stride 32 w.r.t. the input (He et al., 2016); however, using 209 multiscale features, OSMA can produce a feature map F_o ('o' refers to OSMA) of stride greater or 210 lesser than 32 which directly controls the encoder computations. With this functionality, OSMA 211 offers better features or tokens enriched with multiscale information to be utilized by the encoder.

212 OSMA has three main steps: First, Local aggregation of multiscale features, Second, performing one-213 step attention on the aggregated features, and *Third*, broadcasting the output into desired resolution. 214

Local Aggregation of Multiscale Features. This step aggregates *n* multiscale features $F^i \in \mathbb{R}^{C \times H_i \times W_i}$ 215 obtained from the backbone, where i = 0, 1, ..., n - 1 is the scale index (Figure 3). In this step, all the

Figure 4: One step attention and Output Broadcasting in OSMA. 'One-Step' should not be confused with 'single layer'; instead, it refers to all the multiscale features being attended simultaneously through 1×1 layer. Output broadcasting infers the shape of the output feature based on the value of *P*.

features are first divided into non-overlapping grids of size $g_i = 2^i * g$. The total number of grids for all the scales is equal and is given by $N_g = H_i/g_i \times W_i/g_i$.

Let $F_i^i \in \mathbb{R}^{C \times g \times g}$ denotes feature map of j^{th} -grid of F^i , where $j = \{0, 1, ..., N_g\}$. For each scale, we 230 flatten F_i^i in spatial dimension and stack all its features with its corresponding grids in other scales, as 231 232 shown in colors in Figure 3. This results in a matrix $M_i \in \mathbb{R}^{C \times T}$ for each grid. This step is repeated 233 for each non-overlapping grid, and we obtain N_g matrices, referring as $Q \in \mathbb{R}^{N_g \times T \times C}$. The number 234 N_g is absorbed into the batch dimension while processing batched training or inference so that all 235 the matrices can be processed in parallel. This step requires all the feature maps F^{i} to be an integral 236 multiple of g. Hence, we align all the feature maps w.r.t. g_i before performing local aggregation via 237 bilinear interpolation. See Table 5 for ablation on g'.

One Step Multiscale Attention. We perform information fusion by attending to all the multiscale features and their channels simultaneously, thus calling it a one-step multiscale attention. We describe this process for a single grid or M_j , also depicted in Figure 4.

In the attention process, M_j undergoes a 1 × 1 convolution operation with weights $\mathbf{W} \in \mathbb{R}^{d \times T \times 1 \times 1}$ which projects M_j in a *d* dimensional latent space. The 1 × 1 layer combines information from all the multiscale features with a unit stride. In this way, the output of multiscale information becomes richer. Then, we perform a layer normalization over the columns of M_j followed by SiLU activation (See Sec. 4.2).

This process is repeated for feature refinement, and the penultimate 1×1 layer produces P channels. The final layer is a linear projection applied over the columns of \hat{M}_j to refine the column information because a column becomes a feature $\in \mathbb{R}^C$ in the output feature map F_o . The output of this step $\hat{Q} \in \mathbb{R}^{N_g \times P \times C}$ is broadcasted into a feature based on the requirements, as discussed next. **Output Feature Broadcasting.** This step broadcasts \hat{Q} into a feature map F_o of the desired resolution based on the pair $\{g_0, P\}$. For example, for $g_0 = 1$, if we aim to produce a feature F_o of size $H_0 \times W_0$, the value of P will be set to 1, or if a feature F_o of size $2H_0 \times 2W_0$ is required, P can be set to 4, or if a feature F_o of size $H_0/2 \times W_0/2$ is needed, $\{g_0 = 2, P = 1\}$ can be used.

This flexibility of generating F_o of desired resolution allows controlling the encoder's input resolution and, hence, its computations. Meanwhile, multiscale information infusion helps improve DETR accuracy with a slight computational overhead. We have studied various $\{g_0, P\}$ combinations, and our empirical results show that $\{g_0 = 1, P = 1\}$ is the best combination for keeping the resolution of F_o equals to F_0 whereas $\{g_0 = 2, P = 1\}$ is best for reducing the resolution. See Table 5.

260 261

262

223

224

225

226 227

228

229

3.3 CONFIGURING CRED FOR DETRS

We mainly test two important configurations. In all the configurations, OSMA feeds the encoder, whereas CRAM feeds the decoder. OSMA is always fed with $\{F^3, F^4, F^5\}$ output of backbone (Figure 1). Each configuration has different settings for OSMA and input sources for CRAM.

Default. OSMA operates at $\{g_0 = 1, P = 1\}$ i.e. its output resolution is same as F^5 . Whereas CRAM is fed with F^4 . In other words, the encoder runs at half the resolution of the decoder.

Configuration: 'DC×0.25'. OSMA operates at $\{g_0 = 2, P = 1\}$, producing feature map of half the resolution of F^5 or quarter of F^4 or quarter of DC. Whereas CRAM is fed with F^4 . This configuration

evaluates the capability of our Cross-Resolution Encoding-Decoding when the encoder is fed with a very low-resolution map. For comparison, the encoder in baselines is fed with F^5 downsampled to half of its resolution. Since this resolution is 1/4 of DC resolution, it is named as DC×0.25.

Additional Configuration: 'OO' This is similar to the default configuration except the input source to CRAM. We use two instances of OSMA: OSMA_e and OSMA_c. The former operates at $\{g_0 = 1, P = 1\}$, feeding the encoder, while the latter operates at $\{g_0 = 1, P = 4\}$, i.e. in upsampling mode, feeding CRAM with a resolution equivalent to F^4 . This configuration tests the capability of OSMA to fuse multiscale information while increasing the output resolution.

Based on the computational budget, F^3 can also be fed to CRAM in the above configurations. However, due to the large feature resolution, decoder computations also come into play, given highresolution images in the MS-COCO benchmark (Lin et al., 2014). Hence, although we have analyzed its effect (See supplement), we do not use this configuration.

283

4 EXPERIMENTS

284 285

Dataset and Evaluation Metric. Following (Carion et al., 2020; Liu et al., 2022; Li et al., 2022), we
use MS-COCO 2017 benchmark (Lin et al., 2014) for evaluation, having 117k training images and
5k validation images. We use MS-COCO's standard evaluation metric of Average Precision (AP) at
different thresholds and different object scales.

Implementation Details. We plug the proposed CRED into state-of-the-art DN-DETR (Li et al., 2022) for all our experimental evaluations, including ablations. However, to showcase generality, we also adapt CRED into other prominent DETR methods, e.g., Conditional DETR (Meng et al., 2021), DAB-DETR (Liu et al., 2022) etc.

We perform experiments in the 50-epoch setting, widely used for DETRs (Liu et al., 2022; Li et al., 2022; Zhang et al., 2023a). We also show results for the 12-epoch or $1 \times$ (Li et al., 2022) setting to demonstrate accelerated convergence due to the improved DETR design in this paper. The base learning for the backbone is set to 1×10^{-5} while for the transformer, it is set to 1×10^{-4} . For the 12-epoch schedule, the learning rate is dropped by 0.1 at 11^{th} epoch, whereas it is dropped at 40^{th} epoch for the 50-epoch schedule. We use $8 \times$ NVIDIA A40 with a batch size of 16 (2 per-GPU) for training. All the ablations are performed at the $1 \times$ setting.

301 302

4.1 MAIN RESULTS

303
304CRED in DETR: 50-Epoch Setting. We plug CRED into representative DETR frameworks. Table 1
shows that CRED boosts the AP in each DETR pipeline. Compared with the baseline without DC,
CRED introduces a slight overhead of 9G FLOPs with only 1 - 2FPS drop. Compared with recent
IMFA (Zhang et al., 2023a), the overhead is ~ 36% less with an improvement of +4FPS. Our CRED
also performs better regarding runtime than the recent sparse sampling-based method (Zhang et al.,
2023a).

Further, when plugged in Conditional-DETR (Meng et al., 2021), CRED delivers the same accuracy at 50% more FPS than the advanced high-resolution DAB-DETR-DC5-R50 (Li et al., 2022). Similarly, with DN-DETR, CRED delivers the same accuracy at 76% more FPS and 50% fewer FLOPs.

CRED in DETR: 12-Epoch Setting. Evaluations in this setting show that CRED speeds up convergence with slight overhead. From Table 1, CRED with DAB-DETR (Liu et al., 2022) is better than vanilla DAB-DETR by 3.3AP with only a drop of 2FPS. Compared with high-resolution DAB-DETR-DC5, CRED is accurate by 0.4AP at 50% fewer FLOPs and 76% more FPS.

With DN-DETR (Li et al., 2022), CRED achieves beyond 41AP in just 103G FLOPs, implying that CRED improves the convergence speed (See Figure 5), i.e. DN-DETR via CRED achieves the performance of its DC counterpart in just 12 epochs at 50% fewer FLOPs and 76% higher FPS. From the table, CRED can improve the performance of smaller backbones like ResNet18 while delivering real-time performance (> 30FPS). This indicates the utility of CRED where smaller backbones are used due to resource constraints. Hence, by using CRED, detection performance can be boosted

323 **CRED in DETR: DC** \times 0.25 **Setting.** This setting is crucial to show the utility of CRED in DETRs for real-time performance. From Table 1, when encoder resolution is dropped to half (DC \times 0.25)

Table 1: Comprehensive evaluation of CRED when applied in prominent DETR models under different settings,
 i.e., training duration, encoder resolution. Our CRED performs better and faster than DETRs operating at high
 resolution (DC). 'R50: ResNet50 and 'R18: ResNet18 He et al. (2016). Refer to Sec.3.3 for DC×0.25.

Method	#Epochs	#Params	#FLOPs	#FPS	AP	AP ₅₀	AP ₇₅	APs	AP_M	APL
Conditional DETR-R50 Meng et al. (2021)	50	44M	90G	26	40.9	61.8	43.3	20.8	44.6	59.2
Conditional DETR-DC5-R50 Meng et al. (2021)	50	44M	195G	15	43.8	64.4	46.7	24.0	47.6	60.7
Conditional DETR-R50 Meng et al. (2021) + CRED	50	45M	100G	25	44.4	64.6	47.8	25.2	46.9	60.7
DAB DETR-R50 Liu et al. (2022)	50	44M	94G	25	42.2	63.1	44.7	21.5	45.7	60.3
DAB DETR-DC-R50 Liu et al. (2022)	50	45M	202G	13	44.5	65.1	47.7	25.3	48.2	62.3
DAB DETR-R50 Liu et al. (2022) + IMFA Zhang et al. (202 DAB DETR R50 Liu et al. (2022) + CRED	(3a) 50	53M 45M	108G	18	45.5	65.0	49.3	27.3	48.3	61.0
DAB DETR-R50 Elu et al. (2022) + CRED	50	4,5101	1050	23	43.4	04.7	47.4	27.0	40.5	02.2
DN-DETR-R50 Li et al. (2022)	50	44M	94G	25	44.1	64.4	46.7	22.9	48.0	63.4
DN-DETR-DC5-R50 Liu et al. (2022)	50	44M	202G	13	46.3	66.4	49.7	26.7	50.0	64.3
DN-DETR-R50 LT et al. (2022) + CRED	30	43141	1056	23	40.2	05.0	49.0	20.0	50.0	03.2
	12 Epoch Schee	lule								
Conditional DETR-R50 Meng et al. (2021)	12	44M	90G	26	32.4	52.1	33.9	14.2	35.2	48.4
Conditional DETR-R50 + CRED	12	45M	100G	25	36.6	56.2	38.7	18.8	39.5	52.6
DAB DETR-R50 Liu et al. (2022)	12	44M	94G	25	35.1	55.5	36.7	16.2	38.1	52.5
DAB DETR-R50-DC5 Liu et al. (2022)	12	44M	202G	13	38.0	60.3	39.8	19.2	40.9	55.4
DAB DETR-R50 Liu et al. (2022) + IMFA Zhang et al. (202	23a) 12	53M	108G	18	37.3	57.9	39.9	20.8	40.7	52.3
DAB DETR-R50 Liu et al. (2022) + CRED	12	45M	103G	23	38.4	58.4	41.0	20.0	41.8	53.9
DN-DETR-R50 Li et al. (2022)	12	44M	94G	25	38.6	59.1	41.0	17.3	42.4	57.7
DN-DETR-DC5-R50 Li et al. (2022)	12	44M	202G	13	41.7	61.4	44.1	21.2	45.0	60.2
DN-DETR-R50 Li et al. (2022) + CRED	12	45M	103G	23	41.1	60.6	44.0	22.2	44.1	58.9
DAB DETR-R18 Liu et al. (2022)	12	31M	49G	38	29.8	49.0	30.5	10.9	32.5	46.9
DAB DETR-R18 Liu et al. (2022) + IMFA Zhang et al. (202	23a) 12	40M	61G	23	31.2	51.5	32.3	13.0	33.2	49.1
DAB-DETR-R18 Liu et al. (2022) + CRED	12	32M	60G	31	33.5	52.0	35.2	16.7	36.0	46.2
DN-DETR-R18 Li et al. (2022)	12	31M	49G	38	32.5	51.6	33.7	13.5	35.1	49.4
DN-DETR-R18 Li et al. (2022) + CRED	12	32M	60G	31	35.0	54.0	36.9	16.3	37.0	51.4
	DC×0.25 Configu	ration								
DAB DETR-R50 Liu et al. (2022)	12	44M	94G	25	35.1	55.5	36.7	16.2	38.1	52.5
DAB DETR-R50 Liu et al. (2022) DC×0.25	12	44M	80G	26	28.4	48.9	30.0	9.8	31.5	47.0
DAB DETR-R50 Liu et al. (2022) + IMFA Zhang et al. (202	23a) DC×0.25 12	44M	96G	18	33.0	54.2	34.5	16.1	35.3	46.5
DAB-DETR-R50 Liu et al. (2022) + CRED DC×0.25	12	45M	94G	24	37.5	57.9	40.1	18.8	40.7	53.0
DN-DETR-R50 Li et al. (2022)	12	44M	94G	25	38.6	59.1	41.0	17.3	42.4	57.7
DN-DETR-R50 Li et al. (2022) DC×0.25	12	44M	80G	26	31.5	52.7	31.5	10.8	33.7	52.0
DN-DETR-R50 Li et al. (2022) + CRED DC×0.25	12	45M	94G	24	40.0	59.4	42.8	20.7	43.1	56.4
DN-DETR-R50 Li et al. (2022)	50	44M	94G	25	44.1	64.4	46.7	22.9	48.0	63.4
DN-DETR-R50 Li et al. (2022) DC×0.25	50	44M	80G	26	39.9	60.1	41.9	19.2	43.5	59.7
DN-DETR-R50 Li et al. (2022) + CRED DC×0.25	50	45M	94G	24	45.8	64.9	49.1	25.9	49.1	62.8
DAB DETR-R18 Liu et al. (2022)	12	31M	49G	38	29.8	49.0	30.5	10.9	32.5	46.9
DAB DETR-R18 Liu et al. (2022) DC×0.25	12	31M	40G	39	24.0	43.8	25.2	4.7	27.9	42.0
DAB DETR-R18 Liu et al. (2022) + IMFA Zhang et al. (202)	23a) DC×0.25 12	40M	50G	25	27.8	46.9	28.8	14.1	29.5	29.3
DAB-DETR-R18 Eld (t al. (2022) + CRED DC×0.25	12	52141	510		32.2	30.7	34.1	15.7	33.2	
DN-DETR-R18 Li et al. (2022)	12	31M	49G	38	32.5	51.6	33.7	13.5	35.1	49.4
DN-DETR-R18 Li et al. (2022) $DC \times 0.25$	12	32M	40G	34	34.2	53.0	36.2	0.5 16.0	26.0 36.1	45.0 50 (
	12	52111	510		04.2	55.0	50.2	10.0	50.1	20.0
			35						1	
	11			1.1	1	1	1 1		//.	



360

361

362

364

365

366

Figure 5: Convergence plots over MS-COCO validation set. (a) It can be seen that despite having 50% fewer FLOPs and 76% higher FPS, CRED converges similarly to the baseline. (b) DETR with smaller backbone and their DC×0.25 variants in 12 epoch setting. Notice that even in the smaller backbone, CRED-enabled model with and without DC×0.25 have similar accuracy, but this gap is noticeable in the baselines with and without DC×0.25. This strengthens the utility of CRED that encoder input resolution can be aggressively dropped to save computations while having better accuracy.

- of any vanilla DETR pipeline, it degrades the performance while also reducing FLOP requirement.
 However, the degradation in the performance supersedes the reduced FLOPs.
- Whereas CRED in this configuration delivers performance better than the vanilla variant. For example, for vanilla DAB-DETR-R50, the AP drops from 35.1 to 28.4 with DC \times 0.25; however, by using CRED in DC \times 0.25 configuration, we achieve +2.4AP than the vanilla DAB-DETR-R50 in same FLOPs and same FPS. A similar case applies to DN-DETR variants with different backbones.

30														
-	Method	MS	SM	DC	#Epochs	#Params	#FLOPs	#FPS	AP	AP ₅₀	AP ₇₅	AP_S	AP_M	AP_L
1	YOLOS-DeiT-S Fang et al. (2021)				150	28M	172G	_	37.6	57.6	39.2	15.9	40.2	57.3
	Faster-RCNN-FPN-R50 Ren et al. (2015); Lin et al. (2017)	1			108	42M	180G	-	42.0	62.1	45.5	26.6	45.5	53.4
	TSP-FCOS-FPN-R50 Sun et al. (2021b)	1			36	52M	189G	_	43.1	62.3	47.0	26.6	46.8	55.9
	TSP-RCNN-FPN-R50 Sun et al. (2021b) Easter PCNN EDN P101 Pen et al. (2015): Lin et al. (2017)	1			108	64M	188G	_	43.8	63.5	48.5	28.6	46.9	55.7 56.0
	Sparse-RCNN-FPN-R50 Sun et al. (2013), Lin et al. (2017)	1			36	106M	166G	_	45.0	64.1	48.9	28.0	47.6	59.5
	DETR-DC5-R50 Carion et al. (2020)			1	500	41M	187G	16	43.3	63.1	45.9	22.5	47.3	61.1
	SAM-DETR-DC5-R50 Zhang et al. (2022)			1	50	58M	210G	-	43.3	64.4	46.2	25.1	46.9	61.0
	SMCA-DETR-R50 Gao et al. (2021a)	1			50	40M	152G	-	43.7	63.6	47.2	24.2	47.0	60.4
	Conditional DETP, DC5, P50 Mang et al. (2021)	1		,	50	40M	1/3G	15	43.8	64.4	4/./	26.4	47.1	58.0
	Anchor-DETR-DC5-R50 Wang et al. (2021)			·	50	37M	172G	-	44.2	64.4	40.7	24.0	48.2	60.6
	Efficient-DETR-R50 Yao et al. (2021)	1		•	36	32M	159G	_	44.2	62.2	48.0	28.4	47.5	56.6
	DAB-DETR-DC5-R50 Liu et al. (2022)			1	50	44M	202G	13	44.5	65.1	47.7	25.3	48.2	62.3
	SAM-DETR-DC5-R50 Zhang et al. (2022) w/ SMCA			1	50	58M	210G	_	45.0	65.4	47.9	26.2	49.0	63.3
	Conditional DETR-DC5-R101 Meng et al. (2021)			1	50	63M	262G	10	45.0	65.5	48.4	26.1	48.9	62.8
	DN-DETR-RIOT LT et al. (2022) Deformable DAB DETR P50 Zhu et al. (2021)	1	1	1	50	0.5M 41M	174G	12	45.2	03.3 64.7	48.5	24.1	49.1	61.7
	IMFA-DAB-DETR-R50 Zhang et al. (2023a)	ľ	ž		50	53M	108G	18	45.5	65.0	49.3	27.3	48.3	61.6
	DAB DETR-DC5-R101 Liu et al. (2022)			1	50	63M	282G	10	45.8	65.9	49.3	27.0	49.8	63.3
	SAP-DETR-DC5-R50 Liu et al. (2023)				50	47M	197G	12	46.0	65.5	48.9	26.4	50.2	62.6
	DN-DETR-DC5-R50 Li et al. (2022)			1	50	44M	202G	13	46.3	66.4	49.7	26.7	50.0	64.3
	Siamese-DETR-R50 Chen et al. (2023)	1	1		50	41M	173G	-	46.3	64.6	50.5	28.1	50.1	61.5
	SAP-DETR-DC5-R101 Liu et al. (2023)	1	~		50	41M 67M	266G	15	46.7	66.1 66.7	50.6 50.5	29.1 27.9	49.7 51.3	62.2 64.3
	• DAB DETR-R50 Liu et al. (2022) + CRED				50	45M	103G	23	45.4	64.9	49.4	27.0	48.5	62.2
	 DN-DETR-R50 Li et al. (2022) + CRED DC×0.25 				50	45M	94G	24	45.8	64.9	49.1	25.9	49.1	62.8
	• DN-DETR-R50 Li et al. (2022) + CRED				50	45M	103G	23	46.2	65.8	49.8	26.8	50.0	63.5
	• DN-DETR-R50 Li et al. (2022) + CRED-OO				50	45M	105G	23	46.8	66.8	50.5	27.4	50.7	64.0
			12 Ep	och S	chedule									
	DETR-R50 Carion et al. (2020) Defermentie DETR D50 Zhu et al. (2021)		,		12	41M	86G	27	15.5	29.4	14.5	4.3	15.1	26.7
	DAB DETR-R50 Line et al. (2021) DAB DETR-R50 Line et al. (2022) + IMEA Zhang et al. (2023a)		1		12	40M 53M	1/3G 108G	12	37.3	57.9	40.5 30.0	21.1	40.7	52.3
	DAB DETR-DC-R101 Carion et al. (2020)		•	1	12	63M	282G	10	40.3	62.6	42.7	22.2	44.0	57.3
	 DN-DETR-R18 Li et al. (2022) + CRED DC×0.25 				12	32M	51G	34	34.2	53.0	36.2	16.0	36.1	50.0
	• DN-DETR-R18 Li et al. (2022) + CRED				12	32M	60G	31	35.0	54.0	36.9	16.3	37.0	51.4
	 DAB-DETR-R50 Liu et al. (2022) + CRED DC×0.25 				12	45M	94G	24	37.5	57.9	40.1	18.8	40.7	53.0
	• DAB DETR-R50 Liu et al. (2022) + CRED				12	45M	103G	23	38.4	58.4	41.0	20.0	41.8	53.9
	 DN-DETR-R50 Li et al. (2022) + CRED DC×0.25 DN-DETR-R50 Li et al. (2022) + CRED 				12	45M 45M	94G 103G	24 23	40.0 41.1	59.4 60.6	42.8 44.0	20.7 22.2	43.1 44.1	56.4 58.9
			-											

Table 2: Comparison with state-of-the-art object detectors on COCO val 2017. 'MS': Multiscale, 'SM': Sparse
 Multiscale Sampling, and 'DC': Dilated Convolution. FPS is reported at (800 × 1280)

407 CRED vs State-of-the-art. We also compare our CRED with state-of-the-art object detection pipelines with multiscale, high-resolution, and sparse sampling approaches. Table 2 shows the results.
409 From the table, it can be seen that CRED-DETR models are better by a large margin (> 50%) in FLOPs and FPS while delivering accuracy comparable with state-of-the-art methods. Even ResNet-18 based models with CRED show competitive performance with Deformable-DETR (Zhu et al., 2021), DAB-DETR (Liu et al., 2022), IMFA (Zhang et al., 2023a) with a stronger backbone ResNet-101.

CRED w/ ResNet-50 performs better than heavy models, even in DC×0.25 configuration and 12epoch setting. For example, DN-DETR-R50 + CRED DC×0.25 is better than multiscale DeformableDETR-R50 (Zhu et al., 2021) by 2.8AP, 45% fewer FLOPs and 50% higher FPS. Similarly, DNDETR-R50 + CRED is better than DAB-DETR-DC5-R101 by 0.8AP, 63% fewer FLOPs, and 130%
higher FPS. Then DN-DETR-R50 + CRED-OO has 60% fewer FLOPs than (Liu et al., 2023) with the same accuracy. Furthermore, the accuracy can be improved using the latest DETR-training techniques of (Zhao et al., 2024a; Hou et al., 2024), which we leave for future work.

Figure 6 further strengthens our results, that CRED, while delivering comparable performance to the state-of-the-art, have far fewer FLOPs and higher FPS. Also, the results indicate that by utilizing the DC \times 0.25 configuration in CRED, DETRs of real-time speed and high accuracy can be constructed, indicating the huge potential of Cross-Resolution Encoding-Decoding in state-of-the-art DETRs.

424 425

4.2 Ablations

We conduct a comprehensive ablation study on CRED design by using the state-of-the-art DN-DETR (Li et al., 2022) framework and provide insight on the design motivations.

Effect of CRAM and OSMA. We analyze the effect of OSMA (Sec. 3.2) and CRAM (Sec. 3.1).
Table 3 shows the analysis. It can be seen that by using any of OSMA or CRAM into the baseline,
the accuracy improves. By using OSMA alone, AP increases by 1.2, indicating that OSMA produces
better input features or tokens for the encoder. While by using only CRAM, AP improves by 1.8AP.



We also change grid size g. Increasing the grid size $g_0 = 2$ reduces the FLOPs by 2G because T (Figure 3) increases and N_g decreases. However, we observed a reduction in the AP. We hypothesize that this happens because the stage-5 (F^5) feature is the smallest resolution. When more than two

486		Table 5: Ablation of OSMA design.											
487													
/188	ablation	$ F^5$	F^4	$F^3 \mid g_0$	Р	d #Params	#FLOPs	AP	AP ₅₀	AP ₇₅	AP_S	AP_M	AP_L
-100	changing	1	1	1	1	21 45M	101G	40.2	59.5	43.2	21.3	43.2	57.9
489	#scales	1	1	✓ 1	1	21 45M	103G	41.1	60.6	44.0	22.2	44.1	58.9
490	Vary ' $\{g_0, P\}$	}' ′	1	$\begin{array}{c c} \checkmark & 2 \\ \checkmark & 1 \end{array}$	2	21 45M 21 45M	101G 103G	39.7 41.1	59.5 60.6	42.3 44.0	21.0 22.2	42.7 44.1	56.8 58.9
491	Vary '{g ₀ , P	}' 🗸	1	✓ 2	1	21 45M	94G	40.0	59.4	42.8	20.7	43.1	56.4
492	DC×0.25.		1	✓ 4	4	21 45M	92G	39.4	59.0	42.1	19.8	42.7	56.0
	vary a	· · ·	~	V 1	1	40 45M	1140	41.5	01.1	44.Z	22.4	44.0	39.0

497 498

499

500 501

features using $g_0 = 2$ are fused with high-res features, the individual feature at low-resolution loses 495 its chance to interact with the high-resolution features individually because these features already 496 have relatively large receptive fields and carry more information.

Although we are interested in keeping the values of d equal to T, we analyze its effect. We observe that it increases the FLOPs while slightly improving the AP. Hence, based on the computational budget requirements, one can change d to achieve the desired performance and runtime.

Table 6: Effect of LayerNorm Ba et al. (2016) and activations in CRED.

LayerNorm	Activation	#Params	#FLOPs	AP	AP ₅₀	AP ₇₅	AP_S	AP_M	AP_L
1	ReLU	45M	103G	40.6	60.4	43.6	21.6	43.6	57.9
1	SiLU	45M	103G	41.1	60.6	44.0	22.2	44.1	58.9
X	SiLU	45M	103G	39.3	59.1	41.9	20.8	42.5	56.6

Table 7: Ablation of CRAM design.

Input	#Params	#FLOPs	AP	AP ₅₀	AP ₇₅	AP_S	AP_M	AP_L
F^4	45M	103G	41.1	60.6	44.0	22.2	44.1	58.9
F^3	45M	147G	42.1	61.4	44.9	23.3	44.8	59.4

511 512 513

510

514 **CRED Design.** Within the CRED design, we study the effect of different activations and the specified 515 use of layer normalization (Ba et al., 2016). Table 6 shows that using ReLU or removing LayerNorm 516 from CRED decreases accuracy. This justifies the configuration described in the paper.

517 Ablation of CRAM. CRAM is studied by changing its input resolution and source. Table 7 shows 518 that despite feeding the encoder with low resolution, CRAM can effectively transfer the encoder 519 knowledge to the high-resolution feature. By default, we feed resolution equal to F^4 to CRAM. When 520 we feed CRAM with F^3 , the computations in the decoder increase mainly in the cross-attention. 521 Although it improves AP, the rise in FLOPs is notable. Hence, we restrict ourselves to feeding CRAM 522 with resolution up to F^4 .

523 524

5 CONCLUSION

525 526

527 In this work, we present a novel Cross-Resolution Encoding-Decoding (CRED) mechanism to 528 improve the accuracy and runtime of DETR methods. CRED is based on its two novel modules 529 Cross Resolution Attention Module (CRAM) and One Step Multiscale Attention (OSMA). CRAM 530 transfers the knowledge of low-resolution encoder output to a high-resolution feature. While OSMA 531 is designed to fuse multiscale features in a single step and produce a feature map of a desired resolution. With the application of CRED into state-of-the-art DETR methods, FLOPs get reduced by 532 50%, and FPS increases by 76% than the high-resolution DETR at equivalent detection performance.

534 Future Scope & Limitations: CRED with its promising results shows huge potential in real-time and 535 affordable DETRs with high accuracy and high-resolution image processing. There is greater scope 536 for improvements, e.g., fusing CRAM and OSMA for even higher performance or adapting CRED to sparse sampling-based DETRs because the current design can not fuse high-resolution features with sparsely sampled encoder embeddings. In addition, CRED has huge scope in Transformer-538 based semantic or instance segmentation by leveraging its attention transfer to improve runtime for processing high-resolution images because semantic segmentation produces high-resolution outputs.

540	References
542	YOLO-v8. In https://github.com/ultralytics/ultralytics.
543 544 545	Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. <i>arXiv preprint arXiv:1607.06450</i> , 2016.
546 547 548	Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In <i>European Conference on Computer Vision</i> , pp. 213–229. Springer, 2020.
549 550 551 552	Zeren Chen, Gengshi Huang, Wei Li, Jianing Teng, Kun Wang, Jing Shao, Chen Change Loy, and Lu Sheng. Siamese detr. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 15722–15731, 2023.
553 554 555	Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. <i>Advances in neural information processing systems</i> , 34:3965–3977, 2021.
556 557 558 559 560	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. <i>arXiv preprint arXiv:2010.11929</i> , 2020.
561 562 563	Yuxin Fang, Bencheng Liao, Xinggang Wang, Jiemin Fang, Jiyang Qi, Rui Wu, Jianwei Niu, and Wenyu Liu. You only look at one sequence: Rethinking transformer in vision through object detection. Advances in Neural Information Processing Systems, 34:26183–26197, 2021.
564 565 566	Peng Gao, Minghang Zheng, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fast convergence of detr with spatially modulated co-attention. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 3621–3630, 2021a.
567 568 569 570	Peng Gao, Minghang Zheng, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fast convergence of detr with spatially modulated co-attention. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 3621–3630, 2021b.
571 572 573	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 770–778, 2016.
574 575 576	Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415, 2016.
577 578 579	Xiuquan Hou, Meiqin Liu, Senlin Zhang, Ping Wei, and Badong Chen. Salience detr: Enhancing de- tection transformer with hierarchical salience filtering refinement. In <i>Proceedings of the IEEE/CVF</i> <i>Conference on Computer Vision and Pattern Recognition</i> , pp. 17574–17583, 2024.
580 581 582 583	Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. DN-DETR: Accelerate DETR training by introducing query denoising. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 13619–13627, 2022.
584 585 586	Feng Li, Ailing Zeng, Shilong Liu, Hao Zhang, Hongyang Li, Lei Zhang, and Lionel M Ni. Lite detr: An interleaved multi-scale encoder for efficient detr. In <i>Proceedings of the IEEE/CVF conference</i> on computer vision and pattern recognition, pp. 18558–18567, 2023.
587 588 589 590 591	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 <i>Computer Vision–</i> <i>ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings,</i> <i>Part V 13</i> , pp. 740–755. Springer, 2014.
592 593	Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In <i>Proceedings of the IEEE conference on computer</i> <i>vision and pattern recognition</i> , pp. 2117–2125, 2017.

619

626

627

594	Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Oi, Hang Su, Jun Zhu, and Lei Zhang, DAB-
595	DETR: Dynamic anchor boxes are better queries for DETR. <i>arXiv preprint arXiv:2201.12329</i> .
596	2022.
597	

- Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and 598 Alexander C Berg. SSD: Single shot multibox detector. In European conference on computer vision, pp. 21-37. Springer, 2016. 600
- 601 Yang Liu, Yao Zhang, Yixin Wang, Yang Zhang, Jiang Tian, Zhongchao Shi, Jianping Fan, and 602 Zhiqiang He. Sap-detr: bridging the gap between salient points and queries-based transformer detector for fast model convergency. In Proceedings of the IEEE/CVF Conference on Computer 603 Vision and Pattern Recognition, pp. 15539–15547, 2023. 604
- 605 Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and 606 Jingdong Wang. Conditional DETR for fast training convergence. In Proceedings of the IEEE/CVF 607 International Conference on Computer Vision, pp. 3651–3660, 2021.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object 609 detection with region proposal networks. Advances in neural information processing systems, 28: 610 91-99, 2015. 611
- 612 Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei 613 Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 614 pp. 14454–14463, 2021a. 615
- 616 Zhiqing Sun, Shengcao Cao, Yiming Yang, and Kris M Kitani. Rethinking transformer-based set 617 prediction for object detection. In Proceedings of the IEEE/CVF international conference on 618 computer vision, pp. 3611-3620, 2021b.
- Yingming Wang, Xiangyu Zhang, Tong Yang, and Jian Sun. Anchor detr: Query design for 620 transformer-based object detection. arXiv preprint arXiv:2109.07107, 2021. 621
- 622 Zhuyu Yao, Jiangbo Ai, Boxun Li, and Chi Zhang. Efficient detr: improving end-to-end object detector with dense prior. arXiv preprint arXiv:2104.01318, 2021. 623
- 624 Gongjie Zhang, Zhipeng Luo, Yingchen Yu, Kaiwen Cui, and Shijian Lu. Accelerating detr conver-625 gence via semantic-aligned matching. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 949–958, 2022.
- Gongjie Zhang, Zhipeng Luo, Zichen Tian, Jingyi Zhang, Xiaoqin Zhang, and Shijian Lu. Towards 628 efficient use of multi-scale features in transformer-based object detectors. In Proceedings of the 629 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6206–6216, 2023a. 630
- 631 Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel Ni, and Heung-Yeung 632 Shum. DINO: DETR with improved denoising anchor boxes for end-to-end object detection. 633 In The Eleventh International Conference on Learning Representations, 2023b. URL https: 634 //openreview.net/forum?id=3mRwyG5one.
- 635 Chuyang Zhao, Yifan Sun, Wenhao Wang, Qiang Chen, Errui Ding, Yi Yang, and Jingdong Wang. Ms-636 detr: Efficient detr training with mixed supervision. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 17027-17036, 2024a. 638
- Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing 639 network. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 640 2881-2890, 2017. 641
- 642 Yian Zhao, Wenyu Lv, Shangliang Xu, Jinman Wei, Guanzhong Wang, Qingqing Dang, Yi Liu, 643 and Jie Chen. Detrs beat yolos on real-time object detection. In Proceedings of the IEEE/CVF 644 Conference on Computer Vision and Pattern Recognition, pp. 16965–16974, 2024b.
- 645 Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable {detr}: 646 Deformable transformers for end-to-end object detection. In International Conference on Learning 647 *Representations*, 2021. URL https://openreview.net/forum?id=qZ9hCDWe6ke.

DETECTION VISUALIZATIONS ON MS-COCO VALIDATION SET

A QUALITATIVE COMPARISON WITH DN-DETR-R50 (LI ET AL., 2022) AND DN-DETR-DC5-R50 LI ET AL. (2022)

To put the qualitative results in context, we have compared CRED-DETR with baseline DN-DETR with low resolution and DC5 setting. Please refer to Figure 7.



(AP = 46.8, FLOPs = 105G, FPS = 23)

DN-DETR (Li et al., 2022) (AP = 44.4, FLOPs = 100G, FPS = 25)

DN-DETR-DC5 (Li et al., 2022) (AP = 46.3, FLOPs = 202G, FPS = 13)

Figure 7: It can be seen that the proposed CRED-DETR can detect objects with high accuracy. *Top:* CRED-DETR detects almost all of the persons standing in the top (4th row) in the image, including the tie. In contrast, baselines struggle to achieve the same as indicated in rows 1-3 in the top image. *Middle:* CRED-DETR can detect the instances of person visible with very small field of view. However, the DC5 variant of the baseline can not be detected. *bottom:* CRED-DETR detects more number of persons. The baseline DN-DETR misclassifies the front region of the train as a suitcase, while CRED-DETR does not. The same is verified with empirical results provided in the paper.

B QUALITATIVE ANALYSIS

Below we have visualized detection results from CRED-DETR from MS-COCO 2017 validation set. The images covers wide range of object from tiny, small to large. It can be seen that CRED-DETR, despite running at low resolution in the encoder, is able to detect all categories of object.



