

Transformer-based method with Deformable Convolution and Central Difference Convolution for Object Detection

基於 Transformer 網路及可變形卷積與中心差分卷積
的物件偵測方法

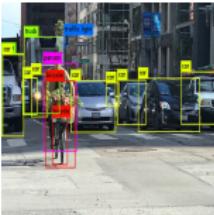
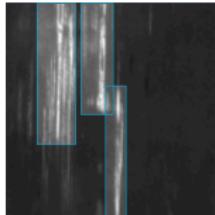
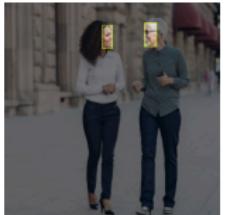
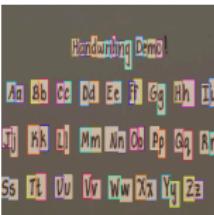
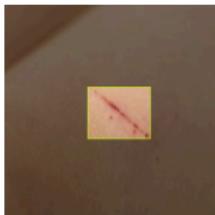
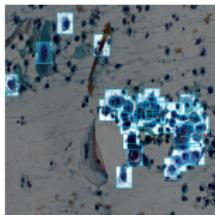
劉律奇

June 21, 2024

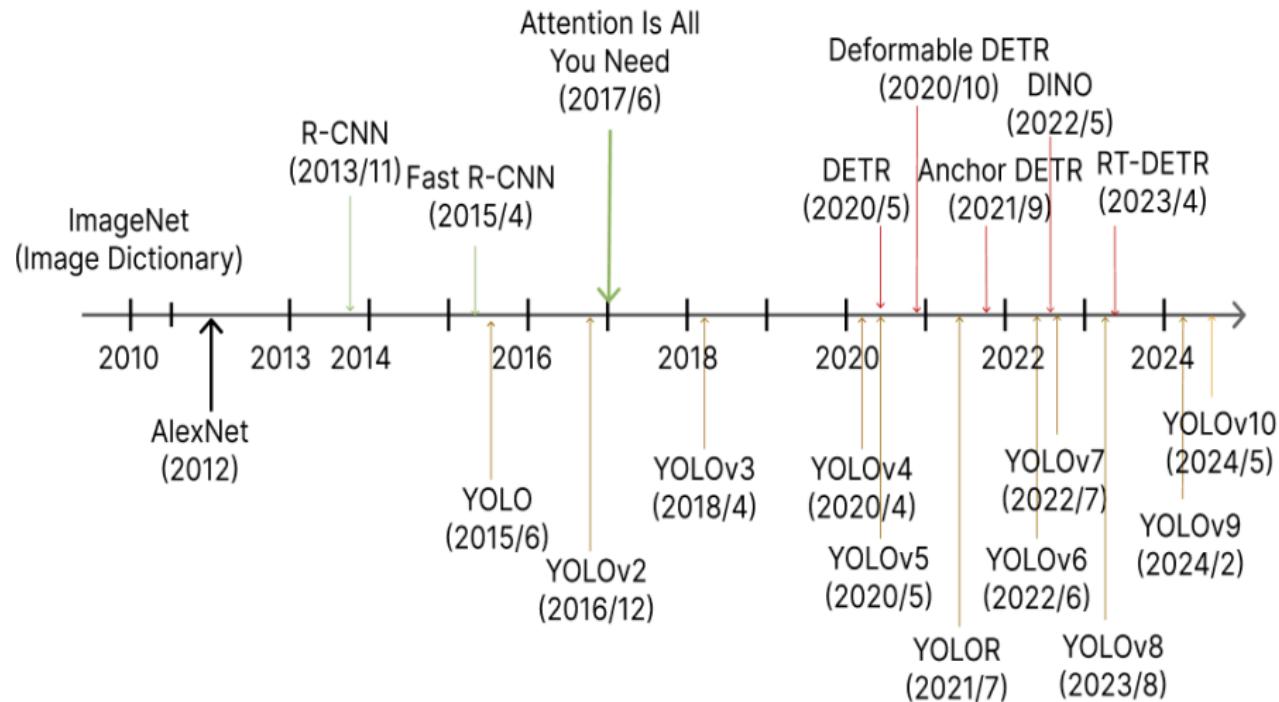
Table of contents

1. Introduction of Object Detection
2. CNN-based Detector YOLO
3. Transformer-based Detector DETR
4. RT-DETR Structure
5. DeformConv and CDC
6. Experimental Result

Object Detection Application

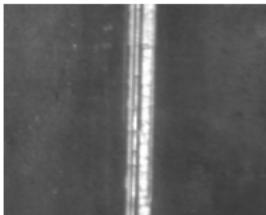
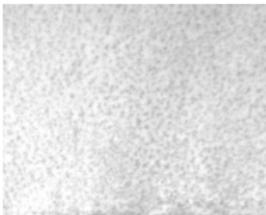
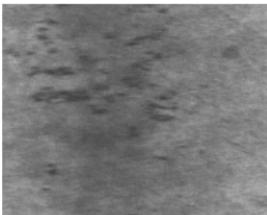
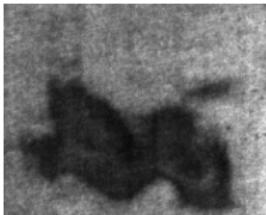
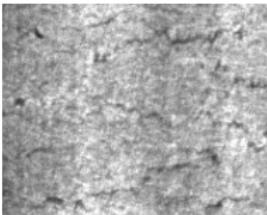
			
Vehicle	Defect	Safety Helmet	Face
			
Hand Writting	Wound	Pap-smear	Pills

History of Object Detection



Research Objective

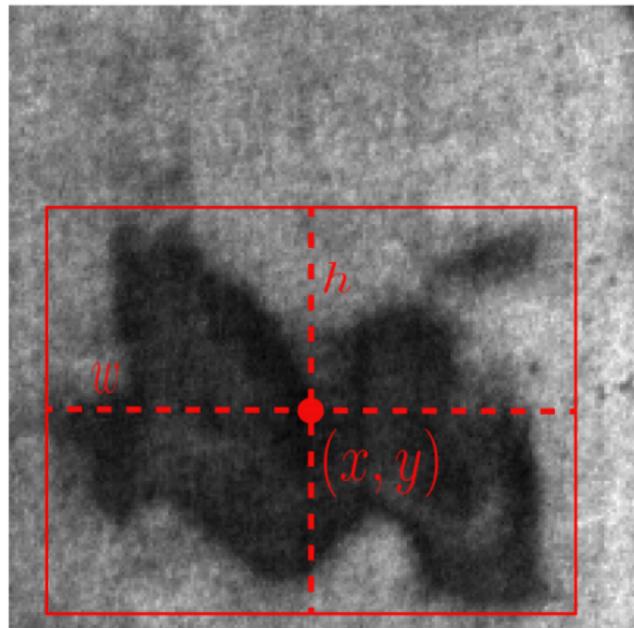
- Improve the detection capabilities of real-time detection transformer (RT-DETR).
- Performing defect detection¹with the RT-DETR model.

		
Scratches	Pitted surface	Rolled in scale
		
Patches	Inclusion	Crazing

¹Kechen Song and Yunhui Yan, "A noise robust method based on completed local binary patterns for hot-rolled steel strip transformer-based detector

Introduction of Object Detection

Object detection is a branch of computer vision that empowers computers to identify and locate objects within images or videos.



$$\text{model}(\text{Image}) = \begin{bmatrix} \vdots \\ P_{\text{Patches}} \\ \vdots \\ x \\ y \\ w \\ h \end{bmatrix}$$

CNN-based Detector YOLO

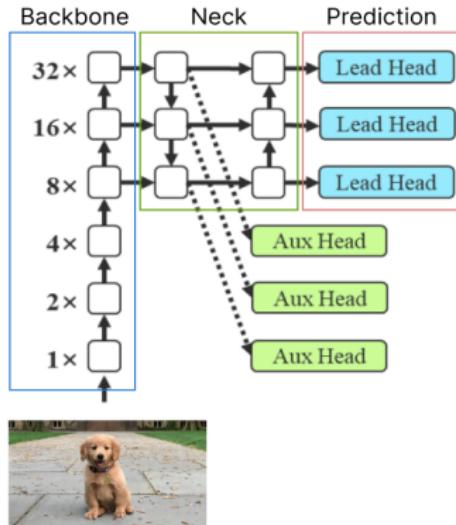


Figure: The main structure of YOLO.³

- **Backbone**
Extracting and encoding features from the image
- **Neck**
Further transforming and refining the features
- **Prediction**
Designed to produce the task-specific prediction

³Chien-Yao Wang, I-Hau Yeh, and Hong-Yuan Mark Liao. *YOLOv9: Learning What You Want to Learn Using Programmable Gradient Information*. 2024. arXiv: 2402.13616 [cs.CV].

YOLO Detecting Head

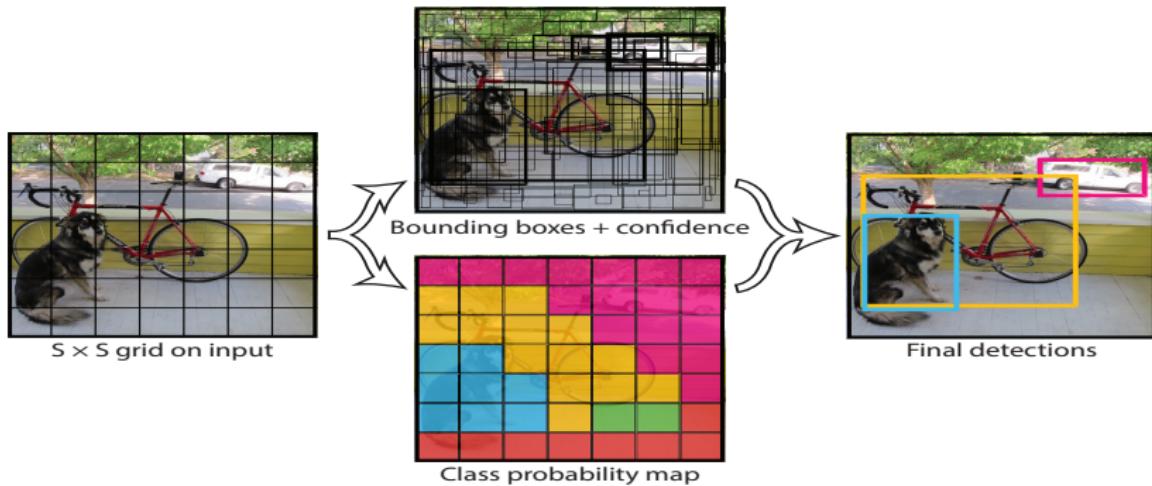


Figure: Detecting Head. YOLO detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.⁵

⁵ Joseph Redmon et al. You Only Look Once: Unified, Real-Time Object Detection. 2016. arXiv: 1506.02640 [cs.CV]. ↗ ↘ ↙ ↛

YOLO Detecting Head



Figure: A prediction grid cell with two bounding boxes.

$$\left\{ \begin{array}{l} \vdots \\ P_{dog} \\ c_1 \\ x_1 \\ y_1 \\ w_1 \\ h_1 \\ \vdots \\ P_{dog} \\ c_2 \\ x_2 \\ y_2 \\ w_2 \\ h_2 \end{array} \right\}$$

Bounding Box1

Bounding Box2

YOLO Detecting Head

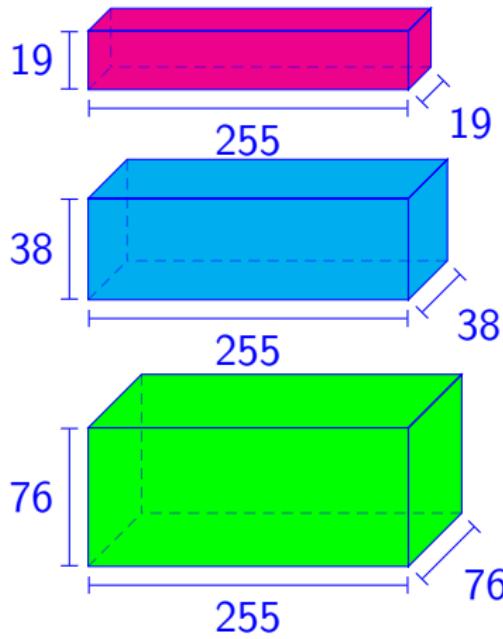
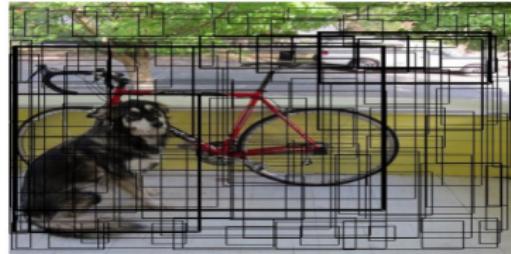


Figure: The Detecting head of YOLOv4⁶.



The cumulative total across all prediction boxes is 22,743.

$$\begin{aligned} 255 &= 3 \times (80 + 1 + 4) \\ &= 3(\text{categories} \\ &\quad + \text{confidence} + \text{box}) \end{aligned}$$

⁶ Alexey Bochkovski, Chien-Yao Wang, and Hong-Yuan Mark Liao. *YOLOv4: Optimal Speed and Accuracy of Object Detection*. 2020. arXiv: 2004.10934 [cs.CV].

Transformer-based Detector DETR

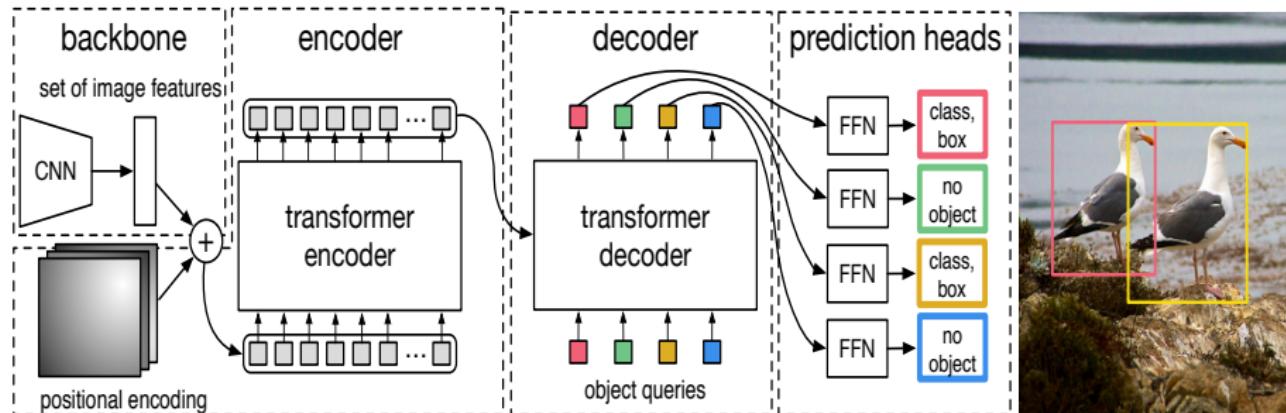


Figure: End-to-End Object Detection with Transformers (DETR)⁷.

- Backbone
- Transformer
- Prediction

⁷Nicolas Carion et al. *End-to-End Object Detection with Transformers*. 2020. arXiv:2005.12872 [cs.IV].

How does DETR predict objects?

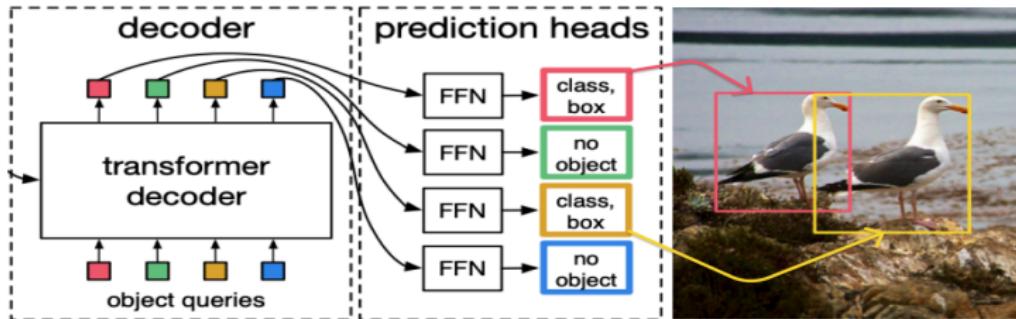


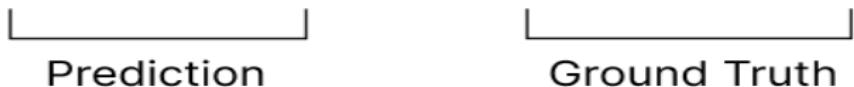
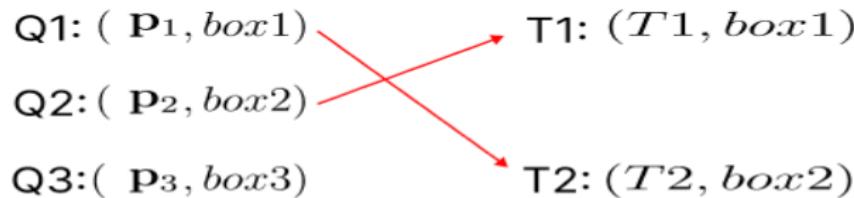
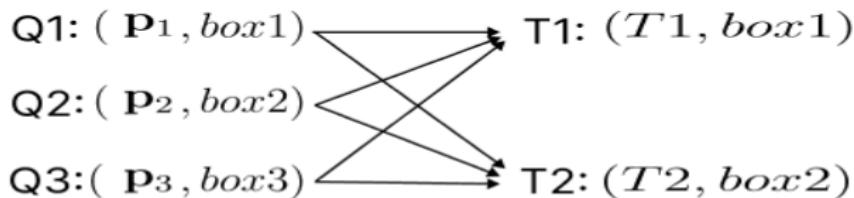
Figure: Each object query is fed into the decoder to predict the objects in the image.

Object query

A small, fixed number of learned positional embeddings⁸.

⁸Carion et al., see n. 7.

Training



Training

no object

The diagram illustrates a training step for a transformer-based detector. It consists of two tables side-by-side, connected by a horizontal arrow pointing from left to right.

Left Table:

	T1	T2	\emptyset
Q1	1.3	1.92	0
Q2	4.5	6.7	0
Q3	9.19	7.41	0

Right Table:

	T1	T2	\emptyset
Q1	1.3	1.92	0
Q2	4.5	6.7	0
Q3	9.19	7.41	0

In the right table, the values 1.92 and 4.5 have been circled in red, indicating they are the target or predicted values for the 'no object' class.

Figure: Find the pair minimize loss cost 6.42

Inference

Set threshold 0.5 as C :

$$Q_1 = (\mathbf{p}_1, \text{box1}) \quad \text{where} \quad \mathbf{p}_1 = (p_{1,1}, p_{1,2}, \dots, p_{1,n})$$

$$Q_2 = (\mathbf{p}_2, \text{box2}) \quad \text{where} \quad \mathbf{p}_2 = (p_{2,1}, p_{2,2}, \dots, p_{2,n})$$

$$Q_3 = (\mathbf{p}_3, \text{box3}) \quad \text{where} \quad \mathbf{p}_3 = (p_{3,1}, p_{3,2}, \dots, p_{3,n})$$

$$Q = \{Q_1, Q_2, Q_3\}$$

The prediction is valid if the maximum probability in any vector exceeds C :

$$\text{Prediction} = \left\{ Q_i \mid \max_j \{p_{i,j}\} > C \quad \text{for} \quad i \in \{1, 2, 3\} \right\}$$

The comparison of YOLO and DETR

YOLO uses a dense prediction approach, while DETR utilizes a set prediction method for object detection. $\text{Head}(S \times S \times B) \gg N$

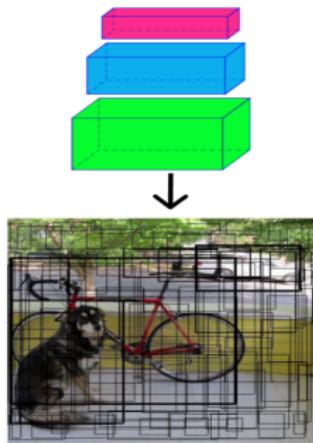


Figure: The number of YOLO Prediction Head($S \times S \times B$).

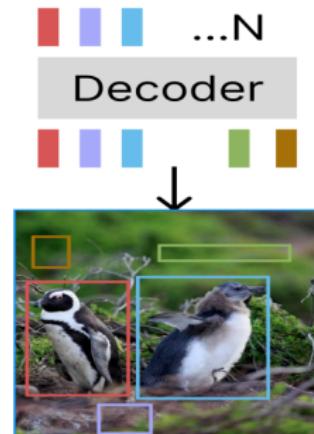


Figure: The number of DETR Prediction N .

RT-DETR

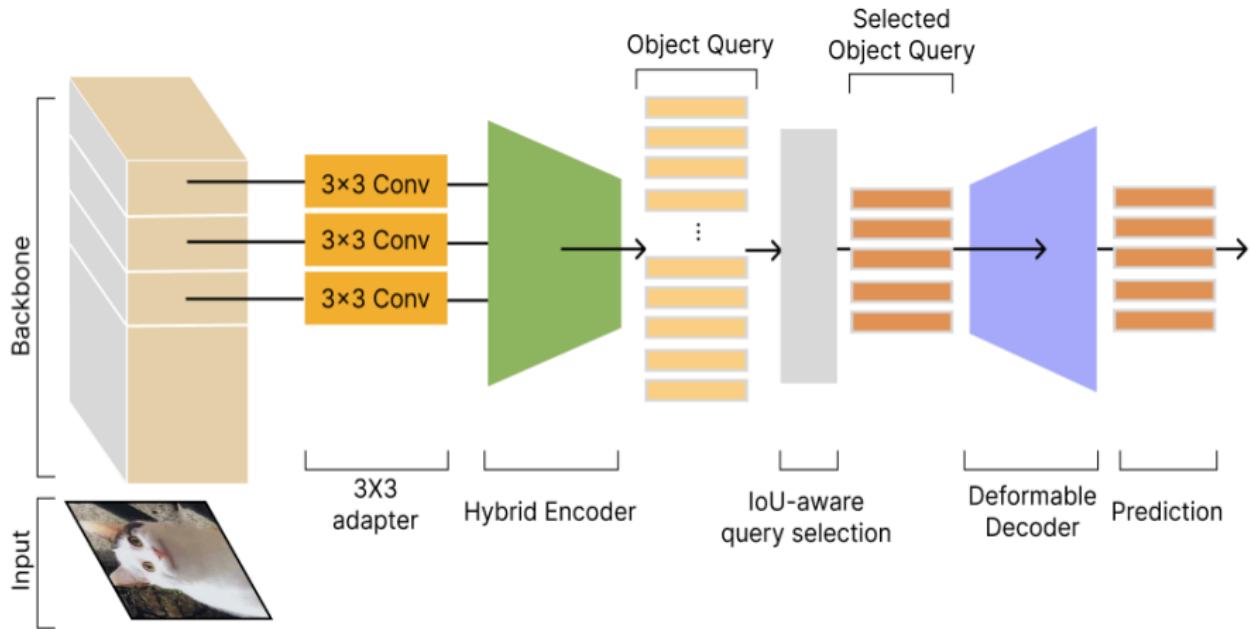


Figure: Real-Time Detection Transformer (RT-DETR)⁹

⁹Yian Zhao et al. "Detrs beat yolos on real-time object detection". In: *arXiv preprint arXiv:2304.08069* (2023).

RT-DETR Ours

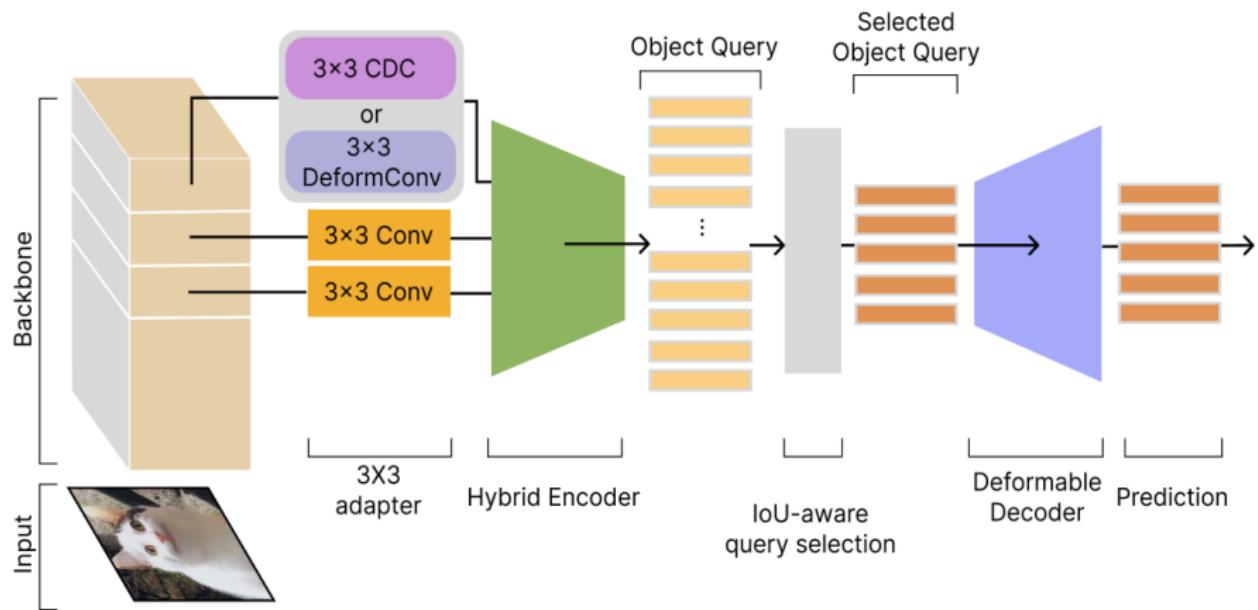
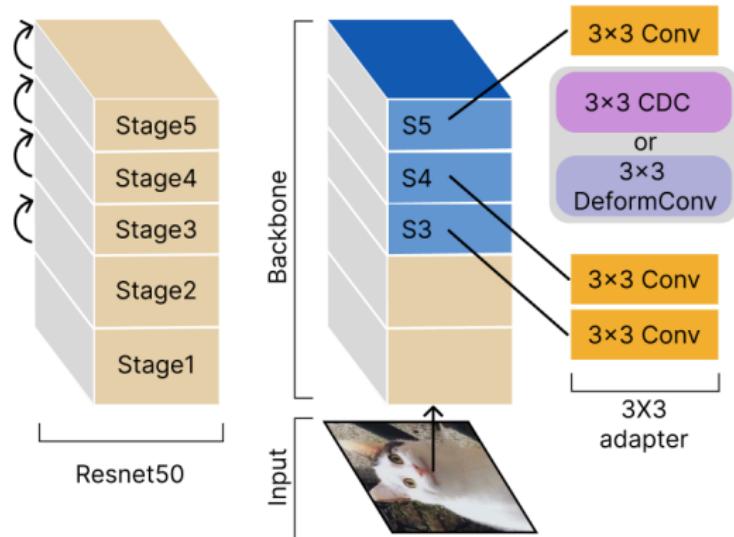


Figure: RT-DETR with Central Difference Convolution (CDC) and Deformable Convolution (DeformConv)

RT-DETR Backbone



- ResNet-50 pretrained on ImageNet.
- CDC and DeformConv at S_5 adapter.

Figure: ResNet-50 as Backbone

Vanilla Convolution

$$y(p_0) = \sum_{p_n \in R} w(p_n) \cdot x(p_0 + p_n)$$

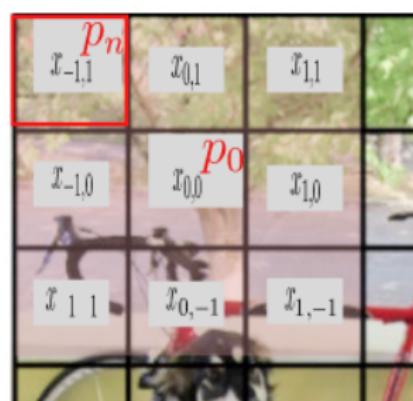
R

p_n	(0, 1)	(1, 1)
(-1, 0)	p_0	(1, 0)
(-1, -1)	(0, -1)	(1, -1)

\mathcal{W}

p_n	$w_{0,1}$	$w_{1,1}$
$w_{-1,0}$	$w_{0,0}$	$w_{1,0}$
$w_{-1,-1}$	$w_{0,-1}$	$w_{1,-1}$

\mathcal{X}



Central Difference Convolution (CDC)

$$y(p_0) = \sum_{p_n \in R} w(p_n) \cdot (x(p_0 + p_n) - x(p_0))$$

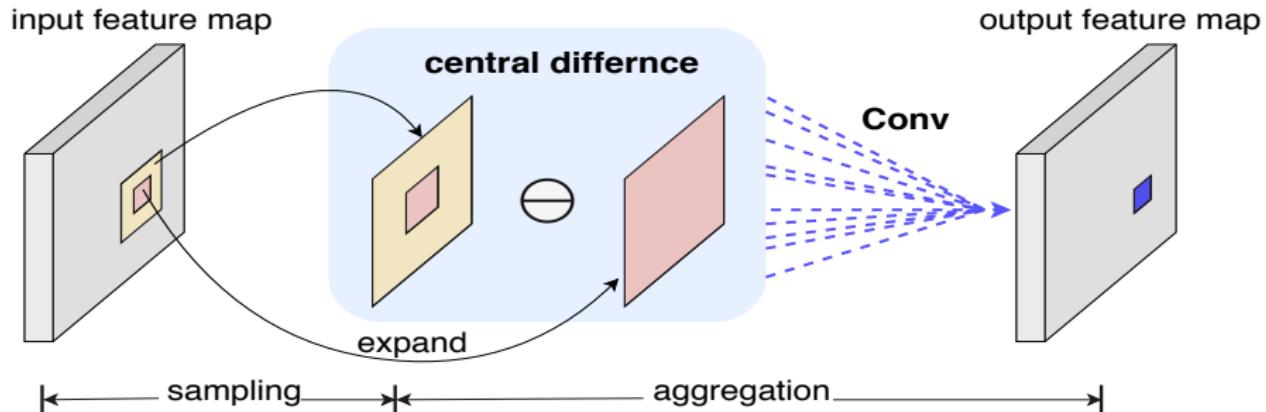


Figure: Central Difference Convolution¹⁰

¹⁰ Zitong Yu et al. *Searching Central Difference Convolutional Networks for Face Anti-Spoofing*. 2020. arXiv: 2003.04092 [cs.CV].

Central Difference Convolution (CDC)

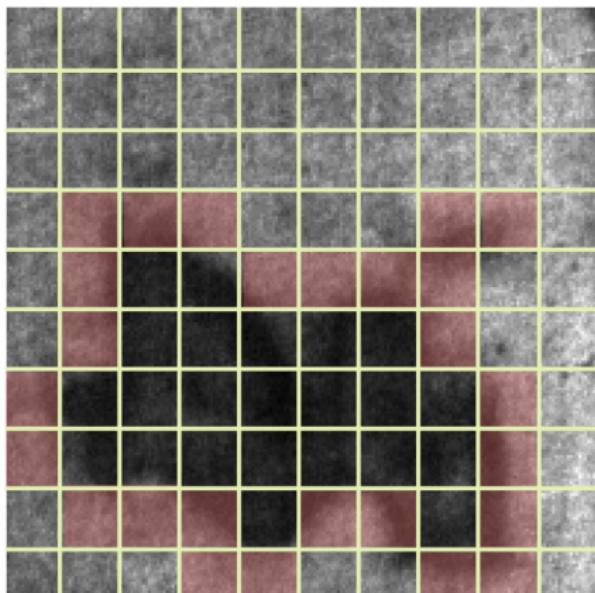
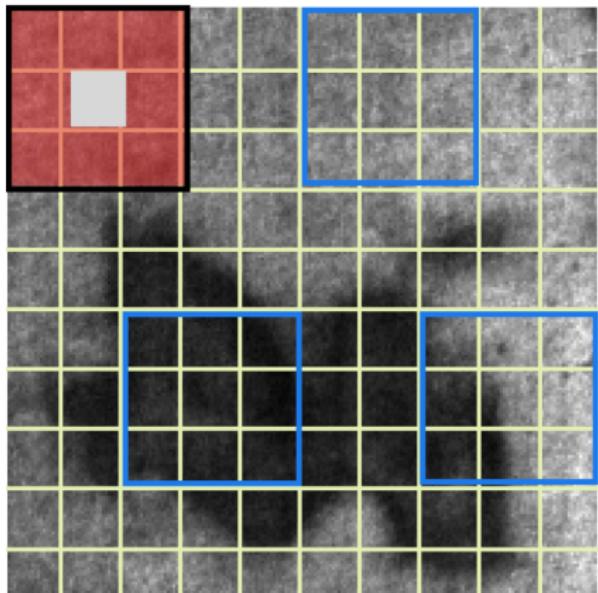


Figure: CDC kernel slid over the defect object.

Deformable Convolution (DeformConv)

$$y(p_0) = \sum_{p_n \in R} w(p_n) \cdot x(p_0 + p_n + \Delta p_n)$$

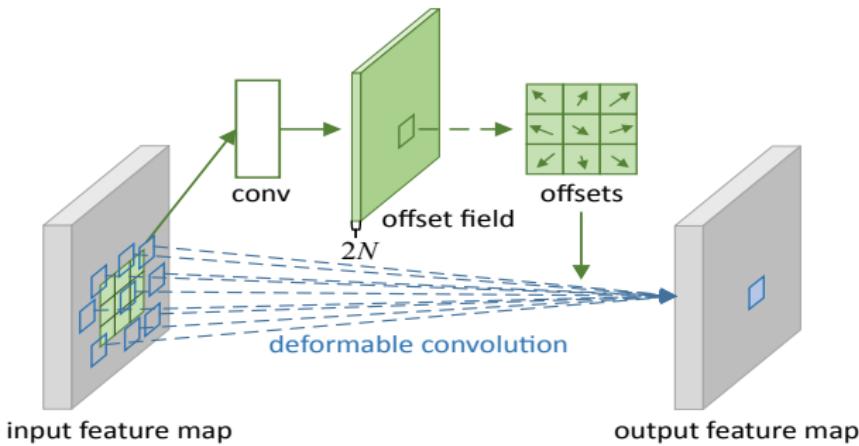


Figure: Deformable Convolution¹¹

¹¹ Jifeng Dai et al. "Deformable convolutional networks". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 764–773.

Deformable Convolution (DeformConv)

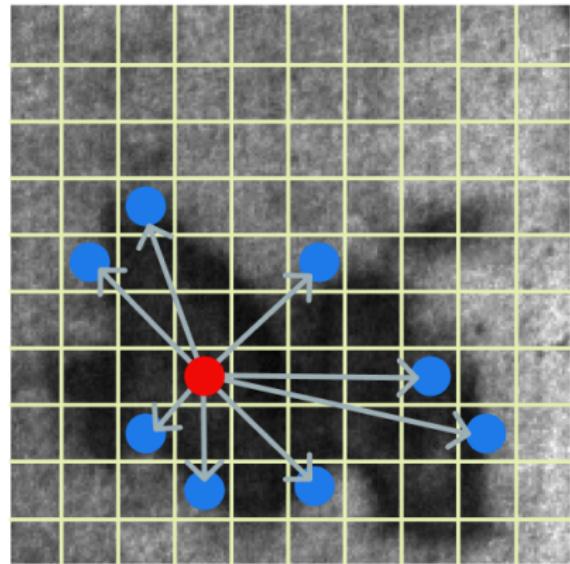
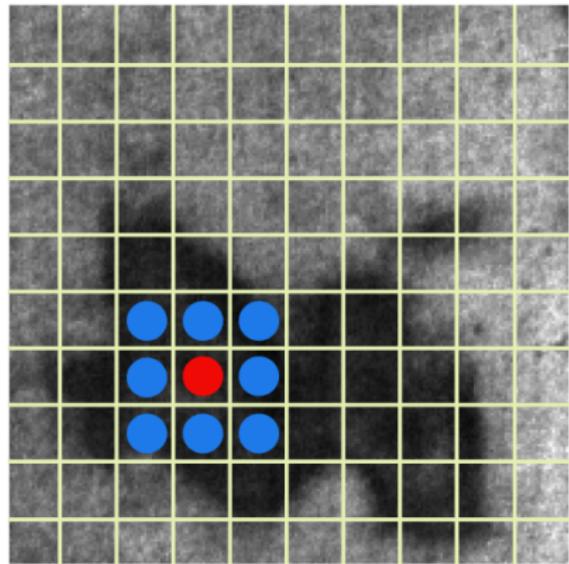


Figure: DeformConv kernel slid over the defect object.

Efficient Hybrid Encoder

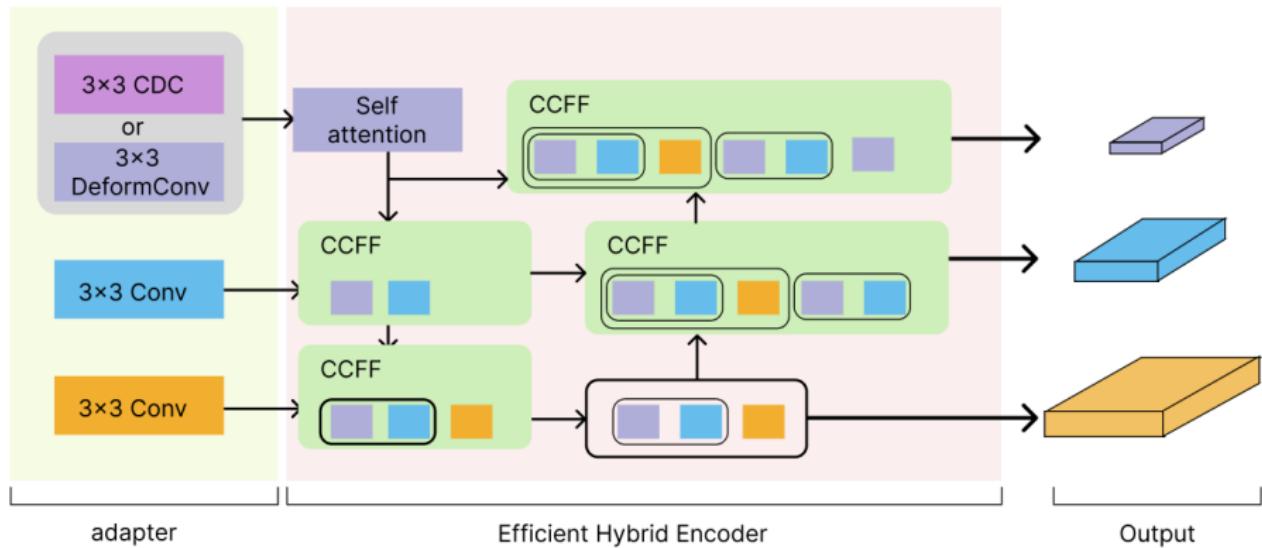
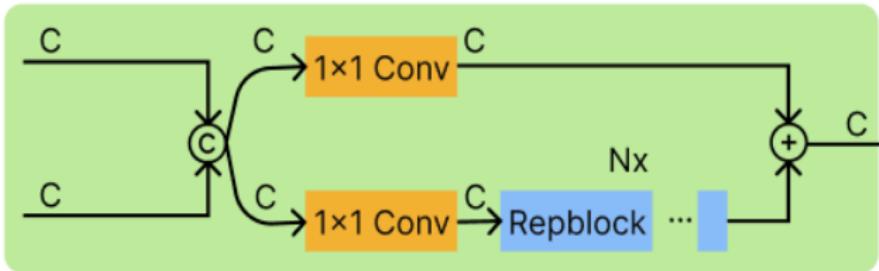


Figure: Efficient Hybrid Encoder¹²

¹²Zhao et al., see n. 9.

Efficient Hybrid Encoder-CCFF



C : Channel © : Concatenate ⊕ : Element-wise add

Figure: Cross-scale Feature-fusion Module (CCFF)¹³

¹³Zhao et al., see n. 9.

IoU-Aware Query Selection

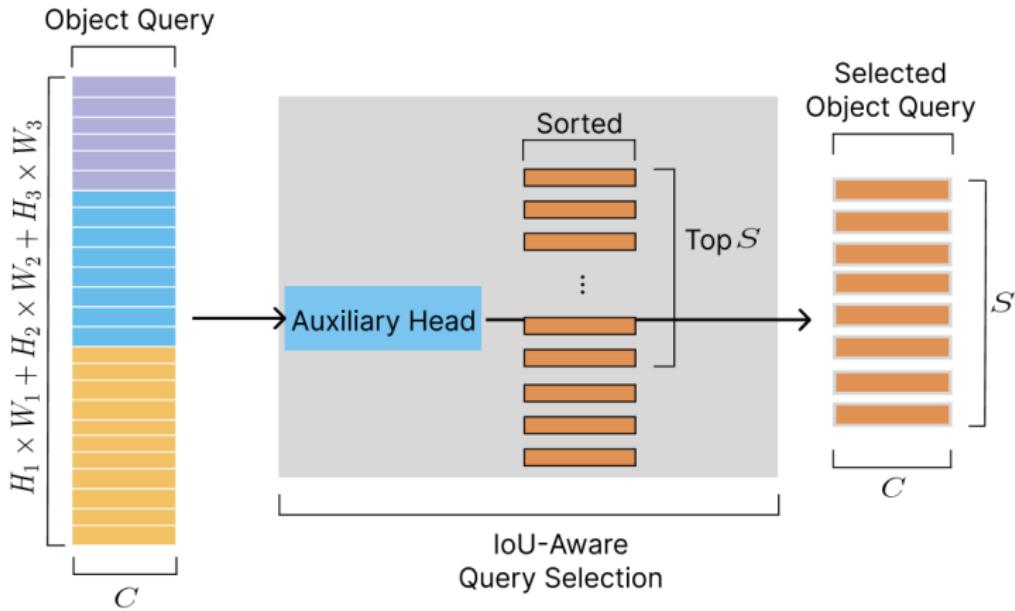
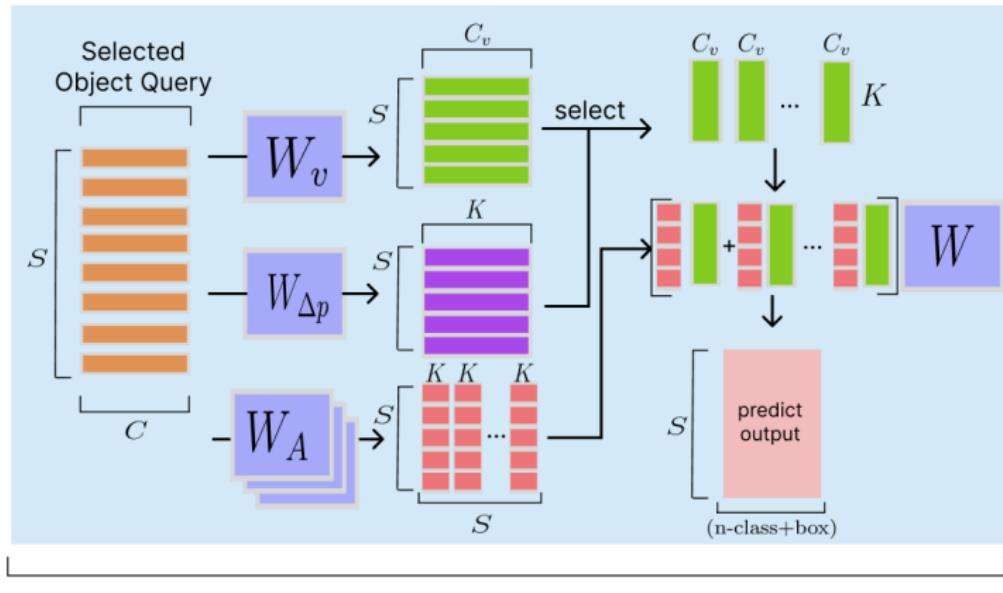


Figure: Selected Top S as Object Queries¹⁴

¹⁴Zhao et al., see n. 9.

Deformable Decoder



Decoder with Deformable Attention

Figure: Decoder with Deformable Attention¹⁵.

¹⁵Xizhou Zhu et al. "Deformable detr: Deformable transformers for end-to-end object detection". In: *arXiv preprint arXiv:2010.04159* (2020).

Loss Function

The loss function employed in this study is identical to that used in the RT-DETR model. The loss function can be decomposed into two components: box loss and classification loss¹⁶, with a weight ratio of 5:1.

$$\begin{aligned}\mathcal{L} &= 5\mathcal{L}_{box} + \mathcal{L}_{class} \\ &= 5\mathcal{L}_1 + VFL\end{aligned}$$

¹⁶ Haoyang Zhang et al. VarifocalNet: An IoU-aware Dense Object Detector. 2021. arXiv: 2008.13367 [cs.CV].



Prediction and Ground Truth

The prediction of a object query and ground truth.

$$\begin{array}{ll} \text{classification} \left\{ \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ \hat{x} \\ \hat{y} \\ \hat{w} \\ \hat{h} \end{bmatrix} \right. & = d \\ \text{box} \left\{ \underbrace{\begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ \hat{x} \\ \hat{y} \\ \hat{w} \\ \hat{h} \end{bmatrix}}_{\text{Prediction}} \right. & \\ \\ \text{classification} \left\{ \begin{bmatrix} 0 \\ 1 \\ 0 \\ x \\ y \\ w \\ h \end{bmatrix} \right. & = t \\ \text{box} \left\{ \underbrace{\begin{bmatrix} 0 \\ 1 \\ 0 \\ x \\ y \\ w \\ h \end{bmatrix}}_{\text{Ground truth}} \right. & \end{array}$$

L1 Loss Function (MAE)

The box loss of a object query.

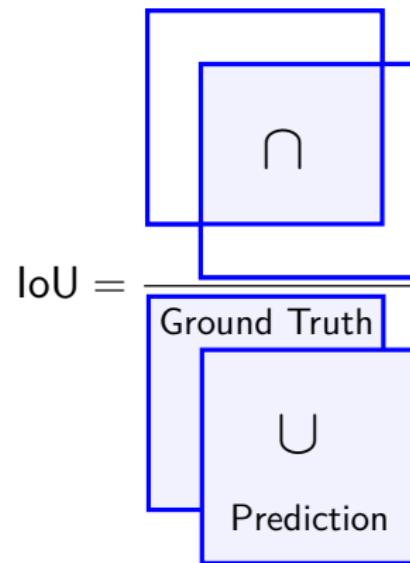
$$d^{box} = \begin{bmatrix} \hat{x} \\ \hat{y} \\ \hat{w} \\ \hat{h} \end{bmatrix} \quad t^{box} = \begin{bmatrix} x \\ y \\ w \\ h \end{bmatrix}$$

$$\mathcal{L}_1 = \sum_{i=1}^4 \frac{|d_i^{box} - t_i^{box}|}{4}$$

Varifocal Loss (VFL)

Preprocess ground truth.

$$t^{class} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \rightarrow t^{class'} = \begin{bmatrix} 0 \\ IoU \\ 0 \end{bmatrix}$$



Varifocal Loss (VFL)

The classification loss of a object query.

$$d^{class} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix} \quad t^{class'} = \begin{bmatrix} 0 \\ \text{IoU} \\ 0 \end{bmatrix} = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix}$$

$$\text{VFL}(p, q) = \begin{cases} -q(q\log(p) + (1-q)\log(1-p)) & q > 0, \\ -\alpha p^\gamma \log(1-p) & q = 0 \end{cases}$$

where the value of $\alpha = 0.2$ and $\gamma = 2$ in the experiment¹⁷.

¹⁷Zhao et al., see n. 9.

Experimental Result NEU-DET

Table: The experiments are conducted on the NEU-DET dataset.

Model	AP0.5-0.95	AP0.5	AP0.75	AP-medium	AP-large
RT-DETR	0.305	0.563	0.300	0.259	0.333
RT-DETR+CDC	0.287	0.532	0.260	0.255	0.311
RT-DETR+DeformConv	0.273	0.515	0.247	0.269*	0.291
Model	AR0.5-0.95	AR0.5	AR0.75	AR-medium	AR-large
RT-DETR	0.253	0.518	0.652	0.602	0.678
RT-DETR+CDC	0.245	0.51	0.646	0.6	0.674
RT-DETR+DeformConv	0.232	0.51	0.646	0.58	0.679*

Experimental Result COCO

Table: The experiments are conducted on the COCO dataset Average Precision.

Model	AP0.5-0.95	AP0.5	AP0.75
RT-DETR	0.517	0.7	0.558
RT-DETR+CDC	0.517	0.699	0.558
RT-DETR+DeformConv	0.464	0.634	0.504
Model	AP-small	AP-medium	AP-large
RT-DETR	0.343	0.562	0.692
RT-DETR+CDC	0.331	0.563*	0.69
RT-DETR+DeformConv	0.295	0.506	0.621

Experimental Result COCO

Table: The experiments are conducted on the COCO dataset Average Recall.

Model	AR0.5-0.95	AR0.5	AR0.75
RT-DETR	0.387	0.648	0.709
RT-DETR+CDC	0.385	0.646	0.711*
RT-DETR+DeformConv	0.366	0.616	0.683
Model	AR-small	AR-medium	AR-large
RT-DETR	0.541	0.752	0.876
RT-DETR+CDC	0.534	0.757*	0.874
RT-DETR+DeformConv	0.484	0.73	0.866

Conclusion

- Research explores CDC and DeformConv on RT-DETR.
- In defect detection tasks, DeformConv might be the better choice due to its higher recall.
- For general object detection, CDC improved RT-DETR's ability to detect medium-sized objects.
- CDC and DeformConv modules enhance RT-DETR for specific object sizes and applications.

Appendix

- Non-Maximum Suppression (NMS)
- Deformable Attention
- Hungarian Algorithm
- CDC
- Auxiliary Head
- Up/Down Sampling

Non-Maximum Suppression

Algorithm 1: Non-Maximum Suppression

Input: $B = \{b_1, \dots, b_N\}, S = \{s_1, \dots, s_N\}, N_t$

B is the list of initial detection boxes

S contains corresponding detection scores

N_t is the NMS threshold

begin

$D \leftarrow \{\};$

while $B \neq \text{empty}$ **do**

$m \leftarrow \arg \max S; M \leftarrow b_m;$

$D \leftarrow D \cup M; B \leftarrow B - M;$

for b_i **to** B **do**

 | $s_i \leftarrow s_i(1 - iou(M, b_i))$

end

end

return D, S

end

Self Attention and Deformable Attention

Let

$$I^{s \times c}, \{W_q, W_v, W_k\} \in \mathbb{R}^{c \times c_v}$$

Self Attention

$$Q = IW_q, K = IW_k, V = IWv \in \mathbb{R}^{s \times c_v}$$

$$A^{c_v \times c_v} = K^T Q$$

$$O^{s \times c_v} = VA$$

Self Attention and Deformable Attention

Deformable Attention

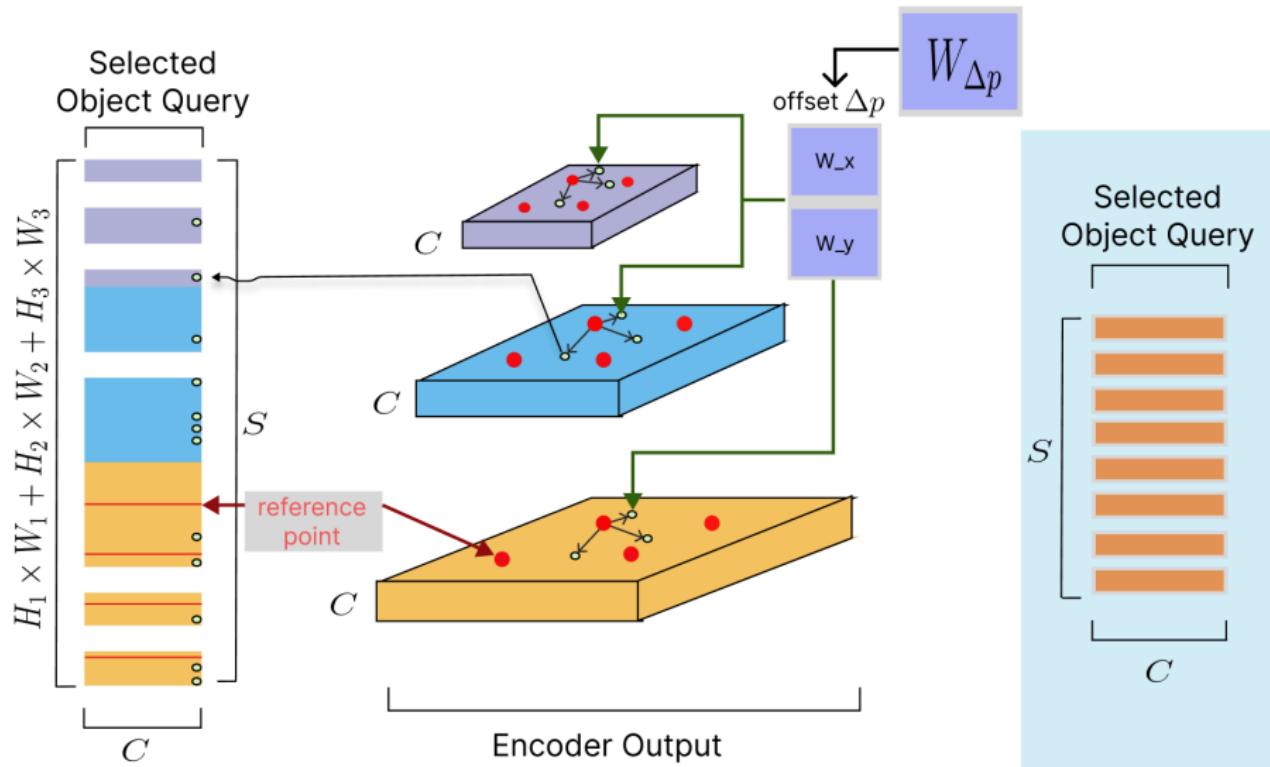
$$V^{s \times c_v} = IW_v^{c \times c_v} \quad \Delta_p = IW_{\Delta_p}^{c \times k} \quad A_i^{s \times k} = IW_{A_i}^{s \times k}$$

$$V \xrightarrow{\text{select by } \Delta_p} V_1, V_2, \dots, V_s^{k \times c_v}$$

$$O^{s \times c_v} = [A_1 V_1 + A_2 V_2 + \dots + A_s V_s] W^{c_v \times c_v}$$

where $i \in [1, 2, 3, \dots, s]$

Deformable Attention More Detail



Deformable Attention More Detail

Deformable Attention

$$V^{s \times c_v} = IW_v^{c \times c_v} \quad \Delta_p = IW_{\Delta_x}^{c \times pk}, IW_{\Delta_y}^{c \times pk} \quad A_i^{s \times pk} = IW_{A_i}^{s \times pk}$$

$$V \xrightarrow{\text{select by } \Delta_p} V_1, V_2, \dots, V_s^{pk \times c_v}$$

$$O^{s \times c_v} = [A_1 V_1 + A_2 V_2 + \dots + A_s V_s] W^{c_v \times c_v}$$

where $i \in [1, 2, 3, \dots, s]$

Hungarian Algorithm

The diagram illustrates a step in the Hungarian Algorithm. On the left, a 3x4 grid represents costs between objects Q1, Q2, Q3 and targets T1, T2, and an empty set \emptyset . An arrow points from this grid to a second, identical grid on the right. In the second grid, the values 1.92 and 4.5 have been circled in red, indicating they are being considered as potential minimum cost pairs.

	T1	T2	\emptyset
Q1	1.3	1.92	0
Q2	4.5	6.7	0
Q3	9.19	7.41	0

no object

→

	T1	T2	\emptyset
Q1	1.3	1.92	0
Q2	4.5	6.7	0
Q3	9.19	7.41	0

Figure: Find the pair minimize cost

Hungarian Algorithm

Preprocess matrix

	T1	T2	\emptyset
Q1	1	2	0
Q2	4	6	0
Q3	9	7	0

subtract row
minima

subtract col
minima

	T1	T2	\emptyset
Q1	0	0	0
Q2	3	4	0
Q3	8	5	0

Hungarian Algorithm

Repeat the followings:

- A Cover all the zeros with a minimum number of lines.
- B Decide whether to stop.
- C Create additional zeros.

Hungarian Algorithm

- Cover all 0 with minimum of lines
- If n lines required then stop else C

	T1	T2	\emptyset
Q1	0	0	0
Q2	3	4	0
Q3	8	5	0



	T1	T2	\emptyset
Q1	0	0	0
Q2	3	4	0
Q3	8	5	0

Hungarian Algorithm

- ① Find smallest element k not cover by lines
- ② Subtract k from all uncovered elements
- ③ add k to all the elements that are covered twice.

The diagram illustrates the Hungarian Algorithm through two stages of a 3x4 grid:

Initial State (Left):

	T1	T2	\emptyset
Q1	0	0	0 +3
Q2	3 3	4 -3	0
Q3	8 -3	5 -3	0

After One Iteration (Right):

	T1	T2	\emptyset
Q1	0	0	3
Q2	0	1	0
Q3	5	2	0

Hungarian Algorithm

	T1	T2	\emptyset
Q1	0	0	3
Q2	0	1	0
Q3	5	2	0

n lines required
stop

	T1	T2	\emptyset
Q1	0	0	3
Q2	0	1	0
Q3	5	2	0

Hungarian Algorithm

	T1	T2	\emptyset
Q1	0	0	3
Q2	0	1	0
Q3	5	2	0

output matrix

Q1 \longrightarrow T2
Q2 \longrightarrow T1
Q3 \longrightarrow \emptyset
minimum cost is 4+2

	T1	T2	\emptyset
Q1	1	2	0
Q2	4	6	0
Q3	9	7	0

original matrix

0	29.9	30.0	30.0	30.1	30.0	30.0	29.9	29.9	30.0	30.0	29.9	30.1	30.0
1	30.0	30.1	30.1	30.0	30.0	30.0	29.9	30.0	30.0	29.8	30.1	30.0	
2	29.9	30.0	30.1	29.9	29.9	30.0	29.9	30.0	29.9	29.8	29.9	30.1	30.0
3	-30.0	29.9	30.1	30.0	29.9	30.0	2.0	2.1	1.9	30.2	30.0	29.9	29.9
4	29.9	30.0	30.0	30.0	2.1	30.1	2.0	2.1	2.0	29.9	30.0	29.9	30.0
5	30.2	29.9	29.9	2.2	2.2	2.1	2.1	2.0	2.1	30.0	29.9	30.1	30.1
6	29.9	30.0	30.0	1.9	2.0	2.0	2.0	1.9	2.1	30.0	30.2	30.1	29.9
7	-30.0	29.9	30.0	2.0	2.0	2.2	2.2	1.9	1.9	29.8	29.8	30.0	30.0
8	29.9	29.8	30.2	30.0	29.9	2.1	2.1	29.9	29.9	30.3	29.9	30.1	30.1
9	-30.1	30.0	29.8	29.9	29.9	2.0	2.0	29.9	30.0	29.9	30.0	30.1	30.0
10	29.9	29.9	30.1	30.2	30.2	30.0	29.8	29.9	29.9	30.2	30.0	30.0	29.9
11	30.0	30.0	29.9	29.9	30.2	29.9	29.9	29.9	29.8	29.9	30.0	30.2	29.9
12	29.9	30.2	29.9	30.0	30.0	30.1	30.1	30.1	29.9	29.9	30.0	30.0	30.1

0	0.3	-1.8	-1.4	-0.2	1.3	1.0	-0.0	0.2	-0.2	-2.4	1.0	0.7	-0.6
1	-0.9	-0.7	-0.4	-0.3	-0.1	-0.6	-0.1	0.5	-0.3	-0.5	1.1	-0.8	-0.4
2	-1.0	0.3	-0.7	0.6	0.7	-28.0	-54.9	-84.7	-55.1	-26.8	0.4	-1.0	0.8
3	-0.2	0.6	-0.5	-27.7	-27.1	-84.2	140.2	282.7	140.6	58.0	-0.5	0.8	0.1
4	-0.1	-0.2	-28.2	-84.1	139.8	168.1	156.2	-0.5	84.0	-83.3	-0.1	0.9	-1.1
5	-0.5	0.6	-55.0	110.9	54.9	27.5	27.6	0.6	83.3	-83.7	0.7	-0.3	-2.0
6	-0.9	-0.0	-84.5	85.1	0.3	1.0	0.5	0.9	83.1	-84.4	-1.7	-0.7	-1.2
7	-0.7	0.4	-55.8	140.0	55.7	26.7	26.2	56.6	140.8	54.6	1.8	0.3	0.6
8	-2.0	1.1	-29.9	-56.8	138.8	55.3	55.8	-139.6	55.6	-31.1	1.0	-1.2	-2.4
9	-0.7	-0.0	1.4	1.1	-54.7	140.3	139.7	-55.2	-0.4	1.2	0.4	-0.6	-0.8
10	-0.1	0.8	-1.5	-1.4	-29.9	-55.8	-55.4	-28.3	0.4	-1.8	0.6	-0.1	0.3
11	-0.0	-0.3	1.1	1.6	-1.4	1.1	0.7	0.2	1.2	-0.1	-0.1	-1.5	1.5
12	-0.4	0.9	1.7	1.1	0.9	-1.3	-1.0	0.5	0.1	1.2	-0.1	-0.4	-0.9

Auxiliary Loss



Figure: GoogLeNet¹⁸

¹⁸Christian Szegedy et al. *Going Deeper with Convolutions*. 2014. arXiv: 1409.4842

Auxiliary Loss

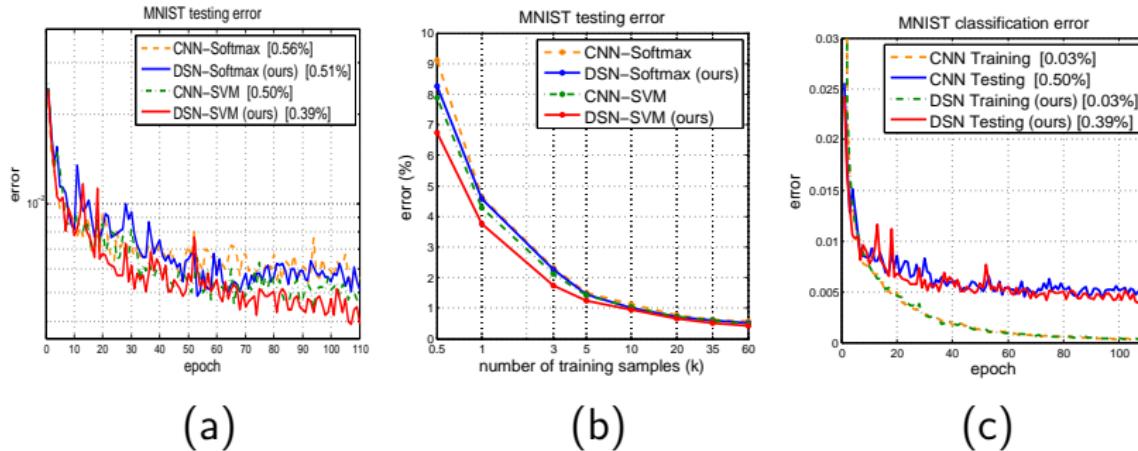


Figure: Classification error on MNIST test. (a) shows test error of competing methods; (b) shows test error w.r.t. the training sample size. (c) training and testing error comparison.²⁰

²⁰Chen-Yu Lee et al. Deeply-Supervised Nets. 2014. arXiv: 1409.5185.

Up/Down Sampling

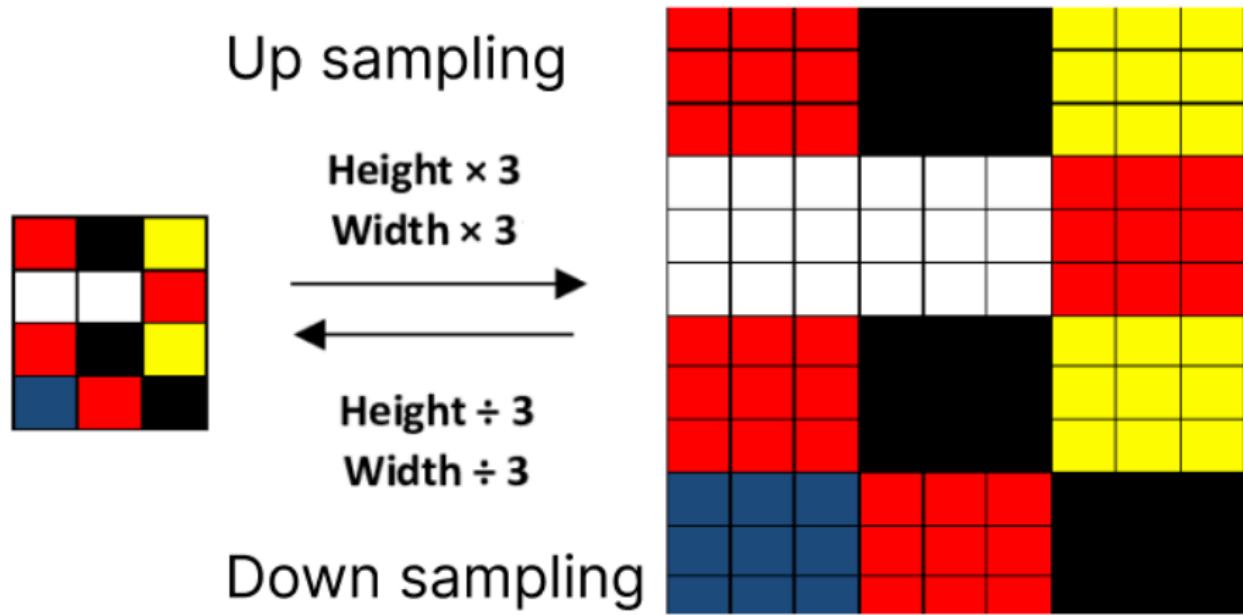


Figure: Image Upsampling and Downsampling²¹

²¹Ding Yang and Bai Jiabao. "An Optimization Method for Video Upsampling and Downsampling Using Interpolation-Dependent Image Downscaling". In: 2021 4th International Conference on Information Communication and Computer Science (ICMS). June 21, 2024 54 / 54