

---

# DAC-DETR: Divide the Attention Layers and Conquer

---

## 1 A Supplementary Material

### 2 A.1 Gathering effect on more decoder layers

3 Section 3.3 (Mechanism Analysis) in the main text shows that DAC-DETR improves the gathering  
4 effect of the cross-attention layer, *i.e.* more and better queries (Fig. 3 in the main text). We supplement  
5 results on more decoder layers (layer-2, layer-3 and layer-4) in Fig. A1.

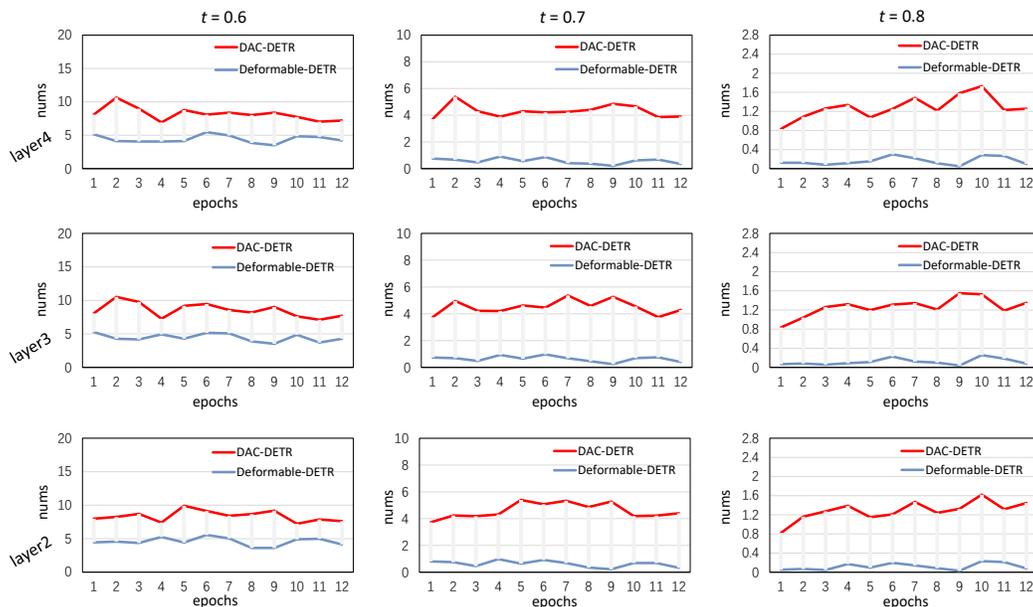


Figure A1: The averaged number of queries that each ground-truth object gathers on the validation set of MS-COCO.

6 The observation in Fig. A1 is consistent with Fig. 3 in the main text: DAC-DETR gathers more  
7 queries for each single object and improves the quality of the best queries. The two corresponding  
8 remarks, *i.e.*, DAC-DETR improves the quantity and quality of the gathered queries, hold across  
9 multiple decoder layers.

### 10 A.2 Comparison on Convergence Speed

11 We investigate the convergence speed of DAC-DETR on three baselines (Deformable-DETR [8],  
12 Deformable-DETR++ [3], and DINO [7]) in Fig. A2. The experiments are conducted on the COCO  
13 2017 [4] detection validation dataset. We adopt ResNet50 [2] backbone and run 12 epochs. It is  
14 observed that DAC-DETR consistently improves the convergence speed over all three baselines.  
15 For example, DAC-DETR outperforms the Deformable-DETR baseline by +8.3 AP and +3.4 AP at  
16 epoch-1 and epoch-12, respectively.

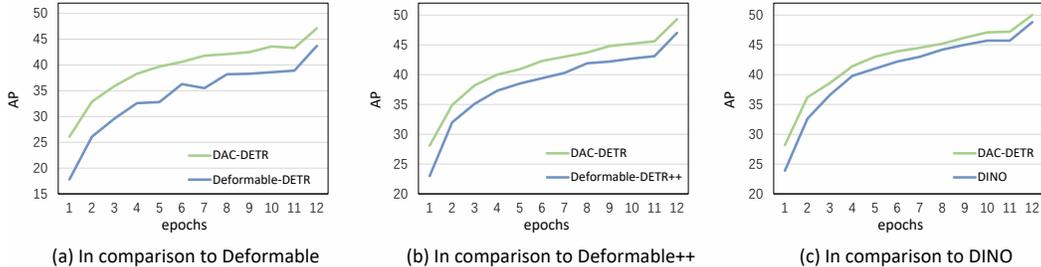


Figure A2: Comparison of convergence speed between DAC-DETR and three baselines.

### 17 A.3 Comparison of training time and inference FPS

18 We compare the average training time per epoch and the inference FPS between DAC-DETR, H-  
 19 DETR [3], and baseline method (Deformable-DETR [8]). For a fair comparison, all the methods  
 20 utilize 8 A100 GPUS for training and a single A100 GPU for inference.

| Method                 | Backbone | Training time (average) | Inference FPS | AP   |
|------------------------|----------|-------------------------|---------------|------|
| Basel (Deformable) [8] | R50      | 58 min                  | 17.8          | 43.7 |
| H-DETR [3]             | R50      | 70 min                  | 17.8          | 45.9 |
| DAC-DETR (ours)        | R50      | 64 min                  | 17.8          | 47.1 |

Table A1: Comparison of average training time on each epoch and inference FPS.

21 From Table A1, we draw two observations: 1) Compared to the baseline, DAC-DETR increases the  
 22 training time per epoch by a small margin (*i.e.*, +6 minutes) while maintaining the same inference  
 23 efficiency. The small increase on training time is because DAC-DETR additionally introduces an  
 24 auxiliary decoder (*i.e.*, C-Decoder) that processes all the queries in parallel. 2) Compared with  
 25 H-DETR (a recent method that employs auxiliary decoder branch), our DAC-DETR is faster to train  
 26 (-6 minutes per epoch). There are two reasons: first, DAC-DETR uses fewer queries than H-DETR.  
 27 Second, the auxiliary C-Decoder in DAC-DETR has fewer attention layers (*i.e.*, no self-attention  
 28 layers).

### 29 A.4 More hyper-parameter analysis

30 In our one-to-many label assignment (Eqn.4 in the main text), we compute the matching score  $m$   
 31 between each query and the object by adding their IoU score and the predicted label score on the  
 32 ground-truth class. To introduce more flexibility, we combine these two scores through a weighted  
 33 sum, which is formulated as:

$$m = (1 - \lambda) \cdot p_{(q)}(\hat{c}) + \lambda \cdot \text{IoU} \langle b_{(q)}, \hat{b} \rangle, \quad (1)$$

34 where  $\lambda$  is a newly-added hyper-parameter for weighting, and all the other variables are the same as  
 35 in the main text (*i.e.*,  $\hat{c}$  and  $\hat{b}$  are the class and bounding box of query  $q$ ,  $p_{(q)}(\hat{c})$  denotes the predicted  
 36 label score on class  $\hat{c}$ .  $b_{(q)}$  denotes the predicted box,  $\langle, \rangle$  denotes the IoU operation between  
 37 predicted box and ground truth  $\hat{b}$ ). We investigate the influence of this hyper-parameter  $\lambda$  in Table A2.

| $\lambda$ | 0.5  | 0.6  | 0.7  | 0.8  | 0.9  | 1.0  |
|-----------|------|------|------|------|------|------|
| AP        | 46.8 | 47.0 | 47.1 | 47.0 | 46.8 | 46.6 |

Table A2: Analysis on the weight  $\lambda$  in the one-to-many label assignment. We adopt Deformable-DETR as the baseline.

38 We observe that DAC-DETR is robust to this hyper-parameter within a large range, and in practice  
 39 use  $\lambda = 0.7$  for all the experiments.

40 **A.5 More Experiments**

41 We evaluate the performance of Align-DETR [1] with the Swin-L [6] backbone on COCO 2017  
 42 detection validation dataset, using the official publicly codes. ( Align-DETR does not report the  
 43 results with Swin-L backbone). The results in Table A3 further confirms the superiority of DAC-  
 44 DETR. After combining an IoU-related loss (Align loss), DAC-DETR surpasses Align-DETR by  
 45 +0.7 AP (12 epochs).

| Method                  | Backbone | epochs | AP   | AP <sub>50</sub> | AP <sub>75</sub> | AP <sub>S</sub> | AP <sub>M</sub> | AP <sub>L</sub> |
|-------------------------|----------|--------|------|------------------|------------------|-----------------|-----------------|-----------------|
| Basel (DINO) [7]        | Swin-L   | 12     | 56.8 | 75.6             | 62.0             | 40.0            | 60.5            | 73.2            |
| Align-DETR † [1]        | Swin-L   | 12     | 57.4 | 75.9             | 62.2             | 40.6            | 61.6            | 73.7            |
| Stable-DINO-4scale [5]  | Swin-L   | 12     | 57.7 | 75.7             | 63.4             | 39.8            | 62.0            | 74.7            |
| DAC-DETR + Align (ours) | Swin-L   | 12     | 58.1 | 76.5             | 63.3             | 40.9            | 62.4            | 75.0            |

Table A3: Evaluation on COCO val2017 with Swin-Transformer Large backbone. †: We evaluate Align-DETR using the official publicly codes.

46 **A.6 Visualization of Object Detection**

47 We visualize some detection results with predicted bounding boxes and label scores in Fig. A3  
 48 and Fig. A4. As shown in Fig. A3, DAC-DETR detects the object "zebra" with limited semantic  
 49 information, whereas Deformable-DETR fails to do so. Compared to Deformable-DETR++, DAC-  
 50 DETR provides more accurate label and box predictions for the object "cat", as shown in Fig. A4.

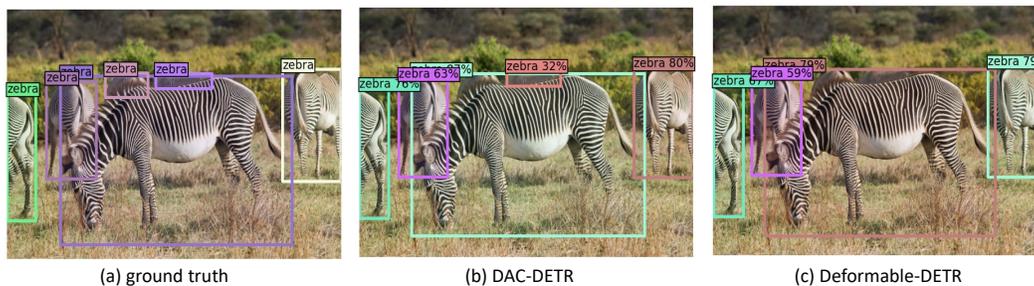


Figure A3: Visualization of the detection results of DAC-DETR and Deformable-DETR.



Figure A4: Visualization of the detection results of DAC-DETR and Deformable-DETR++.

51 **References**

- 52 [1] Zhi Cai, Songtao Liu, Guodong Wang, Zheng Ge, Xiangyu Zhang, and Di Huang. Align-detr:  
53 Improving detr with simple iou-aware bce loss, 2023.
- 54 [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
55 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
56 pages 770–778, 2016.
- 57 [3] Ding Jia, Yuhui Yuan, Haodi He, Xiaopei Wu, Haojun Yu, Weihong Lin, Lei Sun, Chao Zhang,  
58 and Han Hu. Detr with hybrid matching. *arXiv preprint arXiv:2207.13080*, 2022.
- 59 [4] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
60 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer  
61 Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014,  
62 Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- 63 [5] Shilong Liu, Tianhe Ren, Jiayu Chen, Zhaoyang Zeng, Hao Zhang, Feng Li, Hongyang Li, Jun  
64 Huang, Hang Su, Jun Zhu, and Lei Zhang. Detection transformer with stable matching, 2023.
- 65 [6] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
66 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings  
67 of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- 68 [7] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M. Ni, and Heung-Yeung  
69 Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection, 2022.
- 70 [8] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr:  
71 Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*,  
72 2020.