

Supplementary Material for AdaMixer: A Fast-Converging Accurate Query-based Object detector

Ziteng Gao¹

Limin Wang¹✉

Bing Han²

Sheng Guo²

¹State Key Laboratory for Novel Software Technology, Nanjing University, China

²MYbank, Ant Group, China

A. 1. Detection Performance on COCO test set

backbone	AP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l
ResNet-50	47.2	66.3	51.4	28.3	49.3	60.6
ResNet-101	48.1	67.3	52.4	28.3	50.6	62.1
ResNeXt-101-DCN	49.3	68.8	53.7	30.0	51.6	64.0
Swin-S	51.3	71.3	55.9	31.4	53.9	66.3

Table 1. AdaMixer performance on COCO test-dev set with **longer training scheme** (36 epochs training and 300 queries) and single model single scale testing.

The performance of AdaMixer on COCO minival set is reported in Table 5 in the paper. Here we report the performance of these AdaMixer models on COCO test-dev set in Table 1, where labels are not publicly available and evaluation is done on the online server. Note that models here are exactly the same in Table 5 in the paper.

A. 2. More Studies on Adaptive 3D Sampling

We here also validate the effectiveness of our proposed 3D feature space and adaptive 3D sampling in Table 2. Using single feature map leads to inferior results for our AdaMixer. Note that there is no feature pyramid networks (FPN) [5] used. The RoIAlign operator, which extracts features in a single level map according to the bounding box, also performs inferior to our adaptive 3D sampling approach. RoIAlign feature sampling locations are highly stricted inner the bounding box and the sampling level (the z coordinate in our work) is also discretized, whereas locations and levels are adaptive on queries and not stricted in our adaptive 3D sampling. This means that RoIAlign lacks multi-level feature modeling and necessitates multi-scale feature interaction necks.

Relations to other work. Our adaptive 3D sampling is in accord with deformable convolutions [2, 8, 9] to explicitly utilize offsets to sample features spatially to model deformation of objects. But our adaptive 3D sampling extends the adaptive modeling to the scale dimension, namely z -

sampling	AP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l
only C ₃ feature	26.2	42.0	27.3	15.8	28.7	34.1
only C ₄ feature	38.3	57.2	41.0	20.0	41.9	54.1
only C ₅ feature	37.8	58.3	39.5	18.0	41.2	51.7
RoIAlign	37.2	58.5	39.0	19.0	39.3	55.6
A2DS	41.3	61.0	44.4	23.3	43.8	57.8
A3DS	42.7	61.5	45.9	24.7	45.4	59.2

Table 2. **Different feature sampling** performance on MS COCO minival set with $1\times$ training scheme. The “A3DS” is abbreviation for Adaptive 3D Sampling used in AdaMixer. The “A2DS” is the 2D variant of adaptive 3D sampling, where it disables learning adaptive z -axis offsets, namely, $\Delta z_{(\cdot)} = 0$ through training and inference. Note that A2DS is still with multiple feature maps and z of a query still decides how feature maps are aggregated. The RoIAlign is performed to extract features from feature maps individually according to the bounding box indicated by a query. The row “only C_(·) feature” variants use a single feature map to perform adaptive 2D sampling. No FPNs used throughout this table.

axis, to cope with the variation of potential objects. Also, our adaptive 3D sampling is performed by queries in the sparse manner, whereas the deformable convolution or attention are usually performed in the dense manner. Moreover, our method naturally generalizes the scale-equalizing operators [4, 7] on discretized feature maps with different scales to a continuous and interpolable one, enjoying more flexible modeling for scale variations.

A. 3. Model Details

Our AdaMixer is implemented on mmdetection framework [1]. In current implementation, the model is fully based on PyTorch primitives without customized CUDA codes.

Initializations. We initialize all backbones from pre-trained model on ImageNet 1K classification [3], including Swin-S [6]. Please refer to our codebase for more detailed initializations in AdaMixer.

detector	backbone	AP _{val}	AP _{test}	#Param (M)	flops (G)	FPS	training hours
Faster R-CNN	R-50	41.8	-	42	207	17	~16
DETR	R-50-DC5	43.3	-	41	187	11	~204
Deformable DETR++	R-50	46.2	46.9	40	173	12	~106
Sparse R-CNN	R-50	45.0	-	110	174	16	~30
AdaMixer	R-50	47.0	47.2	139	132	15	29
AdaMixer	R-101	48.0	48.1	158	208	12	37
AdaMixer	X-101-DCN	49.5	49.3	160	214	8	65
AdaMixer	Swin-S	51.3	51.3	164	234	10	51

Table 3. Training and inference details about different models. AdaMixer models here are ones in Table 5 of the paper with 300 queries. AP_{val} and AP_{test} are the average precision on COCO minival and test-dev set, respectively. FPS is calculated on a single Nvidia V100 card. Training hours for other detectors are estimated on 8 V100 cards when obtaining the speed of a few training iterations.

Training and inference speed. We evaluate AdaMixer training and inference speed and compare it with other detectors in Table 3. Our AdaMixer shows advantage not only on theoretical flops but on the actual training time and inference FPS. Note that the current implementation, the adaptive 3D sampling procedure is based on PyTorch primitive `grid_sample`, which can be optimized in the future.

A. 4. Visualizations

We also show visualizations of sampling points and final detections of our AdaMixer detector in Figure 1 and 2. These visualizations are done on samples in COCO minival set. The first row shows the input image and detection results and the other rows show sampling points for each stage. Figure 1 and 2 shows two different query results¹. The sampling z-axis coordinate, \tilde{z} , is visualized as the point size and a bigger point corresponds with the larger \tilde{z} . Sampling points of different groups are colored with different colors. We can see that AdaMixer actually sees out of the box and different groups have preference for different semantics. Moreover, when comparing Figure 1 and 2, we can find that sampling point patterns also vary across different queries, indicating the enhanced sampling adaptability and flexibility in AdaMixer. Counterintuitively, sampling points of a query do not gather together in a monotonous manner across stages. Instead, they often stretch wider after focusing on small regions and then focus again (stage 2→3→4→5→6). We suspect that this behavior is beneficial to performing bounding box estimation more accurately. Please see `visualizations` folder in our codebase or supplementary file for more visualizations.

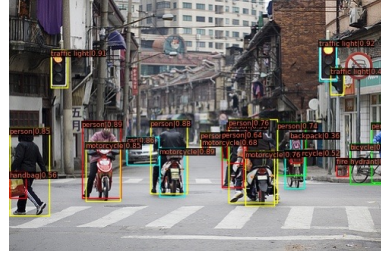
¹In a figure, the query index is consistent through stages.

References

- [1] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Mmdetection: Open mmlab detection toolbox and benchmark. In *arXiv*, 2019. 1
- [2] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, 2017. 1
- [3] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 1
- [4] Yanghao Li, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Scale-aware trident networks for object detection. In *CVPR*, 2019. 1
- [5] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 1
- [6] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 1
- [7] Xinjiang Wang, Shilong Zhang, Zhuoran Yu, Litong Feng, and Wayne Zhang. Scale-equalizing pyramid convolution for object detection. In *CVPR*, 2020. 1
- [8] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *CVPR*, 2019. 1
- [9] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: deformable transformers for end-to-end object detection. In *ICLR*, 2020. 1



input image



detection results



stage 1:
sampling points for person #1



stage 2:
sampling points for person #1



stage 3:
sampling points for person #1



stage 4:
sampling points for person #1



stage 5:
sampling points for person #1

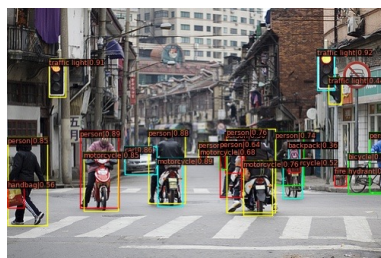


stage 6:
sampling points for person #1

Figure 1. Visualizations of detection results and sampling points of AdaMixer with ResNet-50 as the backbone.



input image



detection results



stage 1:
sampling points for person #2



stage 2:
sampling points for person #2



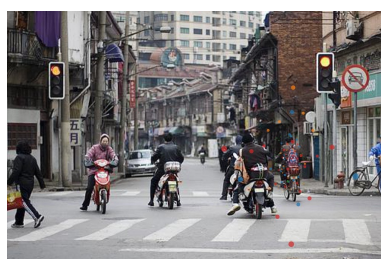
stage 3:
sampling points for person #2



stage 4:
sampling points for person #2



stage 5:
sampling points for person #2



stage 6:
sampling points for person #2

Figure 2. Visualizations of detection results and sampling points of AdaMixer with ResNet-50 as the backbone (cont'd).