# DuET: <u>Du</u>al Incremental Object Detection via <u>E</u>xemplar-Free <u>T</u>ask Arithmetic

## Supplementary Material

## A. Detailed Evaluation Protocol and Metrics

This section provides a detailed discussion of the proposed evaluation protocol and metrics for the DuIOD task. Table S1 outlines the training sequence that is followed for four different DuIOD experiments, along with the detailed evaluation protocol that is used to comprehensively evaluate the performance of different object detectors on the respective DuIOD setting. Unlike existing metrics [1, 21, 24] that focus only on catastrophic forgetting, we used the **Retention-Adaptability Index (RAI)**, which balances both knowledge retention and generalisation to unseen categories across evolving domains. We define RAI as the mean of the Average Retention Index (Avg RI) and Average Generalisation Index (Avg GI), which are discussed in the sections below.

$$RAI = \frac{\text{Avg RI} + \text{Avg GI}}{2} \quad (1)$$

### A.1. Average Retention Index

For each domain $\mathcal{D}_i$ corresponding to task $\mathcal{T}_i$ where $i \in \{1, \dots, T-1\}$, we define the Retention Index $RI_{\mathcal{D}_i}$ as:

$$RI_{\mathcal{D}_i} = \frac{\text{mAP}_{\text{old}}^{\mathcal{T}_T}(\mathcal{D}_i[\mathcal{C}_i])}{\text{mAP}_{\text{new}}^{\mathcal{T}_i}(\mathcal{D}_i[\mathcal{C}_i])} \quad (2)$$

Here, $\text{mAP}_{\text{old}}^{\mathcal{T}_T}(\mathcal{D}_i[\mathcal{C}_i])$ denotes the mean Average Precision (mAP) at IoU threshold = 0.5 of the object detector at the final task $\mathcal{T}_T$ on the classes $\mathcal{C}_i$ which were learned from domain $\mathcal{D}_i$, and $\text{mAP}_{\text{new}}^{\mathcal{T}_i}(\mathcal{D}_i[\mathcal{C}_i])$ is the mAP when classes $\mathcal{C}_i$ from domain $\mathcal{D}_i$ were first encountered and learned, at task $\mathcal{T}_i$. The Avg RI is then calculated as:

$$\text{Avg RI} = \frac{1}{T-1} \sum_{i=1}^{T-1} RI_{\mathcal{D}_i}. \quad (3)$$

To illustrate this, consider the multi-phase experiment (Table S1) with training sequence: *Night Sunny [1:2] → Daytime Sunny [3:4] → Daytime Foggy [5:7].* The Avg RI is computed as the mean of the Retention Index values for Night Sunny (NS) and Daytime Sunny (DS) domains at the final task $\mathcal{T}_3$ as follows:

$$RI_{\text{NS}} = \frac{\text{mAP}_{\text{old}}^{\mathcal{T}_3}(\text{NS}[1:2])}{\text{mAP}_{\text{new}}^{\mathcal{T}_1}(\text{NS}[1:2])} \quad RI_{\text{DS}} = \frac{\text{mAP}_{\text{old}}^{\mathcal{T}_3}(\text{DS}[3:4])}{\text{mAP}_{\text{new}}^{\mathcal{T}_2}(\text{DS}[3:4])} \quad (4)$$

$$\text{Avg RI} = \frac{RI_{\text{NS}} + RI_{\text{DS}}}{2} \quad (5)$$

Hence, in this case, a higher Avg RI indicates how effectively the object detector has retained past knowledge from old Night Sunny and Daytime Sunny domains in the final task $\mathcal{T}_3$. Conversely, a lower value indicates significant catastrophic forgetting.

### A.2. Average Generalization Index

The Generalisation Index ($GI_{\mathcal{D}_i, \mathcal{T}_j}$) quantifies how well the model detects unseen classes from domain $\mathcal{D}_i$ at task $\mathcal{T}_j$. These classes were not part of the training set for task $\mathcal{T}_j$, meaning the model is required to generalise beyond its explicitly trained classes (see Table S1). For a given domain $\mathcal{D}_i$ at task $\mathcal{T}_j$, the Generalization Index is computed as:

$$GI_{\mathcal{D}_i, \mathcal{T}_j} = \frac{\text{mAP}_{\text{unseen}}^{\mathcal{T}_j}(\mathcal{D}_i[\mathcal{C}_{\text{unseen}}])}{\text{mAP}_{\text{ref}}(\mathcal{D}_i[\mathcal{C}_{\text{unseen}}])} \quad (6)$$

Here, $\text{mAP}_{\text{unseen}}^{\mathcal{T}_j}(\mathcal{D}_i[\mathcal{C}_{\text{unseen}}])$ is the mAP of the model at task $\mathcal{T}_j$ on the unseen classes $\mathcal{C}_{\text{unseen}}$ from domain $\mathcal{D}_i$, and $\text{mAP}_{\text{ref}}(\mathcal{D}_i[\mathcal{C}_{\text{unseen}}])$ is the reference mAP obtained by training the object detector solely on these unseen classes on domain $\mathcal{D}_i$. The Average Generalisation Index (Avg GI) over all relevant domain-task pairs is then computed as:

$$\text{Avg GI} = \frac{1}{\mathcal{N}} \sum_{(\mathcal{D}_i, \mathcal{T}_j)} GI_{\mathcal{D}_i, \mathcal{T}_j} \quad (7)$$

where $\mathcal{N}$ is the total number of unseen-class domain-task pairs considered in the evaluation.

Continuing with the same example, the Avg GI is computed as the mean of the Generalization Index values for unseen classes across the Night Sunny (NS), Daytime Sunny (DS), and Daytime Foggy (DF) domains for a total of five domain-task pairs- two from task $\mathcal{T}_2$ and three from $\mathcal{T}_3$:

$$GI_{\text{NS}, \mathcal{T}_2} = \frac{\text{mAP}_{\text{unseen}}^{\mathcal{T}_2}(\text{NS}[3:4])}{\text{mAP}_{\text{ref}}(\text{NS}[3:4])} \quad GI_{\text{DS}, \mathcal{T}_2} = \frac{\text{mAP}_{\text{unseen}}^{\mathcal{T}_2}(\text{DS}[1:2])}{\text{mAP}_{\text{ref}}(\text{DS}[1:2])} \quad (8)$$

$$GI_{\text{NS}, \mathcal{T}_3} = \frac{\text{mAP}_{\text{unseen}}^{\mathcal{T}_3}(\text{NS}[3:4])}{\text{mAP}_{\text{ref}}(\text{NS}[3:4])} \quad GI_{\text{DS}, \mathcal{T}_3} = \frac{\text{mAP}_{\text{unseen}}^{\mathcal{T}_3}(\text{DS}[1:2])}{\text{mAP}_{\text{ref}}(\text{DS}[1:2])}$$

$$GI_{\text{DF}, \mathcal{T}_3} = \frac{\text{mAP}_{\text{unseen}}^{\mathcal{T}_3}(\text{DF}[1:4])}{\text{mAP}_{\text{ref}}(\text{DF}[1:4])} \quad (9)$$

$$\text{Avg GI} = \frac{GI_{\text{NS}, \mathcal{T}_2} + GI_{\text{DS}, \mathcal{T}_2} + GI_{\text{NS}, \mathcal{T}_3} + GI_{\text{DS}, \mathcal{T}_3} + GI_{\text{DF}, \mathcal{T}_3}}{5} \quad (10)$$

Hence, in this case, a higher Avg GI indicates better zero-shot generalisation to unseen categories: NS [3:4], DS [1:2], and DF [1:4] across incremental training. Conversely, a lower value suggests that the model is overfitting to seen classes and fails to generalise.

Table S1. Training Sequence & Evaluation Protocol for different DuIOD experiments.

| DuIOD Experiment | Training Sequence | | Evaluation Protocol | | |
|---|---|---|---|---|---|
| | Task | Class IDs | New Classes | Old Classes | Unseen Classes |
| **Pascal Series Datasets** | | | | | |
| **Two Phase** VOC [1:10] $\rightarrow$ Clipart [11:20] | $\mathcal{T}_1$ | 1-10 from VOC | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_1}(\mathrm{VOC}[1:10])$ | — | — |
| | $\mathcal{T}_2$ | 11-20 from Clipart | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_2}(\mathrm{Clipart}[11:20])$ | $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_2}(\mathrm{VOC}[1:10])$ | $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{VOC}[11:20])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{Clipart}[1:10])$ |
| **Multi Phase** Watercolor [1:3] $\rightarrow$ Comic [4:6] $\rightarrow$ Clipart [7:13] $\rightarrow$ VOC [14:20] | $\mathcal{T}_1$ | 1-3 from Watercolor | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_1}(\mathrm{Watercolor}[1:3])$ | — | — |
| | $\mathcal{T}_2$ | 4-6 from Comic | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_2}(\mathrm{Comic}[4:6])$ | $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_2}(\mathrm{Watercolor}[1:3])$ | $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{Watercolor}[4:6])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{Comic}[1:3])$ |
| | $\mathcal{T}_3$ | 7-13 from Clipart | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_3}(\mathrm{Clipart}[7:13])$ | $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_3}(\mathrm{Watercolor}[1:3])$ $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_3}(\mathrm{Comic}[4:6])$ | $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_3}(\mathrm{Watercolor}[4:6])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_3}(\mathrm{Comic}[1:3])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_3}(\mathrm{Clipart}[1:6])$ |
| | $\mathcal{T}_4$ | 14-20 from VOC | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_4}(\mathrm{VOC}[14:20])$ | $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_4}(\mathrm{Watercolor}[1:3])$ $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_4}(\mathrm{Comic}[4:6])$ $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_4}(\mathrm{Clipart}[7:13])$ | $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_4}(\mathrm{Watercolor}[4:6])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_4}(\mathrm{Comic}[1:3])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_4}(\mathrm{Clipart}[1:6])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_4}(\mathrm{VOC}[1:13])$ |
| **Diverse Weather Series Datasets** | | | | | |
| **Two Phase** Daytime Sunny [1:4] $\rightarrow$ Night Sunny [5:7] | $\mathcal{T}_1$ | 1-4 from Daytime Sunny | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_1}(\mathrm{Daytime\ Sunny}[1:4])$ | — | — |
| | $\mathcal{T}_2$ | 5-7 from Night Sunny | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_2}(\mathrm{Night\ Sunny}[5:7])$ | $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_2}(\mathrm{Daytime\ Sunny}[1:4])$ | $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{Daytime\ Sunny}[5:7])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{Night\ Sunny}[1:4])$ |
| **Multi Phase** Night Sunny [1:2] $\rightarrow$ Daytime Sunny [3:4] $\rightarrow$ Daytime Foggy [5:7] | $\mathcal{T}_1$ | 1-2 from Night Sunny | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_1}(\mathrm{Night\ Sunny}[1:2])$ | — | — |
| | $\mathcal{T}_2$ | 3-4 from Daytime Sunny | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_2}(\mathrm{Daytime\ Sunny}[3:4])$ | $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_2}(\mathrm{Night\ Sunny}[1:2])$ | $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{Night\ Sunny}[3:4])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_2}(\mathrm{Daytime\ Sunny}[1:2])$ |
| | $\mathcal{T}_3$ | 5-7 from Daytime Foggy | $\mathrm{mAP}_{\mathrm{new}}^{\mathcal{T}_3}(\mathrm{Daytime\ Foggy}[5:7])$ | $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_3}(\mathrm{Night\ Sunny}[1:2])$ $\mathrm{mAP}_{\mathrm{old}}^{\mathcal{T}_3}(\mathrm{Daytime\ Sunny}[3:4])$ | $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_3}(\mathrm{Night\ Sunny}[3:4])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_3}(\mathrm{Daytime\ Sunny}[1:2])$ $\mathrm{mAP}_{\mathrm{unseen}}^{\mathcal{T}_3}(\mathrm{Daytime\ Foggy}[1:4])$ |

## B. Loss Function Formulation

To ensure effective incremental learning in case of DuIOD, we employ a combination of standard detector loss $\mathcal{L}_{\mathrm{Detector}}$, a modified distillation loss $\mathcal{L}_{\mathrm{Distill}}^*$ (discussed below), and the Directional Consistency Loss $\mathcal{L}_{\mathrm{DC}}$ (discussed in main paper `Section 3.5`). This section details the formulation of $\mathcal{L}_{\mathrm{total}}$ using these loss components.

Knowledge distillation plays a crucial role in mitigating catastrophic forgetting during incremental learning. In our approach, we extend the standard distillation loss ($\mathcal{L}_{\mathrm{Distill}}$) for incremental learning [13] by incorporating a dynamic thresholding mechanism that filters low-confidence classification outputs and high-variance bounding box predictions from the old (previous task) model.

Let $\mathcal{M}_{\theta_{t-1}}$ represent the previous task model, and $\mathcal{M}_{\theta_t}$ be the current model being trained on $\mathcal{T}_t$. Given input data for the current task $\mathcal{X}_t$, the classification outputs and predicted bounding boxes from both models will be:

$$\mathbf{z}_{\mathrm{curr}} = \mathcal{M}_{\theta_t}(x), \quad \mathbf{z}_{\mathrm{old}} = \mathcal{M}_{\theta_{t-1}}(x) \qquad (11)$$

where $\mathbf{z} = (\mathbf{c}, \mathbf{b})$, with $\mathbf{c}$ being classification logits and $\mathbf{b}$ the predicted bounding box coordinates.

The classification distillation loss is computed as:

$$\mathcal{L}_{\mathrm{Distill}_{\mathrm{cls}}}^* = \frac{1}{|\mathcal{M}_{\mathrm{cls}}^*|} \sum_{i \in \mathcal{M}^*} \left\| \mathbf{c}_{\mathrm{curr}}^{(i)} - \mathbf{c}_{\mathrm{old}}^{(i)} \right\|^2 \qquad (12)$$

where $\mathcal{M}_{\mathrm{cls}}^*$ is the dynamically selected mask that excludes predictions with low confidence scores in $\mathbf{c}_{\mathrm{old}}$:

$$\mathcal{M}_{\mathrm{cls}}^* = \{ i \mid \max(\mathbf{c}_{\mathrm{old}}^{(i)}) \geq \tau_{\mathrm{cls}} \} \qquad (13)$$

where $\tau_{\mathrm{cls}}$ is an adaptive threshold computed as the 75th percentile of $\max(\mathbf{c}_{\mathrm{old}})$ values.

Similarly, the bounding box regression distillation loss is computed by computing the KL divergence between the softmax of bounding box outputs from the current and old models:

$$\mathcal{L}_{\mathrm{Distill}_{\mathrm{bbox}}}^* = \frac{1}{|\mathcal{M}_{\mathrm{bbox}}^*|} \sum_{j \in \mathcal{M}_{\mathrm{bbox}}^*} \mathcal{D}_{\mathrm{KL}}\left( \mathrm{Softmax}(\mathbf{b}_{\mathrm{curr}}^{(j)}) \middle\| \mathrm{Softmax}(\mathbf{b}_{\mathrm{old}}^{(j)}) \right) \quad (14)$$

where $\mathcal{M}_{\mathrm{bbox}}^*$ filters out bounding boxes with high variance in $\mathbf{b}_{\mathrm{old}}$:

$$\mathcal{M}_{\mathrm{bbox}}^* = \{ j \mid \mathrm{Var}(\mathbf{b}_{\mathrm{old}}^{(j)}) \leq \tau_{\mathrm{bbox}} \} \qquad (15)$$

where $\tau_{\mathrm{bbox}}$ is an adaptive threshold computed as the 75th percentile of bounding box variance values.

The final modified distillation loss becomes:

$$\mathcal{L}_{\mathrm{Distill}}^* = \mathcal{L}_{\mathrm{Distill}_{\mathrm{cls}}}^* + \mathcal{L}_{\mathrm{Distill}_{\mathrm{bbox}}}^* \qquad (16)$$

$\mathcal{L}_{\mathrm{Detector}}$ depends on the object detector used. In our framework, we augment the detector losses of YOLO11 [10] and RT-DETR [18] object detectors with $\mathcal{L}_{\mathrm{Distill}}^*$ and $\mathcal{L}_{\mathrm{DC}}$ on the Ultralytics [10] pipeline. In the case of YOLO11, the detection loss consists of classification loss, bounding box regression loss, and Distribution Focal Loss

[10], while in the case of RT-DETR, the detection loss follows a Hungarian matching strategy, consisting of classification, bounding box, and Generalized IoU (GIoU) losses [18].

Hence, the total loss for incremental tasks ($t \geq 2$) is computed as:

$$\mathcal{L}_{\text{Total}} = \mathcal{L}_{\text{Detector}} + \lambda_{\text{Distill}}\mathcal{L}^*_{\text{Distill}} + \lambda_{\text{DC}}\mathcal{L}_{\text{DC}} \qquad (17)$$

where $\lambda_{\text{Distill}}$ and $\lambda_{\text{DC}}$ are scaling coefficients that control the impact of distillation and directional consistency losses, respectively.

## C. Extended Ablation Studies

### C.1. Impact of Loss Components

Table S2. Performance comparison of different model-merging algorithms on VOC [1:10] → Clipart [11:20] depicting the impact of $\mathcal{L}_{\text{DC}}$, with YOLO11n [10] as the base detector. Among columns, best in **bold**, second best _underlined_.

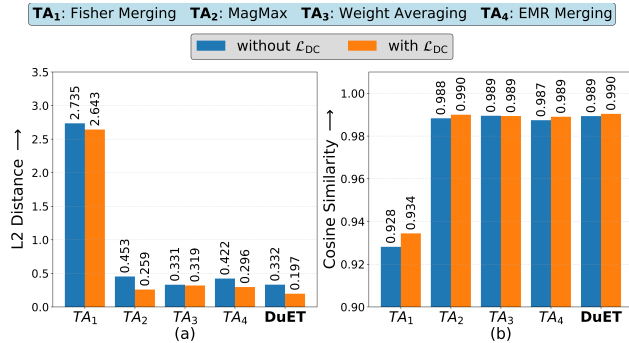| Model-merging algorithm | $\mathcal{L}_{\text{DC}}$ | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|
| Fisher-Merging [20] | ✗ | 20.27 | 17.15 | 18.71 |
| Fisher-Merging [20] | ✓ | 21.64 | 24 | 22.82 (+ 4.11) |
| MagMax [19] | ✗ | 65.05 | 28.09 | 46.57 |
| MagMax [19] | ✓ | 66.79 | 28.28 | 47.54 (+ 0.97) |
| Weight-Averaging [8] | ✗ | 66.42 | 31.42 | 48.92 |
| Weight-Averaging [8] | ✓ | 76.12 | 37.53 | 56.83 (+ 7.91) |
| EMR-Merging [7] | ✗ | 67.66 | 34.4 | 51.03 |
| EMR-Merging [7] | ✓ | 68.03 | 36.46 | 52.25 (+ 1.22) |
| **DuET (Ours)** | ✗ | _87.06_ | _37.75_ | _62.41_ |
| **DuET (Ours)** | ✓ | **87.44** | **44.54** | **65.99** (+ 3.58) |



Figure S1. **Impact of $\mathcal{L}_{DC}$ in (a) reducing L2 Distance and (b) improving cosine similarity.** These results are obtained on the VOC [1:10] → Clipart [11:20] experiment using YOLO11n [10] as the base detector with Incremental Head and Sequential Finetuning.

**Impact of $\mathcal{L}_{DC}$.** Continuing the ablations from the main paper (Section 6), in this section, we further investigate the role of $\mathcal{L}_{DC}$. In Table S2, we compare the performance of different model-merging algorithms, with and without $\mathcal{L}_{DC}$. The results show that $\mathcal{L}_{DC}$ consistently improves the RAI across all methods, with an average RAI improvement of **+3.56%** among all merging methods, with DuET achieving the best performance. Moreover, the bar charts in Figure S1 compare the L2 distance and cosine similarity between the merged model weights with both old and current model weights across different model-merging algorithms. The results show that $\mathcal{L}_{DC}$ significantly reduces L2 distance by **43.46%** averaged across all methods, with DuET achieving the lowest values. Lower L2 distance suggests that after incorporating $\mathcal{L}_{DC}$, the merged model lies closer to the original models, ensuring effective knowledge integration from both. Similarly, incorporation of $\mathcal{L}_{DC}$ consistently improves cosine similarity across all methods by **0.23%** average, with DuET achieving the highest values. Higher cosine similarity suggests that $\mathcal{L}_{DC}$ helps the merged model better align with the original models.

Table S3. Ablation studies of different loss components augmented with detector loss ($\mathcal{L}_{\text{Detector}}$).

| Loss Component | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|
| $\mathcal{L}_{Detector} + \mathcal{L}_{Distill}$ | 72.64 | 33.74 | 53.19 |
| $\mathcal{L}_{Detector} + \mathcal{L}^*_{Distill}$ | 87.06 | 37.75 | 62.41 |
| $\mathcal{L}_{Detector} + \mathcal{L}^*_{Distill} + \mathcal{L}_{DC}$ | **87.44** | **44.54** | **65.99** |

**Impact of $\mathcal{L}^*_{Distill}$.** Table S3 presents the ablation studies of different loss components augmented with detector loss ($\mathcal{L}_{\text{Detector}}$). We observe that the inclusion of $\mathcal{L}^*_{\text{Distill}}$ instead of $\mathcal{L}_{\text{Distill}}$ significantly improves all metrics, with a **+14.42%** increase in Avg RI, **+4.01%** increase in Avg GI, and **+9.22%** increase in RAI. The addition of $\mathcal{L}_{\text{DC}}$ brings in additional improvements, leading to the best performance across all metrics.

### C.2. Sensitivity Analysis for key hyper-parameters

Figure S2 shows the sensitivity analysis for key hyperparameters used in the DuET approach. Base scaling coefficient $\alpha_{base}$ (Figure S2a) effectively controls the contributions from the old (prior task) model and current model; hence, a value of **0.5** ensures a balanced trade-off between past knowledge retention and new adaptation, while extreme values ($\alpha_{base} < 0.3$ or $\alpha_{base} > 0.7$) significantly degrade RAI. The limiting factor $\gamma$ (Figure S2b) impacts task-merging, with $\gamma = 0.1$ giving optimal results for both _Pascal Series_ and _Diverse Weather Series_ datasets. The scaling coefficients $\lambda_{Distill}$ and $\lambda_{DC}$ (Figures S2c and S2d) control the impact of Distillation and Directional Consistency losses, respectively. We observe that, for both of them, a value of **0.01** gives the best results and effectively helps in mitigating catastrophic forgetting by improving retention while preventing sign conflicts; deviations from these values lead to reduced adaptability and degraded performance
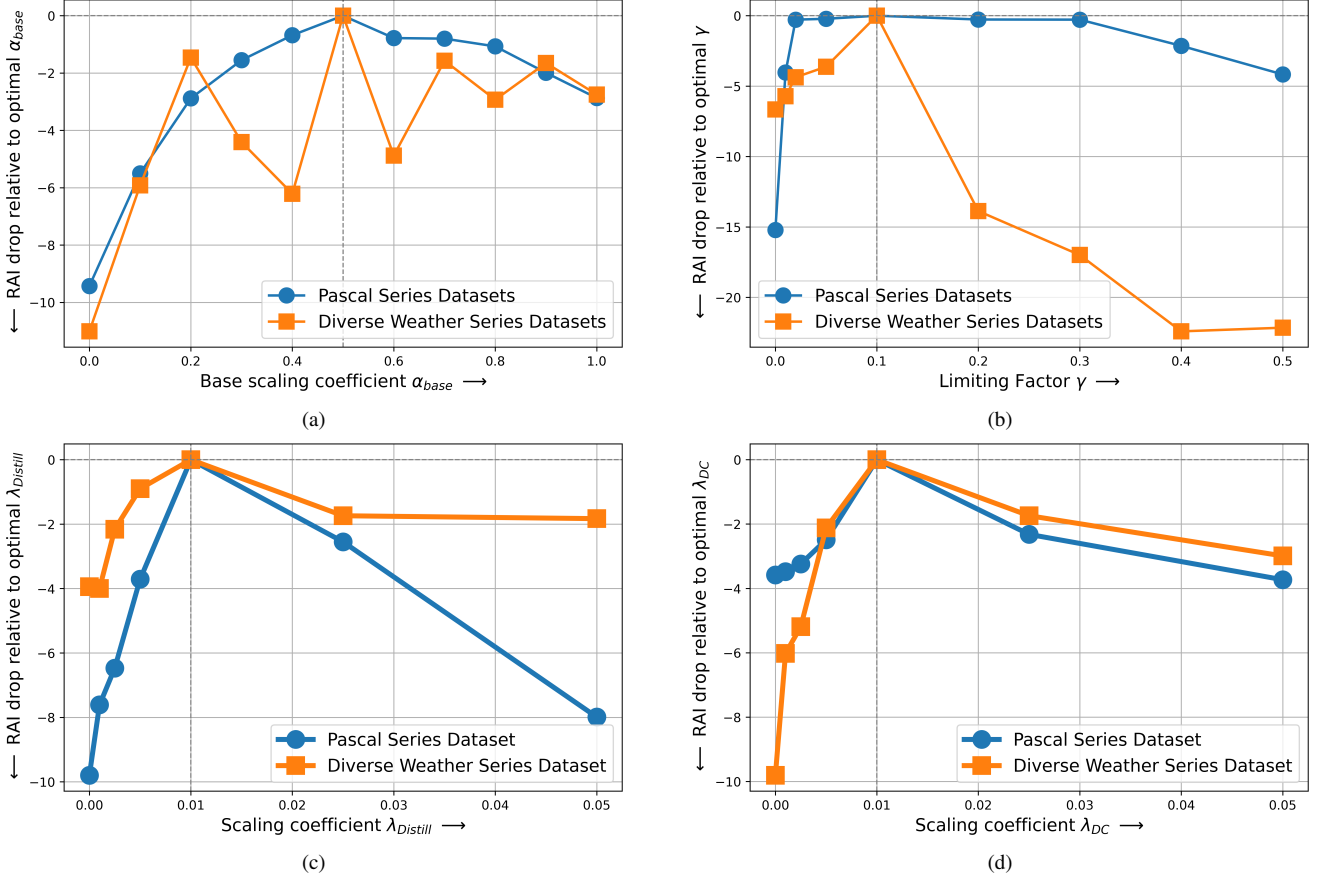
Figure S2. **Sensitivity analysis of DuET approach to key hyperparameters.** (a) shows the effect of varying base scaling coefficient $\alpha_{base}$. (b) illustrates the impact of the limiting factor $\gamma$ on RAI. (c) and (d) depicts the effect of varying scaling coefficients $\lambda_{Distill}$ and $\lambda_{DC}$ on RAI respectively.

on both series of datasets.

Table S4. Influence of random domain and class permutations across three incremental tasks on the *Diverse Weather Series* dataset.

| $\mathcal{T}_1$ | $\mathcal{T}_2$ | $\mathcal{T}_3$ | Avg RI | Avg GI | RAI |
|---|---|---|---|---|---|
| Night Sunny [5:7] | Daytime Sunny [1:2] | Daytime Foggy [3:4] | 83.49 | 51.01 | 67.25 |
| Night Sunny [3:4] | Daytime Sunny [5:7] | Daytime Foggy [1:2] | 80.39 | 51.76 | 66.08 |
| Night Sunny [1:2] | Daytime Sunny [3:4] | Daytime Foggy [5:7] | 88.57 | 41.92 | 65.25 |
| Daytime Foggy [1:2] | Night Sunny [3:4] | Daytime Sunny [5:7] | 78.34 | 50.54 | 64.44 |
| Daytime Sunny [1:2] | Daytime Foggy [3:4] | Night Sunny [5:7] | 88.33 | 35.97 | 62.15 |
| | **Standard Deviation** | | 4.12 | 6.26 | 1.72 |

## C.3. Influence of random class-domain order

In real-world incremental learning scenarios, the sequence in which new classes and domains are introduced can influence knowledge retention and generalisation. To check the sensitivity of the proposed DuET approach to such variations, we conducted experiments with shuffled class orders while keeping the same domain progression (top three rows in Table S4) and shuffled domain orders while keeping the same class sets (bottom three rows). We observe

that Avg RI remains consistently high across all permutations ($> 78\%$), and there are minor variations in RAI and Avg RI with standard deviations of **1.72** and **4.12**, respectively. However, the slight variation in Avg GI, with a standard deviation of **6.26**, stems from domain shifts affecting generalisation. This indicates that the proposed DuET approach maintains a consistent performance irrespective of the randomness in class-domain orders.

## C.4. Complexity Analysis

Table S11 presents a detailed complexity analysis of various methods evaluated in the multi-phase experiment: *Watercolour [1:3] → Comic [4:6] → Clipart [7:13] → VOC [14:20]*. Table S11 compares computational complexity in terms of GFLOPs, trainable parameters (in millions), average inference speed (in milliseconds), and average memory footprint (in gigabytes) across all incremental tasks. While training time (in hours) as evaluated on a single NVIDIA A100-PCIE-40GB on *Daytime Sunny [1:4] → Night Rainy [5:7]* experiment is reported in Table S14. The results

demonstrate that DuET retains the real-time detection capabilities of YOLO11n [10], effectively transforming it into a robust real-time incremental object detector with only a minimal increase in memory footprint (**0.244 GB**) compared to its unaltered counterpart, Sequential FT (**0.235 GB**). Furthermore, since the proposed DuET approach does not modify the base detector architecture, it preserves the same GFLOPs and the number of trainable parameters.

In contrast to other model-merging algorithms [7, 8, 19, 20], which require storing task vectors—and consequently, model weights—for every task, our approach is designed to be more efficient and scalable. DuET maintains only two shared task vectors at any given task: $\tau_{\text{old}}$ (derived from the previous phase's model weights) and $\tau_{\text{curr}}$ (derived from the current phase's model weights), along with the pre-trained model weights. This design utilizes the fact that knowledge from earlier tasks, $\mathcal{T}_1, \mathcal{T}_2, \ldots, \mathcal{T}_{t-2}$, is already encapsulated within the previous phase's weights, $\theta s_{t-1}$. Consequently, DuET avoids the overhead of maintaining a complete history of task vectors, resulting in a consistent memory footprint across all incremental tasks ($\mathcal{T}_t, t \geq 2$), while in case of other TA approaches, the memory footprint grows linearly with the number of tasks (Figure S3).

## D. Implementation Details

Our implementation is primarily based on the Ultralytics framework[1](v8.3.9), with YOLO11n [10] primarily serving as the base detector, also extending to other variants of YOLO11 and RT-DETR [18]. Following the default configuration provided by Ultralytics, we used AdamW [17] optimiser with auto lr find and OneCycleLR scheduler, keeping a batch size of 64. For every task, we trained the detector for 100 epochs, keeping five warm-up epochs with a higher initial learning rate by a factor of 10. For the base task ($t = 1$), we use the default weight decay of 0.0005, while for incremental tasks ($t \geq 2$), we slightly increase it to 0.001 to prevent overfitting to new tasks and help prevent catastrophic forgetting. The same protocol is used while preparing other baselines for a fair comparison. Moreover, unlike LDB [22] and CL-DETR [16], we keep all layers trainable during incremental training to ensure that shared task vectors effectively capture the shift in shared knowledge across incremental tasks.

## E. Comprehensive Results

### E.1. Detailed analysis of Quantitative Results:

Tables S7 to S14 present the comprehensive results of various methods on different DuIOD experiments across multiple base detectors. We conducted a total of seven DuIOD experiments—five two-phase and two multi-phase experi-

ments—three from the *Pascal Series* and the remaining four from the *Diverse Weather Series* datasets.

We provide detailed results for the five two-phase experiments: *VOC[1:10]* → *Clipart[11:20]* (Table S7), *Clipart [1:10]* → *VOC [11:20]* (Table S8), *Daytime Sunny [1:4]* → *Night Sunny [5:7]* (Table S12), *Night Sunny [1:4]* → *Daytime Sunny [5:7]* (Table S9) and *Daytime Sunny [1:4]* → *Night Rainy [5:7]* (Table S14). Meanwhile, the results for the two multi-phase experiments—*Watercolor [1:3]* → *Comic [4:6]* → *Clipart [7:13]* → *VOC [14:20]* and *Night Sunny [1:2]* → *Daytime Sunny [3:4]* → *Daytime Foggy [5:7]*—are presented in Tables S10 and S13, respectively.

We evaluate DuET across all DuIOD experiments using five detection backbones: DeformableDETR [26], YOLO11n & YOLO11x [10], and RTDETR-l & RTDETR-x [18]. Our results show that DuET consistently outperforms the baselines in nearly all DuIOD experiments, demonstrating its effectiveness in addressing the DuIOD task. Notably, DuET outperforms both CL-DETR [16] & LDB [22] on their respective backbones (Deformable DETR [26] & ViTDet [12]) with a **+5.97%** and **+11.98%** RAI gain, preserving **80.3%** vs. **66.85%** and **50.99%** vs. **41.74%** Avg RI respectively (refer Tables S7 to S14). This indicates that **gains are method-specific and not backbone dependent**.

### E.2. Qualitative Visualizations

In Figure S4, we present qualitative visualisations for various methods on the task sequence: *Watercolour [1:3]* → *Comic [4:6]* → *Clipart [7:13]* → *VOC [14:20]*. The detection results are shown for unseen classes: Watercolour [4:6], Comic [1:3], Clipart [1:6], and VOC [1:13] in the final task, $\mathcal{T}_4$. Our observations indicate that DuET consistently outperforms other methods by accurately detecting most objects across different domains and classes. Notably, in the second row (Comic [1:3]), the `bicycle` class, which was learned in $\mathcal{T}_1$ (Watercolor [1:3]), and the `person` class, introduced in $\mathcal{T}_2$ (Comic [4:6]), are examined. We observe that only DuET successfully retains the knowledge from $\mathcal{T}_1$ and correctly detects the `bicycle` class in the unseen Comic [1:3] domain. A similar trend is observed in the fourth row (VOC [1:13]). Additionally, in the third row, the `car` class, introduced in $\mathcal{T}_1$, is not detected by other methods in the unseen Clipart [1:6] domain, whereas DuET consistently identifies it. These qualitative results further emphasise DuET's effectiveness in adapting to unseen classes across different domains in the DuIOD task. A similar trend is observed in the *Diverse Weather Series* (Figure S5), where DuET consistently outperforms other methods by accurately detecting most objects across various domains and classes.

---

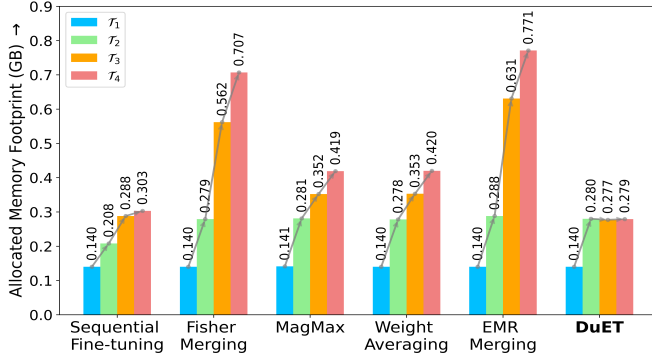[1]https://github.com/ultralytics/ultralytics

Figure S3. Comparison of allocated memory footprint (in GB) for various model-merging approaches on the multi-phase experiment with four tasks: Watercolor [1:3] ($\mathcal{T}_1$) → Comic [4:6] ($\mathcal{T}_2$) → Clipart [7:13] ($\mathcal{T}_3$) → VOC [14:20] ($\mathcal{T}_4$).

Table S5. **Dataset Statistics:** Class-wise distribution across different domains in *Pascal Series* datasets.

| Class ID | Class Name | Watercolor [9] | Comic [9] | Clipart [9] | VOC [4] |
|---|---|---|---|---|---|
| 1 | bicycle | ✓ | ✓ | ✓ | ✓ |
| 2 | bird | ✓ | ✓ | ✓ | ✓ |
| 3 | car | ✓ | ✓ | ✓ | ✓ |
| 4 | cat | ✓ | ✓ | ✓ | ✓ |
| 5 | dog | ✓ | ✓ | ✓ | ✓ |
| 6 | person | ✓ | ✓ | ✓ | ✓ |
| 7 | aeroplane | | | ✓ | ✓ |
| 8 | boat | | | ✓ | ✓ |
| 9 | bottle | | | ✓ | ✓ |
| 10 | bus | | | ✓ | ✓ |
| 11 | chair | | | ✓ | ✓ |
| 12 | cow | | | ✓ | ✓ |
| 13 | diningtable | | | ✓ | ✓ |
| 14 | horse | | | ✓ | ✓ |
| 15 | motorbike | | | ✓ | ✓ |
| 16 | pottedplant | | | ✓ | ✓ |
| 17 | sheep | | | ✓ | ✓ |
| 18 | sofa | | | ✓ | ✓ |
| 19 | train | | | ✓ | ✓ |
| 20 | tvmonitor | | | ✓ | ✓ |
| **Total Classes** | | 6 | 6 | 20 | 20 |
| **Train Images** | | 1000 | 1000 | 500 | 16551 |
| **Val Images** | | 1000 | 1000 | 500 | 4952 |

Table S6. **Dataset Statistics:** Class-wise distribution across different weather conditions in *Diverse Weather Series* datasets.

| Class ID | Class Name | Daytime Sunny [25] | Night Sunny [25] | Daytime Foggy [2, 6] |
|---|---|---|---|---|
| 1 | bike | ✓ | ✓ | ✓ |
| 2 | bus | ✓ | ✓ | ✓ |
| 3 | car | ✓ | ✓ | ✓ |
| 4 | motor | ✓ | ✓ | ✓ |
| 5 | person | ✓ | ✓ | ✓ |
| 6 | rider | ✓ | ✓ | ✓ |
| 7 | truck | ✓ | ✓ | ✓ |
| **Total Classes** | | 7 | 7 | 7 |
| **Train Images** | | 19317 | 25868 | 1829 |
| **Val Images** | | 8289 | 7756 | 688 |

## F. Background Shift

Background shift is a major issue in IOD scenarios [15, 21], where previously learned object categories, if unannotated in subsequent tasks, are treated as background. DuET mitigates this shift by explicitly decomposing model parameters into shared and task-specific components, and then merging these through TA with dynamic, layer-wise retention and adaptation weights (see Section 3.4). This strategy ensures that parameters crucial for previously learned object representations remain stable, thus preventing catastrophic forgetting and minimising the likelihood of previously learned objects being erroneously classified as background when they become unlabeled in subsequent tasks.

## G. Dataset Statistics

Following prior works [3, 11, 22, 23], we evaluate the DuET approach on two dataset series: the *Pascal Series* and the *Diverse Weather Series*, which cover diverse environmental conditions and domain variations, respectively. Table S5 presents the class-wise distribution, capturing cross-domain variations across four different domains: Watercolour, Comic, Clipart, and VOC. Following [11, 14], we combined the PASCAL VOC 2007 and 2012 [4] datasets to form the VOC domain, while the Watercolour, Comic, and Clipart domains were taken from [9]. Watercolour and Comic domains consist of six object categories, forming a subset of the 20 object categories present in Clipart and VOC. We used class splits of $10+10$ and $3+3+7+7$ with Class IDs as mentioned in Table S5 to conduct two-phase and multi-phase DuIOD experiments on the *Pascal Series* datasets, respectively. Similarly, Table S6 presents the class-wise distribution across three different weather conditions: Daytime Sunny, Night Sunny, and Daytime Foggy. Following [11, 22], the datasets are taken from BDD100k [25], Foggy Cityscapes [2], and Adverse Weather [6]. Since each domain contains a common set of seven classes, we used class splits of $4+3$ and $2+2+3$, with Class IDs sorted in alphabetical order (as shown in Table S6), to perform two-phase and multi-phase DuIOD experiments on the *Diverse Weather Series* datasets, respectively.

## H. Limitations and future works

Since DuET is a task vector-based model-merging approach, it inherits the limitations of existing task vector-based methods, and hence it cannot be generalised to models trained from scratch and requires access to pre-trained object detectors to calculate shared task vectors at each incremental task. This is a common limitation of task vector-based methods. Moreover, DuET merges shared task vectors through a weighted linear interpolation mechanism, which may be suboptimal for highly heterogeneous class shifts or extreme domain variations across incremental tasks. Future work could explore more sophisticated non-linear merging approaches to better capture the shared knowledge across tasks.

Table S7. Results of various methods on VOC [1:10] → Clipart [11:20] with different base detectors. Among columns, best in **bold**, second best _underlined_.

| Method | Base Detector | T1 VOC [1:10] | T2: Clipart [11:20] | | | | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|---|---|---|---|---|
| | | | Old VOC [1:10] | New Clipart [11:20] | Unseen Clipart [1:10] | Unseen VOC [11:20] | | | |
| LDB [22] | ViTDet | $74.9_{\pm0.7}$ | $50.10_{\pm0.5}$ | $22.30_{\pm0.8}$ | $8.90_{\pm0.4}$ | $9.60_{\pm0.6}$ | $66.89_{\pm0.5}$ | $18.76_{\pm0.3}$ | $42.83_{\pm0.4}$ |
| **DuET (Ours)** | ViTDet | $74.9_{\pm0.2}$ | $54.20_{\pm0.4}$ | $17.60_{\pm0.2}$ | $17.90_{\pm0.3}$ | $11.80_{\pm0.4}$ | $72.36_{\pm0.2}$ | $32.63_{\pm0.3}$ | $52.50_{\pm0.2}$ |
| CL-DETR [16] | Deformable DETR | $56.1_{\pm0.6}$ | $38.29_{\pm0.7}$ | $9.04_{\pm0.5}$ | $9.22_{\pm0.8}$ | $3.02_{\pm0.4}$ | $68.29_{\pm0.3}$ | $40.72_{\pm0.4}$ | $54.51_{\pm0.3}$ |
| **DuET (Ours)** | Deformable DETR | $56.1_{\pm0.6}$ | $42.32_{\pm0.3}$ | $4.10_{\pm0.2}$ | $15.63_{\pm0.4}$ | $1.68_{\pm0.4}$ | $75.48_{\pm0.6}$ | $72.37_{\pm0.3}$ | _$73.93_{\pm0.5}$_ |
| Sequential FT | RTDETR-l | $87.1_{\pm0.6}$ | $0.00_{\pm0.0}$ | $55.00_{\pm0.7}$ | $0.00_{\pm0.0}$ | $32.20_{\pm0.5}$ | $0.00_{\pm0.0}$ | $18.85_{\pm0.6}$ | $9.43_{\pm0.4}$ |
| LwF [13] | RTDETR-l | $87.1_{\pm0.6}$ | $3.13_{\pm0.2}$ | $25.90_{\pm0.8}$ | $1.30_{\pm0.3}$ | $18.50_{\pm0.4}$ | $3.59_{\pm0.7}$ | $12.30_{\pm0.5}$ | $7.95_{\pm0.6}$ |
| ERD [5] | RTDETR-l | $87.1_{\pm0.6}$ | $1.74_{\pm0.3}$ | $56.00_{\pm0.8}$ | $1.42_{\pm0.2}$ | $37.80_{\pm0.5}$ | $2.00_{\pm0.1}$ | $23.74_{\pm0.6}$ | $12.87_{\pm0.8}$ |
| **DuET (Ours)** | RTDETR-l | $87.1_{\pm0.6}$ | $46.10_{\pm0.1}$ | $68.00_{\pm0.2}$ | $28.10_{\pm0.4}$ | $62.20_{\pm0.4}$ | $52.93_{\pm0.7}$ | $68.20_{\pm0.5}$ | $60.57_{\pm0.6}$ |
| Sequential FT | RTDETR-x | $89.2_{\pm0.5}$ | $0.00_{\pm0.0}$ | $56.52_{\pm0.7}$ | $0.33_{\pm0.2}$ | $30.43_{\pm0.6}$ | $0.00_{\pm0.0}$ | $17.66_{\pm0.4}$ | $8.83_{\pm0.3}$ |
| LwF [13] | RTDETR-x | $89.2_{\pm0.5}$ | $20.40_{\pm0.5}$ | $28.80_{\pm0.3}$ | $17.00_{\pm0.7}$ | $17.90_{\pm0.4}$ | $22.87_{\pm0.6}$ | $28.27_{\pm0.5}$ | $25.57_{\pm0.2}$ |
| ERD [5] | RTDETR-x | $89.2_{\pm0.5}$ | $22.60_{\pm0.7}$ | $55.20_{\pm0.8}$ | $3.83_{\pm0.3}$ | $32.80_{\pm0.4}$ | $25.34_{\pm0.5}$ | $22.73_{\pm0.1}$ | $24.04_{\pm0.8}$ |
| **DuET (Ours)** | RTDETR-x | $89.2_{\pm0.5}$ | $60.50_{\pm0.5}$ | $26.80_{\pm0.4}$ | $49.00_{\pm0.3}$ | $22.50_{\pm0.2}$ | $67.83_{\pm0.4}$ | $64.93_{\pm0.6}$ | $66.38_{\pm0.7}$ |
| Sequential FT | YOLO11n | $80.4_{\pm0.3}$ | $0.60_{\pm0.1}$ | $36.70_{\pm0.8}$ | $1.02_{\pm0.3}$ | $17.40_{\pm0.4}$ | $0.75_{\pm0.2}$ | $12.86_{\pm0.4}$ | $6.81_{\pm0.3}$ |
| LwF [13] | YOLO11n | $80.4_{\pm0.3}$ | $58.40_{\pm0.8}$ | $3.96_{\pm0.3}$ | $28.60_{\pm0.7}$ | $5.00_{\pm0.2}$ | $72.64_{\pm0.3}$ | $33.74_{\pm0.5}$ | $53.19_{\pm0.4}$ |
| ERD [5] | YOLO11n | $80.4_{\pm0.3}$ | $55.20_{\pm0.4}$ | $20.60_{\pm0.7}$ | $30.50_{\pm0.5}$ | $16.70_{\pm0.8}$ | $68.66_{\pm0.4}$ | $43.68_{\pm0.3}$ | $56.17_{\pm0.6}$ |
| **DuET (Ours)** | YOLO11n | $80.4_{\pm0.3}$ | $70.30_{\pm0.3}$ | $8.45_{\pm0.3}$ | $33.80_{\pm0.3}$ | $12.80_{\pm0.3}$ | **$87.44_{\pm0.2}$** | $44.54_{\pm0.1}$ | $65.99_{\pm0.3}$ |
| Sequential FT | YOLO11x | $88.4_{\pm0.5}$ | $0.00_{\pm0.0}$ | $43.50_{\pm0.8}$ | $0.00_{\pm0.0}$ | $16.10_{\pm0.6}$ | $0.00_{\pm0.0}$ | $10.13_{\pm0.7}$ | $5.07_{\pm0.4}$ |
| LwF [13] | YOLO11x | $88.4_{\pm0.5}$ | $57.30_{\pm0.5}$ | $38.00_{\pm0.3}$ | $40.30_{\pm0.7}$ | $30.30_{\pm0.4}$ | $64.82_{\pm0.6}$ | _$74.57_{\pm0.2}$_ | $69.70_{\pm0.3}$ |
| ERD [5] | YOLO11x | $88.4_{\pm0.5}$ | $23.70_{\pm0.7}$ | $46.80_{\pm0.8}$ | $26.00_{\pm0.5}$ | $23.60_{\pm0.3}$ | $26.81_{\pm0.4}$ | $50.66_{\pm0.2}$ | $38.74_{\pm0.8}$ |
| **DuET (Ours)** | YOLO11x | $88.4_{\pm0.5}$ | $74.30_{\pm0.3}$ | $52.40_{\pm0.2}$ | $44.50_{\pm0.1}$ | $46.80_{\pm0.1}$ | _$84.05_{\pm0.5}$_ | **$90.73_{\pm0.3}$** | **$87.39_{\pm0.6}$** |

Table S8. Results of various methods on Clipart [1:10] → VOC [11:20] with different base detectors. Among columns, best in **bold**, second best _underlined_.

| Method | Base Detector | T1 Clipart [1:10] | T2: VOC [11:20] | | | | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|---|---|---|---|---|
| | | | Old Clipart [1:10] | New VOC [11:20] | Unseen VOC [1:10] | Unseen Clipart [11:20] | | | |
| LDB [22] | ViTDet | $36.4_{\pm0.4}$ | $16.30_{\pm0.3}$ | $23.80_{\pm0.5}$ | $7.10_{\pm0.2}$ | $9.10_{\pm0.6}$ | $44.78_{\pm0.7}$ | $15.81_{\pm0.4}$ | $30.30_{\pm0.8}$ |
| **DuET (Ours)** | ViTDet | $36.4_{\pm0.2}$ | $31.20_{\pm0.3}$ | $34.50_{\pm0.3}$ | $24.40_{\pm0.3}$ | $1.60_{\pm0.1}$ | $85.71_{\pm0.1}$ | $18.23_{\pm0.2}$ | $51.97_{\pm0.3}$ |
| CL-DETR [16] | Deformable DETR | $10.5_{\pm0.6}$ | $8.88_{\pm0.7}$ | $27.02_{\pm0.4}$ | $3.88_{\pm0.3}$ | $10.33_{\pm0.8}$ | $84.57_{\pm0.5}$ | **$54.35_{\pm0.2}$** | _$69.46_{\pm0.7}$_ |
| **DuET (Ours)** | Deformable DETR | $10.5_{\pm0.2}$ | $9.54_{\pm0.1}$ | $20.08_{\pm0.1}$ | $3.17_{\pm0.2}$ | $10.23_{\pm0.2}$ | **$90.86_{\pm0.1}$** | _$53.22_{\pm0.2}$_ | **$72.04_{\pm0.2}$** |
| Sequential FT | RTDETR-l | $44.2_{\pm0.5}$ | $0.00_{\pm0.0}$ | $81.50_{\pm0.7}$ | $0.00_{\pm0.0}$ | $30.80_{\pm0.6}$ | $0.00_{\pm0.0}$ | $29.79_{\pm0.8}$ | $14.90_{\pm0.4}$ |
| LwF [13] | RTDETR-l | $44.2_{\pm0.4}$ | $2.81_{\pm0.2}$ | $66.00_{\pm0.7}$ | $0.73_{\pm0.3}$ | $37.20_{\pm0.5}$ | $6.36_{\pm0.6}$ | $36.39_{\pm0.4}$ | $21.38_{\pm0.8}$ |
| ERD [5] | RTDETR-l | $44.2_{\pm0.5}$ | $0.37_{\pm0.3}$ | $81.20_{\pm0.6}$ | $2.81_{\pm0.4}$ | $27.50_{\pm0.7}$ | $0.84_{\pm0.2}$ | $28.21_{\pm0.8}$ | $14.53_{\pm0.5}$ |
| **DuET (Ours)** | RTDETR-l | $44.2_{\pm0.1}$ | $37.80_{\pm0.2}$ | $8.17_{\pm0.1}$ | $29.40_{\pm0.2}$ | $13.20_{\pm0.2}$ | $85.52_{\pm0.1}$ | $29.64_{\pm0.2}$ | _$57.58_{\pm0.1}$_ |
| Sequential FT | RTDETR-x | $47.0_{\pm0.6}$ | $0.00_{\pm0.0}$ | $81.60_{\pm0.7}$ | $0.00_{\pm0.0}$ | $35.70_{\pm0.5}$ | $0.00_{\pm0.0}$ | $37.27_{\pm0.8}$ | $18.64_{\pm0.3}$ |
| LwF [13] | RTDETR-x | $47.0_{\pm0.5}$ | $2.42_{\pm0.3}$ | $64.30_{\pm0.6}$ | $1.25_{\pm0.2}$ | $35.70_{\pm0.8}$ | $5.15_{\pm0.4}$ | $37.97_{\pm0.7}$ | $21.56_{\pm0.5}$ |
| ERD [5] | RTDETR-x | $47.0_{\pm0.4}$ | $0.67_{\pm0.2}$ | $82.00_{\pm0.7}$ | $0.93_{\pm0.3}$ | $34.40_{\pm0.6}$ | $1.43_{\pm0.5}$ | $36.43_{\pm0.8}$ | $18.93_{\pm0.4}$ |
| **DuET (Ours)** | RTDETR-x | $47.0_{\pm0.2}$ | $41.10_{\pm0.2}$ | $4.64_{\pm0.2}$ | $21.30_{\pm0.1}$ | $5.27_{\pm0.2}$ | _$87.45_{\pm0.2}$_ | $17.44_{\pm0.2}$ | $52.45_{\pm0.2}$ |
| Sequential FT | YOLO11n | $47.1_{\pm0.7}$ | $0.00_{\pm0.0}$ | $73.60_{\pm0.6}$ | $0.00_{\pm0.0}$ | $29.10_{\pm0.5}$ | $0.00_{\pm0.0}$ | $30.12_{\pm0.8}$ | $15.06_{\pm0.4}$ |
| LwF [13] | YOLO11n | $47.1_{\pm0.6}$ | $31.40_{\pm0.7}$ | $4.00_{\pm0.5}$ | $20.30_{\pm0.3}$ | $5.36_{\pm0.8}$ | $66.67_{\pm0.4}$ | $18.17_{\pm0.2}$ | $42.42_{\pm0.7}$ |
| ERD [5] | YOLO11n | $47.1_{\pm0.5}$ | $33.20_{\pm0.3}$ | $0.72_{\pm0.6}$ | $20.70_{\pm0.4}$ | $0.63_{\pm0.7}$ | $70.49_{\pm0.2}$ | $13.53_{\pm0.8}$ | $42.01_{\pm0.5}$ |
| **DuET (Ours)** | YOLO11n | $47.1_{\pm0.2}$ | $32.70_{\pm0.1}$ | $44.00_{\pm0.2}$ | $21.70_{\pm0.2}$ | $26.10_{\pm0.1}$ | $69.43_{\pm0.2}$ | $40.51_{\pm0.2}$ | $54.97_{\pm0.1}$ |
| Sequential FT | YOLO11x | $36.3_{\pm0.5}$ | $0.00_{\pm0.0}$ | $77.50_{\pm0.6}$ | $0.00_{\pm0.0}$ | $33.50_{\pm0.8}$ | $0.00_{\pm0.0}$ | $38.15_{\pm0.7}$ | $19.08_{\pm0.4}$ |
| LwF [13] | YOLO11x | $36.3_{\pm0.4}$ | $25.10_{\pm0.3}$ | $0.96_{\pm0.8}$ | $13.00_{\pm0.5}$ | $1.59_{\pm0.6}$ | $69.15_{\pm0.4}$ | $9.16_{\pm0.7}$ | $39.16_{\pm0.3}$ |
| ERD [5] | YOLO11x | $36.3_{\pm0.5}$ | $29.40_{\pm0.7}$ | $0.52_{\pm0.3}$ | $12.90_{\pm0.6}$ | $0.68_{\pm0.2}$ | $80.99_{\pm0.8}$ | $8.07_{\pm0.4}$ | $44.53_{\pm0.5}$ |
| **DuET (Ours)** | YOLO11x | $36.3_{\pm0.2}$ | $19.30_{\pm0.1}$ | $1.25_{\pm0.2}$ | $6.06_{\pm0.2}$ | $3.03_{\pm0.1}$ | $53.17_{\pm0.2}$ | $6.88_{\pm0.1}$ | $30.03_{\pm0.2}$ |

Table S9. Results of various methods on Daytime Sunny [1:4] → Night Sunny [5:7] with different base detectors. Among columns, best in **bold**, second best <u>underlined</u>.

| Method | Base Detector | T1 Daytime Sunny [1:4] | T2: Night Sunny [5:7] Old Daytime Sunny [1:4] | New Night Sunny [5:7] | Unseen Night Sunny [1:4] | Unseen Daytime Sunny [5:7] | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|---|---|---|---|---|
| LDB [22] | VitDet | $45.3_{\pm0.6}$ | $0.50_{\pm0.3}$ | $15.10_{\pm0.4}$ | $0.30_{\pm0.5}$ | $16.90_{\pm0.7}$ | $1.10_{\pm0.2}$ | $22.41_{\pm0.3}$ | $11.76_{\pm0.6}$ |
| **DuET (Ours)** | VitDet | $45.3_{\pm0.2}$ | $12.48_{\pm0.3}$ | $11.60_{\pm0.3}$ | $4.33_{\pm0.2}$ | $9.60_{\pm0.2}$ | $27.55_{\pm0.2}$ | $28.22_{\pm0.1}$ | $27.89_{\pm0.2}$ |
| CL-DETR [16] | Deformable DETR | $46.3_{\pm0.4}$ | $27.41_{\pm0.5}$ | $31.94_{\pm0.6}$ | $19.85_{\pm0.3}$ | $32.55_{\pm0.4}$ | $59.20_{\pm0.2}$ | <u>$54.96_{\pm0.5}$</u> | $57.08_{\pm0.4}$ |
| **DuET (Ours)** | Deformable DETR | $46.3_{\pm0.2}$ | $39.1_{\pm0.1}$ | $15.06_{\pm0.2}$ | $28.17_{\pm0.2}$ | $4.33_{\pm0.1}$ | $84.45_{\pm0.2}$ | $33.45_{\pm0.1}$ | $58.95_{\pm0.2}$ |
| Sequential FT | RTDETR-l | $57.2_{\pm0.5}$ | $0.00_{\pm0.0}$ | $77.40_{\pm0.7}$ | $2.52_{\pm0.3}$ | $39.80_{\pm0.6}$ | $0.00_{\pm0.0}$ | $35.36_{\pm0.8}$ | $17.68_{\pm0.4}$ |
| LwF [13] | RTDETR-l | $57.2_{\pm0.4}$ | $0.15_{\pm0.2}$ | $76.40_{\pm0.7}$ | $0.03_{\pm0.1}$ | $41.50_{\pm0.5}$ | $0.26_{\pm0.2}$ | $35.01_{\pm0.8}$ | $17.64_{\pm0.4}$ |
| ERD [5] | RTDETR-l | $57.2_{\pm0.5}$ | $0.09_{\pm0.1}$ | $80.50_{\pm0.7}$ | $0.04_{\pm0.1}$ | $39.80_{\pm0.6}$ | $0.16_{\pm0.2}$ | $33.59_{\pm0.8}$ | $16.88_{\pm0.4}$ |
| **DuET (Ours)** | RTDETR-l | $57.2_{\pm0.2}$ | $27.30_{\pm0.1}$ | $8.63_{\pm0.2}$ | $20.10_{\pm0.2}$ | $7.88_{\pm0.1}$ | $47.73_{\pm0.2}$ | $21.00_{\pm0.1}$ | $34.37_{\pm0.2}$ |
| Sequential FT | RTDETR-x | $61.0_{\pm0.6}$ | $0.00_{\pm0.0}$ | $84.80_{\pm0.7}$ | $0.00_{\pm0.0}$ | $40.80_{\pm0.5}$ | $0.00_{\pm0.0}$ | $33.77_{\pm0.8}$ | $16.89_{\pm0.3}$ |
| LwF [13] | RTDETR-x | $61.0_{\pm0.5}$ | $0.57_{\pm0.2}$ | $79.10_{\pm0.7}$ | $0.61_{\pm0.3}$ | $40.60_{\pm0.6}$ | $0.93_{\pm0.2}$ | $34.03_{\pm0.8}$ | $17.48_{\pm0.4}$ |
| ERD [5] | RTDETR-x | $61.0_{\pm0.4}$ | $0.81_{\pm0.2}$ | $84.80_{\pm0.7}$ | $0.98_{\pm0.3}$ | $38.90_{\pm0.6}$ | $1.33_{\pm0.2}$ | $32.87_{\pm0.8}$ | $17.10_{\pm0.4}$ |
| **DuET (Ours)** | RTDETR-x | $61.0_{\pm0.2}$ | $34.40_{\pm0.1}$ | $6.51_{\pm0.2}$ | $28.10_{\pm0.2}$ | $6.02_{\pm0.1}$ | $56.39_{\pm0.2}$ | $24.15_{\pm0.1}$ | $40.27_{\pm0.2}$ |
| Sequential FT | YOLO11n | $49.4_{\pm0.3}$ | $0.00_{\pm0.0}$ | $62.20_{\pm0.5}$ | $12.60_{\pm0.4}$ | $35.90_{\pm0.3}$ | $0.00_{\pm0.0}$ | $45.88_{\pm0.6}$ | $22.94_{\pm0.3}$ |
| LwF [13] | YOLO11n | $49.4_{\pm0.2}$ | $27.60_{\pm0.4}$ | $0.34_{\pm0.6}$ | $21.30_{\pm0.3}$ | $0.67_{\pm0.5}$ | $55.87_{\pm0.3}$ | $21.88_{\pm0.7}$ | $38.88_{\pm0.6}$ |
| ERD [5] | YOLO11n | $49.4_{\pm0.5}$ | $33.00_{\pm0.4}$ | $34.00_{\pm0.3}$ | $26.10_{\pm0.6}$ | $29.10_{\pm0.7}$ | $66.80_{\pm0.5}$ | $53.04_{\pm0.3}$ | $59.92_{\pm0.4}$ |
| **DuET (Ours)** | YOLO11n | $49.4_{\pm0.2}$ | $43.50_{\pm0.1}$ | $22.20_{\pm0.3}$ | $31.60_{\pm0.2}$ | $27.40_{\pm0.1}$ | $88.06_{\pm0.2}$ | **$56.95_{\pm0.1}$** | $72.51_{\pm0.2}$ |
| Sequential FT | YOLO11x | $64.2_{\pm0.6}$ | $12.50_{\pm0.4}$ | $68.60_{\pm0.7}$ | $18.80_{\pm0.3}$ | $46.00_{\pm0.8}$ | $19.47_{\pm0.2}$ | $47.60_{\pm0.5}$ | $33.54_{\pm0.4}$ |
| LwF [13] | YOLO11x | $64.2_{\pm0.7}$ | $62.10_{\pm0.6}$ | $0.04_{\pm0.2}$ | $42.40_{\pm0.8}$ | $0.01_{\pm0.1}$ | <u>$96.73_{\pm0.5}$</u> | $27.79_{\pm0.4}$ | $62.26_{\pm0.8}$ |
| ERD [5] | YOLO11x | $64.2_{\pm0.6}$ | $62.20_{\pm0.7}$ | $0.06_{\pm0.2}$ | $42.70_{\pm0.8}$ | $0.07_{\pm0.1}$ | $95.95_{\pm0.5}$ | $28.04_{\pm0.4}$ | $62.46_{\pm0.8}$ |
| **DuET (Ours)** | YOLO11x | $64.2_{\pm0.2}$ | $61.60_{\pm0.1}$ | $12.90_{\pm0.2}$ | $44.70_{\pm0.2}$ | $17.10_{\pm0.1}$ | **$96.88_{\pm0.2}$** | $42.41_{\pm0.1}$ | <u>$69.18_{\pm0.2}$</u> |

Table S10. Results of various methods on Night Sunny [1:4] → Daytime Sunny [5:7] with different base detectors. Among columns, best in **bold**, second best <u>underlined</u>.

| Method | Base Detector | T1 Night Sunny [1:4] | T2: Daytime Sunny [5:7] Old Night Sunny [1:4] | New Daytime Sunny [5:7] | Unseen Daytime Sunny [1:4] | Unseen Night Sunny [5:7] | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|---|---|---|---|---|
| LDB [22] | ViTDet | $37.0_{\pm0.6}$ | $0.40_{\pm0.4}$ | $18.30_{\pm0.7}$ | $0.10_{\pm0.2}$ | $14.30_{\pm0.5}$ | $1.08_{\pm0.3}$ | $20.66_{\pm0.7}$ | $10.87_{\pm0.8}$ |
| **DuET (Ours)** | VitDet | $37.0_{\pm0.2}$ | $5.50_{\pm0.1}$ | $18.33_{\pm0.1}$ | $3.33_{\pm0.5}$ | $14.70_{\pm0.2}$ | $14.86_{\pm0.1}$ | $24.80_{\pm0.2}$ | $19.83_{\pm0.1}$ |
| CL-DETR [16] | Deformable DETR | $48.8_{\pm0.7}$ | $25.90_{\pm0.8}$ | $45.70_{\pm0.5}$ | $19.88_{\pm0.6}$ | $41.33_{\pm0.2}$ | $53.04_{\pm0.7}$ | <u>$55.64_{\pm0.3}$</u> | $54.34_{\pm0.8}$ |
| **DuET (Ours)** | Deformable DETR | $48.8_{\pm0.2}$ | $38.81_{\pm0.1}$ | $4.24_{\pm0.2}$ | $29.09_{\pm0.1}$ | $7.59_{\pm0.2}$ | $79.48_{\pm0.1}$ | $37.69_{\pm0.2}$ | $58.59_{\pm0.2}$ |
| Sequential FT | RTDETR-l | $70.0_{\pm0.6}$ | $0.00_{\pm0.0}$ | $58.90_{\pm0.7}$ | $0.00_{\pm0.0}$ | $40.80_{\pm0.5}$ | $0.00_{\pm0.0}$ | $24.64_{\pm0.8}$ | $12.32_{\pm0.3}$ |
| LwF [13] | RTDETR-l | $70.0_{\pm0.7}$ | $8.73_{\pm0.8}$ | $21.80_{\pm0.5}$ | $6.00_{\pm0.6}$ | $8.37_{\pm0.4}$ | $12.47_{\pm0.7}$ | $10.30_{\pm0.3}$ | $11.39_{\pm0.8}$ |
| ERD [5] | RTDETR-l | $70.0_{\pm0.6}$ | $1.09_{\pm0.7}$ | $57.60_{\pm0.5}$ | $1.86_{\pm0.8}$ | $43.10_{\pm0.4}$ | $1.56_{\pm0.7}$ | $27.65_{\pm0.6}$ | $14.61_{\pm0.3}$ |
| **DuET (Ours)** | RTDETR-l | $70.0_{\pm0.2}$ | $48.70_{\pm0.2}$ | $1.86_{\pm0.1}$ | $43.80_{\pm0.2}$ | $1.03_{\pm0.1}$ | $69.57_{\pm0.2}$ | $38.91_{\pm0.2}$ | $54.24_{\pm0.2}$ |
| Sequential FT | RTDETR-x | $73.3_{\pm0.7}$ | $0.00_{\pm0.0}$ | $58.80_{\pm0.8}$ | $0.00_{\pm0.0}$ | $44.10_{\pm0.6}$ | $0.00_{\pm0.0}$ | $24.39_{\pm0.7}$ | $12.20_{\pm0.3}$ |
| LwF [13] | RTDETR-x | $73.3_{\pm0.6}$ | $10.70_{\pm0.7}$ | $9.25_{\pm0.5}$ | $6.02_{\pm0.6}$ | $7.66_{\pm0.4}$ | $14.60_{\pm0.7}$ | $9.17_{\pm0.3}$ | $11.89_{\pm0.8}$ |
| ERD [5] | RTDETR-x | $73.3_{\pm0.7}$ | $7.03_{\pm0.8}$ | $57.00_{\pm0.5}$ | $0.93_{\pm0.6}$ | $42.50_{\pm0.4}$ | $9.59_{\pm0.7}$ | $24.27_{\pm0.8}$ | $16.93_{\pm0.3}$ |
| **DuET (Ours)** | RTDETR-x | $73.3_{\pm0.2}$ | $66.30_{\pm0.2}$ | $6.40_{\pm0.1}$ | $58.90_{\pm0.2}$ | $5.26_{\pm0.1}$ | <u>$90.45_{\pm0.2}$</u> | $51.19_{\pm0.2}$ | <u>$70.82_{\pm0.2}$</u> |
| Sequential FT | YOLO11n | $50.1_{\pm0.7}$ | $0.12_{\pm0.8}$ | $37.60_{\pm0.5}$ | $0.25_{\pm0.6}$ | $25.8_{\pm0.4}$ | $0.24_{\pm0.7}$ | $19.48_{\pm0.8}$ | $9.86_{\pm0.3}$ |
| LwF [13] | YOLO11n | $50.1_{\pm0.6}$ | $39.00_{\pm0.7}$ | $0.29_{\pm0.2}$ | $33.90_{\pm0.8}$ | $1.51_{\pm0.3}$ | $77.84_{\pm0.5}$ | $35.44_{\pm0.7}$ | $56.64_{\pm0.4}$ |
| ERD [5] | YOLO11n | $50.1_{\pm0.7}$ | $39.60_{\pm0.6}$ | $0.06_{\pm0.2}$ | $34.20_{\pm0.8}$ | $0.04_{\pm0.1}$ | $79.04_{\pm0.5}$ | $34.65_{\pm0.7}$ | $56.85_{\pm0.4}$ |
| **DuET (Ours)** | YOLO11n | $50.1_{\pm0.2}$ | $47.10_{\pm0.2}$ | $20.10_{\pm0.1}$ | $41.30_{\pm0.2}$ | $19.10_{\pm0.1}$ | **$94.01_{\pm0.2}$** | **$56.03_{\pm0.2}$** | **$75.02_{\pm0.2}$** |
| Sequential FT | YOLO11x | $76.3_{\pm0.7}$ | $0.33_{\pm0.8}$ | $50.20_{\pm0.5}$ | $2.62_{\pm0.6}$ | $44.60_{\pm0.4}$ | $0.43_{\pm0.7}$ | $28.51_{\pm0.8}$ | $14.47_{\pm0.3}$ |
| LwF [13] | YOLO11x | $76.3_{\pm0.8}$ | $67.50_{\pm0.7}$ | $0.66_{\pm0.2}$ | $50.10_{\pm0.8}$ | $0.23_{\pm0.1}$ | $88.47_{\pm0.5}$ | $39.16_{\pm0.7}$ | $63.82_{\pm0.4}$ |
| ERD [5] | YOLO11x | $76.3_{\pm0.7}$ | $67.70_{\pm0.8}$ | $0.46_{\pm0.2}$ | $51.30_{\pm0.8}$ | $0.25_{\pm0.1}$ | $88.73_{\pm0.5}$ | $40.10_{\pm0.7}$ | $64.42_{\pm0.4}$ |
| **DuET (Ours)** | YOLO11x | $76.3_{\pm0.2}$ | $22.60_{\pm0.1}$ | $43.70_{\pm0.2}$ | $21.60_{\pm0.2}$ | $40.00_{\pm0.1}$ | $29.62_{\pm0.2}$ | $40.58_{\pm0.2}$ | $35.10_{\pm0.2}$ |

Table S11. Computational complexity analysis of various methods evaluated on the multi-phase experiment: Watercolor [1:3] → Comic [4:6] → Clipart [7:13] → VOC [14:20]. The inference speed and memory footprint are averaged across all incremental tasks.

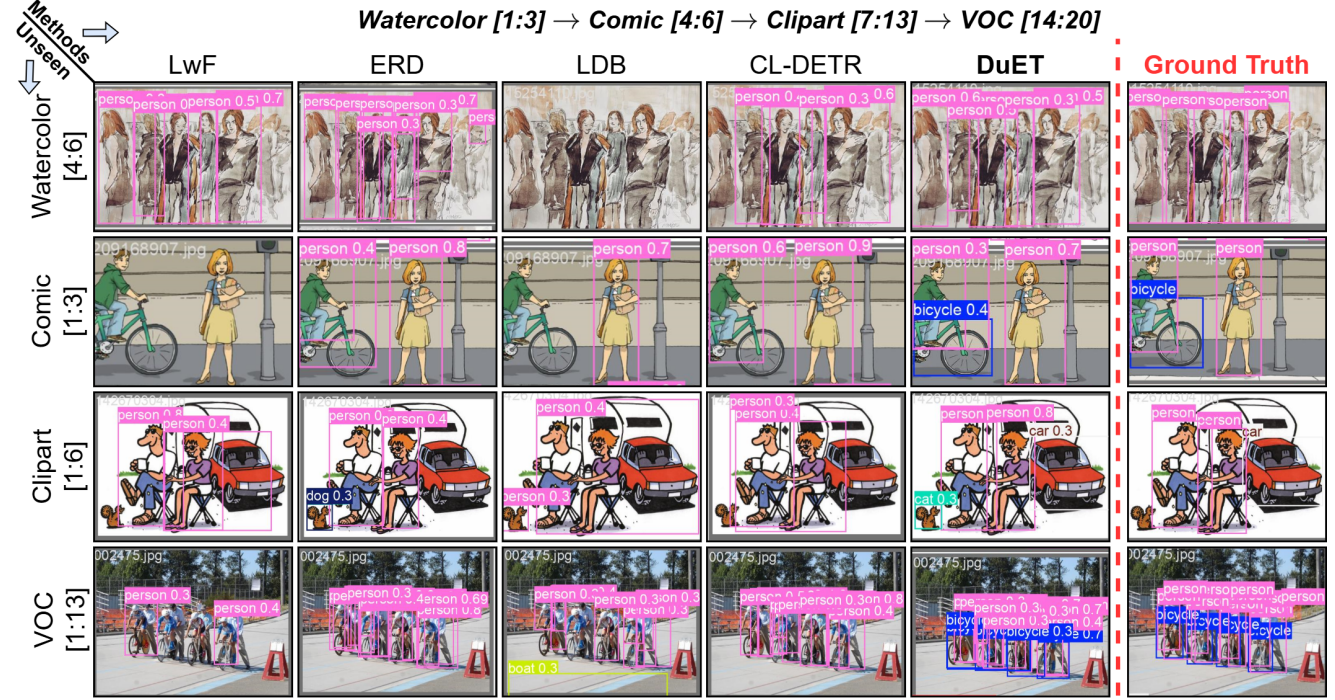| Method | Base Detector | GFLOPs | Trainable Params (M) | Avg. Inference Speed (ms) | Avg. Memory Footprint (GB) |
|---|---|---|---|---|---|
| Sequential FT | YOLO11n [10] | 6.3 | 2.58 | 9.150 | 0.235 |
| LwF [13] | YOLO11n [10] | 6.3 | 2.58 | 9.125 | 0.261 |
| ERD [5] | YOLO11n [10] | 6.3 | 2.58 | 9.075 | 0.257 |
| LDB [22] | ViTDet [12] | 1829.61 | 110.52 | 137.92 | 1.818 |
| CL-DETR [16] | Deformable DETR [26] | 11.77 | 39.85 | 39.075 | 0.789 |
| **DuET (Ours)** | YOLO11n [10] | 6.3 | 2.58 | 4.4 | 0.244 |



Figure S4. Qualitative comparisons on *Pascal Series* multi-phase experiment: Watercolour [1:3] → Comic [4:6] → Clipart [7:13] → VOC [14:20] for different methods on the DuIOD task. The four rows display detection results on unseen classes: Watercolour [4:6], Comic [1:3], Clipart [1:6], and VOC [1:13] on the final task $\mathcal{T}_4$. (zoomed in for best view).
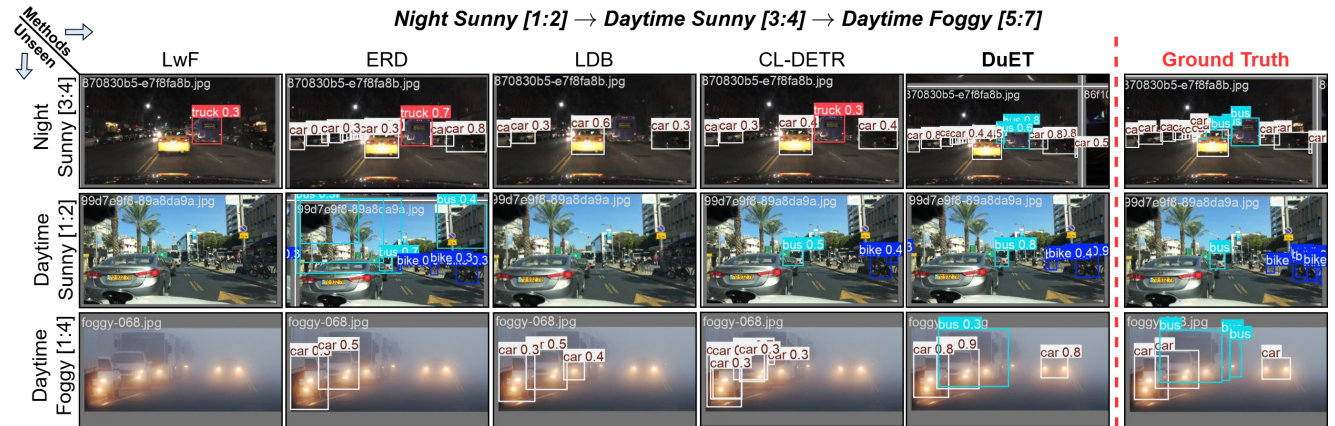


Figure S5. Qualitative comparisons on *Diverse Series* multi-phase experiment: Night Sunny [1:2] → Daytime Sunny [3:4] → Daytime Foggy [5:7] for different methods on the DuIOD task. The four rows display detection results on unseen classes: Night Sunny [3:4], Daytime Sunny [1:2], and Daytime Foggy [1:4] on the final task $\mathcal{T}_3$. (zoomed-in for best view).

Table S12. Results of various methods on Watercolor [1:3] → Comic [4:6] → Clipart [7:13] → VOC [14:20] with different base detectors. Among columns, best in **bold**, second best _underlined_.

| Method | Base Detector | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|
| LDB [22] | ViTDet | $86.08_{\pm0.6}$ | $19.57_{\pm0.5}$ | $52.83_{\pm0.4}$ |
| **DuET (Ours)** | ViTDet | $65.57_{\pm0.2}$ | $40.44_{\pm0.1}$ | $53.01_{\pm0.2}$ |
| CL-DETR [16] | Deformable DETR | $71.73_{\pm0.5}$ | $36.63_{\pm0.6}$ | $54.18_{\pm0.4}$ |
| **DuET (Ours)** | Deformable DETR | $88.54_{\pm0.2}$ | $34.81_{\pm0.1}$ | $61.68_{\pm0.2}$ |
| Sequential FT | RTDETR-l | $0.00_{\pm0.0}$ | $7.76_{\pm0.4}$ | $3.88_{\pm0.3}$ |
| LwF [13] | RTDETR-l | $0.40_{\pm0.2}$ | $13.47_{\pm0.7}$ | $6.94_{\pm0.5}$ |
| ERD [5] | RTDETR-l | $1.05_{\pm0.4}$ | $17.91_{\pm0.6}$ | $9.48_{\pm0.3}$ |
| **DuET (Ours)** | RTDETR-l | $24.75_{\pm0.2}$ | $22.89_{\pm0.2}$ | $23.82_{\pm0.1}$ |
| Sequential FT | RTDETR-x | $0.00_{\pm0.0}$ | $6.35_{\pm0.4}$ | $3.18_{\pm0.2}$ |
| LwF [13] | RTDETR-x | $65.51_{\pm0.8}$ | $16.69_{\pm0.3}$ | $41.10_{\pm0.7}$ |
| ERD [5] | RTDETR-x | $0.05_{\pm0.1}$ | $10.42_{\pm0.2}$ | $5.24_{\pm0.3}$ |
| **DuET (Ours)** | RTDETR-x | $40.80_{\pm0.2}$ | $18.88_{\pm0.2}$ | $29.84_{\pm0.1}$ |
| Sequential FT | YOLO11n | $0.00_{\pm0.0}$ | $11.05_{\pm0.5}$ | $5.53_{\pm0.3}$ |
| LwF [13] | YOLO11n | $52.66_{\pm0.6}$ | $17.01_{\pm0.4}$ | $34.84_{\pm0.3}$ |
| ERD [5] | YOLO11n | $54.76_{\pm0.5}$ | _$41.13_{\pm0.4}$_ | $47.95_{\pm0.7}$ |
| **DuET (Ours)** | YOLO11n | _$89.30_{\pm0.2}$_ | $42.60_{\pm0.1}$ | $65.95_{\pm0.2}$ |
| Sequential FT | YOLO11x | $0.00_{\pm0.0}$ | $10.21_{\pm0.4}$ | $5.11_{\pm0.3}$ |
| LwF [13] | YOLO11x | $10.56_{\pm0.3}$ | $17.46_{\pm0.5}$ | $14.01_{\pm0.2}$ |
| ERD [5] | YOLO11x | $54.19_{\pm0.6}$ | $8.49_{\pm0.2}$ | $31.34_{\pm0.8}$ |
| **DuET (Ours)** | YOLO11x | **$96.72_{\pm0.2}$** | $26.49_{\pm0.1}$ | _$61.61_{\pm0.2}$_ |

Table S13. Results of various methods on Night Sunny [1:2] → Daytime Sunny [3:4] → Daytime Foggy [5:7] with different base detectors. Among columns, best in **bold**, second best _underlined_.

| Method | Base Detector | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|
| LDB [22] | ViTDet | $50.50_{\pm0.6}$ | $5.42_{\pm0.7}$ | $27.96_{\pm0.5}$ |
| **DuET (Ours)** | VitDet | $39.87_{\pm0.2}$ | $17.12_{\pm0.1}$ | $28.50_{\pm0.2}$ |
| CL-DETR [16] | Deformable DETR | _$64.26_{\pm0.3}$_ | **$43.46_{\pm0.5}$** | _$53.86_{\pm0.6}$_ |
| **DuET (Ours)** | Deformable DETR | $62.99_{\pm0.2}$ | $45.11_{\pm0.1}$ | $54.05_{\pm0.2}$ |
| Sequential FT | RTDETR-l | $0.00_{\pm0.0}$ | $14.97_{\pm0.4}$ | $7.49_{\pm0.3}$ |
| LwF [13] | RTDETR-l | $14.76_{\pm0.2}$ | $1.24_{\pm0.3}$ | $8.00_{\pm0.5}$ |
| ERD [5] | RTDETR-l | $5.40_{\pm0.4}$ | $15.83_{\pm0.6}$ | $10.62_{\pm0.3}$ |
| **DuET (Ours)** | RTDETR-l | $20.76_{\pm0.2}$ | $10.55_{\pm0.2}$ | $15.66_{\pm0.1}$ |
| Sequential FT | RTDETR-x | $0.00_{\pm0.0}$ | $22.74_{\pm0.4}$ | $11.37_{\pm0.3}$ |
| LwF [13] | RTDETR-x | $7.62_{\pm0.3}$ | $8.86_{\pm0.2}$ | $8.24_{\pm0.4}$ |
| ERD [5] | RTDETR-x | $3.02_{\pm0.2}$ | $19.16_{\pm0.3}$ | $11.09_{\pm0.2}$ |
| **DuET (Ours)** | RTDETR-x | $27.39_{\pm0.2}$ | $23.65_{\pm0.1}$ | $25.52_{\pm0.2}$ |
| Sequential FT | YOLO11n | $0.00_{\pm0.0}$ | $30.51_{\pm0.6}$ | $15.26_{\pm0.5}$ |
| LwF [13] | YOLO11n | $27.94_{\pm0.7}$ | $23.78_{\pm0.5}$ | $25.86_{\pm0.6}$ |
| ERD [5] | YOLO11n | $44.60_{\pm0.6}$ | $39.40_{\pm0.5}$ | $42.00_{\pm0.3}$ |
| **DuET (Ours)** | YOLO11n | **$88.57_{\pm0.2}$** | _$41.92_{\pm0.1}$_ | **$65.25_{\pm0.2}$** |
| Sequential FT | YOLO11x | $0.00_{\pm0.0}$ | $18.16_{\pm0.4}$ | $9.08_{\pm0.3}$ |
| LwF [13] | YOLO11x | $19.57_{\pm0.3}$ | $28.44_{\pm0.5}$ | $24.01_{\pm0.2}$ |
| ERD [5] | YOLO11x | $45.85_{\pm0.6}$ | $37.38_{\pm0.2}$ | $41.62_{\pm0.8}$ |
| **DuET (Ours)** | YOLO11x | $43.46_{\pm0.2}$ | $38.37_{\pm0.1}$ | $40.92_{\pm0.2}$ |

---

**Algorithm 1:** DuET Training Algorithm

**Input:** Pre-trained model weights: $\theta_0$, Sequence of tasks: $\{T_1, T_2, \ldots, T_T\}$.
**Output:** Final model weights: $\theta_T$.

1   Initialize model with pre-trained weights: $\theta_0$;
2   **for** $t = 1, 2, \ldots, T$ **do**
3     **if** $t = 1$ **then**
4       Train model on task $T_1$ using $\mathcal{L}_{\text{Detector}}$;
5       Update weights: $\theta_1 \leftarrow \theta_0 - \eta \cdot \nabla_\theta \mathcal{L}_{\text{Detector}}$;
6       Decompose weights:
7       $\theta_0 \rightarrow [\theta_{s_0}, \theta_{\tau_0}], \quad \theta_1 \rightarrow [\theta_{s_1}, \theta_{\tau_1}]$;
8       Compute shared task vector: $\tau_{s_1} = \theta_{s_1} - \theta_{s_0}$;
9       Compute task-specific task vector: $\tau_{\tau_1} \leftarrow \theta_{\tau_1}$;
10     **else**
11       Initialize: $\theta_t \leftarrow \theta_{t-1}$ ;
12       Train model on task $T_t$ using $\mathcal{L}_{\text{Total}}$;
13       Update weights: $\theta_t \leftarrow \theta_{t-1} - \eta \cdot \nabla_\theta \mathcal{L}_{\text{Total}}$;
14       Decompose weights:
15       $\theta_{t-1} \rightarrow [\theta_{s_{t-1}}, \theta_{\tau_{t-1}}], \quad \theta_t \rightarrow [\theta_{s_t}, \theta_{\tau_t}]$;
16       Compute shared task vectors:
17       $\tau_{\text{old}} = \tau_{s_{t-1}} = \theta_{s_{t-1}} - \theta_{s_0}$;
18       $\tau_{\text{curr}} = \tau_{s_t} = \theta_{s_t} - \theta_{s_0}$;
19       Compute task-specific task vectors:
20       $\tau_{t_{\text{old}}} = \theta_{\tau_{t-1}}, \quad \tau_{t_{\text{curr}}} = \theta_{\tau_t}$;
21       Update shared weights using DuET:
       $(\theta_{s_t})_{\text{incre}} \leftarrow \textbf{DuET}(\tau_{\text{curr}}, \tau_{\text{old}}, \theta_{s_0})$;
22       Update task-specific weights:
       $(\theta_{\tau_t})_{\text{incre}} \leftarrow [\tau_{t_{\text{old}}}, \tau_{t_{\text{curr}}}]$;
23       Load new updated weights:
       $\theta_t \leftarrow [(\theta_{s_t})_{\text{incre}}, (\theta_{\tau_t})_{\text{incre}}]$;
24     **end**
25   **end**
26   **return** $\theta_T$;

---

**Algorithm 2:** DuET Task Arithmetic Algorithm

**Input: Parameters:** Shared pre-trained weights: $\theta_{s_0}$, Old Task Vector: $\tau_{\text{old}}$, Current Task Vector: $\tau_{\text{curr}}$
      **Hyperparameters:** Limiting Factor: $\gamma$, Base Scaling Coefficient: $\alpha_{\text{base}}$, Numerical Stability constant: $\epsilon$
**Output:** Updated inremental shared weights: $\theta_{s_t}^{\text{incre}}$

1   **for** *Model Layer:* $l = 1, 2, \ldots, L$ **do**
2     $p_l = \dfrac{\|\tau_{\text{old}}^l\| - \|\tau_{\text{curr}}^l\|}{\|\tau_{\text{old}}^l + \tau_{\text{curr}}^l\| + \epsilon}$,
3     $\delta_l = \gamma \cdot \tanh(p_l)$
4     $\alpha_l = \alpha_{\text{base}} + \text{clamp}(\delta_l, -\gamma, \gamma)$
5     $\beta_l = 1 - \alpha_l$
6     $(\theta_{s_t}^l)_{\text{incre}} = \theta_{s_0}^l + \alpha_l \cdot \tau_{\text{old}}^l + \beta_l \cdot \tau_{\text{curr}}^l$
7   **end**
8   $(\theta_{s_t})_{\text{incre}} = \{(\theta_{s_t}^l)_{\text{incre}}\}_{l=1}^L$
9   **return** $(\theta_{s_t})_{\text{incre}}$

Table S14. Results of various methods on Daytime Sunny [1:4] → Night Rainy [5:7] with different base detectors. Among columns, best in **bold**, second best <u>underlined</u>. Training time is reported on single NVIDIA A100-PCIE-40GB.

| Method | Base Detector | Training Time $(\mathcal{T}_1 + \mathcal{T}_2)$ (hours) | T1 Daytime Sunny [1:4] | T2: Night Rainy [5:7] | | | | Avg RI (%) | Avg GI (%) | RAI (%) |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Old Daytime Sunny [1:4] | New Night Rainy [5:7] | Unseen Night Rainy [1:4] | Unseen Daytime Sunny [5:7] | | | |
| LDB [22] | VitDet | 23.73 | $45.3_{\pm0.6}$ | $1.40_{\pm0.3}$ | $8.10_{\pm0.4}$ | $0.02_{\pm0.5}$ | $16.10_{\pm0.7}$ | $3.09_{\pm0.2}$ | $21.02_{\pm0.3}$ | $12.05_{\pm0.6}$ |
| **DuET (Ours)** | VitDet | 24.13 | $45.3_{\pm0.2}$ | $2.02_{\pm0.3}$ | $16.10_{\pm0.1}$ | $1.58_{\pm0.3}$ | $15.20_{\pm0.2}$ | $4.46_{\pm0.2}$ | $22.47_{\pm0.1}$ | $13.47_{\pm0.2}$ |
| CL-DETR [16] | Deformable DETR | 91.96 | $46.3_{\pm0.4}$ | $26.26_{\pm0.4}$ | $9.75_{\pm0.5}$ | $7.41_{\pm0.6}$ | $14.88_{\pm0.3}$ | $56.72_{\pm0.4}$ | <u>$53.93_{\pm0.2}$</u> | $55.33_{\pm0.5}$ |
| **DuET (Ours)** | Deformable DETR | 92.63 | $46.3_{\pm0.2}$ | $26.75_{\pm0.1}$ | $3.33_{\pm0.2}$ | $7.98_{\pm0.2}$ | $13.9_{\pm0.1}$ | <u>$57.78_{\pm0.2}$</u> | **$54.58_{\pm0.1}$** | <u>$56.18_{\pm0.2}$</u> |
| CL-DETR [16] | RT-DETR-l | 33.87 | $57.2_{\pm0.4}$ | $5.11_{\pm0.5}$ | $28.1_{\pm0.6}$ | $5.04_{\pm0.3}$ | $14.3_{\pm0.4}$ | $8.93_{\pm0.2}$ | $20.82_{\pm0.5}$ | $14.88_{\pm0.4}$ |
| **DuET (Ours)** | RT-DETR-l | 33.89 | $57.2_{\pm0.4}$ | $11.8_{\pm0.1}$ | $14.90_{\pm0.2}$ | $12.50_{\pm0.2}$ | $17.30_{\pm0.1}$ | $20.63_{\pm0.2}$ | $36.50_{\pm0.1}$ | $28.57_{\pm0.2}$ |
| LwF [13] | YOLO11n | 11.01 | $49.4_{\pm0.2}$ | $21.50_{\pm0.4}$ | $0.17_{\pm0.6}$ | $9.36_{\pm0.3}$ | $0.55_{\pm0.5}$ | $43.52_{\pm0.3}$ | $17.78_{\pm0.7}$ | $30.65_{\pm0.6}$ |
| ERD [5] | YOLO11n | 11.50 | $49.4_{\pm0.2}$ | $25.60_{\pm0.4}$ | $16.70_{\pm0.3}$ | $12.20_{\pm0.6}$ | $17.60_{\pm0.7}$ | $51.82_{\pm0.5}$ | $38.96_{\pm0.3}$ | $45.39_{\pm0.4}$ |
| **DuET (Ours)** | YOLO11n | 11.02 | $49.4_{\pm0.2}$ | $43.30_{\pm0.1}$ | $2.67_{\pm0.3}$ | $14.00_{\pm0.2}$ | $8.93_{\pm0.1}$ | **$87.65_{\pm0.2}$** | $34.18_{\pm0.1}$ | **$60.92_{\pm0.2}$** |

# References

[1] Li Chen, Chunyan Yu, and Lvcai Chen. A new knowledge distillation for incremental object detection. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–7. IEEE, 2019. 1

[2] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 6

[3] Muhammad Sohail Danish, Muhammad Haris Khan, Muhammad Akhtar Munir, M Saquib Sarfraz, and Mohsen Ali. Improving single domain-generalized object detection: A focus on diversification and alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17732–17742, 2024. 6

[4] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010. 6

[5] Tao Feng, Mang Wang, and Hangjie Yuan. Overcoming catastrophic forgetting in incremental object detection via elastic response distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9427–9436, 2022. 7, 8, 9, 10, 11

[6] Mahmoud Hassaballah, Mourad A Kenk, Khan Muhammad, and Shervin Minaee. Vehicle detection and tracking in adverse weather using a deep learning framework. *IEEE transactions on intelligent transportation systems*, 22(7):4230–4242, 2020. 6

[7] Chenyu Huang, Peng Ye, Tao Chen, Tong He, Xiangyu Yue, and Wanli Ouyang. Emr-merging: Tuning-free high-performance model merging. *arXiv preprint arXiv:2405.17461*, 2024. 3, 5

[8] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*, 2022. 3, 5

[9] Naoto Inoue, Ryosuke Furuta, Toshihiko Yamasaki, and Kiyoharu Aizawa. Cross-domain weakly-supervised object detection through progressive domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5001–5009, 2018. 6

[10] Glenn Jocher and Jing Qiu. Ultralytics yolo11, 2024. 2, 3, 5, 9

[11] Madhu Kiran, Marco Pedersoli, Jose Dolz, Louis-Antoine Blais-Morin, Eric Granger, et al. Incremental multi-target domain adaptation for object detection with efficient domain transfer. *Pattern Recognition*, 129:108771, 2022. 6

[12] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *European conference on computer vision*, pages 280–296. Springer, 2022. 5, 9

[13] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017. 2, 7, 8, 9, 10, 11

[14] Xialei Liu, Hao Yang, Avinash Ravichandran, Rahul Bhotika, and Stefano Soatto. Multi-task incremental learning for object detection. *arXiv preprint arXiv:2002.05347*, 2020. 6

[15] Yuyang Liu, Yang Cong, Dipam Goswami, Xialei Liu, and Joost van de Weijer. Augmented box replay: Overcoming foreground shift for incremental object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11367–11377, 2023. 6

[16] Yaoyao Liu, Bernt Schiele, Andrea Vedaldi, and Christian Rupprecht. Continual detection transformer for incremental object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23799–23808, 2023. 5, 7, 8, 9, 10, 11

[17] Ilya Loshchilov, Frank Hutter, et al. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 5:5, 2017. 5

[18] Wenyu Lv, Shangliang Xu, Yian Zhao, Guanzhong Wang, Jinman Wei, Cheng Cui, Yuning Du, Qingqing Dang, and Yi Liu. Detrs beat yolos on real-time object detection, 2023. 2, 3, 5

[19] Daniel Marczak, Bartłomiej Twardowski, Tomasz Trzciński,

and Sebastian Cygert. Magmax: Leveraging model merging for seamless continual learning. In *European Conference on Computer Vision*, pages 379–395. Springer, 2025. 3, 5

[20] Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. *Advances in Neural Information Processing Systems*, 35:17703–17716, 2022. 3, 5

[21] Angelo G Menezes, Gustavo de Moura, Cézanne Alves, and André CPLF de Carvalho. Continual object detection: a review of definitions, strategies, and challenges. *Neural networks*, 161:476–493, 2023. 1, 6

[22] Xiang Song, Yuhang He, Songlin Dong, and Yihong Gong. Non-exemplar domain incremental object detection via learning domain bias. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 15056–15065, 2024. 5, 6, 7, 8, 9, 10, 11

[23] Aming Wu and Cheng Deng. Single-domain generalized object detection in urban scene via cyclic-disentangled self-distillation. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 847–856, 2022. 6

[24] Dongbao Yang, Yu Zhou, Aoting Zhang, Xurui Sun, Dayan Wu, Weiping Wang, and Qixiang Ye. Multi-view correlation distillation for incremental object detection. *Pattern Recognition*, 131:108863, 2022. 1

[25] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2636–2645, 2020. 6

[26] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 5, 9