

469 Appendix

470 In Sec. A, we present the proofs for the lemmas, theorems, and corollaries presented in the main body
 471 of our work. Sec. B discusses the correspondence of existing methods to specific cases within our
 472 framework. In Sec. C, we provide a detailed presentation of our final algorithm, UDIL, including an
 473 algorithmic description, a visual diagram, and implementation details. We introduce the experimental
 474 settings, including the evaluation metrics and specific training schemes. Finally, in Sec. D, we present
 475 additional empirical results with varying memory sizes and provide more visualization results.

476 A Proofs of Lemmas, Theorems, and Corollaries

477 Before proceeding to prove any lemmas or theorems, we first introduce three crucial additional lemmas
 478 that will be utilized in the subsequent sections. Among these, Lemma A.1 offers a generalization
 479 bound for any weighted summation of ERM losses across multiple domains. Furthermore, Lemma A.2
 480 provides a generalization bound for a weighted summation of *labeling functions* within a given domain.
 481 Lastly, we highlight Lemma 3 in [4] as Lemma A.3, which will be used to establish the upper bound
 482 for Lemma 3.3.

483 **Lemma A.1 (Generalization Bound of α -weighted Domains).** *Let \mathcal{H} be a hypothesis space of*
 484 *VC dimension d . Assume N_j denotes the number of the samples collected from domain j , and*
 485 *$N = \sum_j N_j$ is the total number of the examples collected from all domains. Then for any $\alpha_j > 0$*
 486 *and $\delta \in (0, 1)$, with probability at least $1 - \delta$:*

$$\sum_j \alpha_j \epsilon_{\mathcal{D}_j}(h) \leq \sum_j \alpha_j \widehat{\epsilon}_{\mathcal{D}_j}(h) + \sqrt{\left(\sum_j \frac{\alpha_j^2}{N_j}\right) \left(8d \log\left(\frac{2eN}{d}\right) + 8 \log\left(\frac{2}{\delta}\right)\right)}. \quad (17)$$

487 *Proof.* Suppose each domain \mathcal{D}_j has a deterministic ground-truth labeling function $f_j : \mathbb{R}^n \rightarrow \{0, 1\}$.
 488 Denote as $\widehat{\epsilon}_\alpha \triangleq \sum_j \alpha_j \widehat{\epsilon}_{\mathcal{D}_j}(h)$ the α -weighted empirical loss evaluated on different domains. Hence,

$$\widehat{\epsilon}_\alpha(h) = \sum_j \alpha_j \widehat{\epsilon}_{\mathcal{D}_j}(h) = \sum_j \alpha_j \frac{1}{N_j} \sum_{\mathbf{x} \in \mathcal{X}_j} \mathbb{1}_{h(\mathbf{x}) \neq f_j(\mathbf{x})} = \frac{1}{N} \sum_j \sum_{k=1}^{N_j} R_{j,k}, \quad (18)$$

489 where $R_{j,k} = \left(\frac{\alpha_j N_j}{N}\right) \cdot \mathbb{1}_{h(\mathbf{x}_k) \neq f_j(\mathbf{x}_k)}$ is a random variable that takes the values in $\left\{\frac{\alpha_j N_j}{N}, 0\right\}$. By the
 490 linearity of the expectation, we have $\epsilon_\alpha(h) = \mathbb{E}[\widehat{\epsilon}_\alpha(h)]$. Following [2, 36], we have

$$\mathbb{P}\{\exists h \in \mathcal{H}, \text{ s.t. } |\widehat{\epsilon}_\alpha(h) - \epsilon_\alpha(h)| \geq \epsilon\} \quad (19)$$

$$\leq 2 \cdot \mathbb{P}\left\{\sup_{h \in \mathcal{H}} |\widehat{\epsilon}_\alpha(h) - \widetilde{\epsilon}'_\alpha(h)| \geq \frac{\epsilon}{2}\right\} \quad (20)$$

$$\leq 2 \cdot \mathbb{P}\left\{\bigcup_{R_{j,k}, R'_{j,k}} \frac{1}{N} \left|\sum_j \sum_{k=1}^{N_j} (R_{j,k} - R'_{j,k})\right| \geq \frac{\epsilon}{2}\right\} \quad (21)$$

$$\leq 2\Pi_{\mathcal{H}}(2N) \exp\left\{-\frac{2(N\epsilon/2)^2}{\sum_j (N_j)(2\alpha_j N/N_j)^2}\right\} \quad (22)$$

$$= 2\Pi_{\mathcal{H}}(2N) \exp\left\{-\frac{\epsilon^2}{8 \sum_j (\alpha_j^2/N_j)}\right\} \quad (23)$$

$$\leq 2(2N)^d \exp\left\{-\frac{\epsilon^2}{8 \sum_j (\alpha_j^2/N_j)}\right\}, \quad (24)$$

491 where in Eqn. 20, $\widetilde{\epsilon}'_\alpha(h)$ is the α -weighted empirical loss evaluated on the ‘‘ghost’’ set of examples
 492 $\{\mathcal{X}'_j\}$; Eqn. 22 is yielded by applying Hoeffding’s inequalities [23] and introducing the growth
 493 function $\Pi_{\mathcal{H}}$ [2, 36, 56] at the same time; Eqn. 24 is achieved by using the fact $\Pi_{\mathcal{H}}(2N) \leq$
 494 $(e \cdot 2N/d)^d \leq (2N)^d$, where d is the VC-dimension of the hypothesis set \mathcal{H} . Finally, by setting Eqn. 24
 495 to δ and solve for the error tolerance ϵ will complete the proof. \square

496 **Lemma A.2 (Generalization Bound of β -weighted Labeling Functions).** *Let \mathcal{D} be a single domain*
 497 *and $\mathcal{X} = \{\mathbf{x}_i\}_i^N$ be a collection of samples drawn from \mathcal{D} ; \mathcal{H} is a hypothesis space of VC dimension*

498 *d. Suppose $\{f_j : \mathbb{R}^n \rightarrow \{0, 1\}\}_j$ is a set of different labeling functions. Then for any $\beta_j > 0$ and*
 499 *$\delta \in (0, 1)$, with probability at least $1 - \delta$:*

$$\sum_j \beta_j \epsilon_{\mathcal{D}}(h, f_j) \leq \sum_j \beta_j \widehat{\epsilon}_{\mathcal{D}}(h, f_j) + \left(\sum_j \beta_j\right) \sqrt{\frac{1}{N} \left(8d \log\left(\frac{2eN}{d}\right) + 8 \log\left(\frac{2}{\delta}\right)\right)}. \quad (25)$$

500 *Proof.* Denote as $\epsilon_{\beta}(h) \triangleq \sum_j \beta_j \epsilon_{\mathcal{D}}(h, f_j)$ the β -weighted error on domain \mathcal{D} and $\{f_j\}_j$ the set of
 501 the labeling functions, and $\widehat{\epsilon}_{\beta} \triangleq \sum_j \beta_j \widehat{\epsilon}_{\mathcal{D}}(h, f_j)$ as the β -weighted empirical loss evaluated on
 502 different labeling functions. We have

$$\widehat{\epsilon}_{\beta}(h) = \sum_j \beta_j \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{h(\mathbf{x}_i) \neq f_j(\mathbf{x}_i)} = \frac{1}{N} \sum_{i=1}^N \sum_j \beta_j \mathbb{1}_{h(\mathbf{x}_i) \neq f_j(\mathbf{x}_i)} \triangleq \frac{1}{N} \sum_{i=1}^N R_i, \quad (26)$$

503 where $R_i = \sum_j \beta_j \mathbb{1}_{h(\mathbf{x}_i) \neq f_j(\mathbf{x}_i)} \in [0, \sum_j \beta_j]$ is a new random variable.

504 Then we have

$$\mathbb{P}\{\exists h \in \mathcal{H}, \text{ s.t. } |\widehat{\epsilon}_{\beta}(h) - \epsilon_{\beta}(h)| \geq \epsilon\} \quad (27)$$

$$\leq 2 \cdot \mathbb{P}\left\{\sup_{h \in \mathcal{H}} |\widehat{\epsilon}_{\beta}(h) - \widehat{\epsilon}'_{\beta}(h)| \geq \frac{\epsilon}{2}\right\} \quad (28)$$

$$\leq 2 \cdot \mathbb{P}\left\{\bigcup_{R_i, R'_i} \frac{1}{N} \left|\sum_{i=1}^N (R_i - R'_i)\right| \geq \frac{\epsilon}{2}\right\} \quad (29)$$

$$\leq 2\Pi_{\mathcal{H}}(2N) \exp\left\{\frac{-2(N\epsilon/2)^2}{N \cdot (2\sum_j \beta_j)^2}\right\} \quad (30)$$

$$\leq 2(2N)^d \exp\left\{-\frac{N\epsilon^2}{8(\sum_j \beta_j)^2}\right\}, \quad (31)$$

505 where in Eqn. 28, $\widehat{\epsilon}'_{\beta}(h)$ is the β -weighted empirical loss evaluated on the “ghost” set of examples
 506 \mathcal{X}' ; Eqn. 30 is yielded by applying Hoeffding’s inequalities [23] and introducing the growth function
 507 $\Pi_{\mathcal{H}}$ [2, 36, 56] at the same time; Eqn. 31 is achieved by using the fact $\Pi_{\mathcal{H}}(2N) \leq (e \cdot 2N/d)^d \leq (2N)^d$,
 508 where d is the VC-dimension of the hypothesis set \mathcal{H} . Finally, by setting Eqn. 31 to δ and solve for
 509 the error tolerance ϵ will complete the proof. \square

510 Lemma A.2 asserts that altering or merging multiple target functions does not impact the generaliza-
 511 tion error term, as long as the sum of the weights for each loss $\sum_j \beta_j$ remains constant and the same
 512 dataset \mathcal{X} is used for estimation. Next we highlight the Lemma 3 in [4] again, as it will be utilized for
 513 proving 3.3.

514 **Lemma A.3.** *For any hypothesis $h, h' \in \mathcal{H}$ and any two different domains $\mathcal{D}, \mathcal{D}'$,*

$$|\epsilon_{\mathcal{D}}(h, h') - \epsilon_{\mathcal{D}'}(h, h')| \leq \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}, \mathcal{D}'). \quad (32)$$

515 *Proof.* By definition, we have

$$\begin{aligned} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}, \mathcal{D}') &= 2 \sup_{h, h' \in \mathcal{H}} |\mathbb{P}_{\mathbf{x} \sim \mathcal{D}}[h(\mathbf{x}) \neq h'(\mathbf{x})] - \mathbb{P}_{\mathbf{x} \sim \mathcal{D}'}[h(\mathbf{x}) \neq h'(\mathbf{x})]| \\ &= 2 \sup_{h, h' \in \mathcal{H}} |\epsilon_{\mathcal{D}}(h, h') - \epsilon_{\mathcal{D}'}(h, h')| \\ &\geq 2 |\epsilon_{\mathcal{D}}(h, h') - \epsilon_{\mathcal{D}'}(h, h')|. \end{aligned} \quad \square$$

516 Now we are ready to prove the main lemmas and theorems in the main body of our work.

517 **Lemma 3.1 (ERM-Based Generalization Bound).** *Let \mathcal{H} be a hypothesis space of VC dimension*
 518 *d . When domain t arrives, there are N_t data points from domain t and \tilde{N}_i data points from each*
 519 *previous domain $i < t$ in the memory bank. With probability at least $1 - \delta$, we have:*

$$\sum_{i=1}^t \epsilon_{\mathcal{D}_i}(h) \leq \sum_{i=1}^t \widehat{\epsilon}_{\mathcal{D}_i}(h) + \sqrt{\left(\frac{1}{N_t} + \sum_{i=1}^{t-1} \frac{1}{\tilde{N}_i}\right) \left(8d \log\left(\frac{2eN}{d}\right) + 8 \log\left(\frac{2}{\delta}\right)\right)}. \quad (33)$$

520 *Proof.* Simply using Lemma A.1 and setting $\alpha_i = 1$ for every $i \in [t]$ completes the proof. \square

521 **Lemma 3.2 (Intra-Domain Model-Based Bound).** Let $h \in \mathcal{H}$ be an arbitrary function in the
 522 hypothesis space \mathcal{H} , and H_{t-1} be the model trained after domain $t - 1$. The domain-specific error
 523 $\epsilon_{\mathcal{D}_i}(h)$ on the previous domain i has an upper bound:

$$\epsilon_{\mathcal{D}_i}(h) \leq \epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1}), \quad (34)$$

524 where $\epsilon_{\mathcal{D}_i}(h, H_{t-1}) \triangleq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_i}[h(\mathbf{x}) \neq H_{t-1}(\mathbf{x})]$.

525 *Proof.* By applying the triangle inequality [4] of the 0-1 loss function, we have

$$\begin{aligned} \epsilon_{\mathcal{D}_i}(h) &= \epsilon_{\mathcal{D}_i}(h, f_i) \\ &\leq \epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1}, f_i) \\ &= \epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1}). \end{aligned} \quad \square$$

526 **Lemma 3.3 (Cross-Domain Model-Based Bound).** Let $h \in \mathcal{H}$ be an arbitrary function in the
 527 hypothesis space \mathcal{H} , and H_{t-1} be the function trained after domain $t - 1$. The domain-specific error
 528 $\epsilon_{\mathcal{D}_i}(h)$ evaluated on the previous domain i then has an upper bound:

$$\epsilon_{\mathcal{D}_i}(h) \leq \epsilon_{\mathcal{D}_t}(h, H_{t-1}) + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \epsilon_{\mathcal{D}_i}(H_{t-1}), \quad (35)$$

529 where $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{P}, \mathcal{Q}) = 2 \sup_{h \in \mathcal{H}\Delta\mathcal{H}} |\Pr_{\mathbf{x} \sim \mathcal{P}}[h(\mathbf{x}) = 1] - \Pr_{\mathbf{x} \sim \mathcal{Q}}[h(\mathbf{x}) = 1]|$ denotes the $\mathcal{H}\Delta\mathcal{H}$ -
 530 divergence between distribution \mathcal{P} and \mathcal{Q} , and $\epsilon_{\mathcal{D}_t}(h, H_{t-1}) \triangleq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t}[h(\mathbf{x}) \neq H_{t-1}(\mathbf{x})]$.

531 *Proof.* By the triangle inequality used above and Lemma A.3, we have

$$\begin{aligned} \epsilon_{\mathcal{D}_i}(h) &\leq \epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1}) \\ &= \epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1}) - \epsilon_{\mathcal{D}_t}(h, H_{t-1}) + \epsilon_{\mathcal{D}_t}(h, H_{t-1}) \\ &\leq \epsilon_{\mathcal{D}_i}(H_{t-1}) + |\epsilon_{\mathcal{D}_i}(h, H_{t-1}) - \epsilon_{\mathcal{D}_t}(h, H_{t-1})| + \epsilon_{\mathcal{D}_t}(h, H_{t-1}) \\ &\leq \epsilon_{\mathcal{D}_t}(h, H_{t-1}) + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \epsilon_{\mathcal{D}_i}(H_{t-1}). \end{aligned} \quad \square$$

532 **Theorem 3.4 (Unified Generalization Bound for All Domains).** Let \mathcal{H} be a hypothesis space of
 533 VC dimension d . Let $N = N_t + \sum_{i=1}^{t-1} \tilde{N}_i$ denoting the total number of data points available to the
 534 training of current domain t , where N_t and \tilde{N}_i denote the numbers of data points collected at domain
 535 t and data points from the previous domain i in the memory bank, respectively. With probability at
 536 least $1 - \delta$, we have:

$$\begin{aligned} \sum_{i=1}^t \epsilon_{\mathcal{D}_i}(h) &\leq \left\{ \sum_{i=1}^{t-1} [\gamma_i \hat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \hat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1})] \right\} + \left\{ \hat{\epsilon}_{\mathcal{D}_t}(h) + \left(\sum_{i=1}^{t-1} \beta_i \right) \hat{\epsilon}_{\mathcal{D}_t}(h, H_{t-1}) \right\} \\ &\quad + \frac{1}{2} \sum_{i=1}^{t-1} \beta_i d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \sum_{i=1}^{t-1} (\alpha_i + \beta_i) \epsilon_{\mathcal{D}_i}(H_{t-1}) \\ &\quad + \sqrt{\left(\frac{(1 + \sum_{i=1}^{t-1} \beta_i)^2}{N_t} + \sum_{i=1}^{t-1} \frac{(\gamma_i + \alpha_i)^2}{\tilde{N}_i} \right) (8d \log \left(\frac{2eN}{d} \right) + 8 \log \left(\frac{2}{\delta} \right))} \\ &\triangleq g(h, H_{t-1}, \Omega), \end{aligned} \quad (36)$$

537 where $\hat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1}) = \frac{1}{N_i} \sum_{\mathbf{x} \in \tilde{\mathcal{X}}_i} \mathbb{1}_{h(\mathbf{x}) \neq H_{t-1}(\mathbf{x})}$, $\hat{\epsilon}_{\mathcal{D}_t}(h, H_{t-1}) = \frac{1}{N_t} \sum_{\mathbf{x} \in \mathcal{X}_t} \mathbb{1}_{h(\mathbf{x}) \neq H_{t-1}(\mathbf{x})}$, and
 538 $\Omega \triangleq \{\alpha_i, \beta_i, \gamma_i\}_{i=1}^{t-1}$.

539 *Proof.* By applying Lemma 3.2 and Lemma 3.3 to each of the past domains, we have

$$\begin{aligned} \epsilon_{\mathcal{D}_i}(h) &= (\alpha_i + \beta_i + \gamma_i) \epsilon_{\mathcal{D}_i}(h) \\ &\leq \gamma_i \epsilon_{\mathcal{D}_i}(h) + \alpha_i [\epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1})] \\ &\quad + \beta_i [\epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_t}(h, H_{t-1}) + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t)]. \end{aligned}$$

Table 4: Unification of Existing Methods under UDIL.

	α_i	β_i	γ_i	Transformed Objective	Condition
UDIL (Ours)	$[0, 1]$	$[0, 1]$	$[0, 1]$	-	-
LwF [29]	0	1	0	$\mathcal{L}_{\text{LwF}}(h) = \widehat{\ell}_{\mathcal{X}_t}(h) + \lambda_o \widehat{\ell}_{\mathcal{X}_t}(h, H_{t-1})$	$\lambda_o = t - 1$
ER [46]	0	0	1	$\mathcal{L}_{\text{ER}}(h) = \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{ B'_i /(t-1)}{ B_t } \widehat{\ell}_{B'_i}(h)$	$ B_t = \frac{ B'_1 }{(t-1)}$
DER++ [6]	1/2	0	1/2	$\mathcal{L}_{\text{DER++}}(h) = \widehat{\ell}_{B_t}(h) + \frac{1}{2} \sum_{i=1}^{t-1} \frac{ B'_i /(t-1)}{ B_t } [\widehat{\ell}_{B'_i}(h) + \widehat{\ell}_{B'_i}(h, H_{t-1})]$	$ B_t = \frac{ B'_1 }{(t-1)}$
CLS-ER [3]	$\lambda/\lambda+1$	0	$1/\lambda+1$	$\mathcal{L}_{\text{CLS-ER}}(h) = \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{1}{i-1} \widehat{\ell}_{B'_i}(h) + \sum_{i=1}^{t-1} \frac{\lambda}{i-1} \widehat{\ell}_{B'_i}(h, H_{t-1})$	$\lambda = t - 2$
ESM-ER [50]	$\lambda/\lambda+1$	0	$1/\lambda+1$	$\mathcal{L}_{\text{CLS-ER}}(h) = \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{1}{r(i-1)} \widehat{\ell}_{B'_i}(h) + \sum_{i=1}^{t-1} \frac{\lambda}{r(i-1)} \widehat{\ell}_{B'_i}(h, H_{t-1})$	$\begin{cases} \lambda = -1 + r(t-1) \\ r = 1 - e^{-1} \end{cases}$

540 Re-organizing the terms will give us

$$\begin{aligned}
\sum_{i=1}^t \epsilon_{\mathcal{D}_i}(h) &\leq \left\{ \sum_{i=1}^{t-1} [\gamma_i \epsilon_{\mathcal{D}_i}(h) + \alpha_i \epsilon_{\mathcal{D}_i}(h, H_{t-1})] \right\} + \left\{ \epsilon_{\mathcal{D}_t}(h) + \left(\sum_{i=1}^{t-1} \beta_i \right) \epsilon_{\mathcal{D}_t}(h, H_{t-1}) \right\} \\
&\quad + \frac{1}{2} \sum_{i=1}^{t-1} \beta_i d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \sum_{i=1}^{t-1} (\alpha_i + \beta_i) \epsilon_{\mathcal{D}_i}(H_{t-1}). \tag{37}
\end{aligned}$$

541 Then applying Lemma A.1 and Lemma A.2 jointly to Eqn. 37 will complete the proof. \square

542 B UDIL as a Unified Framework

543 In this section, we will delve into a comprehensive discussion of our UDIL framework, which serves
544 as a unification of numerous existing methods. It is important to note that we incorporate methods
545 designed for task incremental and class incremental scenarios that can be easily adapted to our domain
546 incremental learning. To provide clarity, we will present the corresponding coefficients $\{\alpha_i, \beta_i, \gamma_i\}$
547 of each method within our UDIL framework (refer to Table 4). Furthermore, we will explore the
548 conditions under which these coefficients are included in this unification process.

549 **Learning without Forgetting (LwF)** [29] was initially proposed for task-incremental learning,
550 incorporating a combination of *shared parameters* and *task-specific parameters*. This framework can
551 be readily extended to domain incremental learning by setting all “domain-specific” parameters to be
552 the same in a static model architecture. LwF was designed for the strict continual learning setting,
553 where no data from past tasks is accessible. To overcome this limitation, LwF records the predictions
554 of the history model H_{t-1} on the current data \mathcal{X}_t at the beginning of the new task t . Subsequently,
555 knowledge distillation (as defined in Definition 4.2) is performed to mitigate forgetting:

$$\mathcal{L}_{\text{old}}(h, H_{t-1}) \triangleq -\frac{1}{N_t} \sum_{\mathbf{x} \in \mathcal{X}_t} \sum_{k=1}^K [H_{t-1}(\mathbf{x})]_k \cdot [\log([h(\mathbf{x})]_k)] = \widehat{\ell}_{\mathcal{X}_t}(h, H_{t-1}), \tag{38}$$

556 where $H_{t-1}(\mathbf{x}), h(\mathbf{x}) \in \mathbb{R}^K$ are the class distribution of \mathbf{x} over K classes produced by the history
557 model and current model, respectively. The loss for learning the current task \mathcal{L}_{new} is defined as

$$\mathcal{L}_{\text{new}}(h) \triangleq -\frac{1}{N_t} \sum_{(\mathbf{x}, y) \in \mathcal{S}_t} \sum_{k=1}^K \mathbb{1}_{y=k} \cdot [\log([h(\mathbf{x})]_k)] = \widehat{\ell}_{\mathcal{X}_t}(h). \tag{39}$$

558 LwF uses a “loss balance weight” λ_o to balance two losses, which gives us its final loss for training:

$$\mathcal{L}_{\text{LwF}}(h) \triangleq \mathcal{L}_{\text{new}}(h) + \lambda_o \cdot \mathcal{L}_{\text{old}}(h, H_{t-1}). \tag{40}$$

559 In LwF, the default setting assumes the presence of two domains (tasks) with $\lambda_o = 1$. However, it
560 is possible to learn multiple domains continuously using LwF’s default configuration. To achieve

561 this, the current domain t can be weighed against the number of previous domains (1 versus $t - 1$).
 562 Specifically, if there is no preference for any particular domain, λ_o should be set to $t - 1$. Remarkably,
 563 this is equivalent to setting $\{\beta_i = 1, \alpha_i = \gamma_i = 0\}$ in our UDIL framework (Row 2 in Table 4).

564 **Experience Replay (ER) [46]** serves as the fundamental operation for replay-based continual learning
 565 methods. It involves storing and replaying a subset of examples from past domains during training.
 566 Following the description and implementation provided by [6], ER operates as follows: during each
 567 training iteration on domain t , a mini-batch B_t of examples is sampled from the current domain,
 568 along with a mini-batch B'_t from the memory. These two mini-batches are then concatenated into a
 569 larger mini-batch ($B_t \cup B'_t$), upon which average gradient descent is performed:

$$\mathcal{L}_{\text{ER}}(h) = \widehat{\ell}_{B_t \cup B'_t}(h) \quad (41)$$

$$= \frac{1}{|B_t| + |B'_t|} \sum_{(\mathbf{x}, y) \in B_t \cup B'_t} \sum_{k=1}^K \mathbb{1}_{y=k} \cdot [\log([h(\mathbf{x})]_k)] \quad (42)$$

$$= \frac{|B_t|}{|B_t| + |B'_t|} \widehat{\ell}_{B_t}(h) + \frac{|B'_t|}{|B_t| + |B'_t|} \widehat{\ell}_{B'_t}(h). \quad (43)$$

570 Suppose that each time the mini-batch of past-domain data is perfectly balanced, meaning that each
 571 domain has the same number of examples in B'_t . In this case, Eqn. 43 can be further decomposed as
 572 follows:

$$\mathcal{L}_{\text{ER}}(h) = \frac{|B_t|}{|B_t| + |B'_t|} \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{|B'_t|/(t-1)}{|B_t| + |B'_t|} \widehat{\ell}_{B'_i}(h), \quad (44)$$

573 where $B'_i = \{(\mathbf{x}, y) | (\mathbf{x}, y) \in (B'_t \cap M_i)\}$ is the subset of the mini-batch that belongs to domain i .

574 Now, by dividing both sides of Eqn. 44 by $(|B_t| + |B'_t|/|B_t|)$ and comparing it to Theorem 3.4, we can
 575 include ER in our UDIL framework when the condition $|B_t| = |B'_t|/(t-1)$ is satisfied. In this case,
 576 ER is equivalent to $\{\alpha_i = \beta_i = 0, \gamma_i = 1\}$ in UDIL (Row 3 in Table 4). It is important to note that
 577 this condition is not commonly met throughout the entire process of continual learning. It can be
 578 achieved by linearly scaling up the size of the mini-batch from the memory (which is feasible in the
 579 early domains) or by linearly scaling down the mini-batch from the current-domain data (which may
 580 cause a drop in model performance). It is worth mentioning that this incongruence *highlights the*
 581 *intrinsic bias of the original ER setting towards current domain learning* and cannot be rectified by
 582 adjusting the batch sizes of the current domain or the memory. However, it does not weaken our
 583 claim of unification.

584 **Dark Experience Replay (DER++) [6]** includes an additional dark experience replay, i.e., knowledge
 585 distillation on the past domain exemplars, compared to ER [46]. Now under the same assumptions
 586 (balanced sampling strategy and $|B_t| = |B'_t|/(t-1)$) as discussed for ER, we can utilize Eqn. 44 to
 587 transform the DER++ loss as follows:

$$\mathcal{L}_{\text{DER++}}(h) = \frac{|B_t|}{|B_t| + |B'_t|} \widehat{\ell}_{B_t}(h) + \frac{1}{2} \sum_{i=1}^{t-1} \frac{|B'_t|/(t-1)}{|B_t| + |B'_t|} \widehat{\ell}_{B'_i}(h) + \frac{1}{2} \sum_{i=1}^{t-1} \frac{|B'_t|/(t-1)}{|B_t| + |B'_t|} \widehat{\ell}_{B'_i}(h, H_{t-1}). \quad (45)$$

588 In this scenario, DER++ is equivalent to $\{\alpha_i = \gamma_i = 1/2, \beta_i = 0\}$ in UDIL (Row 4 in Table 4).

589 **Complementary Learning System based Experience Replay (CLS-ER) [3]** involves the mainte-
 590 nance of two history models, namely the plastic model $H_{t-1}^{(p)}$ and the stable model $H_{t-1}^{(s)}$, throughout
 591 the continual training process of the working model h . Following each update of the working model,
 592 the two history models are stochastically updated at different rates using exponential moving averages
 593 (EMA) of the working model’s *parameters*:

$$H_{t-1}^{(i)} \leftarrow \alpha^{(i)} \cdot H_{t-1}^{(i)} + (1 - \alpha^{(i)}) \cdot h, \quad i \in \{p, s\}, \quad (46)$$

594 where $\alpha^{(p)} \leq \alpha^{(s)}$ is set such that the plastic model undergoes rapid updates, allowing it to swiftly
 595 adapt to newly acquired knowledge, while the stable model maintains a “long-term memory” spanning
 596 multiple tasks. Throughout training, CLS-ER assesses the certainty generated by both history models
 597 and employs the logits from the more certain model as the target for knowledge distillation.

598 In the general formulation of the UDIL framework, the history model H_{t-1} is not required to be a
 599 single model with the same architecture as the current model h . In fact, if there are no constraints on

600 memory consumption, we have the flexibility to train and preserve a domain-specific model H_i for
 601 each domain i . During testing, we can simply select the prediction with the highest certainty from each
 602 domain-specific model. From this perspective, the “two-history-model system” employed in CLS-ER
 603 can be viewed as a specific and limited version of the all-domain history models. Consequently, we
 604 can combine the two models used in CLS-ER into a single history model H_{t-1} as follows:

$$H_{t-1}(\mathbf{x}) \triangleq \begin{cases} H_{t-1}^{(p)}(\mathbf{x}) & \text{if } [H_{t-1}^{(p)}(\mathbf{x})]_y > [H_{t-1}^{(s)}(\mathbf{x})]_y \\ H_{t-1}^{(s)}(\mathbf{x}) & \text{o.w.} \end{cases} \quad (47)$$

605 where $(\mathbf{x}, y) \in \mathcal{M}$ is an arbitrary exemplar stored in the memory bank.

606 At each iteration of training, CLS-ER samples a mini-batch B_t from the current domain and a
 607 mini-batch B'_t from the episodic memory. It then concatenates B_t and B'_t for the cross entropy loss
 608 minimization with the ground-truth labels, and uses B'_t to minimize the MSE loss between the logits
 609 of h and H_{t-1} . To align the loss formulation of CLS-ER with that of ESM-ER [50], here we consider
 610 the scenarios where the losses evaluated on B_t and B'_t are individually calculated, i.e., we consider
 611 $\widehat{\ell}_{B_t}(h) + \widehat{\ell}_{B'_t}(h)$ instead of $\widehat{\ell}_{B_t \cup B'_t}(h)$. Based on the assumption from [20], the MSE loss on the
 612 logits is equivalent to the cross-entropy loss on the predictions under certain conditions. Therefore,
 613 following the same balanced sampling strategy assumptions as in ER, the original CLS-ER training
 614 objective can be transformed as follows:

$$\mathcal{L}_{\text{CLS-ER}}(h) = \widehat{\ell}_{B_t}(h) + \widehat{\ell}_{B'_t}(h) + \lambda \widehat{\ell}_{B'_t}(h, H_{t-1}) \quad (48)$$

$$= \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{1}{t-1} \widehat{\ell}_{B'_i}(h) + \sum_{i=1}^{t-1} \frac{\lambda}{t-1} \widehat{\ell}_{B'_i}(h, H_{t-1}). \quad (49)$$

615 Therefore, by imposing the constraint $\alpha_i + \beta_i + \gamma_i = 1$, we find that $\lambda = t - 2$. Substituting this value
 616 back into $\lambda/t-1$ yields the equivalence that CLS-ER corresponds to $\{\alpha_i = \lambda/\lambda+1, \beta_i = 0, \gamma_i = 1/\lambda+1\}$
 617 in UDIL, where $\lambda = t - 2$ (Row 5 in Table 4).

618 **Error Sensitivity Modulation based Experience Replay (ESM-ER)** [50] builds upon CLS-ER by
 619 incorporating an additional error sensitivity modulation module. The primary goal of ESM-ER is to
 620 mitigate sudden representation drift caused by excessively large loss values during current-domain
 621 learning. Let’s consider $(\mathbf{x}, y) \sim \mathcal{D}_t$, which represents a sample from the current domain batch. In
 622 ESM-ER, the cross-entropy loss value of this sample is evaluated using the stable model $H_{t-1}^{(s)}$ and
 623 can be expressed as:

$$\ell(\mathbf{x}, y) = -\log([H_{t-1}^{(s)}(\mathbf{x})]_y). \quad (50)$$

624 To screen out those samples with a high loss value, ESM-ER assigns each sample a weight λ by
 625 comparing the loss with their expectation value, for which ESM-ER uses a running estimate μ as its
 626 replacement. This can be formulated as follows:

$$\lambda(\mathbf{x}) = \begin{cases} 1 & \text{if } \ell(\mathbf{x}, y) \leq \beta \cdot \mu \\ \frac{\mu}{\ell(\mathbf{x}, y)} & \text{o.w.} \end{cases} \quad (51)$$

627 where β is a hyperparameter that determines the margin for a sample to be classified as low-loss.
 628 For the sake of analysis, we make the following assumptions: (i) $\beta = 1$; (ii) the actual expected loss
 629 value $\mathbb{E}_{\mathbf{x}, y}[\ell(\mathbf{x}, y)]$ is used instead of the running estimate μ ; (iii) a *hard screening mechanism* is
 630 employed instead of the current re-scaling approach. Based on these assumptions, we determine the
 631 sample-wise weights λ^* according to the following rule:

$$\lambda^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \ell(\mathbf{x}, y) \leq \mathbb{E}_{\mathbf{x}, y}[\ell(\mathbf{x}, y)] \\ 0 & \text{o.w.} \end{cases} \quad (52)$$

632 Under the assumption that the loss value $\ell(\mathbf{x}, y)$ follows an exponential distribution, denoted
 633 as $\ell(\mathbf{x}, y) \sim \text{Exp}(\lambda_0)$, where the probability density function is given by $f(\ell(\mathbf{x}, y), \lambda_0) =$
 634 $\lambda_0 e^{-\lambda_0 \ell(\mathbf{x}, y)}$, we can calculate the expectation of the loss as $\mathbb{E}_{\mathbf{x}, y}[\ell(\mathbf{x}, y)] = 1/\lambda_0$. Based on

635 this, we can now determine the expected ratio r of the *unscreened samples* in a mini-batch using the
 636 following equation:

$$r = \int_0^{\frac{1}{\lambda_0}} 1 \cdot \lambda_0 e^{-\lambda_0 \ell(\mathbf{x}, y)} d\ell(\mathbf{x}, y) = \int_0^1 e^{-y} dy = (1 - e^{-1}). \quad (53)$$

637 The ratio r represents the proportion of effective samples in the current-domain batch, as the weights
 638 $\lambda^*(\mathbf{x})$ of the remaining samples are set to 0 due to their high loss value. Consequently, the original
 639 training loss of ESM-ER can be transformed as follows:

$$\mathcal{L}_{\text{ESM-ER}}(h) = r \cdot \widehat{\ell}_{B_t}(h) + \widehat{\ell}_{B'_t}(h) + \lambda \widehat{\ell}_{B'_t}(h, H_{t-1}) \quad (54)$$

$$= r \cdot \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{1}{t-1} \widehat{\ell}_{B'_i}(h) + \sum_{i=1}^{t-1} \frac{\lambda}{t-1} \widehat{\ell}_{B'_i}(h, H_{t-1}). \quad (55)$$

640 After applying the constraint of $\alpha_i + \beta_i + \gamma_i = 1$, we obtain $\lambda = r \cdot (t - 1) - 1$. Substituting this
 641 value back into $\lambda/r(t-1)$, we find that ESM-ER is equivalent to $\{\alpha_i = \lambda/\lambda+1, \beta_i = 0, \gamma_i = 1/\lambda+1\}$ in
 642 UDIL, where $\lambda = r \cdot (t - 1) - 1 = (1 - e^{-1})(t - 1) - 1$ should be set (Row 6 in Table 4).

Algorithm 1 Unified Domain Incremental Learning (UDIL) for Domain t Training

Require: history model $H_{t-1} = P_{t-1} \circ E_{t-1}$, current model $h_\theta = p \circ e$, discriminator model d_ϕ ;

Require: dataset from the current domain \mathcal{S}_t , memory bank $\mathcal{M} = \{M_i\}_{i=1}^{t-1}$;

Require: training steps S , batch size B , learning rate η ;

Require: domain alignment strength coefficient λ_d , hyperparameter for generalization effect C .

- 1: $h_\theta \leftarrow H_{t-1}$ ▷ Initialization of the current model.
 - 2: $\Omega \triangleq \{\alpha_i, \beta_i, \gamma_i\} \leftarrow \{1/3, 1/3, 1/3\}$, for $\forall i \in [t - 1]$ ▷ Initialization of the replay coefficient Ω .
 - 3: **for** $s = 1, \dots, S$ **do**
 - 4: $B_t \sim \mathcal{S}_t; B_i \sim M_i, \forall i \in [t - 1]$ ▷ Sample a mini-batch of data from all domains.
 - 5: $\phi \leftarrow \phi - \eta \cdot \lambda_d \cdot \nabla_\phi V_d(d, e, \overset{\circ}{\Omega})$ ▷ Discriminator training with Eqn. 16.
 - 6: $\Omega \leftarrow \Omega - \eta \cdot \nabla_\Omega V_{0-1}(h, \overset{\circ}{\Omega})$ ▷ Find a tighter bound with Eqn. 15.
 - 7: $\theta \leftarrow \theta - \eta \cdot \nabla_\theta (V_l(h_\theta, \overset{\circ}{\Omega}) - \lambda_d V_d(d, e, \overset{\circ}{\Omega}))$ ▷ Model training with Eqn. 14 and Eqn. 16.
 - 8: **end for**
 - 9: $H_t \leftarrow h$
 - 10: $\mathcal{M} \leftarrow \text{BalancedSampling}(\mathcal{M}, \mathcal{S}_t)$
 - 11: **return** H_t ▷ For training on domain $t + 1$.
-

643 C Implementation Details of UDIL

644 This section delves into the implementation details of the UDIL algorithm. The algorithmic descrip-
 645 tion of UDIL is presented in Alg. 1 and a diagram is presented in Fig. 2. However, there are several
 646 practical issues to be further addressed here, including (i) how to exert the constraints of probability
 647 simplex ($[\alpha_i, \beta_i, \gamma_i] \in \mathbb{S}^2$) and (ii) how the memory is maintained. These two problems will be
 648 addressed in Sec. C.1 and Sec. C.2. Next, Sec. C.3 will cover the evaluation metrics used in this
 649 paper. Finally, Sec. C.4 and Sec. C.5 will present a detailed introduction to the main baselines and
 650 the specific training schemes we follow for empirical evaluation.

651 C.1 Modeling the Replay Coefficients $\Omega = \{\alpha_i, \beta_i, \gamma_i\}$

652 Instead of directly modeling Ω in a way such that it can be updated by gradient descent and satisfies
 653 the constraints that $\alpha_i + \beta_i + \gamma_i = 1$ and $\alpha_i, \beta_i, \gamma_i \geq 0$ at the same time, we use a set of logit variables
 654 $\{\bar{\alpha}_i, \bar{\beta}_i, \bar{\gamma}_i\} \in \mathbb{R}^3$ and the *softmax* function to indirectly calculate Ω during training. Concretely, we
 655 have:

$$\begin{bmatrix} \alpha_i \\ \beta_i \\ \gamma_i \end{bmatrix} = \text{softmax} \left(\begin{bmatrix} \bar{\alpha}_i \\ \bar{\beta}_i \\ \bar{\gamma}_i \end{bmatrix} \right) = \begin{bmatrix} \exp(\bar{\alpha}_i)/Z_i \\ \exp(\bar{\beta}_i)/Z_i \\ \exp(\bar{\gamma}_i)/Z_i \end{bmatrix}, \quad (56)$$

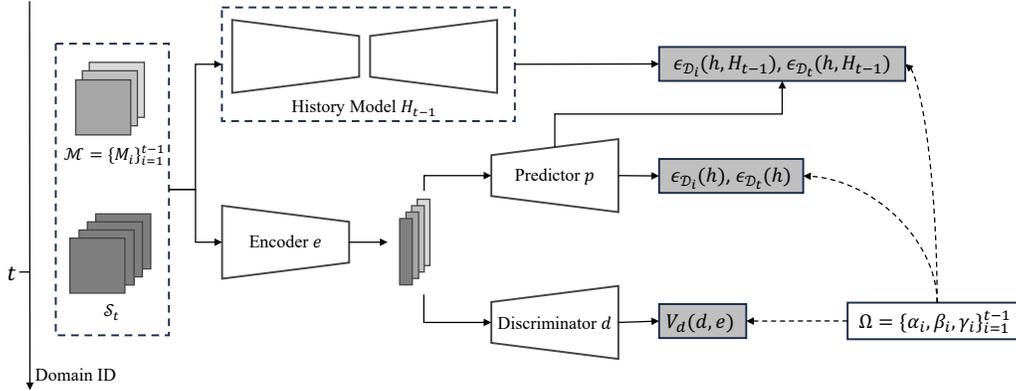


Figure 2: **Diagram of UDIL.** $\mathcal{M} = M_{i=1}^{t-1}$ represents the memory bank that stores all the past exemplars. \mathcal{S}_t corresponds to the dataset from the current domain t , and the current model $h = p \circ e$ is depicted separately in the diagram. The three different categories of losses are illustrated in the dark rectangles, while the weighting effect of the learned replay coefficient $\Omega = \{\alpha_i, \beta_i, \gamma_i\}_{i=1}^{t-1}$ is depicted using dashed lines.

656 where $Z_i = \exp(\bar{\alpha}_i) + \exp(\bar{\beta}_i) + \exp(\bar{\gamma}_i)$ is the normalizing constant. At the beginning of training
 657 on domain t , the logit variables $\{\bar{\alpha}_i, \bar{\beta}_i, \bar{\gamma}_i\} = \{0, 0, 0\}$ are initialized to all zeros, since we do not
 658 have any bias towards any upper bound. During training, they are updated in the same way as the
 659 other parameters with gradient descent.

660 C.2 Memory Maintenance with Balanced Sampling

661 Different from DER++ [6] and its following work [3, 50] that use reservoir sampling [58] to maintain
 662 the episodic memory, UDIL adopts a random balanced sampling after training on each domain. To
 663 be more concrete, given a memory bank with fixed size $|\mathcal{M}|$, after domain t 's training is complete,
 664 we will assign each domain a quota of $|\mathcal{M}|/t$. For the current domain t , we will randomly sample
 665 $\lceil |\mathcal{M}|/t \rceil$ exemplars from its dataset; for all the previous domains $i \in [t-1]$, we will randomly swap
 666 out $\lceil |\mathcal{M}|/t-1 - |\mathcal{M}|/t \rceil$ exemplars from the memory to make sure each domain has roughly the same
 667 number of exemplars. To ensure a fair comparison, we use the same random balanced sampling
 668 strategy for all the other baselines. The following Alg. 2 shows the detailed procedure of random
 balanced sampling.

Algorithm 2 Balanced Sampling for UDIL

Require: memory bank $\mathcal{M} = \{M_i\}_{i=1}^{t-1}$, current domain dataset \mathcal{S}_t , domain ID t .

```

1: for  $i = 1, \dots, t-1$  do
2:   for  $j = 1, \dots, \lceil |\mathcal{M}|/t-1 - |\mathcal{M}|/t \rceil$  do
3:      $(x, y) \leftarrow \text{RandomSample}(M_i)$ 
4:      $(x', y') \leftarrow \text{RandomSample}(\mathcal{S}_t)$ 
5:     Swap  $(x', y')$  into  $\mathcal{M}$ , replacing  $(x, y)$ 
6:   end for
7: end for
8: return  $\mathcal{M}$ 

```

669

670 C.3 Evaluation Metrics

671 In continual learning, many evaluation metrics are based on the **Accuracy Matrix** $R \in \mathbb{R}^{T \times T}$,
 672 where T represents the total number of tasks (domains). In the accuracy matrix R , the entry $R_{i,j}$
 673 corresponds to the accuracy of the model when evaluated on task j after training on task i . With this
 674 definition in mind, we primarily focus on the following specific metrics:

675 **Average Accuracy (Avg. Acc.)** up until domain t represents the average accuracy of the first t
 676 domains after training on these domains. We denote it as A_t and define it as follows:

$$A_t \triangleq \frac{1}{t} \sum_{i=1}^t R_{t,i}. \quad (57)$$

677 In most of the continual learning literature, the final average accuracy A_T is usually reported. In our
 678 paper, this metric is reported in the column labeled “**overall**”. The average accuracy of a model is a
 679 crucial metric as it directly corresponds to the primary optimization goal of minimizing the error on
 680 all domains, as defined in Eqn. 3.

681 Additionally, to better illustrate the learning (and forgetting) process of a model across multiple
 682 domains, we propose the use of the “**Avg. of Avg. Acc.**” metric $A_{t_1:t_2}$, which represents the average
 683 of average accuracies for a consecutive range of domains starting from domain t_1 and ending at
 684 domain t_2 . Specifically, we define this metric as follows:

$$A_{t_1:t_2} \triangleq \frac{1}{t_2-t_1+1} \sum_{i=t_1}^{t_2} A_i. \quad (58)$$

685 This metric provides a condensed representation of the trend in accuracy variation compared to
 686 directly displaying the entire series of average accuracies $\{A_1, A_2, \dots, A_T\}$. We report this Avg. of
 687 Avg. Acc. metric in all tables (except in Table 2 due to the limit of space).

688 **Average Forgetting (i.e., ‘Forgetting’ in the main paper)** defines the average of the largest drop of
 689 accuracy for each domain up till domain t . We denote this metric as F_t and define it as follows:

$$F_t \triangleq \frac{1}{t-1} \sum_{j=1}^{t-1} f_t(j), \quad (59)$$

690 where $f_t(j)$ is the forgetting on domain j after the model completes the training on domain t , which
 691 is defined as:

$$f_t(j) \triangleq \max_{l \in [t-1]} \{R_{l,j} - R_{t,j}\}. \quad (60)$$

692 Typically, the average forgetting is reported after training on the last domain T . Measuring forgetting
 693 is of great practical significance, especially when two models have similar average accuracies. It
 694 indicates how a model balances *stability* and *plasticity*. If a model \mathcal{P} achieves a reasonable final
 695 average accuracy across different domains but exhibits high forgetting, we can conclude that this
 696 model has high plasticity and low stability. It quickly adapts to new domains but at the expense of
 697 performance on past domains. On the other hand, if another model \mathcal{S} has a similar average accuracy
 698 to \mathcal{P} but significantly lower average forgetting, we can infer that the model \mathcal{S} has high stability and
 699 low plasticity. It sacrifices performance on recent domains to maintain a reasonable performance
 700 on past domains. Hence, to gain a comprehensive understanding of model performance, we focus
 701 on evaluating two key metrics: *Avg. Acc.* and *Forgetting*. These metrics provide insights into
 702 how models balance stability and plasticity and allow us to assess their overall performance across
 703 different domains.

704 **Forward Transfer** W_t quantifies the extent to which learning from past $t - 1$ domains contributes to
 705 the performance on the next domain t . It is defined as follows:

$$W_t \triangleq \frac{1}{t-1} \sum_{i=2}^t R_{i-1,i} - r_i, \quad (61)$$

706 where r_i is the accuracy of a randomly initialized model evaluated on domain i . For domain
 707 incremental learning, where the model does not have access to future domain data and does not
 708 explicitly optimize for higher Forward Transfer, the results of this metric are *typically random*.
 709 *Therefore, we do not report this metric in the complete tables presented in this section.*

710 C.4 Introduction to Baselines

711 We compare UDIL with the state-of-the-art continual learning methods that are either specifically
 712 designed for domain incremental learning or can be easily adapted to the domain incremental learning

713 setting. Exemplar-free baselines include online Elastic Weight Consolidation (**oEWC**) [51], Synaptic
 714 Intelligence (**SI**) [66], and Learning without Forgetting (**LwF**) [29]. Memory-based domain incremen-
 715 tal learning baselines include Gradient Episodic Memory (**GEM**) [34], Averaged Gradient Episodic
 716 Memory (**A-GEM**) [8], Experience Replay (**ER**) [46], Dark Experience Replay (**DER++**) [6], and
 717 two recent methods, Complementary Learning System based Experience Replay (**CLS-ER**) [3] and
 718 Error Sensitivity Modulation based Experience Replay (**ESM-ER**) [50]. In addition, we implement
 719 the fine-tuning (**Fine-tune**) [29] and joint-training (**Joint**) as the performance lower bound and upper
 720 bound (Oracle). Here we provide a short description of the primary idea of the memory-based domain
 721 incremental learning baselines.

- 722 • **GEM** [34]: The baseline method that uses the memory to provide additional optimization
 723 constraints during learning the current domain. Specifically, the update of the model cannot
 724 point towards the direction at which the loss of any exemplar increases.
- 725 • **A-GEM** [8]: The improved baseline method where the constraints of GEM are averaged as
 726 one, which shortens the computational time significantly.
- 727 • **ER** [46]: The fundamental memory-based domain incremental learning framework where
 728 the mini-batch of the memory is regularly replayed with the current domain data.
- 729 • **DER++** [6]: A simple yet effective replay-based method where an additional logits distilla-
 730 tion (dubbed “dark experience replay”) is applied compared to the vanilla ER.
- 731 • **CLS-ER** [3]: A complementary learning system inspired replay method, where two expo-
 732 nential moving average models are used to serve as the semantic memory, which provides
 733 the logits distillation target during training.
- 734 • **ESM-ER** [50]: An improved version of CLS-ER, where the effect of large errors when
 735 learning the current domain is reduced, dubbed “error sensitivity modulation”.

736 C.5 Training Schemes

737 **Training Process.** For each group of experiments, we run three rounds with different seeds and
 738 report the mean and standard deviation of the results. We follow the optimal configurations (epochs
 739 and learning rate) stated in [6, 50] for the baselines in *P-MNIST* and *R-MNIST* dataset. For *HD-Balls*
 740 and *Seq-CORe50*, we first search for the optimal training configuration for the joint learning, and then
 741 grid-search the configuration in a small range near it for the baselines listed above. For our UDIL
 742 framework, as it involves adversarial training for the domain embedding alignment, we typically need
 743 a configuration that has larger number of epochs and smaller learning rate. We use a simple grid
 744 search to achieve the optimal configuration for it as well.

745 **Model Architectures.** For the baseline methods and UDIL in the same dataset, we adopt the same
 746 backbone neural architectures to ensure fair comparison. In *HD-Balls*, we adopt the same multi-layer
 747 perceptron with the same separation of encoder and decoder as in CIDA [59], where the hidden
 748 dimension is set to 800. In *P-MNIST* and *R-MNIST*, we adopt the same multi-layer perceptron
 749 architecture as in DER++ [6] with hidden dimension set to 800 as well. In *Seq-CORe50*, we use the
 750 ResNet18 [19] as our backbone architecture for all the methods, where the layers before the final
 751 average pooling are treated as the encoder e , and the remaining part is treated as the predictor p .

752 **Hyperparameter Setting.** For setting the hyper-parameter embedding alignment strength coeffi-
 753 cient λ_d and parameter C that models the combined effect of VC-dimension d and error tolerance δ ,
 754 we use grid search for each dataset, where the range $\lambda_d \in [0.01, 100]$ and $C \in [0, 1000]$ are used.

755 D Additional Empirical Results

756 This section presents additional empirical results of the UDIL algorithm. Sec. D.1 will show the
 757 additional results on different constraints with varying memory sizes. Sec. D.2 provides additional
 758 qualitative results: visualization of embedding distributions, to showcase the importance of the
 759 embedding alignment across domains.

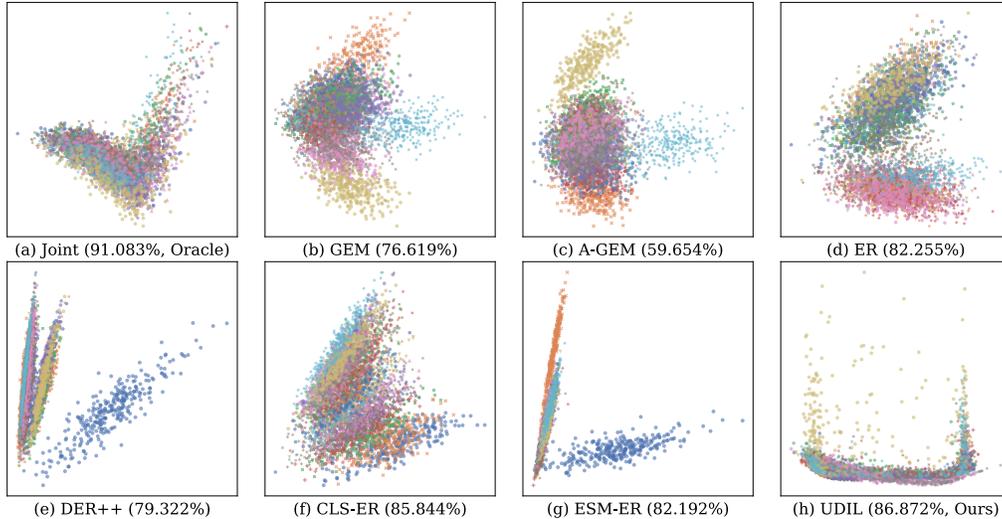


Figure 3: More Results on *HD-Balls*. Data is colored according to domain ID. All data is plotted after PCA [5]. **(a-h)** Accuracy and embeddings learned by Joint (oracle), UDIL, and six baselines with memory. Joint, as the *oracle*, naturally aligns different domains, and UDIL outperforms all baselines in terms of embedding alignment and accuracy.

760 D.1 Empirical Results on Varying Memory Sizes

761 Here we present additional empirical results to validate the effectiveness of our UDIL framework
 762 using varying memory sizes. The evaluation is conducted on three real-world datasets, as shown
 763 in Table 5, Table 6, and Table 7. By increasing the memory size from 400 to 800 in Table 5 and 6
 764 and from 500 to 1000 in Table 7, we can investigate the impact of having access to a larger pool of
 765 past experiences on the continual learning performance, which might occur when the constraint on
 766 memory capacity is relaxed. This allows us to study the benefits of a more extensive memory in terms
 767 of knowledge retention and performance improvement. On the other hand, by further decreasing the
 768 memory size to the extreme of 200 in Table 5 and 6, we can explore the consequences of severely
 769 limited memory capacity. This scenario simulates situations where memory constraints are extremely
 770 tight, and the model can only retain a small fraction of past domain data, for example, a model
 771 deployed on edge devices. To ensure a fair comparison, here we use the same best configuration
 772 found in the main body of this work.

773 The results in all three tables demonstrate a clear advantage of our UDIL framework when the
 774 memory size is limited. In *P-MNIST* and *R-MNIST*, when the memory size $|\mathcal{M}| = 200$, the overall
 775 performance of UDIL reaches 91.483% and 82.796% respectively, which outperforms the second
 776 best model DER++ by 0.757% and a remarkable 6.125%. In *Seq-CORE50*, when the memory
 777 size $|\mathcal{M}| = 500$ is set, UDIL holds a 3.474% lead compared to the second best result. When the
 778 memory size is larger, the gap between UDIL and the baseline models is smaller. This is because
 779 when the memory constraint is relaxed, all the continual learning models should be at least closer
 780 to the performance upper bound, i.e., joint learning or ‘Joint (Oracle)’ in the tables, causing the
 781 indistinguishable results among each other. Apparently, DER++ favors larger memory more than
 782 UDIL, while UDIL can still maintain a narrow lead in the large scale dataset *Seq-CORE50*.

783 D.2 Visualization of Embedding Spaces

784 Here we provide more embedding space visualization results for the baselines with the utilization
 785 of memory, shown in Fig. 3. As one of the primary objectives of our algorithm, embedding space
 786 alignment across multiple domains naturally follows the pattern shown in the joint learning and
 787 therefore leads to a higher performance.

Table 5: **Performances (%) evaluated on *P*-MNIST.** Average Accuracy (Avg. Acc.) and Forgetting are reported to measure the methods’ performance. “ \uparrow ” and “ \downarrow ” mean higher and lower numbers are better, respectively. We use **boldface** and underlining to denote the best and the second-best performance, respectively. We use “-” to denote “not applicable”.

Method	Buffer	$\mathcal{D}_{1:5}$	$\mathcal{D}_{6:10}$	$\mathcal{D}_{11:15}$	$\mathcal{D}_{16:20}$	Overall	
		Avg. Acc (\uparrow)				Avg. Acc (\uparrow)	Forgetting (\downarrow)
Fine-tune	-	92.506 \pm 2.062	87.088 \pm 1.337	81.295 \pm 2.372	72.807 \pm 1.817	70.102 \pm 2.945	27.522 \pm 3.042
oEWC [51]	-	92.415 \pm 0.816	87.988 \pm 1.607	83.098 \pm 1.843	78.670 \pm 0.902	78.476 \pm 1.223	18.068 \pm 1.321
SI [66]	-	92.282 \pm 0.862	87.986 \pm 1.622	83.698 \pm 1.220	79.669 \pm 0.709	79.045 \pm 1.357	17.409 \pm 1.446
LwF [29]	-	95.025 \pm 0.487	91.402 \pm 1.546	83.984 \pm 2.103	76.046 \pm 2.004	73.545 \pm 2.646	24.556 \pm 2.789
GEM [34]	200	93.310 \pm 0.374	91.900 \pm 0.456	89.813 \pm 0.914	87.251 \pm 0.524	86.729 \pm 0.203	9.430 \pm 0.156
A-GEM [8]		93.326 \pm 0.363	91.466 \pm 0.605	89.048 \pm 1.005	86.518 \pm 0.604	85.712 \pm 0.228	10.485 \pm 0.196
ER [46]		94.087 \pm 0.762	92.397 \pm 0.464	89.999 \pm 1.060	87.492 \pm 0.448	86.963 \pm 0.303	9.273 \pm 0.255
DER++ [6]		94.708 \pm 0.451	<u>94.582\pm0.158</u>	<u>93.271\pm0.585</u>	90.980 \pm 0.610	90.333 \pm 0.587	6.110 \pm 0.545
CLS-ER [3]		94.761 \pm 0.340	93.943 \pm 0.197	<u>92.725\pm0.566</u>	<u>91.150\pm0.357</u>	<u>90.726\pm0.218</u>	<u>5.428\pm0.252</u>
ESM-ER [50]		<u>95.198\pm0.236</u>	94.029 \pm 0.427	91.710 \pm 1.056	88.181 \pm 1.021	86.851 \pm 0.858	10.007 \pm 0.864
UDIL (Ours)		95.747\pm0.486	94.695\pm0.256	93.756\pm0.343	92.254\pm0.564	91.483\pm0.270	4.399\pm0.314
GEM [34]		400	93.557 \pm 0.225	92.635 \pm 0.306	91.246 \pm 0.492	89.565 \pm 0.342	89.097 \pm 0.149
A-GEM [8]	93.432 \pm 0.333		92.064 \pm 0.439	90.038 \pm 0.726	87.988 \pm 0.335	87.560 \pm 0.087	8.577 \pm 0.053
ER [46]	93.525 \pm 1.101		91.649 \pm 0.362	90.426 \pm 0.456	88.728 \pm 0.353	88.339 \pm 0.044	7.180 \pm 0.029
DER++ [6]	94.952 \pm 0.403		95.089\pm0.075	94.458\pm0.328	93.257\pm0.249	92.950\pm0.361	<u>3.378\pm0.245</u>
CLS-ER [3]	94.262 \pm 0.649		93.195 \pm 0.148	92.623 \pm 0.195	91.839 \pm 0.187	91.598 \pm 0.117	3.795 \pm 0.144
ESM-ER [50]	<u>95.413\pm0.139</u>		94.654 \pm 0.314	93.353 \pm 0.588	91.022 \pm 0.781	89.829 \pm 0.698	6.888 \pm 0.738
UDIL (Ours)	95.992\pm0.349		<u>95.026\pm0.250</u>	<u>94.212\pm0.280</u>	<u>93.094\pm0.326</u>	<u>92.666\pm0.108</u>	2.853\pm0.107
GEM [34]	800		93.717 \pm 0.177	93.116 \pm 0.206	92.166 \pm 0.335	91.076 \pm 0.342	90.609 \pm 0.364
A-GEM [8]		93.612 \pm 0.241	92.523 \pm 0.375	90.718 \pm 0.739	88.543 \pm 0.391	88.020 \pm 0.851	8.081 \pm 0.867
ER [46]		93.827 \pm 0.871	92.457 \pm 0.217	91.688 \pm 0.277	90.617 \pm 0.289	90.252 \pm 0.056	5.188 \pm 0.045
DER++ [6]		95.295 \pm 0.317	95.539\pm0.041	95.099\pm0.187	94.423\pm0.151	94.227\pm0.261	<u>2.106\pm0.161</u>
CLS-ER [3]		94.463 \pm 0.537	93.567 \pm 0.093	93.182 \pm 0.137	92.744 \pm 0.112	92.578 \pm 0.152	2.803 \pm 0.183
ESM-ER [50]		<u>95.567\pm0.150</u>	95.136 \pm 0.202	94.301 \pm 0.347	92.981 \pm 0.397	92.408 \pm 0.387	4.170 \pm 0.357
UDIL (Ours)		96.082\pm0.313	<u>95.207\pm0.196</u>	<u>94.642\pm0.156</u>	<u>93.997\pm0.194</u>	<u>93.724\pm0.043</u>	1.633\pm0.035
Joint (Oracle)		∞	-	-	-	-	96.368 \pm 0.042

Table 6: **Performances (%) evaluated on *R-MNIST***. Average Accuracy (Avg. Acc.) and Forgetting are reported to measure the methods’ performance. “ \uparrow ” and “ \downarrow ” mean higher and lower numbers are better, respectively. We use **boldface** and underlining to denote the best and the second-best performance, respectively. We use “-” to denote “not applicable”.

Method	Buffer	$\mathcal{D}_{1:5}$	$\mathcal{D}_{6:10}$	$\mathcal{D}_{11:15}$	$\mathcal{D}_{16:20}$	Overall	
		Avg. Acc (\uparrow)				Avg. Acc (\uparrow)	Forgetting (\downarrow)
Fine-tune	-	92.961 \pm 2.683	76.617 \pm 8.011	60.212 \pm 3.688	49.793 \pm 1.552	47.803 \pm 1.703	52.281 \pm 1.797
oEWC [51]	-	91.765 \pm 2.286	76.226 \pm 7.622	60.320 \pm 3.892	50.505 \pm 1.772	48.203 \pm 0.827	51.181 \pm 0.867
SI [66]	-	91.867 \pm 2.272	76.801 \pm 7.391	60.956 \pm 3.504	50.301 \pm 1.538	48.251 \pm 1.381	51.053 \pm 1.507
LwF [29]	-	95.174 \pm 1.154	83.044 \pm 5.935	65.899 \pm 4.061	55.980 \pm 1.296	54.709 \pm 0.515	45.473 \pm 0.565
GEM [34]	200	93.441 \pm 0.610	88.620 \pm 2.381	81.034 \pm 2.704	73.112 \pm 1.922	70.545 \pm 0.623	27.684 \pm 0.645
A-GEM [8]		92.667 \pm 1.352	82.772 \pm 5.503	70.579 \pm 4.028	60.462 \pm 2.001	57.958 \pm 0.579	40.969 \pm 0.580
ER [46]		94.705 \pm 0.790	89.171 \pm 2.883	79.962 \pm 3.365	71.787 \pm 1.608	69.627 \pm 0.911	28.749 \pm 0.993
DER++ [6]		94.904 \pm 0.414	91.637 \pm 1.871	84.915 \pm 2.315	78.373 \pm 1.244	76.671 \pm 0.391	21.743 \pm 0.409
CLS-ER [3]		95.131 \pm 0.523	91.421 \pm 1.732	84.773 \pm 2.665	77.733 \pm 1.480	75.609 \pm 0.418	22.483 \pm 0.456
ESM-ER [50]		95.378\pm0.531	90.800 \pm 2.528	83.438 \pm 2.581	76.987 \pm 1.219	75.203 \pm 0.143	23.564 \pm 0.157
UDIL (Ours)		95.097 \pm 0.447	93.101\pm1.305	89.194\pm1.472	84.704\pm1.722	82.796\pm1.882	12.971\pm2.389
GEM [34]		400	93.842 \pm 0.313	90.663 \pm 1.856	85.392 \pm 1.856	79.061 \pm 1.578	76.619 \pm 0.581
A-GEM [8]	92.820 \pm 1.274		83.564 \pm 5.024	72.616 \pm 3.865	62.223 \pm 2.081	59.654 \pm 0.122	39.196 \pm 0.171
ER [46]	94.916 \pm 0.457		91.491 \pm 1.878	86.029 \pm 2.176	78.688 \pm 1.323	76.794 \pm 0.696	20.696 \pm 0.744
DER++ [6]	95.246 \pm 0.228		93.627 \pm 1.147	90.011 \pm 1.289	85.601 \pm 0.982	84.258 \pm 0.544	13.692 \pm 0.560
CLS-ER [3]	95.233 \pm 0.271		92.740 \pm 1.268	89.111 \pm 1.305	83.678 \pm 1.388	81.771 \pm 0.354	15.455 \pm 0.356
ESM-ER [50]	95.825\pm0.303		93.378 \pm 1.480	89.290 \pm 1.604	83.868 \pm 1.163	82.192 \pm 0.164	16.195 \pm 0.150
UDIL (Ours)	95.274 \pm 0.469		94.043\pm0.759	91.511\pm0.990	87.809\pm0.849	86.635\pm0.686	8.506\pm1.181
GEM [34]	800		94.212 \pm 0.322	92.482 \pm 1.125	89.191 \pm 1.346	84.866 \pm 1.317	82.772 \pm 1.079
A-GEM [8]		92.902 \pm 1.194	84.611 \pm 4.451	75.150 \pm 3.421	64.510 \pm 2.437	61.240 \pm 1.026	37.528 \pm 1.089
ER [46]		95.144 \pm 0.281	92.997 \pm 1.195	89.319 \pm 1.365	84.352 \pm 1.681	81.877 \pm 1.157	15.285 \pm 1.196
DER++ [6]		95.496 \pm 0.261	94.960\pm0.568	93.013\pm0.689	90.820\pm0.687	89.746\pm0.356	7.821 \pm 0.371
CLS-ER [3]		95.462 \pm 0.174	93.927 \pm 0.881	91.275 \pm 0.930	87.816 \pm 0.988	86.418 \pm 0.215	10.598 \pm 0.228
ESM-ER [50]		96.086\pm0.361	94.746 \pm 0.915	92.393 \pm 0.974	89.745 \pm 0.712	88.662 \pm 0.263	9.409 \pm 0.255
UDIL (Ours)		95.354 \pm 0.480	94.711 \pm 0.563	92.776 \pm 0.695	90.399 \pm 0.755	89.191 \pm 0.685	6.351\pm1.304
Joint (Oracle)		∞	-	-	-	-	97.150 \pm 0.036

Table 7: **Performances (%) evaluated on Seq-CORE50.** Avg. Acc. and Forgetting are reported to measure the methods’ performance. “ \uparrow ” and “ \downarrow ” mean higher and lower numbers are better, respectively. We use **boldface** and underlining to denote the best and the second-best performance, respectively. We use “-” to denote “not applicable” and “*” to denote out-of-memory (*OOM*) error when running the experiments.

Method	Buffer	$\mathcal{D}_{1:3}$	$\mathcal{D}_{4:6}$	$\mathcal{D}_{7:9}$	$\mathcal{D}_{10:11}$	Overall	
		Avg. Acc (\uparrow)				Avg. Acc (\uparrow)	Forgetting (\downarrow)
Fine-tune	-	73.707 \pm 13.144	34.551 \pm 1.254	29.406 \pm 2.579	28.689 \pm 3.144	31.832 \pm 1.034	73.296 \pm 1.399
oEWC [51]	-	74.567 \pm 13.360	35.915 \pm 0.260	30.174 \pm 3.195	28.291 \pm 2.522	30.813 \pm 1.154	74.563 \pm 0.937
SI [66]	-	74.661 \pm 14.162	34.345 \pm 1.001	30.127 \pm 2.971	28.839 \pm 3.631	32.469 \pm 1.315	73.144 \pm 1.588
LwF [29]	-	80.383 \pm 10.190	28.357 \pm 1.143	31.386 \pm 0.787	28.711 \pm 2.981	31.692 \pm 0.768	72.990 \pm 1.350
GEM [34]	500	79.852 \pm 6.864	38.961 \pm 1.718	39.258 \pm 2.614	36.859 \pm 0.842	37.701 \pm 0.273	22.724 \pm 1.554
A-GEM [8]		80.348 \pm 9.394	41.472 \pm 3.394	43.213 \pm 1.542	39.181 \pm 3.999	43.181 \pm 2.025	33.775 \pm 3.003
ER [46]		90.838 \pm 2.177	79.343 \pm 2.699	68.151 \pm 0.226	65.034 \pm 1.571	66.605 \pm 0.214	32.750 \pm 0.455
DER++ [6]		92.444 \pm 1.764	88.652 \pm 1.854	80.391 \pm 0.107	78.038 \pm 0.591	78.629 \pm 0.753	21.910 \pm 1.094
CLS-ER [3]		89.834 \pm 1.323	78.909 \pm 1.724	70.591 \pm 0.322	*	*	*
ESM-ER [50]		84.905 \pm 6.471	51.905 \pm 3.257	53.815 \pm 1.770	50.178 \pm 2.574	52.751 \pm 1.296	25.444 \pm 0.580
UDIL (Ours)		98.152\pm1.665	89.814\pm2.302	83.052\pm0.151	81.547\pm0.269	82.103\pm0.279	19.589\pm0.303
GEM [34]	1000	78.717 \pm 4.831	43.269 \pm 3.419	40.908 \pm 2.200	40.408 \pm 1.168	41.576 \pm 1.599	18.537 \pm 1.237
A-GEM [8]		78.917 \pm 8.984	41.172 \pm 4.293	44.576 \pm 1.701	38.960 \pm 3.867	42.827 \pm 1.659	33.800 \pm 1.847
ER [46]		90.048 \pm 2.699	84.668 \pm 1.988	77.561 \pm 1.281	72.268 \pm 0.720	72.988 \pm 0.566	25.997 \pm 0.694
DER++ [6]		89.510 \pm 5.726	92.492 \pm 0.902	88.883 \pm 0.794	86.108 \pm 0.284	86.392 \pm 0.714	13.128 \pm 0.474
CLS-ER [3]		92.004 \pm 0.894	85.044 \pm 1.276	*	*	*	*
ESM-ER [50]		85.120 \pm 4.339	54.852 \pm 5.511	61.714 \pm 1.840	55.098 \pm 3.834	58.932 \pm 0.959	20.134 \pm 0.643
UDIL (Ours)		98.648\pm1.174	93.447\pm1.111	90.545\pm0.705	87.923\pm0.232	88.155\pm0.445	12.882\pm0.460
Joint (Oracle)	∞	-	-	-	-	99.137 \pm 0.049	-

788 **References**

- 789 [1] R. Aljundi, F. Babiloni, M. Elhoseiny, M. Rohrbach, and T. Tuytelaars. Memory aware synapses: Learning
790 what (not) to forget. In *Proceedings of the European conference on computer vision (ECCV)*, pages
791 139–154, 2018.
- 792 [2] M. Anthony, P. L. Bartlett, P. L. Bartlett, et al. *Neural network learning: Theoretical foundations*, volume 9.
793 cambridge university press Cambridge, 1999.
- 794 [3] E. Arani, F. Sarfraz, and B. Zonooz. Learning fast, learning slow: A general continual learning method
795 based on complementary learning system. *arXiv preprint arXiv:2201.12604*, 2022.
- 796 [4] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning
797 from different domains. *Machine learning*, 79:151–175, 2010.
- 798 [5] C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- 799 [6] P. Buzzega, M. Boschini, A. Porrello, D. Abati, and S. Calderara. Dark experience for general continual
800 learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930,
801 2020.
- 802 [7] H. Cha, J. Lee, and J. Shin. Co2l: Contrastive continual learning. In *Proceedings of the IEEE/CVF*
803 *International conference on computer vision*, pages 9516–9525, 2021.
- 804 [8] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny. Efficient lifelong learning with a-gem. *arXiv*
805 *preprint arXiv:1812.00420*, 2018.
- 806 [9] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. Torr, and M. Ranzato. On tiny
807 episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.
- 808 [10] Z. Chen, J. Zhuang, X. Liang, and L. Lin. Blending-target domain adaptation by adversarial meta-adaptation
809 networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
810 pages 2248–2257, 2019.
- 811 [11] S. Dai, K. Sohn, Y.-H. Tsai, L. Carin, and M. Chandraker. Adaptation across extreme variations using
812 unlabeled domain bridges. *arXiv preprint arXiv:1906.02238*, 2019.
- 813 [12] D. Deng, G. Chen, J. Hao, Q. Wang, and P.-A. Heng. Flattening sharpness for dynamic gradient projection
814 memory benefits continual learning. *Advances in Neural Information Processing Systems*, 34:18710–18721,
815 2021.
- 816 [13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers
817 for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 818 [14] P. Dhar, R. V. Singh, K.-C. Peng, Z. Wu, and R. Chellappa. Learning without memorizing. In *Proceedings*
819 *of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5138–5146, 2019.
- 820 [15] J. Gallardo, T. L. Hayes, and C. Kanan. Self-supervised training enhances online continual learning. *arXiv*
821 *preprint arXiv:2103.14010*, 2021.
- 822 [16] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky.
823 Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–
824 2030, 2016.
- 825 [17] P. Garg, R. Saluja, V. N. Balasubramanian, C. Arora, A. Subramanian, and C. Jawahar. Multi-domain
826 incremental learning for semantic segmentation. In *Proceedings of the IEEE/CVF Winter Conference on*
827 *Applications of Computer Vision*, pages 761–771, 2022.
- 828 [18] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.
829 Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- 830 [19] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the*
831 *IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 832 [20] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint*
833 *arXiv:1503.02531*, 2015.
- 834 [21] G. E. Hinton, S. Osindero, and Y.-W. Teh. A fast learning algorithm for deep belief nets. *Neural*
835 *computation*, 18(7):1527–1554, 2006.

- 836 [22] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- 837 [23] W. Hoeffding. Probability inequalities for sums of bounded random variables. *The collected works of*
838 *Wassily Hoeffding*, pages 409–426, 1994.
- 839 [24] C.-Y. Hung, C.-H. Tu, C.-E. Wu, C.-H. Chen, Y.-M. Chan, and C.-S. Chen. Compacting, picking and
840 growing for unforgetting continual learning. *Advances in Neural Information Processing Systems*, 32,
841 2019.
- 842 [25] T. Kalb, M. Roschani, M. Ruf, and J. Beyerer. Continual learning for class-and domain-incremental
843 semantic segmentation. In *2021 IEEE Intelligent Vehicles Symposium (IV)*, pages 1345–1351. IEEE, 2021.
- 844 [26] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ra-
845 malho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings*
846 *of the national academy of sciences*, 114(13):3521–3526, 2017.
- 847 [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural
848 networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 849 [28] Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available:
850 <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 851 [29] Z. Li and D. Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine*
852 *intelligence*, 40(12):2935–2947, 2017.
- 853 [30] Z. Li, L. Zhao, Z. Zhang, H. Zhang, D. Liu, T. Liu, and D. N. Metaxas. Steering prototype with prompt-
854 tuning for rehearsal-free continual learning. *arXiv preprint arXiv:2303.09447*, 2023.
- 855 [31] V. Lomonaco and D. Maltoni. Core50: a new dataset and benchmark for continuous object recognition. In
856 S. Levine, V. Vanhoucke, and K. Goldberg, editors, *Proceedings of the 1st Annual Conference on Robot*
857 *Learning*, volume 78 of *Proceedings of Machine Learning Research*, pages 17–26. PMLR, 13–15 Nov
858 2017.
- 859 [32] V. Lomonaco, D. Maltoni, and L. Pellegrini. Rehearsal-free continual learning over small non-iid batches.
860 In *CVPR Workshops*, volume 1, page 3, 2020.
- 861 [33] M. Long, Z. CAO, J. Wang, and M. I. Jordan. Conditional adversarial domain adaptation. In S. Bengio,
862 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural*
863 *Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- 864 [34] D. Lopez-Paz and M. Ranzato. Gradient episodic memory for continual learning. *Advances in neural*
865 *information processing systems*, 30, 2017.
- 866 [35] M. J. Mirza, M. Masana, H. Possegger, and H. Bischof. An efficient domain-incremental learning approach
867 to drive in all weather conditions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
868 *Pattern Recognition*, pages 3001–3011, 2022.
- 869 [36] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- 870 [37] L. T. Nguyen-Meidine, A. Belal, M. Kiran, J. Dolz, L.-A. Blais-Morin, and E. Granger. Unsupervised
871 multi-target domain adaptation through knowledge distillation. In *Proceedings of the IEEE/CVF Winter*
872 *Conference on Applications of Computer Vision*, pages 1339–1347, 2021.
- 873 [38] Z. Ni, H. Shi, S. Tang, L. Wei, Q. Tian, and Y. Zhuang. Revisiting catastrophic forgetting in class
874 incremental learning. *arXiv preprint arXiv:2107.12308*, 2021.
- 875 [39] Z. Ni, L. Wei, S. Tang, Y. Zhuang, and Q. Tian. Continual vision-language representation learning with
876 off-diagonal information, 2023.
- 877 [40] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang. Domain adaptation via transfer component analysis. *IEEE*
878 *transactions on neural networks*, 22(2):199–210, 2010.
- 879 [41] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data*
880 *engineering*, 22(10):1345–1359, 2010.
- 881 [42] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang. Moment matching for multi-source
882 domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages
883 1406–1415, 2019.

- 884 [43] Q. Pham, C. Liu, and S. Hoi. Dualnet: Continual learning, fast and slow. *Advances in Neural Information*
885 *Processing Systems*, 34:16131–16144, 2021.
- 886 [44] R. Ramesh and P. Chaudhari. Model zoo: A growing" brain" that learns continually. *arXiv preprint*
887 *arXiv:2106.03027*, 2021.
- 888 [45] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert. icarl: Incremental classifier and representation
889 learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages
890 2001–2010, 2017.
- 891 [46] M. Riemer, I. Cases, R. Ajemian, M. Liu, I. Rish, Y. Tu, and G. Tesauro. Learning to learn without
892 forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- 893 [47] G. Saha, I. Garg, and K. Roy. Gradient projection memory for continual learning. *arXiv preprint*
894 *arXiv:2103.09762*, 2021.
- 895 [48] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada. Maximum classifier discrepancy for unsupervised
896 domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
897 pages 3723–3732, 2018.
- 898 [49] S. Sankaranarayanan, Y. Balaji, C. D. Castillo, and R. Chellappa. Generate to adapt: Aligning domains
899 using generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and*
900 *pattern recognition*, pages 8503–8512, 2018.
- 901 [50] F. Sarfraz, E. Arani, and B. Zonooz. Error sensitivity modulation based experience replay: Mitigating
902 abrupt representation drift in continual learning. *arXiv preprint arXiv:2302.11344*, 2023.
- 903 [51] J. Schwarz, W. Czarnecki, J. Luketina, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and R. Hadsell.
904 Progress & compress: A scalable framework for continual learning. In *International conference on machine*
905 *learning*, pages 4528–4537. PMLR, 2018.
- 906 [52] J. Serra, D. Suris, M. Miron, and A. Karatzoglou. Overcoming catastrophic forgetting with hard attention
907 to the task. In *International Conference on Machine Learning*, pages 4548–4557. PMLR, 2018.
- 908 [53] B. Sun and K. Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *Computer*
909 *Vision–ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8–10 and 15–16, 2016, Proceedings,*
910 *Part III 14*, pages 443–450. Springer, 2016.
- 911 [54] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for
912 domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- 913 [55] G. M. van de Ven, T. Tuytelaars, and A. S. Tolias. Three types of incremental learning. *Nature Machine*
914 *Intelligence*, pages 1–13, 2022.
- 915 [56] V. N. Vapnik and A. Y. Chervonenkis. On the uniform convergence of relative frequencies of events to
916 their probabilities. *Measures of complexity: festschrift for alexey chervonenkis*, pages 11–30, 2015.
- 917 [57] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin.
918 Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- 919 [58] J. S. Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*,
920 11(1):37–57, 1985.
- 921 [59] H. Wang, H. He, and D. Katabi. Continuously indexed domain adaptation. *arXiv preprint arXiv:2007.01807*,
922 2020.
- 923 [60] L. Wang, X. Zhang, Q. Li, J. Zhu, and Y. Zhong. Coscl: Cooperation of small continual learners is stronger
924 than a big one. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October*
925 *23–27, 2022, Proceedings, Part XXVI*, pages 254–271. Springer, 2022.
- 926 [61] Y. Wang, Z. Huang, and X. Hong. S-prompts learning with pre-trained transformers: An occam’s razor for
927 domain incremental learning. *arXiv preprint arXiv:2207.12819*, 2022.
- 928 [62] Y. Wu, Y. Chen, L. Wang, Y. Ye, Z. Liu, Y. Guo, and Y. Fu. Large scale incremental learning. In
929 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 374–382,
930 2019.
- 931 [63] Z. Xu, G.-Y. Hao, H. He, and H. Wang. Domain-indexing variational bayes: Interpretable domain index
932 for domain adaptation. In *International Conference on Learning Representations*, 2023.

- 933 [64] Z. Xu, G.-H. Lee, Y. Wang, H. Wang, et al. Graph-relational domain adaptation. *arXiv preprint*
934 *arXiv:2202.03628*, 2022.
- 935 [65] J. Yoon, E. Yang, J. Lee, and S. J. Hwang. Lifelong learning with dynamically expandable networks. *arXiv*
936 *preprint arXiv:1708.01547*, 2017.
- 937 [66] F. Zenke, B. Poole, and S. Ganguli. Continual learning through synaptic intelligence. In *International*
938 *conference on machine learning*, pages 3987–3995. PMLR, 2017.
- 939 [67] Y. Zhang, T. Liu, M. Long, and M. Jordan. Bridging theory and algorithm for domain adaptation. In
940 *International conference on machine learning*, pages 7404–7413. PMLR, 2019.
- 941 [68] M. Zhao, S. Yue, D. Katabi, T. S. Jaakkola, and M. T. Bianchi. Learning sleep stages from radio signals: A
942 conditional adversarial architecture. In *International Conference on Machine Learning*, pages 4100–4109.
943 PMLR, 2017.