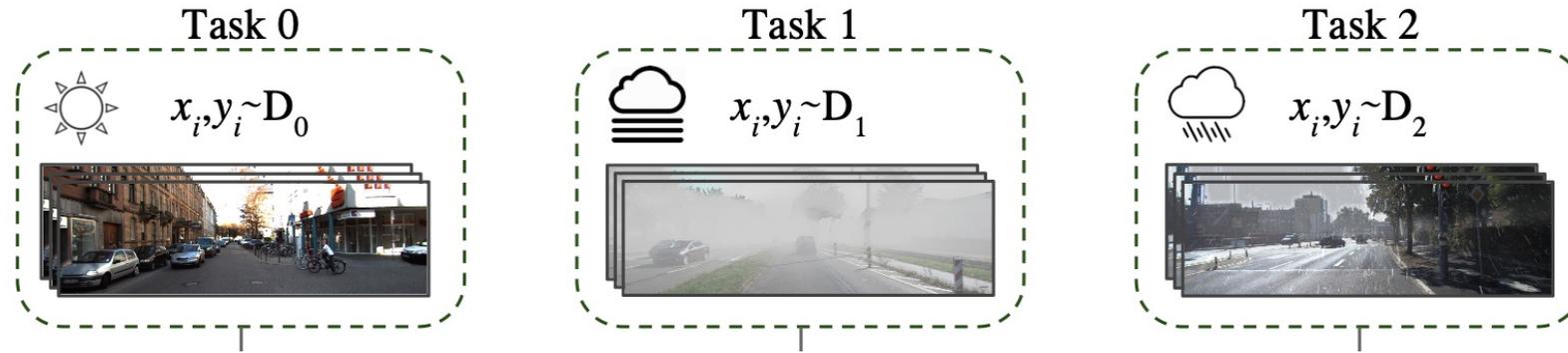


A Unified Approach to Domain Incremental Learning with Memory: Theory and Algorithm

Haizhou Shi, Hao Wang

Computer Science Department, Rutgers University

- Domain Incremental Learning (DIL)
 - Machine learning models are required to incrementally learn the evolving data distributions.
 - E.g., autonomous driving under different weather conditions.

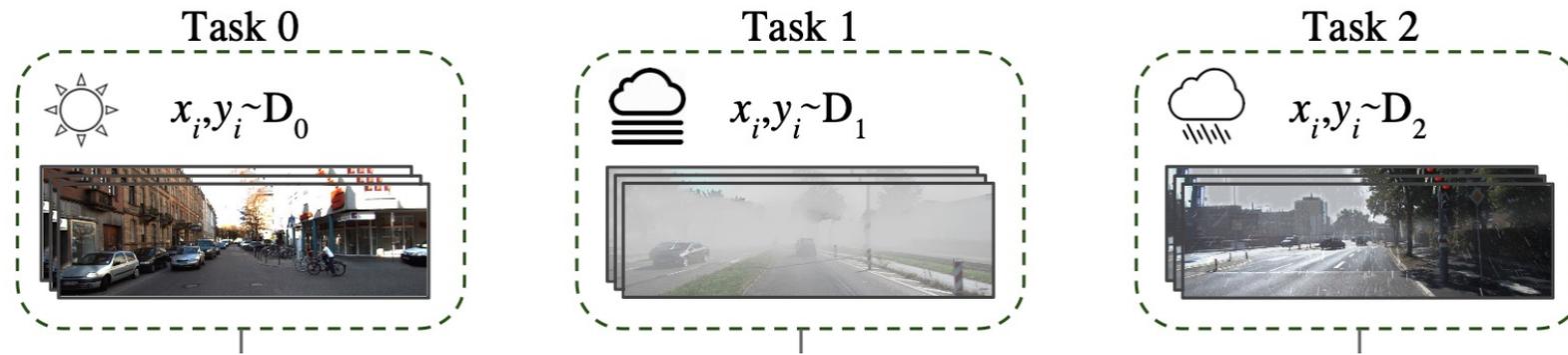


- Memory constraint: no (or very limited size of) the past data can be stored during training.



- Domain Incremental Learning (DIL)

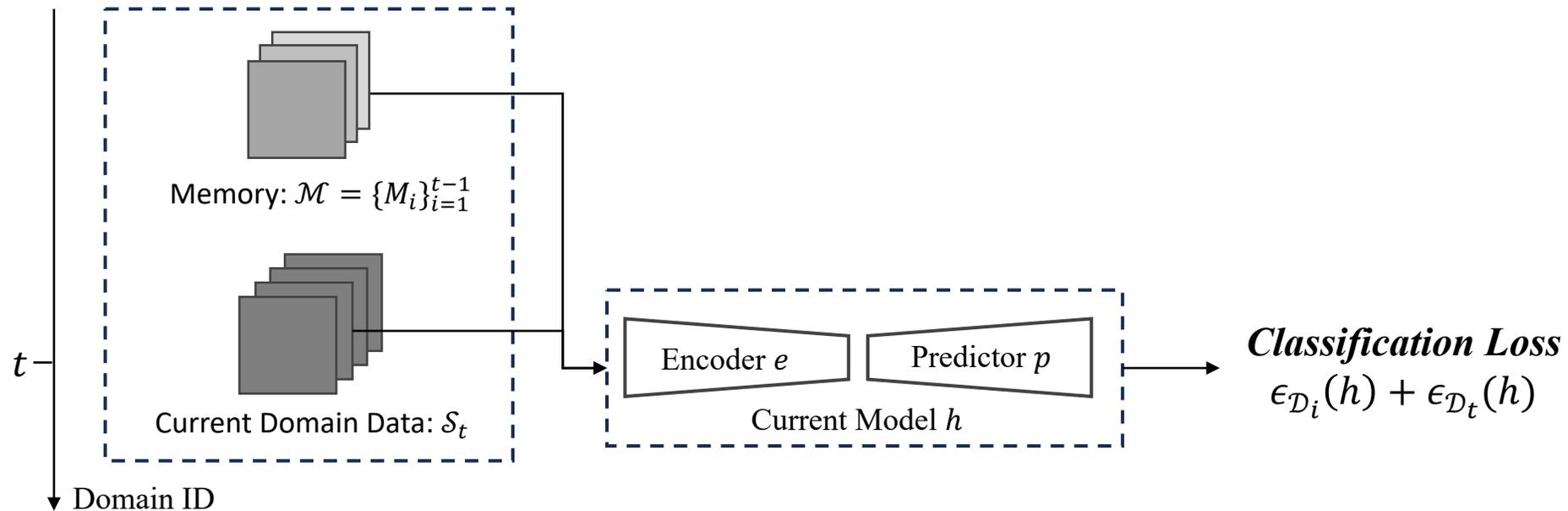
- Machine learning models are required to incrementally learn the evolving data distributions.
- E.g., autonomous driving under different weather conditions.



- Memory constraint: no (or very limited size of) the past data can be stored during training.
- Goal of DIL: minimize the model's risk over *all domains*.

$$\mathcal{L}^*(\theta) = \mathcal{L}_t(\theta) + \mathcal{L}_{1:t-1}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_t} [\ell(y, h_\theta(x))] + \sum_{i=1}^{t-1} \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(y, h_\theta(x))]$$

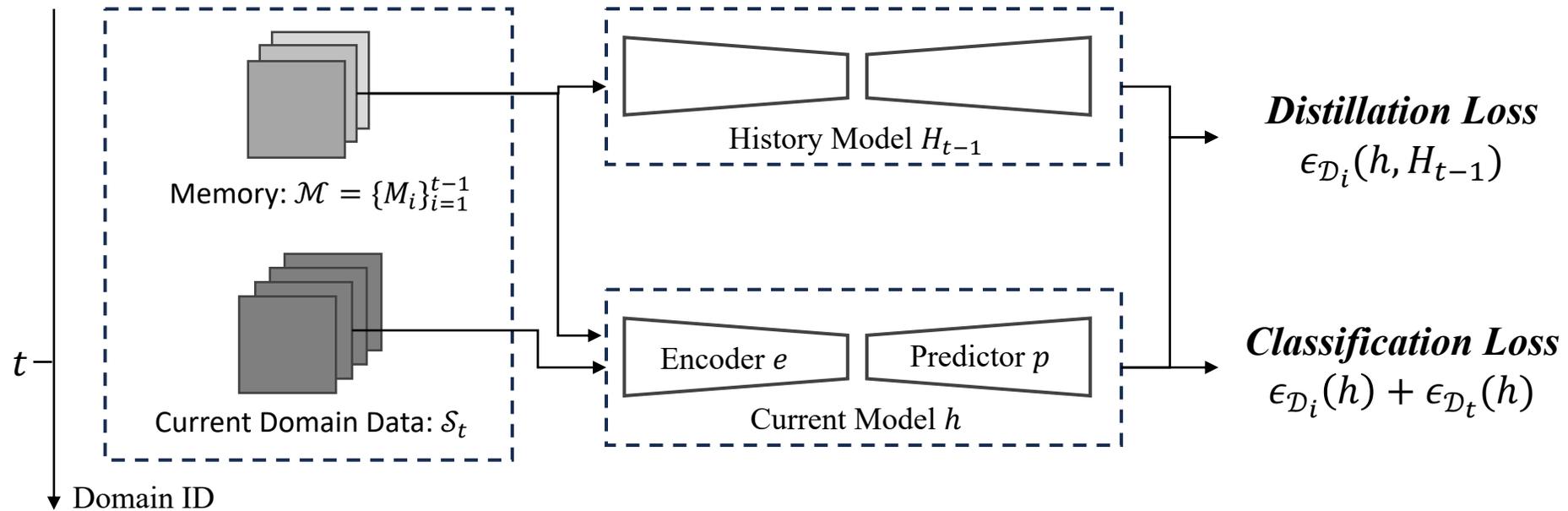
- Empirical Risk Minimization (ERM) via Experience Replay (ER)



- [Lemma 3.1] Trivially replaying the memory will cause *a loose generalization bound*.

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

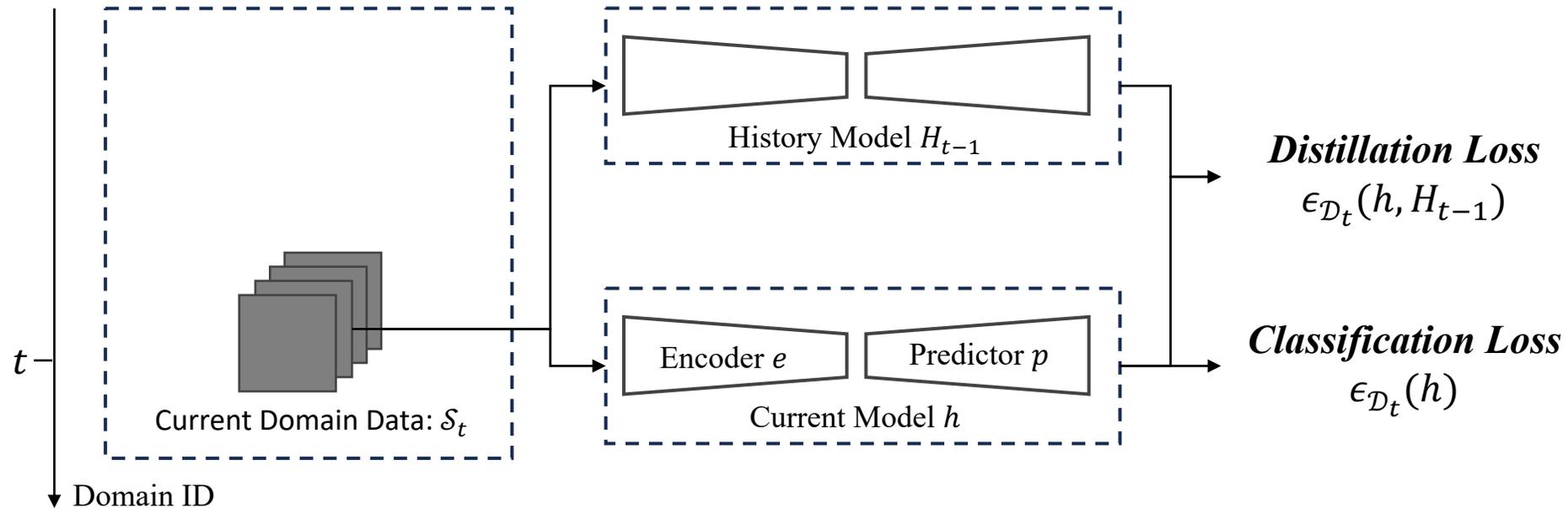
- Dark Experience Replay (DER++)



- [Lemma 3.2] Intra-Domain Model-Based Bound

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

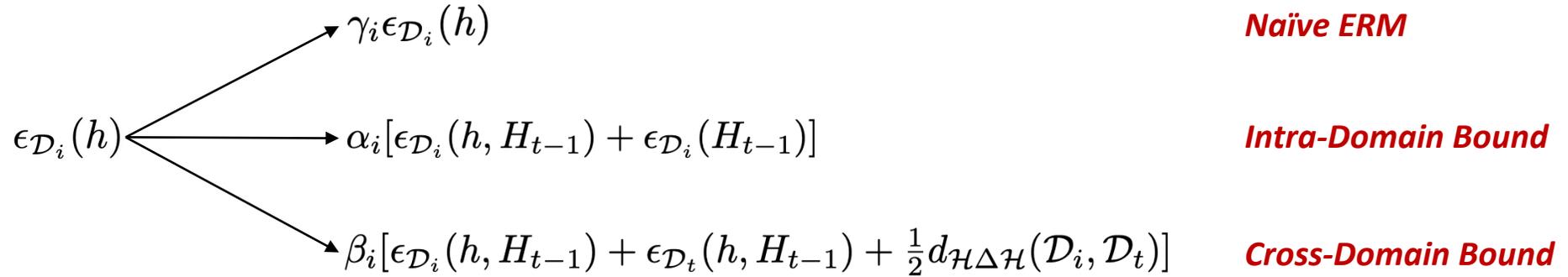
- Learning without Forgetting (LwF)



- [Lemma 3.3] Cross-Domain Model-Based Bound

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

- A set of coefficients $\{\alpha_i, \beta_i, \gamma_i\}_{i=1}^{t-1}$ (with $\alpha_i + \beta_i + \gamma_i = 1$) integrates them into one unified bound.



- [Theorem 3.4] Unified Generalization Bound for all domains

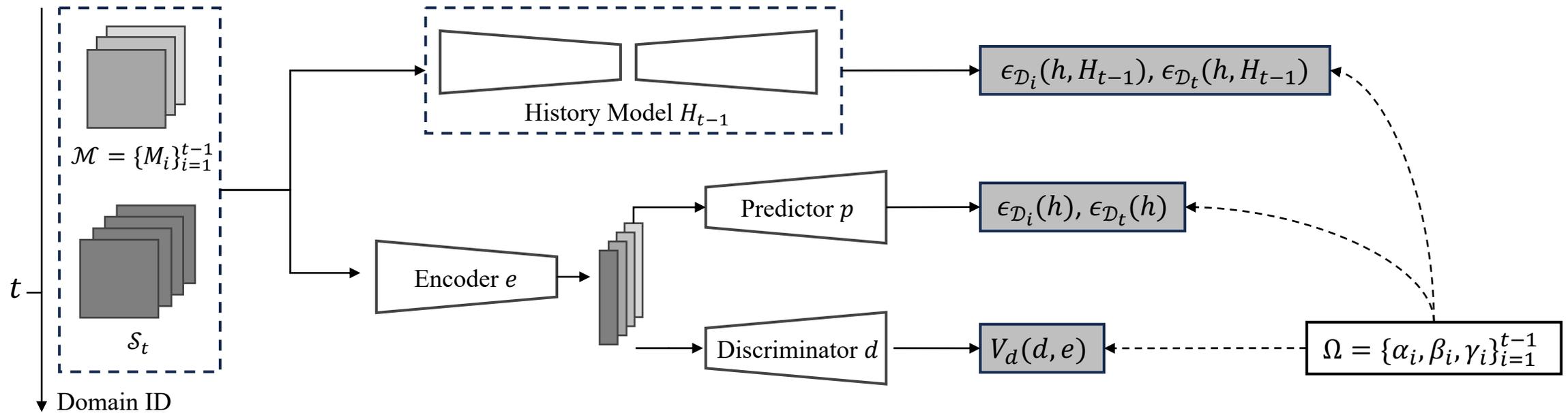
$$\begin{aligned} \sum_{i=1}^t \epsilon_{\mathcal{D}_i}(h) &\leq \left\{ \sum_{i=1}^{t-1} [\gamma_i \hat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \hat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1})] \right\} + \left\{ \hat{\epsilon}_{\mathcal{D}_t}(h) + \left(\sum_{i=1}^{t-1} \beta_i \right) \hat{\epsilon}_{\mathcal{D}_t}(h,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

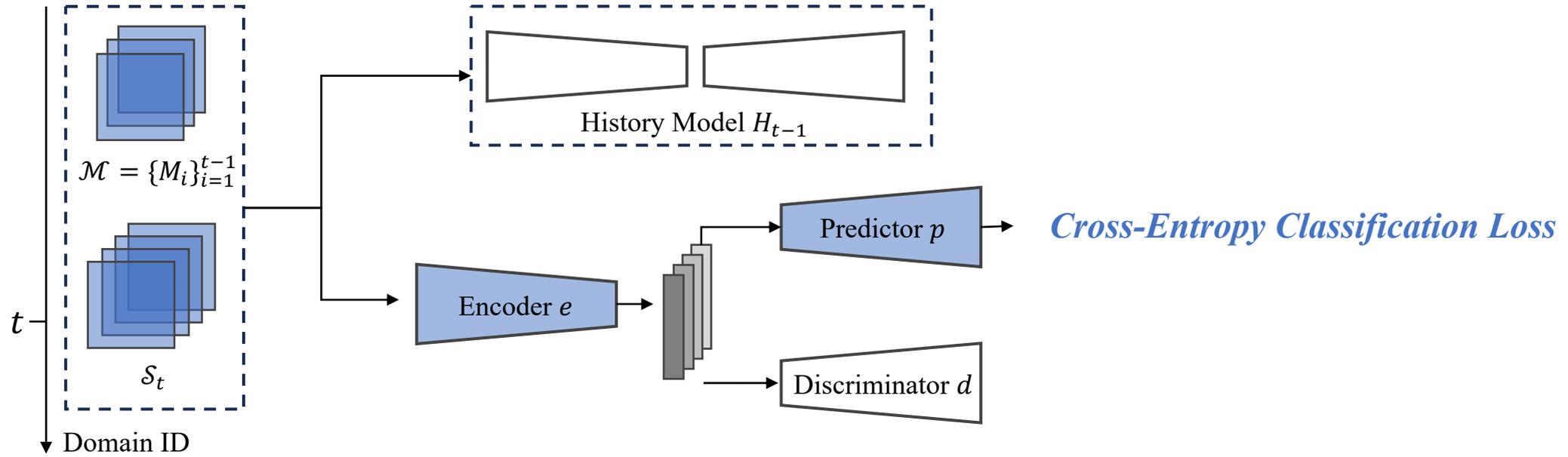
- UDIL *unifies* multiple existing methods under certain conditions.

	α_i	β_i	γ_i	Transformed Objective	Condition
UDIL (Ours)	$[0, 1]$	$[0, 1]$	$[0, 1]$	-	-
LwF [52]	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 [75]	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'_t }{(t-1)}$
DER++ [8]	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'_t }{(t-1)}$
iCaRL [74]	1	0	0	$\mathcal{L}_{\text{iCaRL}}(h) = \widehat{\ell}'_{B_t}(h) + \sum_{i=1}^{t-1} \frac{ B'_i /(t-1)}{ B_t } \widehat{\ell}'_{B'_i}(h, H_{t-1})$	$ B_t = \frac{ B'_t }{(t-1)}$
CLS-ER [4]	$\frac{\lambda}{\lambda+1}$	0	$\frac{1}{\lambda+1}$	$\mathcal{L}_{\text{CLS-ER}}(h) = \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{1}{t-1} \widehat{\ell}_{B'_i}(h) + \sum_{i=1}^{t-1} \frac{\lambda}{t-1} \widehat{\ell}_{B'_i}(h, H_{t-1})$	$\lambda = t - 2$
ESM-ER [80]	$\frac{\lambda}{\lambda+1}$	0	$\frac{1}{\lambda+1}$	$\mathcal{L}_{\text{ESM-ER}}(h) = \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{1}{r(t-1)} \widehat{\ell}_{B'_i}(h) + \sum_{i=1}^{t-1} \frac{\lambda}{r(t-1)} \widehat{\ell}_{B'_i}(h, H_{t-1})$	$\begin{cases} \lambda = -1 + r(t-1) \\ r = 1 - e^{-1} \end{cases}$
BiC [100]	$\frac{t-1}{2t-1}$	$\frac{t-1}{2t-1}$	$\frac{1}{2t-1}$	$\mathcal{L}_{\text{BiC}}(h) = \widehat{\ell}_{B_t}(h) + \sum_{i=1}^{t-1} \frac{(t-1) B_i }{ B_t } \widehat{\ell}_{B'_i}(h, H_{t-1}) + (t-1) \widehat{\ell}_{B_t}(h, H_{t-1}) + \sum_{i=1}^{t-1} \frac{ B_i }{ B_t } \widehat{\ell}_{B'_i}(h)$	$ B_i = B_t $

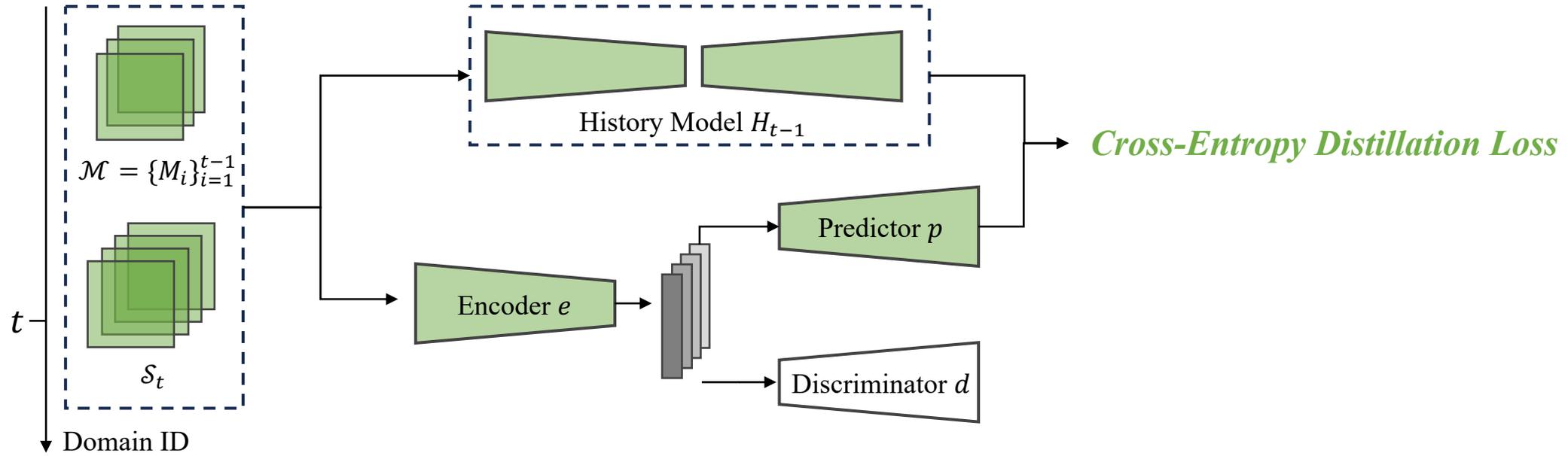
UDIL: An *Adaptive* Bound for DIL

- UDIL can *adaptively* adjust the coefficients based on the data and the history model H_{t-1} .
- It will, ideally, minimize the *tightest bound* in the family of all the generalization bounds.

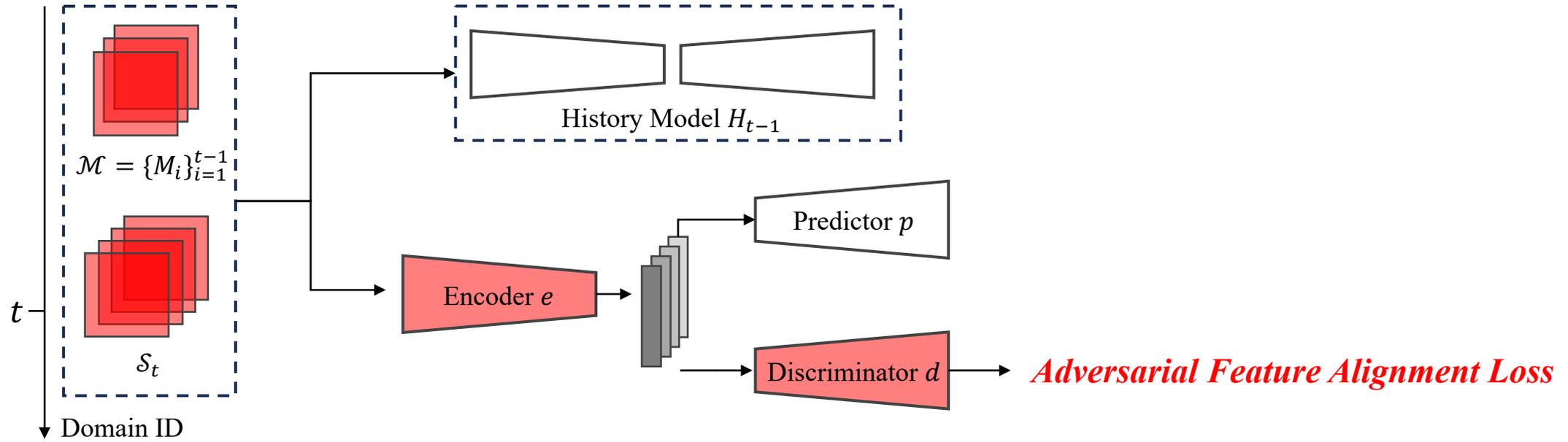




$$\begin{aligned}
 \sum_{i=1}^t \epsilon_{\mathcal{D}_i}(h) &\leq \left\{ \sum_{i=1}^{t-1} [\gamma_i \hat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \hat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1})] \right\} + \left\{ \hat{\epsilon}_{\mathcal{D}_t}(h) + \left(\sum_{i=1}^{t-1} \beta_i \right) \hat{\epsilon}_{\mathcal{D}_t}(h,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

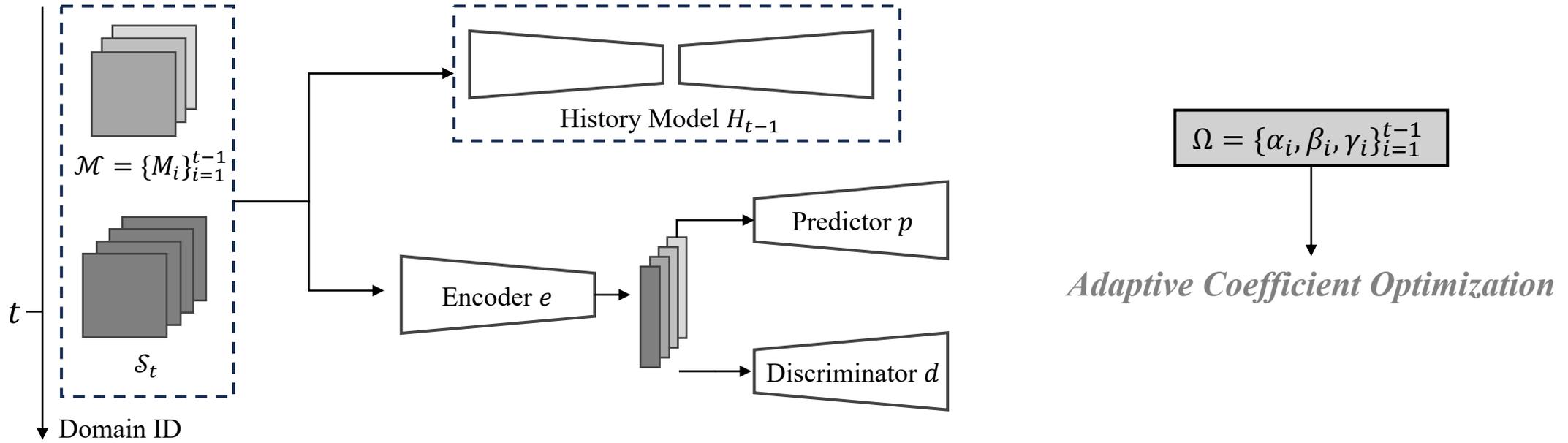


$$\begin{aligned}
 \sum_{i=1}^t \epsilon_{\mathcal{D}_i}(h) &\leq \left\{ \sum_{i=1}^{t-1} [\gamma_i \hat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \hat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1})] \right\} + \left\{ \hat{\epsilon}_{\mathcal{D}_t}(h) + \left(\sum_{i=1}^{t-1} \beta_i \right) \hat{\epsilon}_{\mathcal{D}_t}(h,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



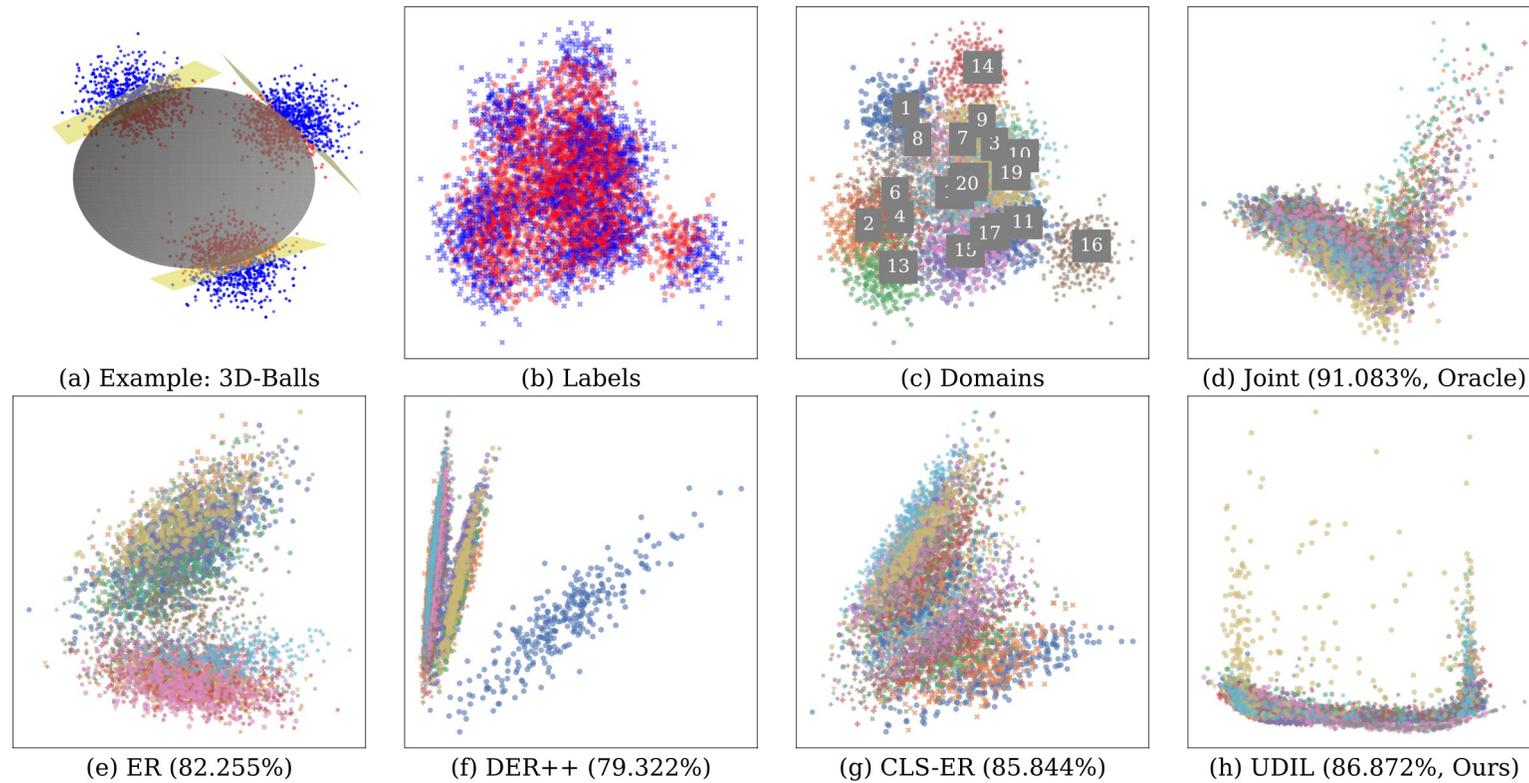
$$\begin{aligned}
 \sum_{i=1}^t \epsilon_{\mathcal{D}_i}(h) &\leq \left\{ \sum_{i=1}^{t-1} [\gamma_i \hat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \hat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1})] \right\} + \left\{ \hat{\epsilon}_{\mathcal{D}_t}(h) + \left(\sum_{i=1}^{t-1} \beta_i \right) \hat{\epsilon}_{\mathcal{D}_t}(h,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

UDIL: An *Adaptive* Bound for DIL



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

- UDIL's representation distribution on synthetic dataset (high-dimensional balls)



- UDIL evaluated on realistic datasets.

HD-Balls, Permuted-MNIST, Rotated-MNIST

Method	Buffer	<i>HD-Balls</i>		<i>P-MNIST</i>		<i>R-MNIST</i>	
		Avg. Acc (\uparrow)	Forgetting (\downarrow)	Avg. Acc (\uparrow)	Forgetting (\downarrow)	Avg. Acc (\uparrow)	Forgetting (\downarrow)
Fine-tune	-	52.319 \pm 0.024	43.520 \pm 0.079	70.102 \pm 2.945	27.522 \pm 3.042	47.803 \pm 1.703	52.281 \pm 1.797
oEWC [47]	-	54.131 \pm 0.193	39.743 \pm 1.388	78.476 \pm 1.223	18.068 \pm 1.321	48.203 \pm 0.827	51.181 \pm 0.867
SI [60]	-	52.303 \pm 0.037	43.175 \pm 0.041	79.045 \pm 1.357	17.409 \pm 1.446	48.251 \pm 1.381	51.053 \pm 1.507
LwF [26]	-	51.523 \pm 0.065	25.155 \pm 0.264	73.545 \pm 2.646	24.556 \pm 2.789	54.709 \pm 0.515	45.473 \pm 0.565
GEM [31]		69.747 \pm 0.656	13.591 \pm 0.779	89.097 \pm 0.149	6.975 \pm 0.167	76.619 \pm 0.581	21.289 \pm 0.579
A-GEM [7]		62.777 \pm 0.295	12.878 \pm 1.588	87.560 \pm 0.087	8.577 \pm 0.053	59.654 \pm 0.122	39.196 \pm 0.171
ER [42]		82.255 \pm 1.552	9.524 \pm 1.655	88.339 \pm 0.044	7.180 \pm 0.029	76.794 \pm 0.696	20.696 \pm 0.744
DER++ [5]	400	79.332 \pm 1.347	13.762 \pm 1.514	92.950\pm0.361	3.378 \pm 0.245	84.258 \pm 0.544	13.692 \pm 0.560
CLS-ER [2]		85.844 \pm 0.165	5.297 \pm 0.281	91.598 \pm 0.117	3.795 \pm 0.144	81.771 \pm 0.354	15.455 \pm 0.356
ESM-ER [46]		71.995 \pm 3.833	13.245 \pm 5.397	89.829 \pm 0.698	6.888 \pm 0.738	82.192 \pm 0.164	16.195 \pm 0.150
UDIL (Ours)		86.872\pm0.195	3.428\pm0.359	92.666\pm0.108	2.853\pm0.107	86.635\pm0.686	8.506\pm1.181
Joint (Oracle)	∞	91.083 \pm 0.332	-	96.368 \pm 0.042	-	97.150 \pm 0.036	-

- UDIL evaluated on realistic datasets.

Sequential C_{ORE}-50

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
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

- Proposed a principled framework, UDIL, for domain incremental learning with memory to *unify various existing methods*.
- Theoretical analysis shows that different existing methods are equivalent to minimizing the same error bound with different *fixed* coefficients.
- UDIL allows *adaptive* coefficients during training, thereby always achieving the tightest bound and improving the performance.

