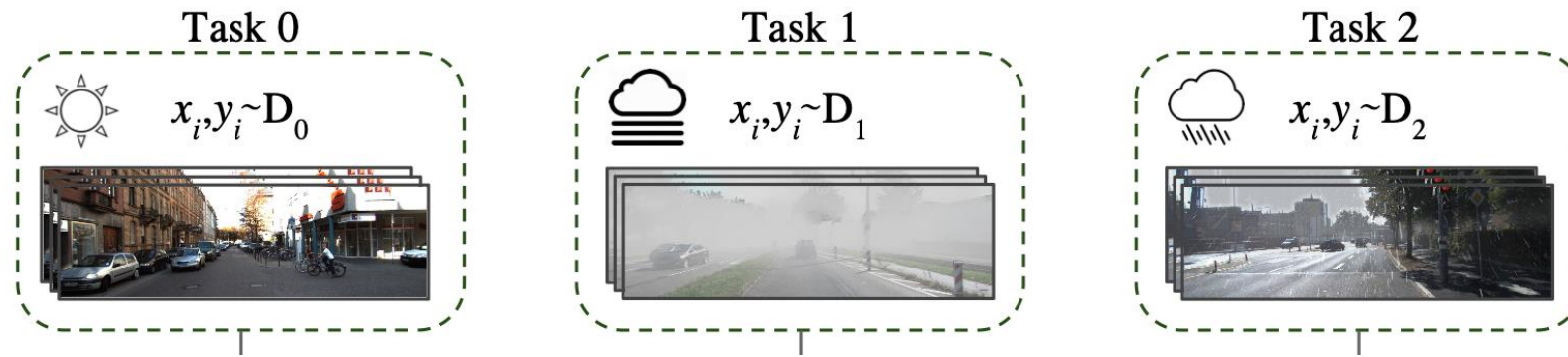# A Unified Approach to Domain Incremental Learning with Memory: Theory and Algorithm
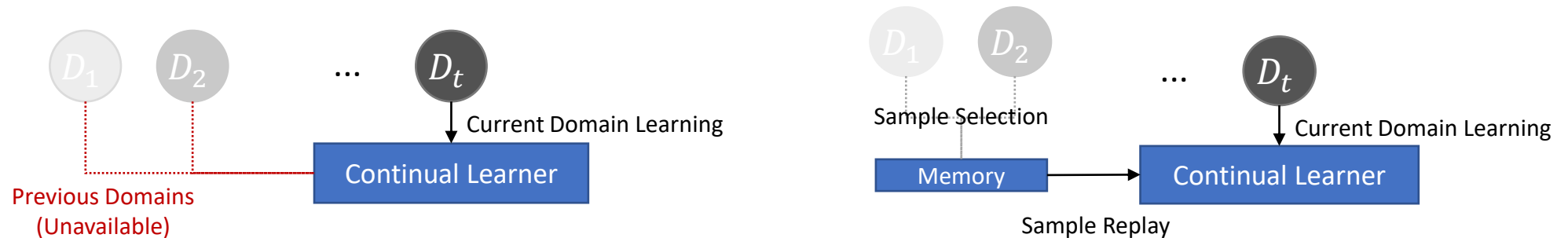
Haizhou Shi, Hao Wang

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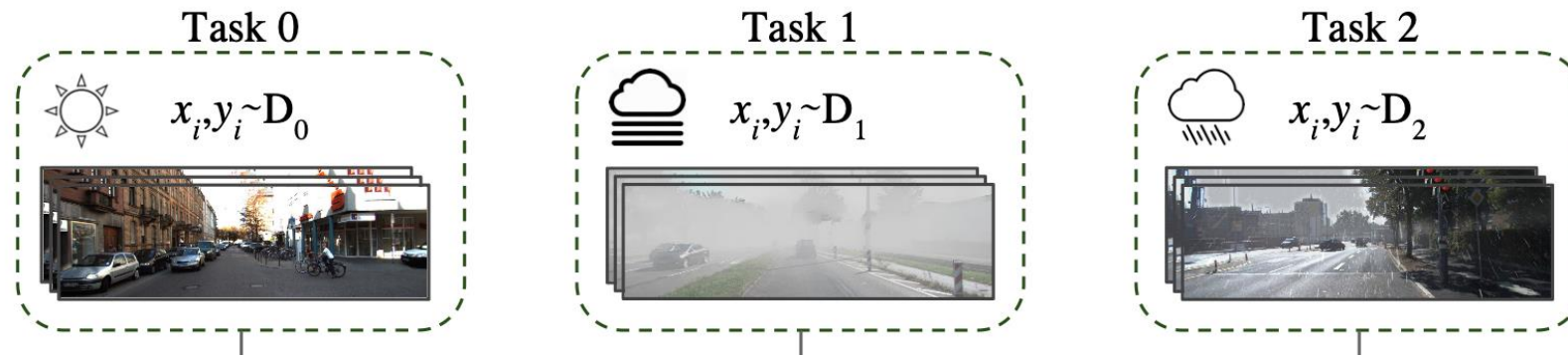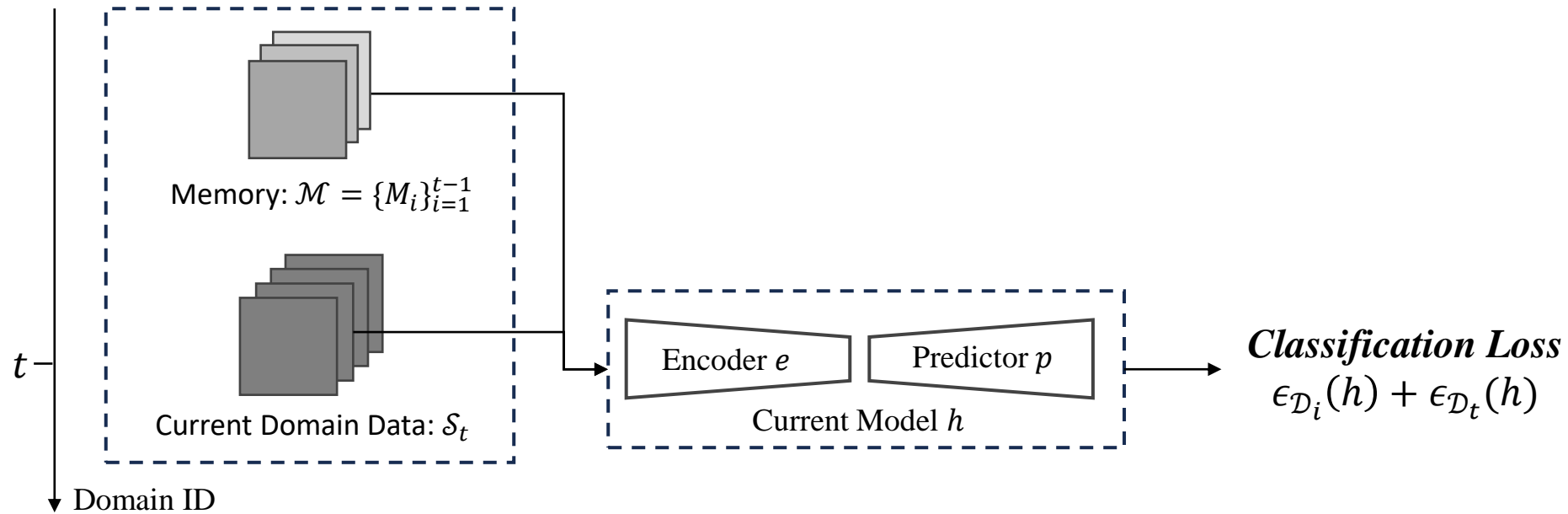
- Domain Incremental Learning (DIL)
  - Machine learning models are required to incrementally learn the evolving data distributions.
  - E.g., autonomous driving under different weather conditions.



- 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***.

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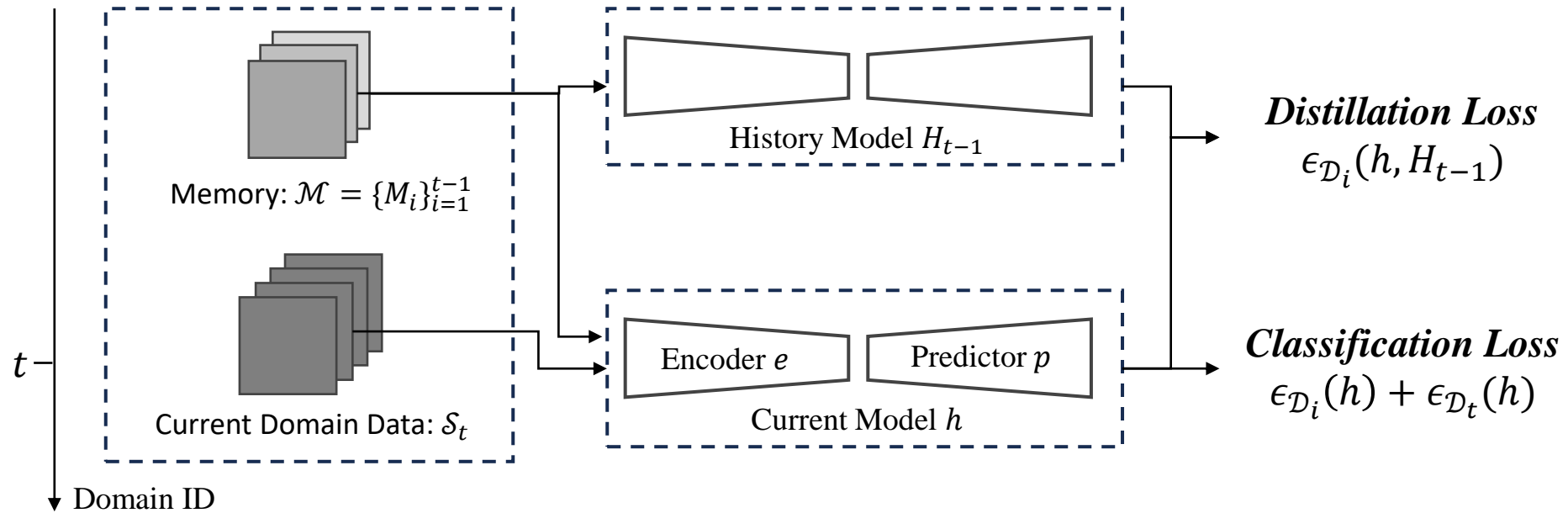
- Empirical Risk Minimization (ERM) via Experience Replay (ER)



Memory: $\mathcal{M} = \{M_i\}_{i=1}^{t-1}$

Current Domain Data: $\mathcal{S}_t$

Encoder $e$   Predictor $p$

Current Model $h$

**Classification Loss**
$\epsilon_{\mathcal{D}_i}(h) + \epsilon_{\mathcal{D}_t}(h)$

$t$

Domain ID

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

$$\sum_{i=1}^{t} \epsilon_{\mathcal{D}_i}(h) \le \sum_{i=1}^{t} \widehat{\epsilon}_{\mathcal{D}_i}(h) + \sqrt{\left(\frac{1}{N_t} + \sum_{i=1}^{t-1} \frac{1}{\widetilde{N}_i}\right)\left(8d \log\left(\frac{2eN}{d}\right) + 8 \log\left(\frac{2}{\delta}\right)\right)}.$$
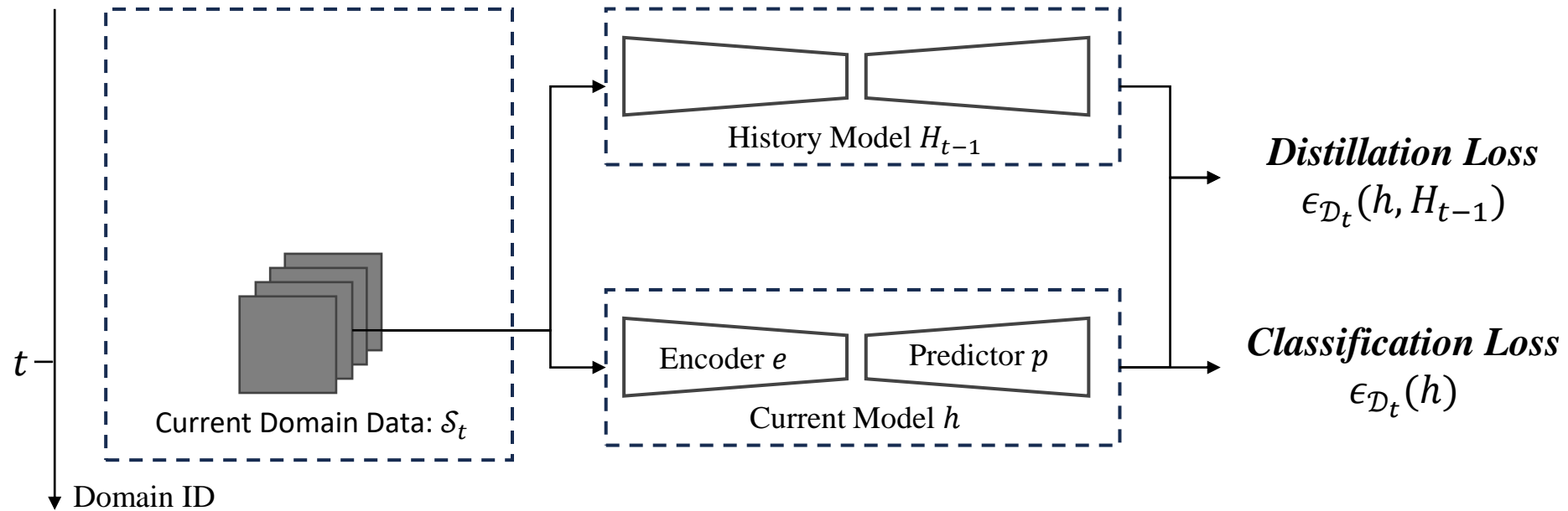
- Dark Experience Replay (DER++)



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

$$\epsilon_{\mathcal{D}_i}(h) \leq \epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1}),$$
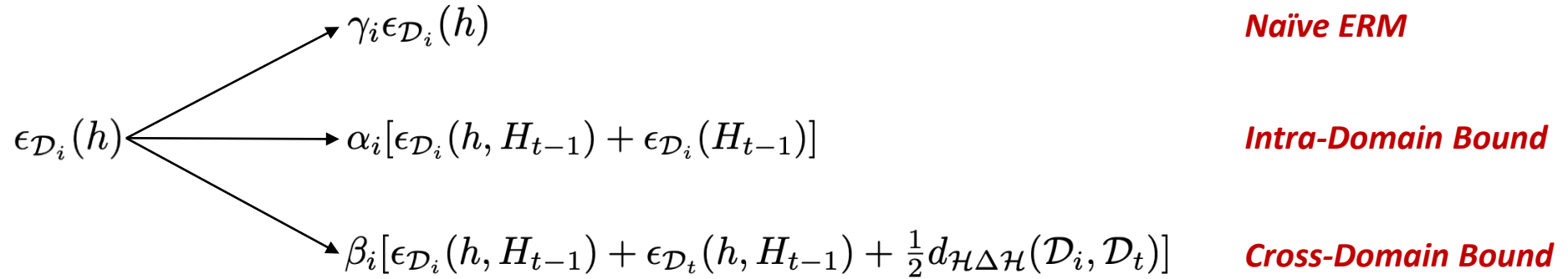
- Learning without Forgetting (LwF)



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

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

- 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.
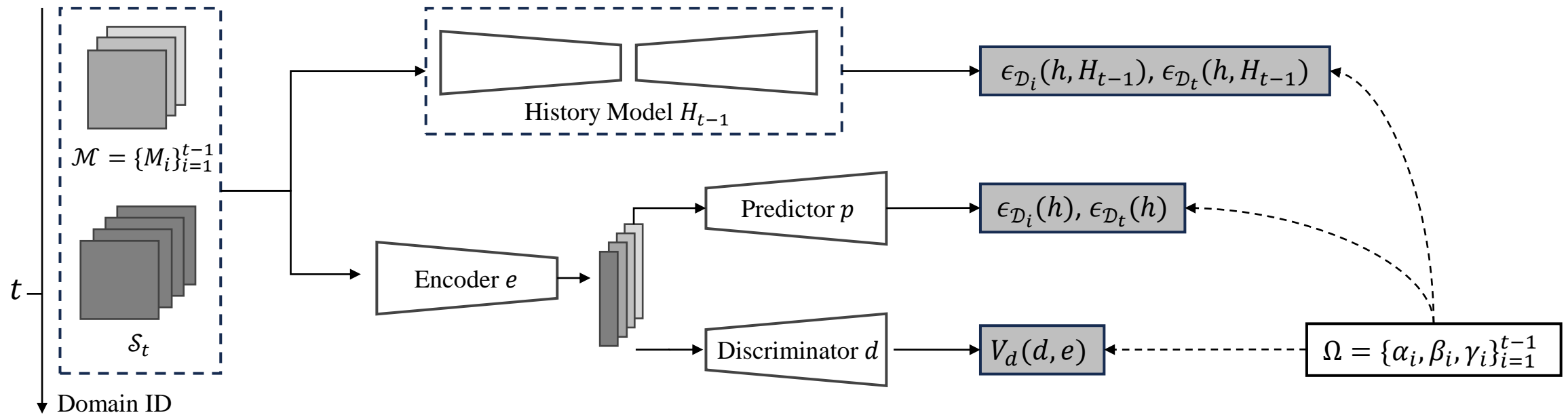
$$\epsilon_{\mathcal{D}_i}(h) \begin{cases} \nearrow \gamma_i \epsilon_{\mathcal{D}_i}(h) & \textit{Naïve ERM} \\ \rightarrow \alpha_i[\epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_i}(H_{t-1})] & \textit{Intra-Domain Bound} \\ \searrow \beta_i[\epsilon_{\mathcal{D}_i}(h, H_{t-1}) + \epsilon_{\mathcal{D}_t}(h, H_{t-1}) + \tfrac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t)] & \textit{Cross-Domain Bound} \end{cases}$$
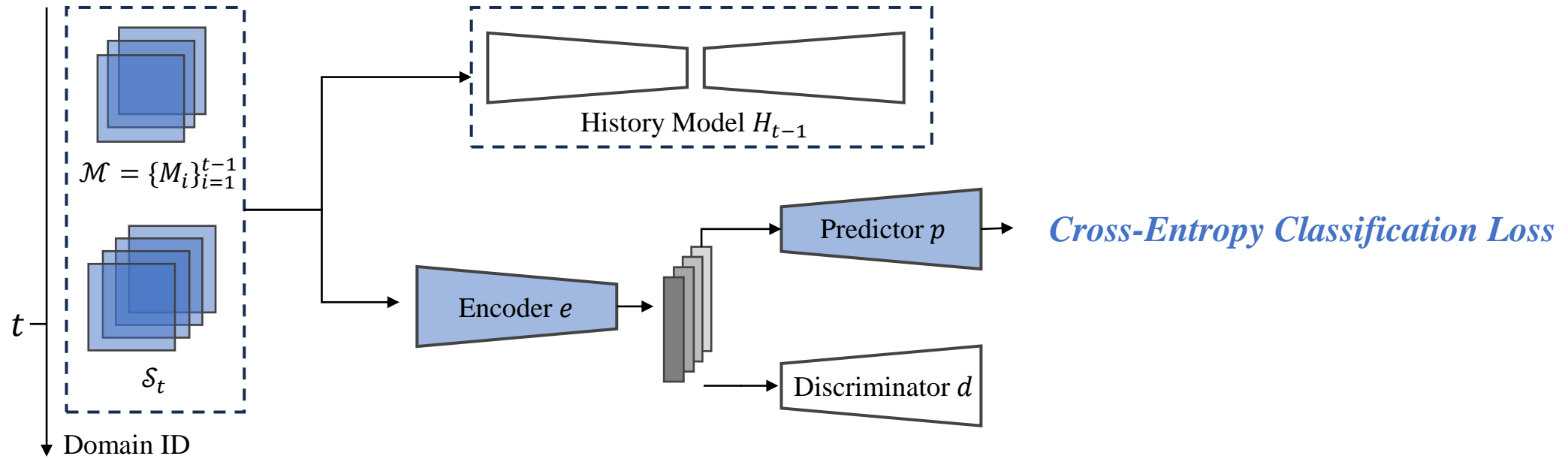
- [Theorem 3.4] Unified Generalization Bound for all domains

$$\sum_{i=1}^{t} \epsilon_{\mathcal{D}_i}(h) \leq \left\{ \sum_{i=1}^{t-1} [\gamma_i \widehat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \widehat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1})] \right\} + \left\{ \widehat{\epsilon}_{\mathcal{D}_t}(h) + (\sum_{i=1}^{t-1} \beta_i) \widehat{\epsilon}_{\mathcal{D}_t}(h, H_{t-1}) \right\}$$

$$+ \frac{1}{2} \sum_{i=1}^{t-1} \beta_i d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \sum_{i=1}^{t-1} (\alpha_i + \beta_i) \epsilon_{\mathcal{D}_i}(H_{t-1})$$

$$+ \sqrt{\left( \frac{(1+\sum_{i=1}^{t-1} \beta_i)^2}{N_t} + \sum_{i=1}^{t-1} \frac{(\gamma_i + \alpha_i)^2}{\widetilde{N}_i} \right) \left( 8d \log\left(\frac{2eN}{d}\right) + 8\log\left(\frac{2}{\delta}\right) \right)}$$

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

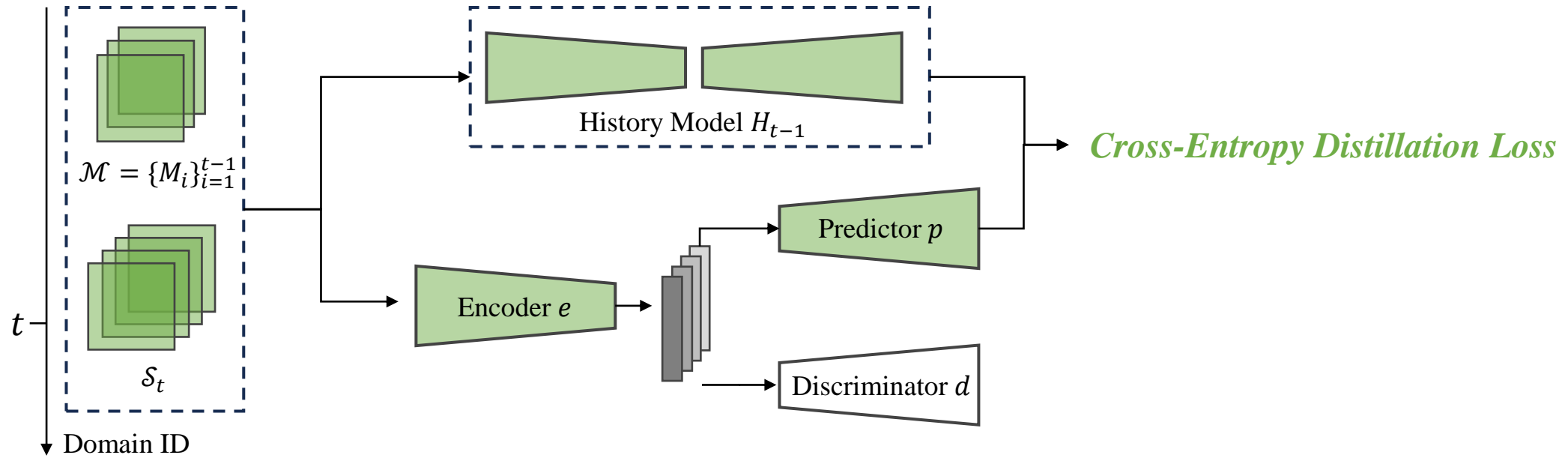| | $\alpha_i$ | $\beta_i$ | $\gamma_i$ | Transformed Objective | Condition |
|---|---|---|---|---|---|
| UDIL (Ours) | $[0,1]$ | $[0,1]$ | $[0,1]$ | - | - |
| LwF [52] | $0$ | $1$ | $0$ | $\mathcal{L}_{\mathrm{LwF}}(h) = \widehat{\ell}_{\mathcal{X}_t}(h) + \lambda_o \widehat{\ell}_{\mathcal{X}_t}(h,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

- UDIL 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.
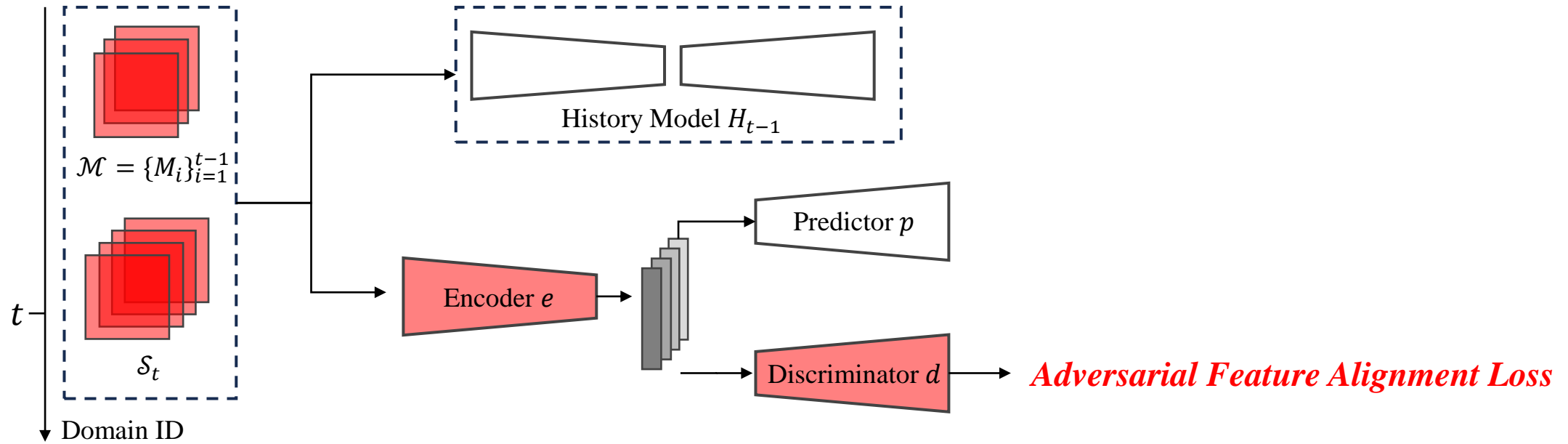
$$\sum_{i=1}^{t} \epsilon_{\mathcal{D}_i}(h) \leq \left\{ \sum_{i=1}^{t-1} \left[ \boxed{\gamma_i \widehat{\epsilon}_{\mathcal{D}_i}(h)} + \alpha_i \widehat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1}) \right] \right\} + \left\{ \boxed{\widehat{\epsilon}_{\mathcal{D}_t}(h)} + \left( \sum_{i=1}^{t-1} \beta_i \right) \widehat{\epsilon}_{\mathcal{D}_t}(h, H_{t-1}) \right\}$$

$$+ \frac{1}{2} \sum_{i=1}^{t-1} \beta_i d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \sum_{i=1}^{t-1} (\alpha_i + \beta_i) \epsilon_{\mathcal{D}_i}(H_{t-1})$$

$$+ \sqrt{ \left( \frac{(1 + \sum_{i=1}^{t-1} \beta_i)^2}{N_t} + \sum_{i=1}^{t-1} \frac{(\gamma_i + \alpha_i)^2}{\widetilde{N}_i} \right) \left( 8d \log\left( \frac{2eN}{d} \right) + 8 \log\left( \frac{2}{\delta} \right) \right) }$$

$$\sum_{i=1}^{t} \epsilon_{\mathcal{D}_i}(h) \leq \left\{ \sum_{i=1}^{t-1} \left[ \gamma_i \widehat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \widehat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1}) \right] \right\} + \left\{ \widehat{\epsilon}_{\mathcal{D}_t}(h) + \left( \sum_{i=1}^{t-1} \beta_i \right) \widehat{\epsilon}_{\mathcal{D}_t}(h, H_{t-1}) \right\}$$

$$+ \frac{1}{2} \sum_{i=1}^{t-1} \beta_i d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \sum_{i=1}^{t-1} (\alpha_i + \beta_i) \epsilon_{\mathcal{D}_i}(H_{t-1})$$

$$+ \sqrt{ \left( \frac{(1 + \sum_{i=1}^{t-1} \beta_i)^2}{N_t} + \sum_{i=1}^{t-1} \frac{(\gamma_i + \alpha_i)^2}{\widetilde{N}_i} \right) \left( 8d \log\left( \frac{2eN}{d} \right) + 8 \log\left( \frac{2}{\delta} \right) \right) }$$
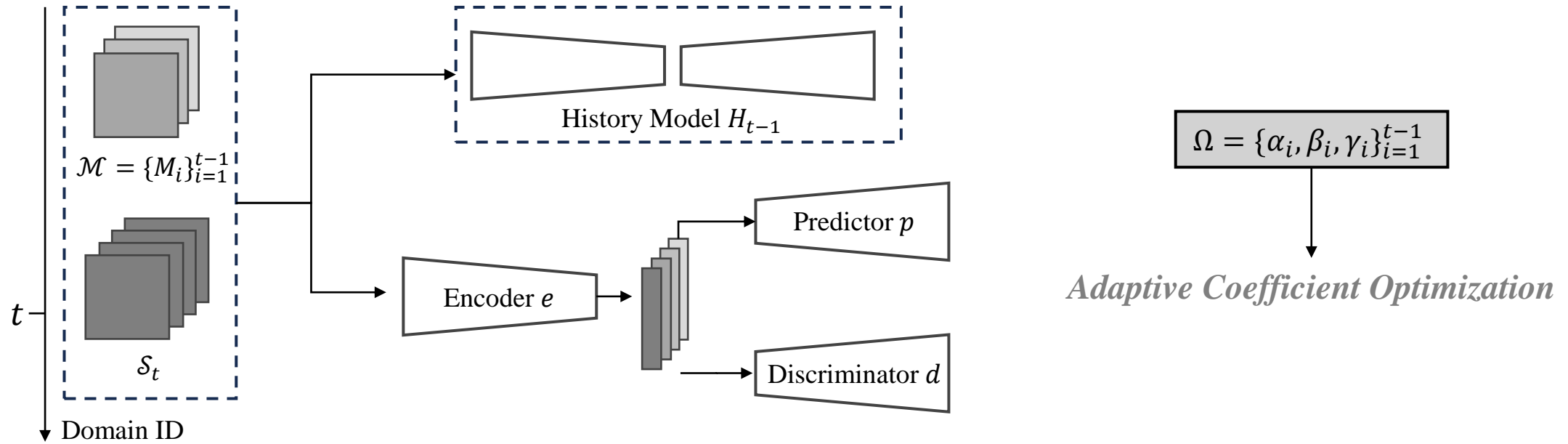
RUTGERS UNIVERSITY



$$\sum_{i=1}^{t} \epsilon_{\mathcal{D}_i}(h) \leq \left\{ \sum_{i=1}^{t-1} \left[ \gamma_i \widehat{\epsilon}_{\mathcal{D}_i}(h) + \alpha_i \widehat{\epsilon}_{\mathcal{D}_i}(h, H_{t-1}) \right] \right\} + \left\{ \widehat{\epsilon}_{\mathcal{D}_t}(h) + \left( \sum_{i=1}^{t-1} \beta_i \right) \widehat{\epsilon}_{\mathcal{D}_t}(h, H_{t-1}) \right\}$$

$$+ \frac{1}{2} \sum_{i=1}^{t-1} \beta_i d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_t) + \sum_{i=1}^{t-1} (\alpha_i + \beta_i) \epsilon_{\mathcal{D}_i}(H_{t-1})$$

$$+ \sqrt{ \left( \frac{(1 + \sum_{i=1}^{t-1} \beta_i)^2}{N_t} + \sum_{i=1}^{t-1} \frac{(\gamma_i + \alpha_i)^2}{\widetilde{N}_i} \right) \left( 8d \log\left( \frac{2eN}{d} \right) + 8 \log\left( \frac{2}{\delta} \right) \right) }$$
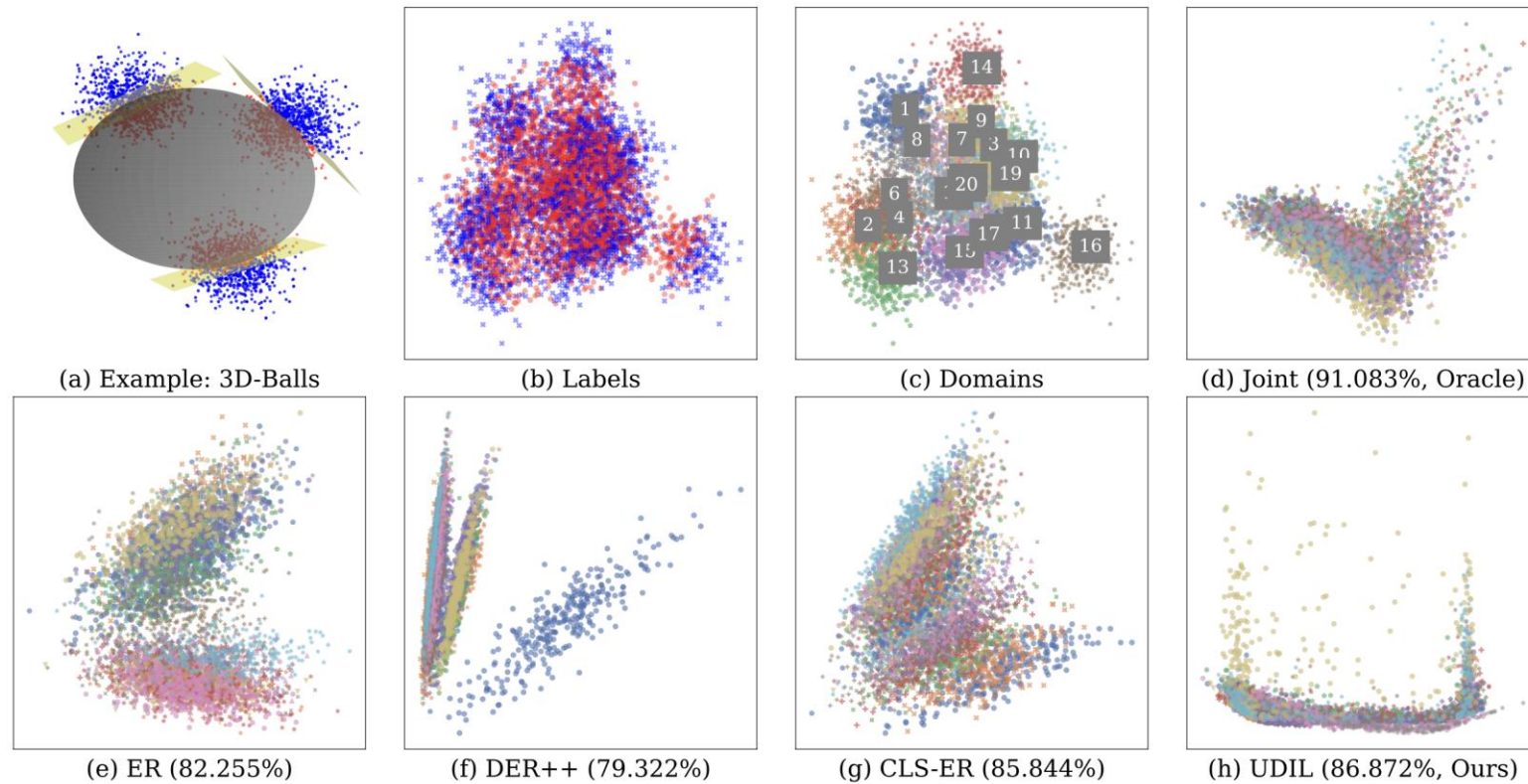
- UDIL's representation distribution on synthetic dataset (high-dimensional balls)



(a) Example: 3D-Balls    (b) Labels    (c) Domains    (d) Joint (91.083%, Oracle)

(e) ER (82.255%)    (f) DER++ (79.322%)    (g) CLS-ER (85.844%)    (h) UDIL (86.872%, Ours)

- UDIL evaluated on realistic datasets.

### *HD-Balls, Permuted-MNIST, Rotated-MNIST*

| Method | Buffer | HD-Balls | | P-MNIST | | R-MNIST | |
|--------|--------|----------|--|---------|--|---------|--|
| | | Avg. Acc (↑) | Forgetting (↓) | Avg. Acc (↑) | Forgetting (↓) | Avg. Acc (↑) | Forgetting (↓) |
| 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 | $79.332_{\pm 1.347}$ | $13.762_{\pm 1.514}$ | $\mathbf{92.950_{\pm 0.361}}$ | $3.378_{\pm 0.245}$ | $84.258_{\pm 0.544}$ | $13.692_{\pm 0.560}$ |
| 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) | | $\mathbf{86.872_{\pm 0.195}}$ | $\mathbf{3.428_{\pm 0.359}}$ | $92.666_{\pm 0.108}$ | $\mathbf{2.853_{\pm 0.107}}$ | $\mathbf{86.635_{\pm 0.686}}$ | $\mathbf{8.506_{\pm 1.181}}$ |
| Joint (Oracle) | ∞ | $91.083_{\pm 0.332}$ | - | $96.368_{\pm 0.042}$ | - | $97.150_{\pm 0.036}$ | - |

- UDIL evaluated on realistic datasets.

### Sequential CORe-50

| Method | Buffer | $\mathcal{D}_{1:3}$ | $\mathcal{D}_{4:6}$ | $\mathcal{D}_{7:9}$ | $\mathcal{D}_{10:11}$ | Overall | |
|---|---|---|---|---|---|---|---|
| | | Avg. Acc (↑) | | | | Avg. Acc (↑) | Forgetting (↓) |
| Fine-tune | - | $73.707_{\pm13.144}$ | $34.551_{\pm1.254}$ | $29.406_{\pm2.579}$ | $28.689_{\pm3.144}$ | $31.832_{\pm1.034}$ | $73.296_{\pm1.399}$ |
| oEWC [51] | - | $74.567_{\pm13.360}$ | $35.915_{\pm0.260}$ | $30.174_{\pm3.195}$ | $28.291_{\pm2.522}$ | $30.813_{\pm1.154}$ | $74.563_{\pm0.937}$ |
| SI [66] | - | $74.661_{\pm14.162}$ | $34.345_{\pm1.001}$ | $30.127_{\pm2.971}$ | $28.839_{\pm3.631}$ | $32.469_{\pm1.315}$ | $73.144_{\pm1.588}$ |
| LwF [29] | - | $80.383_{\pm10.190}$ | $28.357_{\pm1.143}$ | $31.386_{\pm0.787}$ | $28.711_{\pm2.981}$ | $31.692_{\pm0.768}$ | $72.990_{\pm1.350}$ |
| GEM [34] | | $79.852_{\pm6.864}$ | $38.961_{\pm1.718}$ | $39.258_{\pm2.614}$ | $36.859_{\pm0.842}$ | $37.701_{\pm0.273}$ | $22.724_{\pm1.554}$ |
| A-GEM [8] | | $80.348_{\pm9.394}$ | $41.472_{\pm3.394}$ | $43.213_{\pm1.542}$ | $39.181_{\pm3.999}$ | $43.181_{\pm2.025}$ | $33.775_{\pm3.003}$ |
| ER [46] | | $90.838_{\pm2.177}$ | $79.343_{\pm2.699}$ | $68.151_{\pm0.226}$ | $65.034_{\pm1.571}$ | $66.605_{\pm0.214}$ | $32.750_{\pm0.455}$ |
| DER++ [6] | 500 | $92.444_{\pm1.764}$ | $88.652_{\pm1.854}$ | $80.391_{\pm0.107}$ | $78.038_{\pm0.591}$ | $78.629_{\pm0.753}$ | $21.910_{\pm1.094}$ |
| CLS-ER [3] | | $89.834_{\pm1.323}$ | $78.909_{\pm1.724}$ | $70.591_{\pm0.322}$ | ⋆ | ⋆ | ⋆ |
| ESM-ER [50] | | $84.905_{\pm6.471}$ | $51.905_{\pm3.257}$ | $53.815_{\pm1.770}$ | $50.178_{\pm2.574}$ | $52.751_{\pm1.296}$ | $25.444_{\pm0.580}$ |
| UDIL (Ours) | | $\mathbf{98.152}_{\pm1.665}$ | $\mathbf{89.814}_{\pm2.302}$ | $\mathbf{83.052}_{\pm0.151}$ | $\mathbf{81.547}_{\pm0.269}$ | $\mathbf{82.103}_{\pm0.279}$ | $\mathbf{19.589}_{\pm0.303}$ |
| GEM [34] | | $78.717_{\pm4.831}$ | $43.269_{\pm3.419}$ | $40.908_{\pm2.200}$ | $40.408_{\pm1.168}$ | $41.576_{\pm1.599}$ | $18.537_{\pm1.237}$ |
| A-GEM [8] | | $78.917_{\pm8.984}$ | $41.172_{\pm4.293}$ | $44.576_{\pm1.701}$ | $38.960_{\pm3.867}$ | $42.827_{\pm1.659}$ | $33.800_{\pm1.847}$ |
| ER [46] | | $90.048_{\pm2.699}$ | $84.668_{\pm1.988}$ | $77.561_{\pm1.281}$ | $72.268_{\pm0.720}$ | $72.988_{\pm0.566}$ | $25.997_{\pm0.694}$ |
| DER++ [6] | 1000 | $89.510_{\pm5.726}$ | $92.492_{\pm0.902}$ | $88.883_{\pm0.794}$ | $86.108_{\pm0.284}$ | $86.392_{\pm0.714}$ | $13.128_{\pm0.474}$ |
| CLS-ER [3] | | $92.004_{\pm0.894}$ | $85.044_{\pm1.276}$ | ⋆ | ⋆ | ⋆ | ⋆ |
| ESM-ER [50] | | $85.120_{\pm4.339}$ | $54.852_{\pm5.511}$ | $61.714_{\pm1.840}$ | $55.098_{\pm3.834}$ | $58.932_{\pm0.959}$ | $20.134_{\pm0.643}$ |
| UDIL (Ours) | | $\mathbf{98.648}_{\pm1.174}$ | $\mathbf{93.447}_{\pm1.111}$ | $\mathbf{90.545}_{\pm0.705}$ | $\mathbf{87.923}_{\pm0.232}$ | $\mathbf{88.155}_{\pm0.445}$ | $\mathbf{12.882}_{\pm0.460}$ |
| Joint (Oracle) | ∞ | - | - | - | - | $99.137_{\pm0.049}$ | - |

# Conclusion

- 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.


Paper


Code