

Table 6: Comparison of different MTL methods on CIFAR100 and VOC07+12 datasets. The 1<sup>st</sup>/2<sup>nd</sup> best results are indicated in red/blue.

Method	Before evolution (%) $\uparrow$			After evolution (%) $\uparrow$		
	CLS	DET	AVG	CLS	DET	AVG
Expert(CLS)+MAD	70.68 $\pm$ 0.00	0.00 $\pm$ 0.00	35.34 $\pm$ 0.00	70.52 $\pm$ 0.17	82.58 $\pm$ 0.13	<b>76.55<math>\pm</math>0.15</b>
Expert(CLS)+Pix2SeqV2	70.68 $\pm$ 0.00	0.00 $\pm$ 0.00	35.34 $\pm$ 0.00	70.24 $\pm$ 0.15	80.14 $\pm$ 0.09	75.19 $\pm$ 0.12
Expert(CLS)+Frozen	70.68 $\pm$ 0.00	0.00 $\pm$ 0.00	35.34 $\pm$ 0.00	<b>70.68<math>\pm</math>0.00</b>	59.92 $\pm$ 0.06	65.30 $\pm$ 0.03
Expert(DET)+MAD	70.68 $\pm$ 0.00	0.00 $\pm$ 0.00	35.34 $\pm$ 0.00	70.52 $\pm$ 0.17	82.58 $\pm$ 0.13	76.55 $\pm$ 0.15
Expert(DET)+Pix2SeqV2	70.68 $\pm$ 0.00	0.00 $\pm$ 0.00	35.34 $\pm$ 0.00	70.24 $\pm$ 0.15	80.14 $\pm$ 0.09	75.19 $\pm$ 0.12
Expert(DET)+Frozen	0.47 $\pm$ 0.00	83.89 $\pm$ 0.00	42.18 $\pm$ 0.00	43.05 $\pm$ 0.08	<b>83.89<math>\pm</math>0.00</b>	63.47 $\pm$ 0.04
DISC	70.68 $\pm$ 0.00	0.00 $\pm$ 0.00	35.34 $\pm$ 0.00	<b>72.28<math>\pm</math>0.08</b> <sub>(+1.60)</sub>	<b>84.92<math>\pm</math>0.04</b> <sub>(+1.03)</sub>	<b>78.60<math>\pm</math>0.06</b> <sub>(+2.05)</sub>

Table 7: The relationship from human society to machine society to method design.

Human Society	Machine Society	Method Design
Social hierarchy	Hierarchical organizational structures	Layer-wise hierarchical structures
Sequential progression	Progressive interaction modes	Dynamic hierarchical collaboration
Cultural learning	Strong-guided communication mechanisms	Dynamic selective collaboration

Table 8: Comparison of different backbones, training strategies and tasks on CIFAR100 and VOC07+12 datasets. The 1<sup>st</sup>/2<sup>nd</sup> best results are indicated in red/blue.

Method	CLS (%) $\uparrow$	DET (%) $\uparrow$	SEG (%) $\uparrow$	AVG (%) $\uparrow$
ResNet-50 (Scratch)	61.15 $\pm$ 0.08	40.11 $\pm$ 0.06	50.83 $\pm$ 0.12	50.69 $\pm$ 0.09
ResNet-50 (ImageNet)	64.02 $\pm$ 0.03	67.74 $\pm$ 0.09	55.63 $\pm$ 0.04	62.46 $\pm$ 0.05
ResNet-101 (ImageNet)	<b>64.73<math>\pm</math>0.09</b>	<b>69.66<math>\pm</math>0.05</b>	<b>56.29<math>\pm</math>0.06</b>	<b>63.56<math>\pm</math>0.07</b>
ResNet-152 (ImageNet)	<b>64.97<math>\pm</math>0.06</b> <sub>(+0.24)</sub>	<b>71.18<math>\pm</math>0.07</b> <sub>(+1.52)</sub>	<b>57.68<math>\pm</math>0.03</b> <sub>(+1.39)</sub>	<b>64.61<math>\pm</math>0.05</b> <sub>(+1.05)</sub>

Table 9: Comparison of different knowledge acquisition methods on CIFAR100 and VOC07+12 datasets. The 1<sup>st</sup>/2<sup>nd</sup> best results are indicated in red/blue.

Method	CLS (%) $\uparrow$	DET (%) $\uparrow$	AVG (%) $\uparrow$
Data (Direct Experience)	69.84 $\pm$ 0.07	82.96 $\pm$ 0.11	76.40 $\pm$ 0.09
Data-augmentation (Direct Experience)	70.68 $\pm$ 0.09	83.89 $\pm$ 0.07	77.29 $\pm$ 0.08
General-Model (Indirect Experience)	70.95 $\pm$ 0.07	83.91 $\pm$ 0.08	77.43 $\pm$ 0.08
Specialist-Model (Indirect Experience)	<b>71.13<math>\pm</math>0.15</b>	<b>84.01<math>\pm</math>0.12</b>	<b>77.57<math>\pm</math>0.14</b>
DISC (Direct and Indirect Experience)	<b>72.28<math>\pm</math>0.08</b> <sub>(+1.15)</sub>	<b>84.92<math>\pm</math>0.04</b> <sub>(+0.91)</sub>	<b>78.60<math>\pm</math>0.06</b> <sub>(+1.03)</sub>

Table 10: Comparison of performance across different tasks on CIFAR10, Food-101 and WIDER FACE(E, M, H) datasets. The 1<sup>st</sup>/2<sup>nd</sup> best results are indicated in red/blue.

Method	CLS (CIFAR10) $\uparrow$	CLS (Food-101) $\uparrow$	DET (WIDER FACE (E)) $\uparrow$	DET (WIDER FACE (M)) $\uparrow$	DET (WIDER FACE (H)) $\uparrow$	AVG (%) $\uparrow$
LSKD(CLS)+MTL	89.03 $\pm$ 0.05	77.16 $\pm$ 0.17	66.28 $\pm$ 0.11	66.53 $\pm$ 0.16	50.22 $\pm$ 0.14	69.84 $\pm$ 0.13
LSKD(DET)+MTL	68.26 $\pm$ 0.09	56.02 $\pm$ 0.15	88.04 $\pm$ 0.16	87.15 $\pm$ 0.12	71.18 $\pm$ 0.11	74.13 $\pm$ 0.13
CrossKD(CLS)+MTL	<b>89.72<math>\pm</math>0.08</b>	<b>78.01<math>\pm</math>0.17</b>	68.15 $\pm$ 0.12	67.31 $\pm$ 0.17	51.25 $\pm$ 0.14	70.89 $\pm$ 0.14
CrossKD(DET)+MTL	69.04 $\pm$ 0.05	57.28 $\pm$ 0.12	<b>88.92<math>\pm</math>0.16</b>	<b>87.99<math>\pm</math>0.12</b>	<b>71.78<math>\pm</math>0.11</b>	<b>75.00<math>\pm</math>0.11</b>
PPAL(CLS)+MTL	89.15 $\pm$ 0.09	77.35 $\pm$ 0.13	67.24 $\pm$ 0.16	67.05 $\pm$ 0.11	50.72 $\pm$ 0.17	70.30 $\pm$ 0.13
PPAL(DET)+MTL	68.62 $\pm$ 0.07	56.63 $\pm$ 0.17	88.31 $\pm$ 0.12	87.63 $\pm$ 0.15	71.53 $\pm$ 0.16	74.54 $\pm$ 0.13
DISC	<b>90.98<math>\pm</math>0.05</b> <sub>(+1.26)</sub>	<b>79.25<math>\pm</math>0.12</b> <sub>(+1.24)</sub>	<b>90.82<math>\pm</math>0.15</b> <sub>(+1.90)</sub>	<b>89.35<math>\pm</math>0.13</b> <sub>(+1.36)</sub>	<b>73.02<math>\pm</math>0.16</b> <sub>(+1.24)</sub>	<b>84.68<math>\pm</math>0.12</b> <sub>(+9.68)</sub>

Table 11: Comparison of computational cost.

Method	FLOPs (G)	Params (M)
LSKD(CLS)	0.14	0.47
LSKD(DET)	52.96	22.59
CrossKD(CLS)	0.16	26.91
CrossKD(DET)	83.32	32.22
PPAL(CLS)	0.17	29.76
PPAL(DET)	85.05	36.45
MAD(CLS)	0.26	90.62
MAD(DET)	194.39	107.85
Pix2SeqV2(CLS)	0.29	116.89
Pix2SeqV2(DET)	204.81	132.62
DISC(CLS)	0.24	67.69
DISC(DET)	152.47	85.42

*Proof.* (Dynamic over static) Let  $D_{\text{train}} = \{x_i, y_i\}_{i=1}^N$  be a training dataset of  $N$  samples, and  $\hat{\mathbb{E}}(f^m)$  be the empirical error of the  $m$ -th model  $f^m$  on  $D_{\text{train}}$ . For any hypothesis  $f \in \mathcal{H}$  (i.e.,  $\mathcal{H} : \mathcal{X} \rightarrow \{-1, 1\}$ ), with probability at least  $1 - \delta$ , the generalization error is bounded by:

$$\text{GError}(f) \leq \underbrace{\sum_{m=1}^M \mathbb{E}(w^m) \hat{\mathbb{E}}(f^m)}_{\text{Term-L (empirical loss)}} + \underbrace{\sum_{m=1}^M \mathbb{E}(w^m) \mathfrak{R}_m(f^m)}_{\text{Term-C (complexity)}} + \underbrace{\sum_{m=1}^M \text{Cov}(w^m, \ell^m)}_{\text{Term-Cov (covariance)}} + M \sqrt{\frac{\ln(1/\delta)}{2N}}, \quad (1)$$

where  $\mathbb{E}(w^m)$  is the expected collaboration weight,  $\mathfrak{R}_m(f^m)$  is the Rademacher complexity of model  $f^m$ , and  $\text{Cov}(w^m, \ell^m)$  is the covariance between the weight and the loss.

In *static collaboration*, the weights  $w_{\text{static}}^m$  are constant, hence

$$\text{Cov}(w_{\text{static}}^m, \ell^m) = 0. \quad (2)$$

In *dynamic collaboration*, the collaboration weight  $w_{\text{dynamic}}^m$  increases as the model loss  $\ell^m$  decreases. Thus,

$$\text{Cov}(w_{\text{dynamic}}^m, \ell^m) < 0, \quad (3)$$

which effectively reduces the Term-Cov and thereby lowers the generalization bound.

Based on the principle of convexity, it can be concluded that:

$$\sum \mathbb{E}(w_{\text{dynamic}}^m) \hat{\mathbb{E}}(f^m) = \sum w_{\text{static}}^m \hat{\mathbb{E}}(f^m). \quad (4)$$

$$\sum \mathbb{E}(w_{\text{dynamic}}^m) \mathfrak{R}_m(f^m) \leq \sum w_{\text{static}}^m \mathfrak{R}_m(f^m). \quad (5)$$

Since the confidence term  $M\sqrt{\frac{\ln(1/\delta)}{2N}}$  is independent of the collaboration strategy, it remains the same. Therefore, suppose the hypothesis space is  $\mathcal{H} : \mathcal{X} \rightarrow \{-1, 1\}$ . Then for any  $f_{\text{dynamic}}, f_{\text{static}} \in \mathcal{H}$ , and for  $1 > \delta > 0$ , it holds that

$$\mathcal{O}(\text{GError}_{\text{dynamic}}) \leq \mathcal{O}(\text{GError}_{\text{static}}), \quad (6)$$

proving that dynamic collaboration yields a tighter generalization error bound than static collaboration.  $\square$

*Proof.* (Dynamic weights depend on both itself and the collaborators) Given the dynamic collaboration formula delineated in  $f(x) = \sum_{m=1}^{|\mathcal{M}|} \omega^m \cdot f^m(x^m)$ , consider  $\ell$  to be the convex logistic loss function. The softmax function is utilized to normalize  $\omega^m = \frac{e^{w^m}}{\sum_{j=1}^{|\mathcal{M}|} e^{w^j}}$ .

Considering the property of convex function, we have:

$$\ell(f(x), y) = \ell\left(\sum_{m=1}^{|\mathcal{M}|} \omega^m f^m(x^m), y\right) \leq \sum_{m=1}^{|\mathcal{M}|} \omega^m \ell(f^m(x^m), y). \quad (7)$$

When computing the expectation of Equation (7) and leveraging the properties of expectation, the subsequent equation is satisfied. To simplify notation,  $\ell(f^m), y$  can be denoted as  $\ell^m$  and  $D$  is an unknown dataset:

$$\begin{aligned} GE(f) &= \mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(f(x), y) \\ &\leq \sum_{m=1}^{|\mathcal{M}|} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^m \ell^m] \\ &= \frac{1}{|\mathcal{M}|} \left( |\mathcal{M}| \cdot \sum_{m=1}^{|\mathcal{M}|} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^m \ell^m] \right) \\ &= \frac{1}{|\mathcal{M}|} \cdot \sum_{m=1}^{|\mathcal{M}|} \left( \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^m \ell^m] + (|\mathcal{M}| - 1) \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \left( 1 - \sum_{j \neq m} \omega^j \right) \ell^m \right] \right) \\ &= \frac{1}{|\mathcal{M}|} \cdot \sum_{m=1}^{|\mathcal{M}|} \left( \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^m] \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^m] + \text{Cov}(\omega^m, \ell^m) \right. \\ &\quad \left. - (|\mathcal{M}| - 1) \cdot \sum_{j \neq m} (\mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^j] \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^m] + \text{Cov}(\omega^j, \ell^m) - \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^m]) \right) \\ &= \frac{1}{|\mathcal{M}|} \cdot \sum_{m=1}^{|\mathcal{M}|} \left[ \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^m] \left( \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^m] + (|\mathcal{M}| - 1) \cdot \left[ 1 - \sum_{j \neq m} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^j] \right] \right) \right. \\ &\quad \left. + \text{Cov}(\omega^m, \ell^m) - (|\mathcal{M}| - 1) \sum_{j \neq m} \text{Cov}(\omega^j, \ell^m) \right] \\ &= \frac{1}{|\mathcal{M}|} \cdot \sum_{m=1}^{|\mathcal{M}|} \left[ |\mathcal{M}| \cdot \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^m] \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^m] + \text{Cov}(\omega^m, \ell^m) - (|\mathcal{M}| - 1) \sum_{j \neq m} \text{Cov}(\omega^j, \ell^m) \right] \\ &= \sum_{m=1}^{|\mathcal{M}|} \left( \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^m] \mathbb{E}_{(x,y) \sim \mathcal{D}} [\omega^m] + \frac{1}{|\mathcal{M}|} \cdot \left[ \text{Cov}(\omega^m, \ell^m) - (|\mathcal{M}| - 1) \sum_{j \neq m} \text{Cov}(\omega^m, \ell^j) \right] \right) \\ &\leq \sum_{m=1}^{|\mathcal{M}|} \left( \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^m] + \frac{1}{|\mathcal{M}|} \text{Cov}(\omega^m, \ell^m) - \frac{|\mathcal{M}| - 1}{|\mathcal{M}|} \sum_{j \neq m} \text{Cov}(\omega^m, \ell^j) \right) \end{aligned} \quad (8)$$

---

To simplify Equation (8), we invoke Rademacher complexity theory, which establishes that with a confidence level of  $1 - \Delta$  where  $0 < \Delta < 1$ , the following holds:

$$\mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell^m] \leq \hat{err}[f^m] + R_N(\mathcal{H}) + \sqrt{\frac{\ln(1/\Delta)}{2N}}. \quad (9)$$

In this context,  $\hat{err}(f^m)$  represents the empirical error of the unimodal function  $f^m$ , and  $\mathcal{H}$  denotes the hypothesis set, defined as  $\mathcal{H} : X \rightarrow \{-1, +1\}$ , which includes  $f$  as a member. The Rademacher complexity is denoted by  $R_N(\mathcal{H})$ . Consequently, we assert that with a confidence level of  $1 - \Delta$ , where  $0 < \Delta < 1$ , the following relationship is upheld:

$$\begin{aligned} GE(f) \leq & |\mathcal{M}| \left( R_N(\mathcal{H}) + \sqrt{\frac{\ln(1/\Delta)}{2N}} \right) + \sum_{m=1}^{|\mathcal{M}|} \hat{err}(f^m) \\ & + \sum_{m=1}^{|\mathcal{M}|} \left[ \frac{1}{|\mathcal{M}|} \underbrace{Cov(\omega^m, \ell^m)}_{\text{Negative self-correlation in DSC}} - \frac{|\mathcal{M}| - 1}{|\mathcal{M}|} \sum_{j \neq m} \underbrace{Cov(\omega^m, \ell^j)}_{\text{Positive collaborative correlation in DSC}} \right], \end{aligned} \quad (10)$$

proving that dynamic weights depend on both itself and the collaborator, and that dynamic collaboration leads to a tighter generalization error bound.

□