

# Overcoming Multi-model Forgetting Supplementary Material

## A. Proofs

**Lemma 1.** Given a dataset  $\mathcal{D}$  and two architectures with shared parameters  $\theta_s$  and private parameters  $\theta_1$  and  $\theta_2$ , and provided that  $p(\theta_1, \theta_2 \mid \theta_s, \mathcal{D}) = p(\theta_1 \mid \theta_s, \mathcal{D})p(\theta_2 \mid \theta_s, \mathcal{D})$ , we have

$$p(\theta_1, \theta_2, \theta_s \mid \mathcal{D}) \propto \frac{p(\mathcal{D} \mid \theta_2, \theta_s)p(\theta_1, \theta_s \mid \mathcal{D})p(\theta_2, \theta_s)}{\int p(\mathcal{D} \mid \theta_1, \theta_s)p(\theta_1, \theta_s)d\theta_1}. \quad (1)$$

*Proof.* Using Bayes' theorem and ignoring constants, we have

$$\begin{aligned} p(\theta \mid \mathcal{D}) &= \frac{p(\theta_1, \theta_2, \theta_s, \mathcal{D})}{p(\mathcal{D})} \\ &\propto p(\theta_1 \mid \theta_2, \theta_s, \mathcal{D})p(\theta_2, \theta_s, \mathcal{D}) \\ &= p(\theta_1 \mid \theta_s, \mathcal{D})p(\mathcal{D} \mid \theta_2, \theta_s)p(\theta_2, \theta_s) \\ &\propto \frac{p(\theta_1, \theta_s, \mathcal{D})p(\mathcal{D} \mid \theta_2, \theta_s)p(\theta_2, \theta_s)}{p(\mathcal{D}, \theta_s)} \\ &\propto \frac{p(\theta_1, \theta_s, \mathcal{D})p(\mathcal{D} \mid \theta_2, \theta_s)p(\theta_2, \theta_s)}{\int p(\mathcal{D} \mid \theta_1, \theta_s)p(\theta_s, \theta_1)d\theta_1} \\ &\propto \frac{p(\theta_1, \theta_s \mid \mathcal{D})p(\mathcal{D} \mid \theta_2, \theta_s)p(\theta_2, \theta_s)}{\int p(\mathcal{D} \mid \theta_1, \theta_s)p(\theta_s, \theta_1)d\theta_1}, \end{aligned}$$

where we used the conditional independence assumption  $p(\theta_1 \mid \theta_2, \theta_s, \mathcal{D}) = p(\theta_1 \mid \theta_s, \mathcal{D})$  in the third line.  $\square$

We now derive a closed-form expression for the denominator of equation (1).

**Lemma 2.** Suppose we have the maximum likelihood estimate  $(\hat{\theta}_1, \hat{\theta}_s)$  for the first model, write  $\text{Card}(\theta_1) + \text{Card}(\theta_s) = p_1 + p_s = p$ , and let the negative Hessian  $\mathbf{H}_p(\hat{\theta}_1, \hat{\theta}_s)$  of the log posterior probability distribution  $\log p(\theta_1, \theta_s \mid \mathcal{D})$  evaluated at  $(\hat{\theta}_1, \hat{\theta}_s)$  be partitioned into four blocks corresponding to  $(\theta_1, \theta_s)$  as

$$\mathbf{H}_p(\hat{\theta}_1, \hat{\theta}_s) = \begin{bmatrix} \mathbf{H}_{11} & \mathbf{H}_{1s} \\ \mathbf{H}_{s1} & \mathbf{H}_{ss} \end{bmatrix}.$$

If the parameters of each model follow Normal distributions, i.e.,  $(\theta_1, \theta_s) \sim \mathcal{N}_p(\mathbf{0}, \sigma^2 \mathbf{I}_p)$ , with  $\mathbf{I}_p$  the  $p$ -dimensional identity matrix, then the denominator of equation (1),  $A =$

$\int p(\mathcal{D} \mid \theta_1, \theta_s)p(\theta_s, \theta_1)d\theta_1$  can be written as

$$A = \exp \left\{ l_p(\hat{\theta}_1, \hat{\theta}_s) - \frac{1}{2} \mathbf{v}^\top \boldsymbol{\Omega} \mathbf{v} \right\} \times (2\pi)^{p_1/2} |\det(\mathbf{H}_{11}^{-1})|^{1/2}, \quad (2)$$

where  $\mathbf{v} = \theta_s - \hat{\theta}_s$  and  $\boldsymbol{\Omega} = \mathbf{H}_{ss} - \mathbf{H}_{1s}^\top \mathbf{H}_{11}^{-1} \mathbf{H}_{1s}$ .

*Proof.* We have

$$\begin{aligned} p(\mathcal{D} \mid \theta_1, \theta_s)p(\theta_s, \theta_1) &\propto e^{l(\theta_1, \theta_s) - (\theta_1, \theta_s)^\top (\theta_1, \theta_s) / 2\sigma^2} \\ &\propto e^{l_p(\theta_1, \theta_s)}, \end{aligned}$$

where  $l(\theta_1, \theta_s) = \log p(\mathcal{D} \mid \theta_1, \theta_s)$ , and  $l_p(\theta_1, \theta_s) = l(\theta_1, \theta_s) - (\theta_1, \theta_s)^\top (\theta_1, \theta_s) / 2\sigma^2$ .

Let  $\mathbf{H}_p(\theta_1, \theta_s) = \mathbf{H}(\theta_1, \theta_s) + \sigma^{-2} \mathbf{I}_p$  be the negative Hessian of  $l_p(\theta_1, \theta_s)$ , with  $\mathbf{I}_p$  the  $p$ -dimensional identity matrix and  $\mathbf{H}(\theta_1, \theta_s)$  the negative Hessian of  $l(\theta_1, \theta_s)$ .

Using the second-order Taylor expansion of  $l_p(\theta_1, \theta_s)$  around its maximum likelihood estimate  $(\hat{\theta}_1, \hat{\theta}_s)$ , we have

$$l_p(\theta_1, \theta_s) = l_p(\hat{\theta}_1, \hat{\theta}_s) - \frac{1}{2} (\theta_1', \theta_s')^\top \mathbf{H}_p(\hat{\theta}_1, \hat{\theta}_s) (\theta_1', \theta_s'); \quad (3)$$

where  $(\theta_1', \theta_s') = (\theta_1, \theta_s) - (\hat{\theta}_1, \hat{\theta}_s)$ . The first derivative is zero since it is evaluated at the maximum likelihood estimate. We now partition our negative Hessian matrix as

$$\mathbf{H}_p(\hat{\theta}_1, \hat{\theta}_s) = \begin{bmatrix} \mathbf{H}_{11} & \mathbf{H}_{1s} \\ \mathbf{H}_{s1} & \mathbf{H}_{ss} \end{bmatrix},$$

which gives

$$\begin{aligned} \mathbf{B} &= [(\theta_1, \theta_s) - (\hat{\theta}_1, \hat{\theta}_s)]^\top \mathbf{H}_p(\hat{\theta}_1, \hat{\theta}_s) [(\theta_1, \theta_s) - (\hat{\theta}_1, \hat{\theta}_s)] \\ &= (\theta_1 - \hat{\theta}_1)^\top \mathbf{H}_{11} (\theta_1 - \hat{\theta}_1) + (\theta_s - \hat{\theta}_s)^\top \mathbf{H}_{ss} (\theta_s - \hat{\theta}_s) \\ &\quad + (\theta_s - \hat{\theta}_s)^\top \mathbf{H}_{s1} (\theta_1 - \hat{\theta}_1) \\ &\quad + (\theta_1 - \hat{\theta}_1)^\top \mathbf{H}_{1s} (\theta_s - \hat{\theta}_s) \\ &= (\theta_1 - \hat{\theta}_1)^\top \mathbf{H}_{11} (\theta_1 - \hat{\theta}_1) + (\theta_s - \hat{\theta}_s)^\top \mathbf{H}_{ss} (\theta_s - \hat{\theta}_s) \\ &\quad + (\theta_1 - \hat{\theta}_1)^\top (\mathbf{H}_{1s} + \mathbf{H}_{s1}^\top) (\theta_s - \hat{\theta}_s). \end{aligned}$$

Let us define  $\mathbf{u} = \theta_1 - \hat{\theta}_1$ ,  $\mathbf{v} = \theta_s - \hat{\theta}_s$  and  $\mathbf{w} =$

$\mathbf{H}_{11}^{-1}\mathbf{H}_{1s}\mathbf{v}$ . We then have,

$$\begin{aligned}
 \mathbf{C} &= (\mathbf{u} + \mathbf{w})^T \mathbf{H}_{11} (\mathbf{u} + \mathbf{w}) \\
 &= \mathbf{u}^T \mathbf{H}_{11} \mathbf{u} + \mathbf{u}^T \mathbf{H}_{11} \mathbf{w} + \mathbf{w}^T \mathbf{H}_{11} \mathbf{u} + \mathbf{w}^T \mathbf{H}_{11} \mathbf{w} \\
 &= (\boldsymbol{\theta}_1 - \hat{\boldsymbol{\theta}}_1)^T \mathbf{H}_{11} (\boldsymbol{\theta}_1 - \hat{\boldsymbol{\theta}}_1) \\
 &\quad + (\boldsymbol{\theta}_1 - \hat{\boldsymbol{\theta}}_1)^T \mathbf{H}_{11} \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} (\boldsymbol{\theta}_s - \hat{\boldsymbol{\theta}}_s) \\
 &\quad + \mathbf{v}^T \mathbf{H}_{1s}^T \mathbf{H}_{11}^{-1} \mathbf{H}_{11} \mathbf{H}_{1s} \mathbf{v} \\
 &\quad + \mathbf{v}^T \mathbf{H}_{1s}^T \mathbf{H}_{11}^{-1} \mathbf{H}_{11} (\boldsymbol{\theta}_1 - \hat{\boldsymbol{\theta}}_1) \\
 &= \mathbf{B} - \mathbf{v}^T \mathbf{H}_{ss} \mathbf{v} + \mathbf{v}^T \mathbf{H}_{1s}^T \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} \mathbf{v} \\
 &= \mathbf{B} - \mathbf{v}^T (\mathbf{H}_{ss} - \mathbf{H}_{1s}^T \mathbf{H}_{11}^{-1} \mathbf{H}_{1s}) \mathbf{v} \\
 &= \mathbf{B} - \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v},
 \end{aligned}$$

with  $\boldsymbol{\Omega} = \mathbf{H}_{ss} - \mathbf{H}_{1s}^T \mathbf{H}_{11}^{-1} \mathbf{H}_{1s}$ .

Thus

$$\mathbf{B} = (\mathbf{u} + \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} \mathbf{v})^T \mathbf{H}_{11} (\mathbf{u} + \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} \mathbf{v}) + \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v}. \quad (4)$$

Given equation (4), we are now able to prove Lemma 2, as integral

$$\begin{aligned}
 D &= \int e^{l_p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_s)} d\boldsymbol{\theta}_1 = \int e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s) - \frac{1}{2} \mathbf{B}} d\boldsymbol{\theta}_1 \\
 &= \int e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s)} e^{-\frac{1}{2} \mathbf{B}} d\boldsymbol{\theta}_1 = e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s)} \int e^{-\frac{1}{2} \mathbf{B}} d\boldsymbol{\theta}_1 \\
 &= \int e^{-\frac{1}{2} ((\mathbf{u} + \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} \mathbf{v})^T \mathbf{H}_{11} (\mathbf{u} + \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} \mathbf{v}) + \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v})} d\boldsymbol{\theta}_1 \\
 &\quad \times e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s)} \\
 &= \int e^{-\frac{1}{2} ((\mathbf{u} + \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} \mathbf{v})^T \mathbf{H}_{11} (\mathbf{u} + \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} \mathbf{v}))} e^{-\frac{1}{2} \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v}} d\boldsymbol{\theta}_1 \\
 &\quad \times e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s)} \\
 &= e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s) - \frac{1}{2} \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v}} \int e^{-\frac{1}{2} (\boldsymbol{\theta}_1 - \mathbf{z})^T \mathbf{H}_{11} (\boldsymbol{\theta}_1 - \mathbf{z})} d\boldsymbol{\theta}_1 \\
 &= e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s) - \frac{1}{2} \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v}} (2\pi)^{\frac{p_1}{2}} |\det(\mathbf{H}_{11}^{-1})|^{\frac{1}{2}} \\
 &\quad \times (2\pi)^{-\frac{p_1}{2}} |\det(\mathbf{H}_{11}^{-1})|^{-\frac{1}{2}} \\
 &\quad \times \int e^{-\frac{1}{2} (\boldsymbol{\theta}_1 - \mathbf{z})^T \mathbf{H}_{11} (\boldsymbol{\theta}_1 - \mathbf{z})} d\boldsymbol{\theta}_1 \\
 &= e^{l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s) - \frac{1}{2} \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v}} (2\pi)^{\frac{p_1}{2}} |\det(\mathbf{H}_{11}^{-1})|^{\frac{1}{2}},
 \end{aligned}$$

where we re-arranged the terms so that the integral is over a normal distribution with mean  $\mathbf{z} = \hat{\boldsymbol{\theta}}_1 - \mathbf{H}_{11}^{-1} \mathbf{H}_{1s} (\boldsymbol{\theta}_s - \hat{\boldsymbol{\theta}}_s)$  and covariance matrix  $\mathbf{H}_{11}^{-1}$ , which can be computed in closed form.  $\square$

From Lemma 1 and Lemma 2, we can obtain equation 3 by replacing the denominator with the closed form above and

taking the log on both size of equation (1). This yields

$$\begin{aligned}
 \log p(\boldsymbol{\theta} | \mathcal{D}) &\propto \log p(\mathcal{D} | \boldsymbol{\theta}_2, \boldsymbol{\theta}_s) + \log p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_s) \\
 &\quad + \log p(\boldsymbol{\theta}_2, \boldsymbol{\theta}_s) \\
 &\quad - \log \left\{ \int p(\mathcal{D} | \boldsymbol{\theta}_1, \boldsymbol{\theta}_s) p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_s) d\boldsymbol{\theta}_1 \right\} \\
 &= \log p(\mathcal{D} | \boldsymbol{\theta}_2, \boldsymbol{\theta}_s) + \log p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_s) \\
 &\quad + \log p(\boldsymbol{\theta}_2, \boldsymbol{\theta}_s) - l_p(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_s) + \frac{1}{2} \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v} \\
 &\quad + \log \left\{ (2\pi)^{\frac{p_1}{2}} |\det(\mathbf{H}_{11}^{-1})|^{\frac{1}{2}} \right\} \\
 &\propto \log p(\mathcal{D} | \boldsymbol{\theta}_2, \boldsymbol{\theta}_s) + \log p(\boldsymbol{\theta}_2, \boldsymbol{\theta}_s) \\
 &\quad + \log p(\boldsymbol{\theta}_1, \boldsymbol{\theta}_s | \mathcal{D}) + \frac{1}{2} \mathbf{v}^T \boldsymbol{\Omega} \mathbf{v}.
 \end{aligned}$$

## B. Plots for CNN Search

In our CNN search experiment, we search for a ‘‘micro’’ cell as in (Pham et al., 2018). We employ the hyper-parameters available in the released ENAS code. The plots depicting error difference as a function of training epochs as provided in Figure 1 (a), (b) and (c). Note that here again the original ENAS is subject to multi-model forgetting, and our WPL helps reducing it. In Figure 1 (d), we show the mean reward as training progresses. While the shape of the reward curve is different from the RNN case, because of a different formulation of the reward function, the general trend is the same; Our approach initially produces lower rewards, but is better at maintaining good models until the end of the search, as indicated by higher rewards in the second half of training.

## C. Best architectures found by the search

In Figure 2, we show the best architectures found by our neural architecture search for the RNN and CNN cases.

## References

Pham, H., Guan, M. Y., Zoph, B., Le, Q. V., and Dean, J. Efficient Neural Architecture Search via Parameter Sharing. *International Conference on Machine Learning (ICML)*, 2018.

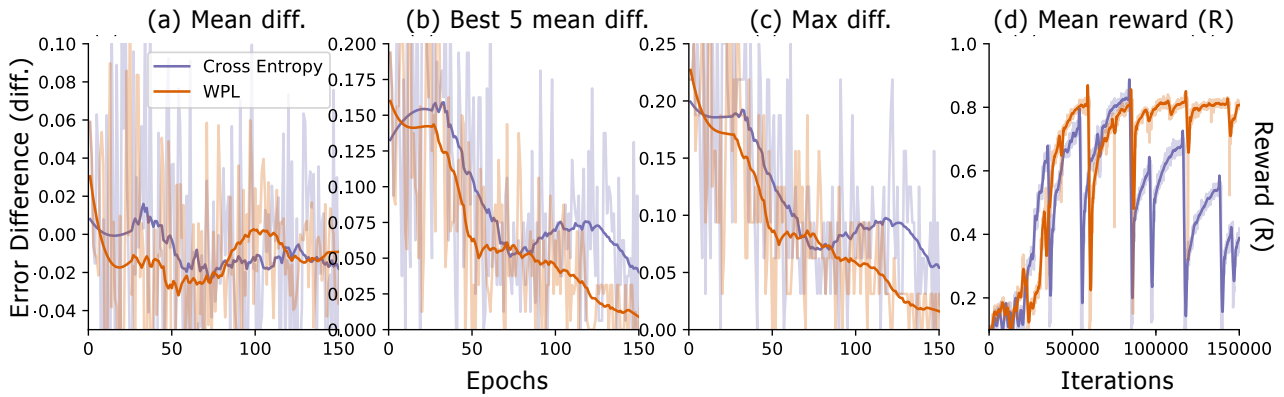


Figure 1. Error differences when searching for CNN architectures. Quantitatively, the multi-model forgetting effect is reduced by up to 99% for (a), 96% for (b), and 98% for (c).

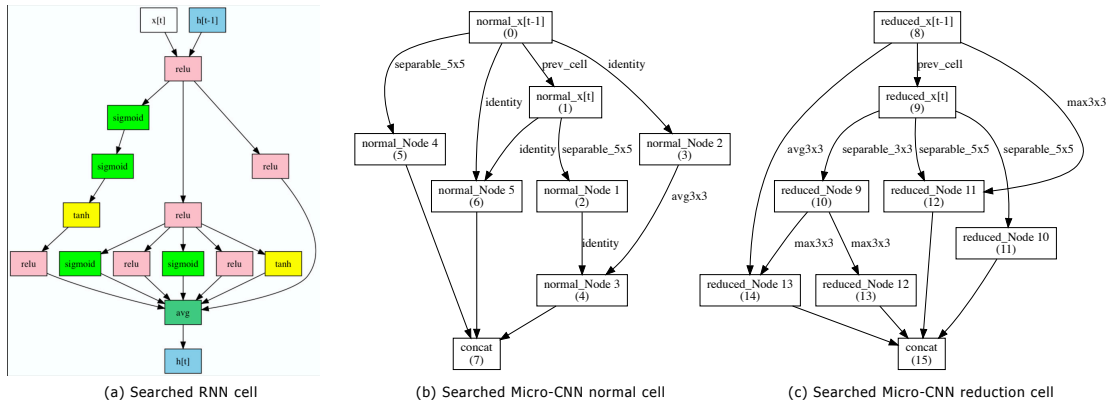


Figure 2. Best architectures found for RNN and CNN. We display the best architecture found by ENAS+WPL, in (a) for the RNN cell, and in (b) and (c) for the CNN normal and reduction cells.