

A Review from Neurips 2020

Insufficient evaluation for the task selection, e.g. adding a baseline, varying the number of ways and the number of training examples. Originally, we only ran the experiment shown in Figure 3 without the baseline Task2Vec and fixing the setting for 5-way 5-shot. In the current version, we add more experiments as shown in Figure 4 with several variations in term of number of classes within a task, and the number of training tasks. The baseline Task2Vec is also included in the comparison.

The proposed approach adds an inefficient step. The proposed approach is not about meta-learning, but to learn a generative model to represent tasks in an embedding space. The obtained representation can be used either for task similarity (as presented in the paper) or further downstream transfer-learning works, such as task augmentation. Here, we use it as a pre-processing step before performing meta-learning for demonstration purpose without claiming its efficiency. If considering such applications, other baselines, such as Task2Vec, also share the same computational burden when calculating task-to-task distances.

Is the proposed approach restricted to gradient based meta-learning algorithms? Meta-learning is used to demonstrate the concept of task similarity quantification proposed in our paper, so our method is not restricted to any particular meta-learning algorithm. To demonstrate that, the submitted paper shows the results using the metric-based Prototypical Networks (we replicate the paper results in Fig. 4b and add the 10-way result in Fig. 4d). The extension to other meta-learning algorithms is trivial.

Experimenting on only 2 data sets. The reason for using those two data sets is that the our proposed approach works with low-dimensional data (Omniglot) or extracted features (mini-ImageNet). We will demonstrate our method on other data sets by embedding images into a lower-dimensional embedding space before performing LDCC. The training at that stage will be to learn both the image embedding and LDCC simultaneously.

B Detail derivation of each term in ELBO

B.1 $\mathbb{E}_q [\ln p(\mathbf{z}|\boldsymbol{\theta})]$

$$\begin{aligned}
 \mathbb{E}_q [\ln p(\mathbf{z}|\boldsymbol{\theta})] &= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K q(z_{dcnk} = 1; \mathbf{r}_{dcn}) \int q(\boldsymbol{\theta}_{dc}; \boldsymbol{\gamma}_{dc}) \ln p(z_{dcnk} = 1 | \boldsymbol{\theta}_{dc}) d\boldsymbol{\theta}_{dc} \\
 &= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K r_{dcnk} \int \text{Dir}_K(\boldsymbol{\theta}_{dc}; \boldsymbol{\gamma}_{dc}) \ln \theta_{dck} d\boldsymbol{\theta}_{dc} \\
 &= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K r_{dcnk} \ln \tilde{\theta}_{dck},
 \end{aligned} \tag{17}$$

where:

$$\ln \tilde{\theta}_{dck} = \psi(\gamma_{dck}) - \psi\left(\sum_{j=1}^K \gamma_{dcj}\right). \tag{18}$$

B.2 $\mathbb{E}_q [\ln p(\boldsymbol{\theta}|\mathbf{y}, \boldsymbol{\alpha})]$

$$\begin{aligned}
 \mathbb{E}_q [\ln p(\boldsymbol{\theta}|\mathbf{y}, \boldsymbol{\alpha})] &= \sum_{d=1}^M \sum_{c=1}^C \sum_{l=1}^L q(y_{dcl} = 1; \boldsymbol{\eta}_{dc}) \int q(\boldsymbol{\theta}_{dc}; \boldsymbol{\gamma}_{dc}) \ln p(\boldsymbol{\theta}_{dc} | y_{dcl} = 1, \boldsymbol{\alpha}) d\boldsymbol{\theta}_{dc} \\
 &= \sum_{d=1}^M \sum_{c=1}^C \sum_{l=1}^L \eta_{dcl} \int \text{Dir}_K(\boldsymbol{\theta}_{dc}; \boldsymbol{\gamma}_{dc}) \ln \text{Dir}_K(\boldsymbol{\theta}_{dc}; \boldsymbol{\alpha}_l) d\boldsymbol{\theta}_{dc}.
 \end{aligned} \tag{19}$$

408 Note that the cross-entropy between 2 Dirichlet distributions can be expressed as:

$$\begin{aligned}
\mathcal{H}[\text{Dir}(\mathbf{x}; \boldsymbol{\alpha}_0), \text{Dir}(\mathbf{x}; \boldsymbol{\alpha}_1)] &= -\mathbb{E}_{\text{Dir}(\mathbf{x}; \boldsymbol{\alpha}_0)} [\ln \text{Dir}(\mathbf{x}; \boldsymbol{\alpha}_1)] \\
&= -\mathbb{E}_{\text{Dir}(\mathbf{x}; \boldsymbol{\alpha}_0)} \left[-\ln B(\boldsymbol{\alpha}_1) + \sum_{k=1}^K (\alpha_{1k} - 1) \ln x_k \right] \\
&= \ln B(\boldsymbol{\alpha}_1) - \sum_{k=1}^K (\alpha_{1k} - 1) \left[\psi(\alpha_{0k}) - \psi \left(\sum_{k'=1}^K \alpha_{0k'} \right) \right], \quad (20)
\end{aligned}$$

409 where:

$$\ln B(\boldsymbol{\alpha}_1) = \sum_{k=1}^K \ln \Gamma(\alpha_{1k}) - \ln \Gamma \left(\sum_{j=1}^K \alpha_{1j} \right). \quad (21)$$

410 Hence:

$$\mathbb{E}_q [\ln p(\boldsymbol{\theta}|\mathbf{y}, \boldsymbol{\alpha})] = \sum_{d=1}^M \sum_{c=1}^C \sum_{l=1}^L \eta_{dcl} \left[-\ln B(\boldsymbol{\alpha}_l) + \sum_{k=1}^K (\alpha_{lk} - 1) \ln \tilde{\theta}_{dck} \right]. \quad (22)$$

411 **B.3** $\mathbb{E}_q [\ln p(\mathbf{y}|\boldsymbol{\phi})]$

$$\begin{aligned}
\mathbb{E}_q [\ln p(\mathbf{y}|\boldsymbol{\phi})] &= \sum_{d=1}^M \sum_{c=1}^C \sum_{l=1}^L q(y_{dcl} = 1; \boldsymbol{\eta}_{dc}) \int q(\boldsymbol{\phi}_d; \boldsymbol{\lambda}_d) \ln p(y_{dcl} = 1|\boldsymbol{\phi}_d) d\boldsymbol{\phi}_d \\
&= \sum_{d=1}^M \sum_{c=1}^C \sum_{l=1}^L \eta_{dcl} \int \text{Dir}(\boldsymbol{\phi}_d; \boldsymbol{\lambda}_d) \ln \phi_{dl} d\boldsymbol{\phi}_d \\
&= \sum_{d=1}^M \sum_{c=1}^C \sum_{l=1}^L \eta_{dcl} \ln \tilde{\phi}_{dl}, \quad (23)
\end{aligned}$$

412 where:

$$\ln \tilde{\phi}_{dl} = \psi(\lambda_{dl}) - \psi \left(\sum_{j=1}^K \lambda_{dj} \right) \quad (24)$$

413 **B.4** $\mathbb{E}_q [\ln p(\boldsymbol{\phi}|\boldsymbol{\delta})]$

$$\begin{aligned}
\mathbb{E}_q [\ln p(\boldsymbol{\phi}|\boldsymbol{\delta})] &= \sum_{d=1}^M \int q(\boldsymbol{\phi}_d; \boldsymbol{\lambda}_d) \ln p(\boldsymbol{\phi}_d|\boldsymbol{\delta}) d\boldsymbol{\phi}_d \\
&= \sum_{d=1}^M \int \text{Dir}_L(\boldsymbol{\phi}_d; \boldsymbol{\lambda}_d) \ln \text{Dir}_L(\boldsymbol{\phi}_d|\boldsymbol{\delta}) d\boldsymbol{\phi}_d \\
&= \sum_{d=1}^M \ln \frac{\Gamma \left(\sum_{l=1}^L \delta_l \right)}{\prod_{l=1}^L \Gamma(\delta_l)} + \sum_{l=1}^L (\delta_l - 1) \ln \tilde{\phi}_{dl}. \quad (25)
\end{aligned}$$

414 **B.5** $\mathbb{E}_q [\ln p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})]$

$$\begin{aligned}
& \mathbb{E}_q [\ln p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] \\
&= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K q(z_{dcnk} = 1; \mathbf{r}_{dcn}) \\
&\quad \times \int q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k; \mathbf{m}_k, \kappa_k, \mathbf{W}_k, \nu_k) \ln p(\mathbf{x}_{dcn} | z_{dcnk} = 1, \boldsymbol{\mu}, \boldsymbol{\Lambda}) d\boldsymbol{\mu} d\boldsymbol{\Lambda} \\
&= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K r_{dcnk} \\
&\quad \times \int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \ln \mathcal{N}(\mathbf{x}_{dcn} | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k \\
&= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K \frac{r_{dcnk}}{2} \int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \times \\
&\quad \times [\ln |\boldsymbol{\Lambda}_k| - (\mathbf{x}_{dcn} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \boldsymbol{\mu}_k) - D \ln(2\pi)] d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k \tag{26}
\end{aligned}$$

415 Note that:

$$\ln \tilde{\Lambda}_k = \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \ln |\boldsymbol{\Lambda}_k| d\boldsymbol{\Lambda}_k = \sum_{i=1}^D \psi\left(\frac{\nu_k - i + 1}{2}\right) + D \ln(2) + \ln |\mathbf{W}_k|, \tag{27}$$

416 and

$$\begin{aligned}
& \mathbb{E}_{\boldsymbol{\mu}_k} [(\mathbf{x}_{dcn} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \boldsymbol{\mu}_k)] \\
&= \mathbb{E}_{\boldsymbol{\mu}_k} [(\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{m}_k - \boldsymbol{\mu}_k) \\
&\quad + (\mathbf{m}_k - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + (\mathbf{m}_k - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{m}_k - \boldsymbol{\mu}_k)] \\
&= (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + \mathbb{E}_{\boldsymbol{\mu}_k} [(\mathbf{m}_k - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{m}_k - \boldsymbol{\mu}_k)] \\
&= (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + \text{tr} [\boldsymbol{\Lambda}_k \mathbb{E}_{\boldsymbol{\mu}_k} [(\mathbf{m}_k - \boldsymbol{\mu}_k)^\top (\mathbf{m}_k - \boldsymbol{\mu}_k)]] \\
&= (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + \text{tr} [\boldsymbol{\Lambda}_k (\kappa_k \boldsymbol{\Lambda}_k)^{-1}] \\
&= (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + D \kappa_k^{-1} \tag{28}
\end{aligned}$$

417

$$\begin{aligned}
\Rightarrow \mathbb{E}_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k} [(\mathbf{x}_{dcn} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \boldsymbol{\mu}_k)] &= \mathbb{E}_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k} [(\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + D \kappa_k^{-1}] \\
&= \nu_k (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \mathbf{W}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) + D \kappa_k^{-1}. \tag{29}
\end{aligned}$$

418 Hence:

$$\begin{aligned}
\mathbb{E}_q [\ln p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] &= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K \frac{r_{dcnk}}{2} \left[\ln \tilde{\Lambda}_k - D \kappa_k^{-1} - D \ln(2\pi) - \right. \\
&\quad \left. - \nu_k (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \mathbf{W}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) \right] \tag{30}
\end{aligned}$$

419 **B.6** $\mathbb{E}_q [\ln p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{m}_0, \kappa_0, \mathbf{W}_0, \nu_0)]$

$$\begin{aligned}
& \mathbb{E}_q [\ln p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{m}_0, \kappa_0, \mathbf{W}_0, \nu_0)] \\
&= \sum_{k=1}^K \int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \\
&\quad \times \ln \left[\mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_0, (\kappa_0 \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_0, \nu_0) \right] d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k \\
&= \sum_{k=1}^K \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \left[\int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \ln \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_0, (\kappa_0 \boldsymbol{\Lambda}_k)^{-1}) d\boldsymbol{\mu}_k \right] d\boldsymbol{\Lambda}_k \\
&\quad + \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \ln \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_0, \nu_0) d\boldsymbol{\Lambda}_k.
\end{aligned} \tag{31}$$

420 Note that:

$$\begin{aligned}
& \int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \ln \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_0, (\kappa_0 \boldsymbol{\Lambda}_k)^{-1}) d\boldsymbol{\mu}_k \\
&= -\frac{1}{2} \int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \\
&\quad \times [D \ln(2\pi) - D \ln \kappa_0 - \ln |\boldsymbol{\Lambda}_k| + \kappa_0 (\boldsymbol{\mu}_k - \mathbf{m}_0)^\top \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \mathbf{m}_0)] d\boldsymbol{\mu}_k \\
&= \frac{1}{2} \left[D \ln \left(\frac{\kappa_0}{2\pi} \right) + \ln |\boldsymbol{\Lambda}_k| - \kappa_0 (\mathbf{m}_k - \mathbf{m}_0)^\top \boldsymbol{\Lambda}_k (\mathbf{m}_k - \mathbf{m}_0) - \frac{D \kappa_0}{\kappa_k} \right]
\end{aligned} \tag{32}$$

421

$$\begin{aligned}
& \Rightarrow \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \left[\int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \ln \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_0, (\kappa_0 \boldsymbol{\Lambda}_k)^{-1}) d\boldsymbol{\mu}_k \right] d\boldsymbol{\Lambda}_k \\
&= \frac{1}{2} \left[D \ln \left(\frac{\kappa_0}{2\pi} \right) + \ln \tilde{\Lambda}_k - \nu_k \kappa_0 (\mathbf{m}_k - \mathbf{m}_0)^\top \mathbf{W}_k (\mathbf{m}_k - \mathbf{m}_0) - \frac{D \kappa_0}{\kappa_k} \right]
\end{aligned} \tag{33}$$

422 and the cross-entropy between 2 Wishart distributions can be written as:

$$\begin{aligned}
& \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \ln \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_0, \nu_0) d\boldsymbol{\Lambda}_k = -\mathcal{H}[(\mathbf{W}_k, \nu_k), (\mathbf{W}_0, \nu_0)] \\
&= \frac{\nu_0}{2} \ln |\mathbf{W}_0^{-1} \mathbf{W}_k| - \frac{D+1}{2} \ln |\mathbf{W}_k| - \frac{\nu_k}{2} \text{tr}(\mathbf{W}_0^{-1} \mathbf{W}_k) - \log \Gamma_D \left(\frac{\nu_0}{2} \right) + \\
&\quad + \frac{\nu_0 - D - 1}{2} \psi_D \left(\frac{\nu_k}{2} \right) - \frac{D(D+1)}{2} \ln(2) \\
&= -\frac{\nu_0}{2} \ln |\mathbf{W}_0| + \frac{\nu_0 - D - 1}{2} \left[\psi_D \left(\frac{\nu_k}{2} \right) + \ln |\mathbf{W}_k| + D \ln(2) \right] - \frac{\nu_k}{2} \text{tr}(\mathbf{W}_0^{-1} \mathbf{W}_k) - \\
&\quad - \log \Gamma_D \left(\frac{\nu_0}{2} \right) - \frac{\nu_0 D}{2} \ln(2) \\
&= -\frac{\nu_0}{2} \ln |\mathbf{W}_0| + \frac{\nu_0 - D - 1}{2} \ln \tilde{\Lambda}_k - \frac{\nu_k}{2} \text{tr}(\mathbf{W}_0^{-1} \mathbf{W}_k) - \log \Gamma_D \left(\frac{\nu_0}{2} \right) - \frac{\nu_0 D}{2} \ln(2), \tag{34}
\end{aligned}$$

423 where $\Gamma_D(\cdot)$ and $\psi_D(\cdot)$ are the multivariate gamma and digamma functions, respectively.

424 Hence, the expectation of interest can be expressed as:

$$\begin{aligned}
& \mathbb{E}_q [\ln p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{m}_0, \kappa_0, \mathbf{W}_0, \nu_0)] \\
&= \sum_{k=1}^K \frac{1}{2} \left[D \ln \left(\frac{\kappa_0}{2\pi} \right) + \ln \tilde{\Lambda}_k - \nu_k \kappa_0 (\mathbf{m}_k - \mathbf{m}_0)^\top \mathbf{W}_k (\mathbf{m}_k - \mathbf{m}_0) - \frac{D \kappa_0}{\kappa_k} \right] - \\
&\quad - \frac{\nu_0}{2} \ln |\mathbf{W}_0| + \frac{\nu_0 - D - 1}{2} \ln \tilde{\Lambda}_k - \frac{\nu_k}{2} \text{tr}(\mathbf{W}_0^{-1} \mathbf{W}_k) - \log \Gamma_D \left(\frac{\nu_0}{2} \right) - \frac{\nu_0 D}{2} \ln(2). \tag{35}
\end{aligned}$$

425 The result is identical to the previous derivation in the literature Bishop, 2006, Eq. (10.74).

426 **B.7** $\mathbb{E}_q [\ln q(\mathbf{z})]$

$$\mathbb{E}_q [\ln q(\mathbf{z})] = \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K r_{dcnk} \ln r_{dcnk}. \quad (36)$$

427 **B.8** $\mathbb{E}_q [\ln q(\boldsymbol{\theta})]$

$$\mathbb{E}_q [\ln q(\boldsymbol{\theta})] = - \sum_{d=1}^M \sum_{c=1}^C \ln B(\gamma_{dc}) + \sum_{j=1}^K (\gamma_{dcj} - 1) \ln \tilde{\theta}_{dcj}. \quad (37)$$

428 **B.9** $\mathbb{E}_q [\ln q(\mathbf{y})]$

$$\mathbb{E}_q [\ln q(\mathbf{y})] = \sum_{d=1}^M \sum_{c=1}^C \sum_{l=1}^L \eta_{dcl} \ln \eta_{dcl}. \quad (38)$$

429 **B.10** $\mathbb{E}_q [\ln q(\boldsymbol{\phi})]$

$$\mathbb{E}_q [\ln q(\boldsymbol{\phi})] = - \sum_{d=1}^M \ln B(\boldsymbol{\lambda}_d) - \sum_{l=1}^L (\lambda_{dl} - 1) \ln \tilde{\phi}_{dl}. \quad (39)$$

430 **B.11** $\mathbb{E}_q [\ln q(\boldsymbol{\mu}, \boldsymbol{\Lambda})]$

$$\begin{aligned} & \mathbb{E}_q [\ln q(\boldsymbol{\mu}, \boldsymbol{\Lambda})] \\ &= \sum_{k=1}^K \mathbb{E}_{q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)} [\ln q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)] \\ &= \sum_{k=1}^K \int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) [\ln \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \\ & \quad + \ln \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k)] d\boldsymbol{\mu}_k d\boldsymbol{\Lambda}_k \\ &= \sum_{k=1}^K \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \left[\int \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) \ln \mathcal{N}(\boldsymbol{\mu}_k; \mathbf{m}_k, (\kappa_k \boldsymbol{\Lambda}_k)^{-1}) d\boldsymbol{\mu}_k \right] d\boldsymbol{\Lambda}_k \\ & \quad + \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \ln \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) d\boldsymbol{\Lambda}_k \\ &= - \sum_{k=1}^K \int \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \frac{1}{2} \ln |2\pi e(\kappa_k \boldsymbol{\Lambda}_k)^{-1}| d\boldsymbol{\Lambda}_k + \mathcal{H}[\mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k), \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k)] \\ &= \sum_{k=1}^K \left\{ \frac{1}{2} \ln \tilde{\Lambda}_k + \frac{D}{2} \ln \left(\frac{\kappa_k}{2\pi} \right) - \frac{D}{2} - \mathcal{H}[\mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k), \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k)] \right\}. \end{aligned} \quad (40)$$

431 **C Optimise ELBO**

432 **C.1 Variational parameter for categorical distribution of \mathbf{z}**

433 Note that \mathbf{r}_{dcn} is a K -dimensional vector parameterized for a categorical distribution. Hence, one
434 constrain for \mathbf{r}_{dcn} is:

$$\sum_{k=1}^K r_{dcnk} = 1. \quad (41)$$

435 We form the Lagrangian which consists of the lower-bound with isolating terms relating r_{dcnk} and
 436 add the appropriate Lagrange multiplier:

$$\begin{aligned} \mathcal{L}[r_{dcnk}] &= \mathbb{E}_q [\ln p(\mathbf{z}|\boldsymbol{\theta}) + \ln p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) - \ln q(\mathbf{z})] + \zeta \left(\sum_{k=1}^K r_{dcnk} - 1 \right) \\ &= r_{dcnk} \ln \tilde{\theta}_{dck} + \frac{r_{dcnk}}{2} \left[\ln \tilde{\Lambda}_k - \nu_k(\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \mathbf{W}_k(\mathbf{x}_{dcn} - \mathbf{m}_k) - D\kappa_k^{-1} - D \ln(2\pi) \right] \\ &\quad - r_{dcnk} \ln r_{dcnk} + \zeta \left(\sum_{k=1}^K r_{dcnk} - 1 \right), \end{aligned} \quad (42)$$

437 where: $\ln \tilde{\theta}_{dck}$ and $\ln \tilde{\Lambda}_k$ are defined in Eqs. (18) and (27).

438 Taking the derivative w.r.t. r_{dcnk} gives:

$$\frac{\partial \mathcal{L}}{\partial r_{dcnk}} = \ln \tilde{\theta}_{dck} + \frac{1}{2} \left[\ln \tilde{\Lambda}_k - \nu_k(\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \mathbf{W}_k(\mathbf{x}_{dcn} - \mathbf{m}_k) - D\kappa_k^{-1} - D \ln(2\pi) \right] - \ln r_{dcnk} - 1 + \zeta. \quad (43)$$

439 Setting this derivative to zero and solving for r_{dcnk} gives:

$$r_{dcnk} \propto \exp \left\{ \ln \tilde{\theta}_{dck} + \frac{1}{2} \left[\ln \tilde{\Lambda}_k - \nu_k(\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \mathbf{W}_k(\mathbf{x}_{dcn} - \mathbf{m}_k) - D\kappa_k^{-1} \right] \right\}. \quad (44)$$

440 C.2 Variational parameter for Dirichlet distribution of θ

441 The lower-bound with isolating terms relating to only γ_{dck} is:

$$\begin{aligned} \mathcal{L}[\gamma_{dck}] &= \mathbb{E}_q [\ln p(\mathbf{z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}|\mathbf{y}, \boldsymbol{\alpha}) - \ln q(\boldsymbol{\theta})] \\ &= \sum_{n=1}^N \sum_{k=1}^K r_{dcnk} \ln \tilde{\theta}_{dck} + \sum_{l=1}^L \eta_{dcl} \sum_{k=1}^K (\alpha_{lk} - 1) \ln \tilde{\theta}_{dck} + \ln B(\boldsymbol{\gamma}_{dc}) - \sum_{k=1}^K (\gamma_{dck} - 1) \ln \tilde{\theta}_{dck} \\ &= \sum_{k=1}^K \ln \tilde{\theta}_{dck} \left[\sum_{n=1}^N r_{dcnk} + \sum_{l=1}^L \eta_{dcl} (\alpha_{lk} - 1) - \gamma_{dck} + 1 \right] + \ln B(\boldsymbol{\gamma}_{dc}), \end{aligned} \quad (45)$$

442 where $B(\mathbf{u})$ is defined in Eq. (21) – the normalizing constant of the Dirichlet distribution parameter-
 443 ized by \mathbf{u} . Note that: the sum notations are applicable only on their own single lines.

444 Taking derivative w.r.t. γ_{dck} gives:

$$\begin{aligned} \frac{\partial \mathcal{L}[\gamma_{dck}]}{\partial \gamma_{dck}} &= \Psi(\gamma_{dck}) \left[\sum_{n=1}^N r_{dcnk} + \sum_{l=1}^L \eta_{dcl} (\alpha_{lk} - 1) - \gamma_{dck} + 1 \right] \\ &\quad - \Psi \left(\sum_{j=1}^K \gamma_{dcj} \right) \sum_{j=1}^K \left[\sum_{n=1}^N r_{dcnj} + \sum_{l=1}^L \eta_{dcl} (\alpha_{lj} - 1) - \gamma_{dcj} + 1 \right]. \end{aligned} \quad (46)$$

445 Setting the derivative to zero and solve for γ_{dck} yields:

$$\gamma_{dck} = 1 + \sum_{n=1}^N r_{dcnk} + \sum_{l=1}^L \eta_{dcl} (\alpha_{lk} - 1). \quad (47)$$

446 C.3 Variational parameter for \mathbf{y}

447 Note that the L -dimensional vector $\boldsymbol{\eta}_{dc}$ is the parameter of a categorical distribution for \mathbf{y}_{dc} , it
 448 satisfies the following constrain:

$$\sum_{l=1}^L \eta_{dcl} = 1. \quad (48)$$

449 The Lagrangian can be expressed as:

$$\mathbb{L}[\mathbf{y}_{dc}] = \sum_{l=1}^L \eta_{dcl} \left[-\ln B(\boldsymbol{\alpha}_l) + \sum_{k=1}^K (\alpha_{lk} - 1) \ln \tilde{\theta}_{dck} \right] + \sum_{l=1}^L \eta_{dcl} \ln \tilde{\phi}_{dl} - \sum_{l=1}^L \eta_{dcl} \ln \eta_{dcl} + \xi \left(\sum_{l=1}^L \eta_{dcl} - 1 \right). \quad (49)$$

450 Taking the derivative w.r.t. η_{dcl} gives:

$$\frac{\partial \mathbb{L}}{\partial \eta_{dcl}} = -\ln B(\boldsymbol{\alpha}_l) + \sum_{k=1}^K (\alpha_{lk} - 1) \ln \tilde{\theta}_{dck} + \psi(\lambda_{dl}) - \psi \left(\sum_{j=1}^K \lambda_{dj} \right) - \ln \eta_{dcl} - 1 + \xi. \quad (50)$$

451 Setting the derivative to zero and solve for η_{dcl} yields:

$$\eta_{dcl} \propto \exp \left[\ln \tilde{\phi}_{dl} - \ln B(\boldsymbol{\alpha}_l) + \sum_{k=1}^K (\alpha_{lk} - 1) \ln \tilde{\theta}_{dck} \right]. \quad (51)$$

452 C.4 Variational parameter for ϕ

$$\begin{aligned} \mathbb{L}[\boldsymbol{\lambda}_d] &= \sum_{c=1}^C \sum_{l=1}^L \eta_{dcl} \ln \tilde{\phi}_{dl} + \sum_{l=1}^L (\delta_l - 1) \ln \tilde{\phi}_{dl} + \ln B(\boldsymbol{\lambda}_d) - \sum_{l=1}^L (\lambda_{dl} - 1) \ln \tilde{\phi}_{dl} \\ &= \ln B(\boldsymbol{\lambda}_d) + \sum_{l=1}^L \ln \tilde{\phi}_{dl} \left(\delta_l - \lambda_{dl} + \sum_{c=1}^C \eta_{dcl} \right). \end{aligned} \quad (52)$$

453 Taking derivative gives:

$$\frac{\partial \mathbb{L}}{\partial \lambda_{dl}} = \Psi(\lambda_{dl}) \left(\delta_l - \lambda_{dl} + \sum_{c=1}^C \eta_{dcl} \right) - \Psi \left(\sum_{j=1}^L \lambda_{dj} \right) \sum_{l=1}^L \left(\delta_l - \lambda_{dl} + \sum_{c=1}^C \eta_{dcl} \right). \quad (53)$$

454 Setting to zero and solving for λ_{dl} gives:

$$\lambda_{dl} = \delta_l + \sum_{c=1}^C \eta_{dcl}. \quad (54)$$

455 C.5 Variational parameter for word-topics

456 The lower-bound with the terms relating only to $\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k$ is:

$$\begin{aligned} \mathbb{L}[\boldsymbol{\mu}, \boldsymbol{\Lambda}] &= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \sum_{k=1}^K \frac{r_{dcnk}}{2} \left[\ln \tilde{\Lambda}_k - \nu_k (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \mathbf{W}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) - \frac{D}{\kappa_k} \right] \\ &\quad - \sum_{k=1}^K \frac{1}{2} \left[\nu_k \kappa_0 (\mathbf{m}_k - \mathbf{m}_0)^\top \mathbf{W}_k (\mathbf{m}_k - \mathbf{m}_0) + \frac{D \kappa_0}{\kappa_k} + D \ln \kappa_k \right] \\ &\quad - D_{\text{KL}} [\mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \parallel \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_0, \nu_0)]. \end{aligned} \quad (55)$$

457 Note that:

$$\frac{\partial}{\partial \mathbf{W}_k} \ln \tilde{\Lambda}_k = (\mathbf{W}_k^{-1})^\top = \mathbf{W}_k^{-1} \quad (56)$$

$$\frac{\partial}{\partial \nu_k} \ln \tilde{\Lambda}_k = \frac{1}{2} \sum_{i=1}^D \Psi \left(\frac{\nu_k - i + 1}{2} \right), \quad (57)$$

458 and:

$$\begin{aligned} &D_{\text{KL}} [\mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k) \parallel \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_0, \nu_0)] \\ &= \mathcal{H} [\mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k), \mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_0, \nu_0)] - \mathcal{H} [\mathcal{W}(\boldsymbol{\Lambda}_k; \mathbf{W}_k, \nu_k)] \\ &= -\frac{\nu_0}{2} \ln |\mathbf{W}_0^{-1} \mathbf{W}_k| + \frac{\nu_k}{2} [\text{tr}(\mathbf{W}_0^{-1} \mathbf{W}_k) - D] + \ln \Gamma_D \left(\frac{\nu_0}{2} \right) - \ln \Gamma_D \left(\frac{\nu_k}{2} \right) \\ &\quad + \frac{\nu_k - \nu_0}{2} \psi_D \left(\frac{\nu_k}{2} \right), \end{aligned} \quad (58)$$

where: Γ_D and ψ_D are the multivariate gamma and digamma function.

Therefore:

$$\begin{cases} \frac{\partial}{\partial \mathbf{W}_k} D_{\text{KL}} &= -\frac{\nu_0}{2} \mathbf{W}_k^{-1} + \frac{\nu_k}{2} \mathbf{W}_0^{-1} \\ \frac{\partial}{\partial \nu_k} D_{\text{KL}} &= \text{tr}(\mathbf{W}_0^{-1} \mathbf{W}_k) - D + \frac{\nu_k - \nu_0}{4} \Psi_D\left(\frac{\nu_k}{2}\right). \end{cases} \quad (59)$$

Taking derivative of the lower-bound w.r.t. the variational parameters of μ_k, Λ_k gives:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{m}_k} = \nu_k \kappa_0 \mathbf{W}_k (\mathbf{m}_0 - \mathbf{m}_k) + \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N r_{dcnk} \nu_k \mathbf{W}_k (\mathbf{x}_{dcn} - \mathbf{m}_k). \quad (60)$$

$$\frac{\partial \mathcal{L}}{\partial \kappa_k} = \frac{D \kappa_0}{2 \kappa_k^2} - \frac{D}{2 \kappa_k} + \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \frac{D r_{dcnk}}{2 \kappa_k^2} = \frac{D}{2 \kappa_k^2} \left[\kappa_0 - \kappa_k + \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N r_{dcnk} \right]. \quad (61)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{W}_k} &= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \frac{r_{dcnk}}{2} [\mathbf{W}_k^{-1} - \nu_k (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top (\mathbf{x}_{dcn} - \mathbf{m}_k)] \\ &\quad - \frac{1}{2} \nu_k \kappa_0 (\mathbf{m}_k - \mathbf{m}_0)^\top (\mathbf{m}_k - \mathbf{m}_0) + \frac{\nu_0}{2} \mathbf{W}_k^{-1} - \frac{\nu_k}{2} \mathbf{W}_0^{-1}. \end{aligned} \quad (62)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \nu_k} &= \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N \frac{r_{dcnk}}{2} \left[\frac{1}{2} \sum_{i=1}^D \Psi\left(\frac{\nu_k - i + 1}{2}\right) - (\mathbf{x}_{dcn} - \mathbf{m}_k)^\top \mathbf{W}_k (\mathbf{x}_{dcn} - \mathbf{m}_k) \right] \\ &\quad - \frac{1}{2} \kappa_0 (\mathbf{m}_k - \mathbf{m}_0)^\top \mathbf{W}_k (\mathbf{m}_k - \mathbf{m}_0) \\ &\quad - \text{tr}(\mathbf{W}_0^{-1} \mathbf{W}_k) + D - \frac{\nu_k - \nu_0}{4} \Psi_D\left(\frac{\nu_k}{2}\right). \end{aligned} \quad (63)$$

Setting these partial derivatives to zero and solving for the 4 parameters of interest give:

$$\begin{aligned} \mathbf{m}_k &= \frac{\kappa_0}{\kappa_0 + N_k} \mathbf{m}_0 + \frac{N_k}{\kappa_0 + N_k} \bar{\mathbf{x}}_k \\ \kappa_k &= \kappa_0 + N_k \\ \mathbf{W}_k^{-1} &= \mathbf{W}_0^{-1} + N_k \mathbf{S}_k + \frac{\kappa_0 N_k}{\kappa_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^\top \\ \nu_k &= \nu_0 + N_k, \end{aligned} \quad (64)$$

where:

$$N_k = \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N r_{dcnk} \quad (65)$$

$$\bar{\mathbf{x}}_k = \frac{1}{N_k} \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N r_{dcnk} \mathbf{x}_{dcn} \quad (66)$$

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{d=1}^M \sum_{c=1}^C \sum_{n=1}^N r_{dcnk} (\mathbf{x}_{dcn} - \bar{\mathbf{x}}_k)(\mathbf{x}_{dcn} - \bar{\mathbf{x}}_k)^\top. \quad (67)$$

D Experiments on Omniglot and mini-ImageNet

For Omniglot, we follow the pre-processing steps as in few-shot image classification without any data augmentation, and use the standard train-test split in the original paper to prevent information leakage. For mini-ImageNet, we follow the common train-test split with 80 classes for training and 20 classes for testing Ravi and Larochelle, 2017. Since the dimension of raw images in mini-ImageNet is large,

we employ the 640-dimensional features extracted from a wide-residual-network Rusu et al., 2019 to ease the calculation.

We follow Algorithm 1 to obtain the parameters of the image-topics posterior using tasks in training set. We use $L = 8$ task-topics and $K = 16$ image-topics for Omniglot, and $L = 3$ and $K = 10$ for mini-ImageNet. The Dirichlet distribution for task-topic proportion is chosen to be symmetric with $\delta = 0.01$. The hyper-parameter α is uniformly chosen in $[0.01, 0.1]$ and held fixed in each experiment. The training for LDCC is carried out with 20 images per class to fit into the memory of a Nvidia 1080 Ti GPU, while the inference of the variational parameter λ is done on all available labelled data in each class (20 for Omniglot and 600 for mini-ImageNet). Note that this is for the inference of LDCC used in the correlation diagram. For the task selection, this number matches the number of shots in the few-shot learning setting.

D.1 Pre-processing

For Omniglot, the grey-scale images of each character are resized to 28-by-28 pixels, resulting in a 728-dimensional vector. We use the original train-test split, where 30 alphabets are used for training, and the other 20 alphabets are used for testing. No rotation is applied to augment classes.

For mini-ImageNet, we use the extracted features that are presented by 640-dimensional vectors. In addition, we normalise those features to be within $[0, 1]$ by multiplying with a factor of 3.

D.2 Hyper-parameters in LDCC

The prior parameters of image-topics, $\mathbf{m}_0, \kappa_0, \mathbf{W}_0, \nu_0$, are selected as follows:

- The mean \mathbf{m}_0 is selected as the mean of all the images in the training of each data set,
- κ_0 is set to 0.01,
- The scale matrix \mathbf{W}_0 is selected as the covariance matrix of all the images in the training set,
- ν_0 is set to the dimension of images (728 for Omniglot and 640 for mini-ImageNet) adding 2.

The hyper-parameters for the learning rate used in online learning are $\tau_0 = 100$ and $\tau_1 = 0.5$. The mini-batch used in both data sets for LDCC is 500.

D.3 Network architectures used in meta-learning

We use two different network architectures on the two data sets. For Omniglot, we follow the “standard” 4 module CNN network, where each module consists of 64 3-by-3 filter convolutional layer, followed by batch normalisation, activated by ReLU and 2-by-2 max-pooling. For mini-ImageNet, we use a fully connected network with 1 hidden layer consisting of 128 hidden units, activated by ReLU.