
Massively Parallel Reweighted Wake-Sleep (Supplementary Material)

Thomas Heap¹

Gavin Leech¹

Laurence Aitchison¹

¹Department of Computer Science, University of Bristol, Bristol

1 PROOF OF EQUIVALENCE OF THE DIFFERENT FORMS OF THE GLOBAL RWS UPDATES

We start with the RWS P update (Eq. 8a), then use $\nabla_{\theta} \log \mathcal{P}_{\text{global}}(z) = (\nabla_{\theta} \mathcal{P}_{\text{global}}(z)) / \mathcal{P}_{\text{global}}(z)$,

$$\mathbb{E} [\Delta\theta_{\text{RWS}}] = \mathbb{E}_{Q_{\phi}(z|x)} \left[\frac{\nabla_{\theta} \mathcal{P}_{\text{global}}(z)}{\mathcal{P}_{\text{global}}(z)} \right]. \quad (1)$$

Using the definition of $\mathcal{P}_{\text{global}}(z)$ (Eq. 4),

$$\mathbb{E} [\Delta\theta_{\text{RWS}}] = \mathbb{E}_{Q_{\phi}(z|x)} \left[\frac{\nabla_{\theta} \frac{1}{K} \sum_k r_k(z)}{\mathcal{P}_{\text{global}}(z)} \right] \quad (2)$$

Substituting for $r_k(z)$ (Eq. 5) in the numerator,

$$\mathbb{E} [\Delta\theta_{\text{RWS}}] = \mathbb{E}_{Q_{\phi}(z|x)} \left[\frac{\frac{1}{K} \sum_k \frac{\nabla_{\theta} P_{\theta}(z^k, x)}{Q_{\phi}(z^k|x)}}{\mathcal{P}_{\text{global}}(z)} \right] \quad (3)$$

substituting $\nabla_{\theta} P_{\theta}(z^k, x) = P_{\theta}(z^k, x) \nabla_{\theta} \log P_{\theta}(z^k, x)$,

$$\mathbb{E} [\Delta\theta_{\text{RWS}}] = \mathbb{E}_{Q_{\phi}(z|x)} \left[\frac{1}{K} \sum_k \frac{P_{\theta}(z^k, x)}{Q_{\phi}(z^k|x) \mathcal{P}_{\text{global}}(z)} \nabla_{\theta} \log P_{\theta}(z^k, x) \right] \quad (4)$$

Noticing that the ratio of $P_{\theta}(z^k, x)$ and $Q_{\phi}(z^k|x)$ in the numerator is equal to $r_k(z)$ (Eq. 5), we get back to Eq. (7a), as required.

The RWS Q update is very similar. Again, we start with Eq. (8b), then use $\nabla_{\theta} \log \mathcal{P}_{\text{global}}(z) = (\nabla_{\theta} \mathcal{P}_{\text{global}}(z)) / \mathcal{P}_{\text{global}}(z)$,

$$\mathbb{E} [\Delta\phi_{\text{RWS}}] = -\mathbb{E}_{Q_{\phi}(z|x)} \left[\frac{\nabla_{\phi} \mathcal{P}_{\text{global}}(z)}{\mathcal{P}_{\text{global}}(z)} \right] \quad (5)$$

Using the definition of $\mathcal{P}_{\text{global}}(z)$ (Eq. 4),

$$\mathbb{E} [\Delta\phi_{\text{RWS}}] = -\mathbb{E}_{Q_{\phi}(z|x)} \left[\frac{\nabla_{\phi} \frac{1}{K} \sum_k r_k(z)}{\mathcal{P}_{\text{global}}(z)} \right] \quad (6)$$

Substituting for $r_k(z)$ (Eq. 5) in the numerator,

$$\mathbb{E} [\Delta\phi_{\text{RWS}}] = -\mathbb{E}_{Q_{\phi}(z|x)} \left[\frac{\frac{1}{K} \sum_k \nabla_{\phi} \frac{P_{\theta}(z^k, x)}{Q_{\phi}(z^k|x)}}{\mathcal{P}_{\text{global}}(z)} \right] \quad (7)$$

Computing the derivative,

$$\mathbb{E} [\Delta\phi_{\text{RWS}}] = \mathbb{E}_{\mathcal{Q}_\phi(z|x)} \left[\frac{\frac{1}{K} \sum_k \frac{P_\theta(z^k, x)}{(\mathcal{Q}_\phi(z^k|x))^2} \nabla_\phi \mathcal{Q}_\phi(z^k|x)}{\mathcal{P}_{\text{global}}(z)} \right]. \quad (8)$$

Noticing that $(\nabla_\phi \mathcal{Q}_\phi(z^k|x)) / \mathcal{Q}_\phi(z^k|x) = \nabla_\phi \log \mathcal{Q}_\phi(z^k|x)$,

$$\mathbb{E} [\Delta\phi_{\text{RWS}}] = \mathbb{E}_{\mathcal{Q}_\phi(z|x)} \left[\frac{\frac{1}{K} \sum_k \frac{P_\theta(z^k, x)}{\mathcal{Q}_\phi(z^k|x)} \nabla_\phi \log \mathcal{Q}_\phi(z^k|x)}{\mathcal{P}_{\text{global}}(z)} \right]. \quad (9)$$

Finally, noticing that the ratio of $P_\theta(z^k, x)$ and $\mathcal{Q}_\phi(z^k|x)$ in the numerator is equal to $r_k(z)$ (Eq. 5), we get back to Eq. (7b), as required.

Both of these derivations may be straightforwardly repeated for the massively parallel setting, simply by replacing $k \in \mathcal{K}$ with $\mathbf{k} \in \mathcal{K}^n$, and by replacing $1/K$ with $1/K^n$.

2 TMC VS MASSIVELY PARALLEL APPROXIMATE POSTERiors

TMC approximate posteriors draw the K samples of the i th latent variable IID,

$$\mathcal{Q}_{\text{TMC}}(z_i | z_j \text{ for all } j \in \text{qa}(i)) = \prod_{k_i \in \mathcal{K}} \mathcal{Q}_{\text{TMC}}(z_i^{k_i} | z_j \text{ for all } j \in \text{qa}(i)) \quad (10)$$

Specifically, TMC draws each sample from an equally weighted mixture over all parent particles,

$$\mathcal{Q}_{\text{TMC}}(z_i^{k_i} | z_j \text{ for all } j \in \text{qa}(i)) = \frac{1}{K^{|\text{qa}(i)|}} \sum_{\mathbf{k}_{\text{qa}(i)}} \mathcal{Q}_{\text{global}}(z_i^{k_i} | z_j^{k_j} \text{ for all } j \in \text{qa}(i)). \quad (11)$$

In contrast, massively parallel methods do not force us to sample particles IID. The key issue with IID sampling is that it introduces the risk of particle degeneracy [Carpenter et al., 1999, Li et al., 2012, 2014, Zhou et al., 2016, Wang et al., 2017]. In particle degeneracy, some of the parent samples (e.g. z_j^1 where $j \in \text{qa}(i)$) might have multiple children, in the sense that multiple z_i^k are sampled from the mixture component arising from z_j^1 . At the same time, some of the parents, (e.g. z_j^2) might have no children, in the sense that no z_i^k are sampled from a mixture component arising from z_j^2 . This is problematic because it reduces diversity in the population of samples, $z_i = (z_i^1, \dots, z_i^K)$, and this reduction in diversity can be especially problematic in models with long chains of latent variables, such as timeseries models. To reduce the risk of particle degeneracy, the massively parallel methods considered here couple the distribution over each of the K particles,

$$\mathcal{Q}_{\text{MP}}(z_i | z_j \text{ for all } j \in \text{qa}(i)) \neq \prod_{k_i \in \mathcal{K}} \mathcal{Q}_{\text{TMC}}(z_i^{k_i} | z_j \text{ for all } j \in \text{qa}(i)). \quad (12)$$

However, we do ensure that the marginal for a single particle is the same as for TMC,

$$\mathcal{Q}_{\text{MP}}(z_i^{k_i} | z_j \text{ for all } j \in \text{qa}(i)) = \frac{1}{K^{|\text{qa}(i)|}} \sum_{\mathbf{k}_{\text{qa}(i)}} \mathcal{Q}_{\text{global}}(z_i^{k_i} | z_j^{k_j} \text{ for all } j \in \text{qa}(i)). \quad (13)$$

To achieve this, we sample a permutation, π for each latent variable, and the permutation tells us which parent particle to consider. To give an example for one parent,

$$\mathcal{Q}_{\text{MP}}(z_i | \pi, z_j) = \prod_{k_i} \mathcal{Q}_{\text{MP}}(z_i^{k_i} | \pi, z_j) \quad (14)$$

$$\mathcal{Q}_{\text{MP}}(z_i^{k_i} | \pi, z_j) = \frac{1}{K} \sum_{k_i} \mathcal{Q}_\phi(z_i^{k_i} | z_j^{\pi k_i}) \quad (15)$$

Critically, if we marginalise over the permutation, the distribution over a single $z_i^{k_i}$ has the same density as that from a uniform mixture,

$$Q_{\text{MP}}\left(z_i^{k_i} \mid z_j\right) = \sum_{\pi} Q_{\text{MP}}\left(z_i^{k_i} \mid \pi, z_j\right) \quad (16)$$

$$Q_{\text{MP}}\left(z_i^{k_i} \mid z_j\right) = \frac{1}{K} \sum_{k_j} Q_{\phi}\left(z_i^{k_i} \mid z_j^{k_j}\right). \quad (17)$$

Finally, if we have multiple parent latent variables, we independently sample a permutation for each latent variable.

3 MASSIVELY PARALLEL IWAE AND RWS

Before getting started, it will prove useful to define some briefer notation than that used in the main text. Specifically, we use,

$$z_{\text{qa}(i)} = \{z_j \text{ for all } j \in \text{qa}(i)\}, \quad (18)$$

$$z_{\text{qa}(i)}^{k_i} = \left\{z_j^{k_i} \text{ for all } j \in \text{qa}(i)\right\}, \quad (19)$$

$$z_{\text{pa}(i)}^{\mathbf{k}_{\text{pa}(i)}} = \left\{z_j^{k_j} \text{ for all } j \in \text{pa}(i)\right\}, \quad (20)$$

so,

$$Q\left(z_i^{k_i} \mid x, z_{\text{qa}(i)}\right) = Q\left(z_i^{k_i} \mid x, z_j \text{ for all } j \in \text{qa}(i)\right), \quad (21)$$

$$P_{\theta}\left(z_i^{k_i} \mid z_{\text{pa}(i)}^{\mathbf{k}_{\text{pa}(i)}}\right) = P_{\theta}\left(z_i^{k_i} \mid z_j^{k_j} \text{ for all } j \in \text{pa}(i)\right) \quad (22)$$

Note that in Eq. (21), we allow for the possibility of a slightly more general form for the approximate posterior, where the distribution over $z_i^{k_i}$ may depend on any of the parent samples. This generalisation ensures that the subsequent derivations generalise to other possible forms for the approximate posterior, such as those for TMC (Eq. 11).

In addition, it is useful to introduce notation to describe the ‘‘non-indexed’’ latent variables (i.e. everything in z that is not $z^{\mathbf{k}}$). The i th non-indexed latents are, $z_i^{/\mathbf{k}_i}$,

$$z_i^{/\mathbf{k}_i} = \left(z_i^1, \dots, z_i^{k_i-1}, z_i^{k_i+1}, \dots, z_i^K\right) \in \mathcal{Z}_i^{K-1}. \quad (23)$$

and $z^{/\mathbf{k}}$ are all non-indexed latents,

$$z^{/\mathbf{k}} = \left(z_1^{/k_1}, z_2^{/k_2}, \dots, z_n^{/k_n}\right) \in \mathcal{Z}^{K-1}. \quad (24)$$

3.1 IWAE

3.1.1 Single-Sample VI

We begin by building intuition by looking at the derivation for the ELBO in the standard single-sample VAE. We start by writing the marginal likelihood as an integral,

$$P_{\theta}(x) = \int dz' P_{\theta}(x, z'). \quad (25)$$

Here, we use $z' \in \mathcal{Z}$ to denote a single sample from the full joint state space; we use z' instead of z because z is reserved for K samples (Eq. 2). Next, we divide and multiply by the approximate posterior probability, $Q_{\phi}(z'|x)$,

$$P_{\theta}(x) = \int dz' Q_{\phi}(z'|x) \frac{P_{\theta}(x, z')}{Q_{\phi}(z'|x)}. \quad (26)$$

Now, we can rewrite the integral as an expectation under the approximate posterior,

$$P_\theta(x) = E_{Q_\phi(z'|x)} \left[\frac{P_\theta(x, z')}{Q_\phi(z'|x)} \right]. \quad (27)$$

Now we take the logarithm on both sides and apply Jensen’s inequality,

$$\log P_\theta(x) = \log E_{Q_\phi(z'|x)} \left[\frac{P_\theta(x, z')}{Q_\phi(z'|x)} \right] \geq E_{Q_\phi(z'|x)} \left[\log \frac{P_\theta(x, z')}{Q_\phi(z'|x)} \right] = \mathcal{L}_{\text{VAE}} \quad (28)$$

Of course, this derivation is specific to the single-sample VAE. But we can pull out an underlying strategy that generalises to the multi-sample setting. In particular, we first come up with an unbiased estimator of the marginal likelihood. In our VAE, this is,

$$\mathcal{P}_{\text{VAE}}(z') = \frac{P_\theta(x, z')}{Q_\phi(z'|x)} \quad (29)$$

Following Eq. (27) we can see that this quantity is an unbiased estimator of the marginal likelihood if z' is sampled from $Q_\phi(z'|x)$,

$$P_\theta(x) = E_{Q_\phi(z'|x)} [\mathcal{P}_{\text{VAE}}(z')] \quad (30)$$

Then we apply Jensen’s inequality (Eq. 28),

$$\log P_\theta(x) \geq \mathcal{L}_{\text{VAE}} = E_{Q_\phi(z'|x)} [\log \mathcal{P}_{\text{VAE}}(z')]. \quad (31)$$

However, this approach highlights key issues with the usual single-sample bound. In particular, the single-sample estimator, $\mathcal{P}_{\text{VAE}}(z')$ can be very high-variance, and variance in the unbiased estimator causes the Jensen bound to become looser.

3.1.2 Global IWAE

To reduce variance in the unbiased estimator, a natural approach is to average K independent samples, and this is exactly what global IWAE does,

$$\mathcal{P}_{\text{global}}(z) = \frac{1}{K} \sum_{k=1}^K r_k(z) = \frac{1}{K} \sum_{k=1}^K \mathcal{P}_{\text{VAE}}(z^k) \quad (32)$$

This is of course an unbiased estimator, as it is the average of K unbiased estimators,

$$P_\theta(x) = E_{Q_\phi(z|x)} [\mathcal{P}_{\text{global}}(z)]. \quad (33)$$

Therefore, applying Jensen’s inequality gives a new lower-bound on the log-marginal likelihood,

$$\log P_\theta(x) = \log E_{Q_\phi(z|x)} [\mathcal{P}_{\text{MP}}(z)] \geq E_{Q_\phi(z|x)} [\log \mathcal{P}_{\text{global}}(z)] = \mathcal{L}_{\text{IWAE}} \quad (34)$$

which is tighter than the usual single-sample ELBO [Burda et al., 2015], and which matches Eq. (6) in the main text.

3.1.3 Massively Parallel IWAE

Our proposed $\mathcal{P}_{\text{MP}}(z)$ (Eq. 15) is the average of K^n terms, rather than K terms in global IWAE. To prove that our massively parallel strategy is valid, our strategy is to show that every term in this average is an unbiased estimator of $\log P_\theta(x)$, in which case the average is also an unbiased estimator, and we can again apply Jensen.

Each term in the average $\mathcal{P}_{\text{MP}}(z)$ (Eq. 15) is of the form $r_{\mathbf{k}}(z)$ (Eq. 16). The expectation of each term is,

$$E_{Q_\phi(z|x)} [r_{\mathbf{k}}(z)] = E_{Q_\phi(z|x)} \left[\frac{P_\theta(x, z^{\mathbf{k}})}{\prod_i Q_\phi(z_i^{k_i} | x, z_{\text{qa}(i)})} \right]. \quad (35)$$

We can rewrite the expectation as an integral,

$$\mathbb{E}_{\mathbb{Q}_\phi(z|x)} [r_{\mathbf{k}}(z)] = \int dz P_\theta(x, z^{\mathbf{k}}) \prod_i \frac{\mathbb{Q}_\phi(z_i|x, z_{\text{qa}(i)})}{\mathbb{Q}_\phi(z_i^{k_i}|x, z_{\text{qa}(i)})}. \quad (36)$$

Bayes theorem tells us,

$$\frac{\mathbb{Q}_\phi(z_i|x, z_{\text{qa}(i)})}{\mathbb{Q}_\phi(z_i^{k_i}|x, z_{\text{qa}(i)})} = \frac{\mathbb{Q}_\phi(z_i^{k_i}, z_i^{/k_i}|x, z_{\text{qa}(i)})}{\mathbb{Q}_\phi(z_i^{k_i}|x, z_{\text{qa}(i)})} = \mathbb{Q}(z_i^{/k_i}|x, z_i^{k_i}, z_{\text{qa}(i)}), \quad (37)$$

Applying Bayes theorem,

$$\mathbb{E}_{\mathbb{Q}_\phi(z|x)} [r_{\mathbf{k}}(z)] = \int dz P_\theta(x, z^{\mathbf{k}}) \prod_i \mathbb{Q}(z_i^{/k_i}|x, z_i^{k_i}, z_{\text{qa}(i)}). \quad (38)$$

Importantly, the integrand is a valid joint distribution over x and z , or equivalently over x , $z^{\mathbf{k}}$ and $z^{/k}$. Thus, integrating over $z^{/k}$ then $z^{\mathbf{k}}$, we find,

$$\mathbb{E}_{\mathbb{Q}_\phi(z|x)} [r_{\mathbf{k}}(z)] = P_\theta(x). \quad (39)$$

As such, each of the $r_{\mathbf{k}}(z)$ terms is an unbiased estimator of the marginal likelihood. As $\mathcal{P}_{\text{MP}}(z)$ (Eq. 15) is just an average of K^n $r_{\mathbf{k}}(z)$ terms, it is also an unbiased estimator. Applying Jensen's inequality to this unbiased estimator,

$$\log P_\theta(x) \geq \mathbb{E}_{\mathbb{Q}_\phi(z|x)} [\log \mathcal{P}_{\text{MP}}(z)] = \mathcal{L}_{\text{MP}}, \quad (40)$$

which mirrors Eq. (17) in the main text.

3.2 RWS

3.2.1 Global RWS

To build intuition, we first give a derivation of the standard RWS updates. Ideally the updates would use samples drawn from the true posterior, $P_\theta(z|x)$,

$$\Delta\theta_{\text{post}} = \mathbb{E}_{P_\theta(z^k|x)} [\nabla_\theta \log P_\theta(z, x)] \quad (41a)$$

$$\Delta\phi_{\text{post}} = \mathbb{E}_{P_\theta(z^k|x)} [\nabla_\phi \log \mathbb{Q}_\phi(z|x)] \quad (41b)$$

The P update is exactly the M-step in EM, and the Q step trains $\mathbb{Q}_\phi(z|x)$ using maximum likelihood based on samples from the true posterior. To simplify the derivations, we note that both of these updates can be understood as computing a moment under the true posterior,

$$\Delta_{\text{post}} = \mathbb{E}_{P_\theta(z^k|x)} [\Delta(z^k)]. \quad (42)$$

For the P update, we have $\Delta_{\text{post}} = \Delta\theta_{\text{post}}$ and $\Delta(z^k) = \nabla_\theta \log P_\theta(z, x)$. For the Q update, we have $\Delta_{\text{post}} = \Delta\phi_{\text{post}}$ and $\Delta(z^k) = \nabla_\phi \log \mathbb{Q}_\phi(z, x)$. Of course, in practice, the true posterior is intractable, so instead we must use some form of importance weighting. We begin by writing the generic form for the updates as an integral,

$$\Delta_{\text{post}} = \int dz^k P(z^k|x) \Delta(z^k). \quad (43)$$

We then multiply and divide by an approximate posterior, $\mathbb{Q}(z^k|x)$,

$$\Delta_{\text{post}} = \int dz^k \mathbb{Q}(z^k|x) \frac{P(z^k|x)}{\mathbb{Q}(z^k|x)} \Delta(z^k). \quad (44)$$

We can rewrite the integral as expectation over the approximate posterior, $Q(z^k|x)$,

$$\Delta_{\text{post}} = \mathbb{E}_{Q(z^k|x)} \left[\frac{P(z^k|x)}{Q(z^k|x)} \Delta(z^k) \right]. \quad (45)$$

This quantity is difficult to use directly, because computing the posterior, $P(z^k|x)$ involves the marginal likelihood, $P_\theta(x)$, which is intractable,

$$P_\theta(z^k|x) = \frac{P_\theta(z^k, x)}{P_\theta(x)} \quad P_\theta(x) = \int dz^k P_\theta(z^k, x). \quad (46)$$

As the true marginal likelihood is intractable, we instead use $\mathcal{P}_{\text{global}}(z)$ (Eq. 4), which is an unbiased estimator of $P(x)$, and is correct in the limit as $K \rightarrow \infty$ [Burda et al., 2015]. This gives updates of the form,

$$\Delta_{\text{global}} = \mathbb{E}_{Q(z^k|x)} \left[\frac{\frac{P(z^k, x)}{Q(z^k|x)}}{\mathcal{P}_{\text{global}}(z)} \Delta(z^k) \right]. \quad (47)$$

Remembering the definition of $r_k(z)$ (Eq. 5), this can be written,

$$\Delta_{\text{global}} = \mathbb{E}_{Q_\phi(z|x)} \left[\frac{r_k(z)}{\mathcal{P}_{\text{global}}(z)} \Delta(z^k) \right]. \quad (48)$$

Finally as the expectation is the same for all k , we can average over k , which gives the expression in the main text (Eq. 7)

3.2.2 Massively Parallel RWS

Now, we can move on to massively parallel RWS. In the previous derivation for global RWS, we showed that each sample, z^k , individually constituted an unbiased estimator. In the massively parallel setting, the key difference is that instead of having K samples z^k , we have K^n samples, $z^{\mathbf{k}}$. In particular,

$$\Delta_{\text{post}} = \mathbb{E}_{P(z^{\mathbf{k}}|x)} [\Delta(z^{\mathbf{k}})] = \int dz^{\mathbf{k}} P(z^{\mathbf{k}}|x) \Delta(z^{\mathbf{k}}). \quad (49)$$

Now, we multiply and divide by $\prod_i Q(z_i^{k_i}|x, z_{\text{qa}(i)})$,

$$\Delta_{\text{post}} = \int dz^{\mathbf{k}} \left(\prod_i Q(z_i^{k_i}|x, z_{\text{qa}(i)}) \right) \frac{P(z^{\mathbf{k}}|x)}{\prod_i Q(z_i^{k_i}|x, z_{\text{qa}(i)})} \Delta(z^{\mathbf{k}}). \quad (50)$$

Now, we introduce and integrate out a distribution over the non-indexed latent variables, $\prod_i Q(z_i^{/k_i}|x, z_i^{k_i}, z_{\text{qa}(i)})$

$$1 = \int dz^{/k} \prod_i Q(z_i^{/k_i}|x, z_i^{k_i}, z_{\text{qa}(i)}), \quad (51)$$

Multiplying Eq. (50) by 1 (Eq. 51),

$$\Delta_{\text{post}} = \int dz^{\mathbf{k}} \left(\prod_i Q(z_i^{k_i}|x, z_{\text{qa}(i)}) \right) \frac{P(z^{\mathbf{k}}|x)}{\prod_i Q(z_i^{k_i}|x, z_{\text{qa}(i)})} \Delta(z^{\mathbf{k}}) \int dz^{/k} \prod_i Q(z_i^{/k_i}|x, z_i^{k_i}, z_{\text{qa}(i)}). \quad (52)$$

Combining the integrals over $z^{\mathbf{k}}$ and $z^{/k}$ into a single integral over z ,

$$\Delta_{\text{post}} = \int dz Q(z|x) \frac{P(z^{\mathbf{k}}|x)}{\prod_i Q(z_i^{k_i}|x, z_{\text{qa}(i)})} \Delta(z^{\mathbf{k}}). \quad (53)$$

Writing the integral as an expectation,

$$\Delta_{\text{post}} = \mathbb{E}_{Q_\phi(z|x)} \left[\frac{P(z^{\mathbf{k}}|x)}{\prod_i Q(z_i^{k_i}|x, z_{\text{pa}(i)})} \Delta(z^{\mathbf{k}}) \right]. \quad (54)$$

Again, the posterior can be written,

$$P_\theta(z^{\mathbf{k}}|x) = \frac{P_\theta(z^{\mathbf{k}}, x)}{P_\theta(x)} \quad P_\theta(x) = \int dz^{\mathbf{k}} P_\theta(z^{\mathbf{k}}, x). \quad (55)$$

Again, the marginal likelihood, $P_\theta(x)$ is intractable. Instead, we use the massively parallel estimate of the marginal likelihood, which was shown to be unbiased in Sec. 3.1.3,

$$\Delta_{\text{MP}} = \mathbb{E}_{Q_\phi(z|x)} \left[\frac{\frac{P(z^{\mathbf{k}}, x)}{\prod_i Q(z_i^{k_i}|x, z_{\text{pa}(i)})}}{\mathcal{P}_{\text{MP}}(z)} \Delta(z^{\mathbf{k}}) \right]. \quad (56)$$

Remembering the definition of $r_{\mathbf{k}}(z)$ (Eq. 16), this can be written,

$$\Delta_{\text{MP}} = \mathbb{E}_{Q_\phi(z|x)} \left[\frac{r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)} \Delta(z^{\mathbf{k}}) \right]. \quad (57)$$

Finally, note that the expectation is the same for every value of \mathbf{k} . Averaging over all K^n values of \mathbf{k} , we get the form in the main text (Eq. 18).

3.3 MOVIELENS GRAPHICAL MODEL

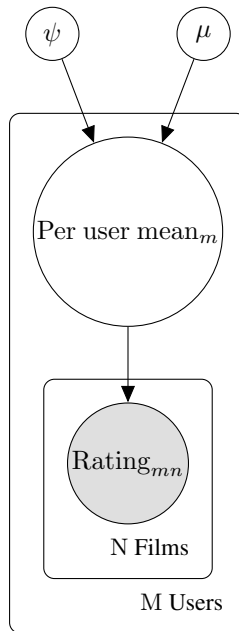


Figure 1: Graphical model for the MovieLens dataset

3.4 BUS DELAY MODEL SPECIFICATION

$$\begin{aligned}
 \text{YearVariance} &\sim \text{Cat}([0.1, 0.5, 0.4, 0.05, 0.05]) \\
 \text{YearMean} &\sim \mathcal{N}(0, 10^{-4}) \\
 \text{BoroughMean}_m &\sim \mathcal{N}(\text{YearMean}, \exp(\text{YearVariance})), m = 1, \dots, M \\
 \text{BoroughVariance}_j &\sim \text{Cat}([0.1, 0.4, 0.05, 0.5, 0.05]), j = 1, \dots, J \\
 \text{IdMean}_{mj} &\sim \mathcal{N}(\text{BoroughMean}_m, \text{BoroughVariance}_j), j = 1, \dots, J, m = 1, \dots, M \\
 \text{WeightVariance}_i &\sim \text{Cat}([0.1, 0.4, 0.5, 0.05, 0.05]), i = 1, \dots, I \\
 \mathbf{C}_i &\sim \mathcal{N}(\mathbf{0}_{\#\text{BusCo.s}}, \text{WeightVariance}_i), i = 1, \dots, I \\
 \mathbf{J}_i &\sim \mathcal{N}(\mathbf{0}_{\#\text{JourneyTypes}}, \text{WeightVariance}_i), i = 1, \dots, I \\
 \text{logits}_{mji} &= \text{IdMean}_{mj} + \mathbf{C}_i * \text{Bus company name}_{mji} + \mathbf{J}_i * \text{Journey type}_{mji} \\
 \text{Delay}_{mji} &\sim \text{NegativeBinomial}(\text{total count} = 130, \text{logits}_{mji}), i = 1, \dots, I, j = 1, \dots, J, \mu = 1, \dots, M
 \end{aligned} \tag{58}$$

Where Bus company name_{mji} is a one-hot encoded indicator variable indicating which bus company was running that route, and Journey type_{mji} similarly indicates which kind of bus journey was being undertaken. A total county of 130 is chosen as this is the largest recorded delay in the dataset.

3.5 BUS BREAKDOWN GRAPHICAL MODEL

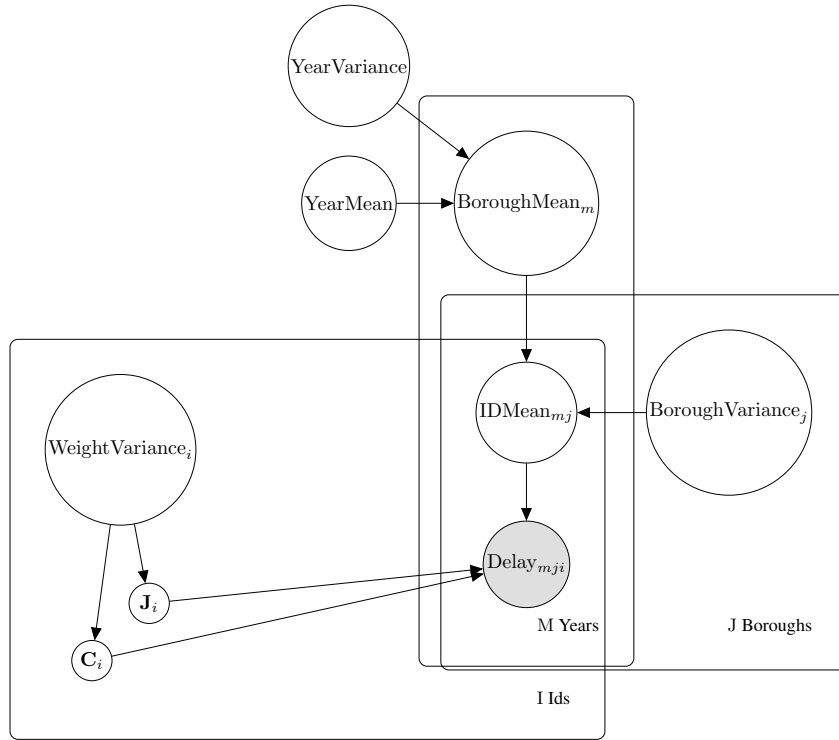


Figure 2: Graphical model for the bus breakdown dataset