

MASSIVELY PARALLEL REWEIGHTED WAKE-SLEEP

Thomas Heap, Gavin Leech, Laurence Aitchison

Department of Computer Science, University of Bristol, UK



Premise

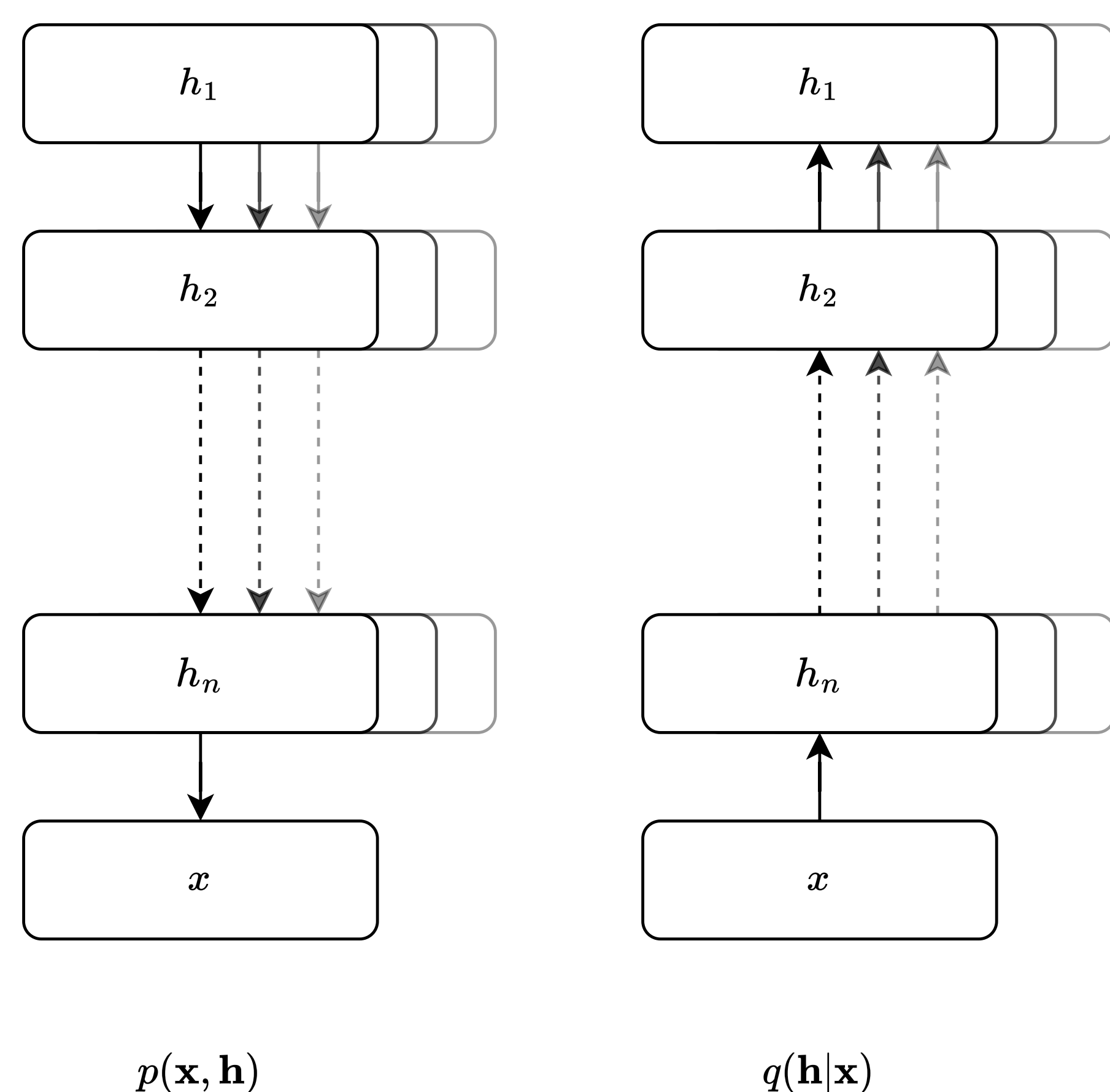
Reweighted wake-sleep (RWS) can perform Bayesian inference in a very general class of models. However, the number of samples required for effective importance weighting is exponential in the number of latents; getting this many importance samples is intractable in all but the smallest models. Here, we develop *massively parallel RWS*, which circumvents this issue by drawing K samples of all n latent variables and reasoning about all K^n possible combinations of samples.

Reasoning about K^n combinations might seem intractable, but the required computations can be performed in polynomial time by exploiting conditional independencies. Our algorithm and implementation of MPRWS show considerable improvements over standard RWS, which draws K samples from the full joint.

Background

We seek the parameters, ϕ , of an approximate posterior $Q_\phi(\mathbf{h}|\mathbf{x})$ with latents \mathbf{h} and data \mathbf{x} .

RWS [2] draws K samples from an underlying approximate posterior, then uses importance weighting to provide a better estimate of the true posterior P_θ . RWS then updates its approximate posterior towards the importance-weighted estimate of the true posterior.



Standard reweighted wake-sleep. Draws K samples from the joint latent space.

However, naively importance weighting, as we do in methods like RWS, does not provide enough samples for accurate estimation of the posterior [3], making practical applications of this method difficult.

In past work [1] we circumvent this problem with a massively parallel scheme drawing K samples for each latent variable, then *effectively* obtaining K^n samples by considering all combinations of K samples for each of the n latents. Let each individual sample for each separate latent variable be h_i^k , where k indexes the sample and i indexes the latent variable.

To sample all K copies of the full joint latent space, TMC [1] uses an IID distribution over the K samples, z_1^1, \dots, z_i^K ,

$$Q_{\text{TMC}}(h|x) = \prod_{i=1}^n \prod_{k \in \mathcal{K}} Q_{\text{TMC}}(h_i^k | h_j \text{ for all } j \in \text{qa}(i)). \quad (1)$$

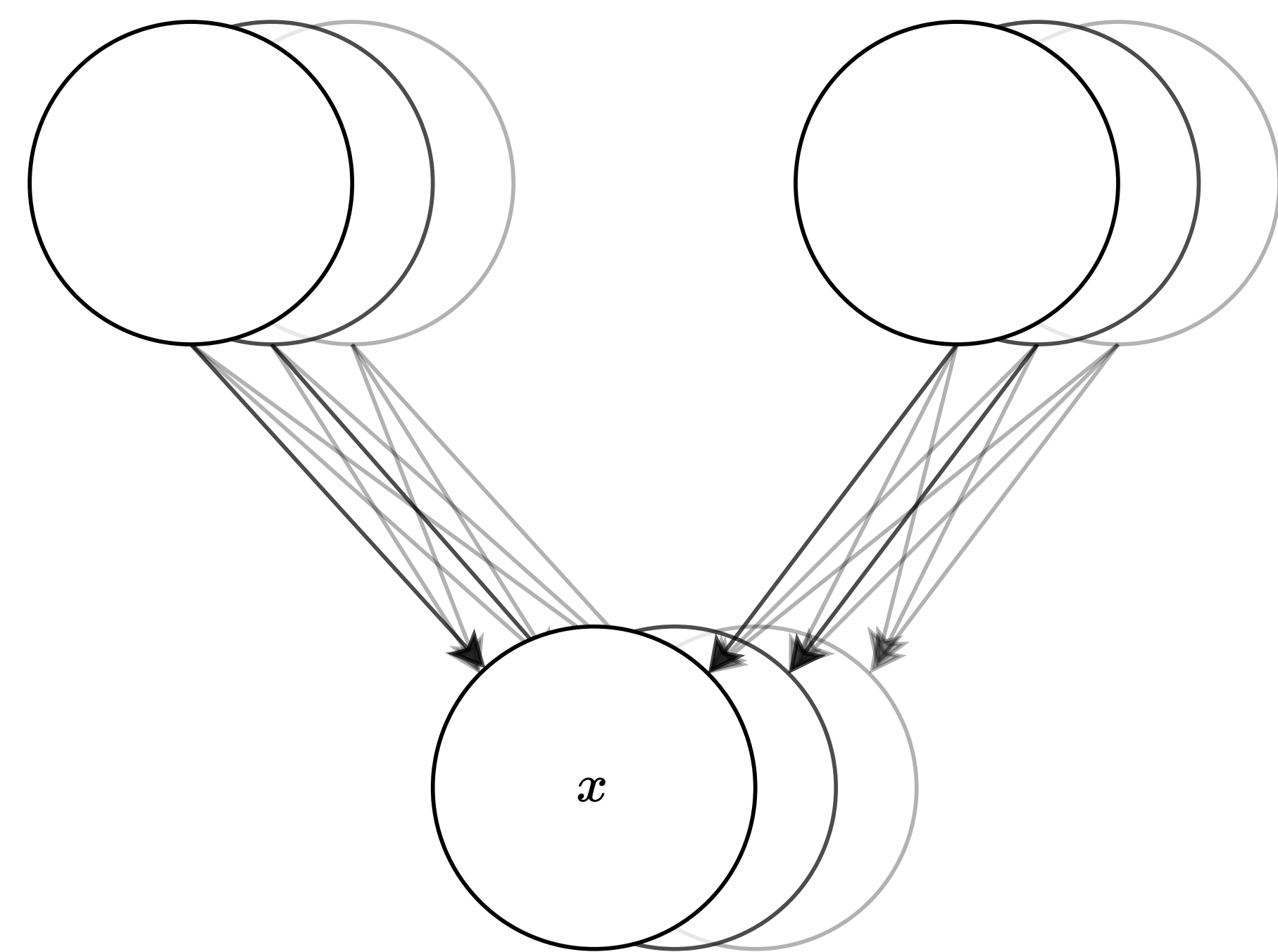
$\text{qa}(i)$ = indices of parents of the i th latent variable under the approximate posterior.

The present work generalises this, allowing dependencies (in Q) between the K samples for a single latent, and deriving parallel update rules for an RWS which averages over combinations of samples.

$$Q_{\text{MP}}(h|x) = \prod_{i=1}^n Q_{\text{MP}}(h_i | h_j \text{ for all } j \in \text{qa}(i)). \quad (2)$$

There are no formal constraints on these dependencies. However, there are practical constraints, namely that we need to be able to efficiently compute the single-particle marginals

Methods



Here, $K = 3$ and $n = 2$, so we reason over $K^n = 9$ combinations of samples.

Algorithm 1: Massively Parallel RWS

Require: Data x , Prior P_θ , Proposal Q_{MP} , $K \geq 1$

for $i \leftarrow 1$ to n **do**

 Sample $z_i \sim Q_{\text{MP}}(z_i | z_j \text{ for all } j \in \text{qa}(i))$

$z \leftarrow \{z_1, \dots, z_{i-1}\} \cup z_i$

$f_{k_i, k_{\text{pa}(i)}}^i(z) \leftarrow \frac{P_\theta(z_i^k | z_j^{k_j} \text{ for all } j \in \text{pa}(i))}{Q_{\text{MP}}(z_i^k | x, z_j \text{ for all } j \in \text{qa}(i))}$

end for

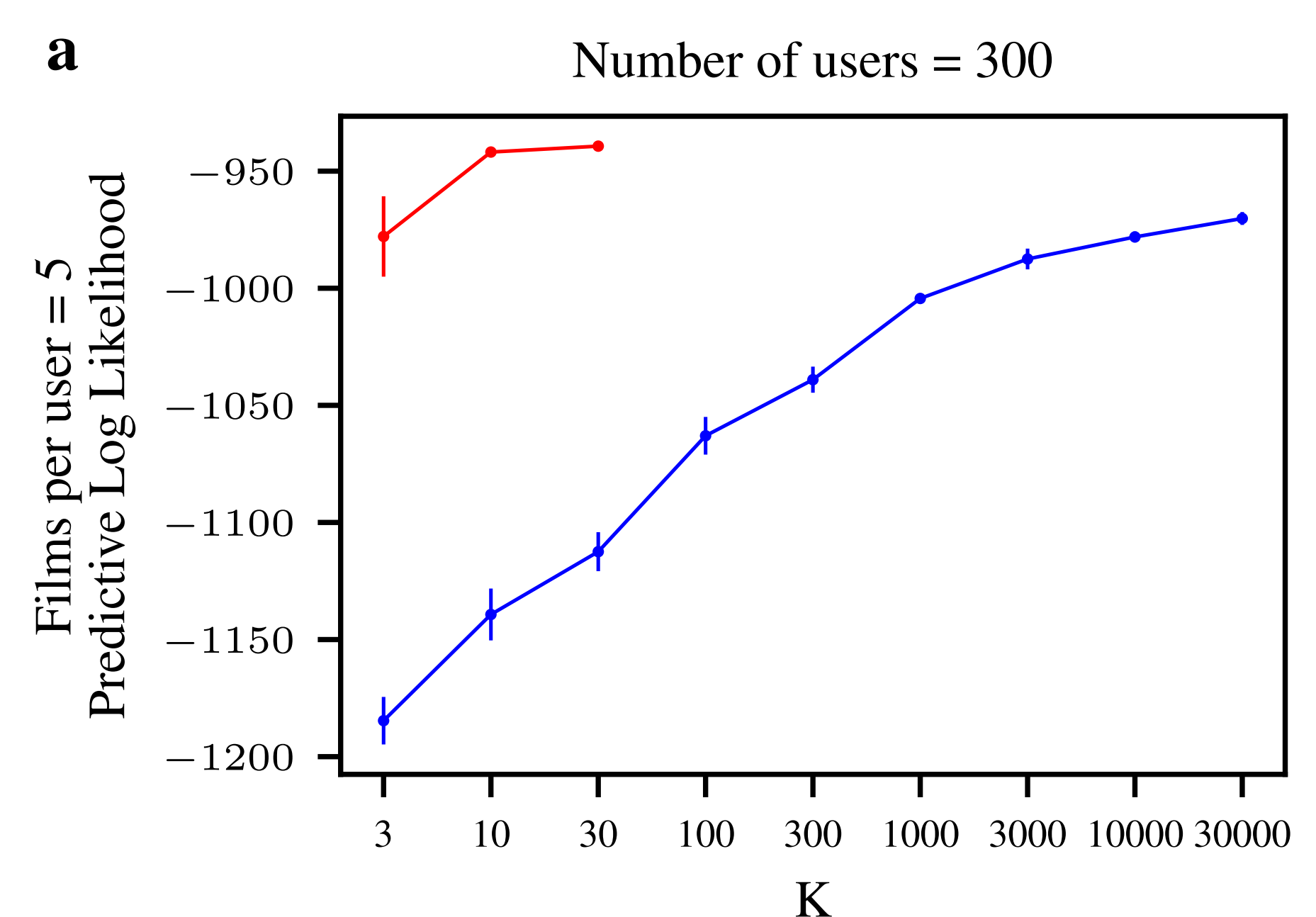
$f_{k_{\text{pa}(x)}}^x(z) \leftarrow P_\theta(x | z_j^{k_j} \text{ for all } j \in \text{pa}(x))$

$\mathcal{P}_{\text{MP}}(z) \leftarrow \frac{1}{K^n} \sum_{k^n} f_{k_{\text{pa}(x)}}^x(z) \prod_i f_{k_i, k_{\text{pa}(i)}}^i(z)$

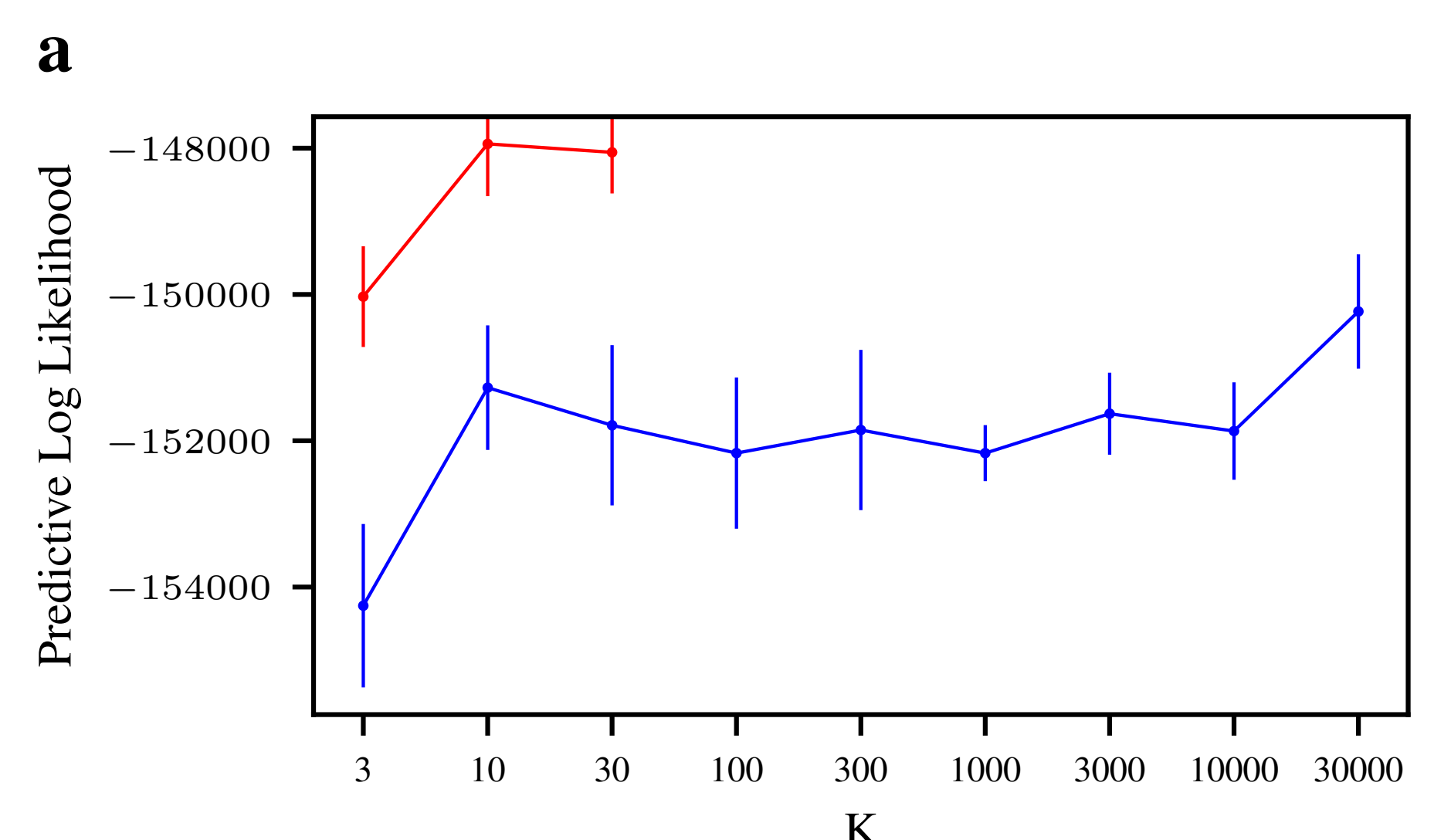
$\Delta \theta_{\text{MP}} \leftarrow \nabla_\theta \log \mathcal{P}_{\text{MP}}(z)$

$\Delta \phi_{\text{MP}} \leftarrow \nabla_\phi (-\log \mathcal{P}_{\text{MP}}(z))$

Results



MPRWS vs global RWS on the movieLens dataset.



MPRWS vs global RWS on the NYC bus breakdown dataset.

References

- [1] Laurence Aitchison. "Tensor Monte Carlo: particle methods for the GPU era". In: *Advances in Neural Information Processing Systems* 32 (2019).
- [2] Jörg Bornschein and Yoshua Bengio. "Reweighted wake-sleep". In: *arXiv preprint arXiv:1406.2751* (2014).
- [3] Sourav Chatterjee and Persi Diaconis. "The sample size required in importance sampling". In: *The Annals of Applied Probability* 28.2 (2018), pp. 1099–1135.