

A Relationship Between Sampling Methods

We organize the sampling methods described in this work in Table 1.

Table 1: Comparison of Sampling Method Designs

Algorithm	Origin	Proposal	M-H Test	Acceptance Ratio	Multiple Proposals	Reference
RWMH ¹	MCMC	Dependent	✓	M-H	✗	Duane <i>et al.</i> (1987)
HMC	MCMC	Dependent	✓	M-H	✗	
SNIS	IS	Independent	✗		✓	
IMH	MCMC	Independent	✓	M-H	✗	Naesseth <i>et al.</i> 2020
CIS	PMCMC ⁴	Independent	✓	Barker	✓	
En. MCMC ²	MCMC	Both	✓	Barker	✓	
PMP MCMC ³	MCMC	Dependent	✓	Barker	✓	Austad 2007

¹ Random-walk Metropolis-Hastings

² Ensemble MCMC

³ Parallel multiple proposals MCMC

⁴ Particle MCMC

In this paper, we designated kernels that use independent proposals and perform a Metropolis-Hastings (M-H) test as “IMH type” kernels. While the original paper of CIS does not mention it as an IMH type, we have shown in Section 3.1 that it is indeed an IMH type kernel that uses Barker’s acceptance ratio and multiple proposals per transition. This, in turn, reveals close connections with ensemble MCMC by Neal (2011a). While parallel multiple proposals MCMC by Austad (2007) also uses Barker’s acceptance ratio and multiple proposals, it only considers dependent proposals, unlike ensemble MCMC. Although in principle, it should work with independent proposals without modification.

B Additional Experimental Results

B.1 Experimental Environment

All of our experiments presented in this paper were executed on a server with 20 Intel Xeon E5-2640 CPUs and 64GB RAM. Each of the CPUs has 20 logical threads with 32k L1 cache, 256k L2 cache, and 25MB L3 cache. All of our experiments can be executed within a few days on a system with similar computational capabilities.

B.2 Additional Results of Logistic Regression Experiments

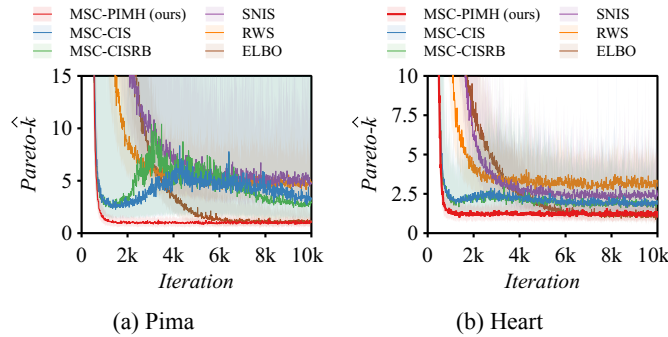


Figure 5: Pareto- \hat{k} results of logistic regression problems. The solid lines are the median of 100 repetitions while the colored regions are the 80% empirical percentiles.

549 B.3 Isotropic Gaussian Experiments

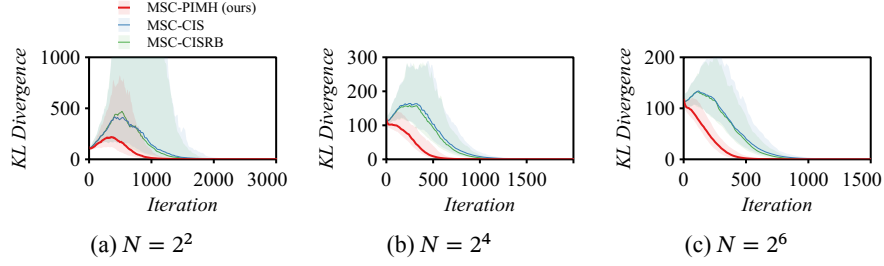


Figure 6: 100-D isotropic Gaussian example with a varying computational budget N . MSC-PIMH converges faster than MSC-CIS and MSC-CISRB regardless of N . Also, the convergence of MSC-PIMH becomes more stable/monotonic as N increases. The solid lines and colored regions are the medians and 80% percentiles computed from 100 repetitions.

550 We perform experiments with a 100-D isotropic multivariate Gaussian distribution. With Gaussian
 551 distributions, convergence can be evaluated exactly since their KL divergence is available in a closed
 552 form. We compare the performance of MSC-PIMH, MSC-CIS, and MSC-CISRB with respect to the
 553 N (number of proposals for MSC-CIS, MSC-CISRB; number of parallel chains for MSC-PIMH).
 554 The results are shown in Figure 6. While MSC-PIMH shows some level of overshoot with $N = 4$,
 555 it shows monotonic convergence with larger N . On the other hand, both MSC-CIS and MSC-
 556 CISRB overshoots even with $N = 64$. This clearly shows that PIMH enjoys better gradient estimates
 557 compared to the CIS kernel.

558 C Numerical Simulation

559 We present numerical simulations of our analyses in Section 3.3 and Section 3.4. In particular, we
 560 visualize the fact that the variance of the CIS kernel can increase with the number of proposals N
 561 when the KL divergence is large, as described in (13).

562 **Experimental Setup** We first set the target distribution as $p(z | x) = \mathcal{N}(0, 1)$ and the proposal dis-
 563 tribution as $q(z; \mu) = \mathcal{N}(\mu, 2)$ with varying mean. We measure the variance of estimating the score
 564 function $s(z, \mu) = \frac{\partial q(z; \mu)}{\partial \mu}$ using the CIS, CISRB, and PIMH kernels, given the previous Markov-
 565 chain denoted by state z_{t-1} and computational budget N . For CIS and CISRB, we set a fixed z_{t-1} ,
 566 while for PIMH, we randomly sample N samples from $z_{t-1} \sim p(z | z)$ (we obtained similar trends
 567 regardless of the distribution of z_{t-1}). The variance is estimated using 2^{14} samples from $K(z_{t-1}, \cdot)$.
 568 We report the variance across varying N and varying KL divergence between $q_\lambda(z)$ and $p(z | x)$.
 569 The latter is performed by varying the difference between the mean of the proposal and the target
 570 distributions denoted by $\Delta\mu = \mathbb{E}_{p(z|x)}[z] - \mathbb{E}_{q_\lambda}[z]$.

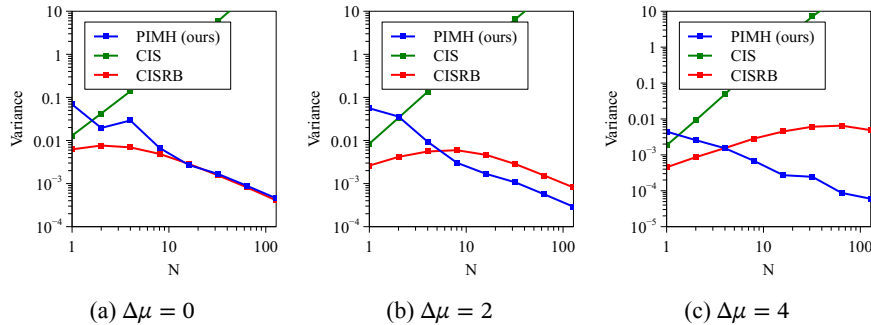


Figure 7: Conditional variance of different MCMC kernels with varying N and varying difference between the mean of the target and proposal distributions.

Results Summary The results are presented in Figure 7. We can see that, when the difference of the mean of the p and q is large, the variance of CISRB *increases* with N . This increasing trend becomes stronger as the KL divergence between p and q increases. While this simulation suggests that CISRB has much smaller variance compared to CIS, our realistic experiments in Section 4 did not reveal such levels of performance gains. Is also visible that PIMH has a slightly larger variance compared to CIS in the small N regime. This is due to the higher acceptance rate of the Metropolis-Hastings acceptance ratio used by PIMH compared to Barker’s acceptance ratio used by CIS (Peskun, 1973; Minh & Minh, 2015).

D Probabilistic Models Considered in Section 4

D.1 Hierarchical Logistic Regression

The hierarchical logistic regression used in Section 4.2 is

$$\begin{aligned}\sigma_\beta &\sim \mathcal{N}^+(0, 1.0) \\ \sigma_\alpha &\sim \mathcal{N}^+(0, 1.0) \\ \beta &\sim \mathcal{N}(\mathbf{0}, \sigma_\beta^2 \mathbf{I}) \\ \alpha &\sim \mathcal{N}(0, \sigma_\alpha^2) \\ p &\sim \mathcal{N}(\mathbf{x}_i^\top \beta + \alpha, \sigma_\alpha^2) \\ y_i &\sim \text{Bernoulli-Logit}(p)\end{aligned}$$

where \mathbf{x}_i and y_i are the predictors and binary target variable of the i th datapoints.

D.2 Stochastic Volatility

The stochastic volatility model used in Section 4.3 is

$$\begin{aligned}\mu &\sim \text{Cauchy}(0, 10) \\ \phi &\sim \text{Uniform}(-1, 1) \\ \sigma &\sim \text{Cauchy}^+(0, 5) \\ h_1 &\sim \mathcal{N}\left(0, \frac{\sigma^2}{1 - \phi^2}\right) \\ h_{t+1} &\sim \mathcal{N}(\mu + \phi(h_t - \mu), \sigma^2) \\ y_t &\sim \mathcal{N}(0, \exp(h_t))\end{aligned}$$

where y_t is the stock price at the t th point in time. We used the reparameterized version where h_t is sampled from a white multivariate Gaussian described by the Stan Development Team (2020).

D.3 Radon Hierarchical Regression

The partially pooled linear regression model used in Section 4.3 is

$$\begin{aligned}\sigma_{a_1} &\sim \text{Gamma}(\alpha = 1, \beta = 0.02) \\ \sigma_{a_2} &\sim \text{Gamma}(\alpha = 1, \beta = 0.02) \\ \sigma_y &\sim \text{Gamma}(\alpha = 1, \beta = 0.02) \\ \mu_{a_1} &\sim \mathcal{N}(0, 1) \\ \mu_{a_2} &\sim \mathcal{N}(0, 1) \\ a_{1,c} &\sim \mathcal{N}(\mu_{a_1}, \sigma_{a_1}^2) \\ a_{2,c} &\sim \mathcal{N}(\mu_{a_2}, \sigma_{a_2}^2) \\ y_i &\sim \mathcal{N}(a_{1,c_i} + a_{2,c_i} x_i, \sigma_y^2)\end{aligned}$$

where $a_{1,c}$ is the intercept at the county c , $a_{2,c}$ is the slope at the county c , c_i is the county of the i th datapoint, x_i and y_i are the floor predictor of the measurement and the measured radon level of the i th datapoint, respectively. The model pools the datapoints into their respective counties, which complicates the posterior geometry (Betancourt, 2020).

E Proofs

Detailed derivation of Equation (6)

$$\mathbb{E}_{K(\mathbf{z}_{t-1}, \mathbf{z})} [f(\mathbf{z})] \quad (16)$$

$$= \mathbb{E}_{q_\lambda} \left[\sum_{i=0}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)}) / \sum_{i=0}^N w(\mathbf{z}^{(i)}) \right] \quad (17)$$

(with a slight abuse of notation, $\mathbf{z}_{t-1} = \mathbf{z}^{(0)}$)

$$= \mathbb{E}_{q_\lambda} \left[\left(\sum_{i=1}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)}) + w(\mathbf{z}_{t-1}) f(\mathbf{z}_{t-1}) \right) / \sum_{i=0}^N w(\mathbf{z}^{(i)}) \right] \quad (18)$$

$$= \mathbb{E}_{q_\lambda} \left[\frac{\sum_{i=1}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})} + \frac{w(\mathbf{z}_{t-1})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})} f(\mathbf{z}_{t-1}) \right] \quad (19)$$

$$= \mathbb{E}_{q_\lambda} \left[\frac{\sum_{i=1}^N w(\mathbf{z}^{(i)})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})} \frac{\sum_{i=1}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)})}{\sum_{i=1}^N w(\mathbf{z}^{(i)})} + \frac{w(\mathbf{z}_{t-1})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})} f(\mathbf{z}_{t-1}) \right] \quad (20)$$

$$= \mathbb{E}_{q_\lambda} \left[\left(1 - \frac{w(\mathbf{z}_{t-1})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})} \right) \frac{\sum_{i=1}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)})}{\sum_{i=1}^N w(\mathbf{z}^{(i)})} + \frac{w(\mathbf{z}_{t-1})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})} f(\mathbf{z}_{t-1}) \right] \quad (21)$$

$$= \mathbb{E}_{q_\lambda} \left[\alpha(\mathbf{z}_{t-1}, \mathbf{z}^{(1:N)}) \frac{\sum_{i=1}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)})}{\sum_{i=1}^N w(\mathbf{z}^{(i)})} + r(\mathbf{z}_{t-1} | \mathbf{z}^{(1:N)}) f(\mathbf{z}_{t-1}) \right] \quad (22)$$

where we denote $\alpha(\mathbf{z}_{t-1}, \mathbf{z}^{(1:N)}) = \frac{\sum_{i=1}^N w(\mathbf{z}^{(i)})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})}$, $r(\mathbf{z}_{t-1} | \mathbf{z}^{(1:N)}) = \frac{w(\mathbf{z}_{t-1})}{\sum_{i=0}^N w(\mathbf{z}^{(i)})}$, and thus,

$$= \mathbb{E}_{q_\lambda} \left[\alpha(\mathbf{z}_{t-1}, \mathbf{z}^{(1:N)}) \frac{\sum_{i=1}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)})}{\sum_{i=1}^N w(\mathbf{z}^{(i)})} \right] + \mathbb{E}_{q_\lambda} [r(\mathbf{z}_{t-1} | \mathbf{z}^{(1:N)})] f(\mathbf{z}_{t-1}) \quad (23)$$

$$= \mathbb{E}_{q_\lambda} \left[\alpha(\mathbf{z}_{t-1}, \mathbf{z}^{(1:N)}) \frac{\sum_{i=1}^N w(\mathbf{z}^{(i)}) f(\mathbf{z}^{(i)})}{\sum_{i=1}^N w(\mathbf{z}^{(i)})} \right] + r(\mathbf{z}_{t-1}) f(\mathbf{z}_{t-1}) \quad (24)$$

□

Proposition 1. Assuming $w(\mathbf{z}_{t-1})$ is large enough to make $r(\mathbf{z} | \mathbf{z}^{(1:N)})$ independent of $\mathbf{z}^{(1:N)}$, the variance can be approximated by

$$\mathbb{V}_{q_\lambda} [\mathbb{E}[f | \mathbf{z}_{t-1}, \mathbf{z}^{(1:N)}]] \approx (1 - r(\mathbf{z}_{t-1}))^2 \mathbb{V}_{q_\lambda} [f_{IS} | \mathbf{z}_{t-1}]. \quad (11)$$

Proof of Proposition 1. We evaluate the variance by approximating the rejection probability as an independent constant. First,

$$\mathbb{V}_{q_\lambda} [\mathbb{E}[f | \mathbf{z}_{t-1}, \mathbf{z}^{(1:N)}]] \quad (25)$$

600 Applying (22),

$$= \mathbb{V}_{q_\lambda} \left[(1 - r(\mathbf{z}_{t-1} \mid \mathbf{z}^{(1:N)})) f_{\text{IS}} + r(\mathbf{z}_{t-1} \mid \mathbf{z}^{(1:N)}) f(\mathbf{z}_{t-1}) \mid \mathbf{z}_{t-1} \right] \quad (26)$$

$$\approx \mathbb{V}_{q_\lambda} \left[(1 - r(\mathbf{z}_{t-1})) f_{\text{IS}} + r(\mathbf{z}_{t-1}) f(\mathbf{z}_{t-1}) \mid \mathbf{z}_{t-1} \right] \quad (27)$$

$$= \mathbb{V}_{q_\lambda} \left[(1 - r(\mathbf{z}_{t-1})) f_{\text{IS}} \mid \mathbf{z}_{t-1} \right] \quad (28)$$

$$= (1 - r(\mathbf{z}_{t-1}))^2 \mathbb{V}_{q_\lambda} [f_{\text{IS}} \mid \mathbf{z}_{t-1}]. \quad (29)$$

601 The equality of (28) follows from the fact that $r(\mathbf{z}_{t-1}) f(\mathbf{z}_{t-1})$ is a constant. \square

602 **Proposition 2.** *The rejection rate $r(\mathbf{z}_{t-1})$ of a CIS sampler with N proposals is bounded below such*
 603 *that*

$$r(\mathbf{z}_{t-1}) \geq \frac{1}{1 + \frac{NZ}{w(\mathbf{z}_{t-1})}}$$

604 where $Z = \mathbb{E}_{q_\lambda(\mathbf{z})} [p(\mathbf{z}, \mathbf{x})/q_\lambda(\mathbf{z})] = \int p(\mathbf{z}, \mathbf{x}) d\mathbf{z}$ is the normalizing constant.

605 *Proof of Proposition 2.* The rejection rate $r(\mathbf{z}_{t-1})$ is given by

$$r(\mathbf{z}_{t-1}) = \mathbb{E}_{q_\lambda} \left[\frac{w(\mathbf{z}_{t-1})}{\sum_{k=1}^N w(\mathbf{z}^{(k)}) + w(\mathbf{z}_{t-1})} \right] \quad (30)$$

$$= \mathbb{E}_{q_\lambda} \left[\left(\frac{\sum_{k=1}^N w(\mathbf{z}^{(k)})}{w(\mathbf{z}_{t-1})} + 1 \right)^{-1} \right]. \quad (31)$$

606 At this point, we apply Jensen's inequality subject to the convex function $f(x) = 1/(1+x)$,

$$\geq \frac{1}{1 + \mathbb{E}_{q_\lambda} \left[\frac{\sum_{k=1}^N w(\mathbf{z}^{(k)})}{w(\mathbf{z}_{t-1})} \right]} \quad (32)$$

$$= \frac{1}{1 + \frac{1}{w(\mathbf{z}_{t-1})} \mathbb{E}_{q_\lambda} \left[\sum_{k=1}^N w(\mathbf{z}^{(k)}) \right]}. \quad (33)$$

607 From the independence of the N proposals, we obtain

$$= \frac{1}{1 + \frac{1}{w(\mathbf{z}_{t-1})} N \mathbb{E}_{q_\lambda} [w(\mathbf{z})]} \quad (34)$$

$$= \frac{1}{1 + \frac{1}{w(\mathbf{z}_{t-1})} NZ}. \quad (35)$$

608 \square

609 **Theorem 1.** *Assuming $\sup_{\mathbf{z}} p(\mathbf{z}|\mathbf{x})/q_\lambda(\mathbf{z}) = M < \infty$, the average rejection rate $r = \int r(\mathbf{z}_{t-1}) p(\mathbf{z}_{t-1} \mid$
 610 $\mathbf{x}) d\mathbf{z}_{t-1}$ of a CIS kernel with N proposals is bounded below such that*

$$r \geq \frac{1}{1 + \frac{N}{\exp(D_{\text{KL}}(p\|q_\lambda))}} - \delta,$$

611 where the sharpness of the bound is given as $0 \leq \delta \leq \frac{M}{\exp^2(D_{\text{KL}}(p\|q_\lambda))}$.

612 *Proof of Theorem 1.* We first show a simple Lemma that relates the rejection weight $w(\mathbf{z}_{t-1})$ with
 613 the KL divergence.

Lemma 1. *The average unnormalized weight of the rejection states is bounded below by the KL divergence such as*

$$Z \exp(D_{\text{KL}}(p \parallel q_\lambda)) \leq \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [w(\mathbf{z}_{t-1})].$$

Proof. By the definition of the inclusive KL divergence,

$$D_{\text{KL}}(p \parallel q_\lambda) = \int p(\mathbf{z} | \mathbf{x}) \log \frac{p(\mathbf{z} | \mathbf{x})}{q_\lambda(\mathbf{z})} d\mathbf{z} \leq \log \mathbb{E}_{p(\mathbf{z}|\mathbf{x})} \left[\frac{p(\mathbf{z} | \mathbf{x})}{q_\lambda(\mathbf{z})} \right] \quad (36)$$

$$= \log \mathbb{E}_{p(\mathbf{z}|\mathbf{x})} \left[\frac{w(\mathbf{z})}{Z} \right] \quad (37)$$

where the right-hand side follows from Jensen's inequality. By a simple change of notation, we relate (37) with the rejection states \mathbf{z}_{t-1} such as

$$D_{\text{KL}}(p \parallel q_\lambda) \leq \log \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} \left[\frac{w(\mathbf{z}_{t-1})}{Z} \right]. \quad (38)$$

Then,

$$\exp(D_{\text{KL}}(p \parallel q_\lambda)) \leq \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} \left[\frac{w(\mathbf{z}_{t-1})}{Z} \right] \quad (39)$$

$$Z \exp(D_{\text{KL}}(p \parallel q_\lambda)) \leq \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [w(\mathbf{z}_{t-1})]. \quad (40)$$

□

Now, from the result of Proposition 2,

$$\mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [r(\mathbf{z}_{t-1})] \geq \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} \left[\frac{w(\mathbf{z}_{t-1})}{w(\mathbf{z}_{t-1}) + NZ} \right] = \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [\varphi(w(\mathbf{z}_{t-1}))], \quad (41)$$

where $\varphi(x) = x/(x + NZ)$. The lower bound has the following relationship

$$\varphi(\mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [w(\mathbf{z}_{t-1})]) \geq \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [\varphi(w(\mathbf{z}_{t-1}))] \quad (42)$$

by the concavity of φ and Jensen's inequality. From this, we denote the *Jensen gap*

$$\delta = \varphi(\mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [w(\mathbf{z}_{t-1})]) - \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [\varphi(w(\mathbf{z}_{t-1}))], \quad (43)$$

where $\delta \geq 0$. Then, by applying (43) to (41),

$$\mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [r(\mathbf{z}_{t-1})] \geq \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [\varphi(w(\mathbf{z}_{t-1}))] \quad (44)$$

$$= \varphi(\mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [w(\mathbf{z}_{t-1})]) - \delta, \quad (45)$$

and by the monotonicity of φ and Lemma 1,

$$\geq \varphi(Z \exp(D_{\text{KL}}(p \parallel q_\lambda))) - \delta \quad (46)$$

$$= \frac{Z \exp(D_{\text{KL}}(p \parallel q_\lambda))}{Z \exp(D_{\text{KL}}(p \parallel q_\lambda)) + NZ} - \delta \quad (47)$$

$$= \frac{\exp(D_{\text{KL}}(p \parallel q_\lambda))}{\exp(D_{\text{KL}}(p \parallel q_\lambda)) + N} - \delta \quad (48)$$

$$= \frac{1}{1 + \frac{N}{\exp(D_{\text{KL}}(p \parallel q_\lambda))}} - \delta. \quad (49)$$

Now we discuss the Jensen gap δ , which directly gives the sharpness of our lower bound. Liao & Berg (2019, Theorem 1) have shown that, for a random variable X satisfying $P(X \in (a, b)) = 1$, where $-\infty \leq a < b \leq \infty$, and a differentiable function $\tilde{\varphi}(x)$, the following inequality holds:

$$\inf_{x \in (a, b)} h(x; \mu) \sigma^2 \leq \mathbb{E}[\tilde{\varphi}(X)] - \tilde{\varphi}(\mathbb{E}[X]), \quad \text{where} \quad h(x; \nu) = \frac{\tilde{\varphi}(x) - \tilde{\varphi}(\nu)}{(x - \nu)^2} - \frac{\tilde{\varphi}'(\nu)}{x - \nu}, \quad (50)$$

622 μ and σ^2 are the mean and variance of X , respectively. Also, Liao & Berg (2019, Lemma 1) have
 623 shown that, if $\tilde{\varphi}'(x)$ is convex, then $\inf_{x \in (a,b)} h(x; \mu) = \lim_{x \rightarrow a} h(x; \mu)$.

624 In our case, the domain is $(a, b) = (0, \infty)$ since $w(\mathbf{z}_{t-1}) > 0$. Since $\varphi'(x) = NZ/(x + NZ)^2$ is convex,
 625 we have

$$\lim_{x \rightarrow 0} h(x; \mu) = \lim_{x \rightarrow 0} \frac{1}{(x - \mu)^2} (\varphi(x) - \varphi(\mu)) - \frac{1}{x - \mu} \varphi'(\mu) \quad (51)$$

$$= \lim_{x \rightarrow 0} \frac{1}{(x - \mu)^2} \left(\frac{x}{x + NZ} - \frac{\mu}{\mu + NZ} \right) - \frac{1}{x - \mu} \left(\frac{NZ}{(\mu + NZ)^2} \right) \quad (52)$$

$$= -\frac{1}{\mu^2} \left(\frac{\mu}{\mu + NZ} \right) + \frac{1}{\mu} \left(\frac{NZ}{(\mu + NZ)^2} \right) \quad (53)$$

$$= -\frac{1}{\mu(\mu + NZ)} + \frac{NZ}{\mu(\mu + NZ)^2} \quad (54)$$

$$> -\frac{1}{\mu^2}. \quad (55)$$

626 Notice that in the context of the original problem, $\mu = \mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [w(\mathbf{z}_{t-1})]$.

627 We finally discuss the variance term σ^2 in (50). Since we assume $\sup p(\mathbf{z}|\mathbf{x})/q_\lambda(\mathbf{z}) = M < \infty$, $0 <$
 628 $\frac{p(\mathbf{z}|\mathbf{x})}{q_\lambda(\mathbf{z})} < M$ for all $\mathbf{z} \in \mathcal{Z}$. Then,

$$\sigma^2 = \mathbb{E} [w^2(\mathbf{z})] - \mathbb{E} [w(\mathbf{z})]^2 \quad (56)$$

$$= \mathbb{E} \left[\left(\frac{Z p(\mathbf{z}|\mathbf{x})}{q_\lambda(\mathbf{z})} \right)^2 \right] - \mathbb{E} \left[\frac{Z p(\mathbf{z}|\mathbf{x})}{q_\lambda(\mathbf{z})} \right]^2 \quad (57)$$

$$= Z^2 \left(\mathbb{E} \left[\left(\frac{p(\mathbf{z}|\mathbf{x})}{q_\lambda(\mathbf{z})} \right)^2 \right] - \mathbb{E} \left[\frac{p(\mathbf{z}|\mathbf{x})}{q_\lambda(\mathbf{z})} \right]^2 \right) \quad (58)$$

$$= Z^2 \mathbb{V} \left[\frac{p(\mathbf{z}|\mathbf{x})}{q_\lambda(\mathbf{z})} \right], \quad (59)$$

629 and by Bhatia & Davis (2000)'s inequality,

$$0 \leq \sigma^2 = Z^2 \mathbb{V} \left[\frac{p(\mathbf{z}|\mathbf{x})}{q_\lambda(\mathbf{z})} \right] \leq Z^2 (M - \mu) \mu. \quad (60)$$

630 By combining the results, we obtain

$$0 \leq \delta \leq -\inf_{x \in (a,b)} h(x; \mu) \sigma^2 = -\sigma^2 \lim_{x \rightarrow 0} h(x; \mu) < \frac{\sigma^2}{\mu^2} < \frac{Z^2 M}{\mu^2} = \frac{Z^2 M}{\mathbb{E}_{p(\mathbf{z}_{t-1}|\mathbf{x})} [w(\mathbf{z}_{t-1})]^2}, \quad (61)$$

631 and by Lemma 1,

$$0 \leq \delta < \frac{M}{\exp^2(D_{\text{KL}}(p \parallel q_\lambda))}. \quad (62)$$

632 □

633 **Proposition 3.** The rejection rate $r(\mathbf{z}_{t-1})$ of a IMH sampler is bounded below such that

$$r(\mathbf{z}_{t-1}) \geq 1 - \frac{Z}{w(\mathbf{z}_{t-1})} \quad \text{where} \quad Z = \int p(\mathbf{z}, \mathbf{x}) d\mathbf{z}.$$

634 *Proof of Proposition 3.* The rejection rate $r(\mathbf{z}_{t-1})$ is given by

$$r(\mathbf{z}_{t-1}) = 1 - \int \alpha(\mathbf{z}, \mathbf{z}_{t-1}) q_\lambda(\mathbf{z}) d\mathbf{z}. \quad (63)$$

635 For an IMH sampler with the Metropolis-Hastings acceptance function and independent proposals,
 636 the rejection rate is bounded such that

$$r(\mathbf{z}_{t-1}) = 1 - \int \min\left(\frac{w(\mathbf{z})}{w(\mathbf{z}_{t-1})}, 1\right) q_\lambda(\mathbf{z}) d\mathbf{z} \quad (64)$$

$$= 1 - \frac{1}{w(\mathbf{z}_{t-1})} \int \min(w(\mathbf{z}), w(\mathbf{z}_{t-1})) q_\lambda(\mathbf{z}) d\mathbf{z} \quad (65)$$

$$= 1 - \frac{1}{w(\mathbf{z}_{t-1})} \int \min\left(\frac{p(\mathbf{z}, \mathbf{x})}{q_\lambda(\mathbf{z})}, w(\mathbf{z}_{t-1})\right) q_\lambda(\mathbf{z}) d\mathbf{z} \quad (66)$$

$$= 1 - \frac{1}{w(\mathbf{z}_{t-1})} \int \min(p(\mathbf{z}, \mathbf{x}), w(\mathbf{z}_{t-1}) q_\lambda(\mathbf{z})) d\mathbf{z} \quad (67)$$

$$\geq 1 - \frac{1}{w(\mathbf{z}_{t-1})} \int p(\mathbf{z}, \mathbf{x}) d\mathbf{z} \quad (68)$$

$$= 1 - \frac{Z}{w(\mathbf{z}_{t-1})} \quad (69)$$

637 The inequality in Equation (68) follows from $\min(p(\mathbf{z}, \mathbf{x}), \cdot) \leq p(\mathbf{z}, \mathbf{x})$ for $\forall \mathbf{z} \in \mathcal{Z}$. □