

Exact Bayesian Inference (for Big Data)

Single- and Multi- Core Approaches

Murray Pollock

H Dai, P Fearnhead, AM Johansen, GO Roberts

m.pollock@warwick.ac.uk

www.warwick.ac.uk/mpollock





© marketoonist.com

- **Big Data challenge?**
 - Algorithmic 'Scalability'
- Target of interest:

$$\pi(x) \propto \prod_{i=0}^N f_i(x).$$

- Want to use [MCMC].
- Approaches:
 - Single-Core
 - Multi-Core

- Big Data challenge?
 - Algorithmic 'Scalability'
- Target of interest:

$$\pi(x) \propto \prod_{i=0}^N f_i(x).$$

- Want to use [MCMC].
- Approaches:
 - Single-Core
 - Multi-Core

- **Big Data challenge?**
 - Algorithmic 'Scalability'
- Target of interest:

$$\pi(x) \propto \prod_{i=0}^N f_i(x).$$

- **Want to use [MCMC].**
- Approaches:
 - Single-Core
 - Multi-Core

- **Big Data challenge?**
 - Algorithmic 'Scalability'
- Target of interest:

$$\pi(x) \propto \prod_{i=0}^N f_i(x).$$

- **Want to use [MCMC].**
- **Approaches:**
 - Single-Core
 - Multi-Core

- **Big Data challenge?**
 - Algorithmic 'Scalability'
- Target of interest:

$$\pi(x) \propto \prod_{i=0}^N f_i(x).$$

- **Want to use [MCMC].**
- Approaches:
 - Single-Core
 - Multi-Core

- **Big Data challenge?**
 - Algorithmic 'Scalability'
- Target of interest:

$$\pi(x) \propto \prod_{i=0}^N f_i(x).$$

- **Want to use [MCMC].**
- Approaches:
 - Single-Core
 - Multi-Core

■ Single-Core

- **Problem:** Metropolis move from $\theta \rightarrow \phi$ is accepted w.p.,

$$\min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\}$$

- **Goal:** Scalability of iterative cost.
- **Lots of work!** Pseudo-Marginal; Stochastic gradient schemes. . .
- *[1] The Scalable Langevin Exact Algorithm: Bayesian Inference for Big Data* available at: <https://arxiv.org/abs/1609.03436>

- Single-Core

- **Problem:** Metropolis move from $\theta \rightarrow \phi$ is accepted w.p.,

$$\min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\}$$

- **Goal:** Scalability of iterative cost.
- **Lots of work!** Pseudo-Marginal; Stochastic gradient schemes. . .
- *[1] The Scalable Langevin Exact Algorithm: Bayesian Inference for Big Data* available at: <https://arxiv.org/abs/1609.03436>

- Single-Core

- **Problem:** Metropolis move from $\theta \rightarrow \phi$ is accepted w.p.,

$$\min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\}$$

- **Goal:** Scalability of iterative cost.
 - **Lots of work!** Pseudo-Marginal; Stochastic gradient schemes...
- [1] *The Scalable Langevin Exact Algorithm: Bayesian Inference for Big Data* available at: <https://arxiv.org/abs/1609.03436>

- Single-Core

- **Problem:** Metropolis move from $\theta \rightarrow \phi$ is accepted w.p.,

$$\min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\}$$

- **Goal:** Scalability of iterative cost.
- **Lots of work!:** Pseudo-Marginal; Stochastic gradient schemes. . .
- [1] *The Scalable Langevin Exact Algorithm: Bayesian Inference for Big Data* available at: <https://arxiv.org/abs/1609.03436>

- Single-Core

- **Problem:** Metropolis move from $\theta \rightarrow \phi$ is accepted w.p.,

$$\min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\}$$

- **Goal:** Scalability of iterative cost.
- **Lots of work!:** Pseudo-Marginal; Stochastic gradient schemes. . .
- **[1] *The Scalable Langevin Exact Algorithm: Bayesian Inference for Big Data*** available at: <https://arxiv.org/abs/1609.03436>

■ Multi-Core

■ Solution to Single-Core:

- 1 Break data into S 'shards' (of size N/S)
- 2 Separate inferences [MCMC]
- 3 'Recombine' on 'mother-core'

■ Problem: Recombining – How do you do it?

■ Lots of work!: Consensus; Averaging; Kernel methods...

■ Constraints / Assumptions

■ [2] Bayesian Fusion: An exact and parallelisable consensus approach to unifying distributed analyses

■ Multi-Core

■ Solution to Single-Core:

- 1 Break data into S 'shards' (of size N/S)
- 2 Separate inferences [MCMC]
- 3 'Recombine' on 'mother-core'

- Problem: Recombining – How do you do it?
- Lots of work!: Consensus; Averaging; Kernel methods...
- Constraints / Assumptions

- *[2] Bayesian Fusion: An exact and parallelisable consensus approach to unifying distributed analyses*

■ Multi-Core

■ Solution to Single-Core:

- 1 Break data into S 'shards' (of size N/S)
- 2 Separate inferences [MCMC]
- 3 'Recombine' on 'mother-core'

■ Problem: Recombining – How do you do it?

- Lots of work!: Consensus; Averaging; Kernel methods...
- Constraints / Assumptions

- *[2] Bayesian Fusion: An exact and parallelisable consensus approach to unifying distributed analyses*

■ Multi-Core

■ Solution to Single-Core:

- 1 Break data into S 'shards' (of size N/S)
- 2 Separate inferences [MCMC]
- 3 'Recombine' on 'mother-core'

■ Problem: Recombining – How do you do it?

■ Lots of work!: Consensus; Averaging; Kernel methods...

■ Constraints / Assumptions

■ [2] Bayesian Fusion: An exact and parallelisable consensus approach to unifying distributed analyses

■ Multi-Core

■ Solution to Single-Core:

- 1 Break data into S 'shards' (of size N/S)
- 2 Separate inferences [MCMC]
- 3 'Recombine' on 'mother-core'

■ Problem: Recombining – How do you do it?

■ Lots of work!: Consensus; Averaging; Kernel methods...

■ Constraints / Assumptions

- *[2] Bayesian Fusion: An exact and parallelisable consensus approach to unifying distributed analyses*

■ Multi-Core

■ Solution to Single-Core:

- 1 Break data into S 'shards' (of size N/S)
- 2 Separate inferences [MCMC]
- 3 'Recombine' on 'mother-core'

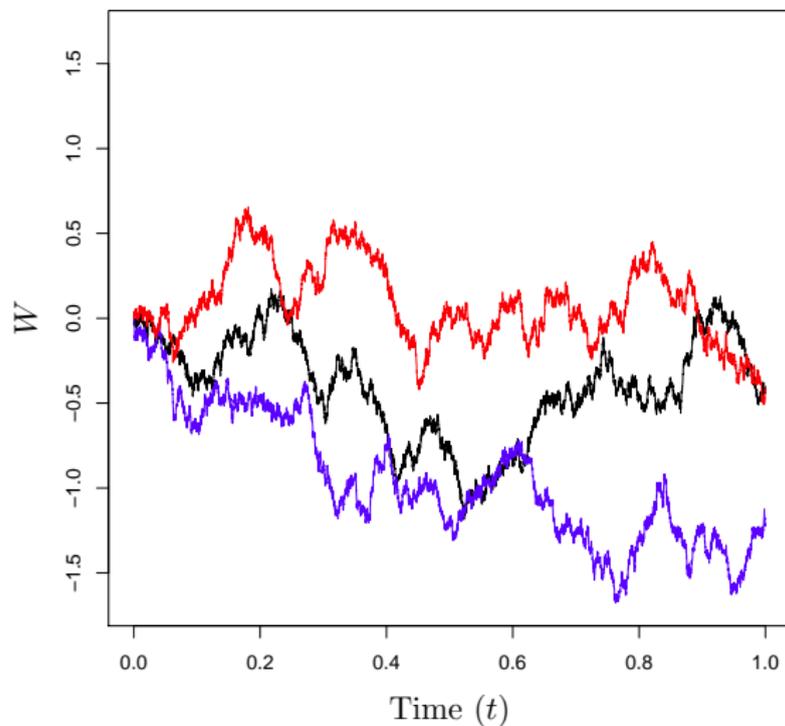
■ Problem: Recombining – How do you do it?

- Lots of work!: Consensus; Averaging; Kernel methods. . .
- Constraints / Assumptions

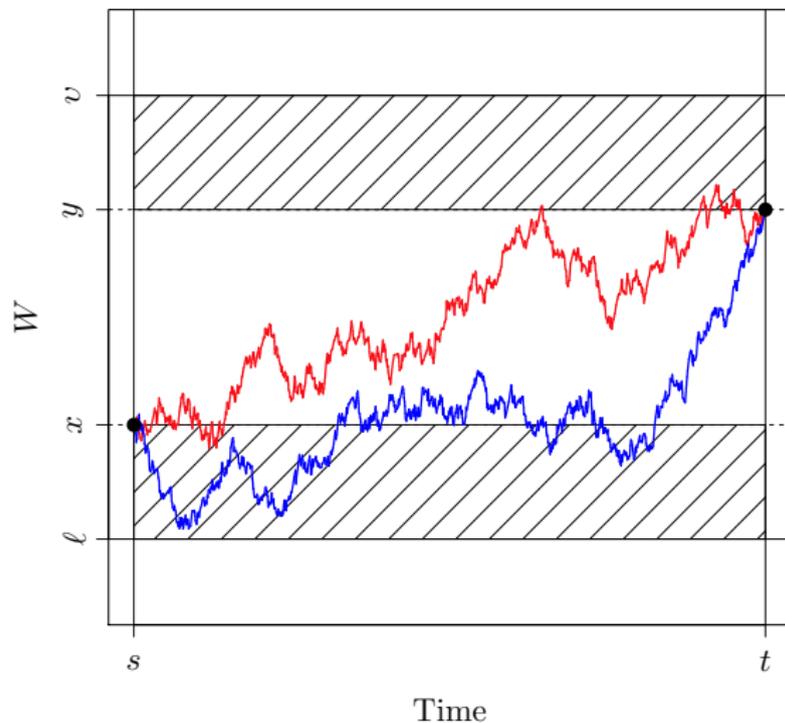
■ *[2] Bayesian Fusion: An exact and parallelisable consensus approach to unifying distributed analyses*

0 - Retrospective Trust / Tricks

Brownian Motion



Brownian Motion



Path-space Rejection Sampling:

- We want $X \sim \mathbb{Q}$ where:

$$\mathbb{Q} : dX_t = \alpha(X_t) dt + \Lambda^{1/2} dB_t, \quad X_0 = x \in \mathbb{R}^d, t \in [0, T]$$

- **Discretisation Free Approach!:** Path-space Rejection Sampler (PRS)
(see [arXiv 1302.6964](https://arxiv.org/abs/1302.6964) for details)
 - 1 $X_T \sim h_T(X_0)$
 - 2 $X^{\text{fin}} \sim \mathbb{P} | X_T^*$ (eg \mathbb{W} or \mathbb{D})
 - 3 (Accept / Reject)** / Assign Weight**

Path-space Rejection Sampling:

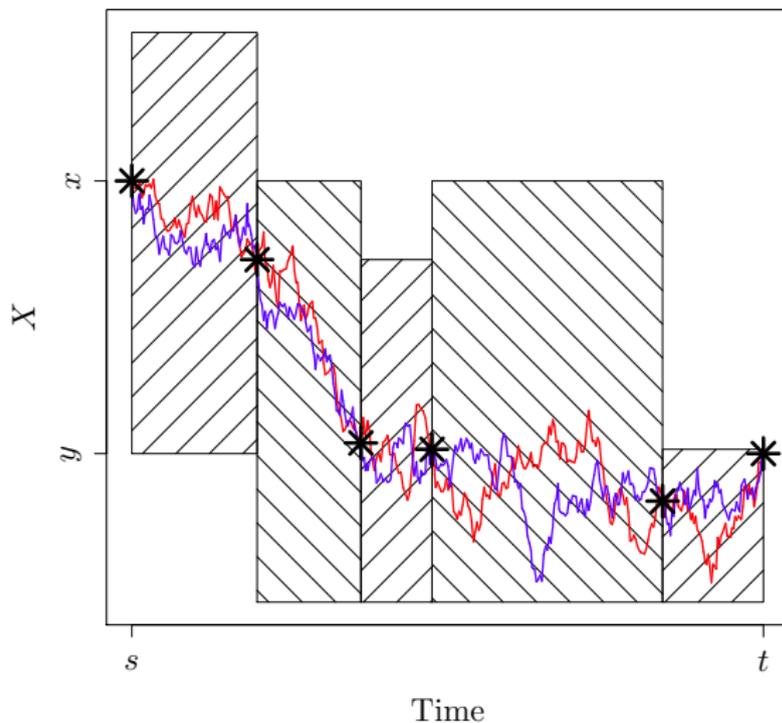
- We want $X \sim \mathbb{Q}$ where:

$$\mathbb{Q} : dX_t = \alpha(X_t) dt + \Lambda^{1/2} dB_t, \quad X_0 = x \in \mathbb{R}^d, t \in [0, T]$$

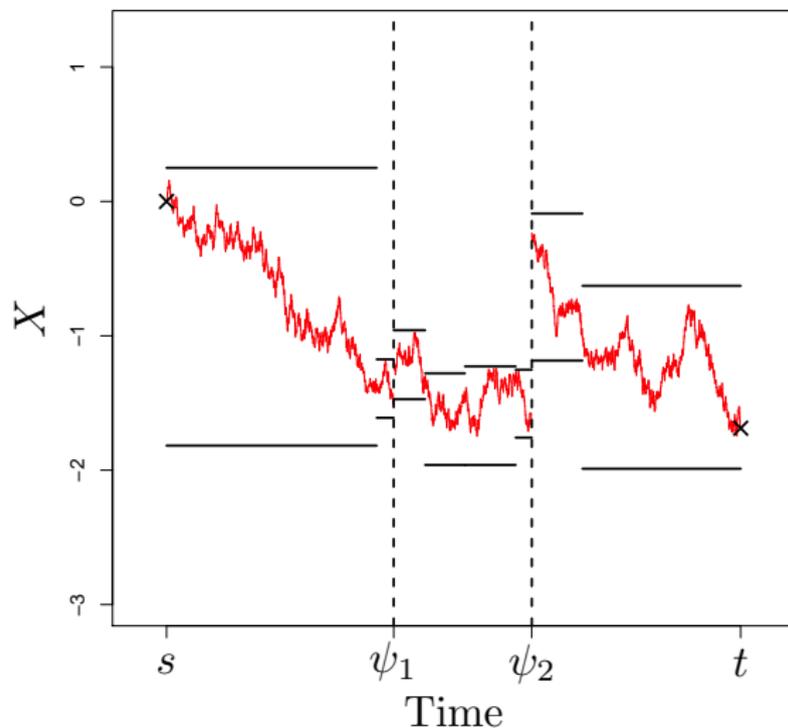
- **Discretisation Free Approach!:** Path-space Rejection Sampler (PRS)
(see [arXiv 1302.6964](https://arxiv.org/abs/1302.6964) for details)

- 1 $X_T \sim h_T(X_0)$
- 2 $X^{\text{fin}} \sim \mathbb{P} | X_T^*$ (eg \mathbb{W} or \mathbb{D})
- 3 (Accept / Reject)** / Assign Weight**

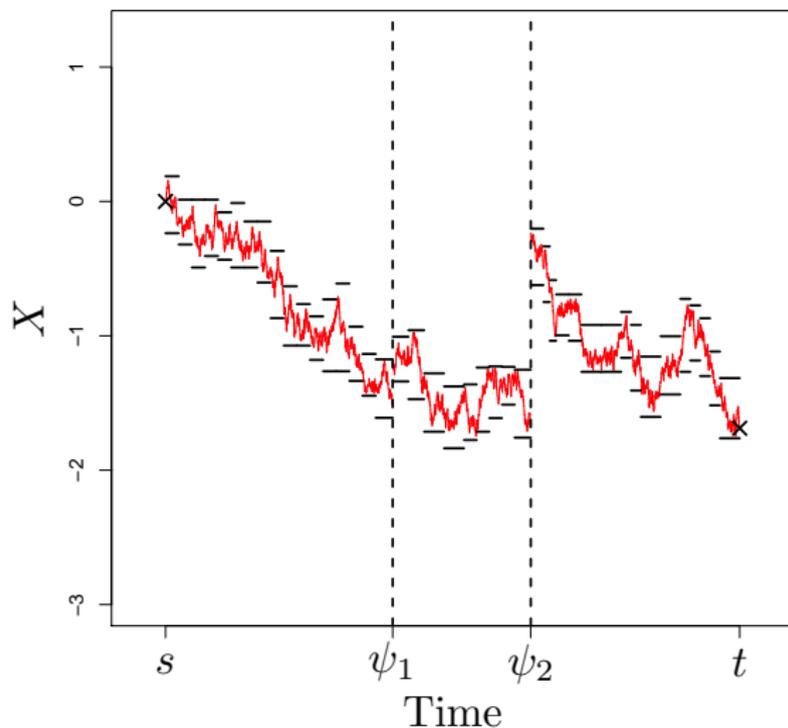
Path-space Rejection Sampling



Path-space Rejection Sampling



Path-space Rejection Sampling



Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution v
 - Direct statistical exploitation... $v \equiv \pi(\mathbb{I})$
 - Langevin + $v \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv v^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y)}_{\propto h} \cdot \underbrace{v^{1/2}(y)}_{\in [0,1]} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow v$
 - $v \equiv \pi^2(\mathbb{D}\mathbb{I})$
- If $X_0 \sim v$, then $\forall t > 0, X_t \sim v$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution v
 - Direct statistical exploitation... $v \equiv \pi(\mathbb{I})$
 - Langevin + $v \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv v^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y)}_{\propto h} \cdot \underbrace{v^{1/2}(y)}_{\in [0,1]} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow v$
 - $v \equiv \pi^2(\mathbb{D}\mathbb{I})$
- If $X_0 \sim v$, then $\forall t > 0, X_t \sim v$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution ν
 - Direct statistical exploitation... $\nu \equiv \pi(\mathbb{I})$
 - Langevin + $\nu \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv \nu^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y)}_{\propto h} \cdot \underbrace{\nu^{1/2}(y)}_{\in [0,1]} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow \nu$
 - $\nu \equiv \pi^2(\mathbb{D}\mathbb{I})$
- If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution ν
 - Direct statistical exploitation... $\nu \equiv \pi(\mathbb{L})$
 - Langevin + $\nu \equiv \pi$ + Discretisation + Correction \implies MALA
 - PRS Class (however [1] $y := X_T \sim h \equiv \nu^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y)}_{\propto h} \cdot \underbrace{\nu^{1/2}(y)}_{\in [0,1]} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow \nu$
 - $\nu \equiv \pi^2(\mathbb{D}\mathbb{L})$
- If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution v
 - Direct statistical exploitation... $v \equiv \pi(\mathbb{L})$
 - Langevin + $v \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv v^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y)}_{\propto h} \cdot \underbrace{v^{1/2}(y)}_{\in [0,1]} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow v$
 - $v \equiv \pi^2(\mathbb{D}\mathbb{L})$
- If $X_0 \sim v$, then $\forall t > 0, X_t \sim v$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution ν
 - Direct statistical exploitation... $\nu \equiv \pi(\mathbb{L})$
 - Langevin + $\nu \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv \nu^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y) \cdot \nu^{1/2}(y)}_{\propto h} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow \nu$
 - $\nu \equiv \pi^2(\mathbb{D}\mathbb{L})$
- If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution v
 - Direct statistical exploitation... $v \equiv \pi(\mathbb{L})$
 - Langevin + $v \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv v^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y) \cdot v^{1/2}(y)}_{\propto h} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow v$
 - $v \equiv \pi^2(\mathbb{D}\mathbb{L})$
- If $X_0 \sim v$, then $\forall t > 0, X_t \sim v$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution v
 - Direct statistical exploitation... $v \equiv \pi(\mathbb{L})$
 - Langevin + $v \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv v^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y) \cdot v^{1/2}(y)}_{\propto h} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow v$
 - $v \equiv \pi^2(\mathbb{D}\mathbb{L})$
- If $X_0 \sim v$, then $\forall t > 0, X_t \sim v$

Langevin Diffusion:

- In \mathbb{Q} set $\alpha(X_t) := \frac{1}{2} \Delta \nabla \log v(X_t)$
- Invariant distribution ν
 - Direct statistical exploitation... $\nu \equiv \pi(\mathbb{L})$
 - Langevin + $\nu \equiv \pi$ + Discretisation + Correction \implies MALA
- PRS Class (however [1] $y := X_T \sim h \equiv \nu^{1/2} \dots$)
 - $\lim_{T \rightarrow \infty} p_T(x, y) = \underbrace{w_T(x, y) \cdot \nu^{1/2}(y)}_{\propto h} \cdot \underbrace{P(X)}_{\in [0,1]} \rightarrow \nu$
 - $\nu \equiv \pi^2(\mathbb{D}\mathbb{L})$
- If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$

1 - Single Core: Quasi-Stationary Monte Carlo

Quasi-Stationary Monte Carlo

- Consider Brownian motion, killed at τ with intensity

$$\kappa(x) = \frac{\|\nabla \log \pi(x)\|^2 + \Delta \log \pi(x)}{2} - \ell \in \mathbb{R}_{\geq 0},$$

and the quasi-limiting distribution

$$\lim_{t \rightarrow \infty} \mathcal{L}(X_t | \tau > t).$$

- Under weak regularity conditions has quasi-stationary distribution π .
- Statistical Interpretation:
 - Big Data? \rightarrow Subsampling
 - Implementation? \rightarrow ScaLE

Quasi-Stationary Monte Carlo

- Consider Brownian motion, killed at τ with intensity

$$\kappa(x) = \frac{\|\nabla \log \pi(x)\|^2 + \Delta \log \pi(x)}{2} - \ell \in \mathbb{R}_{\geq 0},$$

and the quasi-limiting distribution

$$\lim_{t \rightarrow \infty} \mathcal{L}(X_t | \tau > t).$$

- Under weak regularity conditions has quasi-stationary distribution π .
- Statistical Interpretation:
 - Big Data? \rightarrow Subsampling
 - Implementation? \rightarrow ScaLE

Quasi-Stationary Monte Carlo

- Consider Brownian motion, killed at τ with intensity

$$\kappa(x) = \frac{\|\nabla \log \pi(x)\|^2 + \Delta \log \pi(x)}{2} - \ell \in \mathbb{R}_{\geq 0},$$

and the quasi-limiting distribution

$$\lim_{t \rightarrow \infty} \mathcal{L}(X_t | \tau > t).$$

- Under weak regularity conditions has quasi-stationary distribution π .
- Statistical Interpretation:
 - Big Data? \rightarrow Subsampling
 - Implementation? \rightarrow ScaLE

Quasi-Stationary Monte Carlo

- Consider Brownian motion, killed at τ with intensity

$$\kappa(x) = \frac{\|\nabla \log \pi(x)\|^2 + \Delta \log \pi(x)}{2} - \ell \in \mathbb{R}_{\geq 0},$$

and the quasi-limiting distribution

$$\lim_{t \rightarrow \infty} \mathcal{L}(X_t | \tau > t).$$

- Under weak regularity conditions has quasi-stationary distribution π .
- Statistical Interpretation:
 - Big Data? \rightarrow Subsampling
 - Implementation? \rightarrow ScaLE

Quasi-Stationary Monte Carlo

- Consider Brownian motion, killed at τ with intensity

$$\kappa(x) = \frac{\|\nabla \log \pi(x)\|^2 + \Delta \log \pi(x)}{2} - \ell \in \mathbb{R}_{\geq 0},$$

and the quasi-limiting distribution

$$\lim_{t \rightarrow \infty} \mathcal{L}(X_t | \tau > t).$$

- Under weak regularity conditions has quasi-stationary distribution π .
- Statistical Interpretation:
 - Big Data? \rightarrow Subsampling
 - Implementation? \rightarrow ScaLE

1.2 - Subsampling

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- QSMC \equiv Simulating BM + inhomogeneous Poisson Process κ
- Evaluating κ is $O(N)$.
- * If $\forall x, \kappa(x) \leq K$ (requires localisation argument) then:
 - Simulating $PP(\kappa(x)) \equiv$ Simulating $PP(K)$ and accepting w.p. $\kappa(X_t)/K$.
- We can make our algorithm worse (!) by choosing $\tilde{K} \geq K \dots$
- Remark on coins
- Suppose $\exists A \sim \mathcal{A}, \tilde{\kappa}_A(\cdot) \in [0, \tilde{K}]$ such that $\mathbb{E}_{\mathcal{A}}[\tilde{\kappa}_A(x)/\tilde{K}] = \kappa(x)/\tilde{K}$ then:
 - Simulating $A \sim \mathcal{A} PP(\tilde{K})$ and accepting w.p. $\tilde{\kappa}_A(X_t)/\tilde{K} \equiv *$.

- Scalability \equiv Finding $A \sim \mathcal{A}$ and $\tilde{\kappa}_A(\cdot)$ which are $O(1)$ (trivial), such that $\tilde{K}/K \geq 1$ scales well...
- Intuition is the diffusion drift is a sum:

$$\nabla \log \pi(x) = \sum_{i=0}^N \nabla \log f_i(x)$$

- We require control variates for good scaling of \tilde{K}/K ... (omitted)

- Scalability \equiv Finding $A \sim \mathcal{A}$ and $\tilde{\kappa}_A(\cdot)$ which are $O(1)$ (trivial), such that $\tilde{K}/K \geq 1$ scales well...
- Intuition is the diffusion drift is a sum:

$$\nabla \log \pi(x) = \sum_{i=0}^N \nabla \log f_i(x)$$

- We require control variates for good scaling of \tilde{K}/K ... (omitted)

- Scalability \equiv Finding $A \sim \mathcal{A}$ and $\tilde{\kappa}_A(\cdot)$ which are $O(1)$ (trivial), such that $\tilde{K}/K \geq 1$ scales well...
- Intuition is the diffusion drift is a sum:

$$\nabla \log \pi(x) = \sum_{i=0}^N \nabla \log f_i(x)$$

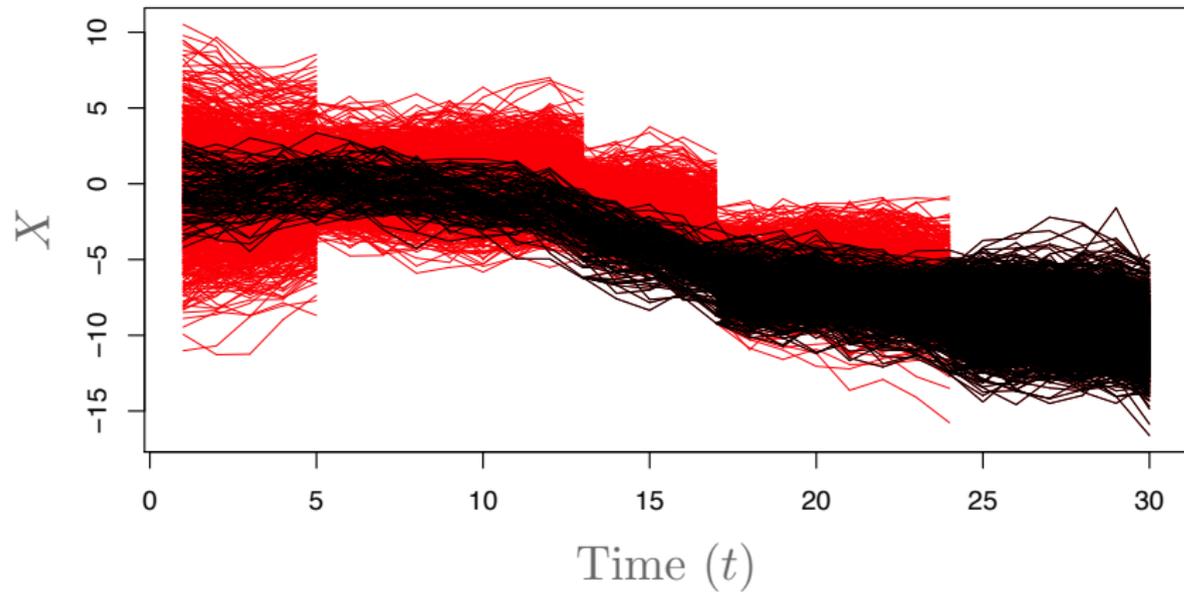
- We require control variates for good scaling of \tilde{K}/K ... (omitted)

1.3 - Single-Core: ScaLE

- **Implementational Problem: Trajectory death!**
- **First Approach: Scalable Langevin Exact Algorithm (ScaLE)**
 - Continuous time multi-level splitting / Importance sampling QSMC + SMC + Resampling

- **Implementational Problem:** Trajectory death!
- **First Approach: Scalable Langevin Exact Algorithm (ScaLE)**
 - Continuous time multi-level splitting / Importance sampling QSMC + SMC + Resampling

- **Implementational Problem:** Trajectory death!
- **First Approach: Scalable Langevin Exact Algorithm (ScaLE)**
 - Continuous time multi-level splitting / Importance sampling QSMC + SMC + Resampling



1.4 - Summary

■ Summary...:

- QSMC: 'Exact' Bayesian Inference
- No intrinsic cost for exactness.
- ScaLE's well!

■ Missing Bits...:

- Localisation
- Theory: QSMC; (SMC-) ScaLE; Re-ScaLE.
- Scaling: Dimensionality; Control-Variate...
- Implementational Details

■ Summary...:

- QSMC: 'Exact' Bayesian Inference
 - No intrinsic cost for exactness.
 - ScaLE's well!

■ Missing Bits...:

- Localisation
- Theory: QSMC; (SMC-) ScaLE; Re-ScaLE.
- Scaling: Dimensionality; Control-Variate...
- Implementational Details

■ Summary...:

- QSMC: 'Exact' Bayesian Inference
- No intrinsic cost for exactness.
- ScaLE's well!

■ Missing Bits...:

- Localisation
- Theory: QSMC; (SMC-) ScaLE; Re-ScaLE.
- Scaling: Dimensionality; Control-Variate...
- Implementational Details

■ Summary...:

- QSMC: 'Exact' Bayesian Inference
- No intrinsic cost for exactness.
- ScaLE's well!

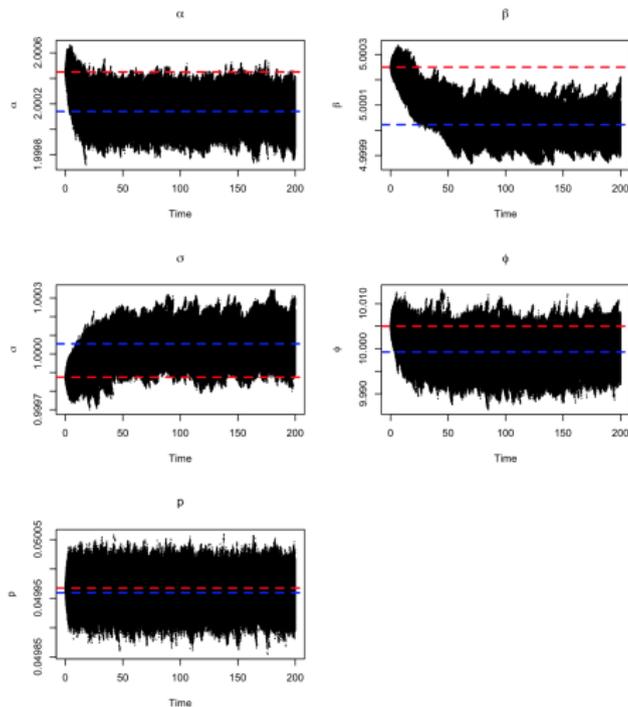
■ Missing Bits...:

- Localisation
- Theory: QSMC; (SMC-) ScaLE; Re-ScaLE.
- Scaling: Dimensionality; Control-Variate...
- Implementational Details

- Summary...:
 - QSMC: 'Exact' Bayesian Inference
 - No intrinsic cost for exactness.
 - ScaLE's well!
- Missing Bits...:
 - Localisation
 - Theory: QSMC; (SMC-) ScaLE; Re-ScaLE.
 - Scaling: Dimensionality; Control-Variate...
 - Implementational Details

Example

2^{27} dataset, contaminated regression model



2 - Multi-Core: Bayesian Fusion

■ Recall Target:

$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

■ C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.

■ Simple Approach... [Think (A)BC]

1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.

2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.

3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).

■ Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:

■ $\mathbb{L}_1, \dots, \mathbb{L}_C, \text{DL}_1, \dots, \text{DL}_C \dots$

- Recall Target:

$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

- C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.

- Simple Approach... [Think (A)BC]

- 1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.

- 2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.

- 3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).

- Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:

- $\mathbb{L}_1, \dots, \mathbb{L}_C, \text{DL}_1, \dots, \text{DL}_C \dots$

- Recall Target:

$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

- C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.

- Simple Approach... [Think (A)BC]

- 1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.

- 2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.

- 3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).

- Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:

- $\mathbb{L}_1, \dots, \mathbb{L}_C, \text{DL}_1, \dots, \text{DL}_C \dots$

■ Recall Target:

$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

■ C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.

■ Simple Approach... [Think (A)BC]

1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.

2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.

3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).

■ Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:

■ $\mathbb{L}_1, \dots, \mathbb{L}_C, \text{DL}_1, \dots, \text{DL}_C \dots$

- Recall Target:

$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

- C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.
- Simple Approach... [Think (A)BC]
 - 1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.
 - 2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.
 - 3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).
- Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:
 - $\mathbb{L}_1, \dots, \mathbb{L}_C, \text{DL}_1, \dots, \text{DL}_C \dots$

- Recall Target:

$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

- C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.

- Simple Approach... [Think (A)BC]

- 1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.

- 2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.

- 3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).

- Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:

- $\mathbb{L}_1, \dots, \mathbb{L}_C, \mathbb{DL}_1, \dots, \mathbb{DL}_C \dots$

- Recall Target:

$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

- C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.

- Simple Approach... [Think (A)BC]

- 1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.
- 2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.
- 3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).

- Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:

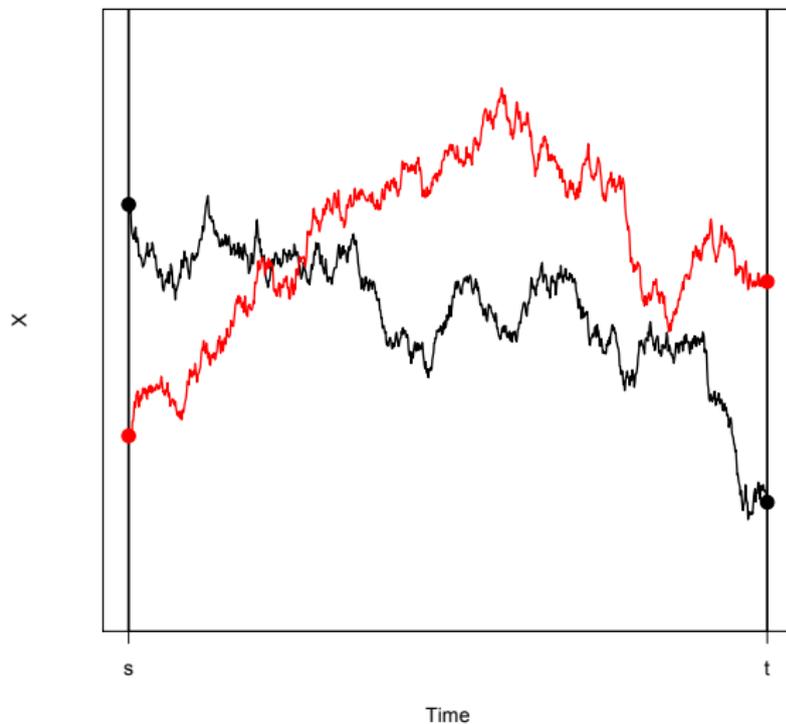
- $\mathbb{L}_1, \dots, \mathbb{L}_C, \text{DL}_1, \dots, \text{DL}_C \dots$

- Recall Target:

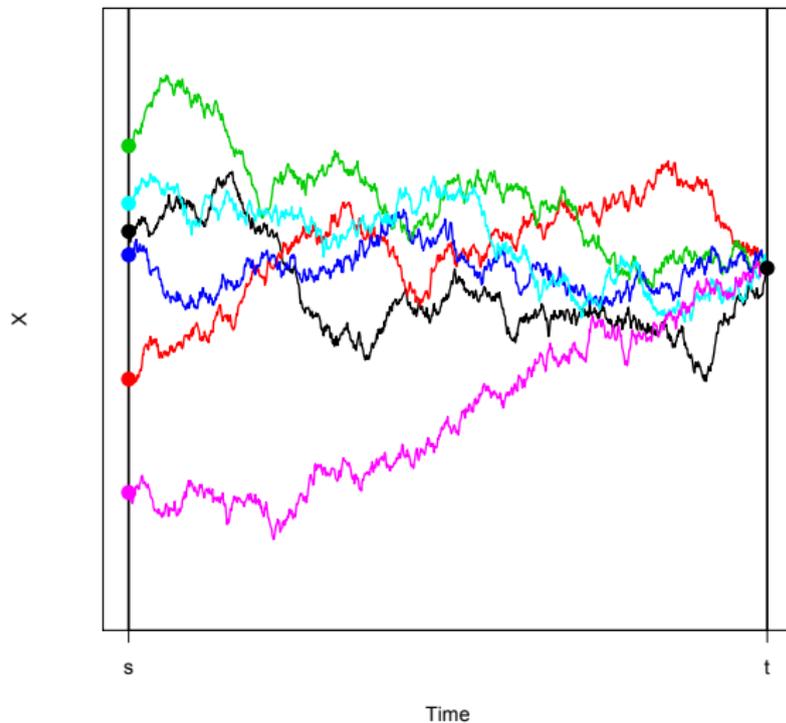
$$\pi(x) \propto \prod_{c=1}^C f_c(x).$$

- C - Number of cores / experts / 'views' ...; f_c - Sub-posterior.
- Simple Approach... [Think (A)BC]
 - 1 Simulate $X^{(1)} \sim f_1, X^{(2)} \sim f_2, \dots, X^{(C)} \sim f_C$.
 - 2 Accept if $X^{(1)} = X^{(2)} = \dots = X^{(C)}$, else go to 1/.
 - 3 Return $X := X^{(1)}$ ($\sim \prod_{i=1}^C f_i \propto \pi$).
- Recall Langevin: If $X_0 \sim \nu$, then $\forall t > 0, X_t \sim \nu$:
 - $\mathbb{L}_1, \dots, \mathbb{L}_C, \mathbb{D}\mathbb{L}_1, \dots, \mathbb{D}\mathbb{L}_C \dots$

Fusion Idea



Fusion Actual



- Fusion Measure ($\mathfrak{x} \in \Omega_0$)

$$d\mathbb{F}(\mathfrak{x}) \propto d\left(\times_{c=1}^C \mathbb{D}\mathbb{L}_c^{\mathbf{x}_0^{(c)}, \mathbf{y}_T}\right)(\mathfrak{x}) \cdot \prod_{c=1}^C \left[f_c^2(\mathbf{x}_0^{(c)}) p_{T,c}^{\text{dl}}(\mathbf{y}_T | \mathbf{x}_0^{(c)}) \cdot \frac{1}{f_c(\mathbf{y}_T)} \right],$$

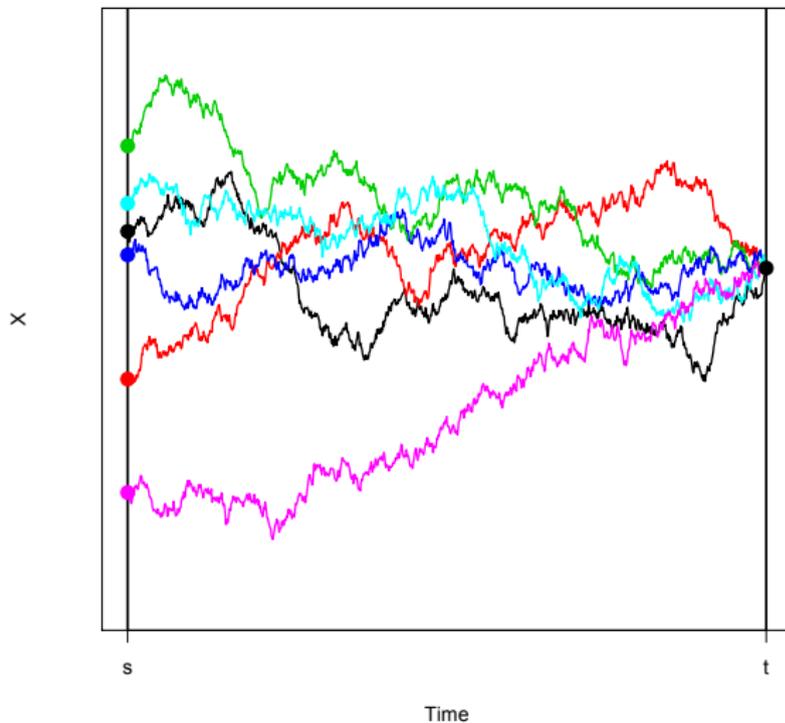
- Key Idea: If $\mathfrak{x} \sim \mathbb{F}$, then $\mathfrak{x}_T \sim \prod_{c=1}^C f_c \propto \pi$ (!)

- Fusion Measure ($\mathfrak{x} \in \Omega_0$)

$$d\mathbb{F}(\mathfrak{x}) \propto d\left(\times_{c=1}^C \mathbb{D}L_c^{\mathbf{x}_0^{(c)}, \mathbf{y}_T}\right)(\mathfrak{x}) \cdot \prod_{c=1}^C \left[f_c^2(\mathbf{x}_0^{(c)}) p_{T,c}^{\text{dl}}(\mathbf{y}_T | \mathbf{x}_0^{(c)}) \cdot \frac{1}{f_c(\mathbf{y}_T)} \right],$$

- Key Idea: If $\mathfrak{x} \sim \mathbb{F}$, then $\mathfrak{x}_T \sim \prod_{c=1}^C f_c \propto \pi$ (!)

Some Details



- 'Standard' Multi-Core Problem $\equiv \mathfrak{X} \sim \mathbb{F}$ (with practical constraints)
 - Rejection Sampling! Possible proposals $\mathfrak{X} \sim \mathbb{P}$, w.p. $P(\mathfrak{X})$:
 - 'Brownian':

$$d\mathbb{P}^{\text{bm}}(\mathfrak{X}) \propto d\left(\times_{c=1}^C \mathbb{W}_c^{\mathbf{X}_0^{(c)}, \mathbf{y}_T}\right)(\mathfrak{X}) \cdot h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T), \quad \mathfrak{X} \in \Omega_0$$

- 'Standard' Multi-Core Problem '≡' $\mathfrak{X} \sim \mathbb{F}$ (with practical constraints)
 - Rejection Sampling! Possible proposals $\mathfrak{X} \sim \mathbb{P}$, w.p. $P(\mathfrak{X})$:
 - 'Brownian':

$$d\mathbb{P}^{\text{bm}}(\mathfrak{X}) \propto d\left(\prod_{c=1}^C \mathbb{W}_c^{\mathbf{X}_0^{(c)}, \mathbf{y}_T}\right)(\mathfrak{X}) \cdot h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T), \quad \mathfrak{X} \in \Omega_0$$

- 'Standard' Multi-Core Problem '≡' $\mathfrak{X} \sim \mathbb{F}$ (with practical constraints)
 - Rejection Sampling! Possible proposals $\mathfrak{X} \sim \mathbb{P}$, w.p. $P(\mathfrak{X})$:
 - 'Brownian':

$$d\mathbb{P}^{\text{bm}}(\mathfrak{X}) \propto d\left(\times_{c=1}^C \mathbb{W}_c^{\mathbf{X}_0^{(c)}, \mathbf{y}_T}\right)(\mathfrak{X}) \cdot h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T), \quad \mathfrak{X} \in \Omega_0$$

■ Simple 'Brownian' Case:

- 'Optimal' $h_T^{\text{bm}}(\cdot, \cdot)$:

$$h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T) \propto \underbrace{\left[\prod_{c=1}^C f_c(\mathbf{X}_0^{(c)}) \right]}_{\text{initial core draws}} \underbrace{\exp\left(-\frac{C \cdot \|\mathbf{y}_T - \bar{\mathbf{X}}_0\|^2}{2T}\right) \cdot \exp\left(-\frac{C\sigma^2}{2T}\right)}_{\text{end point draw}}$$

- Need RS for $h_T^{\text{bm}}(\cdot, \cdot)$ end point.
- Accept with probability

$$P(x) := \exp\left[-\sum_{c=1}^C \int_0^T \kappa_c(\mathbf{X}_t^{(c)}) dt\right] \in [0, 1]$$

- Exact ('Talking') vs. Approximate ('Silent' / 'Lecture') Fusion
- Remark: 'Ornstein-Uhlenbeck' special case

■ Simple 'Brownian' Case:

- 'Optimal' $h_T^{\text{bm}}(\cdot, \cdot)$:

$$h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T) \propto \underbrace{\left[\prod_{c=1}^C f_c(\mathbf{X}_0^{(c)}) \right]}_{\text{initial core draws}} \underbrace{\exp\left(-\frac{C \cdot \|\mathbf{y}_T - \bar{\mathbf{X}}_0\|^2}{2T}\right) \cdot \exp\left(-\frac{C\sigma^2}{2T}\right)}_{\text{end point draw}}$$

- Need RS for $h_T^{\text{bm}}(\cdot, \cdot)$ end point.
- Accept with probability

$$P(x) := \exp\left[-\sum_{c=1}^C \int_0^T \kappa_c(\mathbf{X}_t^{(c)}) dt\right] \in [0, 1]$$

- Exact ('Talking') vs. Approximate ('Silent' / 'Lecture') Fusion
- Remark: 'Ornstein-Uhlenbeck' special case

- Simple 'Brownian' Case:

- 'Optimal' $h_T^{\text{bm}}(\cdot, \cdot)$:

$$h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T) \propto \underbrace{\left[\prod_{c=1}^C f_c(\mathbf{X}_0^{(c)}) \right]}_{\text{initial core draws}} \underbrace{\exp\left(-\frac{C \cdot \|\mathbf{y}_T - \bar{\mathbf{X}}_0\|^2}{2T}\right) \cdot \exp\left(-\frac{C\sigma^2}{2T}\right)}_{\text{end point draw}}$$

- Need RS for $h_T^{\text{bm}}(\cdot, \cdot)$ end point.

- Accept with probability

$$P(x) := \exp\left[-\sum_{c=1}^C \int_0^T \kappa_c(\mathbf{X}_t^{(c)}) dt\right] \in [0, 1]$$

- Exact ('Talking') vs. Approximate ('Silent' / 'Lecture') Fusion

- Remark: 'Ornstein-Uhlenbeck' special case

- Simple 'Brownian' Case:

- 'Optimal' $h_T^{\text{bm}}(\cdot, \cdot)$:

$$h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T) \propto \underbrace{\left[\prod_{c=1}^C f_c(\mathbf{X}_0^{(c)}) \right]}_{\text{initial core draws}} \underbrace{\exp\left(-\frac{C \cdot \|\mathbf{y}_T - \bar{\mathbf{X}}_0\|^2}{2T}\right) \cdot \exp\left(-\frac{C\sigma^2}{2T}\right)}_{\text{end point draw}}$$

- Need RS for $h_T^{\text{bm}}(\cdot, \cdot)$ end point.
- Accept with probability

$$P(\mathfrak{x}) := \exp\left[-\sum_{c=1}^C \int_0^T \kappa_c(\mathbf{X}_t^{(c)}) dt\right] \in [0, 1]$$

- Exact ('Talking') vs. Approximate ('Silent' / 'Lecture') Fusion

- Remark: 'Ornstein-Uhlenbeck' special case

■ Simple 'Brownian' Case:

- 'Optimal' $h_T^{\text{bm}}(\cdot, \cdot)$:

$$h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T) \propto \underbrace{\left[\prod_{c=1}^C f_c(\mathbf{X}_0^{(c)}) \right]}_{\text{initial core draws}} \underbrace{\exp\left(-\frac{C \cdot \|\mathbf{y}_T - \bar{\mathbf{X}}_0\|^2}{2T}\right) \cdot \exp\left(-\frac{C\sigma^2}{2T}\right)}_{\text{end point draw}}$$

- Need RS for $h_T^{\text{bm}}(\cdot, \cdot)$ end point.
- Accept with probability

$$P(\mathfrak{x}) := \exp\left[-\sum_{c=1}^C \int_0^T \kappa_c(\mathbf{X}_t^{(c)}) dt\right] \in [0, 1]$$

- Exact ('Talking') vs. Approximate ('Silent' / 'Lecture') Fusion

■ Remark: 'Ornstein-Uhlenbeck' special case

- Simple 'Brownian' Case:

- 'Optimal' $h_T^{\text{bm}}(\cdot, \cdot)$:

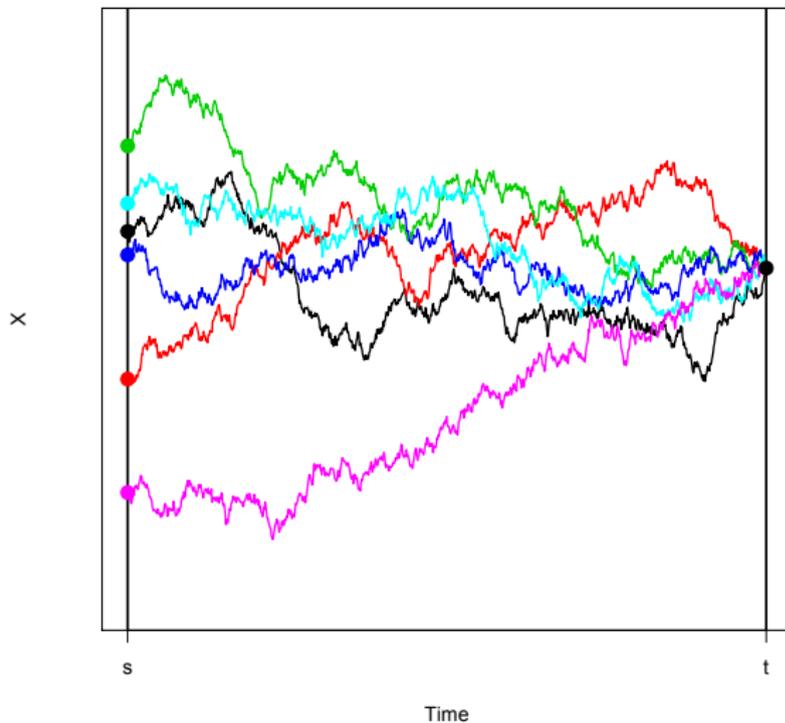
$$h_T^{\text{bm}}(\mathbf{X}_0^{(1:C)}, \mathbf{y}_T) \propto \underbrace{\left[\prod_{c=1}^C f_c(\mathbf{X}_0^{(c)}) \right]}_{\text{initial core draws}} \underbrace{\exp\left(-\frac{C \cdot \|\mathbf{y}_T - \bar{\mathbf{X}}_0\|^2}{2T}\right) \cdot \exp\left(-\frac{C\sigma^2}{2T}\right)}_{\text{end point draw}}$$

- Need RS for $h_T^{\text{bm}}(\cdot, \cdot)$ end point.
- Accept with probability

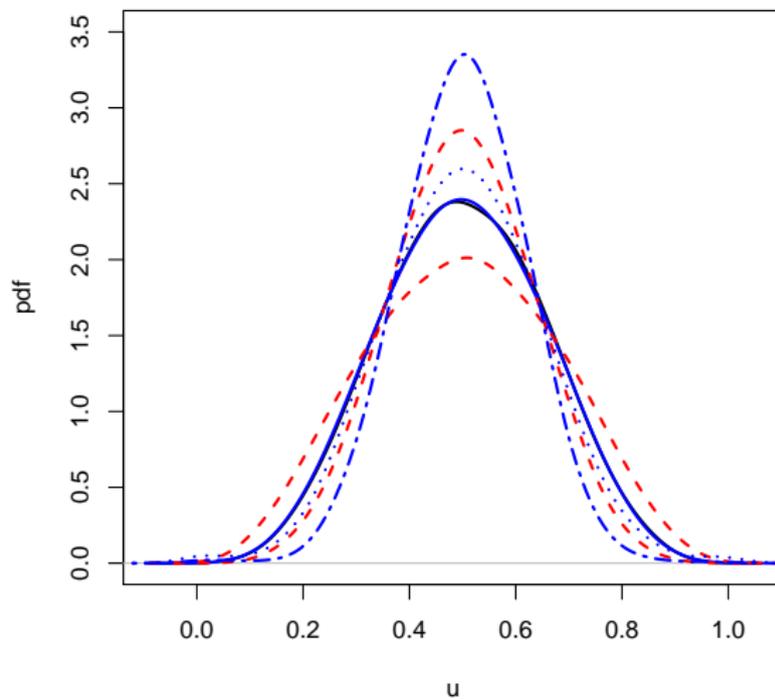
$$P(\mathfrak{x}) := \exp\left[-\sum_{c=1}^C \int_0^T \kappa_c(\mathbf{X}_t^{(c)}) dt\right] \in [0, 1]$$

- Exact ('Talking') vs. Approximate ('Silent' / 'Lecture') Fusion
- Remark: 'Ornstein-Uhlenbeck' special case

Some Details



Beta(5,5) density



Questions?