

# Divide-and-Conquer Fusion: Methods for unifying distributed analyses

Ryan Chan

Murray Pollock (Newcastle), Hongsheng Dai (Essex), Adam Johansen (Warwick), Gareth Roberts (Warwick)  
Poster session: Tuesday 28th June (19:00-22:00)

28 June 2022



**The  
Alan Turing  
Institute**

# Outline

Introduction to Fusion methodologies

- What is the Fusion problem?

- Monte Carlo Fusion

- Limitations of Monte Carlo Fusion

Divide-and-Conquer Generalised Monte Carlo Fusion

Divide-and-Conquer Generalised Bayesian Fusion

Examples

- Logistic regression

- Negative Binomial regression

Concluding remarks and future directions

# Fusion Problem

- Target fusion density:

$$f(\mathbf{x}) \propto \prod_{c=1}^C f_c(\mathbf{x})$$

where each *sub-posterior*,  $f_c(\mathbf{x})$ , is a density representing one of the  $C$  distributed inferences we wish to unify

- No general analytical approach
- Monte Carlo: assume we can sample  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$
- Applications:
  - Expert elicitation: combining views of multiple experts
  - Big Data (by construction)
    - Partitioning large datasets to make them more manageable
  - Inference in privacy settings

# Fusion Problem

- Target fusion density:

$$f(\mathbf{x}) \propto \prod_{c=1}^C f_c(\mathbf{x})$$

where each *sub-posterior*,  $f_c(\mathbf{x})$ , is a density representing one of the  $C$  distributed inferences we wish to unify

- No general analytical approach
  - Monte Carlo: assume we can sample  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$
- Applications:
  - Expert elicitation: combining views of multiple experts
  - Big Data (by construction)
    - Partitioning large datasets to make them more manageable
  - Inference in privacy settings

# Fusion Problem

- Target fusion density:

$$f(\mathbf{x}) \propto \prod_{c=1}^C f_c(\mathbf{x})$$

where each *sub-posterior*,  $f_c(\mathbf{x})$ , is a density representing one of the  $C$  distributed inferences we wish to unify

- No general analytical approach
- Monte Carlo: assume we can sample  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$
- Applications:
  - Expert elicitation: combining views of multiple experts
  - Big Data (by construction)
    - Partitioning large datasets to make them more manageable
  - Inference in privacy settings

# Fusion Problem

- Target fusion density:

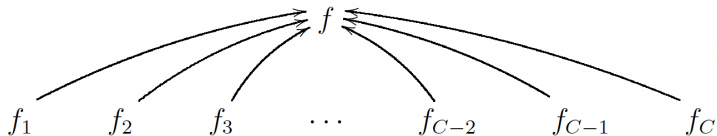
$$f(\mathbf{x}) \propto \prod_{c=1}^C f_c(\mathbf{x})$$

where each *sub-posterior*,  $f_c(\mathbf{x})$ , is a density representing one of the  $C$  distributed inferences we wish to unify

- No general analytical approach
- Monte Carlo: assume we can sample  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$
- Applications:
  - Expert elicitation: combining views of multiple experts
  - Big Data (by construction)
    - Partitioning large datasets to make them more manageable
  - Inference in privacy settings

# Fork-and-join

The **fork-and-join** approach:



## Some Fork-and-Join methods

- Several fork-and-join methods have been developed (typically for Bayesian inference for large dataset applications):
  - Kernel density averaging (KDEMC) [Neiswanger et al., 2014]
  - Weierstrass sampler (WRS) [Wang and Dunson, 2013]
  - Consensus Monte Carlo (CMC) [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
  - CMC is exact if all sub-posteriors are Gaussian
  - All theory is asymptotic in the number of observations
- However, Monte Carlo Fusion [Dai et al., 2019] (and subsequently Bayesian Fusion [Dai et al., 2021]) is **exact** in the sense it targets the correct fusion density



## Some Fork-and-Join methods

- Several fork-and-join methods have been developed (typically for Bayesian inference for large dataset applications):
  - Kernel density averaging (KDEMC) [Neiswanger et al., 2014]
  - Weierstrass sampler (WRS) [Wang and Dunson, 2013]
  - Consensus Monte Carlo (CMC) [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
  - CMC is exact if all sub-posteriors are Gaussian
  - All theory is asymptotic in the number of observations
- However, Monte Carlo Fusion [Dai et al., 2019] (and subsequently Bayesian Fusion [Dai et al., 2021]) is **exact** in the sense it targets the correct fusion density

## Some Fork-and-Join methods

- Several fork-and-join methods have been developed (typically for Bayesian inference for large dataset applications):
  - Kernel density averaging (KDEMC) [Neiswanger et al., 2014]
  - Weierstrass sampler (WRS) [Wang and Dunson, 2013]
  - Consensus Monte Carlo (CMC) [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
  - CMC is exact if all sub-posteriors are Gaussian
  - All theory is asymptotic in the number of observations
- However, Monte Carlo Fusion [Dai et al., 2019] (and subsequently Bayesian Fusion [Dai et al., 2021]) is **exact** in the sense it targets the correct fusion density

## Some Fork-and-Join methods

- Several fork-and-join methods have been developed (typically for Bayesian inference for large dataset applications):
  - Kernel density averaging (KDEMC) [Neiswanger et al., 2014]
  - Weierstrass sampler (WRS) [Wang and Dunson, 2013]
  - Consensus Monte Carlo (CMC) [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
  - CMC is exact if all sub-posteriors are Gaussian
  - All theory is asymptotic in the number of observations
- However, Monte Carlo Fusion [Dai et al., 2019] (and subsequently Bayesian Fusion [Dai et al., 2021]) is **exact** in the sense it targets the correct fusion density

## Some Fork-and-Join methods

- Several fork-and-join methods have been developed (typically for Bayesian inference for large dataset applications):
  - Kernel density averaging (KDEMC) [Neiswanger et al., 2014]
  - Weierstrass sampler (WRS) [Wang and Dunson, 2013]
  - Consensus Monte Carlo (CMC) [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
  - CMC is exact if all sub-posteriors are Gaussian
  - All theory is asymptotic in the number of observations
- However, Monte Carlo Fusion [Dai et al., 2019] (and subsequently Bayesian Fusion [Dai et al., 2021]) is **exact** in the sense it targets the correct fusion density

# An extended target density

## Proposition

Suppose that  $p_c(\mathbf{y}|\mathbf{x}^{(c)})$  is the transition density of a *stochastic process with stationary distribution*  $f_c^2(\mathbf{x})$ . The  $(C + 1)d$ -dimensional (fusion) density proportional to the integrable function

$$g\left(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}\right) \propto \prod_{c=1}^C \left[ f_c^2\left(\mathbf{x}^{(c)}\right) \cdot p_c\left(\mathbf{y}|\mathbf{x}^{(c)}\right) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

admits the *marginal density*  $f$  for  $\mathbf{y}$ .

**Main idea:** If we can sample from  $g$ , then we can obtain a draw from the fusion density ( $\mathbf{y} \sim f$ )

# An extended target density

## Proposition

Suppose that  $p_c(\mathbf{y}|\mathbf{x}^{(c)})$  is the transition density of a *stochastic process with stationary distribution*  $f_c^2(\mathbf{x})$ . The  $(C + 1)d$ -dimensional (fusion) density proportional to the integrable function

$$g\left(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}\right) \propto \prod_{c=1}^C \left[ f_c^2\left(\mathbf{x}^{(c)}\right) \cdot p_c\left(\mathbf{y}|\mathbf{x}^{(c)}\right) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

admits the *marginal density*  $f$  for  $\mathbf{y}$ .

**Main idea:** If we can sample from  $g$ , then we can obtain a draw from the fusion density ( $\mathbf{y} \sim f$ )

## An extended target density

- There are many possible choices for  $p_c(\mathbf{y}|\mathbf{x}^{(c)})$
- Let  $p_c(\mathbf{y}|\mathbf{x}^{(c)}) := p_{T,c}(\mathbf{y}|\mathbf{x}^{(c)})$ , the transition density of the  $d$ -dimensional (double) Langevin (DL) diffusion processes  $\mathbf{X}_t^{(c)}$  from  $\mathbf{x}^{(c)}$  to  $\mathbf{y}$  for  $c = 1, \dots, C$ , for a pre-defined time  $T > 0$  given by

$$d\mathbf{X}_t^{(c)} = \nabla \log f_c \left( \mathbf{X}_t^{(c)} \right) dt + d\mathbf{W}_t^{(c)},$$

(where  $\mathbf{W}_t^{(c)}$  is  $d$ -dimensional Brownian motion and  $\nabla$  is the gradient operator over  $\mathbf{x}$ )

- Has stationary distribution  $f_c^2(\mathbf{x})$
- Sample paths of DL diffusions can be simulated exactly using Path Space Rejection Sampling / Exact Algorithm methodology [Beskos et al., 2005, 2006; Pollock et al., 2016]

## An extended target density

- There are many possible choices for  $p_c(\mathbf{y}|\mathbf{x}^{(c)})$
- Let  $p_c(\mathbf{y}|\mathbf{x}^{(c)}) := p_{T,c}(\mathbf{y}|\mathbf{x}^{(c)})$ , the transition density of the  $d$ -dimensional (double) Langevin (DL) diffusion processes  $\mathbf{X}_t^{(c)}$  from  $\mathbf{x}^{(c)}$  to  $\mathbf{y}$  for  $c = 1, \dots, C$ , for a pre-defined time  $T > 0$  given by

$$d\mathbf{X}_t^{(c)} = \nabla \log f_c \left( \mathbf{X}_t^{(c)} \right) dt + d\mathbf{W}_t^{(c)},$$

(where  $\mathbf{W}_t^{(c)}$  is  $d$ -dimensional Brownian motion and  $\nabla$  is the gradient operator over  $\mathbf{x}$ )

- Has stationary distribution  $f_c^2(\mathbf{x})$
- Sample paths of DL diffusions can be simulated exactly using Path Space Rejection Sampling / Exact Algorithm methodology [Beskos et al., 2005, 2006; Pollock et al., 2016]



## Constructing a rejection sampler for $g$

- Extended target density:

$$g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C \left[ f_c^2(\mathbf{x}^{(c)}) \cdot p_{T,c}(\mathbf{y} | \mathbf{x}^{(c)}) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

- Consider the proposal density  $h$  for the extended target  $g$ :

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

- $\bar{\mathbf{x}} = \frac{1}{C} \sum_{c=1}^C \mathbf{x}^{(c)}$
- $T$  is an arbitrary positive constant

## Constructing a rejection sampler for $g$

- Extended target density:

$$g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C \left[ f_c^2(\mathbf{x}^{(c)}) \cdot p_{T,c}(\mathbf{y} | \mathbf{x}^{(c)}) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

- Consider the proposal density  $h$  for the extended target  $g$ :

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

- $\bar{\mathbf{x}} = \frac{1}{C} \sum_{c=1}^C \mathbf{x}^{(c)}$
- $T$  is an arbitrary positive constant

## Constructing a rejection sampler for $g$

- Simulation from  $h$  is easy:

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

1. Simulate  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$  independently
2. Simulate  $\mathbf{y} \sim \mathcal{N}_d(\bar{\mathbf{x}}, \frac{T}{C}\mathbb{I}_d)$ 
  - This value  $\mathbf{y}$  ends up being our proposal for  $f$

## Constructing a rejection sampler for $g$

- Simulation from  $h$  is easy:

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

1. Simulate  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$  independently
2. Simulate  $\mathbf{y} \sim \mathcal{N}_d(\bar{\mathbf{x}}, \frac{T}{C}\mathbb{I}_d)$ 
  - This value  $\mathbf{y}$  ends up being our proposal for  $f$

## Constructing a rejection sampler for $g$

- Simulation from  $h$  is easy:

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

1. Simulate  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$  independently
2. Simulate  $\mathbf{y} \sim \mathcal{N}_d(\bar{\mathbf{x}}, \frac{T}{C} \mathbb{I}_d)$ 
  - This value  $\mathbf{y}$  ends up being our proposal for  $f$

# Rejection sampling - acceptance probability

- Acceptance probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1$$

where

$$\left\{ \begin{array}{l} \rho_0 := e^{-\frac{C\sigma^2}{2T}}, \quad \sigma^2 = \frac{1}{C} \sum_{c=1}^C \|\mathbf{x}^{(c)} - \bar{\mathbf{x}}\|^2 \\ \rho_1 := \mathbb{E}_{\bar{\mathbb{W}}} \left( \prod_{c=1}^C \left[ \exp \left\{ - \int_0^T \left( \phi_c \left( \mathbf{X}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \right) \end{array} \right\}$$

where  $\bar{\mathbb{W}}$  denotes the law of  $C$  independent Brownian bridges  $\mathbf{X}_t^{(1)}, \dots, \mathbf{X}_t^{(C)}$  with  $\mathbf{X}_0 = \mathbf{x}^{(c)}$  and  $\mathbf{X}_T^{(c)} = \mathbf{y}$

- **Trade-off** with choice of  $T$ : as  $T$  increases,  $\rho_0$  increases, but this results in  $\rho_1$  to be small (might typically decrease exponentially with  $T$ )

## Rejection sampling - acceptance probability

- Acceptance probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1$$

where

$$\left\{ \begin{array}{l} \rho_0 := e^{-\frac{C\sigma^2}{2T}}, \quad \sigma^2 = \frac{1}{C} \sum_{c=1}^C \|\mathbf{x}^{(c)} - \bar{\mathbf{x}}\|^2 \\ \rho_1 := \mathbb{E}_{\bar{\mathbb{W}}} \left( \prod_{c=1}^C \left[ \exp \left\{ -\int_0^T \left( \phi_c \left( \mathbf{X}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \right) \end{array} \right.$$

where  $\bar{\mathbb{W}}$  denotes the law of  $C$  independent Brownian bridges  $\mathbf{X}_t^{(1)}, \dots, \mathbf{X}_t^{(C)}$  with  $\mathbf{X}_0 = \mathbf{x}^{(c)}$  and  $\mathbf{X}_T^{(c)} = \mathbf{y}$

- Trade-off with choice of  $T$ : as  $T$  increases,  $\rho_0$  increases, but this results in  $\rho_1$  to be small (might typically decrease exponentially with  $T$ )

## Rejection sampling - acceptance probability

- Acceptance probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1$$

where

$$\left\{ \begin{array}{l} \rho_0 := e^{-\frac{C\sigma^2}{2T}}, \quad \sigma^2 = \frac{1}{C} \sum_{c=1}^C \|\mathbf{x}^{(c)} - \bar{\mathbf{x}}\|^2 \\ \rho_1 := \mathbb{E}_{\bar{\mathbb{W}}} \left( \prod_{c=1}^C \left[ \exp \left\{ - \int_0^T \left( \phi_c \left( \mathbf{X}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \right) \end{array} \right.$$

where  $\bar{\mathbb{W}}$  denotes the law of  $C$  independent Brownian bridges  $\mathbf{X}_t^{(1)}, \dots, \mathbf{X}_t^{(C)}$  with  $\mathbf{X}_0 = \mathbf{x}^{(c)}$  and  $\mathbf{X}_T^{(c)} = \mathbf{y}$

- **Trade-off** with choice of  $T$ : as  $T$  increases,  $\rho_0$  increases, but this results in  $\rho_1$  to be small (might typically decrease exponentially with  $T$ )



## $\rho_1$ Acceptance Probability

$$\rho_1 := \mathbb{E}_{\bar{\mathbb{W}}} \left( \prod_{c=1}^C \left[ \exp \left\{ - \int_0^T \left( \phi_c \left( \mathbf{x}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \right)$$

where

- $\phi_c(\mathbf{x}) = \frac{1}{2} \left( \|\nabla \log f_c(\mathbf{x})\|^2 + \Delta \log f_c(\mathbf{x}) \right)$
- $\Phi_c$  are constants such that for all  $\mathbf{x}$ ,  $\phi_c(\mathbf{x}) \geq \Phi_c$  for  $c \in \{1, \dots, C\}$
- Events of probability  $\rho_1$  can be simulated using **Poisson thinning** and methodology called **Path-space Rejection Sampling (PSRS)** or the **Exact Algorithm** (Beskos et al. [2005], Beskos et al. [2006], Pollock et al. [2016])

# $\rho_1$ Acceptance Probability

$$\rho_1 := \mathbb{E}_{\bar{\mathbb{W}}} \left( \prod_{c=1}^C \left[ \exp \left\{ - \int_0^T \left( \phi_c \left( \mathbf{x}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \right)$$

where

- $\phi_c(\mathbf{x}) = \frac{1}{2} \left( \|\nabla \log f_c(\mathbf{x})\|^2 + \Delta \log f_c(\mathbf{x}) \right)$
- $\Phi_c$  are constants such that for all  $\mathbf{x}$ ,  $\phi_c(\mathbf{x}) \geq \Phi_c$  for  $c \in \{1, \dots, C\}$
- Events of probability  $\rho_1$  can be simulated using Poisson thinning and methodology called Path-space Rejection Sampling (PSRS) or the Exact Algorithm (Beskos et al. [2005], Beskos et al. [2006], Pollock et al. [2016])

# $\rho_1$ Acceptance Probability

$$\rho_1 := \mathbb{E}_{\bar{\mathbb{W}}} \left( \prod_{c=1}^C \left[ \exp \left\{ - \int_0^T \left( \phi_c \left( \mathbf{x}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \right)$$

where

- $\phi_c(\mathbf{x}) = \frac{1}{2} \left( \|\nabla \log f_c(\mathbf{x})\|^2 + \Delta \log f_c(\mathbf{x}) \right)$
- $\Phi_c$  are constants such that for all  $\mathbf{x}$ ,  $\phi_c(\mathbf{x}) \geq \Phi_c$  for  $c \in \{1, \dots, C\}$
- Events of probability  $\rho_1$  can be simulated using **Poisson thinning** and methodology called **Path-space Rejection Sampling (PSRS)** or the **Exact Algorithm** (Beskos et al. [2005], Beskos et al. [2006], Pollock et al. [2016])

# Monte Carlo Fusion - Summary

- Aim: Sample from  $g$  (admits marginal density  $f$  for  $\mathbf{y}$ )
- Proposal:

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

- Accept  $\mathbf{y}$  as a draw from fusion density  $f$  with probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1$$

# Monte Carlo Fusion - Summary

- Aim: Sample from  $g$  (admits marginal density  $f$  for  $\mathbf{y}$ )
- Proposal:

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

- Accept  $\mathbf{y}$  as a draw from fusion density  $f$  with probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1$$

# Monte Carlo Fusion - Summary

- Aim: Sample from  $g$  (admits marginal density  $f$  for  $\mathbf{y}$ )
- Proposal:

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left(-\frac{C \cdot \|\mathbf{y} - \bar{\mathbf{x}}\|^2}{2T}\right)$$

- Accept  $\mathbf{y}$  as a draw from fusion density  $f$  with probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1$$

# Limitations of Monte Carlo Fusion

- **Robustness:** there is a lack of robustness when:
  - sub-posterior correlation increases
  - $C$  increases
  - $d$  increases
  - combining **conflicting** sub-posteriors
- **Aim:** To construct a fusion algorithm / framework to alleviate some of these limitations (see Dai et al. [2021]; Chan et al. [2021] for full details)

# Limitations of Monte Carlo Fusion

- **Robustness:** there is a lack of robustness when:
  - sub-posterior correlation increases
  - $C$  increases
  - $d$  increases
  - combining **conflicting** sub-posteriors
- **Aim:** To construct a fusion algorithm / framework to alleviate some of these limitations (see Dai et al. [2021]; Chan et al. [2021] for full details)



# The Generalised Monte Carlo Fusion (GMCF) approach

## Problem: Scalability with sub-posterior correlation

- Recall we have the flexibility to choose different  $p_c$  (transition density of stochastic process with  $f_c^2$  invariant density)
- Now, we choose  $p_c$  to be the transition density of the  $d$ -dimensional (double) Langevin (DL) diffusion processes  $\mathbf{X}_t^{(c)}$  with covariance matrix,  $\Lambda_c$  from  $\mathbf{x}^{(c)}$  to  $\mathbf{y}$  for  $c = 1, \dots, C$ , over  $[0, T]$  given by

$$d\mathbf{X}_t^{(c)} = \Lambda_c \nabla \log f_c \left( \mathbf{X}_t^{(c)} \right) dt + \Lambda_c^{1/2} d\mathbf{W}_t^{(c)},$$

- Has stationary density proportional to  $f_c^2(x)$
- $\Lambda_c$  is the *preconditioning matrix* (enables incorporation of covariance / correlation structure into our algorithm)

# The Generalised Monte Carlo Fusion (GMCF) approach

## Problem: Scalability with sub-posterior correlation

- Recall we have the flexibility to choose different  $p_c$  (transition density of stochastic process with  $f_c^2$  invariant density)
- Now, we choose  $p_c$  to be the transition density of the  $d$ -dimensional (double) Langevin (DL) diffusion processes  $\mathbf{X}_t^{(c)}$  with covariance matrix,  $\mathbf{\Lambda}_c$  from  $\mathbf{x}^{(c)}$  to  $\mathbf{y}$  for  $c = 1, \dots, C$ , over  $[0, T]$  given by

$$d\mathbf{X}_t^{(c)} = \mathbf{\Lambda}_c \nabla \log f_c \left( \mathbf{X}_t^{(c)} \right) dt + \mathbf{\Lambda}_c^{1/2} d\mathbf{W}_t^{(c)},$$

- Has stationary density proportional to  $f_c^2(\mathbf{x})$
- $\mathbf{\Lambda}_c$  is the *preconditioning matrix* (enables incorporation of covariance / correlation structure into our algorithm)

## Constructing an importance sampler

- Switch to importance sampler for the extended target density  $g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})$ :
  - Rejection sampling can be wasteful
  - We will subsequently embed this approach within a SMC algorithm
- Consider an **alternative** proposal density  $h$  for the extended target  $g$ :

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{(\mathbf{y} - \tilde{\mathbf{x}})^\top \Lambda^{-1} (\mathbf{y} - \tilde{\mathbf{x}})}{2T} \right\},$$

where

$$\tilde{\mathbf{x}} := \left( \sum_{c=1}^C \Lambda_c^{-1} \right)^{-1} \left( \sum_{c=1}^C \Lambda_c^{-1} \mathbf{x}^{(c)} \right), \quad \Lambda^{-1} := \sum_{c=1}^C \Lambda_c^{-1}.$$

## Constructing an importance sampler

- Switch to importance sampler for the extended target density  $g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})$ :
  - Rejection sampling can be wasteful
  - We will subsequently embed this approach within a SMC algorithm
- Consider an **alternative** proposal density  $h$  for the extended target  $g$ :

$$h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{(\mathbf{y} - \tilde{\mathbf{x}})^\top \mathbf{\Lambda}^{-1} (\mathbf{y} - \tilde{\mathbf{x}})}{2T} \right\},$$

where

$$\tilde{\mathbf{x}} := \left( \sum_{c=1}^C \mathbf{\Lambda}_c^{-1} \right)^{-1} \left( \sum_{c=1}^C \mathbf{\Lambda}_c^{-1} \mathbf{x}^{(c)} \right), \quad \mathbf{\Lambda}^{-1} := \sum_{c=1}^C \mathbf{\Lambda}_c^{-1}.$$

# Importance weights

- Importance weights:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1$$

where

$$\begin{cases} \rho_0 := \exp \left\{ - \sum_{c=1}^C \frac{(\tilde{\mathbf{x}} - \mathbf{x}^{(c)})^\top \boldsymbol{\Lambda}_c^{-1} (\tilde{\mathbf{x}} - \mathbf{x}^{(c)})}{2T} \right\} \\ \rho_1 := \prod_{c=1}^C \mathbb{E}_{\mathbb{W}_{\boldsymbol{\Lambda}_c}} \left[ \exp \left\{ - \int_0^T \left( \phi_c \left( \mathbf{X}_t^{(c)} \right) - \boldsymbol{\Phi}_c \right) dt \right\} \right] \end{cases}$$

where  $\phi_c(\mathbf{x}) := \frac{1}{2} \left( \nabla \log f_c(\mathbf{x})^\top \boldsymbol{\Lambda}_c \nabla \log f_c(\mathbf{x}) + \text{Tr}(\boldsymbol{\Lambda}_c \nabla^2 \log f_c(\mathbf{x})) \right)$ , with  $\mathbb{W}_{\boldsymbol{\Lambda}_c}$  denoting the law of a Brownian bridge  $\{\mathbf{X}_t^{(c)}, t \in [0, T]\}$  with  $\mathbf{X}_0^{(c)} := \mathbf{x}^{(c)}$ ,  $\mathbf{X}_T^{(c)} := \mathbf{y}$  and covariance matrix  $\boldsymbol{\Lambda}_c$

# Scalability with sub-posterior correlation

In our *Generalised* Monte Carlo Fusion setting:

- Able to incorporate covariance / correlation information within our proposals and through  $p_c$  and  $h$  (in MCF  $\mathbf{\Lambda}_c = \mathbb{I}_d$  for  $c = 1, \dots, C$ )
- Unfortunately no longer have i.i.d. draws from  $f$  but now have weighted samples to approximate  $f$  (later embed within divide-and-conquer SMC [Lindsten et al., 2017] framework)

# Scalability with sub-posterior correlation

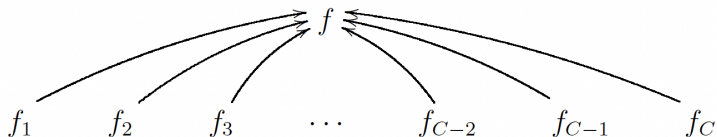
In our *Generalised* Monte Carlo Fusion setting:

- Able to incorporate covariance / correlation information within our proposals and through  $p_c$  and  $h$  (in MCF  $\mathbf{\Lambda}_c = \mathbb{I}_d$  for  $c = 1, \dots, C$ )
- Unfortunately no longer have i.i.d. draws from  $f$  but now have weighted samples to approximate  $f$  (later embed within divide-and-conquer SMC [Lindsten et al., 2017] framework)

# Divide-and-Conquer Monte Carlo Fusion

## Problem: Scalability with $C$

The (Generalised) Monte Carlo Fusion algorithm implies a **fork-and-join** approach:



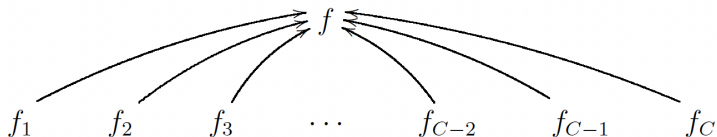
- Not necessarily the most efficient way to combine sub-posteriors
- For MCF, acceptance probabilities typically decrease geometrically with  $C$



# Divide-and-Conquer Monte Carlo Fusion

## Problem: Scalability with $C$

The (Generalised) Monte Carlo Fusion algorithm implies a **fork-and-join** approach:



- Not necessarily the most efficient way to combine sub-posteriors
- For MCF, acceptance probabilities typically decrease geometrically with  $C$

# Divide-and-Conquer Monte Carlo Fusion

## Problem: Scalability with $C$

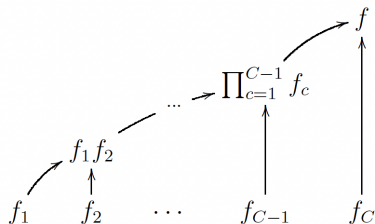
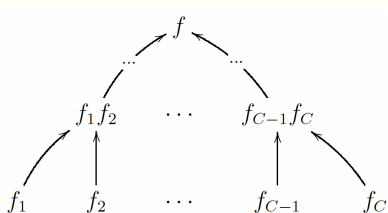
The (Generalised) Monte Carlo Fusion algorithm implies a **fork-and-join** approach:



- Not necessarily the most efficient way to combine sub-posteriors
- For MCF, acceptance probabilities typically decrease geometrically with  $C$

# Divide-and-Conquer Monte Carlo Fusion

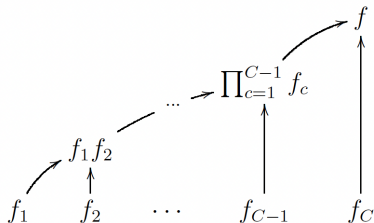
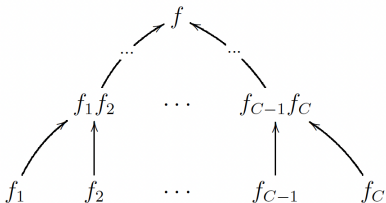
- **Solution:** Divide-and-Conquer Monte Carlo Fusion
  - We could perform fusion in a **proper divide-and-conquer** framework
    - i.e. a fork-and-join method is recursively applied
  - Two possible choices are **balanced-binary** (left) and **progressive** (right) trees



Note: Other trees are possible

# Divide-and-Conquer Monte Carlo Fusion

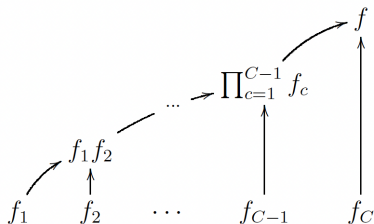
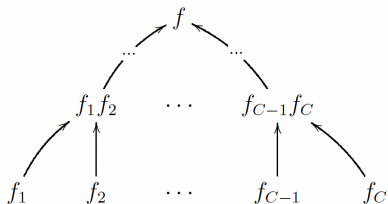
- **Solution:** Divide-and-Conquer Monte Carlo Fusion
  - We could perform fusion in a **proper divide-and-conquer** framework
    - i.e. a fork-and-join method is recursively applied
  - Two possible choices are **balanced-binary** (left) and **progressive** (right) trees



Note: Other trees are possible

# Divide-and-Conquer Monte Carlo Fusion

- **Solution:** Divide-and-Conquer Monte Carlo Fusion
  - We could perform fusion in a **proper divide-and-conquer** framework
    - i.e. a fork-and-join method is recursively applied
  - Two possible choices are **balanced-binary** (left) and **progressive** (right) trees

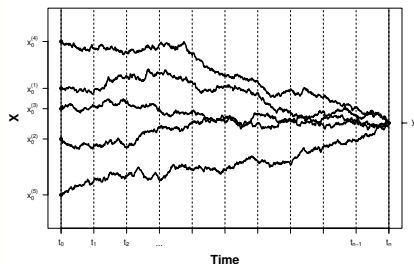


**Note:** Other trees are possible

# Generalised Bayesian Fusion

## Problem: Robustness to conflicting sub-posteriors

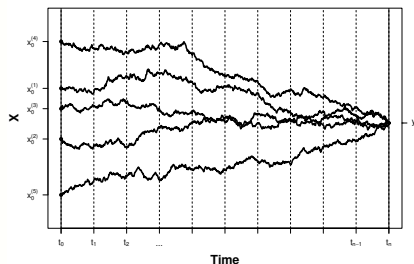
- Generalising the Bayesian Fusion approach of Dai et al. [2021]
- Recall choosing a value  $T > 0$  for MCF can be hard:
  - Want to make  $T$  large so that  $\rho_0$  is large - but this makes  $\rho_1$  smaller (since we have to simulate a diffusion over a longer time horizon  $T$ )
- **Solution:** Introduce temporal partition of  $T$ 
  - Have the flexibility to choose  $T$  large enough for initialisation, while being able to have small intervals in the partition



# Generalised Bayesian Fusion

## Problem: Robustness to conflicting sub-posteriors

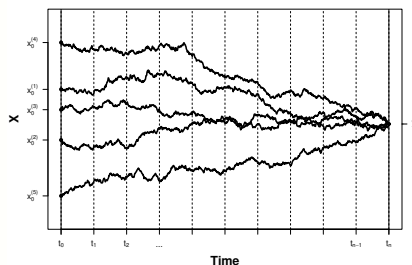
- Generalising the Bayesian Fusion approach of Dai et al. [2021]
- Recall choosing a value  $T > 0$  for MCF can be hard:
  - Want to make  $T$  large so that  $\rho_0$  is large - but this makes  $\rho_1$  smaller (since we have to simulate a diffusion over a longer time horizon  $T$ )
- **Solution:** Introduce temporal partition of  $T$ 
  - Have the flexibility to choose  $T$  large enough for initialisation, while being able to have small intervals in the partition



# Generalised Bayesian Fusion

## Problem: Robustness to conflicting sub-posteriors

- Generalising the Bayesian Fusion approach of Dai et al. [2021]
- Recall choosing a value  $T > 0$  for MCF can be hard:
  - Want to make  $T$  large so that  $\rho_0$  is large - but this makes  $\rho_1$  smaller (since we have to simulate a diffusion over a longer time horizon  $T$ )
- **Solution:** Introduce temporal partition of  $T$ 
  - Have the flexibility to choose  $T$  large enough for initialisation, while being able to have small intervals in the partition





## Examples

- We compare our methodology with the approximate methodologies KDEMC [Neiswanger et al., 2014], WRS [Wang and Dunson, 2013] and CMC [Scott et al., 2016]
- To compare methods we calculate the **integrated absolute distance** metric

$$IAD = \frac{1}{2d} \sum_{j=1}^d \int \left| \hat{f}(x_j) - f(x_j) \right| dx_j \in [0, 1]$$

where  $\hat{f}(x_j)$  is the marginal density for  $x_j$  based on the method applied (computed using a kernel density estimate) and  $f(x_j)$  is target marginal density

- Gives a measure of how accurate our samples are to our target (lower is better)

## Examples

- We compare our methodology with the approximate methodologies KDEMC [Neiswanger et al., 2014], WRS [Wang and Dunson, 2013] and CMC [Scott et al., 2016]
- To compare methods we calculate the **integrated absolute distance** metric

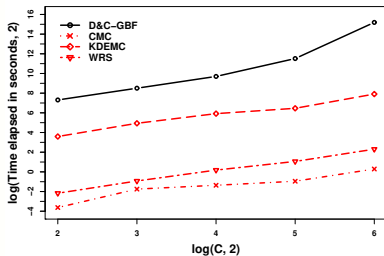
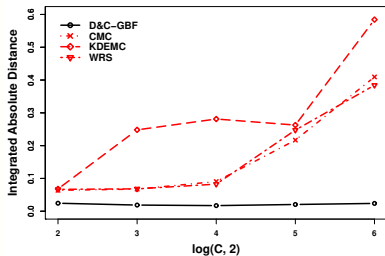
$$IAD = \frac{1}{2d} \sum_{j=1}^d \int \left| \hat{f}(\mathbf{x}_j) - f(\mathbf{x}_j) \right| d\mathbf{x}_j \in [0, 1]$$

where  $\hat{f}(\mathbf{x}_j)$  is the marginal density for  $\mathbf{x}_j$  based on the method applied (computed using a kernel density estimate) and  $f(\mathbf{x}_j)$  is target marginal density

- Gives a measure of how accurate our samples are to our target (lower is better)

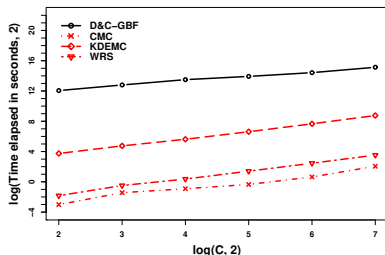
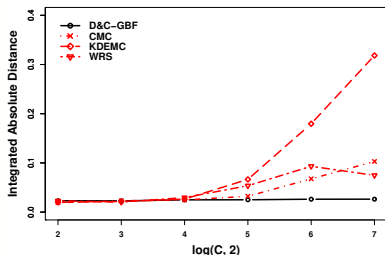
# Logistic regression

- Simulated data example with  $n = 1000$  and  $d = 5$  (and set  $N = 10000$ )
  - Small data size means that large data assumptions will fail
- We split the data into  $C = 4, 8, 16, 32, 64$  and apply D&C-GBF (using a balanced binary tree approach)



# Negative Binomial regression

- Using the **Bike sharing dataset** ( $n = 17379$ ,  $d = 10$ ) (and set  $N = 10000$ )
- We split the data into  $C = 4, 8, 16, 32, 64, 128$  and apply D&C-GBF (using a balanced binary tree approach)



## Ongoing research questions

- Reducing the computational cost of the Fusion approach
  - Exactness comes at a cost
- Practical implementation considerations for specific applications:
  - Big data setting: evaluating  $\phi_c$  has  $\mathcal{O}(m_c)$  cost - can perhaps employ **sub-sampling** methods to reduce this cost
  - **Confidential fusion** (**Con**-fusion): where sharing information/data between cores is **not** permitted
- Scalability with dimension
  - Performance with regards to dimension has improved since MCF, but not been explicitly addressed

## References

- Beskos, A., Papaspiliopoulos, O., Roberts, G. O., and Fearnhead, P. (2006). Exact and computationally efficient likelihood-based estimation for discretely observed diffusion processes (with discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3):333–382.
- Beskos, A., Roberts, G. O., et al. (2005). Exact simulation of diffusions. *The Annals of Applied Probability*, 15(4):2422–2444.
- Chan, R. S., Pollock, M., Johansen, A. M., and Roberts, G. O. (2021). Divide-and-Conquer Monte Carlo Fusion. Statistics e-print 2110.07265, arXiv.
- Dai, H., Pollock, M., and Roberts, G. O. (2019). Monte Carlo Fusion. *Journal of Applied Probability*, 56(1):174–191.
- Dai, H., Pollock, M., and Roberts, G. O. (2021). Bayesian Fusion: Scalable unification of distributed statistical analyses. Statistics e-print 2102.02123, arXiv.
- Lindsten, F., Johansen, A. M., Naesseth, C. A., Kirkpatrick, B., Schön, T. B., Aston, J. A., and Bouchard-Côté, A. (2017). Divide-and-Conquer with Sequential Monte Carlo. *Journal of Computational and Graphical Statistics*, 26(2):445–458.
- Neiswanger, W., Wang, C., and Xing, E. P. (2014). Asymptotically Exact, Embarrassingly Parallel MCMC. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, UAI'14*, page 623–632, Arlington, Virginia, USA. AUAI Press.
- Pollock, M., Johansen, A. M., Roberts, G. O., et al. (2016). On the exact and  $\epsilon$ -strong simulation of (jump) diffusions. *Bernoulli*, 22(2):794–856.
- Scott, S. L., Blocker, A. W., Bonassi, F. V., Chipman, H. A., George, E. I., and McCulloch, R. E. (2016). Bayes and Big Data: The Consensus Monte Carlo Algorithm. *International Journal of Management Science and Engineering Management*, 11(2):78–88.
- Wang, X. and Dunson, D. B. (2013). Parallelizing MCMC via Weierstrass Sampler. Statistics e-print 1312.4605, arXiv.

Poster session: Tuesday 28th June (19:00-22:00)