

# Divide-and-Conquer Generalised Bayesian Fusion

## A method of unifying distributed analyses

Ryan Chan<sup>1,3</sup>, Murray Pollock<sup>1,2</sup>, Adam Johansen<sup>1,3</sup>, Gareth Roberts<sup>1,3</sup>

<sup>1</sup> University of Warwick, <sup>2</sup> Newcastle University, <sup>3</sup> The Alan Turing Institute

### Generalised Monte Carlo Fusion (GMCF)

Provides theory to carry out *exact* inference for the target:

$$f(\mathbf{x}) \propto f_1(\mathbf{x}) \cdots f_C(\mathbf{x}) = \prod_{c=1}^C f_c(\mathbf{x}). \quad (1)$$

Uses *importance sampling* on the extended target distribution:

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

Let  $p_c(\mathbf{y} | \mathbf{x}^{(c)})$  be a transition density of a stochastic process with stationary distribution  $f_c^2(\mathbf{x})$ , then (2) admits  $f$  as a marginal for  $\mathbf{y}$ .

There are many possible choices for  $p_c$ , but we focus on letting  $p_c$  be transition density of a  $d$ -dimensional double Langevin diffusion process for  $t \in [0, T]$ :

$$d\mathbf{X}_t^{(c)} = \Lambda_c \nabla \log f_c(\mathbf{X}_t^{(c)}) dt + \Lambda_c^{\frac{1}{2}} d\mathbf{W}_t^{(c)}, \quad \mathbf{X}_0^{(c)} = \mathbf{x}^{(c)}, \quad \mathbf{X}_T^{(c)} = \mathbf{y}. \quad (3)$$

Proposal density:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{(\mathbf{y} - \tilde{\mathbf{x}})^\top \Lambda^{-1} (\mathbf{y} - \tilde{\mathbf{x}})}{2T} \right\}, \quad (4)$$

where

$$\tilde{\mathbf{x}} := \left( \sum_{c=1}^C \Lambda_c^{-1} \right)^{-1} \left( \sum_{c=1}^C \Lambda_c^{-1} \mathbf{x}^{(c)} \right), \quad \Lambda^{-1} := \sum_{c=1}^C \Lambda_c^{-1}. \quad (5)$$

Under certain mild conditions,

$$\frac{g(\mathbf{x}^{(1:C)}, \mathbf{y})}{h(\mathbf{x}^{(1:C)}, \mathbf{y})} \propto \rho_0 \cdot \rho_1 =: w'(\mathbf{x}^{(1:C)}, \mathbf{y}), \quad (6)$$

where

$$\rho_0(\mathbf{x}^{(1:C)}) := \exp \left\{ -\sum_{c=1}^C \frac{(\tilde{\mathbf{x}} - \mathbf{x}^{(c)})^\top \Lambda_c^{-1} (\tilde{\mathbf{x}} - \mathbf{x}^{(c)})}{2T} \right\}, \quad (7)$$

$$\rho_1(\mathbf{x}^{(1:C)}, \mathbf{y}) := \prod_{c=1}^C \mathbb{E}_{\mathbb{W}_{\Lambda_c}} \left[ \exp \left\{ -\int_0^T \phi_c(\mathbf{X}_t^{(c)}) dt \right\} \right], \quad (8)$$

and  $\phi_c(\mathbf{x}) := \frac{1}{2} (\nabla \log f_c(\mathbf{x})^\top \Lambda_c \nabla \log f_c(\mathbf{x}) + \text{Tr}(\Lambda_c \nabla^2 \log f_c(\mathbf{x})))$ , where  $\mathbb{W}_{\Lambda_c}$  denotes 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  $\Lambda_c$ .

**Algorithm 1** Generalised Monte Carlo Fusion (GMCF)

1. Simulate a proposal  $\mathbf{y}$  from  $h$  (4):
  - a) For  $c = 1, \dots, C$ , simulate  $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$  and calculate  $\tilde{\mathbf{x}}$ .
  - b) Simulate  $\mathbf{y} \sim \mathcal{N}_d(\tilde{\mathbf{x}}, T\Lambda)$ .
2. Assign unnormalised weight  $w'(\mathbf{x}^{(1:C)}, \mathbf{y})$  as per (6).

Interpretation: GMCF corrects a perturbation of a weighted average of Monte Carlo samples from the sub-posteriors.

### Generalised Bayesian Fusion (GBF)

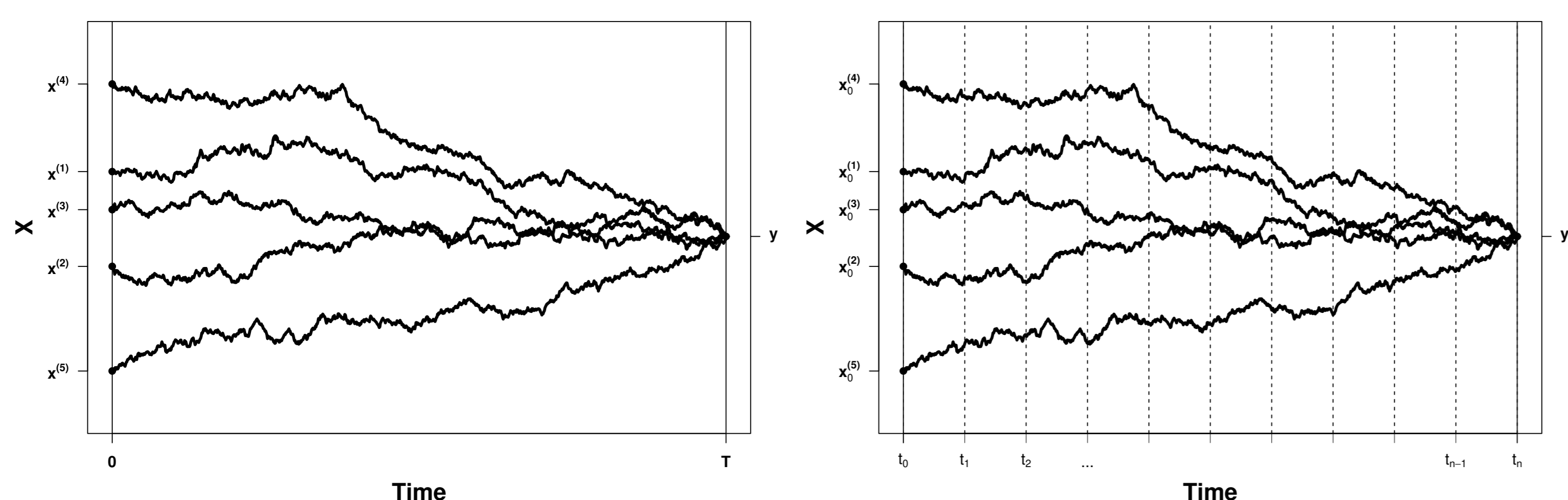


Figure 1: GMCF approach.

Figure 2: GBF approach.

- GMCF simulates  $C$  independent stochastic processes such that they coalesce at time  $T$  and the marginal distribution at time  $T$  has density  $f$  (1).
- GBF replaces IS with a sequential Monte Carlo approach which steps through a sequence of distributions between the initial proposal distribution to the target fusion density by considering a temporal partition of  $[0, T]$ :  
 $\mathcal{P} = \{t_0, t_1, \dots, t_n : 0 = t_0 < t_1 < \dots < t_n := T\}$ .

### Divide-and-Conquer Fusion

Problem: Fusion becomes inefficient as the number of sub-posteriors increases.

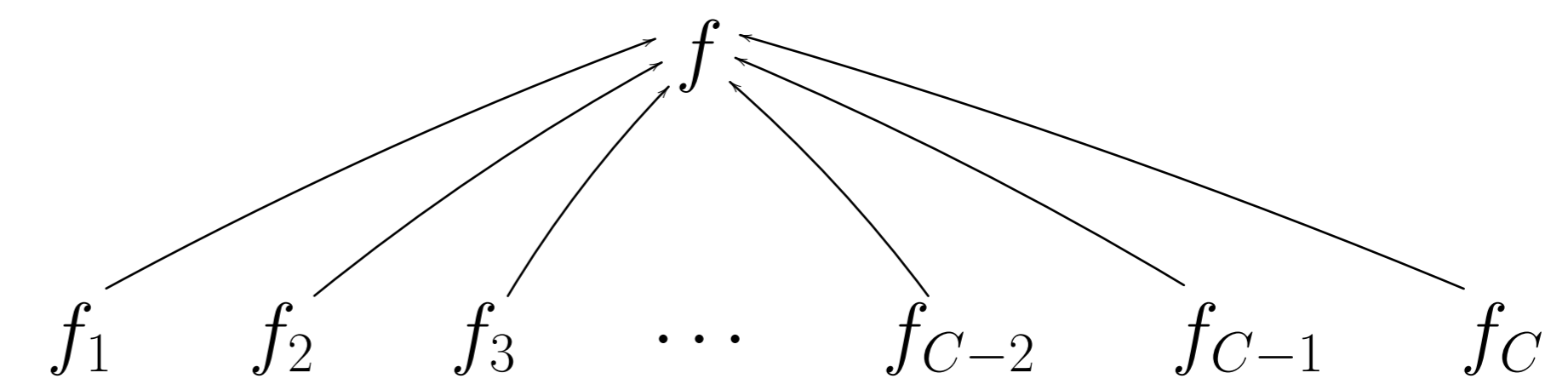


Figure 3: The 'fork-and-join' approach.

However, we can adopt a *divide-and-conquer* approach and embed our Fusion algorithm within a *divide-and-conquer sequential Monte Carlo* framework:

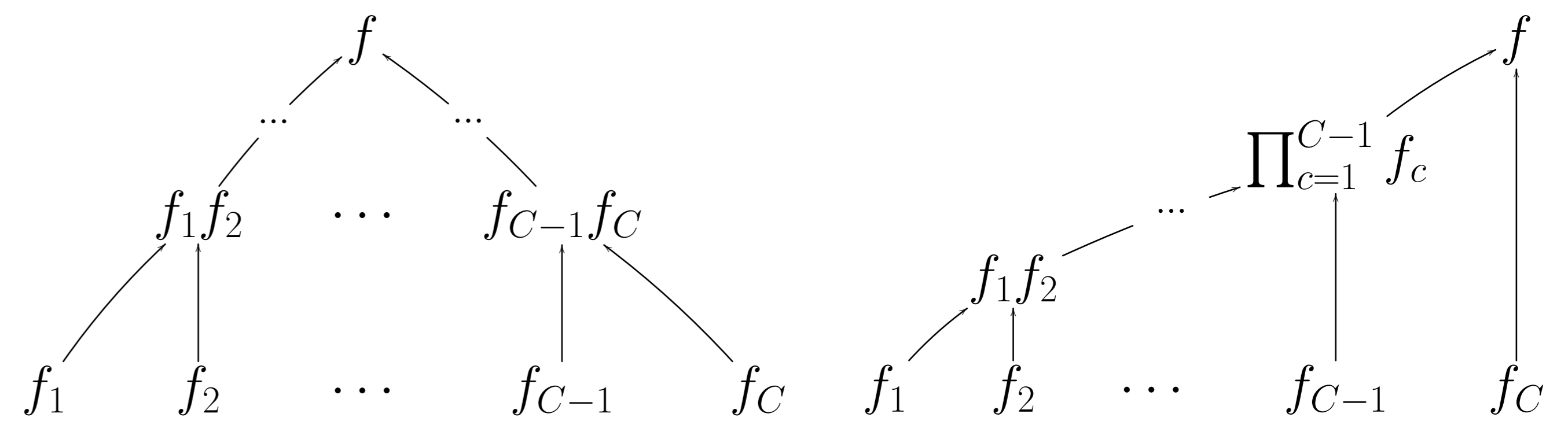


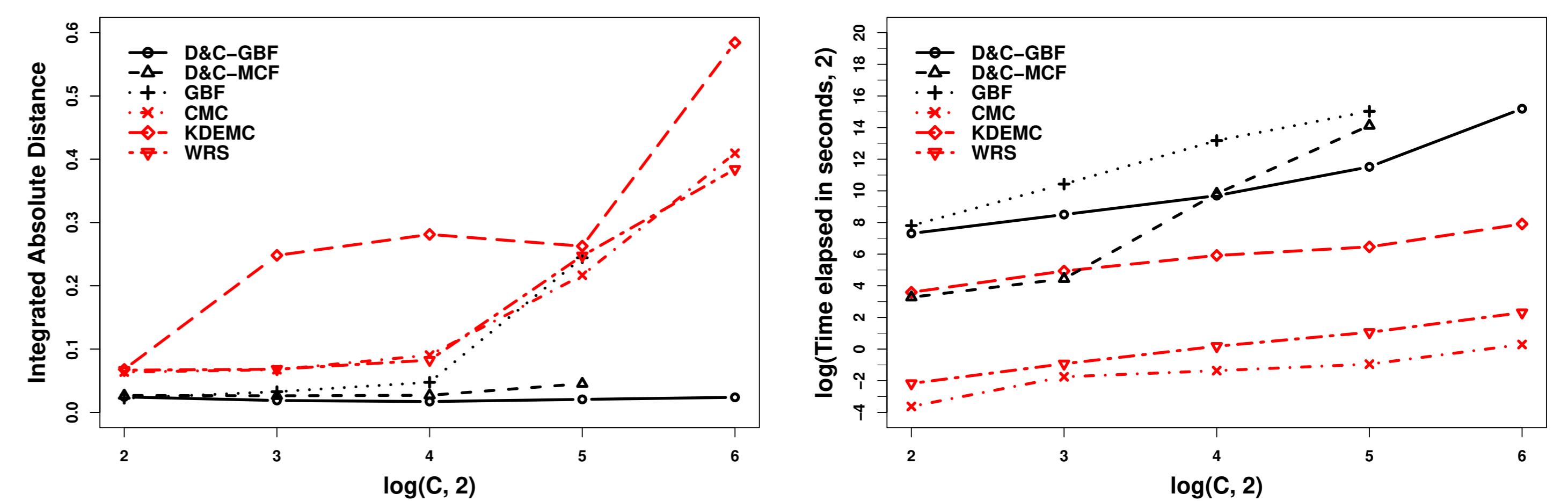
Figure 4: A balanced-binary tree.

Figure 5: A progressive tree.

### Logistic regression

- Simulated data with  $n = 1000$  and  $d = 5$
- To compare methods, use *IAD* (where  $f(\mathbf{x}_j)$  is target density) (lower better):

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



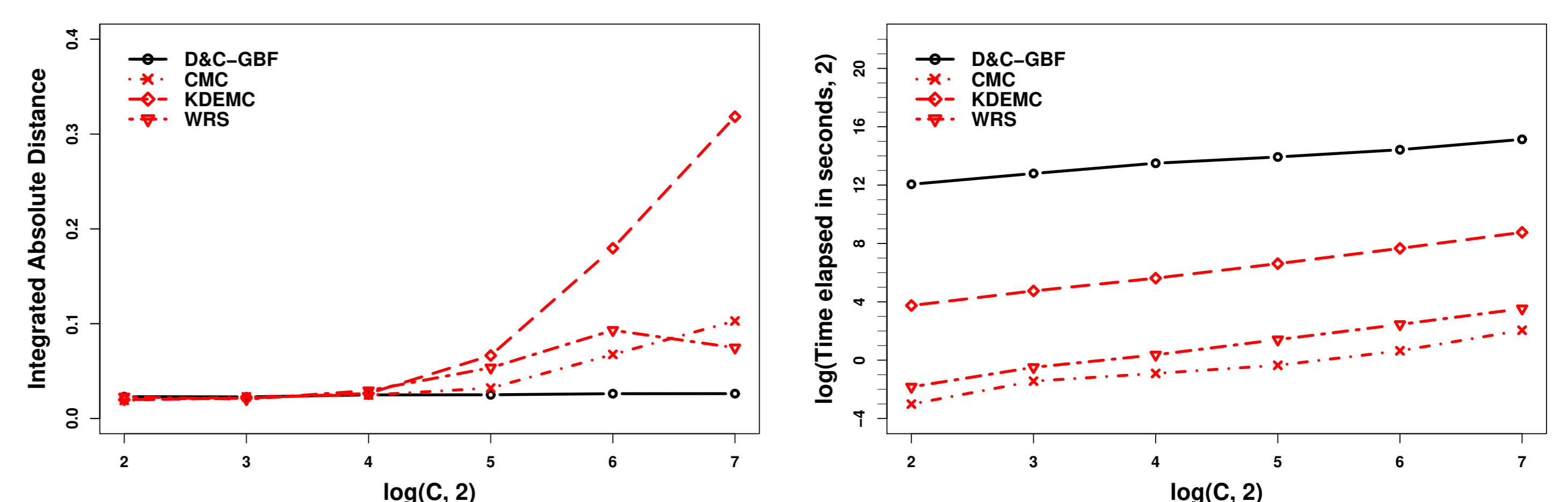
(a): Integrated Absolute Distance.

(b): Computational cost.

Figure 6: Comparison of competing methodologies to Fusion applied to logistic regression.

### Negative Binomial regression

- Bike sharing dataset ( $n = 17379$ ,  $d = 10$ )



(a): Integrated Absolute Distance.

(b): Computational cost.

Figure 7: Comparison of competing methodologies to Fusion applied to NB regression.

### References

- Chan, R.S.Y., Johansen, A.M., Pollock, M., and Roberts, G.O. 2021. *Divide-and-Conquer Monte Carlo Fusion*. Submitted. arXiv:2110.07265.
- Dai, Hongsheng, Murray Pollock, and Gareth Roberts. *Monte Carlo Fusion*. Journal of Applied Probability (56)1. 2019: pp. 174-191.
- Dai, H., Pollock, M. and Roberts, G., 2021. *Bayesian Fusion: Scalable unification of distributed statistical analyses*. Submitted. arXiv:2102.02123.