

Probabilistic reconstruction of cellular differentiation trees from single-cell RNA-seq data



Miriam Shiffman^{1,2,5}, Will Stephenson², Geoffrey Schiebinger^{3,5}, Trevor Campbell^{2,3}, Jonathan Huggins², Aviv Regev^{4,5}, Tamara Broderick^{2,3}



¹ Computational & Systems Biology, ² CSAIL, ³ IDSS, ⁴ Biology @



abstract

- interested in understanding **cellular differentiation** — how one cell type gives rise to others
- but** must infer this unobserved, inherently dynamic process from static, noisy snapshots (e.g. scRNA-seq)
- Bayesian nonparametric** approach provides flexible tree densities & principled uncertainties

Here, **we develop a full generative model and inference** (novel MCMC sampler) to:

- directly model **sources of uncertainty** in single-cell transcriptomic data
- infer **interpretable, probabilistic** insight into cell state & differentiation topology

biological motivation

Question: how does a less-specialized progenitor (stem cell) reliably give rise to many cell fates?

- ubiquitous to multicellular life (development + adulthood), yet not well-understood at systems level
- ↳ implications for **disease + basic science**

- insight into stable differentiation landscape

↔ insight into **master regulatory wiring** / molecular program

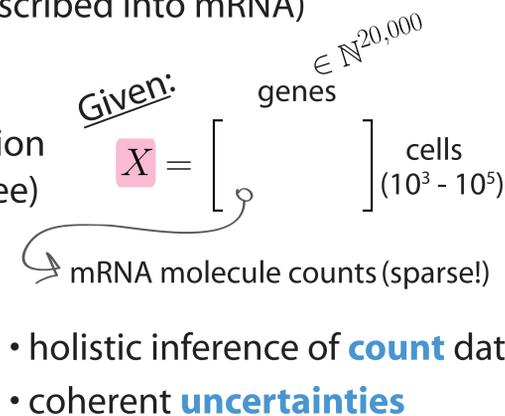
Challenge #1: reconstruct dynamic process from many static snapshots

Challenge #2: do not directly observe cell state, but rather a noisy proxy — e.g. transcriptome (genes recently transcribed into mRNA)

model desiderata

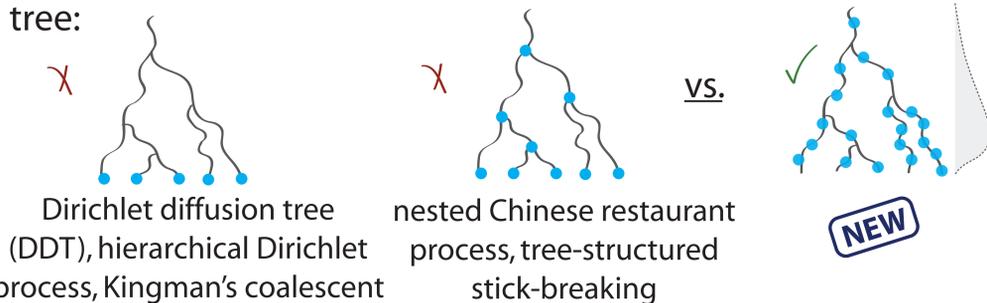
Goal: reconstruct continuous differentiation process (where latent cell states form a tree)

- learn tree **topology** of differentiation
- learn where **cells** fall along tree
- learn how **genes** drive branching, how expression changes along branches



Bayesian tree models

Probabilistic tree models allow for explicit **quantification of uncertainty** & flexible **binary branch topologies**, but existing models do not allow data to live arbitrarily along tree:



our augmented Dirichlet diffusion tree model

Recall: to perform inference, we begin by writing out the generative model, then use Bayes' rule to invert it.

- Generate tree** (node states + times) according to K -leaf DDT

- Gaussian diffusion b/w nodes
- branch times via hazard process
- regularize depth with prior over K

- Generate cells along tree:**

each cell traverses the tree, starting at root & choosing branches w.p.

\propto # of cells down each path.

Stop at time drawn from continuous distribution over $(0,1)$.

For each gene:

- draw latent cell state (λ) according to Brownian bridge defined by latent states of neighbors

- draw mRNA transcript count,

$$x \sim \text{Binom} \left(N_{\text{UMI}}, 1 - e^{-q h(\lambda)} \right)$$

max countable transcripts/gene

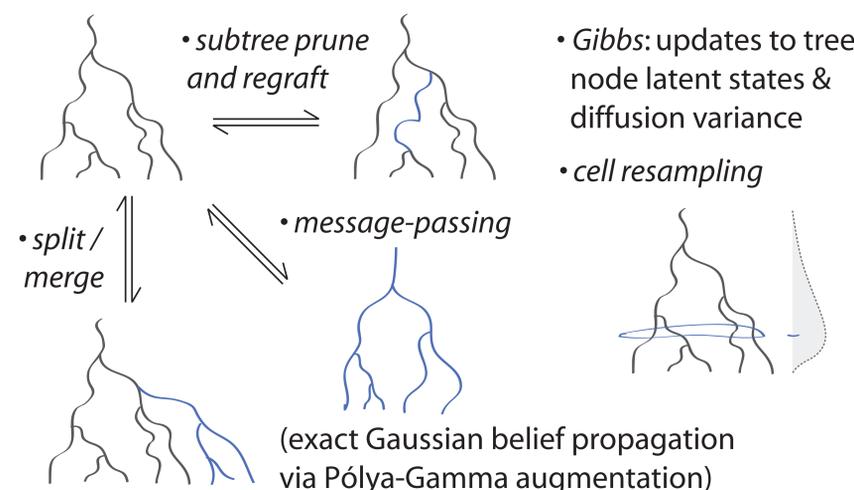
hyperparam related to dropout

positive link, $h: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$

from summing over nonzero Poissons

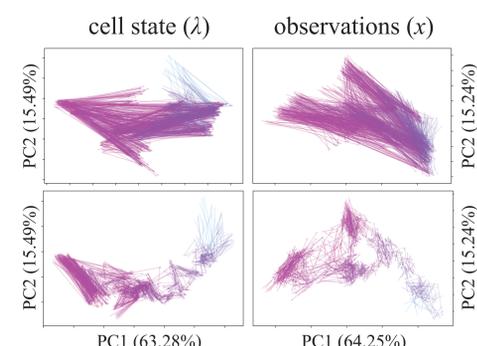
inference

As usual, exact inference intractable. Instead: approximate inference via novel **Markov chain Monte Carlo** sampler...



experiments

Partial **recovery of latent structure** from simulated data, including recovery of the true # of leaves.



convergence to ground truth demonstrated by shortening of arrow from true to inferred value per "cell," from initialized tree (top) to max-likelihood sampled (bottom)

future / ongoing

- beyond toy data
- tree metrics
- cell fate determinism vs. stochasticity
- test predictions by experimental perturbation