# CREATIVE PROBABILISTIC PROGRAMMING FOR BIOLOGY

Miriam Shiffman

shiffman@ { mit.edu
            broadinstitute.org

"MEANINGFULNESS" of a learned representation in biology can only be measured
w.r.t. a particular biological CONTEXT or question.

MODELING is the structure that provides this CONTEXT and endows latent representations
with meaning.

*PROBABILISTIC* MODELING is often the best choice
→ interpretability / decision theory
→ coherent framework for hierarchies, noise
→ biology itself is PROBABILISTIC !

PROBABILISTIC PROGRAMMING LANGUAGES are one tool missing from widespread adoption in biology,
w/ potential to more naturally + holistically meld the modeling process w/ the process of wet lab science.

**how can experimental biology be restructured around probabilistic modeling ?**
(as an ongoing part of data collection & experimental design, beyond *post hoc* analysis)

**how can PPLs be extended to meet the particular challenges of biology ?**
(and promote model-tinkering in new & creative ways)

## PROBABILISTIC PROGRAMMING LANGUAGES (PPLs)...

add `random variables` to the list of
types we expect in a language: `str`, `int`, ...
→ fundamental operations in probability = fundamental
(automated) features: `sample`, `condition`, `infer`

minimize edit distance b/w  | writing down the mathematical / coding up the executable | model

promote EXPERIMENTATION & CREATIVITY in generative modeling
→ complex architectures out of legolike abstractions
→ tweak the model but not the algorithm
→ concise, intuitive, human-readable

...just as differentiable languages have done for neural networks.

## THE FUTURE...?

how to:
→ visualize uncertainty for high-D, multimodal posteriors ?
→ integrate PPLs more intimately into wet lab, like optimizing experimental protocols ?
→ extend support for discrete structures **LIKE TREES**, a common regime in biology ?
*(when posterior not diff'able w.r.t. its params, precluding VI & HMC)*

could:
→ uncertainty quantification inform next gene to perturb or tissue to sequence ?
→ probabilistic programs of biological processes be synthesized from data ?
→ useful biological structures be encoded as PPL primitives ?
*(Gene Ontology, KEGG pathways, genome coordinates, ...)*

MIT CSAIL
BROAD INSTITUTE

## PROBABILISTIC DIFFERENTIATION TREES, w/o PPLs...

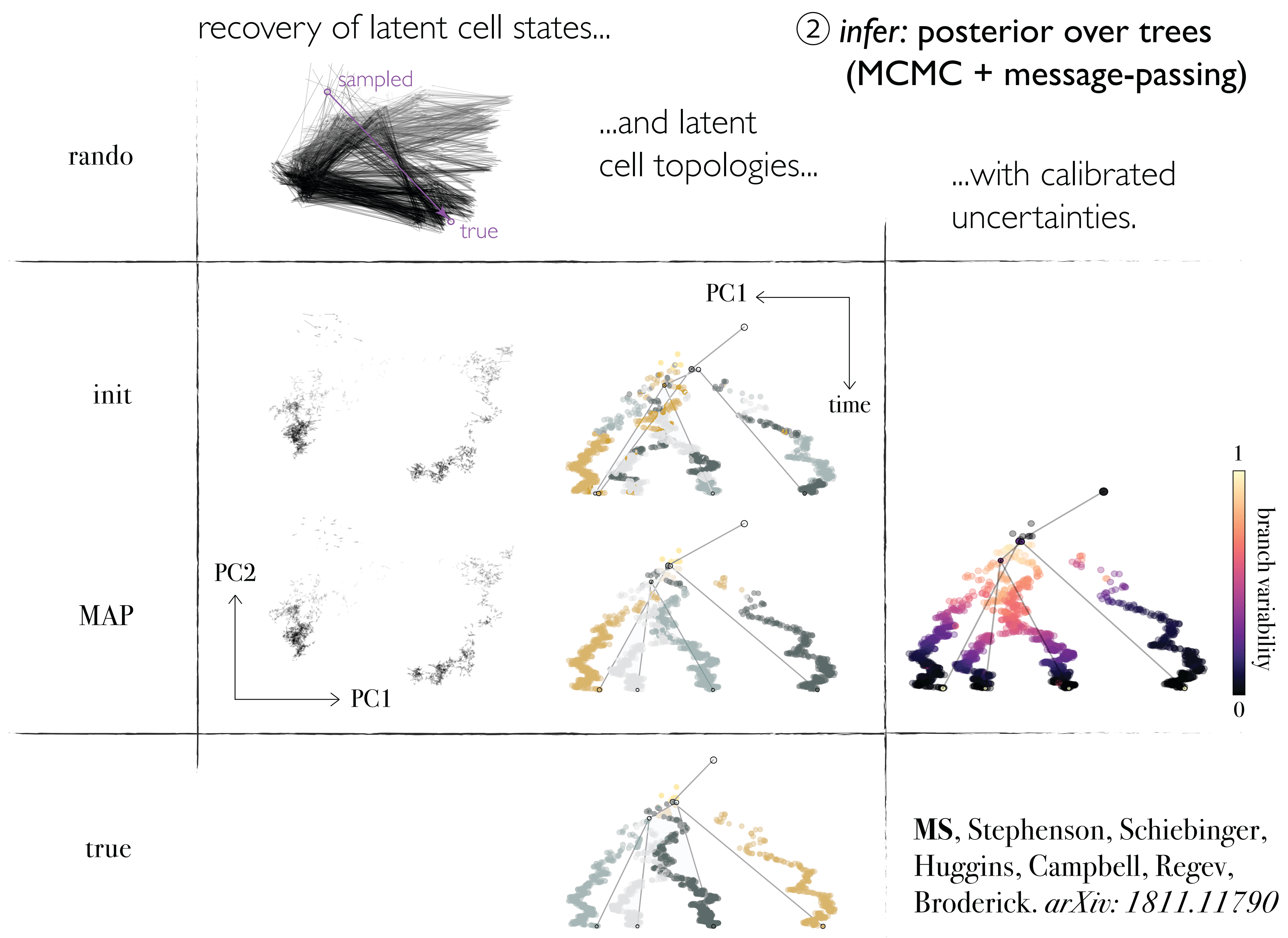how do stem cells give rise to specialized divergent reproducible cell fates ?

*[Waddington '57]*

genes
cells
# transcripts

potency
"cell state"
"cell state"

many static snapshots
(scRNA-seq)
→ branching dynamics
(continuous, probabilistic tree)

$\#leaves \sim \text{Poisson} + 1$
$tree \mid \#leaves \sim \text{DirichletDiffusionTree}$

cell positions

$times \sim \text{Beta}$
$branches \mid tree, times \sim \text{richGetRicher}$
$states \mid branches, tree, times \sim \text{GaussianDiffusion}$

$expressionProfiles \mid states \sim \text{Binom}(N_{UMI}, \sigma(states))$

time
state

① *simulate*: 2000 cells, 10 genes; fixed topology / times ("time course")

recovery of latent cell states...
② *infer*: posterior over trees
(MCMC + message-passing)

...and latent cell topologies...
...with calibrated uncertainties.

rando
sampled
true

init
PC1
time

PC2
PC1
MAP

PC1

true

1
branch variability
0

**MS**, Stephenson, Schiebinger, Huggins, Campbell, Regev, Broderick. *arXiv: 1811.11790*