# Inferring tumor progression in large datasets

M. Neyshabouri, S. Jun, J. Lagergren. *PLOS Computational Biology.*

Presented by Gillian Chu

# Tumor Progression

What is tumor progression?

- Find a set of linearly ordered mutually exclusive driver pathways

# Tumor Progression

What is tumor progression?

- Find a set of linearly ordered mutually exclusive driver pathways

Why is it complex?

- Driver mutations are hard to identify
- Driver gene mutations are mediated by pathways
- Set of genes provide the same effect
- Selective advantage of one gene exhausts others in a mutually exclusive manner

# Overview

- Presents a probabilistic model of mutually exclusive linearly ordered pathways
- Sampling based inference algorithm to train the model
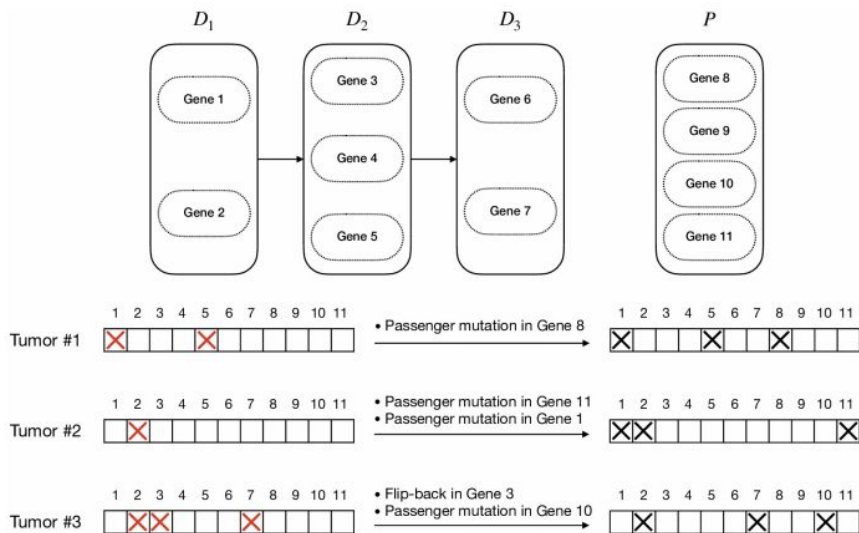- Test it against ILP solutions

# Overview

1. Linear Progression Model
2. Probabilistic Generative Process
3. Model Selection
4. Experiments on Simulated Data
5. Real Data Validation

# Overview

# Linear Progression Model



- Passenger mutations
- Flip-back events

# Notation

- N genes
- $\mathcal{P} = (D_1, D_2, \ldots, D_L, P)$ as an ordered partition of N. We have a linear pathway progression model of length L. D1 - DL are driver pathways, P is passengers.
- $Y \in \{0, 1\}^{M \times N}$ Mutation matrix Y for M tumors
- Noise parameters: δ, $\epsilon$

# Overview

1. Linear Progression Model
2. Probabilistic Generative Process
3. Model Selection
4. Experiments on Simulated Data
5. Real Data Validation
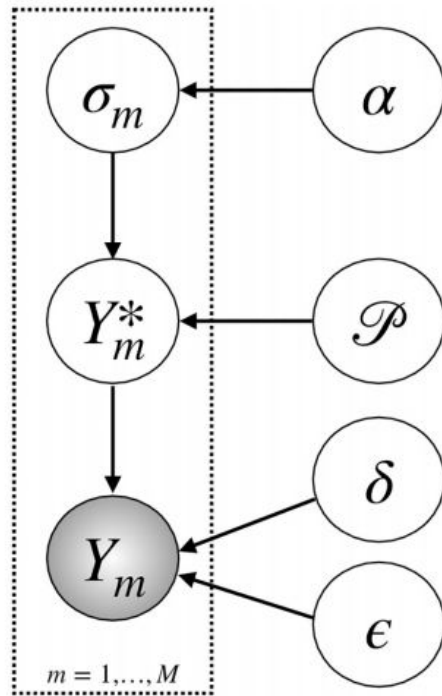
# Probabilistic Generative Model



**Algorithm 1** Generative process

1: **Input:** $\alpha, \mathcal{P}, \epsilon, \delta$
2: **Output:** $Y$
3: **for** $m \in \{1, \ldots, M\}$ **do**
4:      Let $Y_m$ be a vector of $N$ zeros
5:      Draw $\sigma_m \sim \text{Categorical}(\alpha)$
6:      **for** $k \in \{1, \ldots, \sigma_m\}$ **do**
7:          Draw $g_k$ (uniformly) from $D_k$
8:          Set $Y_m[g_k] = 1$
9:      **for** $n = 1$ to $N$ **do**
10:         **if** $Y_m[n] == 0$ **then**
11:             Set $Y_m[n] = 1$ with prob. $\epsilon$
12:         **else**
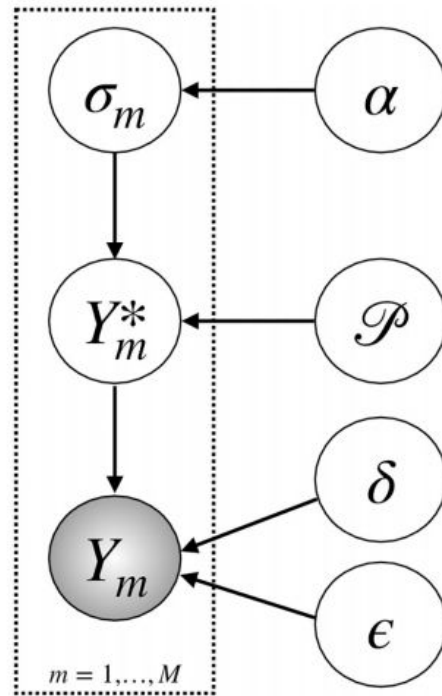13:             Set $Y_m[n] = 0$ with prob. $\delta$

# Probabilistic Generative Model



**Algorithm 1** Generative process

1: **Input:** $\alpha, \mathcal{P}, \epsilon, \delta$
2: **Output:** $Y$
3: **for** $m \in \{1, \ldots, M\}$ **do**
4:     Let $Y_m$ be a vector of $N$ zeros
5:     Draw $\sigma_m \sim \text{Categorical}(\alpha)$
6:     **for** $k \in \{1, \ldots, \sigma_m\}$ **do**
7:         Draw $g_k$ (uniformly) from $D_k$
8:         Set $Y_m[g_k] = 1$
9:     **for** $n = 1$ to $N$ **do**
10:        **if** $Y_m[n] == 0$ **then**
11:            Set $Y_m[n] = 1$ with prob. $\epsilon$
12:        **else**
13:            Set $Y_m[n] = 0$ with prob. $\delta$

# Overview

# Model Likelihood

- Determine the number of driver pathways L by computing marginal likelihood given L
- $Y = \{Y_m : m \in \{1, \ldots, M\}\}$ Mutation matrix for a collection of tumors
- $\mathcal{P} = (D_1, D_2, \ldots, D_L, P)$ Pathway progression model
- Find the $p(Y|\mathcal{P}, \alpha, \epsilon, \delta)$

# Model Selection w/ MCMC

- Initialize $\mathcal{P}^0, \epsilon^0, \delta^0$.

- For $t = 1, \dots, T$:

  - Sample $\mathcal{P}^t \sim \pi(\mathcal{P}|L, \epsilon^{t-1}, \delta^{t-1}, Y, \alpha)$,

  - Sample $\epsilon^t \sim \pi(\epsilon|\mathcal{P}^t, \delta^{t-1}, Y, \alpha)$,

  - Sample $\delta^t \sim \pi(\delta|\mathcal{P}^t, \epsilon^t, Y, \alpha)$.

$$\pi(\mathcal{P}|\epsilon, \delta, Y, L, \alpha) = \frac{p(Y|\epsilon, \delta, \mathcal{P}, \alpha)p(\mathcal{P}|L)}{p(Y|\epsilon, \delta, L, \alpha)} \propto p(Y|\epsilon, \delta, \mathcal{P}, \alpha)$$
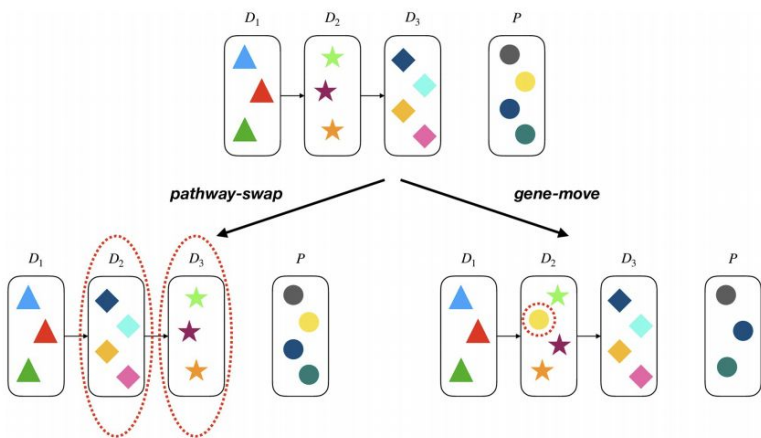
MCMC algorithm to generate samples from the posterior

- Given Y observations for fixed model length L
- Sample the progression model and error parameters for T iter

# Progression Model Proposal



Proposal function randomly chosen via Bernoulli:

- Pathway-swap
- gene-move

# Overview

1. Linear Progression Model
2. Probabilistic Generative Process
3. Model Selection
4. Experiments on Simulated Data
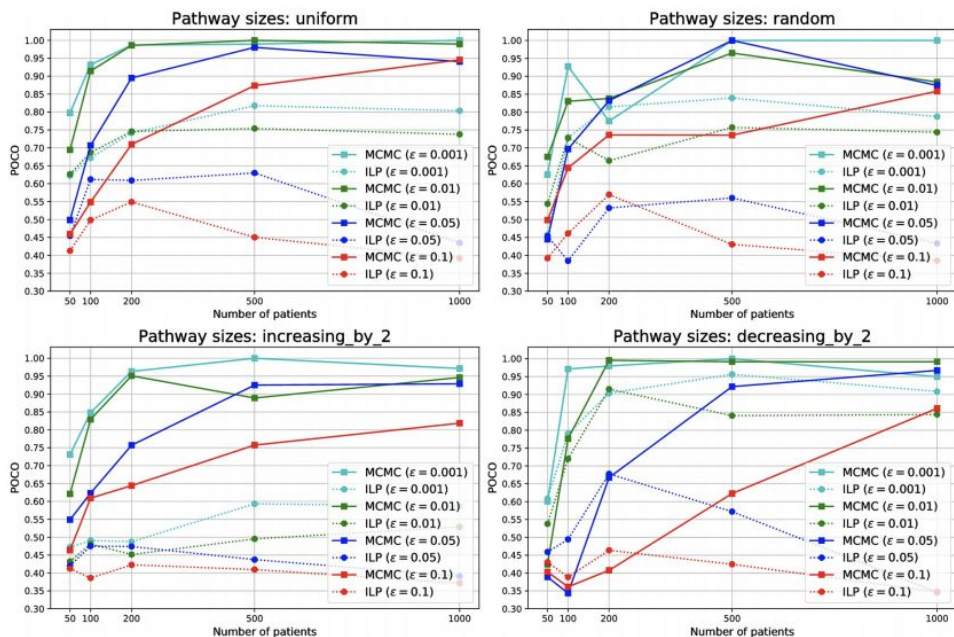5. Real Data Validation

# Metrics

- Model Evidence
- POCO (Percentage of Correct Ordering of Genes)
  - Given an inferred model, for each pair of genes check relative position
  - POCO is the percentage of gene pairs with correct relative position
- F1 Score

$$F1 = 2 \times \frac{Precision * Recall}{Precision + Recall}$$

# Exp 1: Known Driver Genes

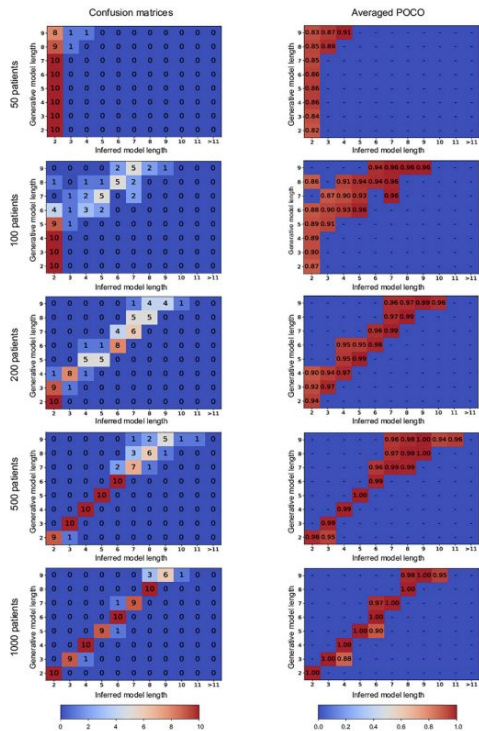- Distributed 25 genes in 5 driver pathways:

# Exp 2: Unknown Driver Genes

- 10 datasets with 500 patients
- ILP requires background mutation rate as input

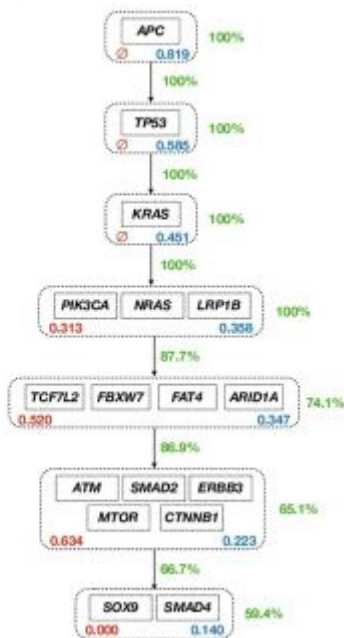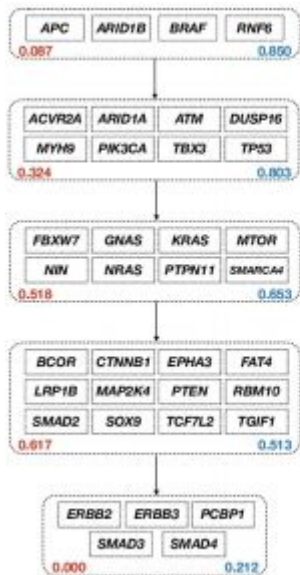| Method | POCO | | | F1 score (driver detection) | | | F1 score (pathway detection) | | |
|---|---|---|---|---|---|---|---|---|---|
| | 5 passengers | 25 passengers | 100 passengers | 5 passengers | 25 passengers | 100 passengers | 5 passengers | 25 passengers | 100 passengers |
| MCMC | 0.942 | 0.936 | 0.951 | 0.974 | 0.960 | 0.935 | 0.916 | 0.877 | 0.868 |
| ILP($\epsilon = 0.01$) | 0.568 | 0.423 | 0.285 | 0.909 | 0.667 | 0.343 | 0.392 | 0.238 | 0.086 |
| ILP($\epsilon = 0.05$) | 0.772 | 0.808 | 0.856 | 0.882 | 0.864 | 0.779 | 0.389 | 0.379 | 0.304 |
| ILP($\epsilon = 0.1$) | 0.551 | 0.494 | 0.713 | 0.493 | 0.428 | 0.400 | 0.380 | 0.333 | 0.276 |

# Exp 3: Model Length Selection



- 10 datasets from length [2, 9] with 50-1000 patients
- POCO is better than length at small patient cohort sizes

# Overview

1. Linear Progression Model
2. Probabilistic Generative Process
3. Model Selection
4. Experiments on Simulated Data
5. Real Data Validation

# Biological Data Analysis



ME scores of the pathways are the red numbers

Orthogonal validation of the driver genes based on domain knowledge

ILP places low mutation rate false positives, vs. MCMC focus on highly mutated genes

$$S_{ME}(g_i, g_j) = S_{ME}(g_j, g_i) = \frac{\text{mutation rate of } g_i \text{ in patients with mutated } g_j}{\text{mutation rate of } g_i \text{ in all patients}}$$

# Conclusion

- Probabilistic generative process for tumor progression
- MCMC process for model selection
- Comparison:
    - ILP places low mutation rate genes, more false positives
    - MCMC places highly mutated genes, more passenger genes