

CIMICE: MARKOV CHAIN INFERENCE METHOD TO IDENTIFY CANCER EVOLUTION

Nicolò Rossi

Nicola Gigante

Carla Piazza

Nicola Vitacolonna

Dip. di Scienze Matematiche, Informatiche e Fisiche - UniUD



PREMISES

THE CONTEXT

- ▶ Investigation on the **mutational history** of a **cancer cell**
- ▶ Relying on **single cell data** at a **single time instant**
- ▶ For the reconstruction a **probabilistic model**

THE AIM

- ▶ Identification of a **minimal set of assumptions** on **models/data**
- ▶ **Detection** of the sources of **uncertainty** in the reconstruction
- ▶ Provision of **suggestions** for further **experiments**

PREMISES

THE CONTEXT

- ▶ Investigation on the **mutational history** of a **cancer cell**
- ▶ Relying on **single cell data** at a **single time instant**
- ▶ For the reconstruction a **probabilistic** model

THE AIM

- ▶ Identification of a **minimal set of assumptions** on **models/data**
- ▶ **Detection** of the sources of **uncertainty** in the reconstruction
- ▶ Provision of **suggestions** for further **experiments**

OUR RESULTS

MODEL RECONSTRUCTION

We find a **minimal set of assumptions** such that:

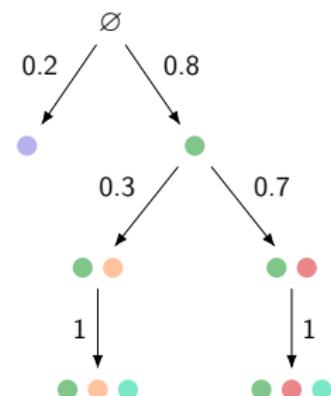
- ▶ without convergent evolutionary paths
 - there is **one** probabilistic underlying model
 - CIMICE infers it in **efficiently** w.r.t. the data size
- ▶ with convergent evolutionary paths
 - there is an **infinite set** of possible models
 - CIMICE heuristically assigns weights on convergences to pick a **preferred one**

OUR RESULTS

MODEL RECONSTRUCTION

We find a **minimal set of assumptions** such that:

- ▶ without convergent evolutionary paths
 - there is **one** probabilistic underlying model
 - CIMICE infers it in **efficiently** w.r.t. the data size
- ▶ with convergent evolutionary paths
 - there is an **infinite set** of possible models
 - CIMICE heuristically assigns weights on convergences to pick a **preferred** one

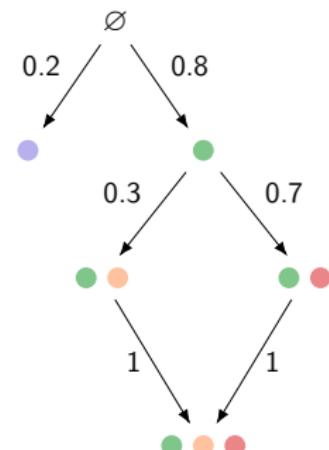


OUR RESULTS

MODEL RECONSTRUCTION

We find a **minimal set of assumptions** such that:

- ▶ without convergent evolutionary paths
 - there is **one** probabilistic underlying model
 - CIMICE infers it efficiently w.r.t. the data size
- ▶ with convergent evolutionary paths
 - there is an **infinite set** of possible models
 - CIMICE heuristically assigns weights on convergences to pick a **preferred one**

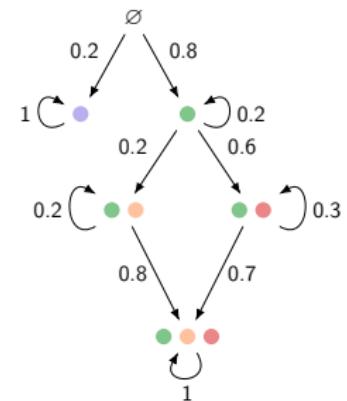


OUR RESULTS

GENERATIVE MODELS

Whenever our **biological assumptions are reasonable**, CIMICE produces **synthetic data** to

- ▶ generate **more data** from an inferred model
- ▶ test different **model inference methods**

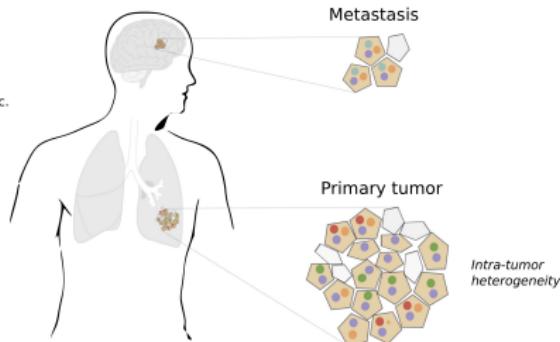
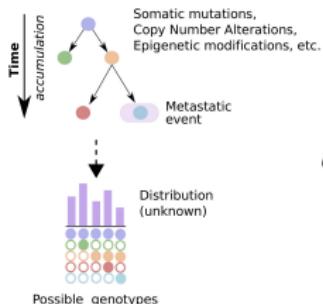


PLAN OF THE TALK

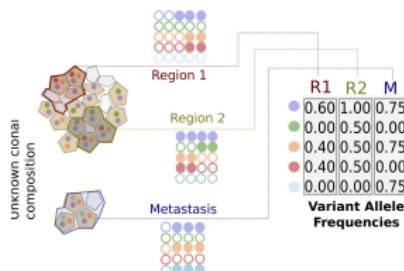
1. Single Cell Data
2. From Biological Assumptions to Models
3. Two Reconstruction Problems
4. Our Inference Algorithm
5. CIMICE Tool
6. Synthetic Models and Tests
7. Real Data
8. Conclusion

SINGLE CELL DATA

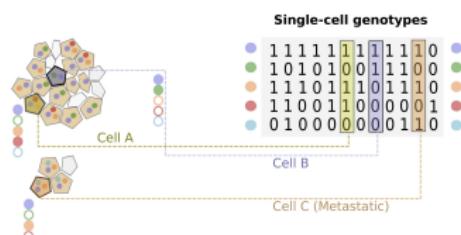
A Tumor phylogeny



B Multi-region bulk sequencing

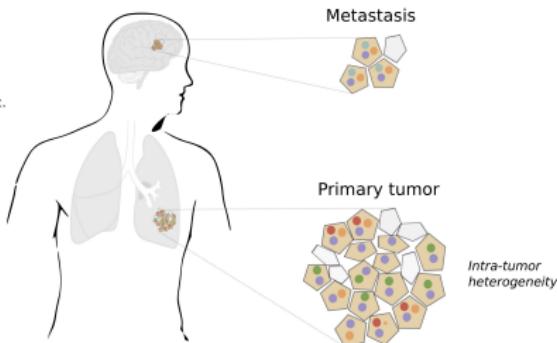
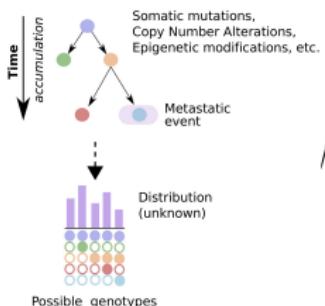


C Single-cell sequencing

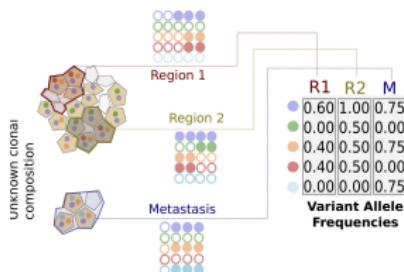


SINGLE CELL DATA

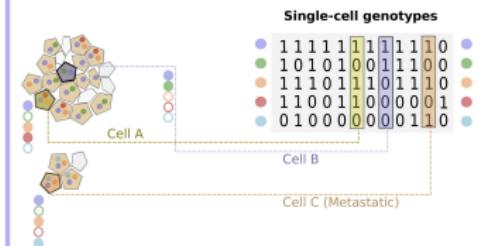
A Tumor phylogeny



B Multi-region bulk sequencing



C Single-cell sequencing



FROM BIOLOGICAL ASSUMPTIONS TO MODELS

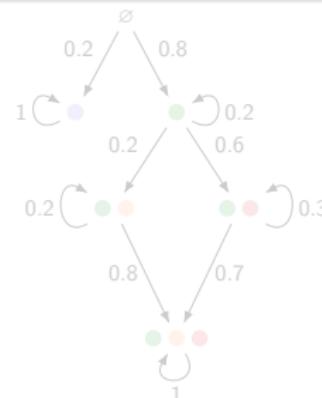
CANCER PROGRESSION MARKOV CHAINS (CPMC)

BIOLOGICAL ASSUMPTIONS

- Along cancer progression cells accumulate gene mutations
 - ∅ Initially there are no mutated genes
- MC New mutations probabilistically depends only on the current genotype
- ▽ Each time a minimal number of mutations is acquired

MODELS: CPMC

- The states are the genotypes
- The healthy genotype is the source
- It is a Discrete Time Markov Chain
- It is acyclic and anti-transitive



FROM BIOLOGICAL ASSUMPTIONS TO MODELS

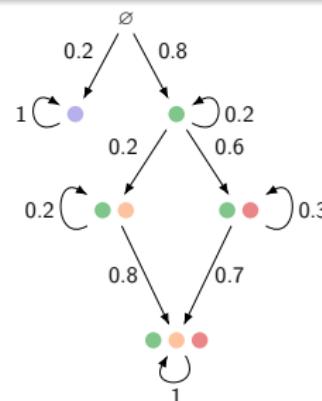
CANCER PROGRESSION MARKOV CHAINS (CPMC)

BIOLOGICAL ASSUMPTIONS

- Along cancer progression cells accumulate gene mutations
 - ∅ Initially there are no mutated genes
- MC** New mutations probabilistically depends only on the current genotype
- ▽ Each time a minimal number of mutations is acquired

MODELS: CPMC

- ▶ The states are the genotypes
- ▶ The healthy genotype is the source
- ▶ It is a Discrete Time Markov Chain
- ▶ It is acyclic and anti-transitive



FROM BIOLOGICAL ASSUMPTIONS TO MODELS

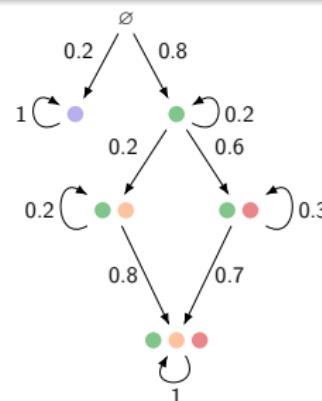
CANCER PROGRESSION MARKOV CHAINS (CPMC)

BIOLOGICAL ASSUMPTIONS

- U Along cancer progression cells accumulate gene mutations
- Ø Initially there are no mutated genes
- MC New mutations probabilistically depends only on the current genotype
- ▽ Each time a minimal number of mutations is acquired

MODELS: CPMC

- The states are the genotypes
- The healthy genotype is the source
- It is a Discrete Time Markov Chain
- It is acyclic and anti-transitive



FROM BIOLOGICAL ASSUMPTIONS TO MODELS

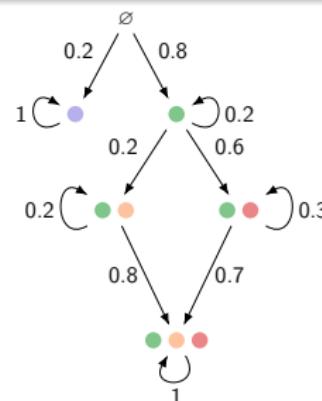
CANCER PROGRESSION MARKOV CHAINS (CPMC)

BIOLOGICAL ASSUMPTIONS

- Along cancer progression cells accumulate gene mutations
 - ∅ Initially there are no mutated genes
- MC New mutations probabilistically depends only on the current genotype
- ▽ Each time a minimal number of mutations is acquired

MODELS: CPMC

- ▶ The states are the genotypes
- ▶ The healthy genotype is the source
- ▶ It is a Discrete Time Markov Chain
- ▶ It is acyclic and anti-transitive



FROM BIOLOGICAL ASSUMPTIONS TO MODELS

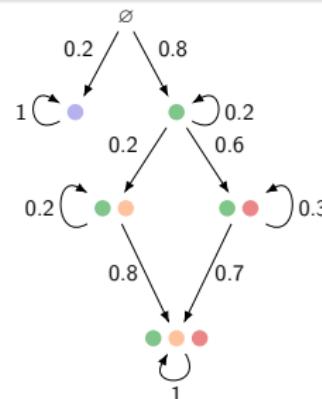
CANCER PROGRESSION MARKOV CHAINS (CPMC)

BIOLOGICAL ASSUMPTIONS

- Along cancer progression cells accumulate gene mutations
 - ∅ Initially there are no mutated genes
- MC New mutations probabilistically depends only on the current genotype
- ▽ Each time a minimal number of mutations is acquired

MODELS: CPMC

- ▶ The states are the genotypes
- ▶ The healthy genotype is the source
- ▶ It is a Discrete Time Markov Chain
- ▶ It is acyclic and anti-transitive



FROM BIOLOGICAL ASSUMPTIONS TO MODELS

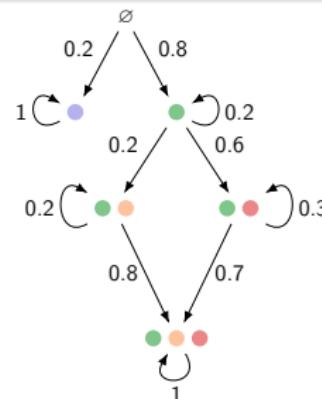
CANCER PROGRESSION MARKOV CHAINS (CPMC)

BIOLOGICAL ASSUMPTIONS

- U Along cancer progression cells accumulate gene mutations
 - Ø Initially there are no mutated genes
- MC New mutations probabilistically depends only on the current genotype
- ▽ Each time a minimal number of mutations is acquired

MODELS: CPMC

- The states are the genotypes
- The healthy genotype is the source
- It is a Discrete Time Markov Chain
- It is acyclic and anti-transitive



WHAT'S NEW?

MAIN DIFFERENCES W.R.T. LITERATURE

We propose a method that:

- ▶ refers to **single cell DNA-seq data**
- ▶ assumes **clean and rich data**
- ▶ points out where **more knowledge is needed**
- ▶ exploits a **deterministic** approach, which can be extended to include **statistical methods** and **expert knowledge**
- ▶ is **patient-driven** and suitable to study **treatment effects**

TWO RECONSTRUCTION PROBLEMS

SINGLE TIME INPUT:

$$D_i$$

a single cell data matrix

MANY TIME INPUT:

$$D_0, D_1, \dots, D_t$$

a time sequence of matrices



OUTPUT:

A CPMC whose simulation would generate the data

In such CPMC time ticks at each mutational event (no self-loops)

REMARK

- ▶ It is not a standard Markov Chain reconstruction problem
- ▶ There can be infinitely many solutions

Appendix

TWO RECONSTRUCTION PROBLEMS

SINGLE TIME INPUT:

$$D_i$$

a single cell data matrix

MANY TIME INPUT:

$$D_0, D_1, \dots, D_t$$

a time sequence of matrices



OUTPUT:

A CPMC whose simulation would generate the data

In such CPMC time ticks at each mutational event (no self-loops)

REMARK

- ▶ It is not a standard Markov Chain reconstruction problem
- ▶ There can be infinitely many solutions

Appendix

OUR INFERENCE ALGORITHM

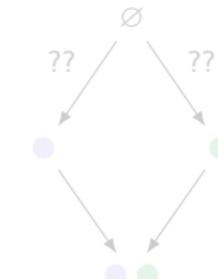
FUNDAMENTALS

- ▶ the **topology** is directly induced by \cup and \triangledown
- ▶ **without convergencies**, the probabilities can be directly computed thanks to the topology, MC, and **Bayes' theorem**
- ▶ with **convergent paths**, we exploit **heuristics** to estimate *backward probabilities* Appendix

CONVERGENCY AMBIGUITIES

s_1	1	0
s_2	0	1
s_3	1	1
s_4	1	1

⇒



OUR INFERENCE ALGORITHM

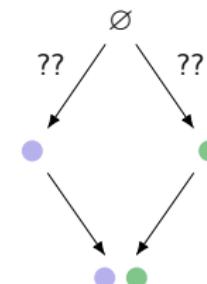
FUNDAMENTALS

- ▶ the **topology** is directly induced by \cup and \triangledown
- ▶ **without convergencies**, the probabilities can be directly computed thanks to the topology, MC, and **Bayes' theorem**
- ▶ with **convergent paths**, we exploit **heuristics** to estimate *backward probabilities* Appendix

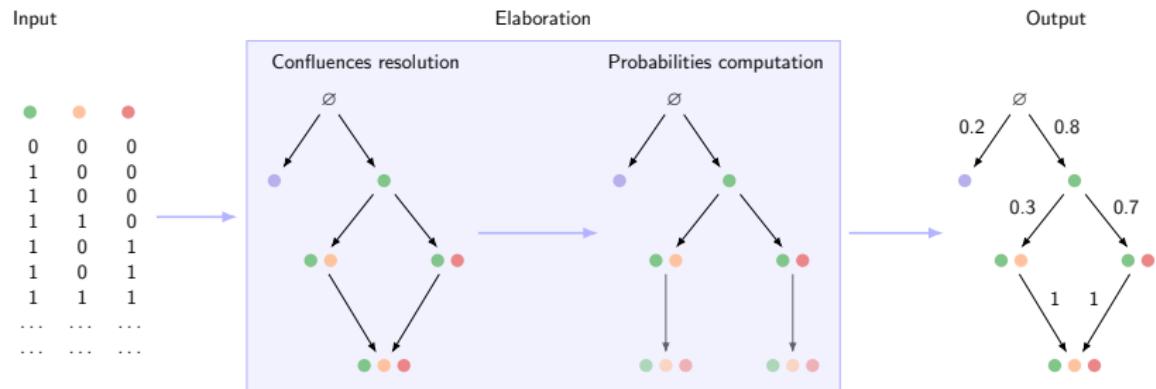
CONVERGENCY AMBIGUITIES

	●	●
s_1	1	0
s_2	0	1
s_3	1	1
s_4	1	1

⇒



CIMICE Tool

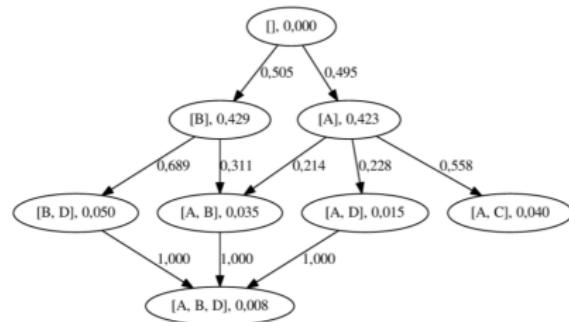
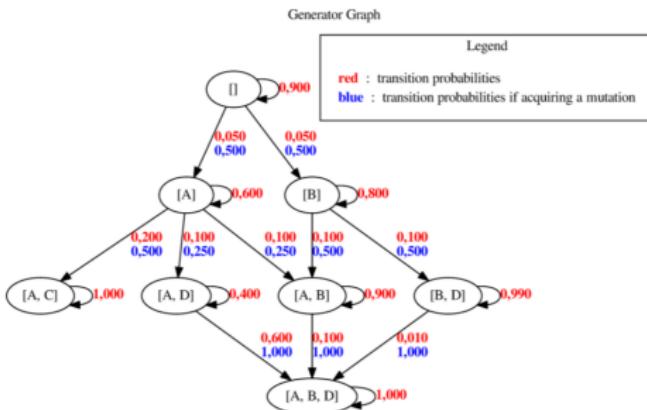


<https://github.com/redsnic/tumorEvolutionWithMarkovChains>

SYNTHETIC MODELS AND TESTS

Random walks on a given CPMC can be used to either:

- ▶ generate more data on a specific case or
- ▶ generate the genotypes of the “artificial” dataset and then test our methodology

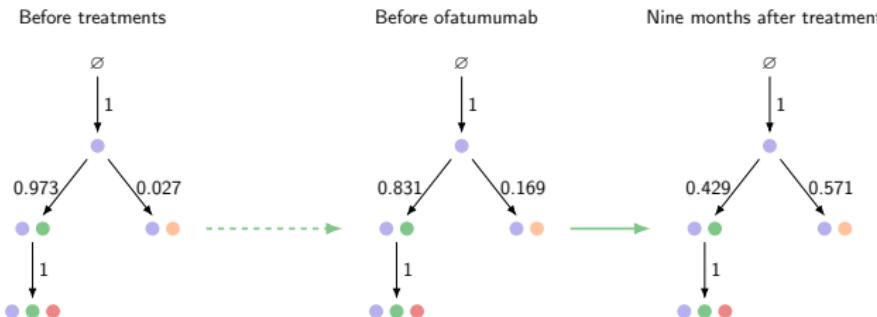


REAL DATA

PRESENT SITUATION

- ▶ small single cell datasets are available
- ▶ in Ogundijo and Wang, BMC Bioinf. 20:6(2019) a method for inferring genotype proportions from is described
- ▶ their output is a suitable input for us

CHRONIC LYMPHOCYTIC LEUKEMIA - PATIENT CLL077



Drugs: chlorambuci, fludarabine and cyclophosphamide, ofatumumab

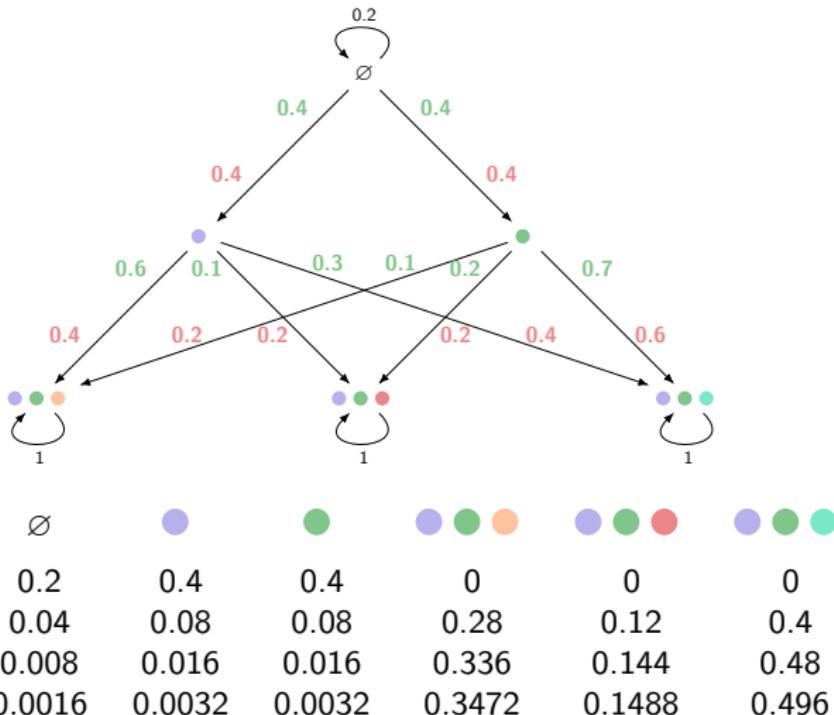
CONCLUSIONS

- ▶ We proposed a method for tumor phylogeny reconstruction using DTMCs
- ▶ The reconstruction algorithm is time efficient, can easily be enhanced with statistical methods, and it is suitable for very large datasets
- ▶ We are planning to include treatments information directly in the model
- ▶ It is possible to use the artificial data produced by CIMICE to test similar tools

RELATED WORKS

- ▶ *The evolution of tumour phylogenetics: principles and practice*
R. Schwartz and A. A. Schäffer, Nature Reviews Genetics
18:213-229, 2017
- ▶ *Learning mutational graphs of individual tumour evolution from single-cell and multi-region sequencing data* D.
Ramazzotti, A. Graudenzi, L. De Sano, M. Antoniotti, and G.
Caravagna, BMC Bioinformatics 20:210, 2019
- ▶ *SeqClone: sequential Monte Carlo based inference of tumor subclones* E. Ogundijo and Xiaodong Wang, BMC
Bioinformatics 20:6, 2019

INFINITELY MANY SOLUTIONS



[Back to Presentation](#)

TECHNICAL DETAILS

- ▶ the data are $P[X(k) = S]$ for each S genotype
- ▶ k is unknown but large enough
- ▶ the probability of the edge (S, T) on the chain with self-loops is

$$P(S, T) = P[X(k) = T | X(k - 1) = S]$$

does not depend on k

- ▶ we are interested in the probability of the edge (S, T) on the chain without self-loops

$$JP(S, T) = \frac{P(S, T)}{\text{norm}(S)}$$

[Back to Presentation](#)

TECHNICAL DETAILS

- ▶ by acyclicity and anti-transitivity of CPMCs

$$P(S, T) = \frac{\sum_{i=1}^k P[f_i(T)] * P[X(i-1) = S | f_i(T)]}{\sum_{i=0}^{k-1} P[X(i) = S]}$$

- ▶ which without self-loops become

$$JP(S, T) = \frac{\sum_{i=1}^k P[f_i(T)] * P[X(i-1) = S | f_i(T)]}{norm(S)}$$

- ▶ in the case **without convergencies** $P[X(i-1) = S | f_i(T)] = 1$ and

$$\sum_{i=1}^k P[f_i(T)] = P[X(\leq k) = T] = P[X(k) = T] + \sum_{Y \in Adj^-[T]} P[X(\leq k) = Y]$$

[Back to Presentation](#)

TECHNICAL DETAILS

- ▶ by acyclicity and anti-transitivity of CPMCs

$$P(S, T) = \frac{\sum_{i=1}^k P[f_i(T)] * P[X(i-1) = S | f_i(T)]}{\sum_{i=0}^{k-1} P[X(i) = S]}$$

- ▶ which without self-loops become

$$JP(S, T) = \frac{\sum_{i=1}^k P[f_i(T)] * P[X(i-1) = S | f_i(T)]}{norm(S)}$$

- ▶ in the case **with convergencies** the heuristics estimate

$$P[X(i-1) = S | f_i(T)] \text{ and } P[X(\leq k) = T]$$

[Back to Presentation](#)