

Mind the gap: Modelling evolution with indels

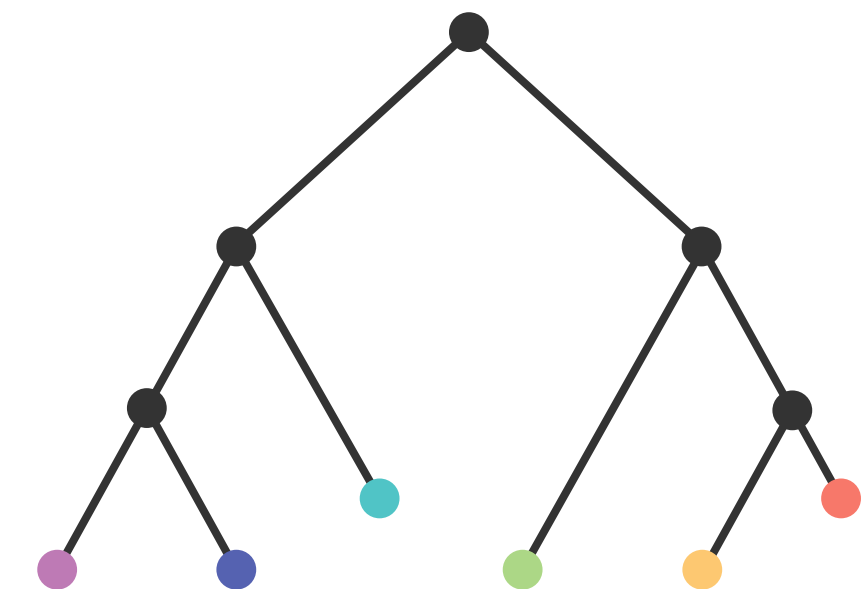
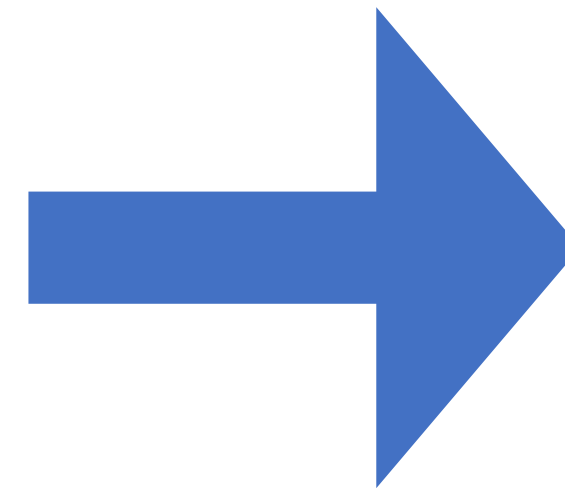
TAMING THE BEAST 2023

PHYLOGENETIC PIPELINES

EXPECTATION



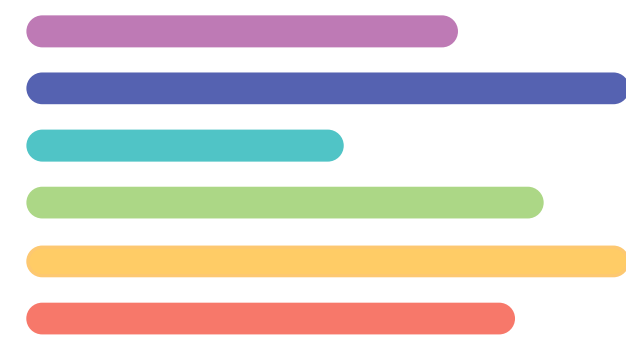
MSA



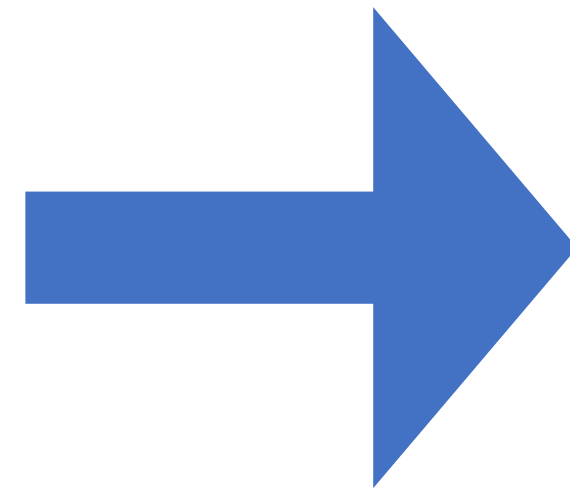
phylogenetic tree

PHYLOGENETIC PIPELINES

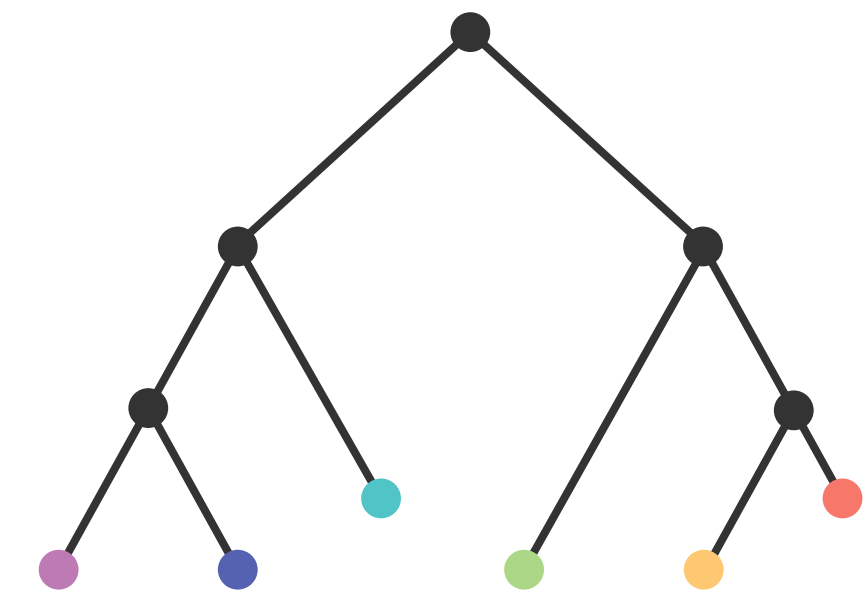
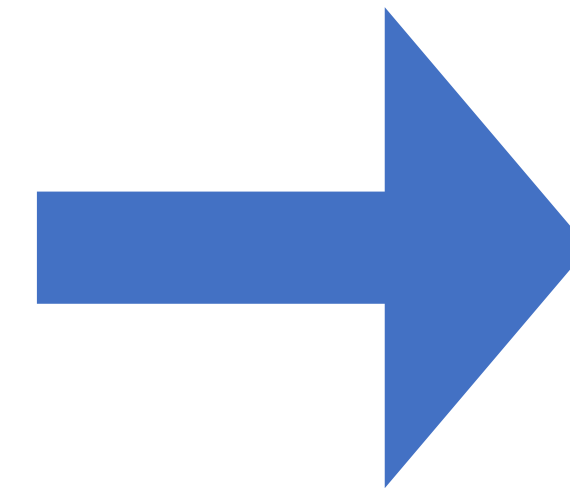
REALITY



unaligned
sequences



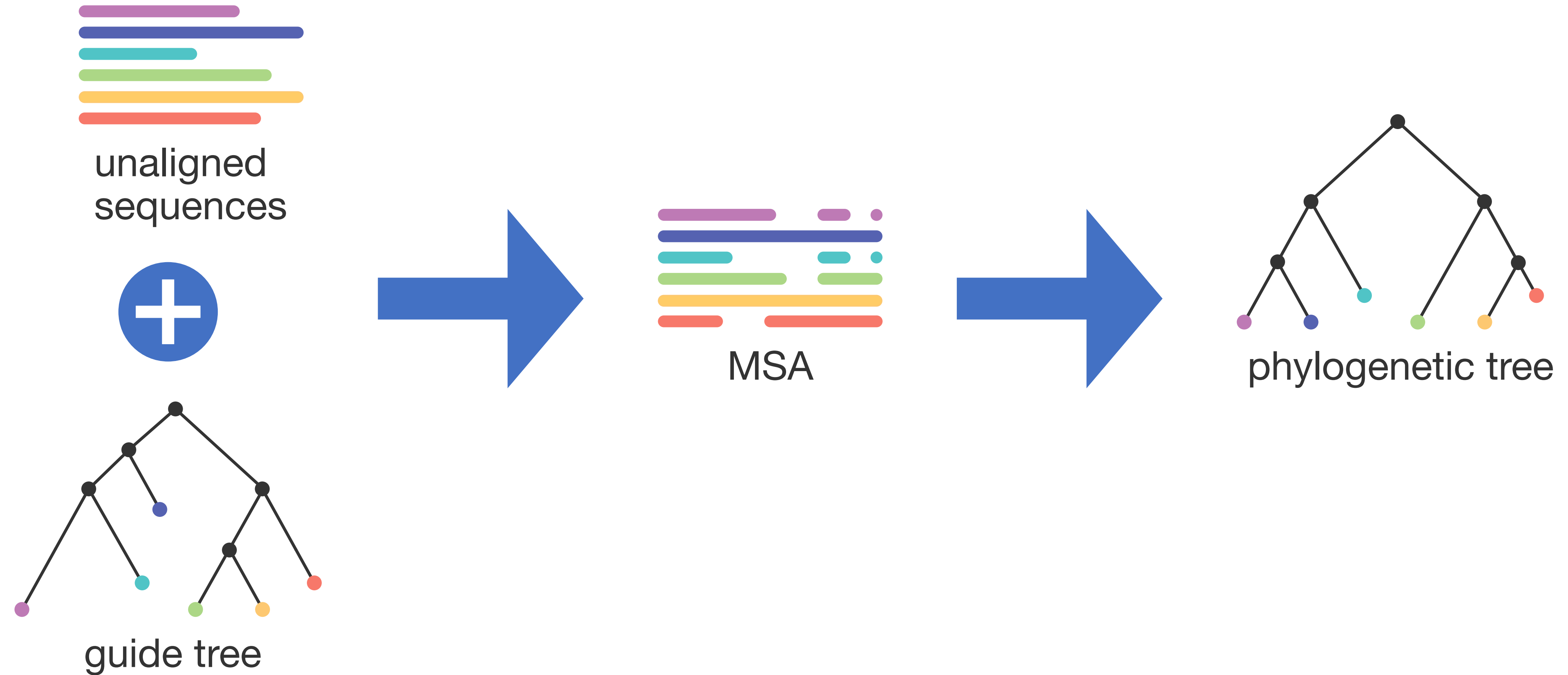
MSA



phylogenetic tree

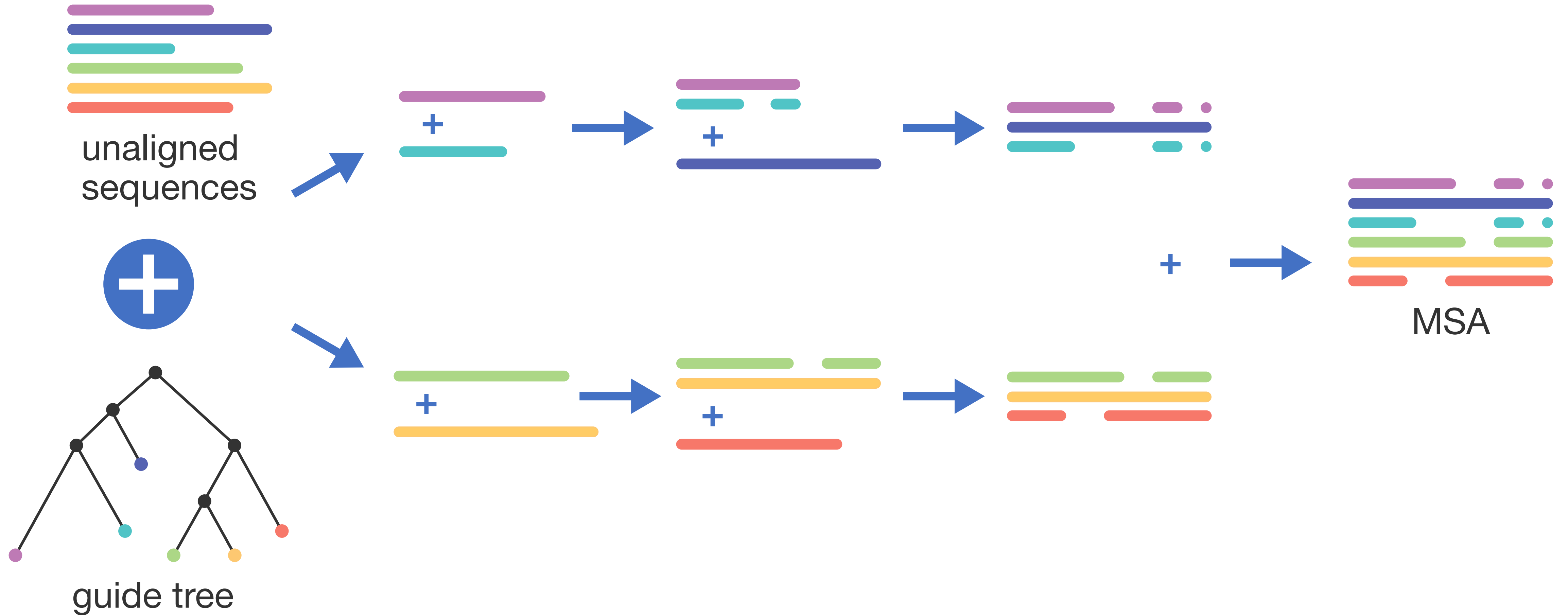
PHYLOGENETIC PIPELINES

REALITY



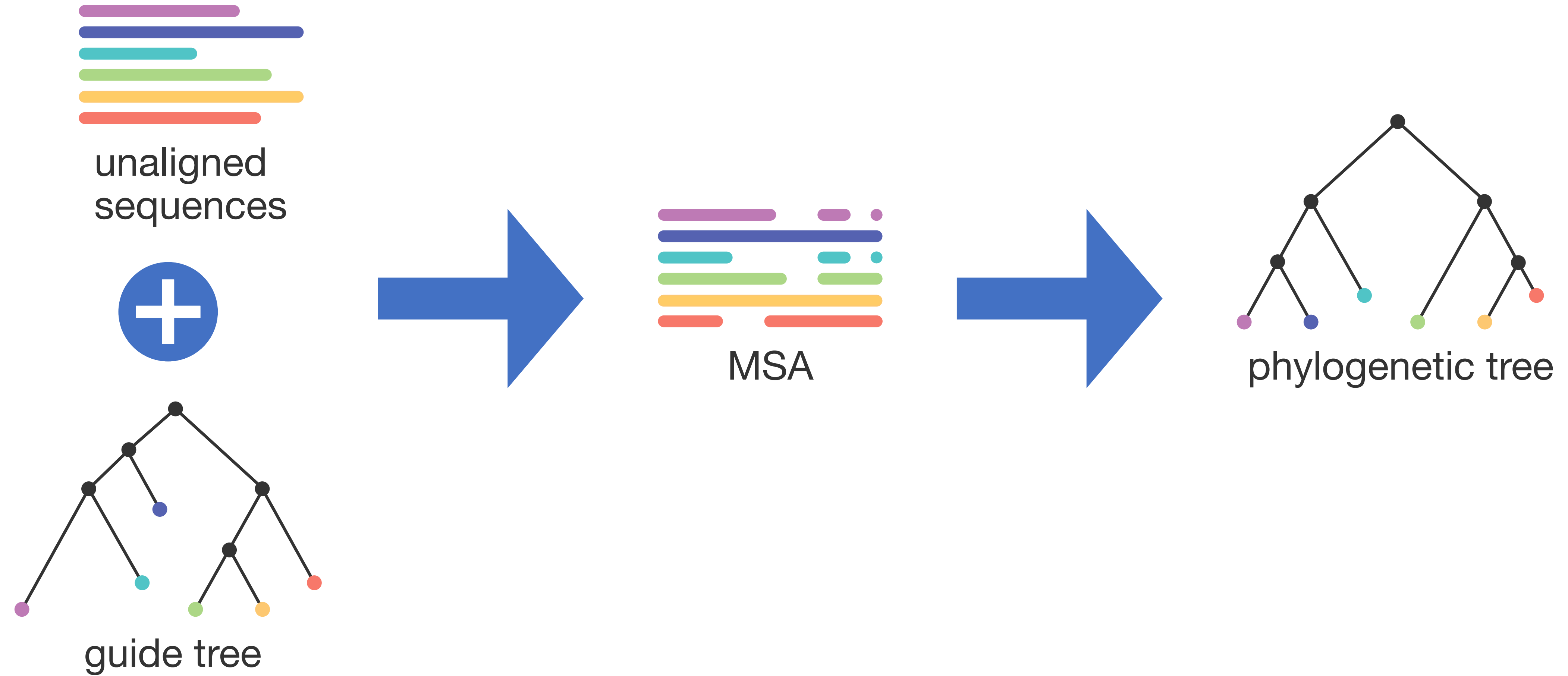
PHYLOGENETIC PIPELINES

REALITY



PHYLOGENETIC PIPELINES

REALITY

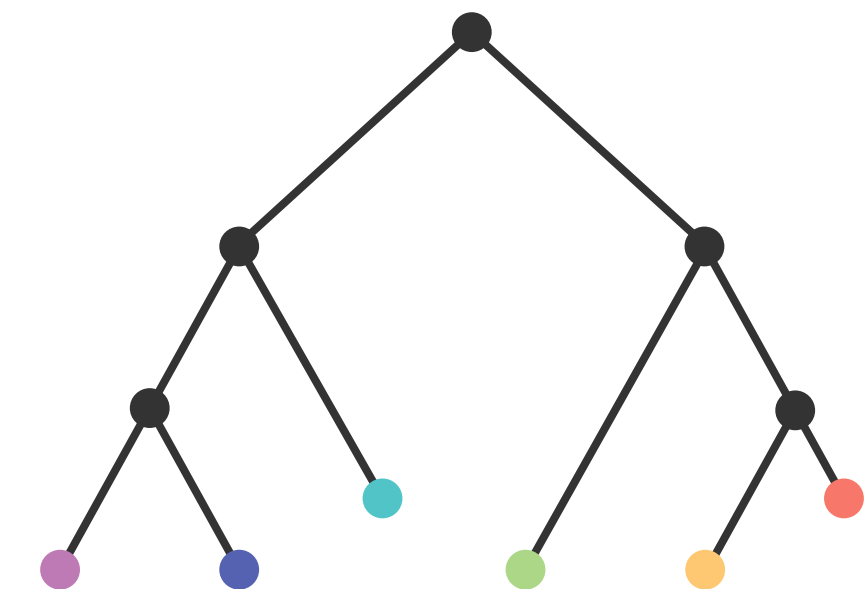
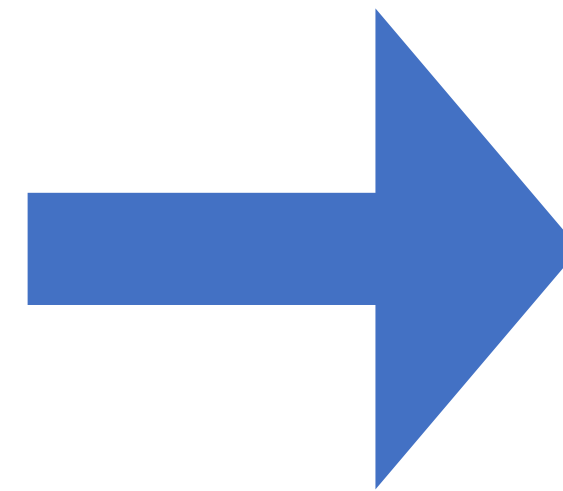


PHYLOGENETIC PIPELINES

REALITY

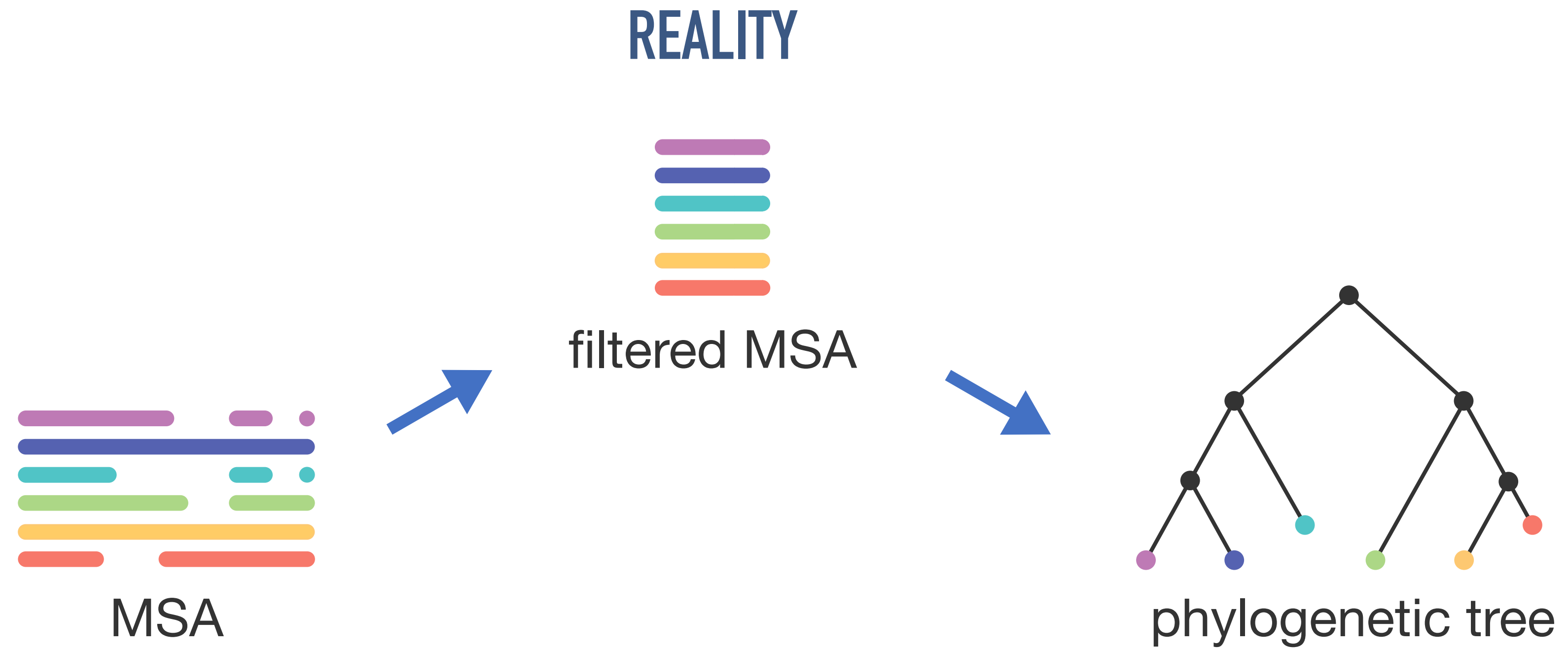


MSA

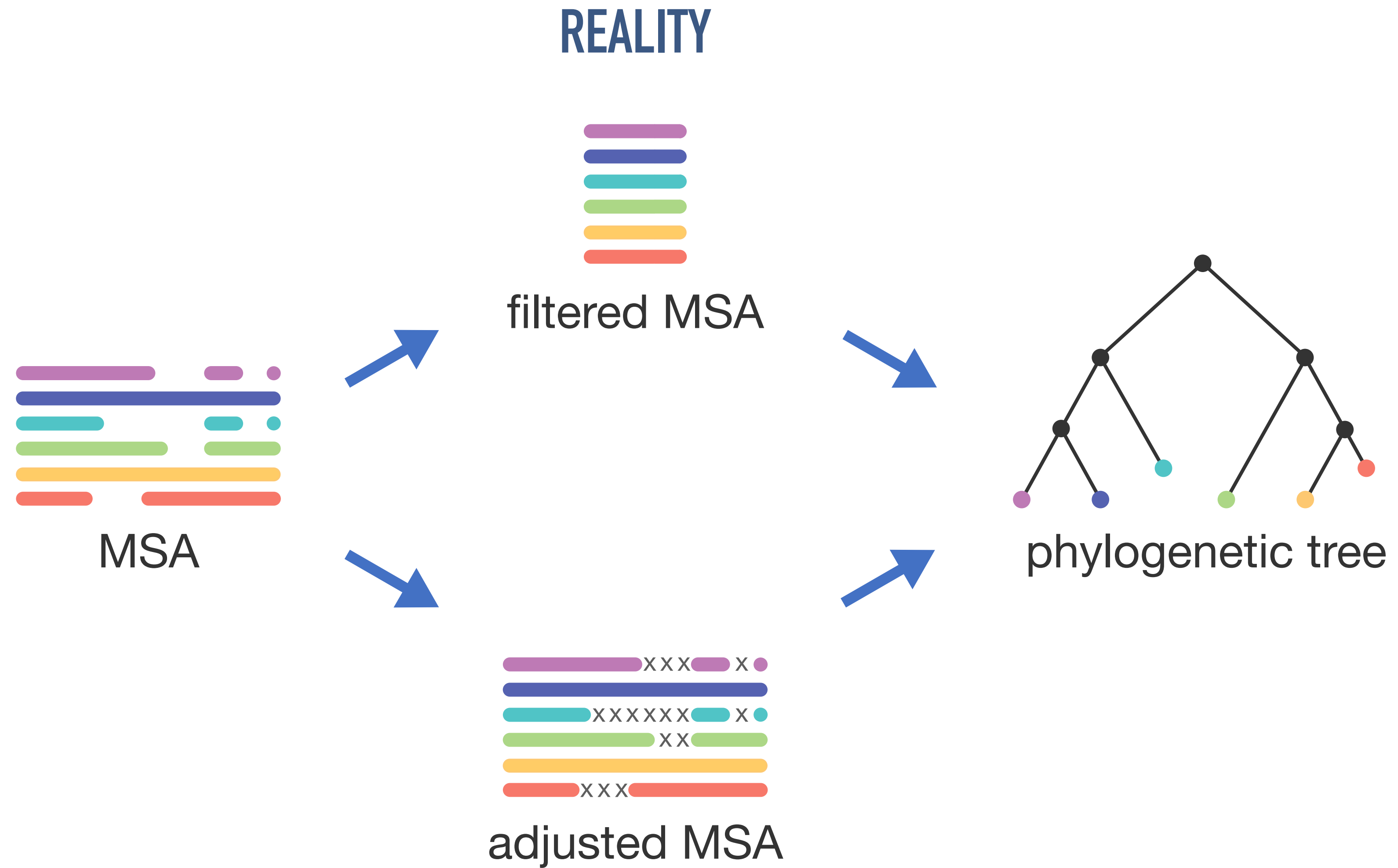


phylogenetic tree

PHYLOGENETIC PIPELINES

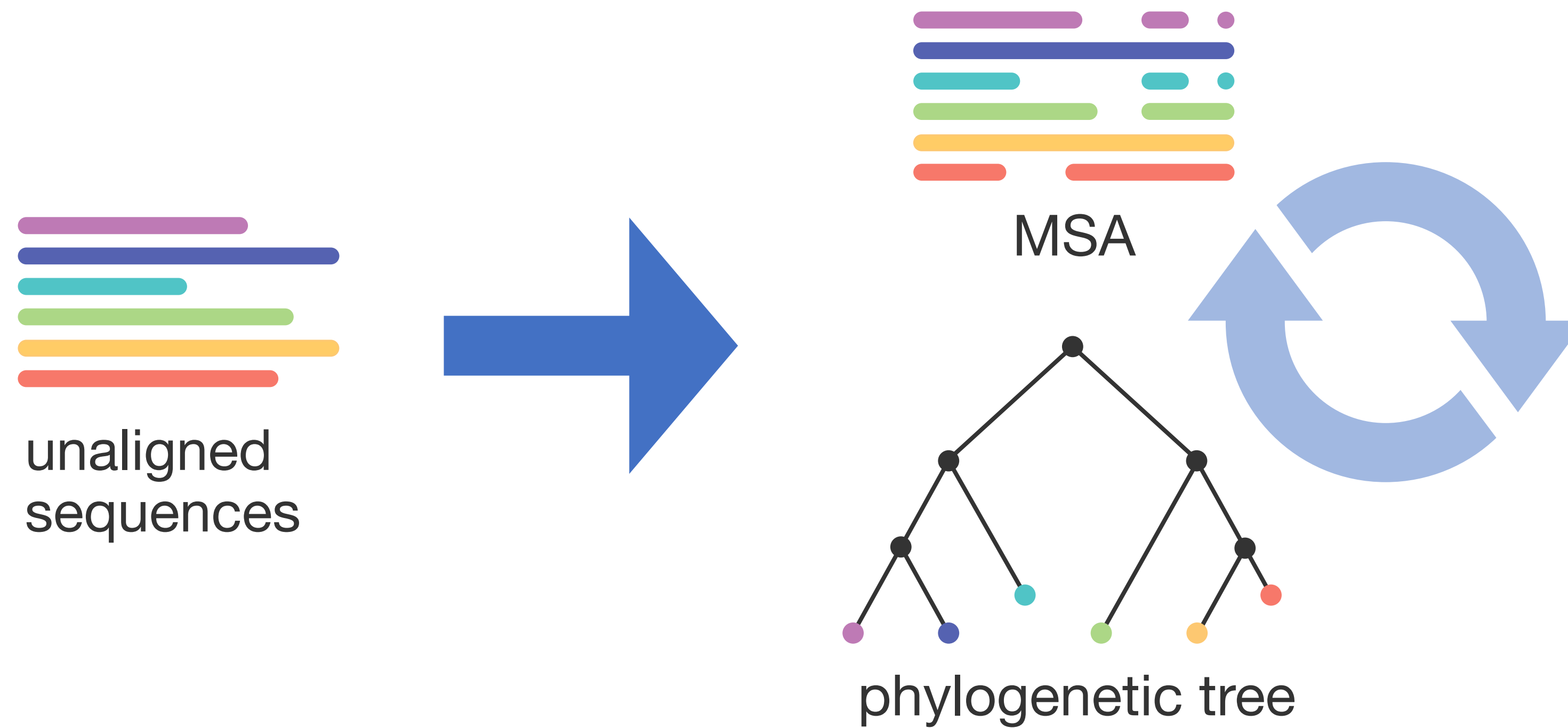


PHYLOGENETIC PIPELINES



PHYLOGENETIC PIPELINES

EXPECTATION 2.0



REMINDER: SUBSTITUTION MODELS

CONTINUOUS TIME MARKOV CHAINS

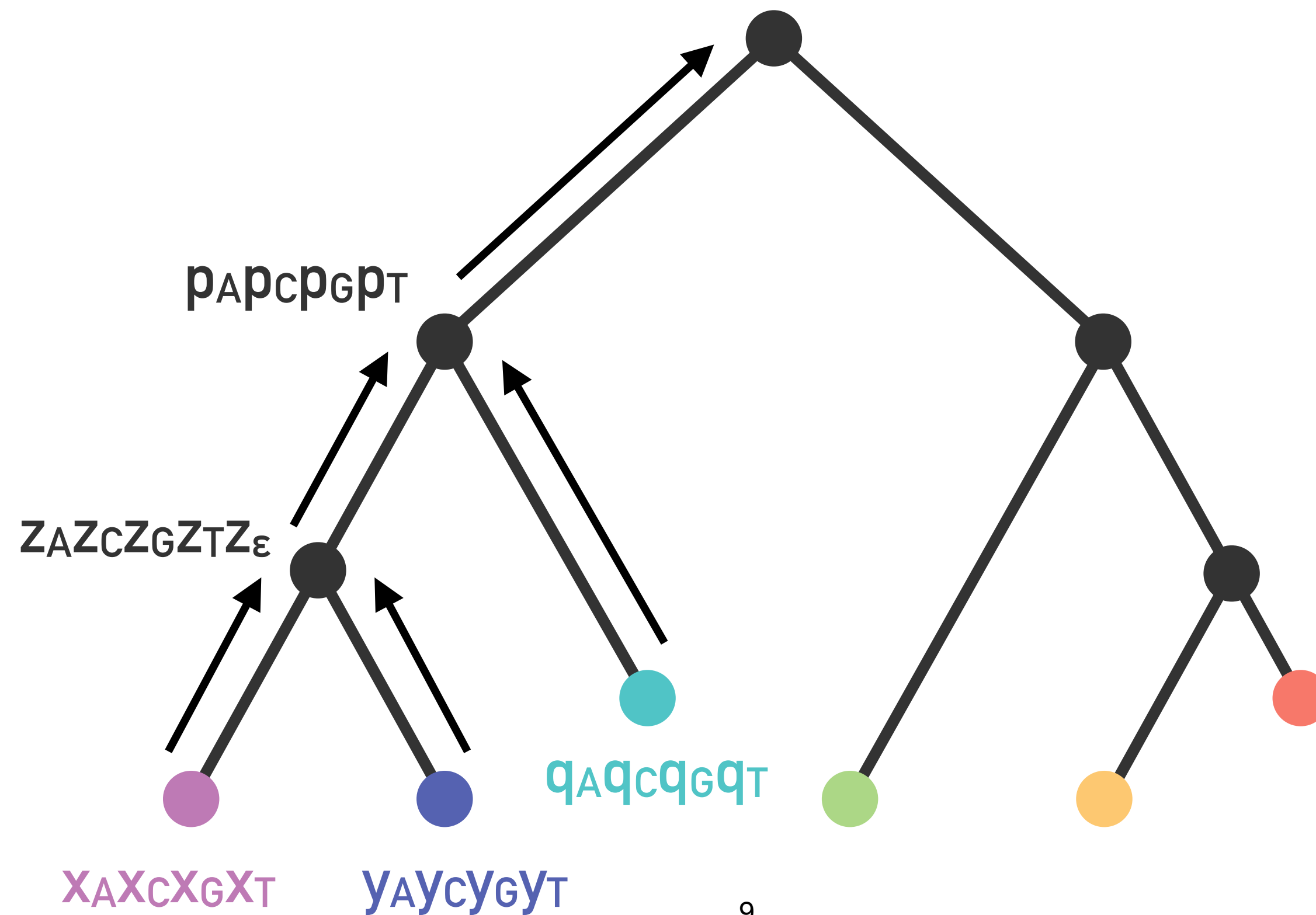
- Q — substitution rate matrix
- $P(t)$ — transition probability matrix
- $P(t) = e^{tQ}$
- $P(t + s) = e^{(t+s)Q} = e^{tQ}e^{sQ} = P(t)P(s)$

	T	C	A	G
T	...	$a\pi_C$	$b\pi_A$	$c\pi_G$
C	$a\pi_T$...	$d\pi_A$	$e\pi_G$
A	$b\pi_T$	$d\pi_C$...	$f\pi_G$
G	$c\pi_T$	$e\pi_C$	$f\pi_A$...

REMINDER: PHYLOGENETIC LIKELIHOOD

USING FELSENSTEIN'S PRUNING ALGORITHM

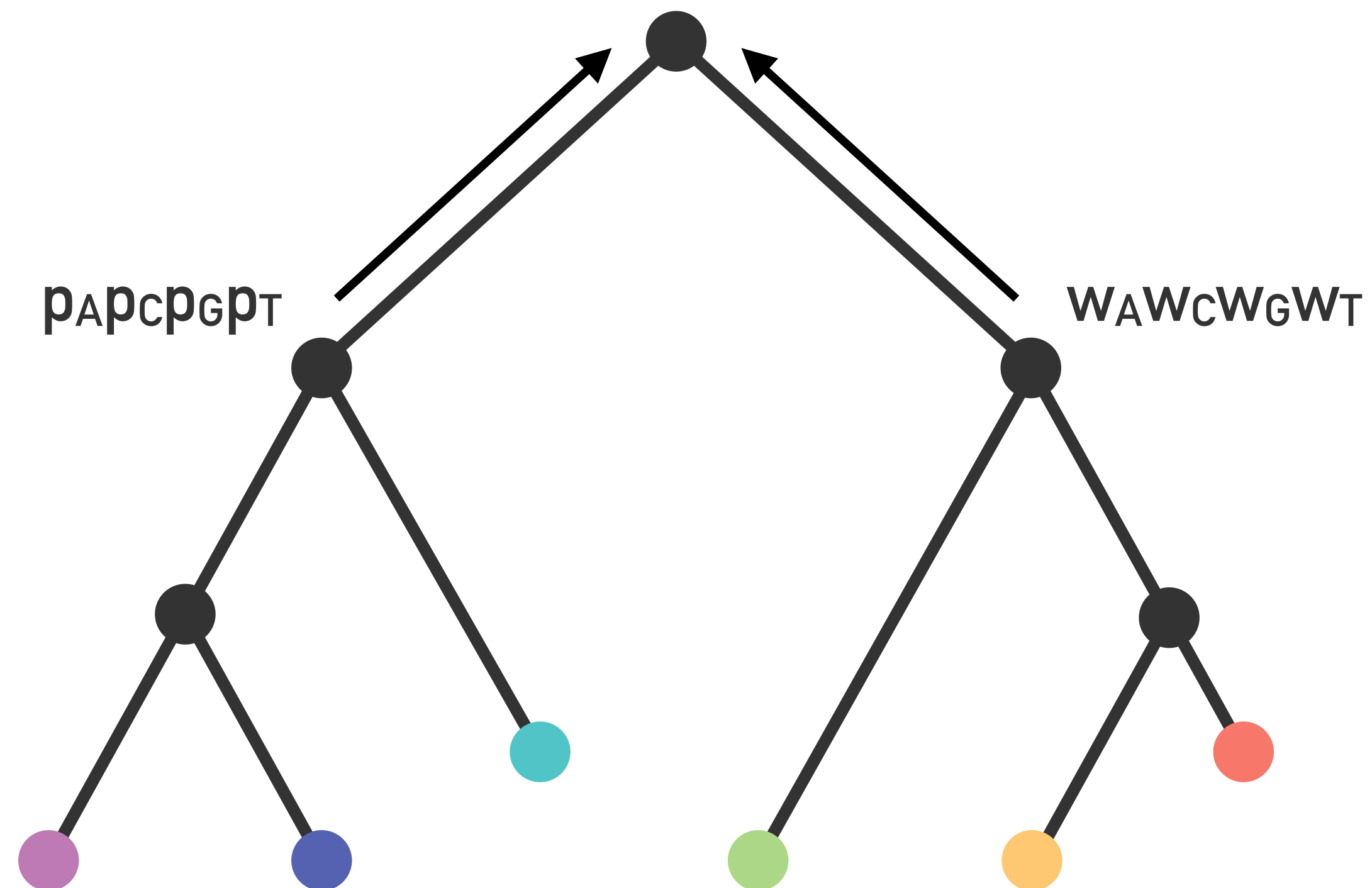
$$L(\text{tree}) = P(\text{matrix} \mid \text{tree}, \text{model})$$



REMINDER: PHYLOGENETIC LIKELIHOOD

USING FELSENSTEIN'S PRUNING ALGORITHM

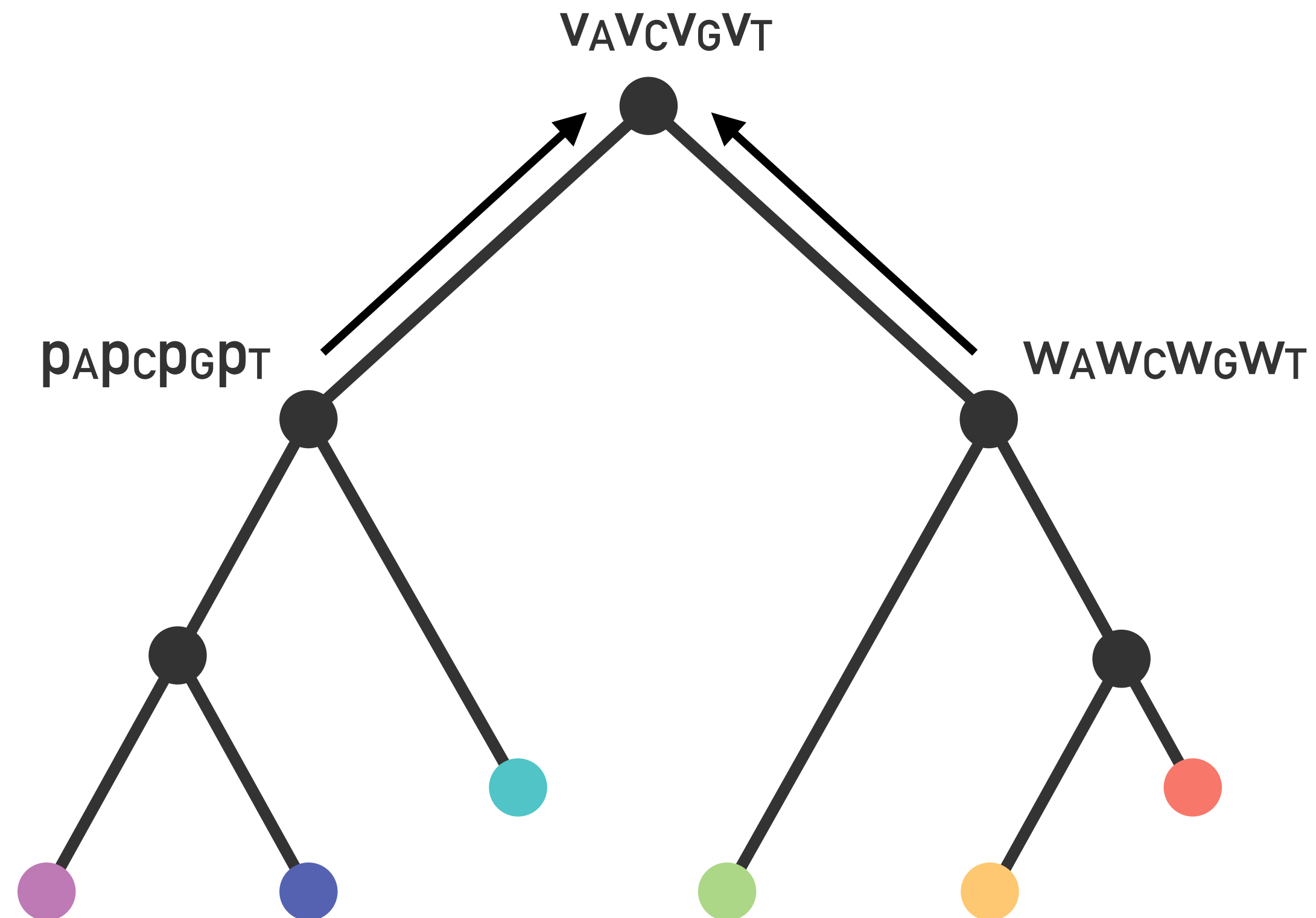
$$L(\text{tree}) = P(\text{data} | \text{tree}, \text{model})$$



REMINDER: PHYLOGENETIC LIKELIHOOD

USING FELSENSTEIN'S PRUNING ALGORITHM

$$L(\text{tree}) = P(\text{sequences} \mid \text{tree}, \text{matrix})$$

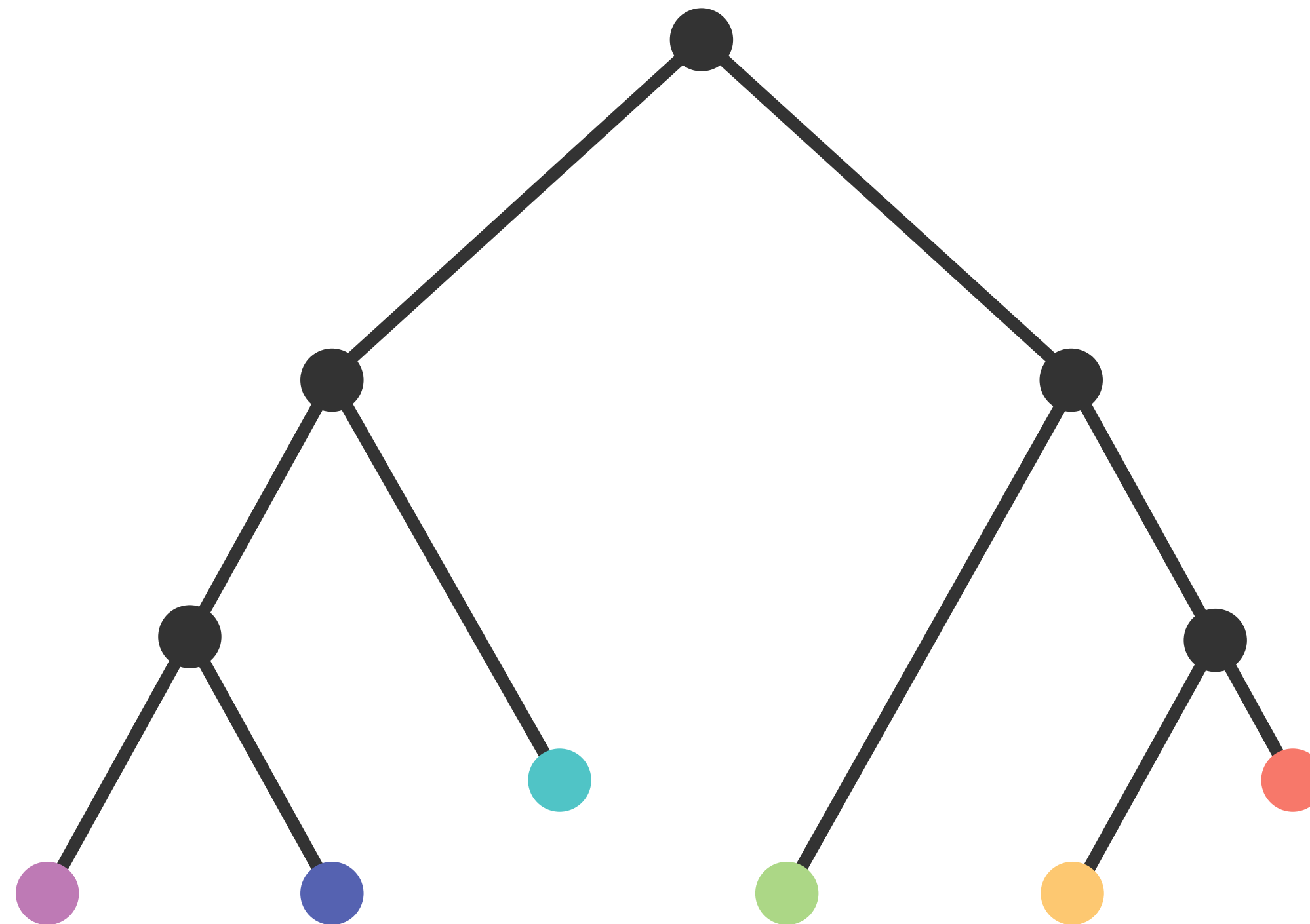


REMINDER: PHYLOGENETIC LIKELIHOOD

USING FELSENSTEIN'S PRUNING ALGORITHM

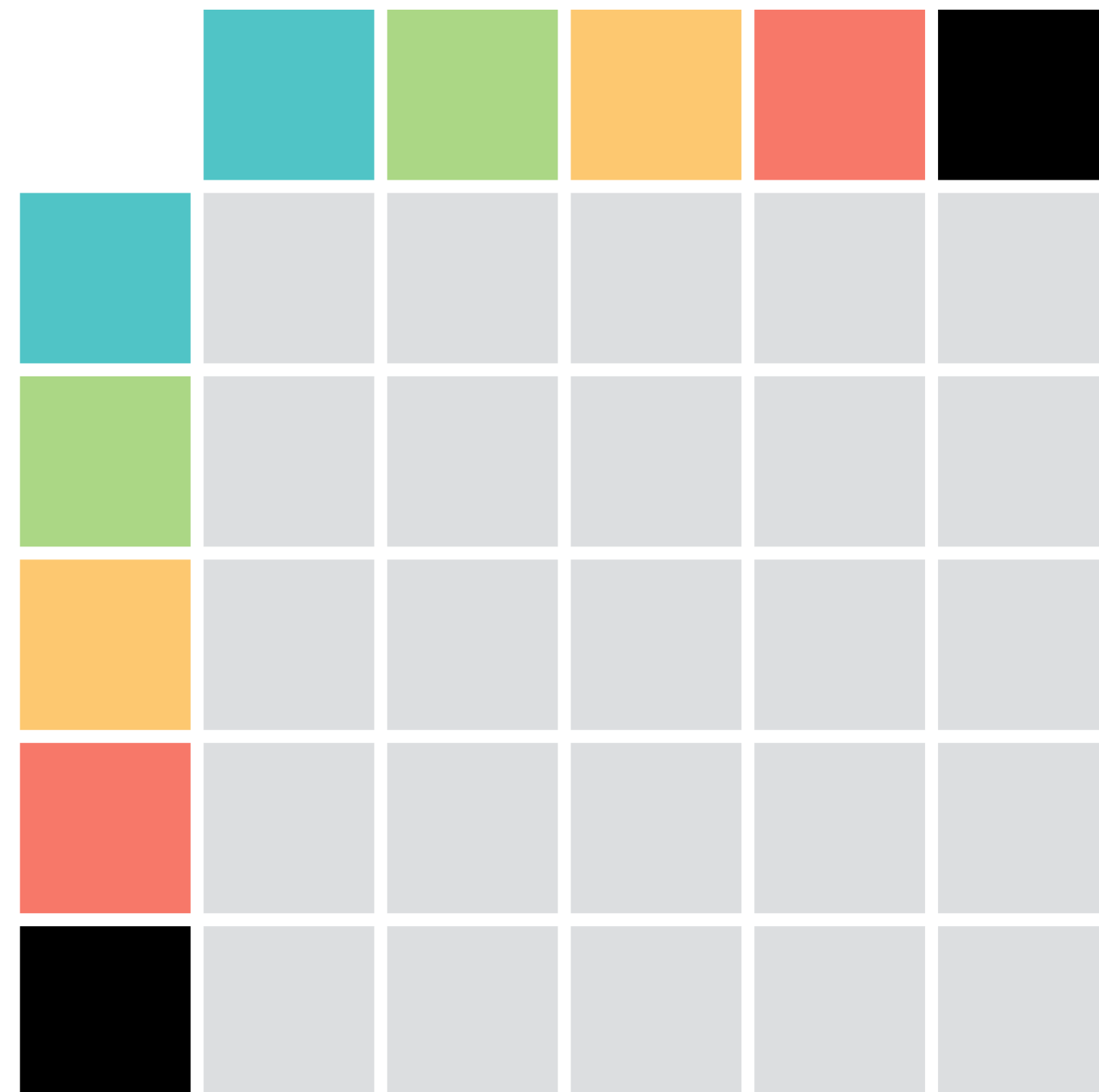
$$L(\text{tree}) = P(\text{sequences} \mid \text{tree}, \text{matrix})$$

Likelihood = $f(v_A v_C v_G v_T)$



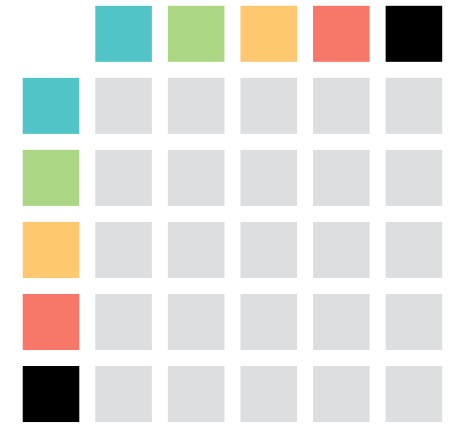
MODELLING INDELS

WHY NOT WITH MARKOV CHAINS?

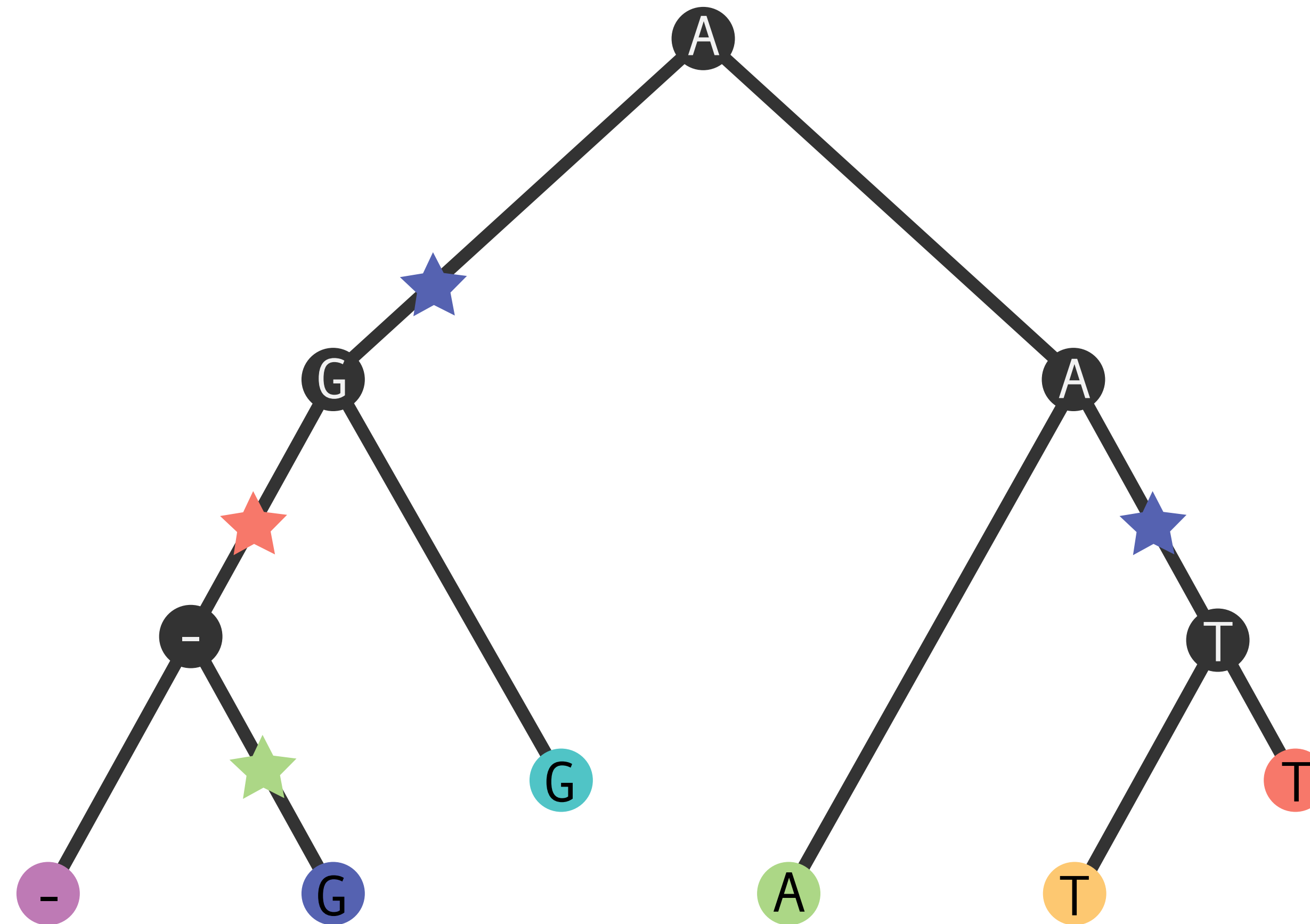


MODELLING INDELS

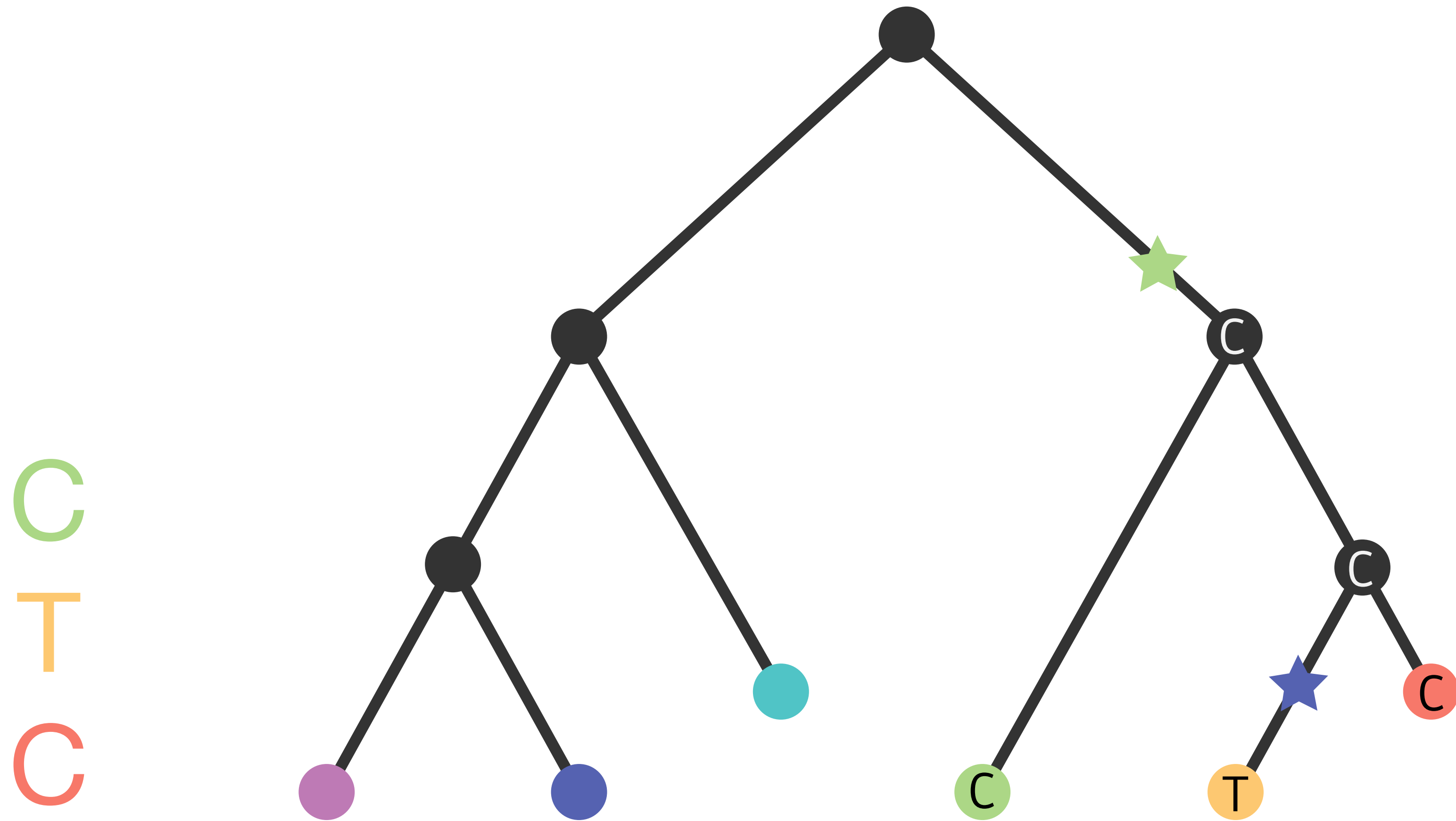
WHY NOT WITH MARKOV CHAINS?



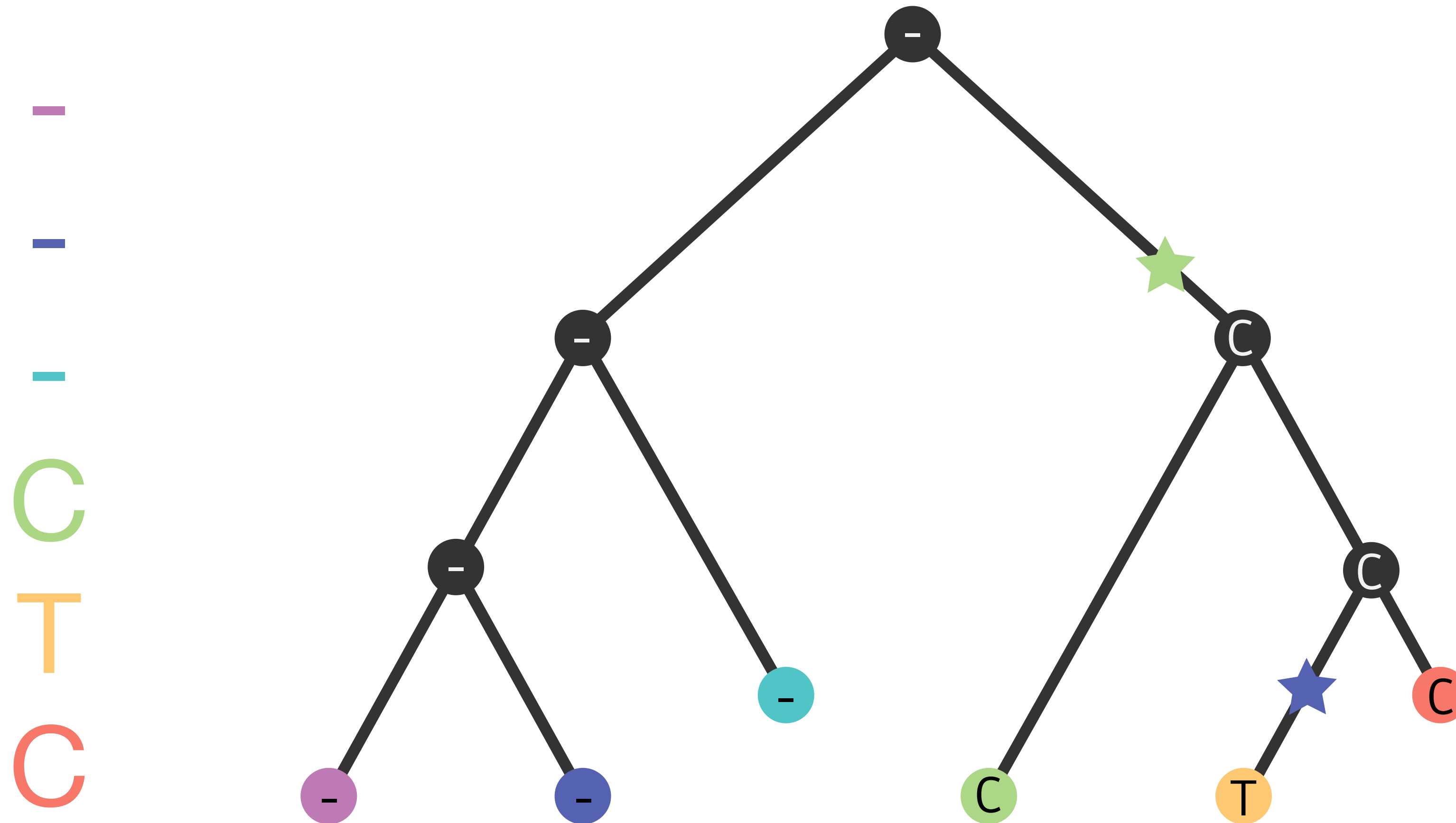
-
G
G
A
T
T



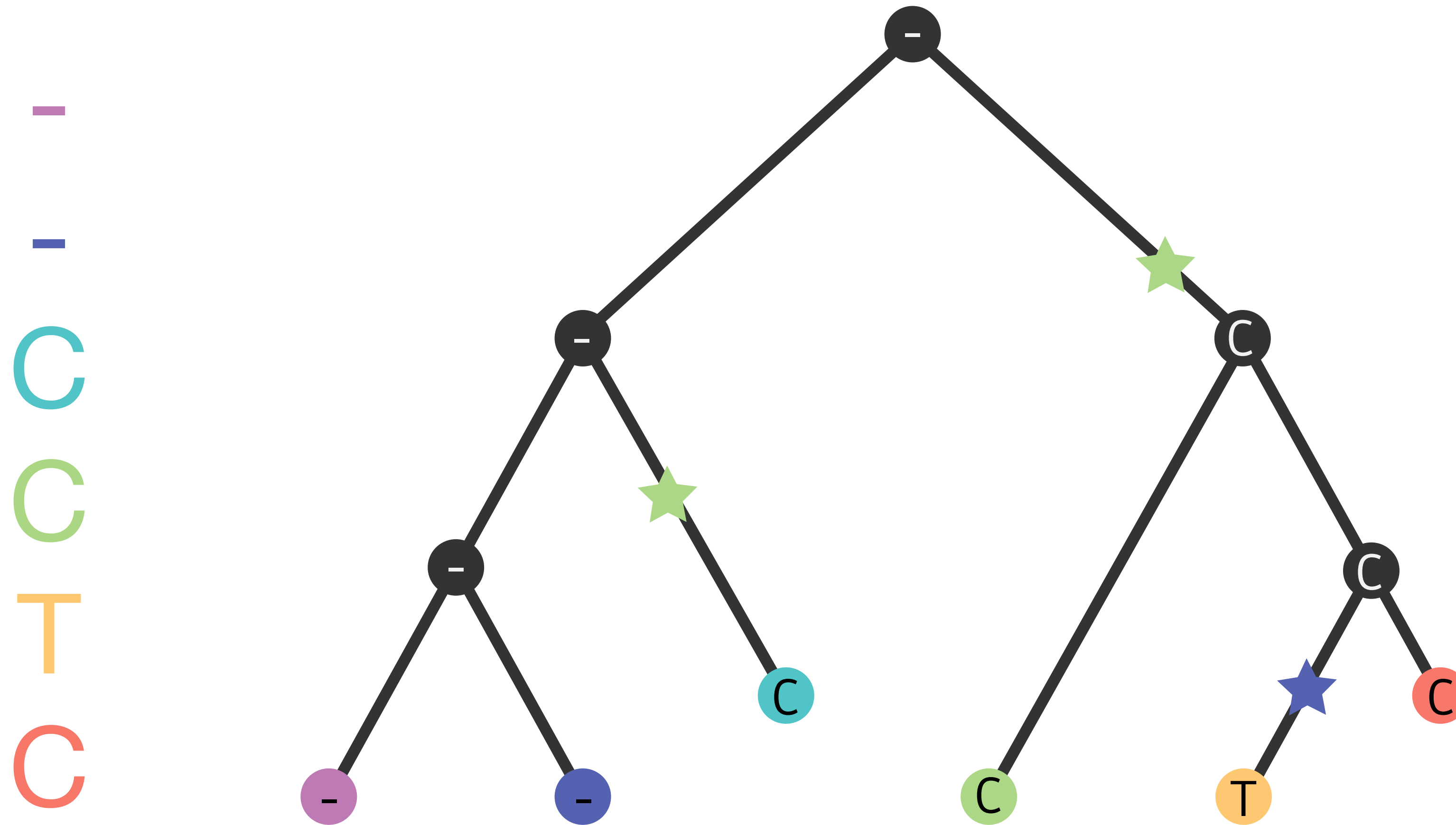
ANOTHER FUNKY INDEL SCENARIO



ANOTHER FUNKY INDEL SCENARIO



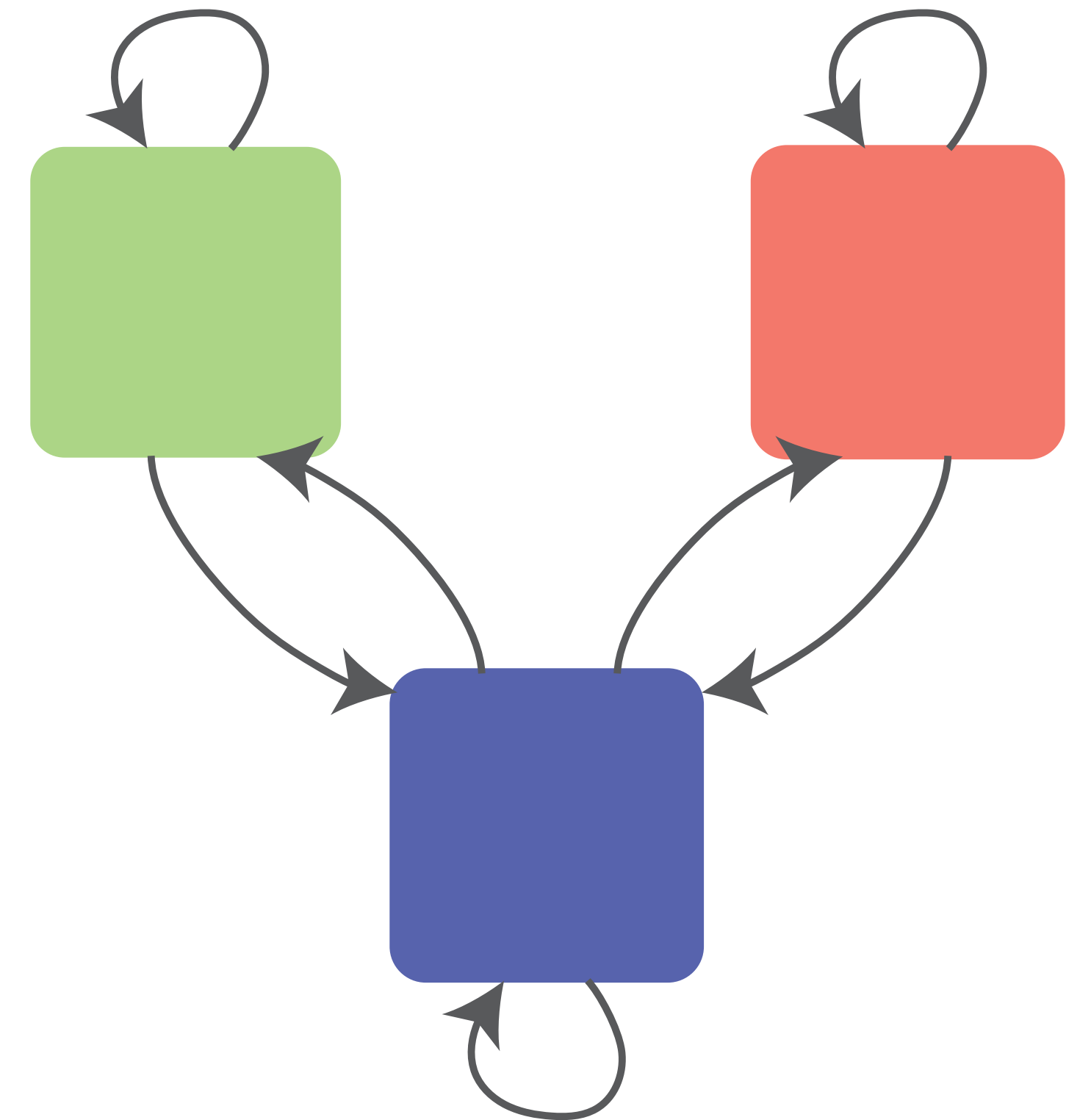
ANOTHER FUNKY INDEL SCENARIO



MODELLING INDELS

USING HMMS

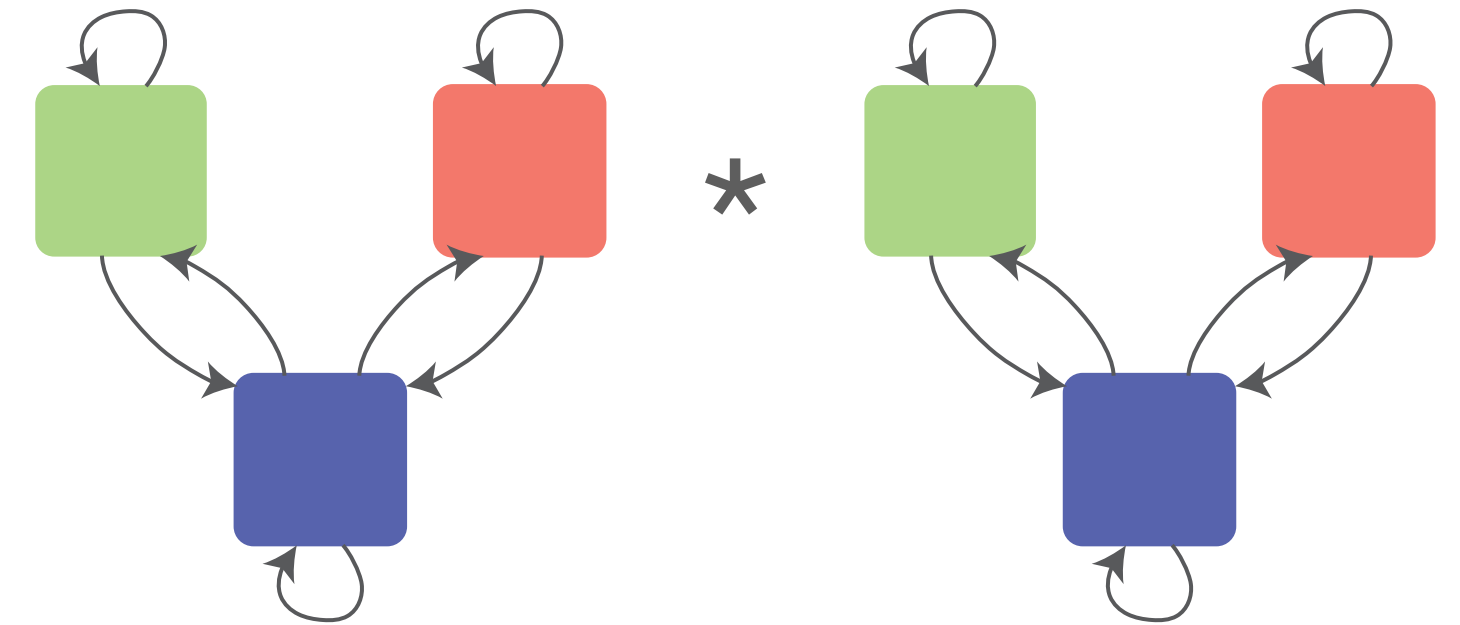
- Process (e.g. TKF91/92) described with an HMM



MODELLING INDELS

USING HMMS

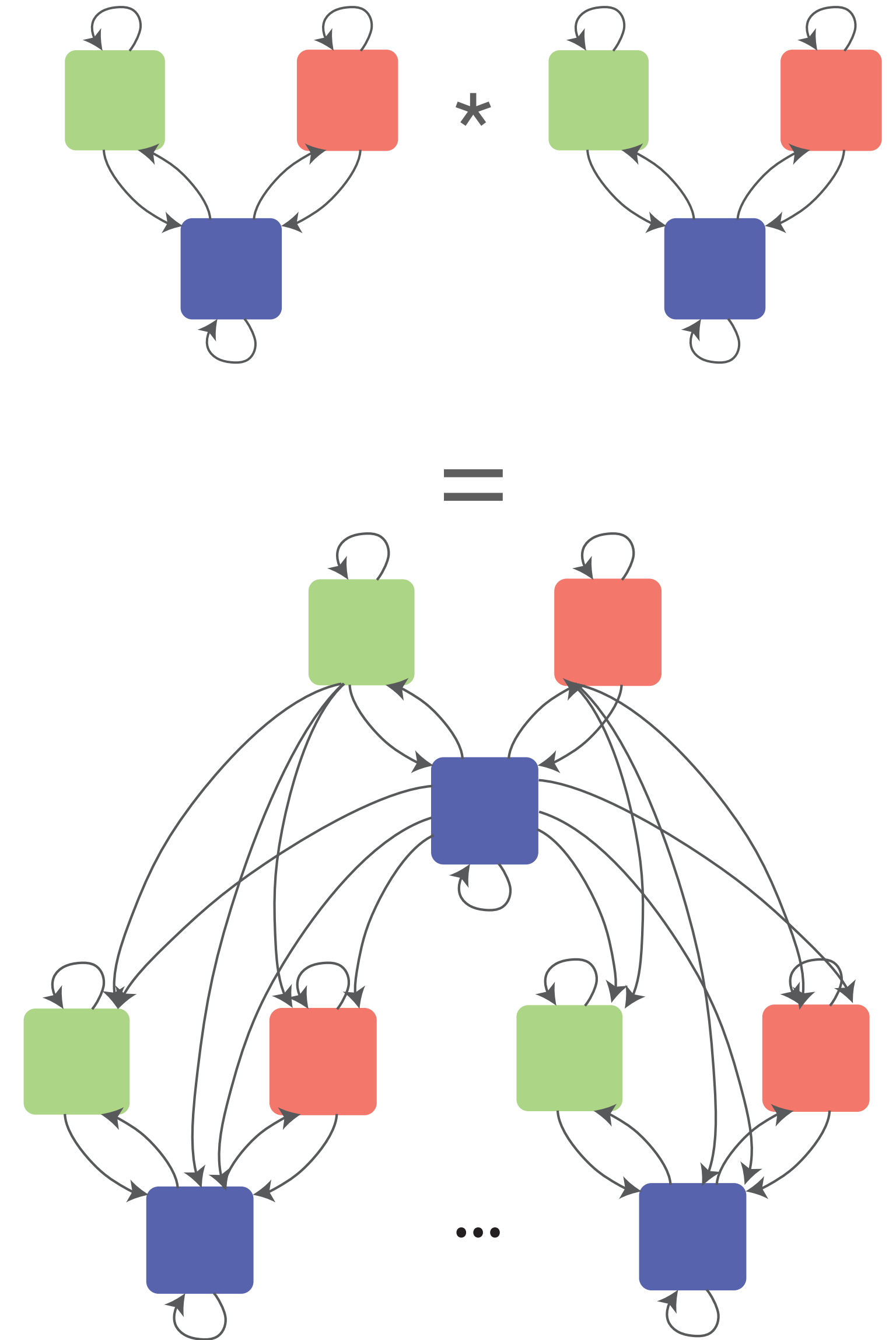
- Process (e.g. TKF91/92) described with an HMM
- Differentiating along a branch explodes the HMM



MODELLING INDELS

USING HMMS

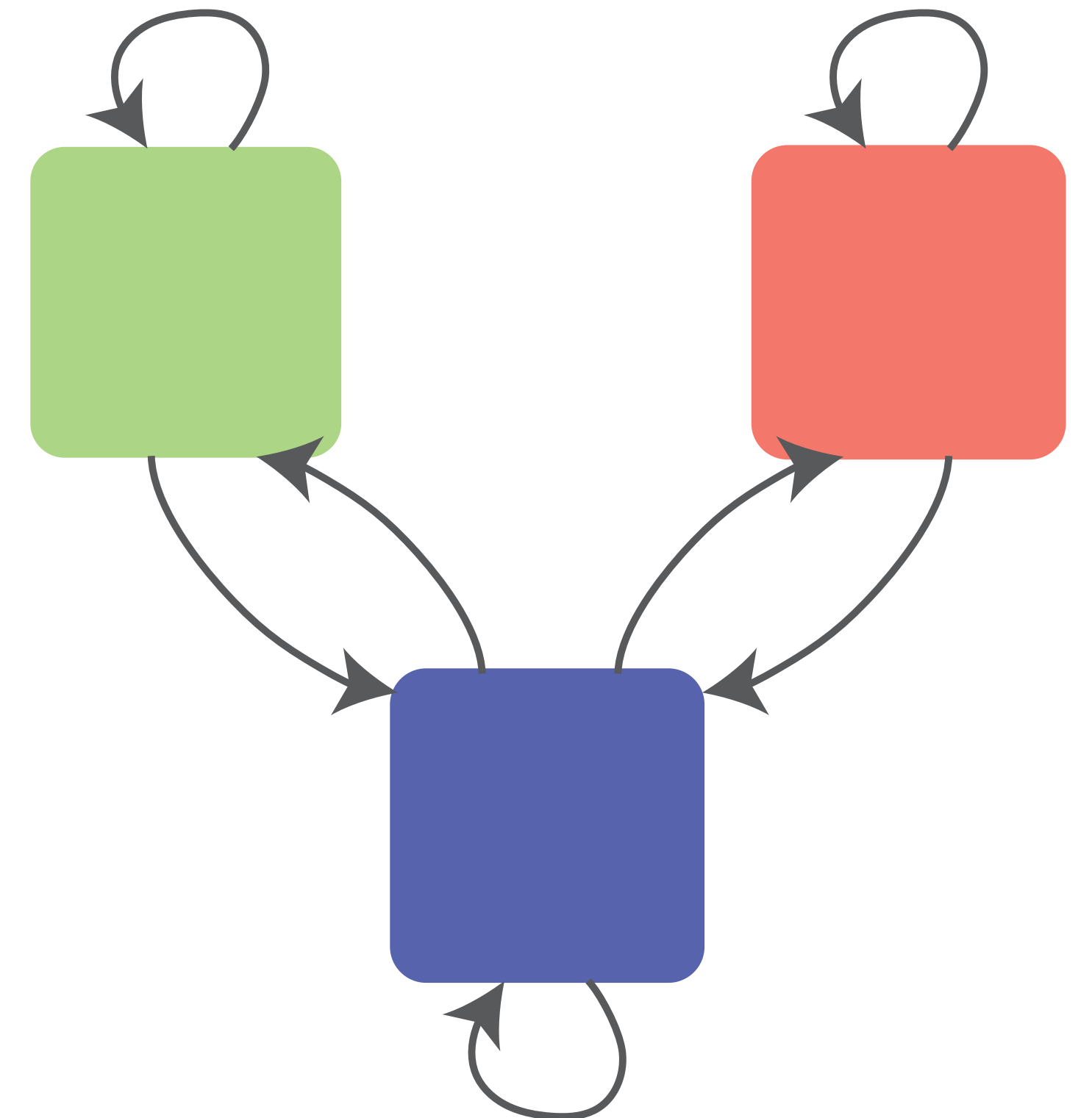
- Process (e.g. TKF91/92) described with an HMM
- Differentiating along a branch explodes the HMM
- Finite branch length \rightarrow infinite states



MODELLING INDELS

USING HMMS

- Process (e.g. TKF91/92) described with an HMM
- Differentiating along a branch explodes the HMM
 - Finite branch length \rightarrow infinite states
- Can be simplified!
 - For local likelihood;
 - Global likelihood complexity is still exponential.



MODELLING INDELS...

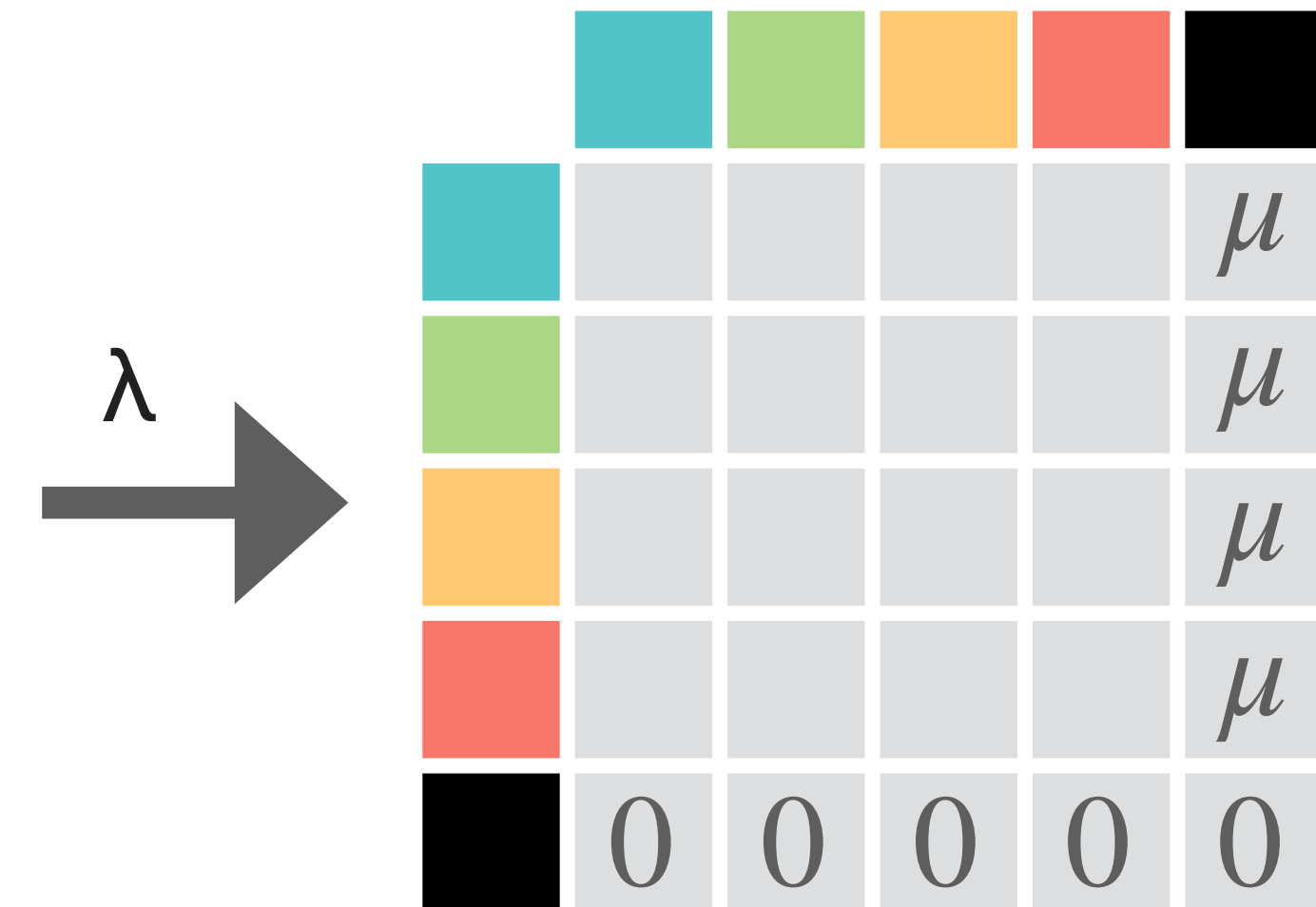
... is very hard

MODELLING INDELS

POISSON INDEL PROCESS (PIP)

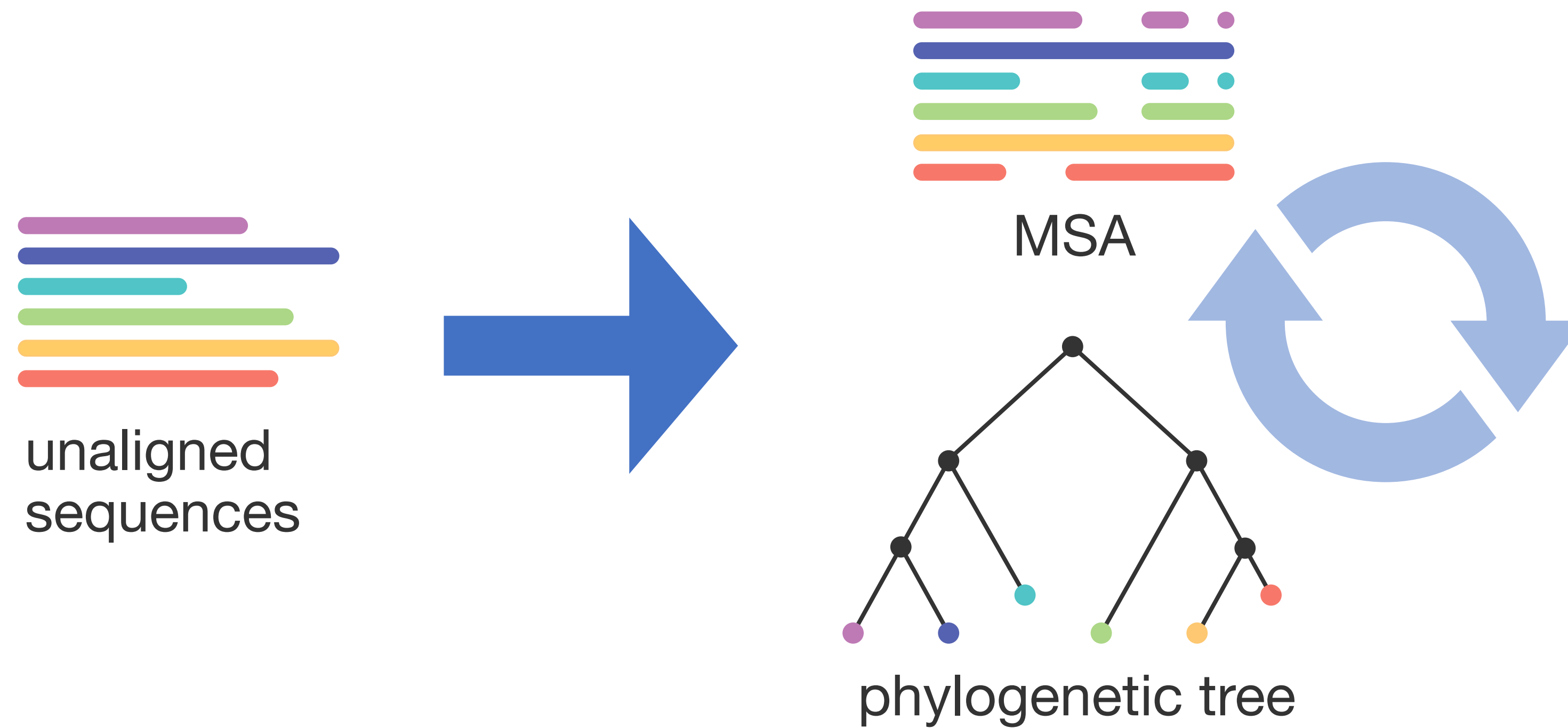
PIP = Poisson insertion process on a tree
+ Substitution/deletion Markov process

- Properties:
 - Single character insertions/deletions
 - Time-reversible
 - Linear time likelihood



JOINT INFERENCE: THE DREAM

EXPECTATION 2.0

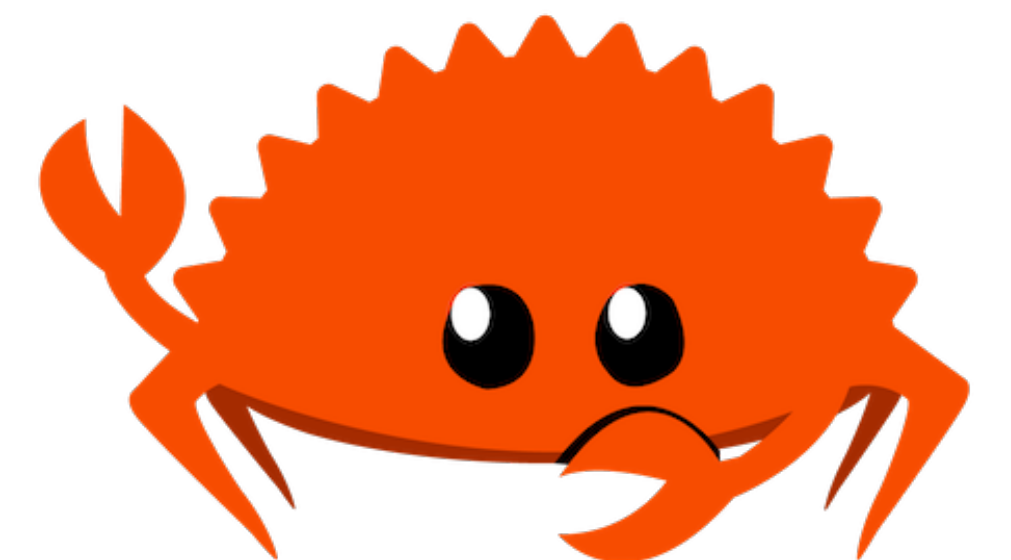


BAli---Phy

JOINT INFERENCE

REALITY

- JATI 1.0 necessary (but not sufficient) conditions:
 - Efficient
 - MSA adjustment in $O(NL^3)$
 - Tree search in $O(N^2L)$
 - Open-source
 - User-friendly
- Implementation:
 - Rust codebase (github.com/acg-team/JATI (private for now))
 - Release information (twitter.com/JulijaPecerska)



TAKEAWAYS

- Indel models are complicated;
- Know your data!!!
- Know your models and priors;
- Know your software and its assumptions.

ACKNOWLEDGEMENTS

- ACGTeam @ ZHAW



- Prof Maria Anisimova



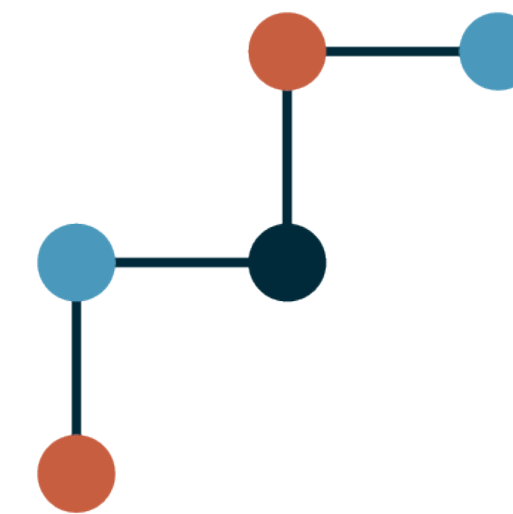
- Dr Manuel Gil



- Clara Iglhaut



- Gholamhossein Jowkar



**Swiss National
Science Foundation**

