# My research progress

# A quick remainder of what I am doing
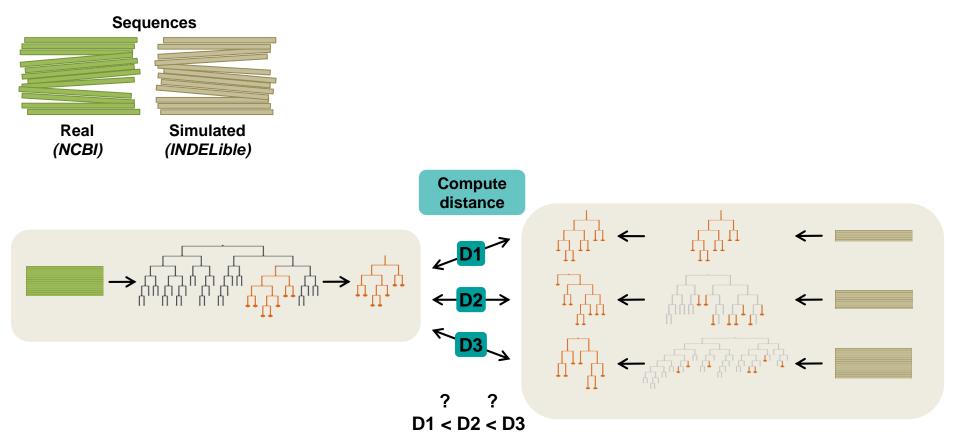
**MSA**
multiple sequence alignment

# Which sequences are simulated?



ITS gene

# GUIDANCE as a measure of aligning difficulty



**Column score: 2/3**

# GUIDANCE as a measure of aligning difficulty



**Residue score: 3/3**

# GUIDANCE as a measure of aligning difficulty

**Real sequences**



**Simulated sequences**



| Gene | Residues score | Column score |
|------|----------------|--------------|
| 18S  | 0.99955        | 0.854151     |
| 26S  | 0.999351       | 0.742893     |
| atpB | 0.999953       | 0.947041     |
| matK | 0.999471       | 0.959492     |
| ITS  | 0.893573       | 0.050596     |

| Gene | Residues score | Column score |
|------|----------------|--------------|
| 18S  | 0.987769       | 0.001852     |
| 26S  | 0.862116       | 0.000057     |
| atpB | 0.995953       | 0.350238     |
| matK | 0.997872       | 0.075541     |
| ITS  | 0.492671       | 0.00093      |

# INDELible – a deeper insight

## Markov chain

### State of the chain = whole sequence

State 1                    State 2

GGTT**A**ATCCCAGT  →  GGTT**G**ATCCCAGT

**Simulating Markov chains using Gillespie's algorithm**

$$\exp(\lambda)$$

EVENT

*Mol. Biol. Evol. 2009*

# INDELible – a deeper insight



**Given:**

GGTTGATCCTGCCAGTAGTCATCT

sequence

branch length (t)

**Denote:**

$$\lambda = I + D + S$$

**I**nsertion rate
**D**eletion rate
**S**ubstitution rate

**Waiting time:**

s1

$s \sim \exp(\lambda)$, mean$=1/\lambda$

# INDELible – a deeper insight

draw an event with probabilities

$s1 > t$                                                         $s1 < t$

$$\frac{I}{\lambda'}, \qquad \frac{D}{\lambda'}, \qquad \frac{S}{\lambda}$$

no event

draw a location

indel                    substitution

$t = 1.2$



uniformly
(L+1 positions)

proportional to
the substitution
rate at each site

$\lambda = I + D + S$

insert, delete or
substitute

update sequence
length L,
time = t-s1

# INDELible – a deeper insight

draw an event with
probabilities

$\dfrac{I}{\lambda'}$ $\qquad$ $\dfrac{D}{\lambda'}$ $\qquad$ $\dfrac{S}{\lambda}$

**s1 < t**

draw a location

indel $\qquad$ substitution

uniformly $\qquad$ proportional to
the substitution
rate at each site

repeated until
s1 + s2 + ... > t

insert, delete or
substitute $\qquad$ update sequence
length L,
time = t-s1

# INDELible – a deeper insight

## Insertion-deletion model

**ir**    **indel rate**        $ir \in [0,1]$

**a**    **indel size**    $a > 1$



**Zeta distribution**

# INDELible – a deeper insight

## Insertion-deletion model

**ir**     **indel rate**                 $ir \in [0,1]$

**a**     **indel size**       $a > 1$



**rl**     **root length**



GGTTGATCCTGCCAGT

# INDELible – a deeper insight

## Insertion-deletion model

**ir**      **indel rate**

**a**      **indel size**

**rl**      **root length**

INDELible – a deeper insight

Gillespie algorithm

True alignment

E: GGTGGATCCTGCT-AGT
D: GGT--ATCCTGCCA-GT
C: AGGTTGATC-GGCCATT
B: GGTTGATCCT---AGT
A: AATTGATTCTGCCATTT

GGTTGATCCTGCCAGT

# INDELible – estimating parameters

average seqs length

alignment length

average gaps number

**rl** =average seqs length

**a** =1.00005, 1.001,1.15…

**ir** =0.0001, 0.00002, 0.0002…

**INDELible**

alignment length

average gaps number

compute relative error

Choose the parameters a, ir which gave the lowest error

# Getting more realistic simulations

Real sequences    Simulated sequences

INDELible

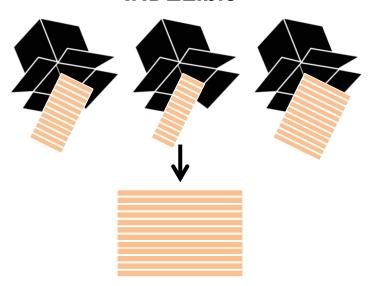ir   a   rl

constant

distributed

block optimized

# Getting more realistic simulations

## INDELible



|  | column score (GUIDANCE) | | |
|---|---|---|---|
|  | **constant** | **distributed** | **block optimized** |
|  | 0.97901696 | 0.97916526 | 0.97672856 |
|  | 0.91552168 | 0.94023076 | 0.91669352 |
|  | 0.9883448 | 0.98878618 | 0.98792382 |
|  | 0.99055344 | 0.99287678 | 0.98774598 |
|  | 0.94787124 | 0.9599763 | 0.9362862 |
|  | 0.76490788 | 0.79192566 | 0.74172564 |
|  | 0.51005486 | 0.52628762 | 0.60550184 |

# Changing route?