



THESIS DEFENSE



Bayesian Assembly of Reads from High Throughput Sequencing

Jonathan Laserson

advised by

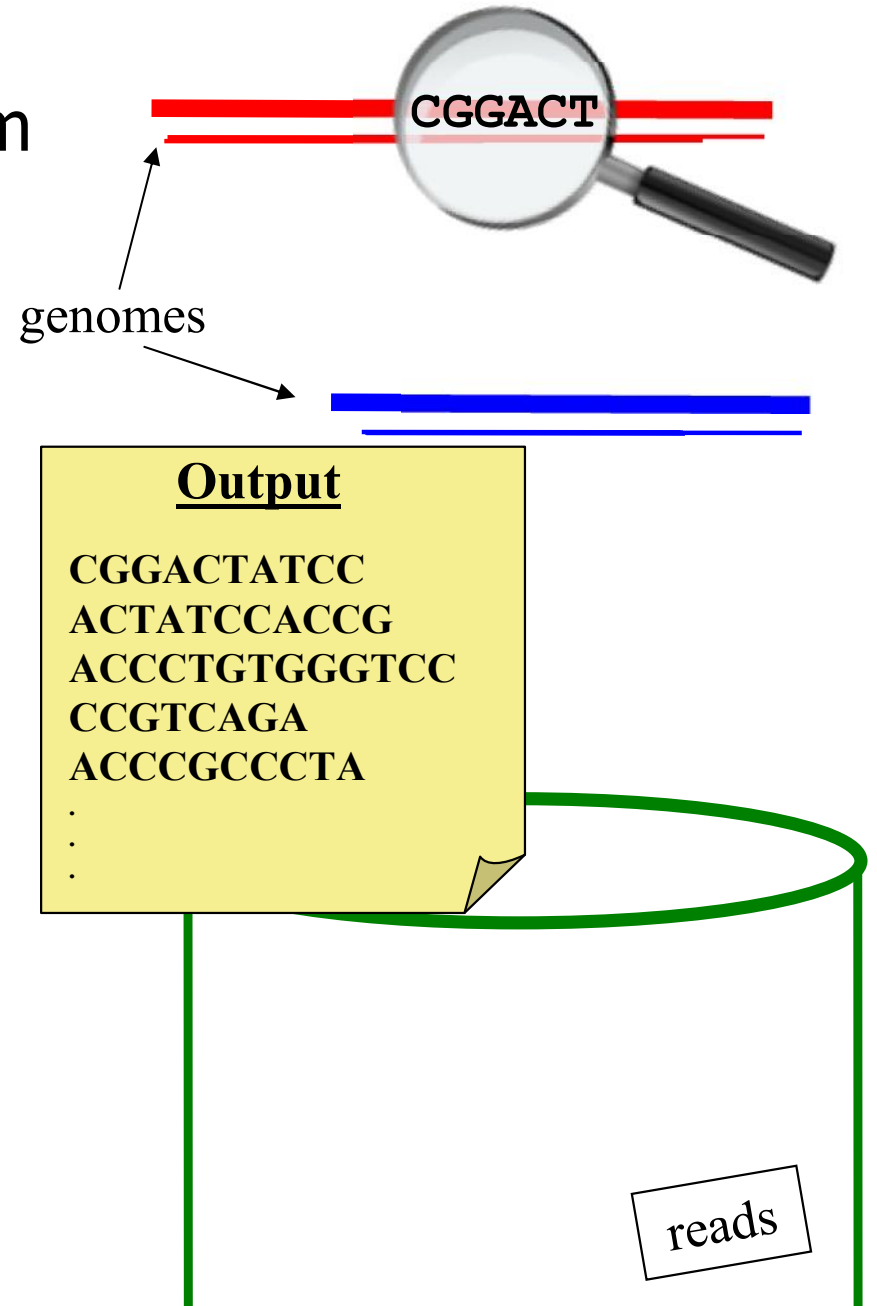
Daphne Koller

October 20th, 2011



High-Throughput Sequencing

- Millions of noisy short reads from a sample
- A **snapshot** of the genomic material in the sample
- Roche 454 GS-FLX Titanium
 - Avg read length 380b
 - reads in a run 500k
 - Base error rate 0.74% (84% indels)
 - Cost per base 0.05\$/kb

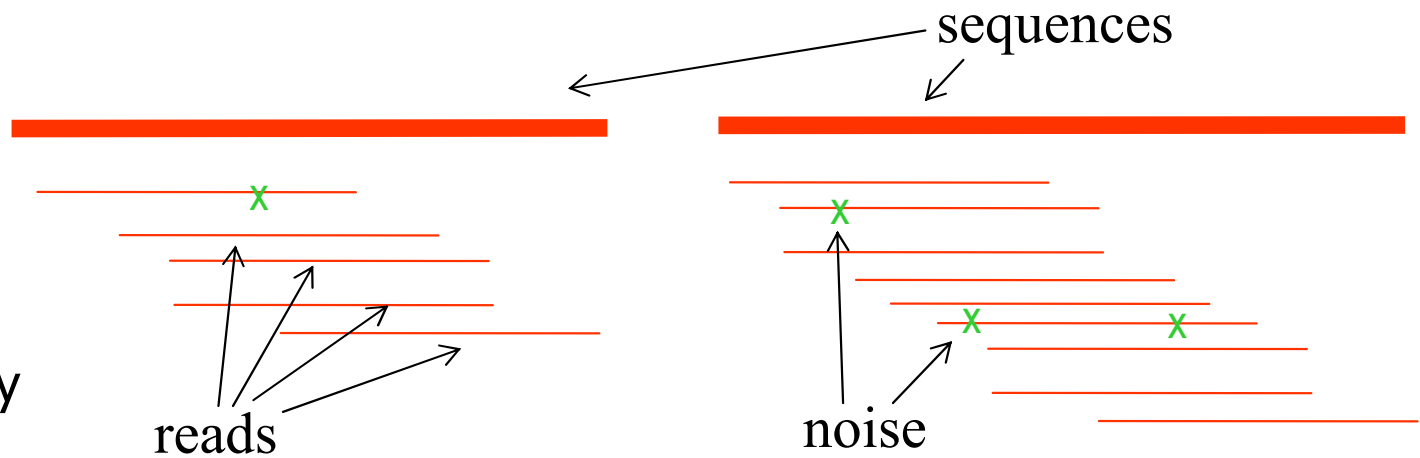




Population Sequencing

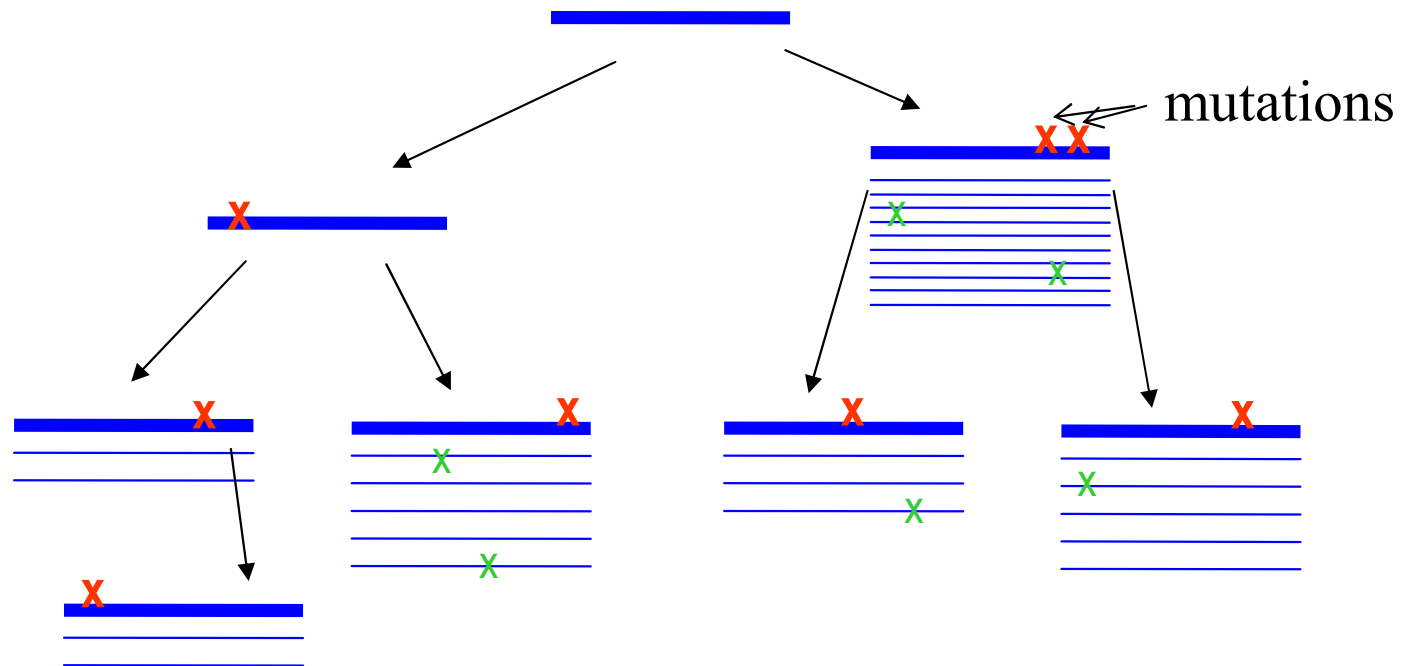
Breadth

- Metagenomics
- *de novo* assembly



Depth

- Phylogenetic trees
- Immune-seq



CTGGGAAAAG AGCCCAGCCCC CGTAGATCA CTGGAAAAG
GCTGACGTAG AAAGCCCAG ACTGACCTG
CCCTCATTTC ACTGCCCTGG GACGTAGATCA
GCCCTGGTGA CGGCTGACGT AGATCATAGA
AGCCCAGCCC CATAGAAACGAT

reads

Input

GENOVO

Output

ACTGCCCTGGGAAAAGCCCAGCCCCTCATTTC

ACTGCCCTGG

ACTGCCCTG

GCCCTGGGA

CTGGGAAAAG

CTGGGAAAAG

AAAGCCCAG

AGCCCAGCCC

AGCCCAGCCCC

CCCTCATTTC

contigs

CGGCTGACGTAGATCATAGAAACGAT

CGGCTGACGT

GCTGACGTAG

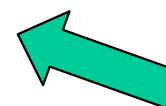
GACGTAGATCA

CGTAGATCA

AGATCATAGA

CATAGAAACGAT

reads



an assembly



Previous Work

- Velvet, Euler:
 - Assume most reads have no error
 - Assume a single genome
 - Trash low-abundance reads
 - Fix reads with one error a-priori
- Newbler: No code / publication

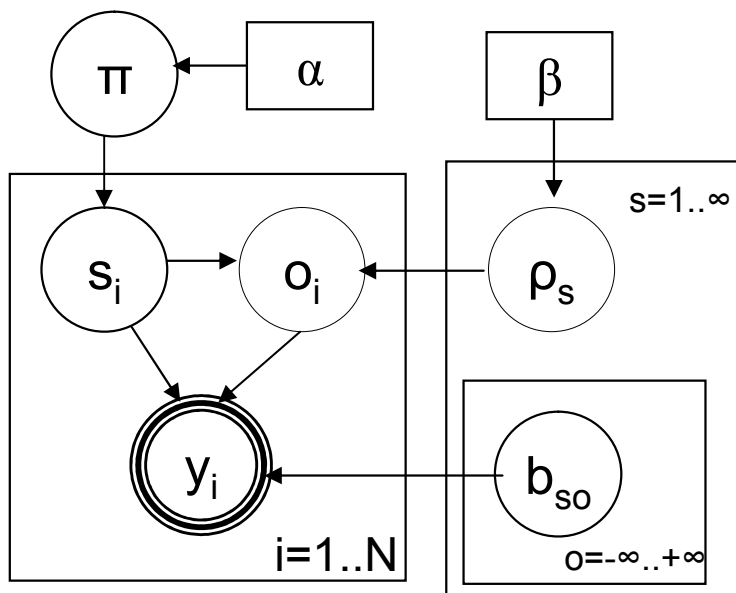


Our Approach: Probabilistic Model

- **Fully Bayesian** model for high-throughput sequencing
- **Jointly** performing 3 tasks:
 - de-noising the reads
 - *de novo* sequence assembly
 - multiple alignment
- **We do not throw away reads**
- Can handle **unbalanced populations** of genomes



Our Approach: Probabilistic Model



$b_{so} \sim U[A, C, G, T]$
 $\mathbf{s} \sim \text{CRP}(\alpha, N)$
 $\rho_s \sim \text{Beta}(1, 1 + \beta)$
 $o_i \sim \mathcal{G}(\rho_{s_i})$
 $y_i \sim \text{read_model}(\mathbf{b}, s_i, o_i)$

$$\log P(\mathbf{y}, \mathbf{b}, \mathbf{s}, \mathbf{o} | \boldsymbol{\rho}) \approx \underbrace{\# \text{errors}}_{\text{min read errors}} \cdot \log(\epsilon) - \underbrace{\# \text{bases}}_{\text{min bases}} \cdot \log(4) + \underbrace{\# \text{seq}}_{\text{min seq}} \cdot \log(\alpha)$$

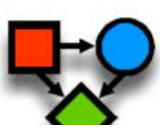
we want to find

$$\underset{\mathbf{b}, \mathbf{s}, \mathbf{o}}{\operatorname{argmax}} P(\mathbf{b}, \mathbf{s}, \mathbf{o} | \mathbf{y})$$

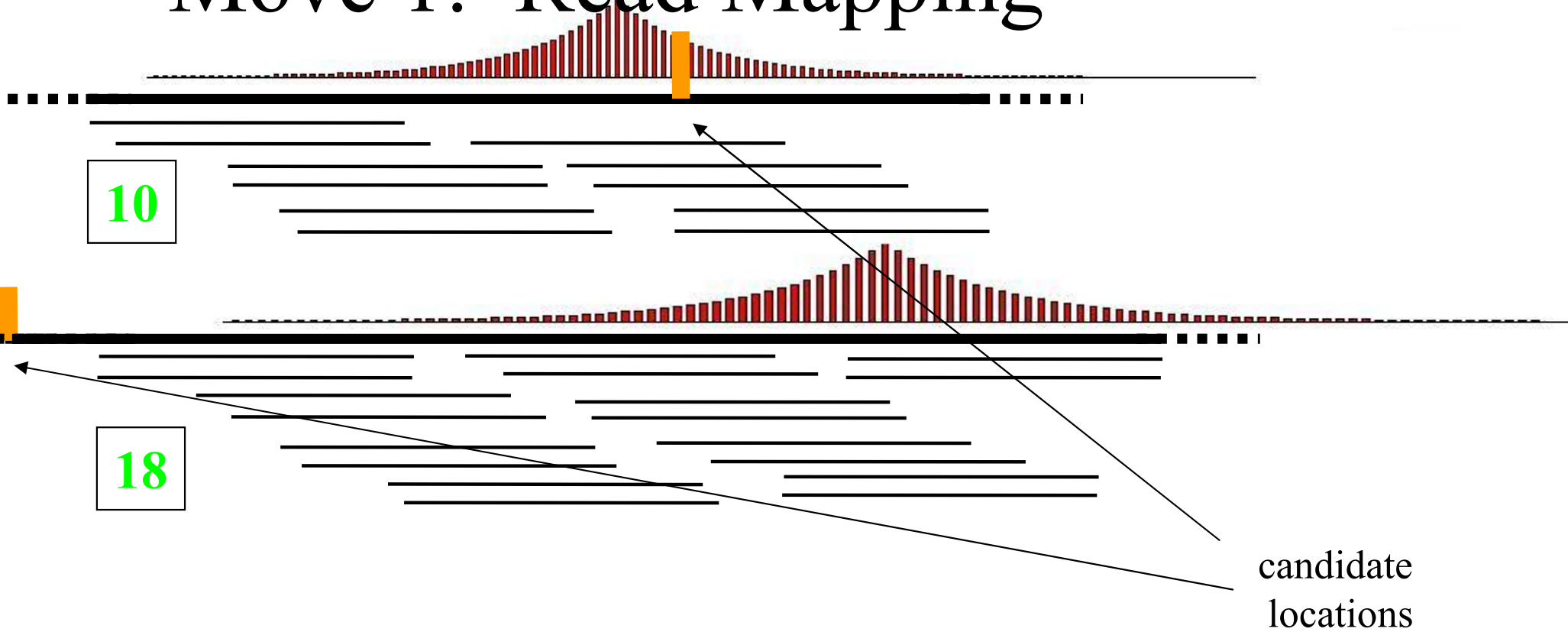


Inference

- MCMC
- Making moves in a probabilistic space
- Stochastic process
 - escape local optimum
 - multiple hypotheses



Move 1: Read Mapping



$$P(s_i=s, o_i=o \mid \mathbf{y}, \mathbf{b}, \boldsymbol{\rho}, \mathbf{s}_{-i}, \mathbf{o}_{-i}) \propto \text{Likelihood Term} \text{ Stickiness Term} \text{ Geometric Term}$$

A READ samples a *sequence* and an *offset* (inside the seq)

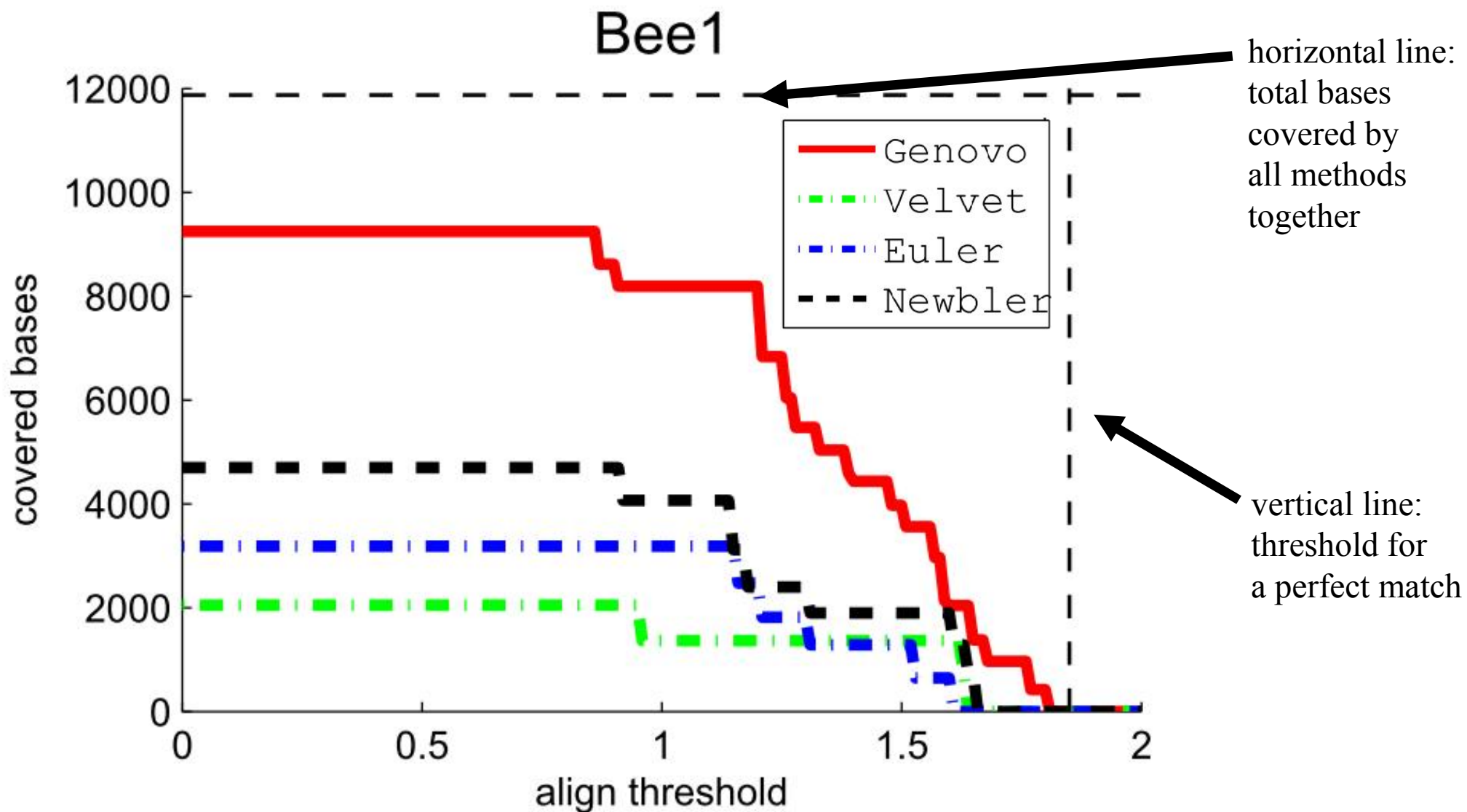
s_i

o_i



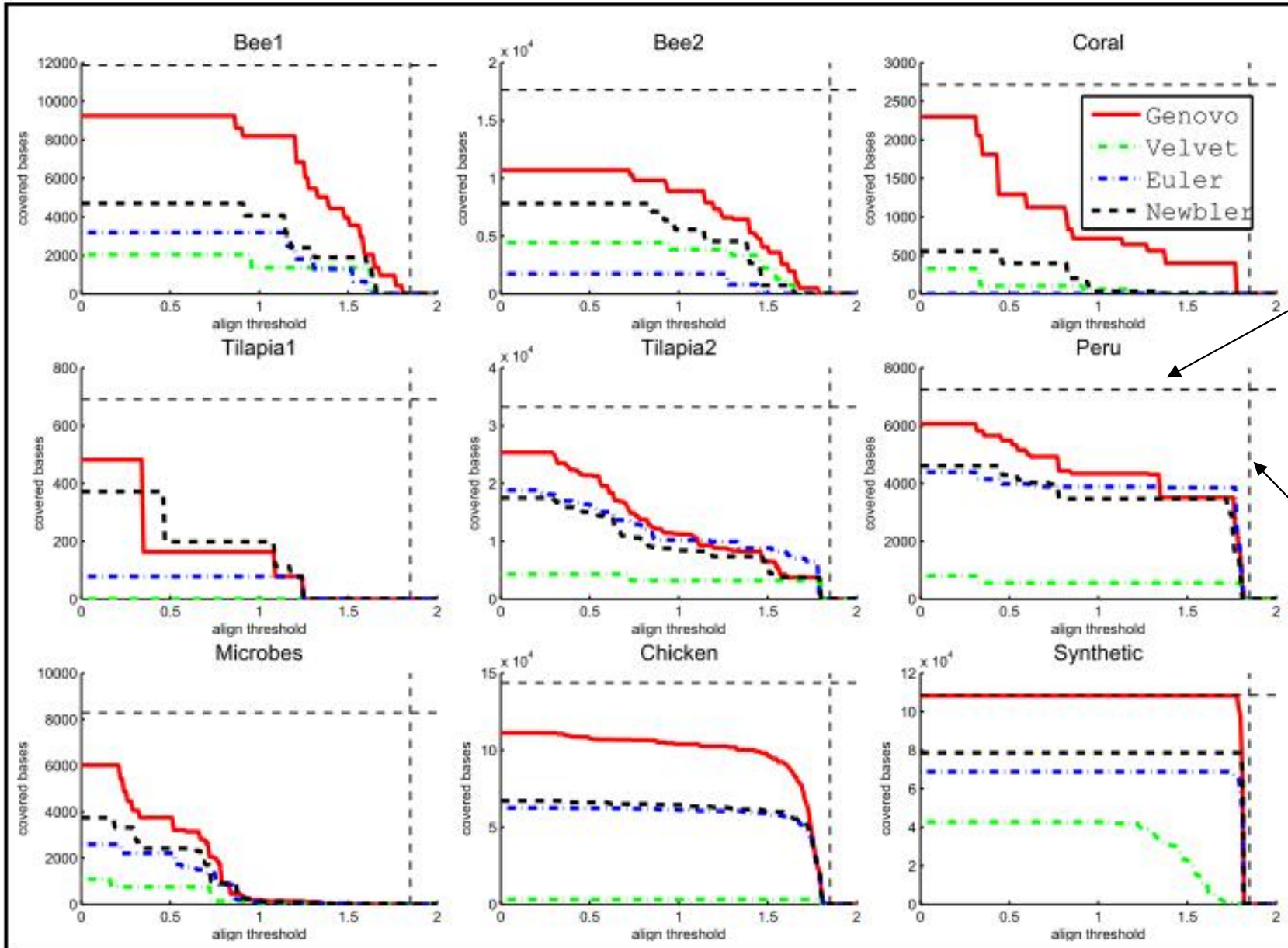
Comparison:

How many bases did you BLAST into?





Comparison: How many bases did you BLAST into?



horizontal line:
total bases
covered by
all methods
together

vertical line:
threshold for
a perfect match



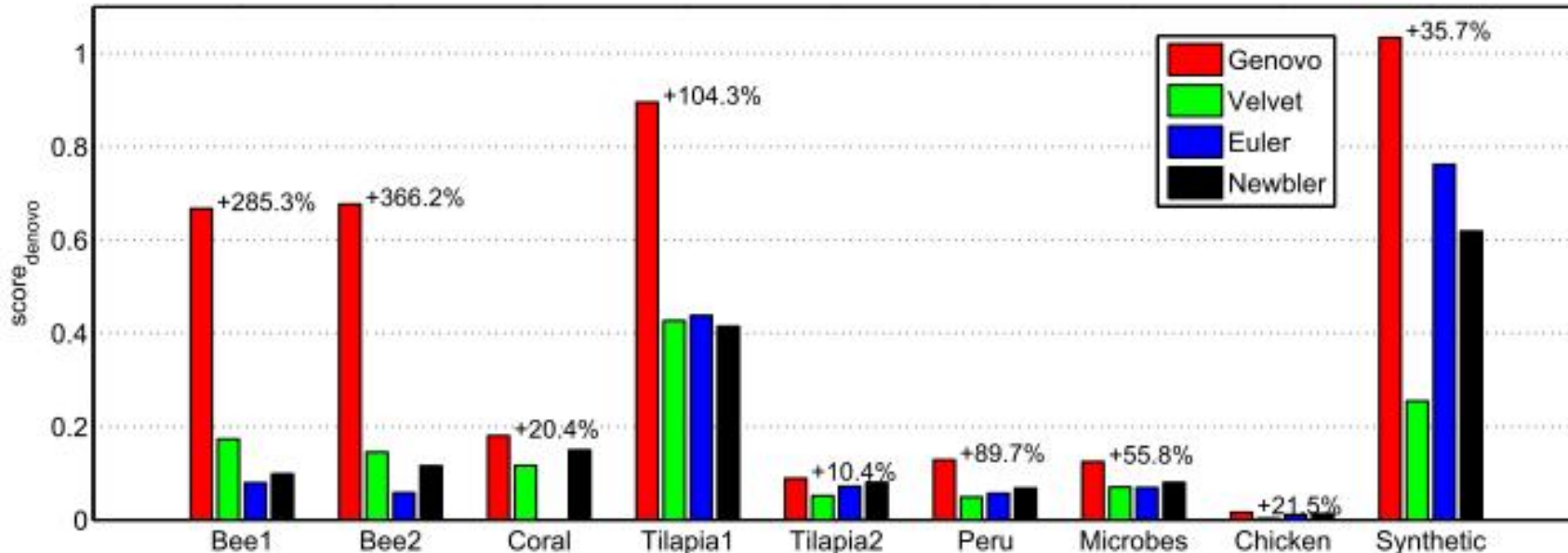
A Score for an Assembly

min read errors

min bases

α controls minimal overlap
between adjacent reads

$$\text{SCORE}_{denovo} = \underbrace{\#errors \cdot \log(\epsilon)}_{\text{min read errors}} - \underbrace{\#bases \cdot \log(4)}_{\text{min bases}} + \underbrace{\#seq \cdot \log(\alpha)}_{\alpha \text{ controls minimal overlap between adjacent reads}}$$





Contributions

- Fully Bayesian model for high-throughput sequencing
- An algorithm jointly performing 3 tasks:
 - de-noising the reads
 - *de novo* sequence Assembly
 - multiple alignment
- A score for *de novo* assembly
- Code freely available

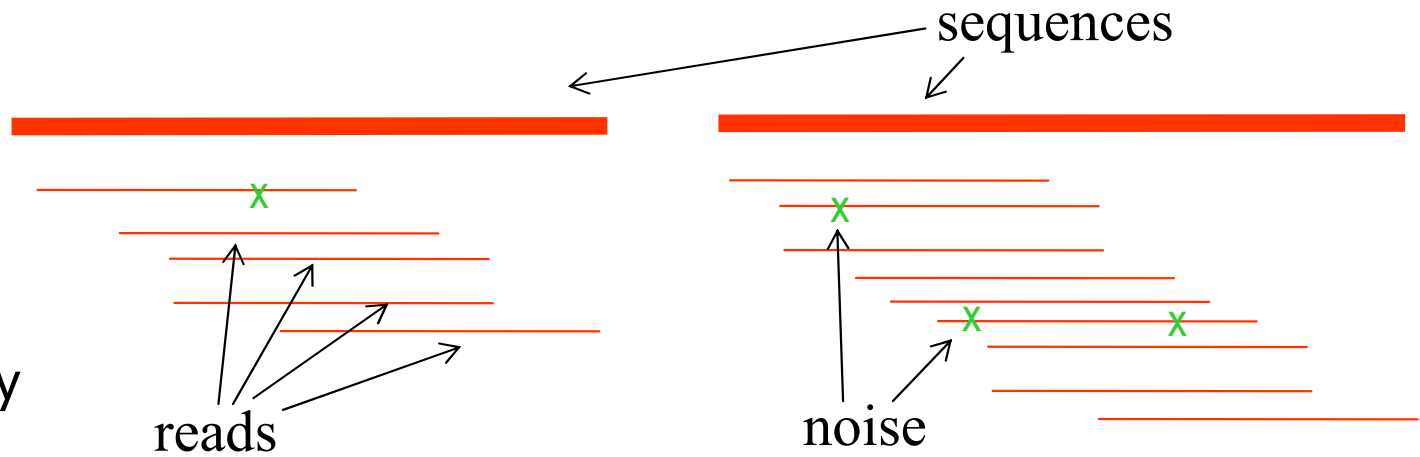
Laserson J, Jojic V, Koller D, **Genovo: *de novo* Assembly for Metagenomes**,
Journal of Computational Biology, 2011.
RECOMB 2010.



Population Sequencing

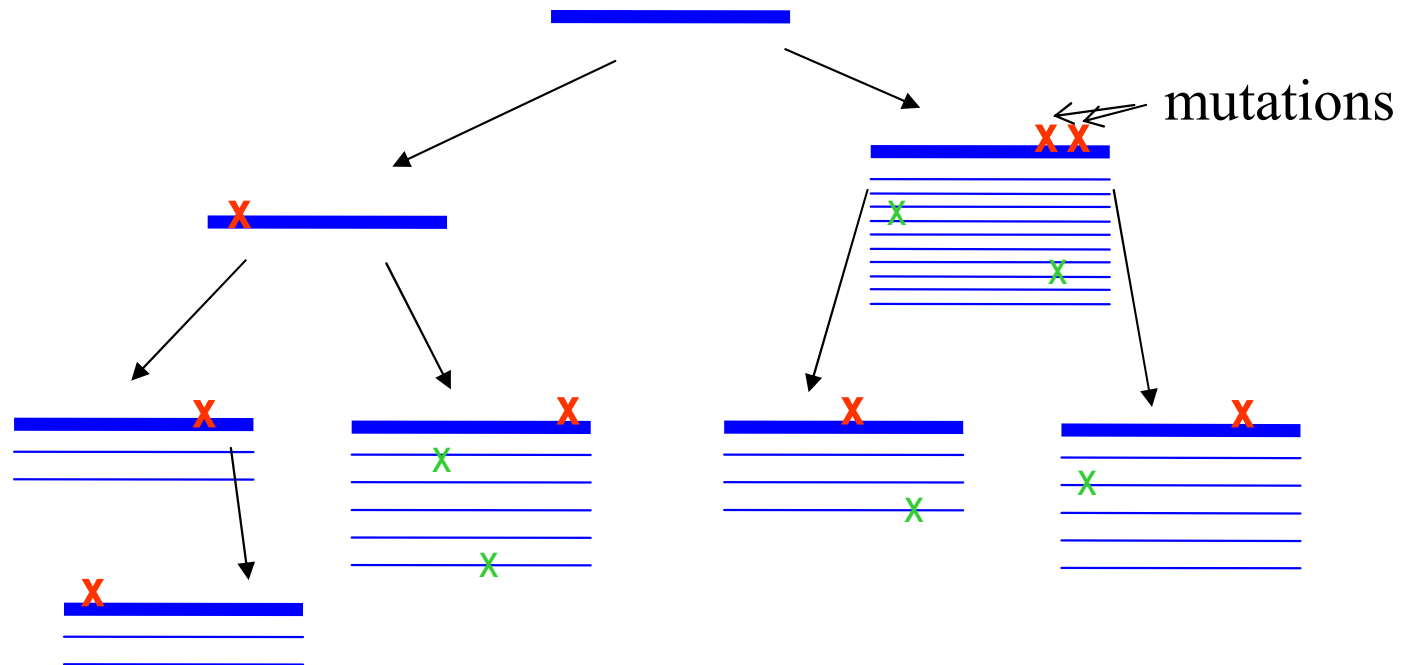
Breadth

- Metagenomics
- *de novo* assembly



Depth

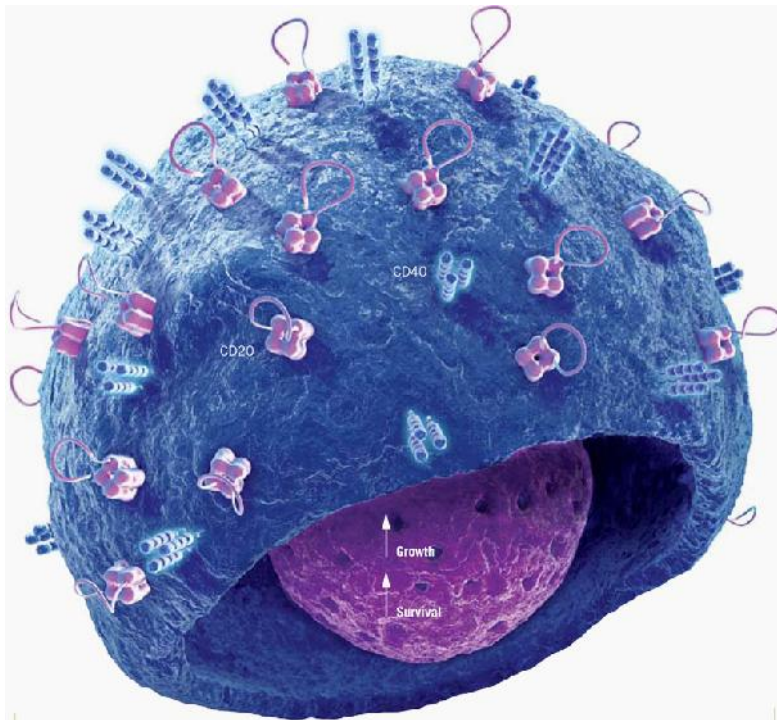
- Phylogenetic trees
- Immune-seq

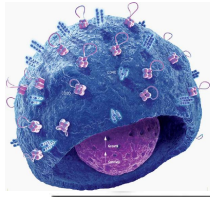




Immune System

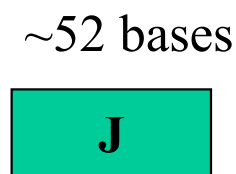
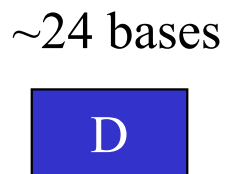
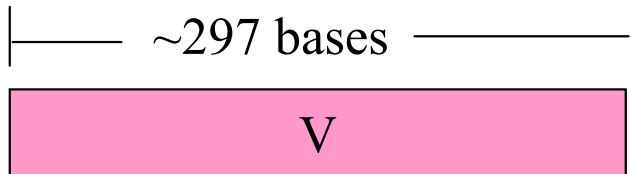
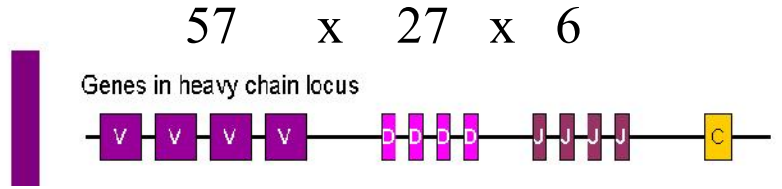
- 10^{10} B-Cells in an individual

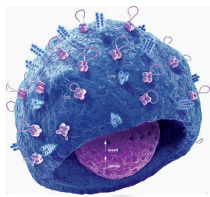




Immune System

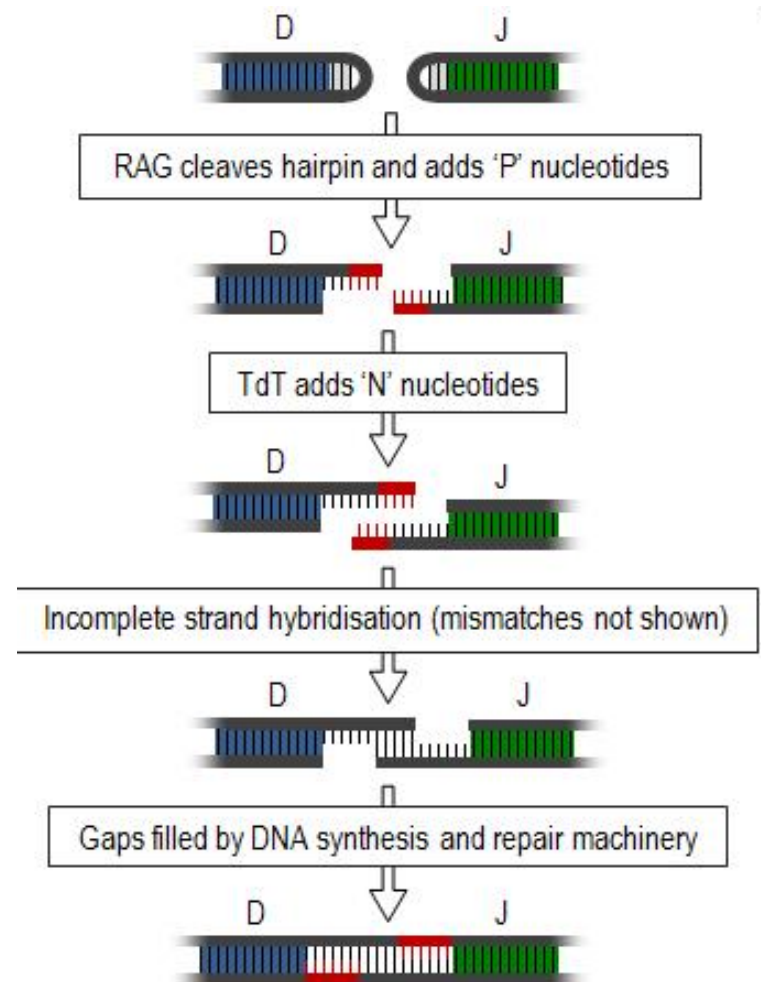
- 10^{10} B-Cells in an individual
- Phase 1: generate basic seq
 - choose V,D,J from catalog





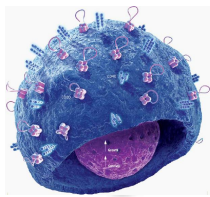
Immune System

- 10^{10} B-Cells in an individual
- Phase 1: generate basic seq
 - choose V,D,J
 - stitch V-D and D-J by:
 - eat up a few letters from ends
 - insert a few letters



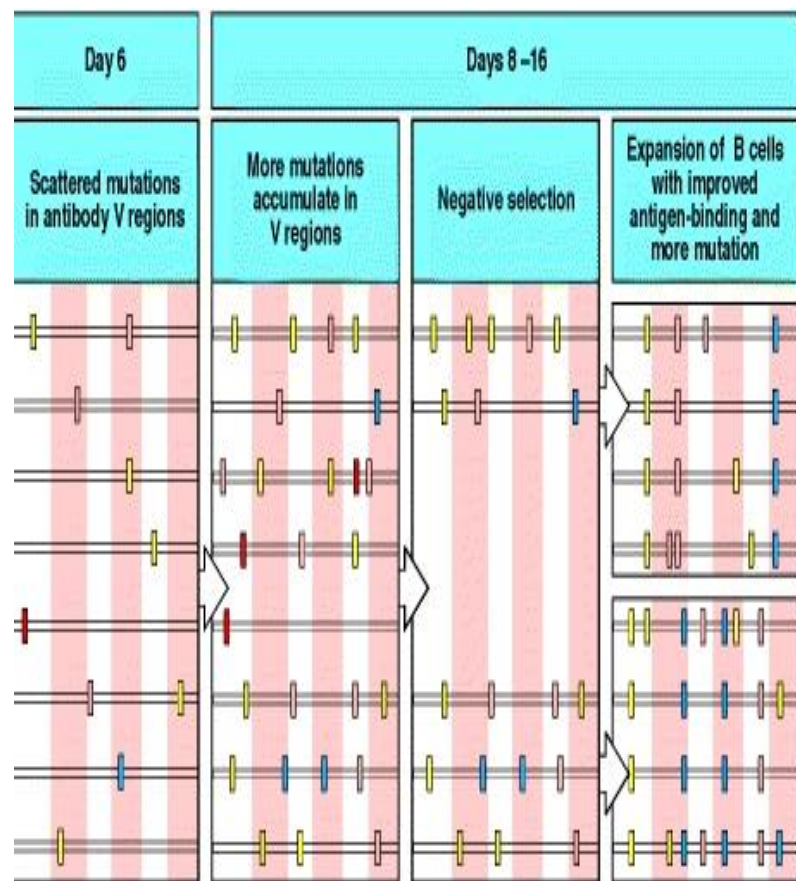
[Illustration: Wikipedia]





Immune System

- 10^{10} B-Cells in an individual
- Phase 1: generate basic seq
- Phase 2: hyper-mutate
 - If detecting an invader, **clone yourself with lots-of-mutations**
 - Even better detection? **keep cloning.**
 - **Remember the successful sequences** for future reference, and keep producing them.
 - This set of sequences is one **clone**.



[Illustration: Janeway, 2001]

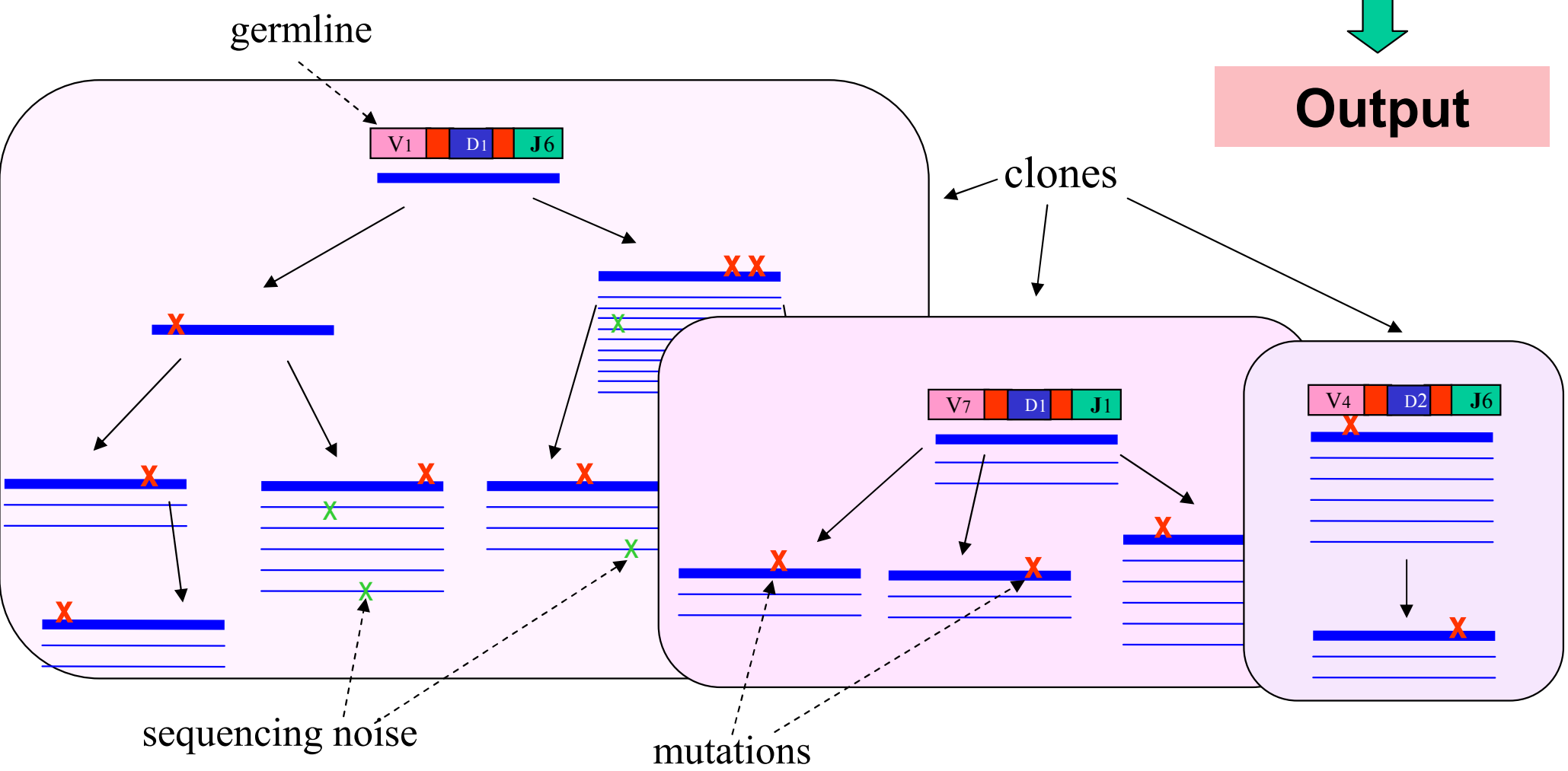




Input

igForest

Output





ReperTree

De Novo Annotation for B-Cells

- Phase 1: Divide reads into clones
 - 1M reads!
 - In a **healthy individual** almost every read is from a different clone.
- Phase 2: Construct a mutation tree for each clone.
 - How to tell a **mutation** from a **sequencing error**?



Overview

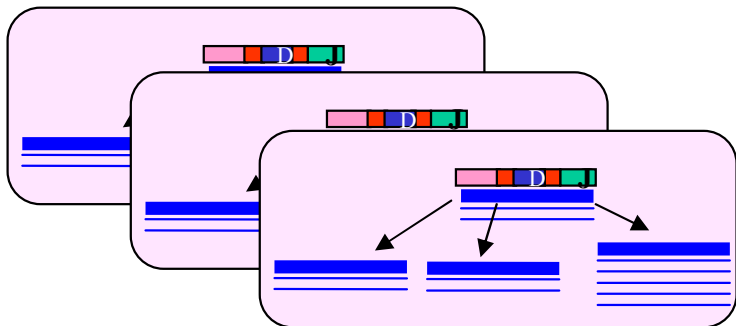
- 10 (healthy) Organ donors
- 4 samples from each individual
 - Blood, spleen, two lymph nodes
- 380k reads total



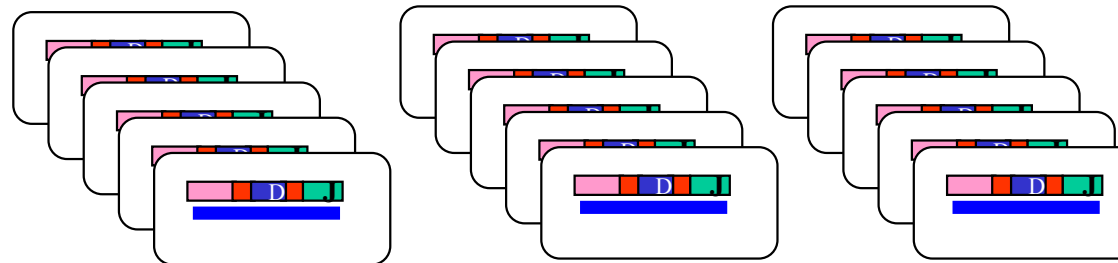
igForest



- 8,000 clones



- 165k other B-cells variants (not clones)





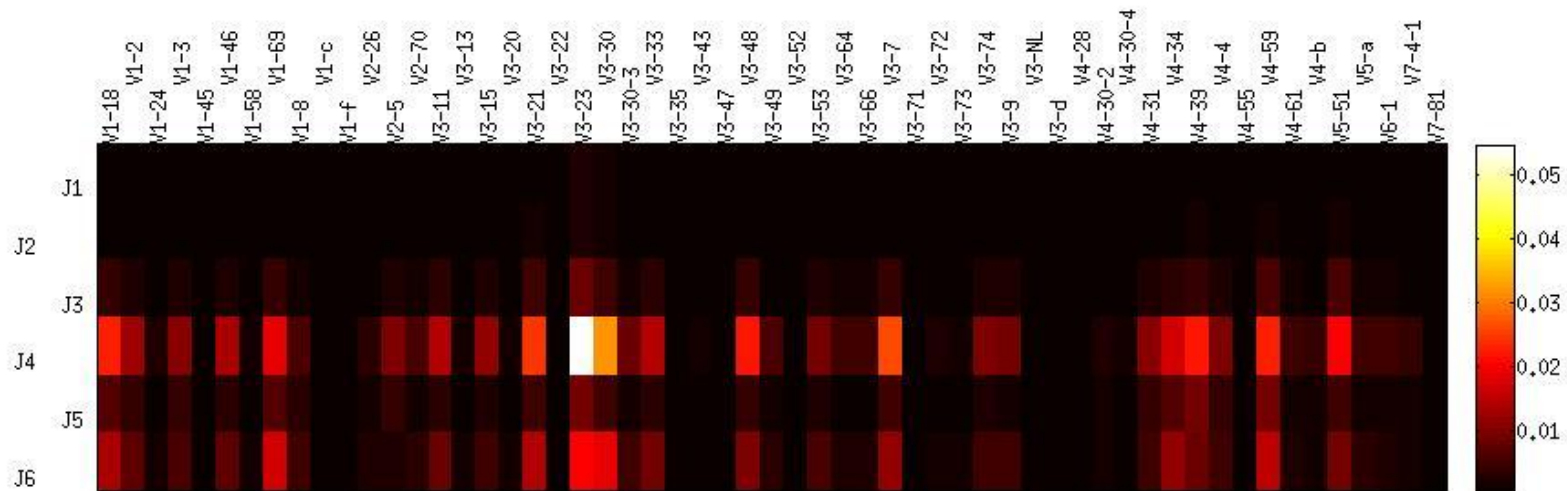
A Repertoire

Genes in heavy chain locus



Vs (57)

Js (6)



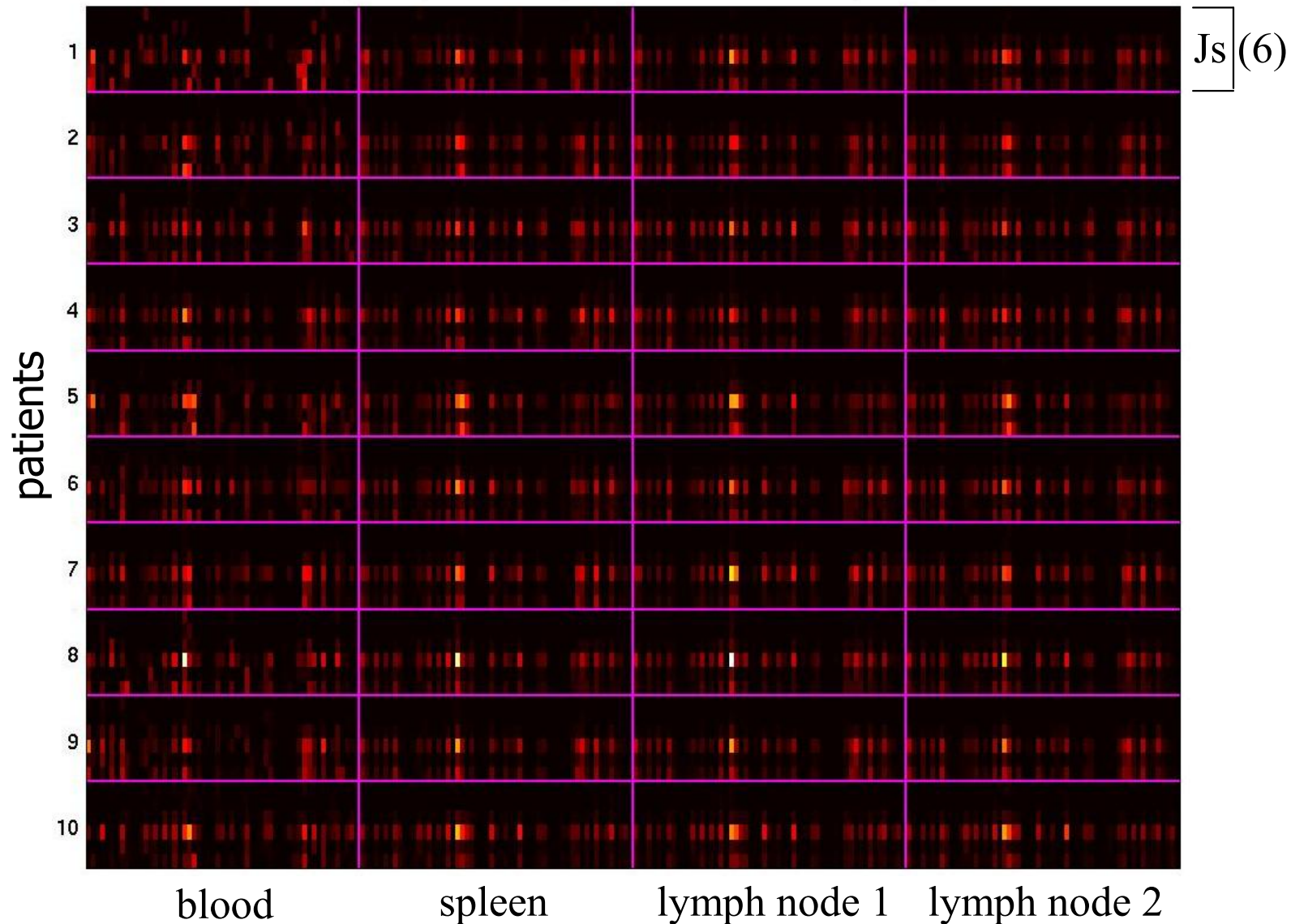
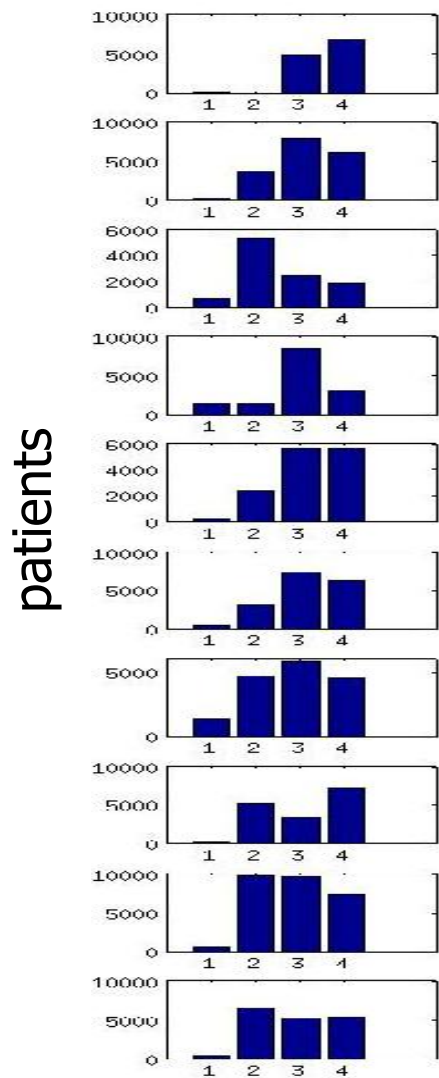
Variants (173k)

Genes in heavy chain locus



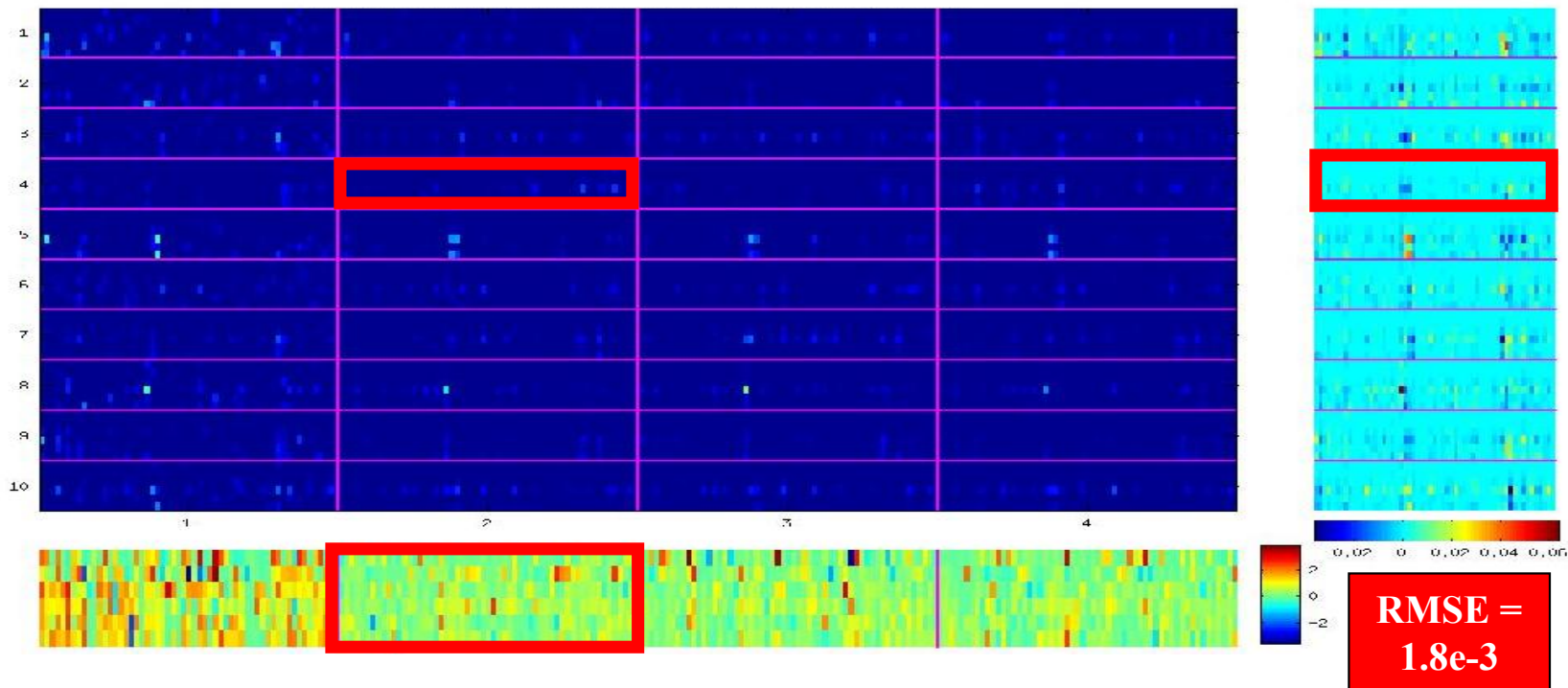
V_s (57)

collapsed repertoire for all patients/sources, normalized per sample





Repertoire Factorization



$X(v,j,donor,source) = a(v,j,source) \times b(v,j,donor)$

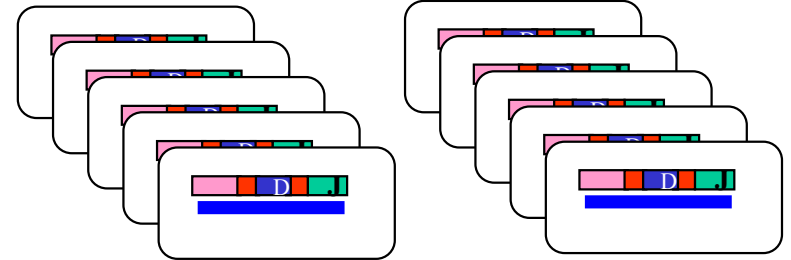
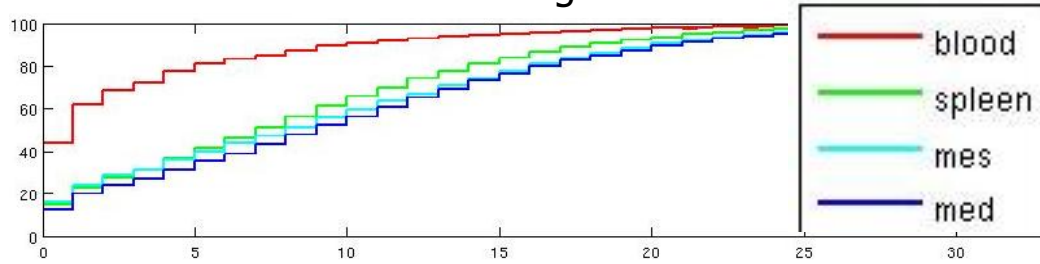
V,J	VJ	VS	JS	VD	JD	VJS	VJD	
X								3
	X							2.9
		X	X					2.9
				X	X			2.2
		X	X	X	X			2.1
	X	X	X	X	X			1.9
						X		2.8
							X	2
						X	X	1.8

*1e-3

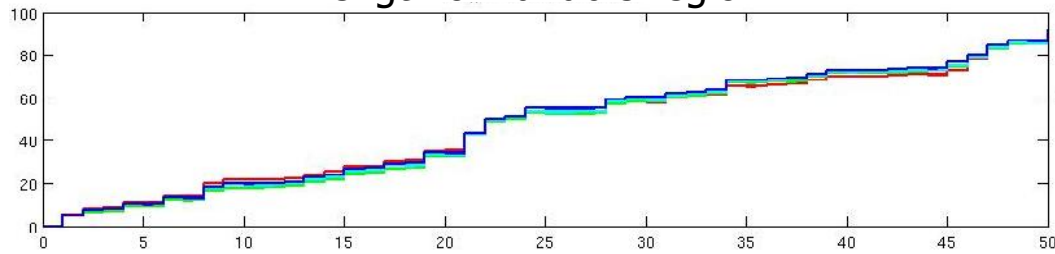


Properties Of Variants

distance to germline

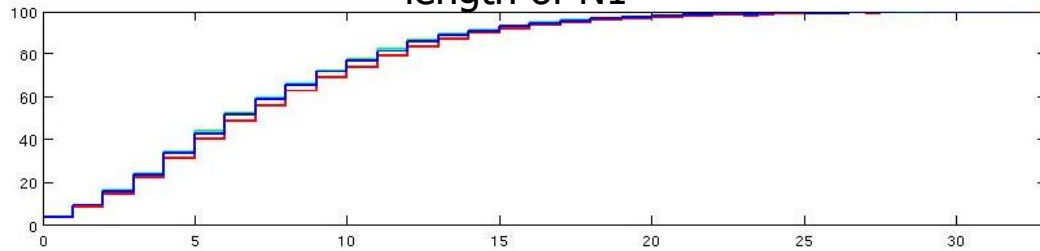


length of variable region

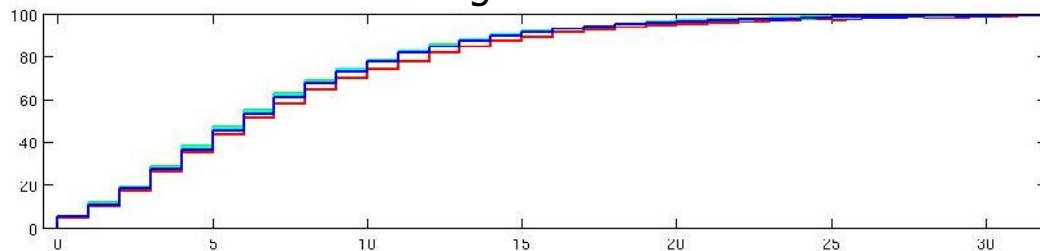


fraction of variants

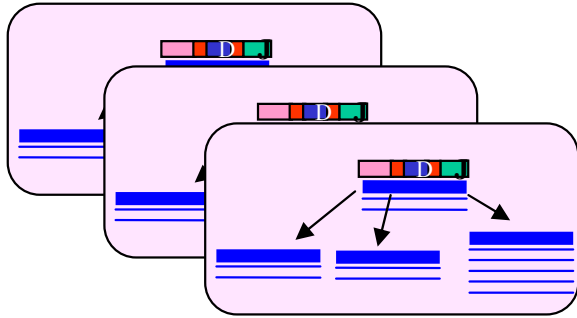
length of N1



length of N2

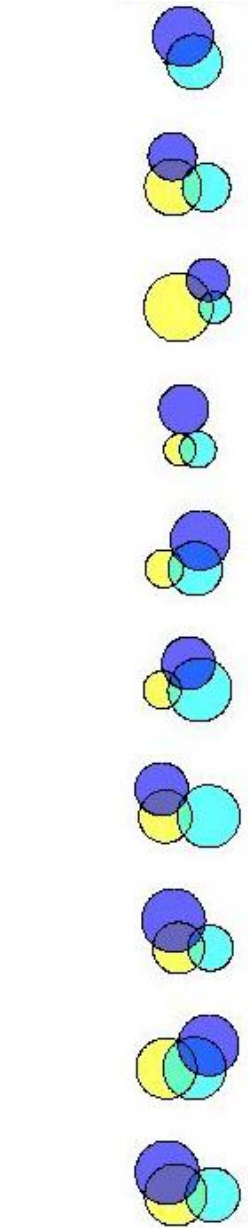
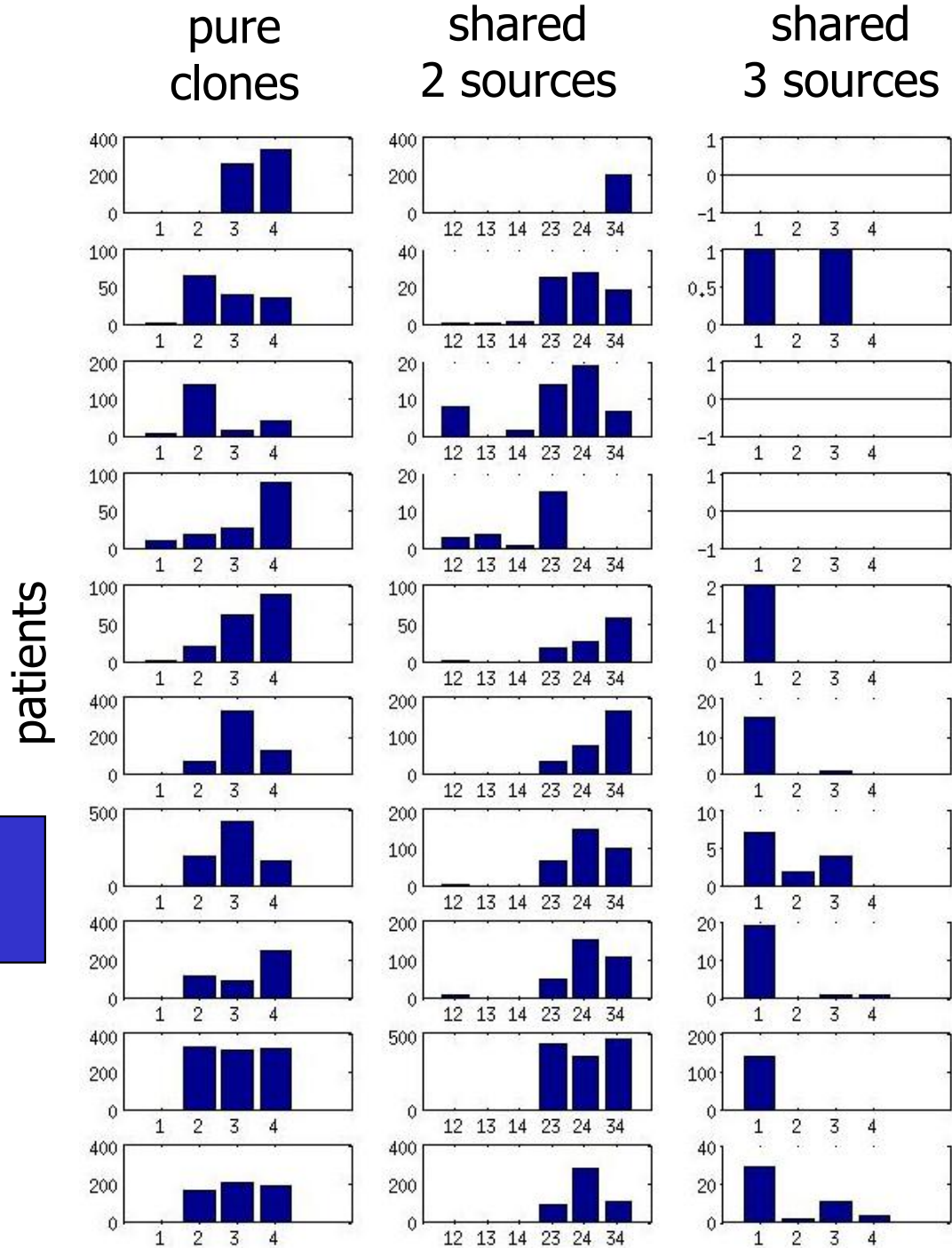
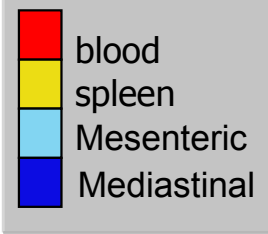


Clones (8,000)



Evidence of clone migration

How do the trees look like?

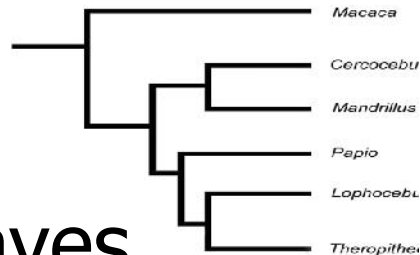




Current Phylogenetic Trees

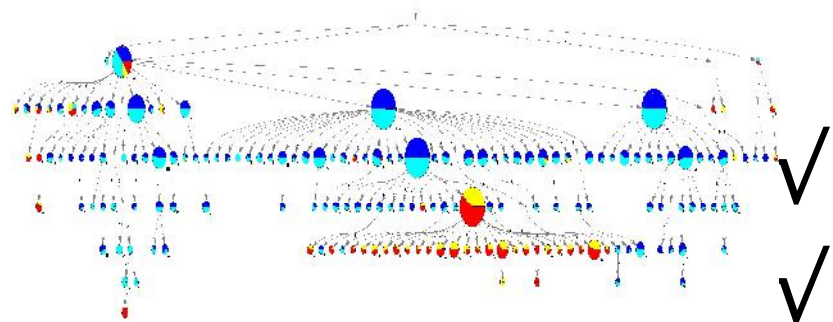
traditional phylo trees

- $X*100$ sequences
- Weeks to run
- No sequencing/PCR errors
- Single solution
- Binary tree
- Data at the leaves
- Continuous mutation model



igForest

- $X*10,000$ ✓
- Hours ✓
- Noise model ✓
- Multiple/statistics ✓
- Discrete ✓





Previous Work on B Cells

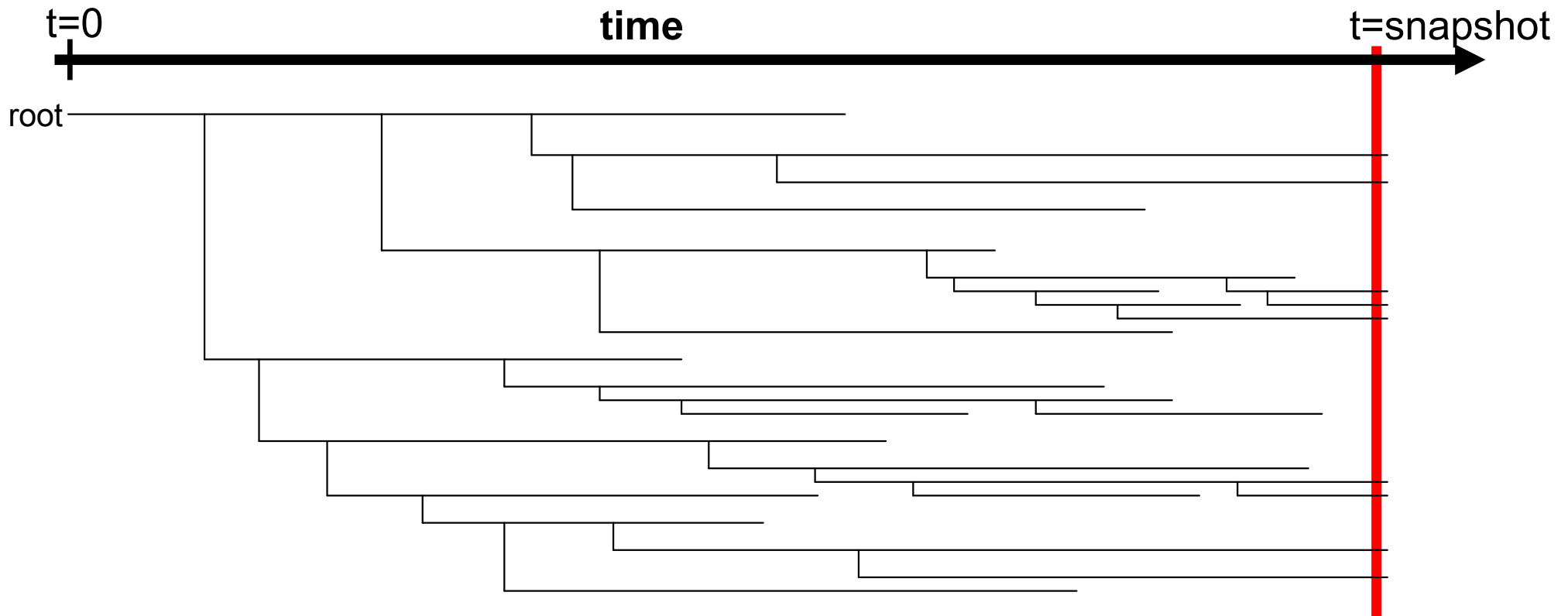
- **ihmmune-align** [Gaëta et al. 2007],
SoDA2 [Munshaw and Kepler 2010]
 - examine each read independently, no tree
- **igTree** [Barak et al. 2008]
 - heuristic search, no high-throughput support
- Our Approach: **igForest**
 - Joint **probabilistic model** for hypermutation *and* sequencing noise



Generate Tree

(Birth/death process)

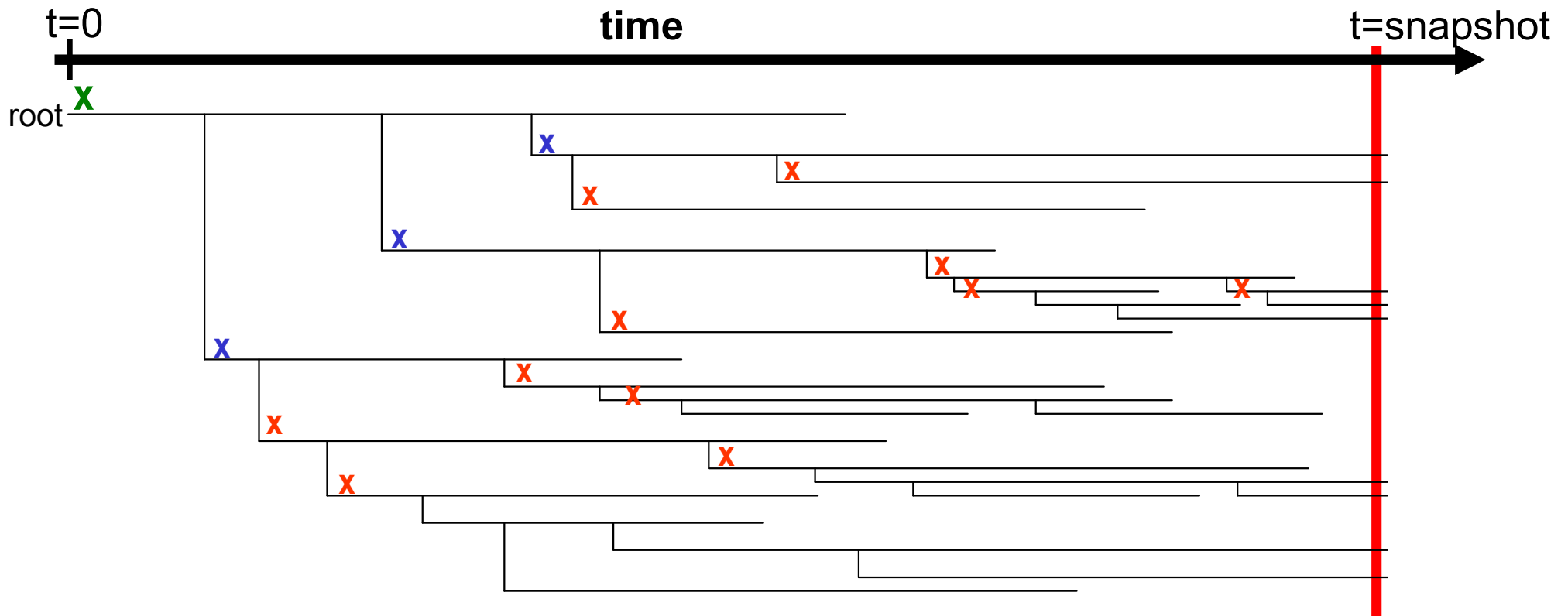
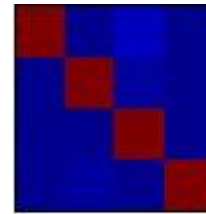
- **Start with one** live cell at $t=0$
- **Give birth** at rate λ
- **Die** at rate δ
- **Stop** when $t = \text{"snapshot time"}$.





Generate Sequences

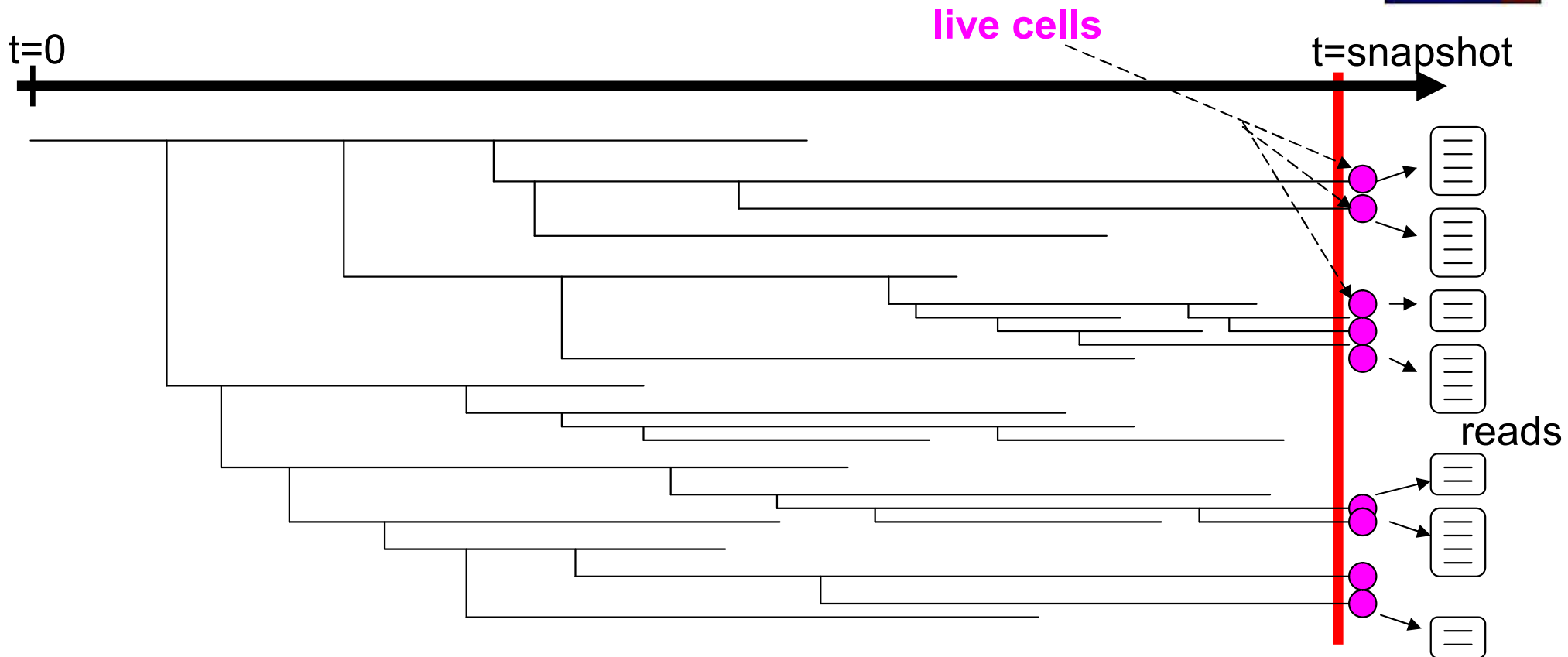
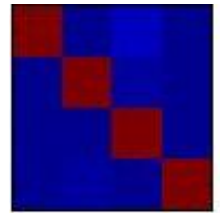
- **Start with** the inferred germline at the root
- **Generate each child sequence** from the sequence of its parent
 - Using a **Mutation Model**
- Continue so, **recursively**.





Generate Reads

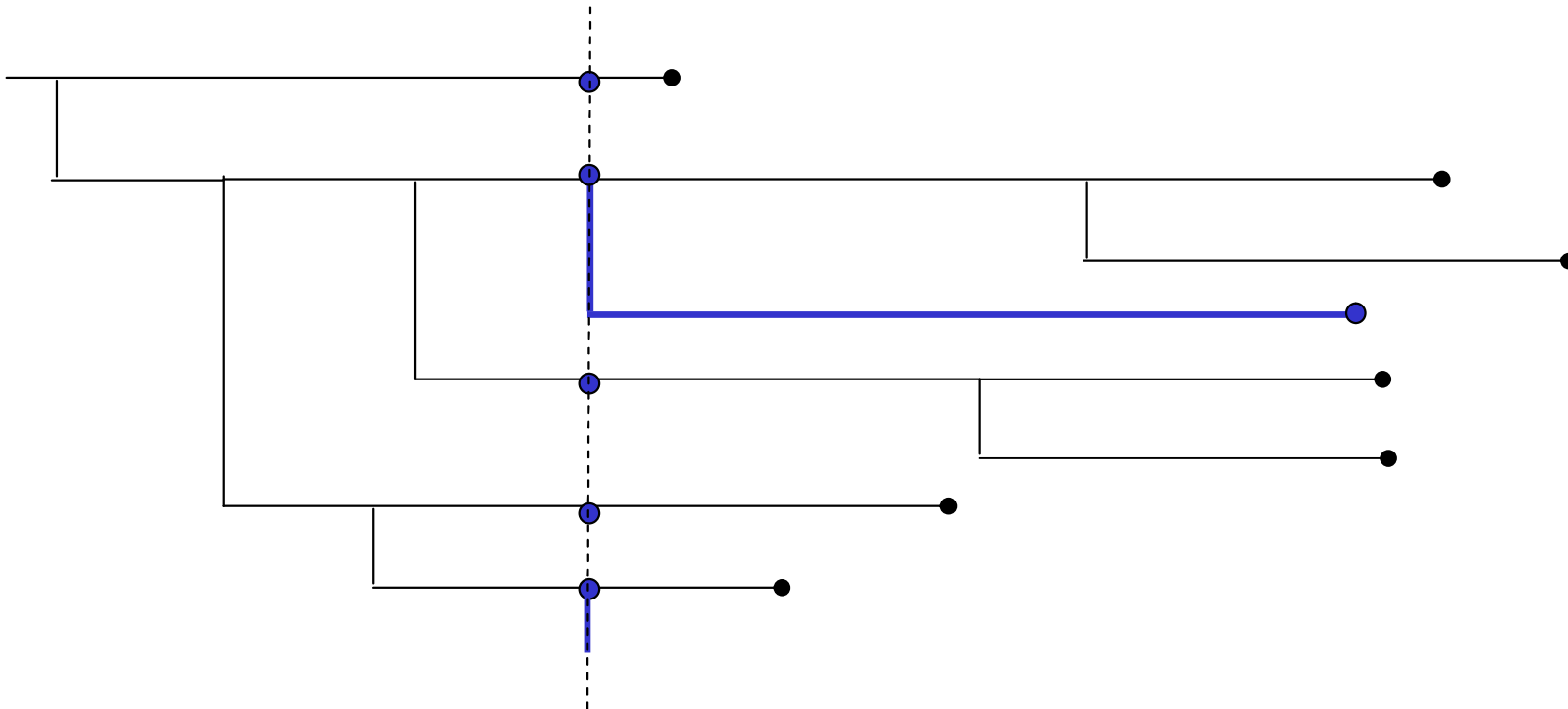
- **Total** N reads.
- **Assign** reads to **live cells** (uniformly).
- **Copy** read sequence from cell, based on **read model**.





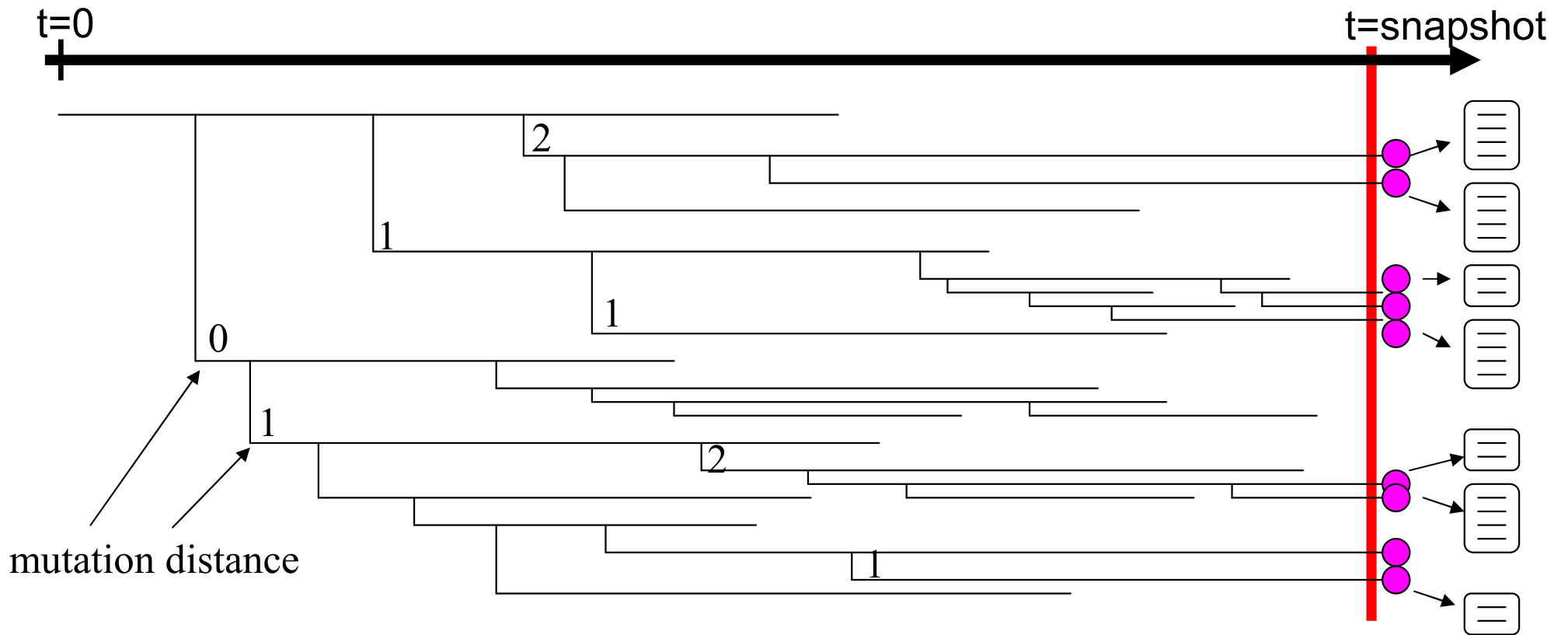
Inference

- MCMC
- Moves to **manipulate tree** and read assignments
- Generating samples of trees along the way



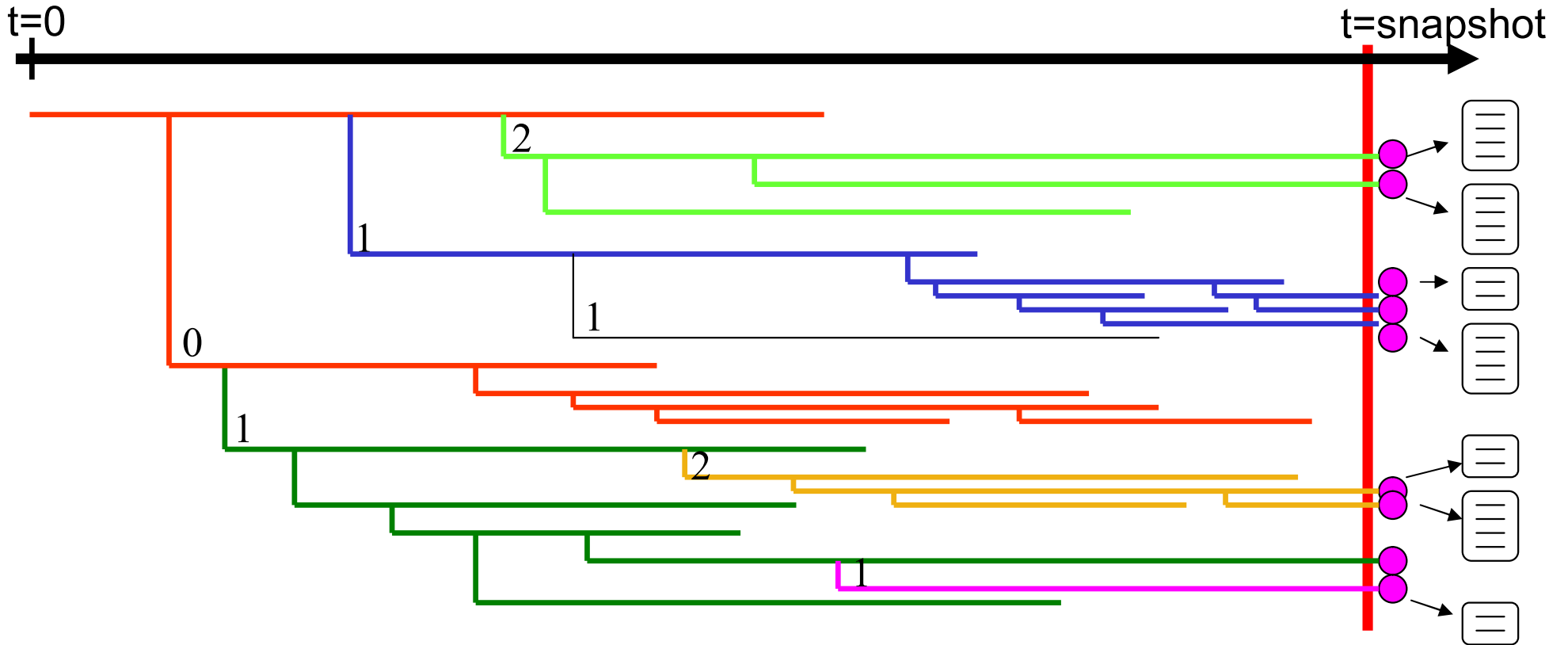
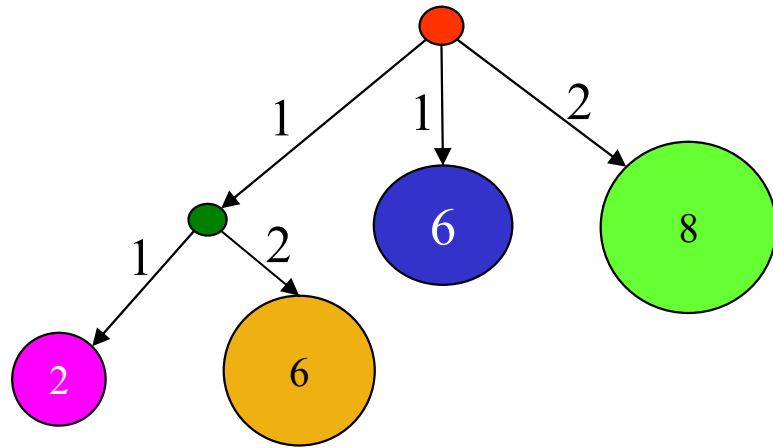


Output Tree



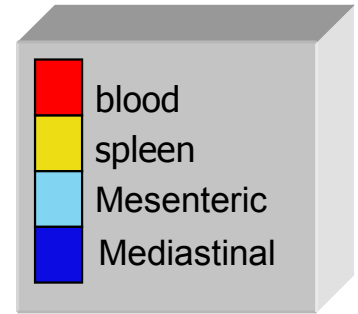
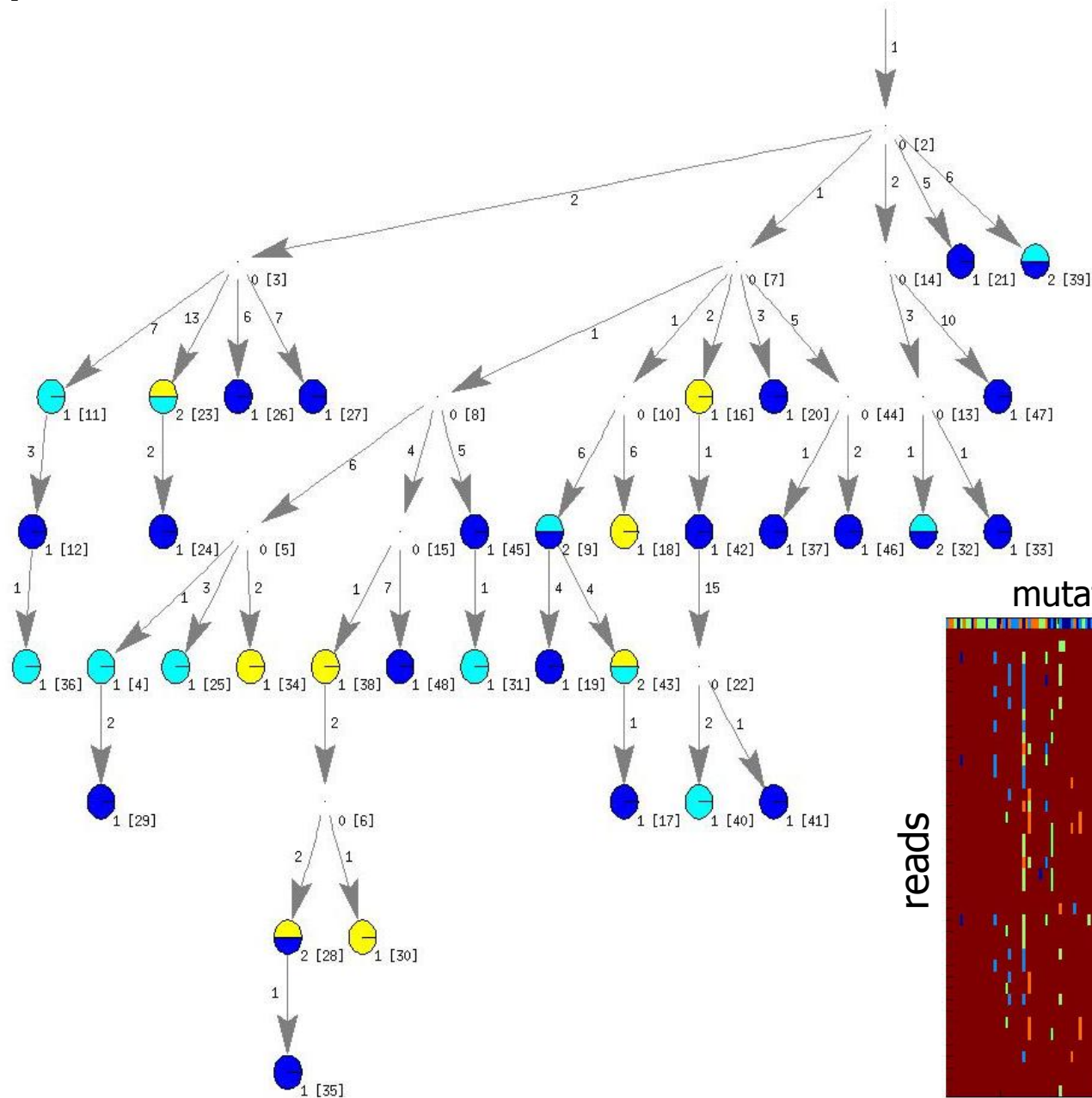


Output Tree

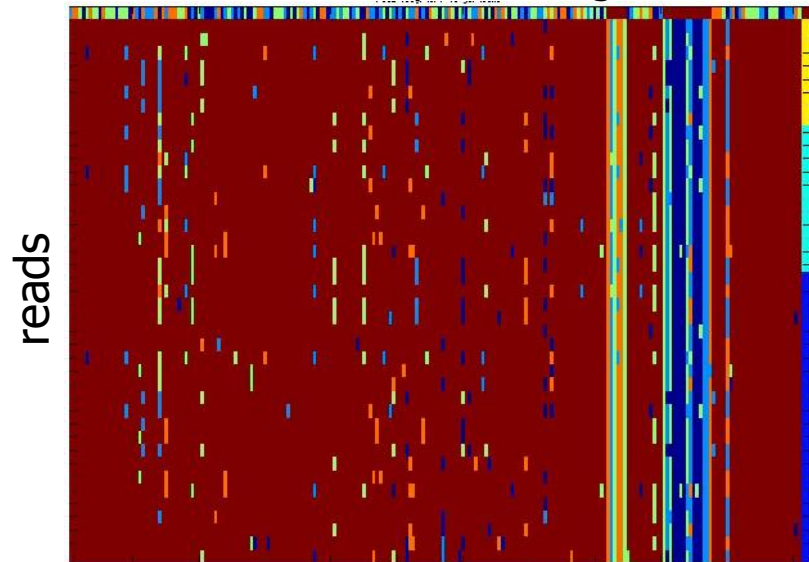




Sample Trees

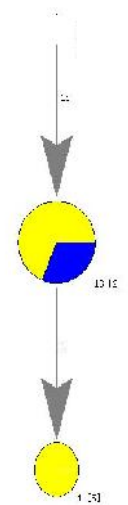
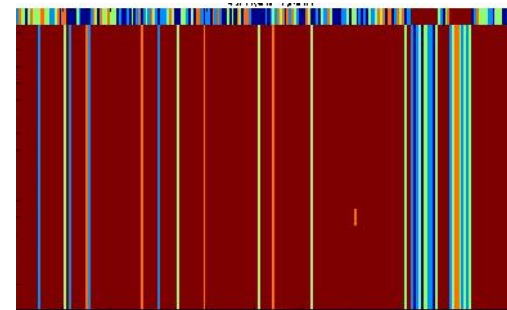
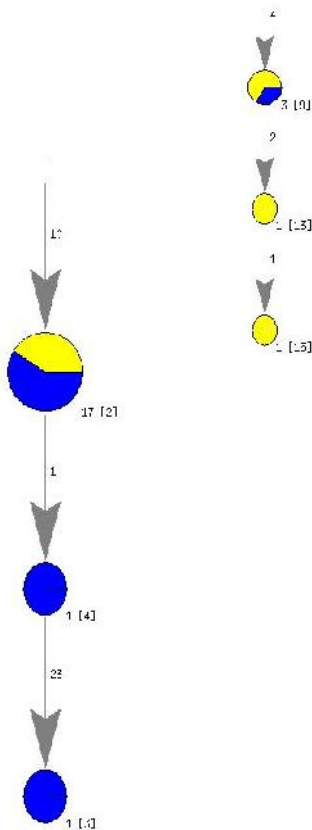
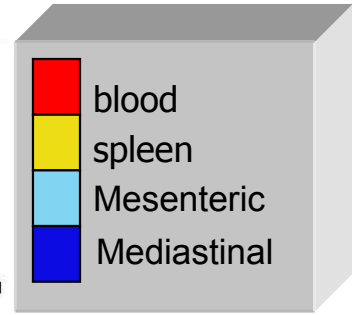
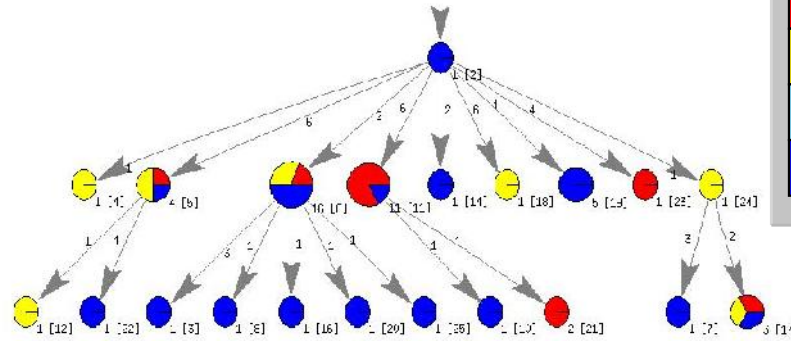
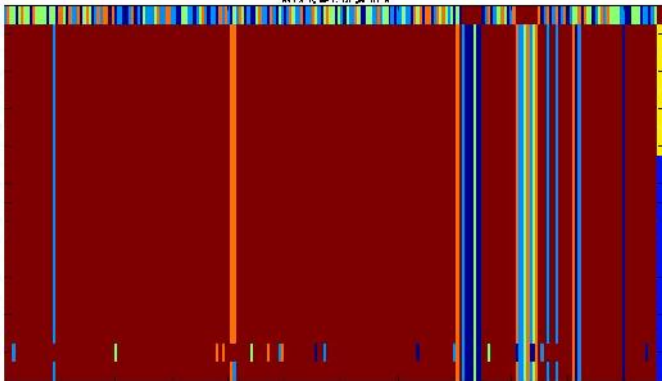
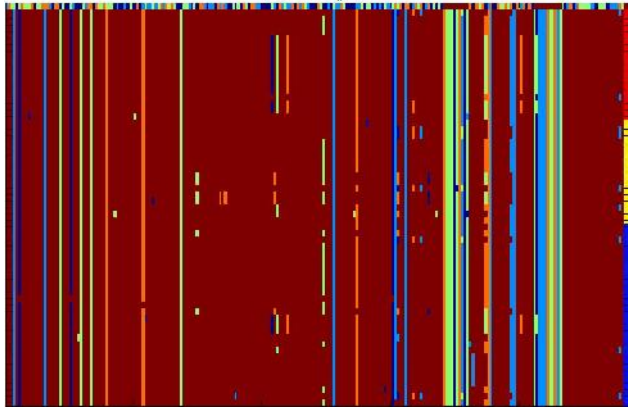


mutations compared to germline



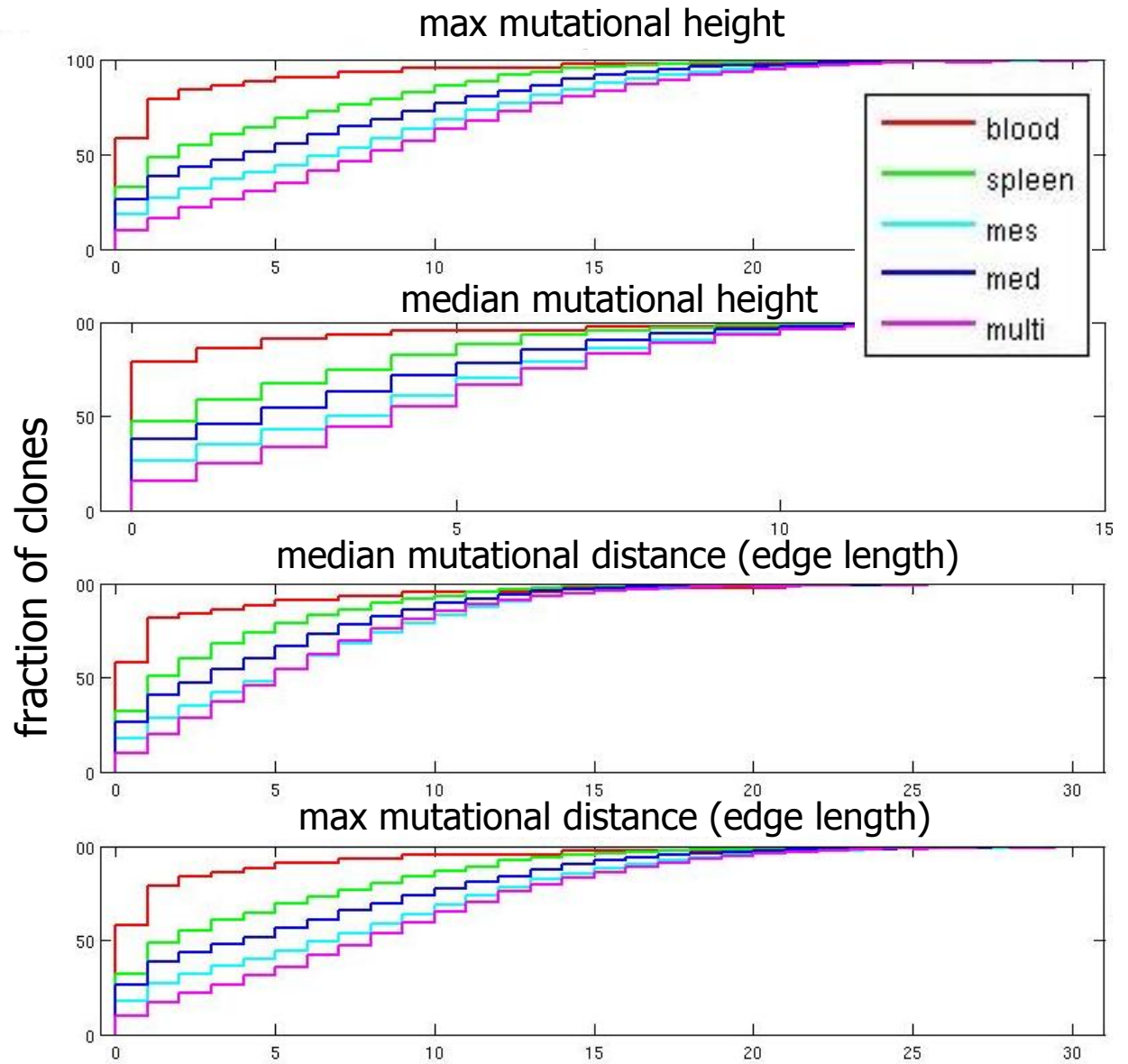
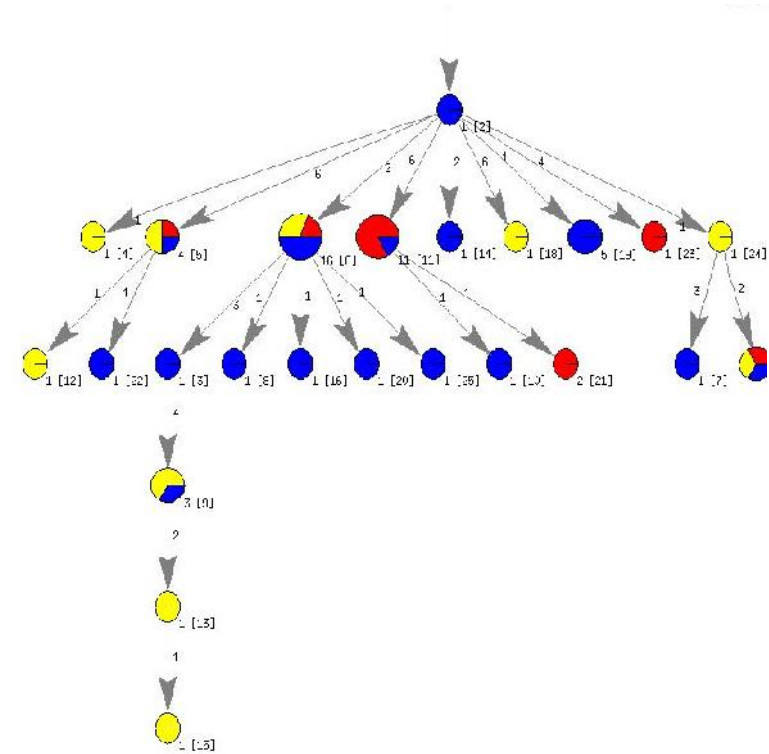


Sample Trees





Properties Of Trees

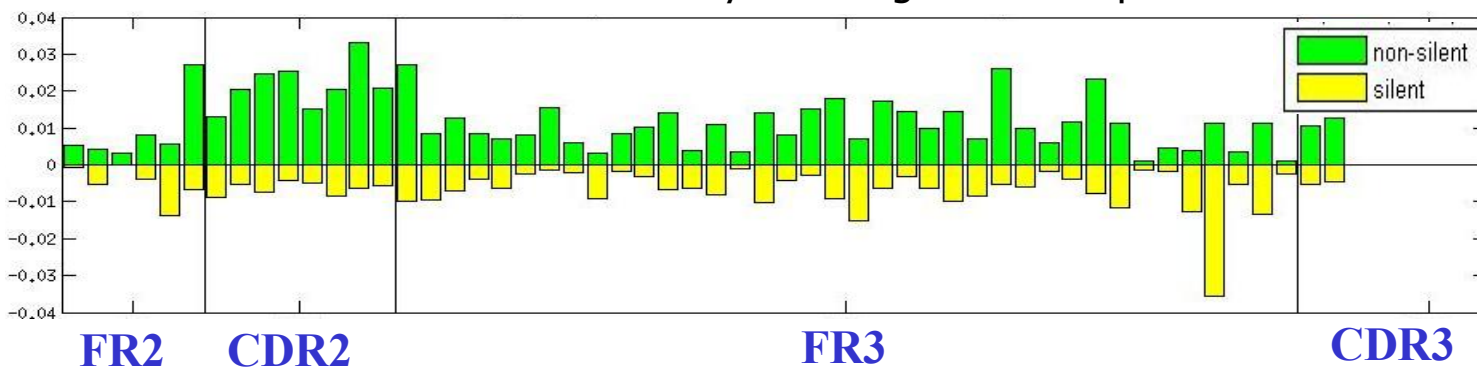




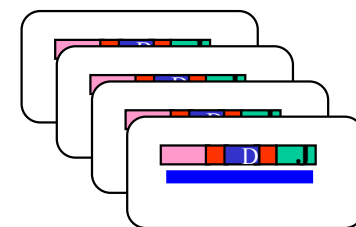
Mutation Analysis – V region

mutations binned by their aligned codon position

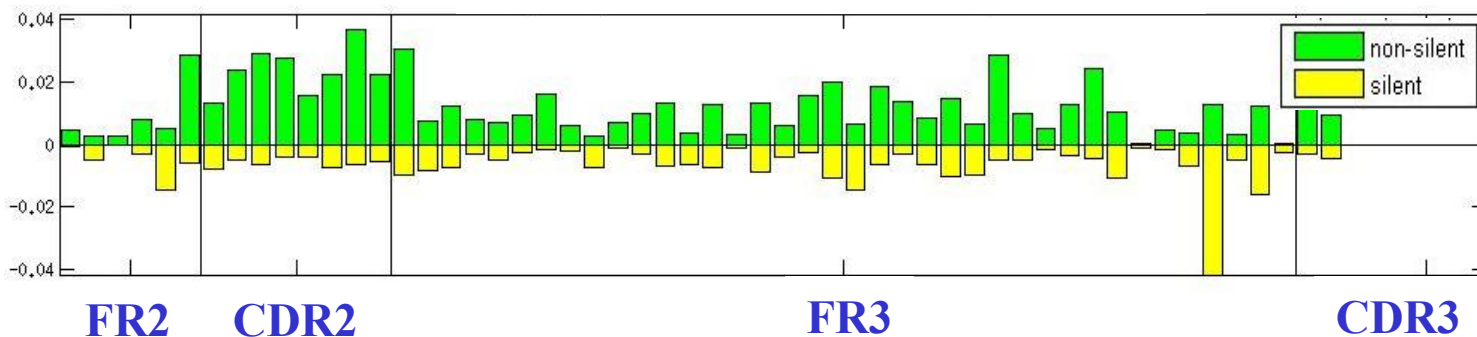
fraction of mutations



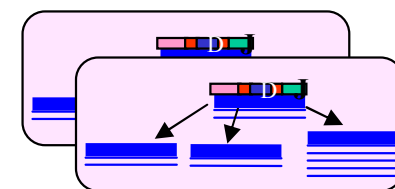
variants



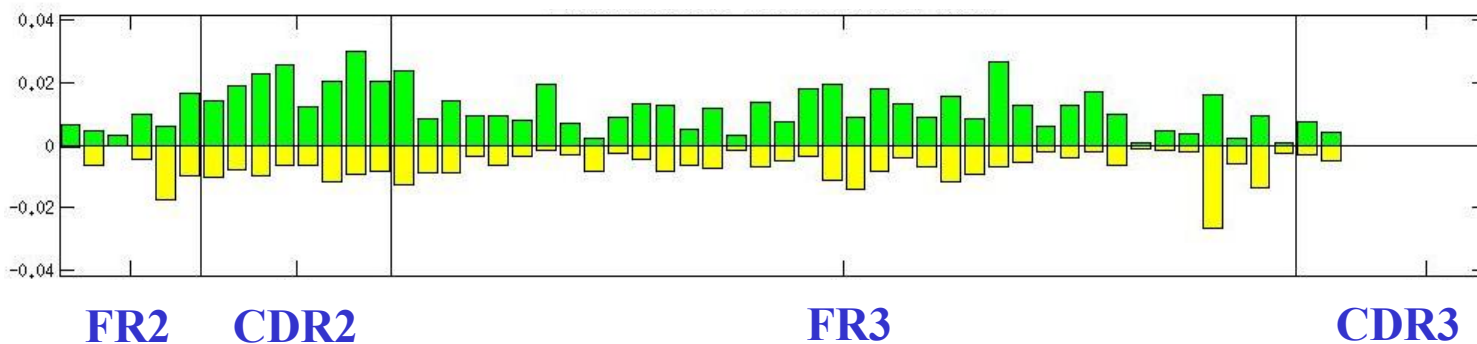
fraction of mutations



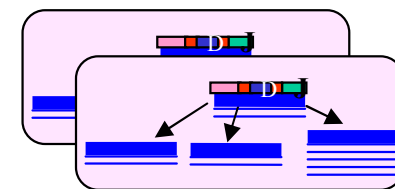
Clones



germline to root



Clones

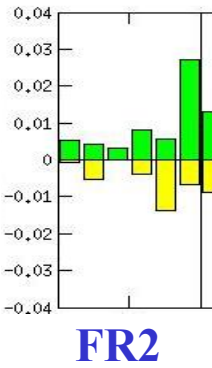


root down the tree

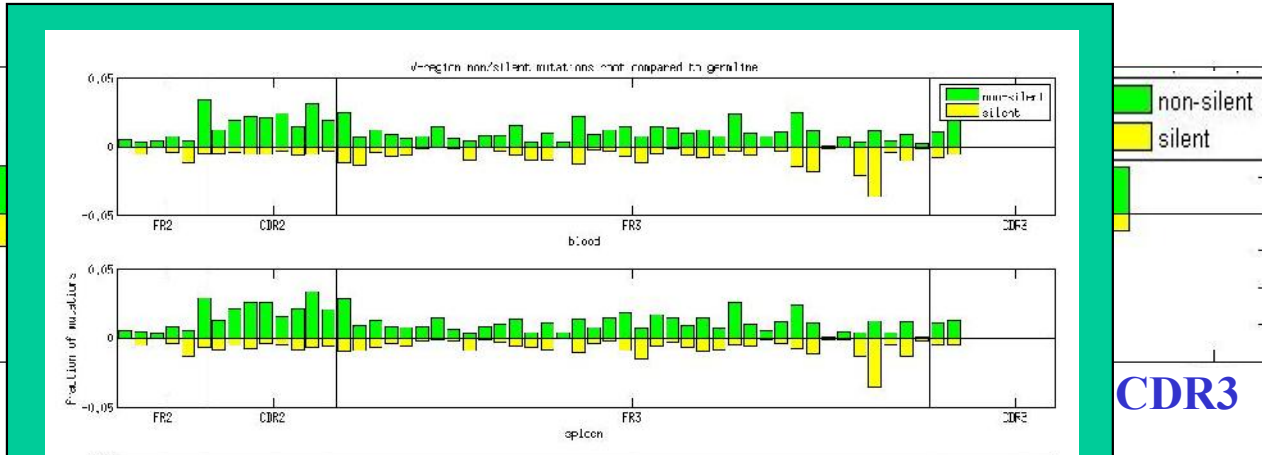


Mutation Analysis – V region

fraction of mutations

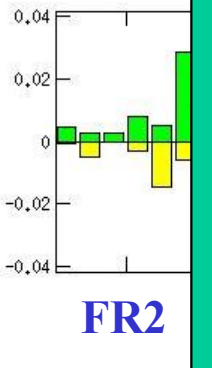


FR2

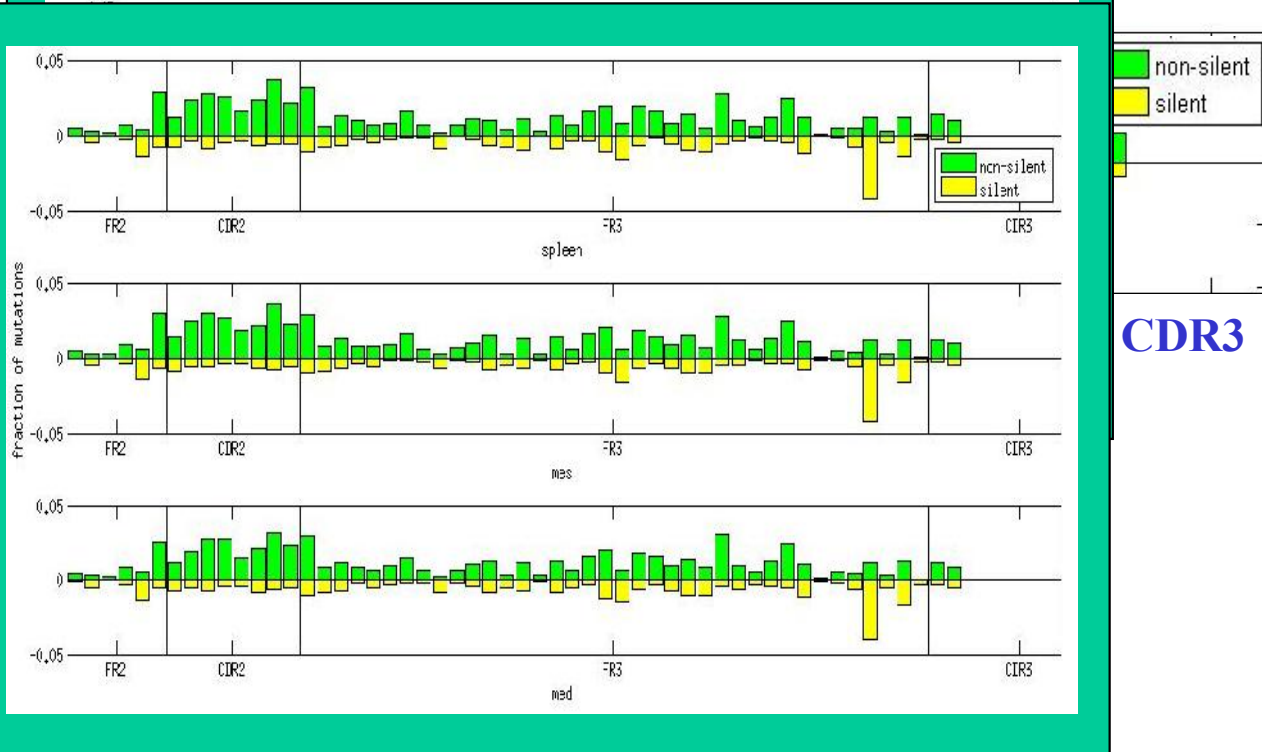


CDR3

fraction of mutations

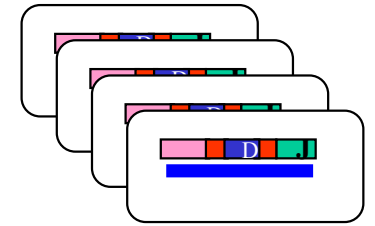


FR2

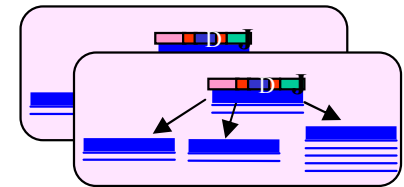


CDR3

variants



Clones



germline to root



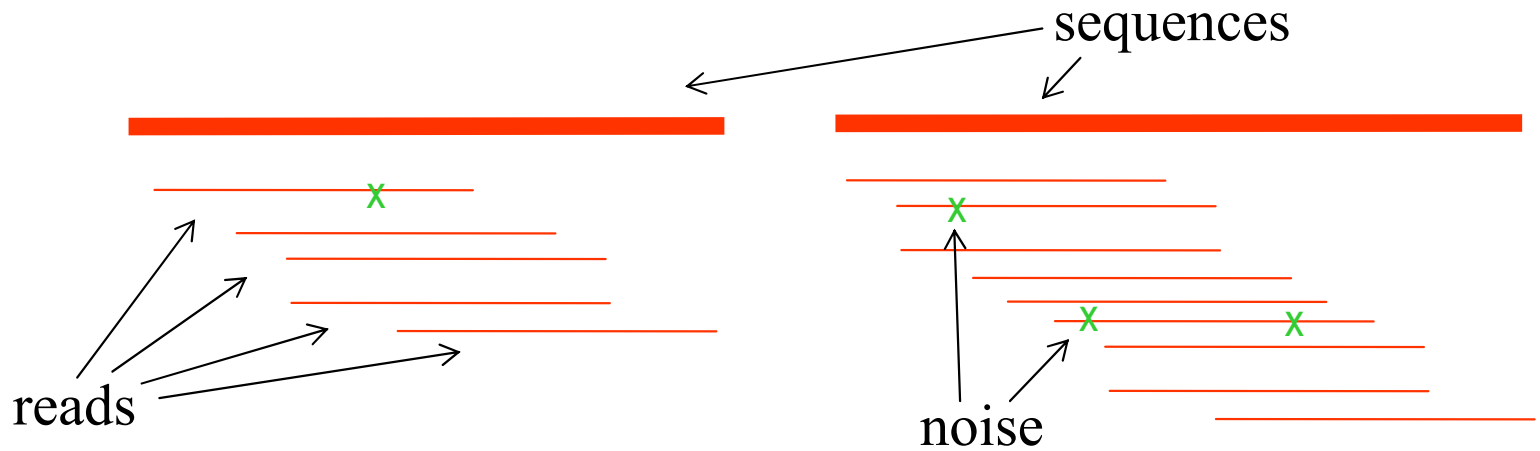
Contributions

- Fully Bayesian phylogenetic analysis
 - $O(10k)$ reads
 - Handles sequencing noise
- Big Picture view of the B cell repertoire
 - Intuitive clone description
- Insight on the Immune System
 - Evidence of clone migration
 - Differences between tissues
 - Deep hypermutation analysis
- Code (will be) freely available



Thank you!

■ Breadth



■ Depth

