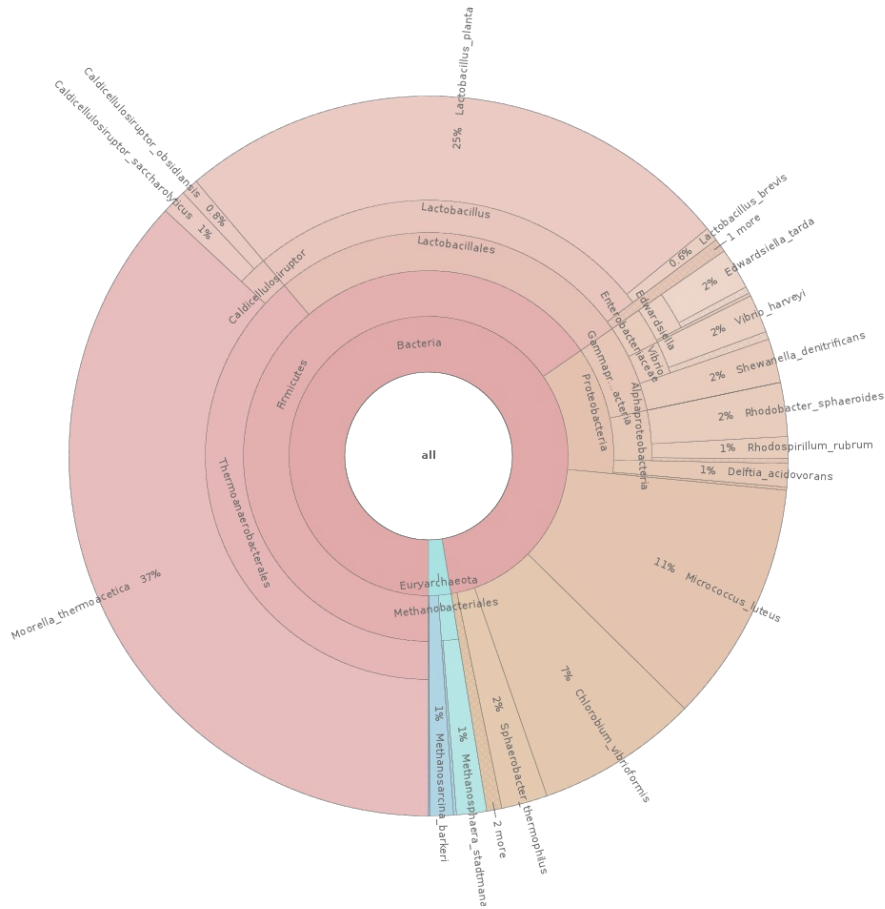


# Estimation of abundances in microbial communities from metagenomic data

September 6<sup>th</sup> 2018

Master 2 Informatique Fondamentale  
Université Paris-Est Marne-la-Vallée

Simone Pignotti

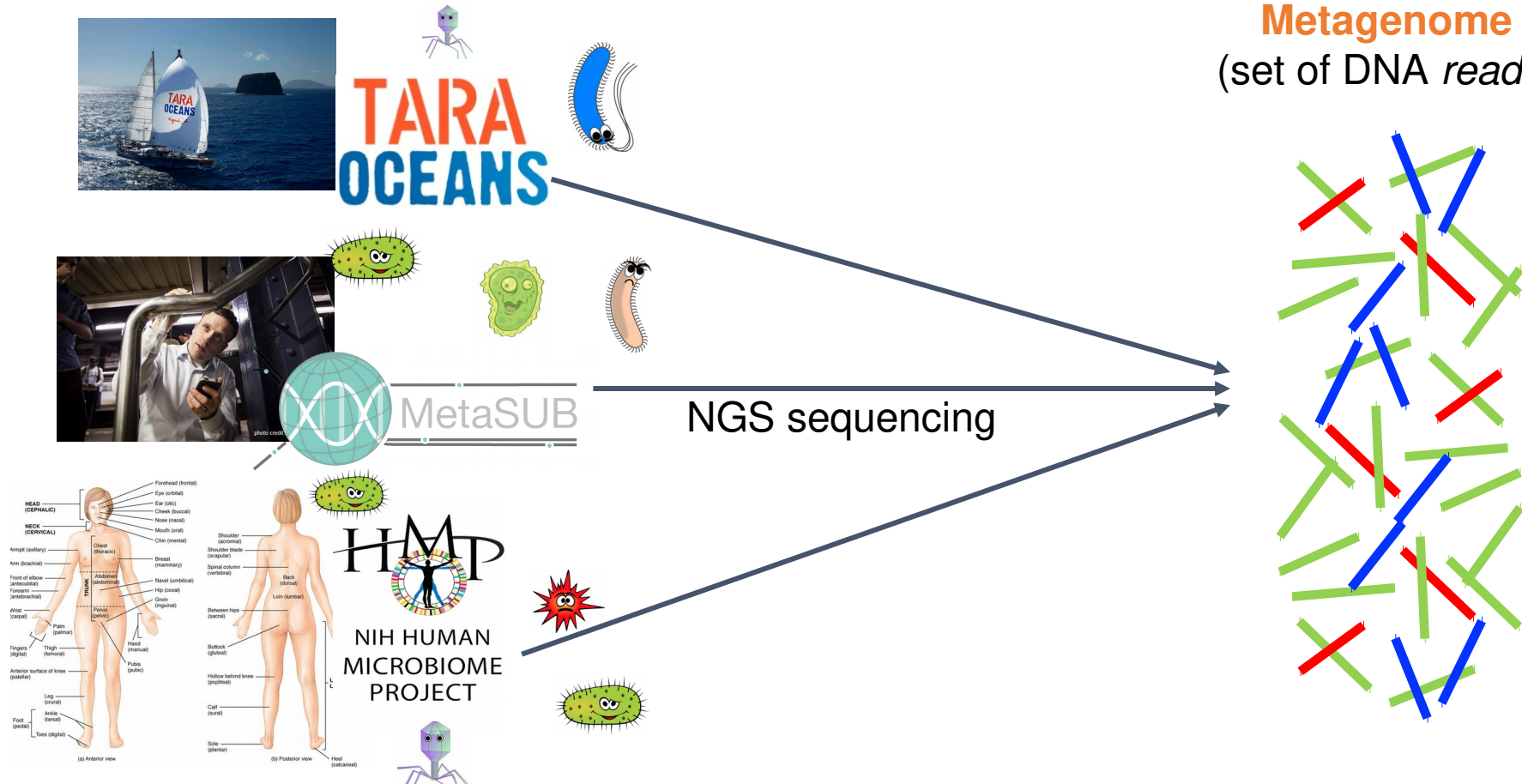


# Overview

- What is Metagenomics
- Sequence classification with ProPhyle
- Estimation of abundances
- Experiments and validation
- Conclusions and perspectives

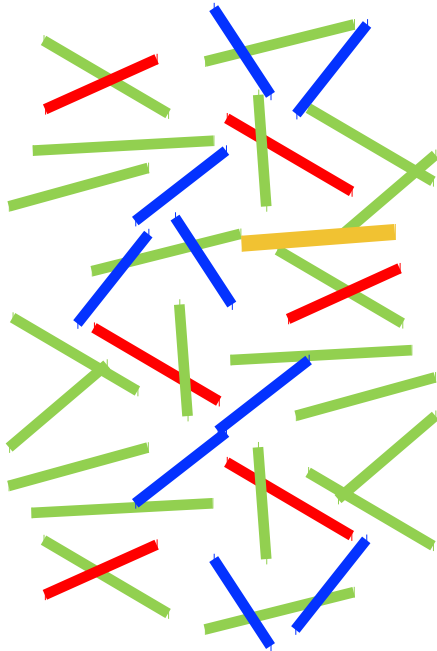
# Metagenomic sequencing

**Metagenome**  
(set of DNA *reads*)

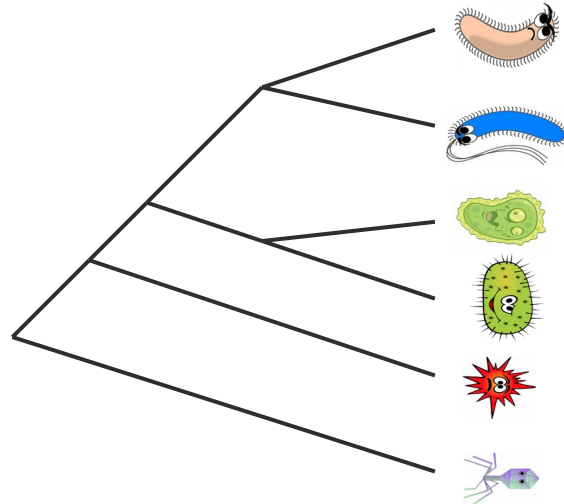


# Phylogeny-based metagenomic classification

**Metagenome** reads  
(billions)



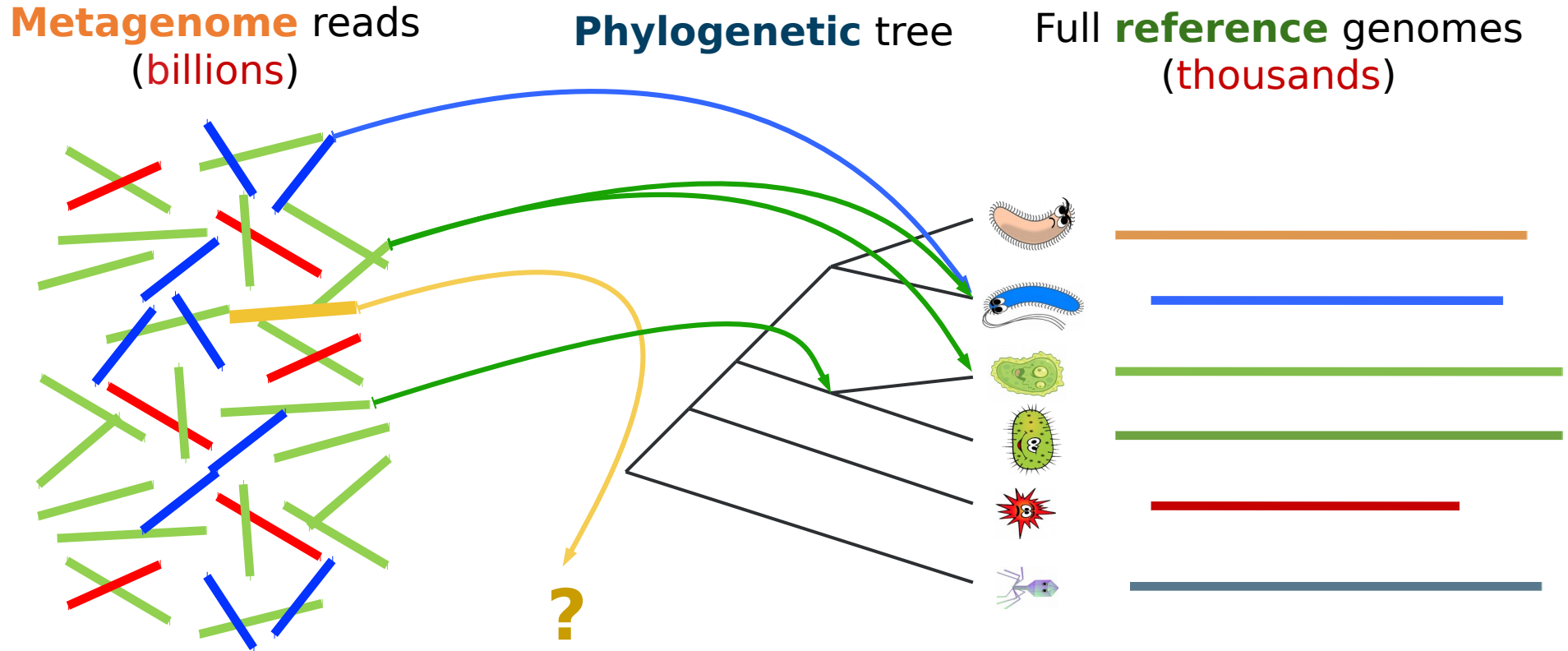
**Phylogenetic** tree



Full **reference** genomes  
(thousands)



# Phylogeny-based metagenomic classification



# Alignment-free methods

**Reference** genome:



**Index** of  $k$ -mers

“Is this  $k$ -mer in the genome?” **yes/no**



**Goal:** estimate the “likelihood” of a read to belong to the genome (assign a score, e.g. number of hits)

**Read:**



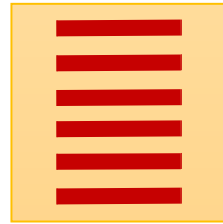
$k$ -mers



# Alignment-free methods

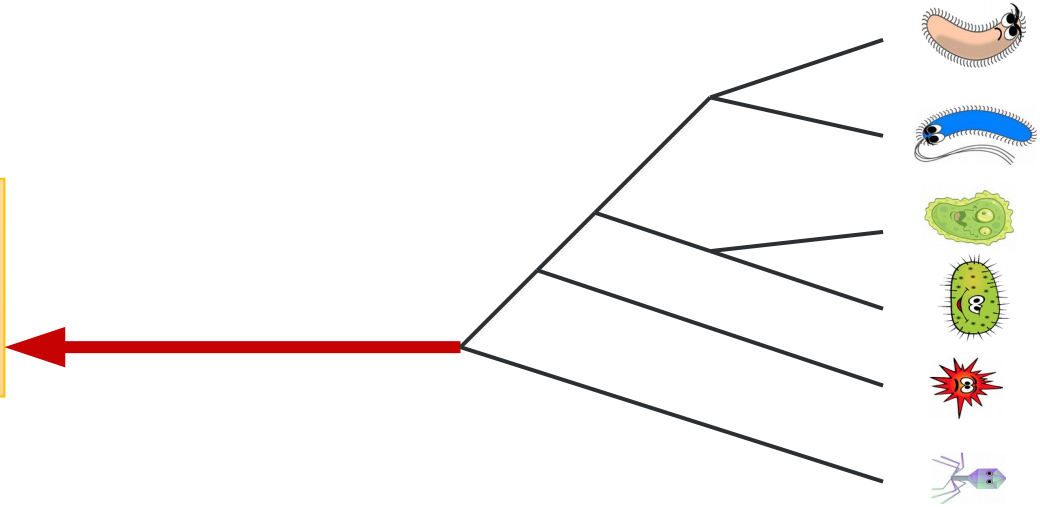
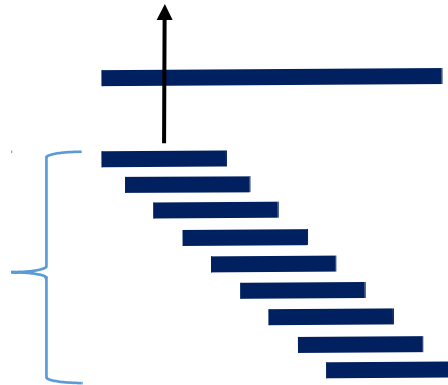
## **Index** of $k$ -mers

- Huge
- Repeated  $k$ -mers



**Read:**

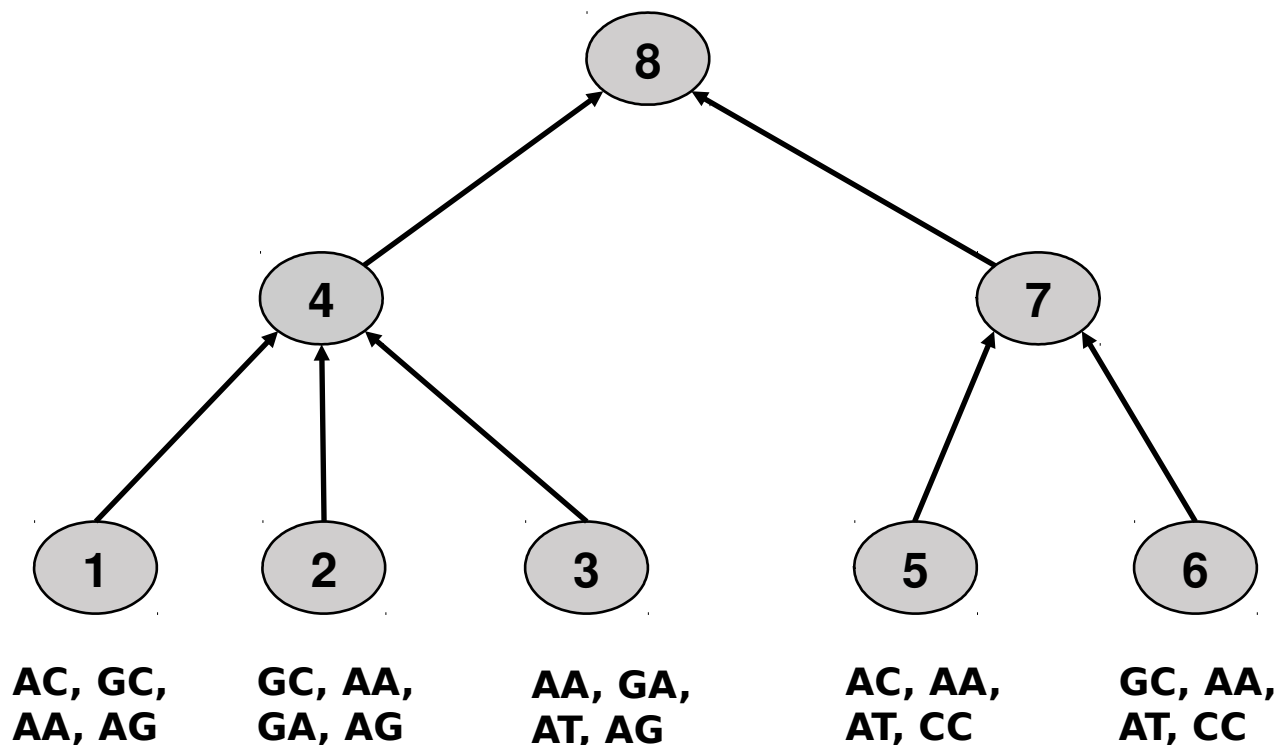
$k$ -mers



Same problem, with tree

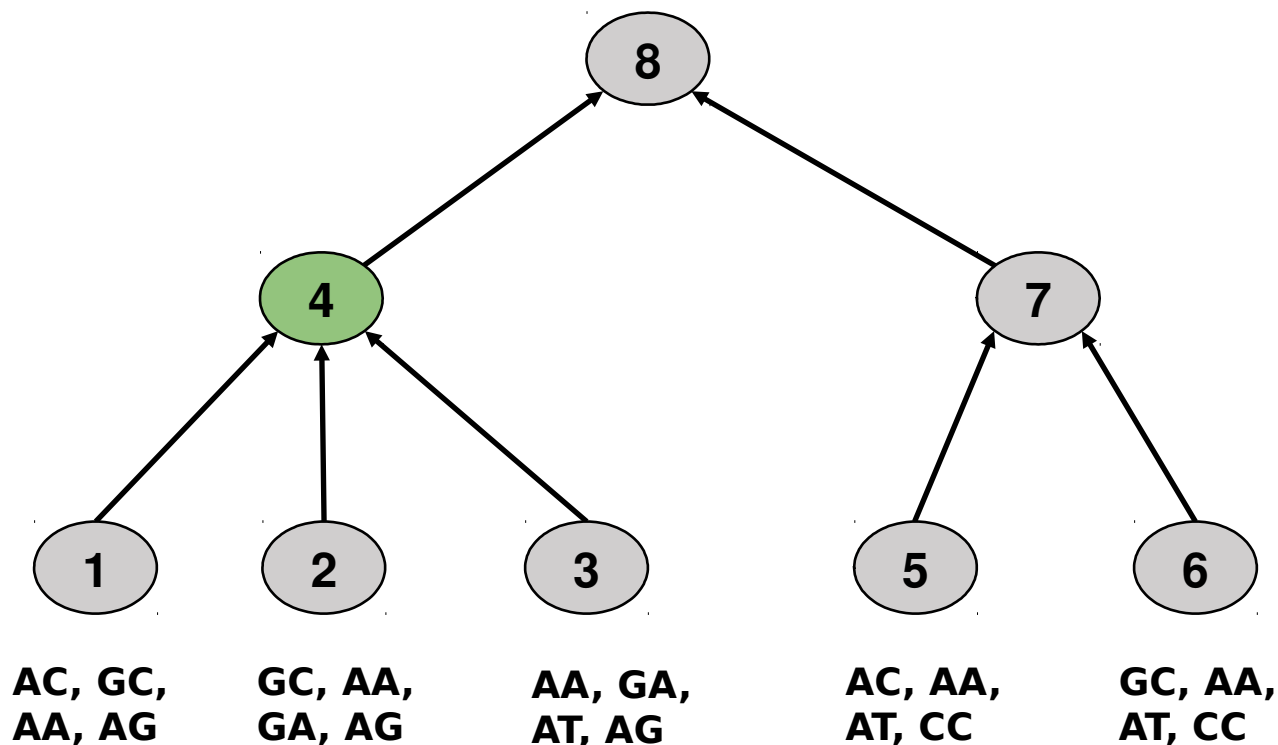
How to **compress** it?

# *k*-mer propagation - ProPhyle

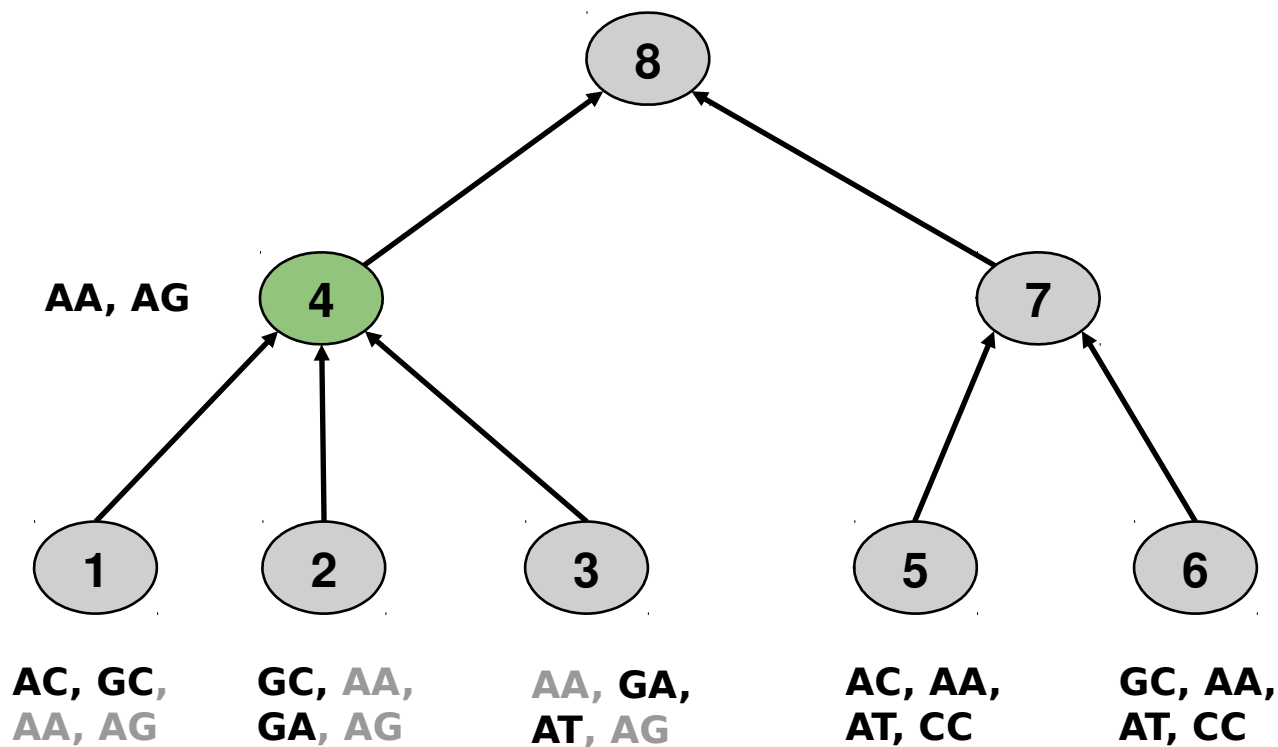




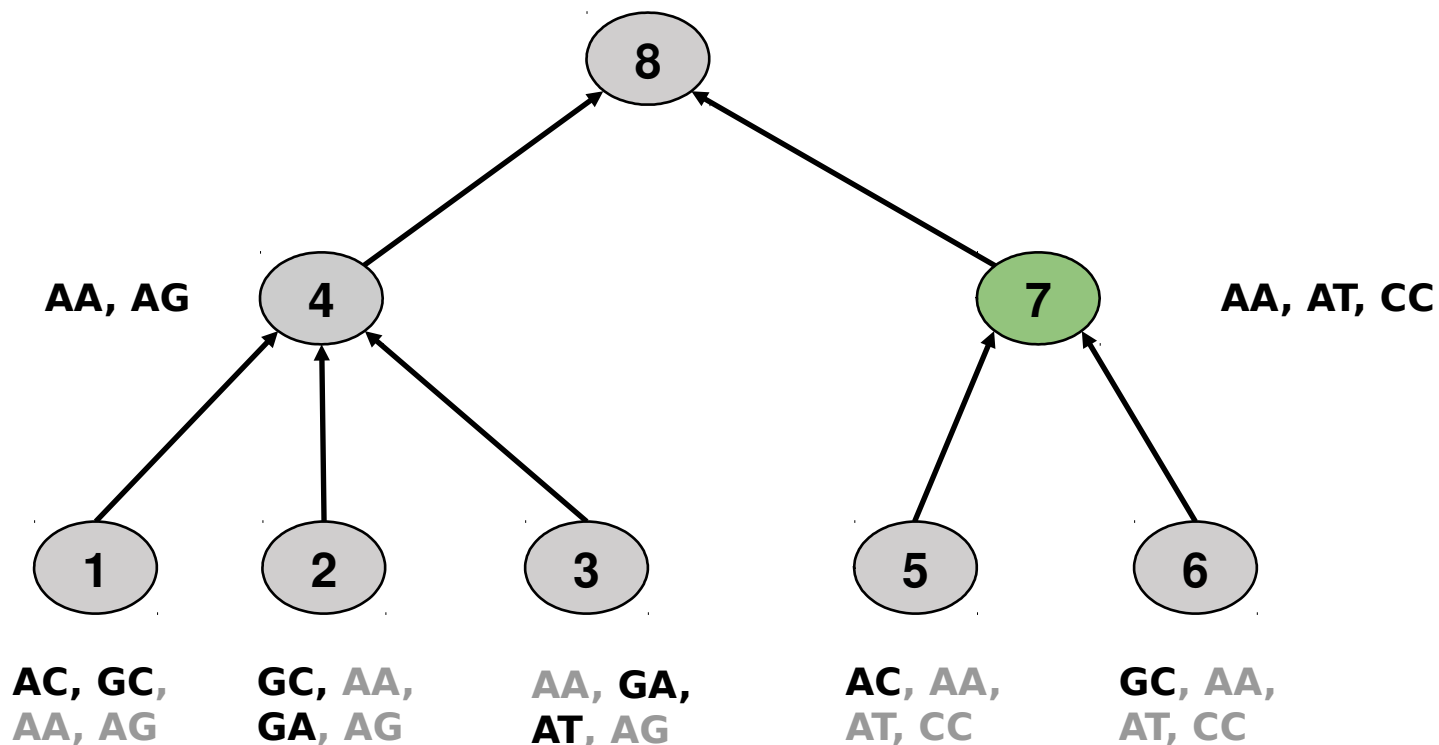
# *k*-mer propagation - ProPhyle



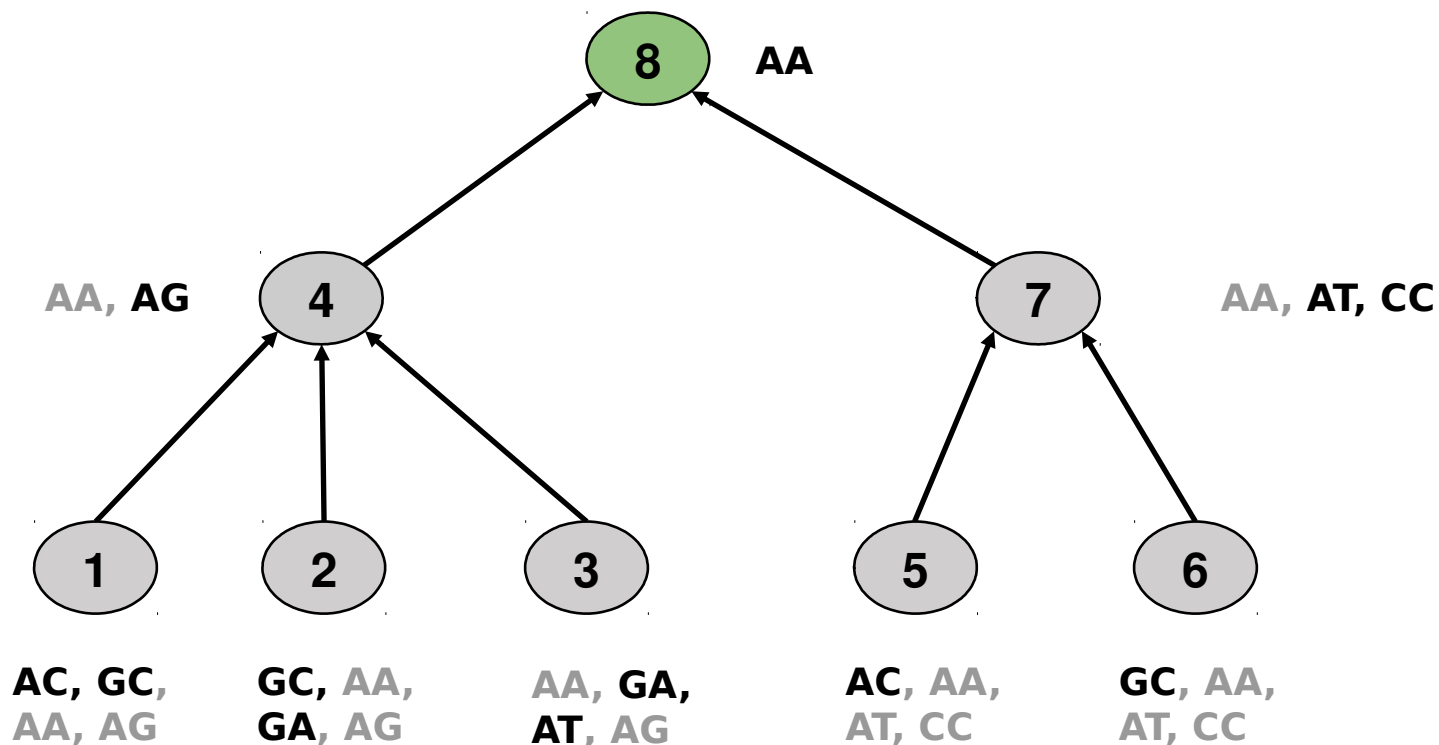
# *k*-mer propagation - ProPhyle



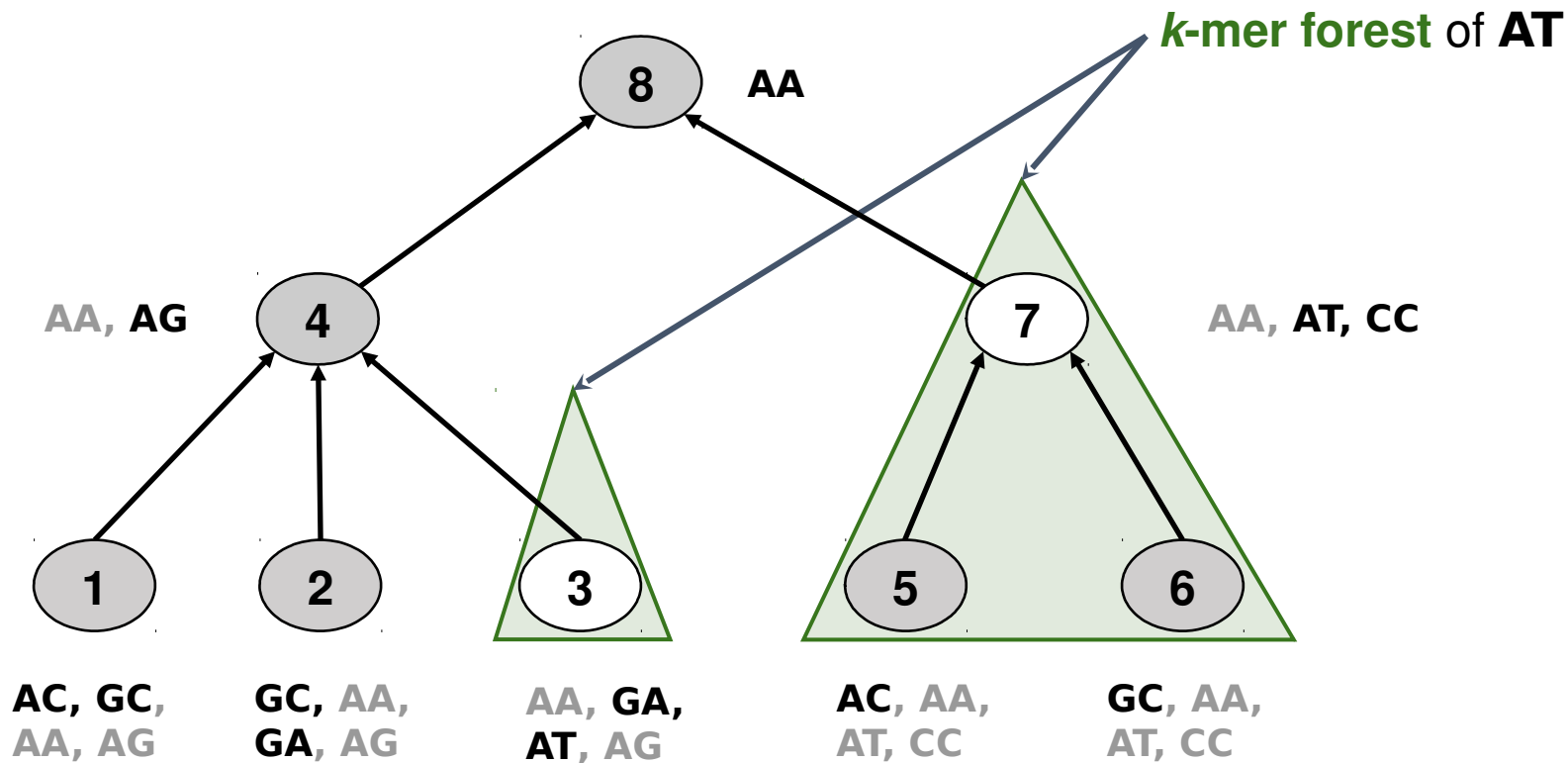
# *k*-mer propagation - ProPhyle



# k-mer propagation - ProPhyle

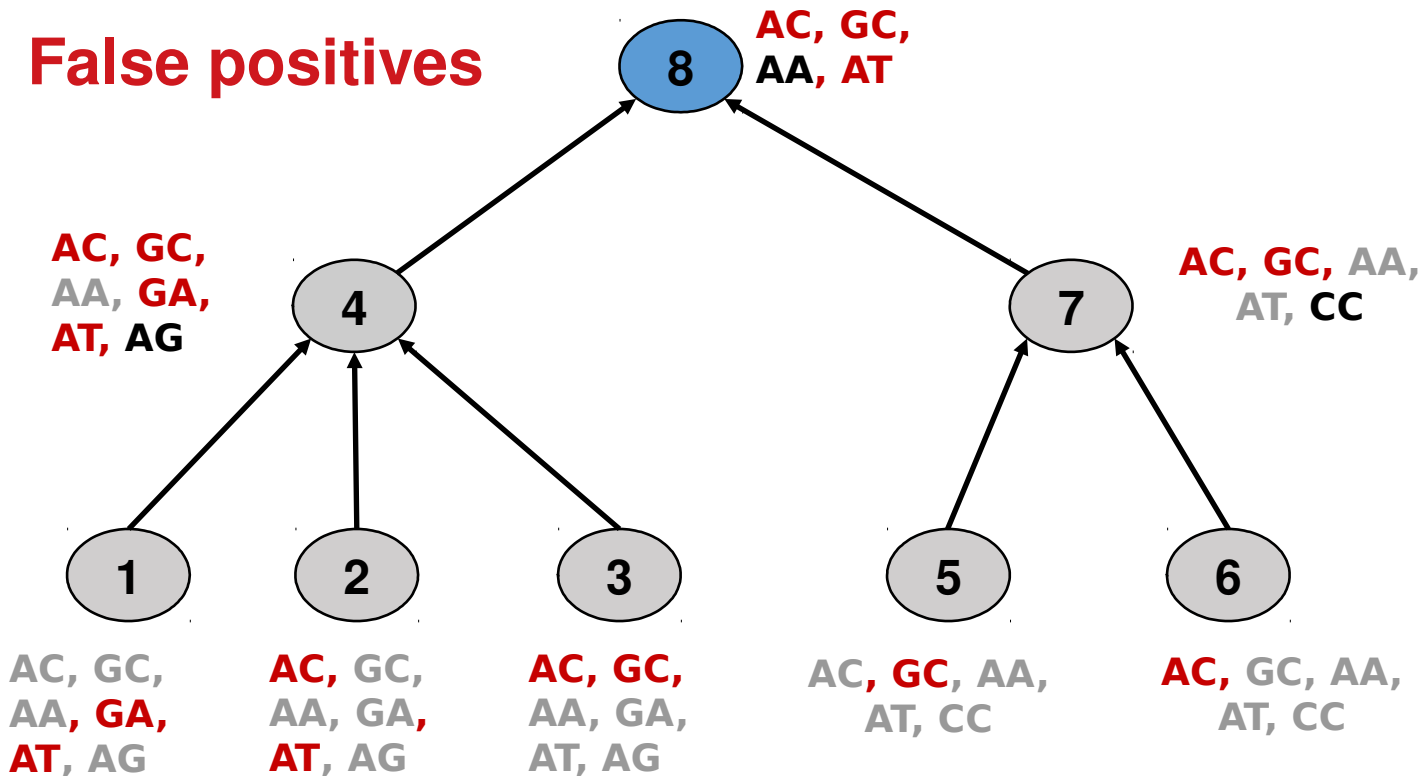


# $k$ -mer propagation - ProPhyle

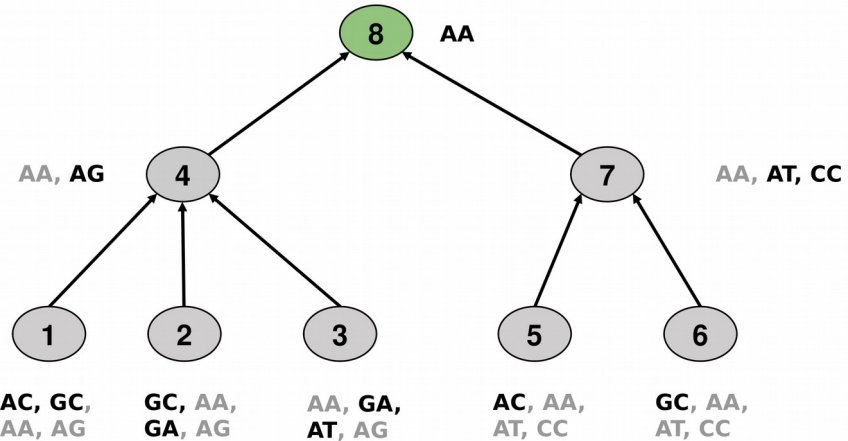


# Kraken's LCA

(lowest common ancestors)

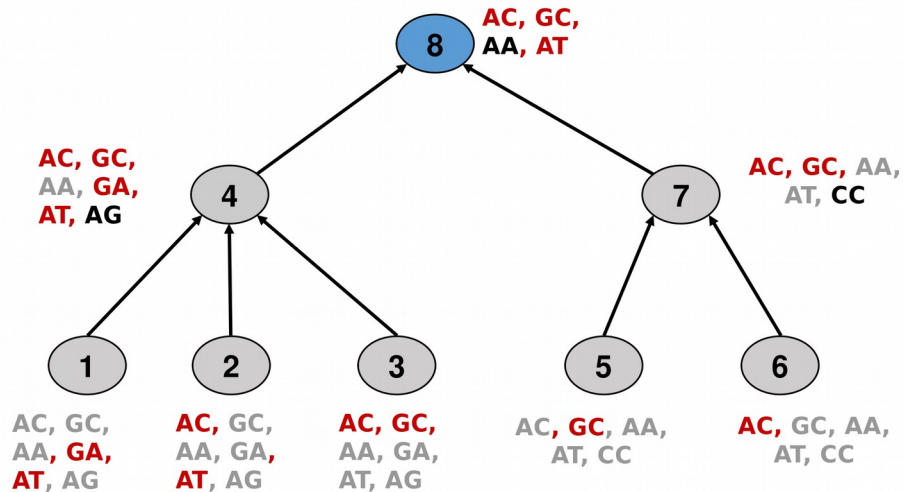


# Assignments – ProPhyle vs Kraken

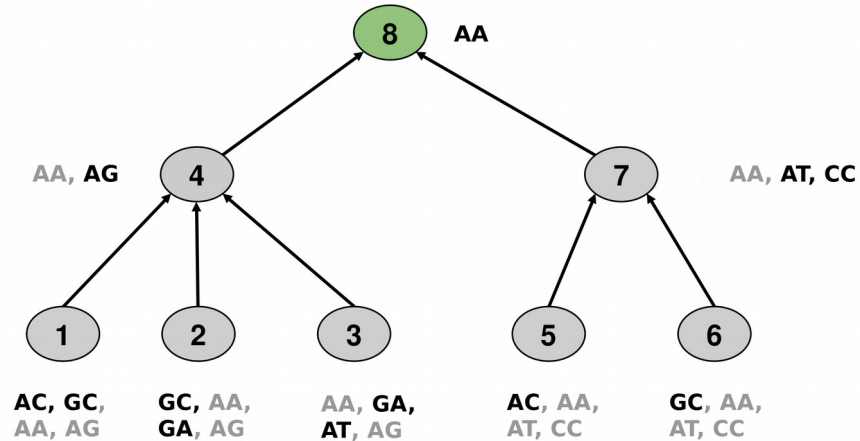


Read: **AGC**

k-mers: **AG**  
**GC**



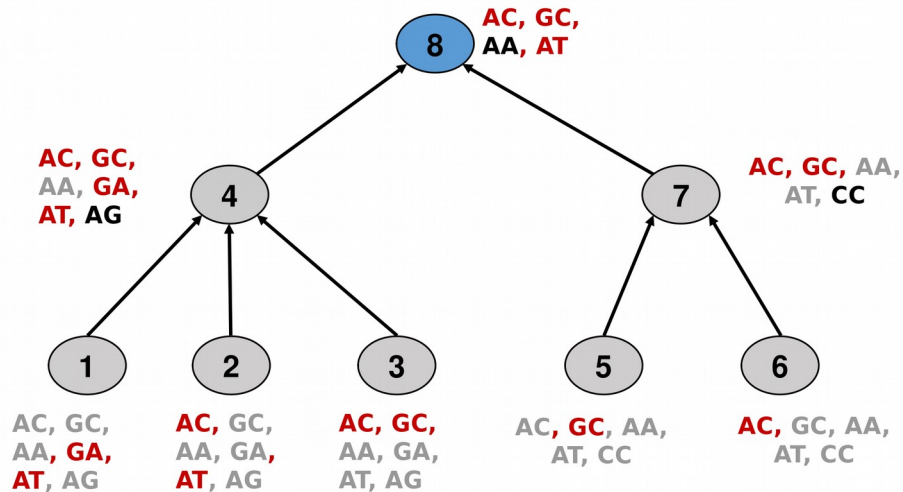
# Assignments – ProPhyle vs Kraken



Read: **AGC**

k-mers: **AG**  
**GC**

ProPhyle: **AG** → {4} → {1, 2, 3}  
**GC** → {1, 2}





# Assignments – ProPhyle vs Kraken

Read: **AGC**

k-mers: **AG**  
**GC**

ProPhyle: **AG** → {4} → {1, 2, 3}  
**GC** → {1, 2, 6}

$s(1) = 2$

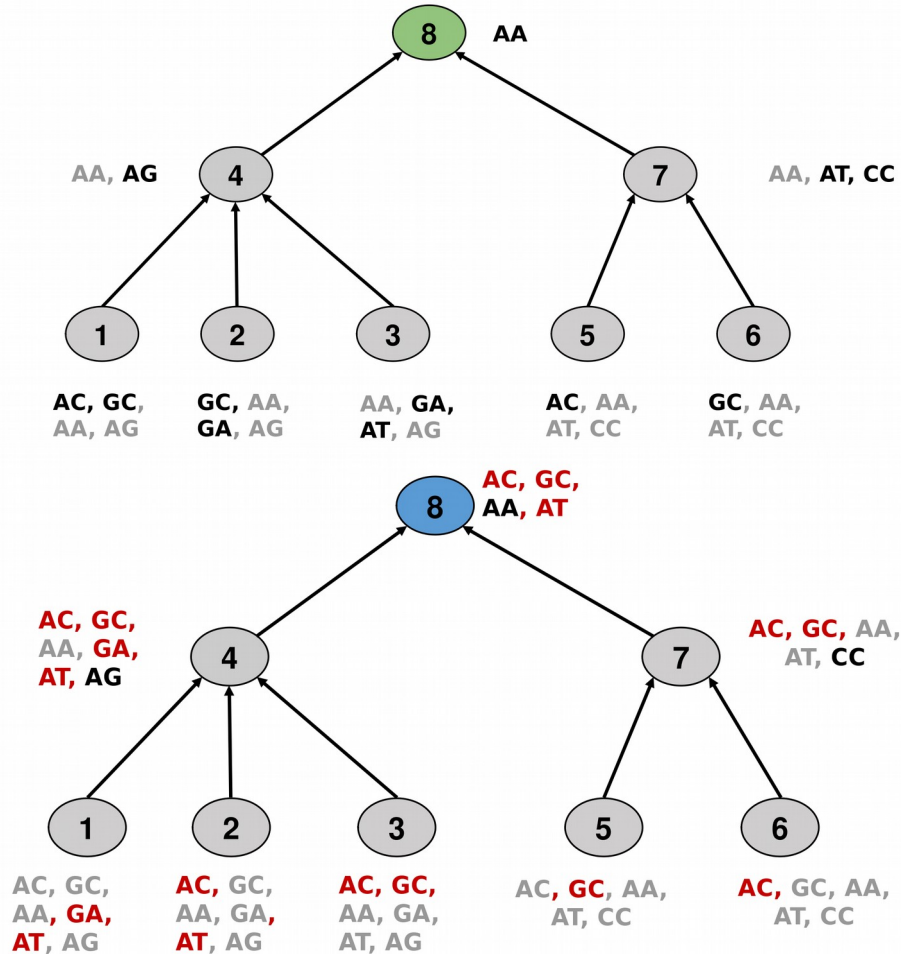
$s(2) = 2$

$s(3) = 1$

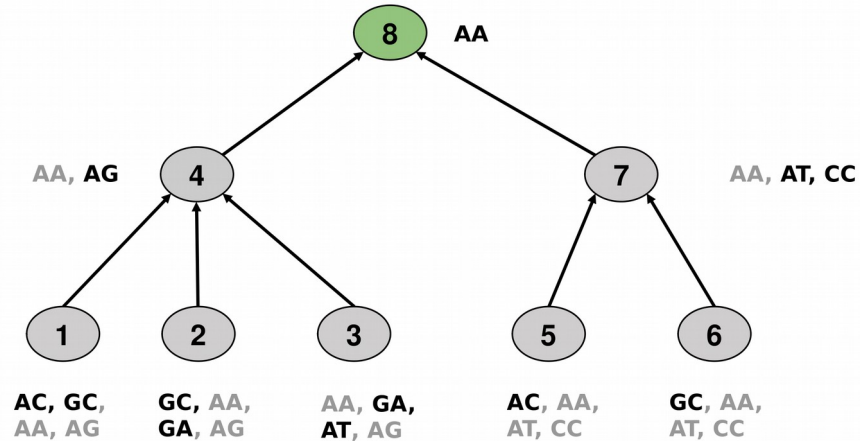
$s(6) = 1$

$s(*) = 0$

read → {1, 2}



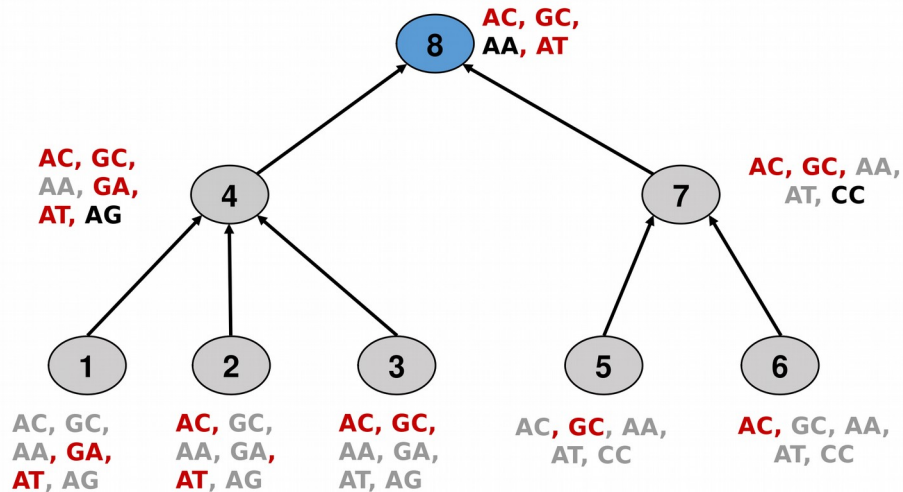
# Assignments – ProPhyle vs Kraken



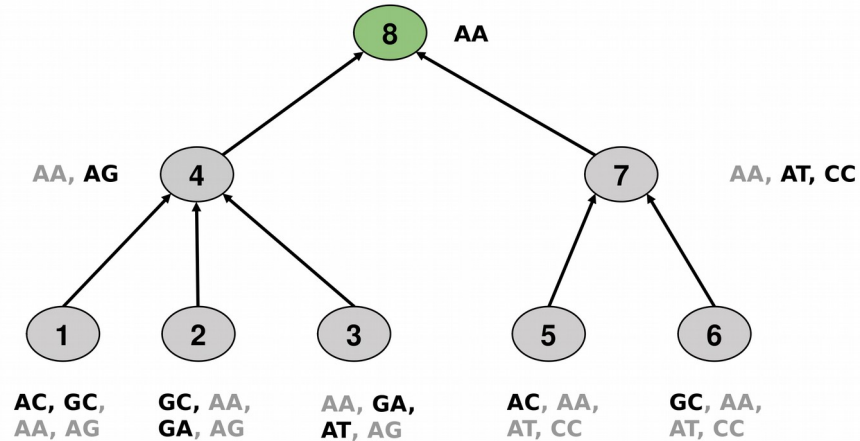
Read: **AGC**

k-mers: **AG**  
**GC**

Kraken: **AG** → 4 → {1, 2, 3}  
**GC** → 8 → {\*}



# Assignments – ProPhyle vs Kraken



Read: **AGC**

k-mers: **AG**  
**GC**

Kraken: **AG** → 4 → {1, 2, 3}  
**GC** → 8 → {\*}

$s(1) = 2$

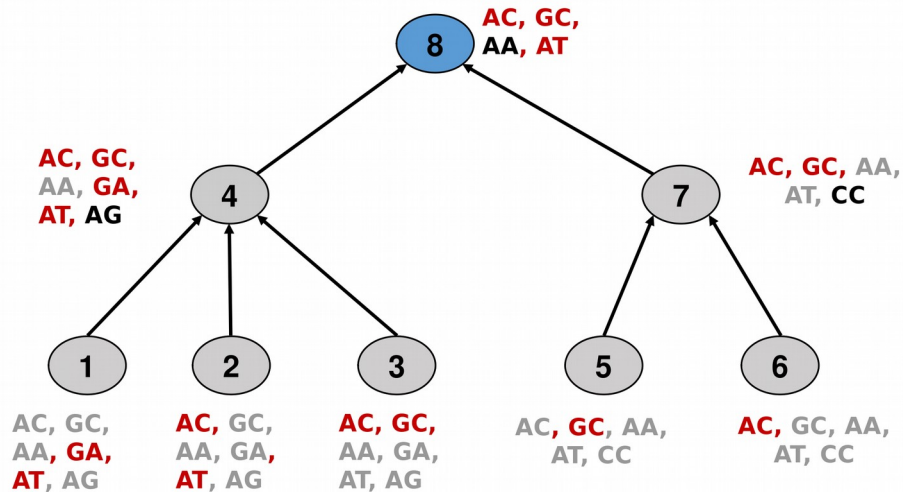
$s(2) = 2$

$s(3) = 2$

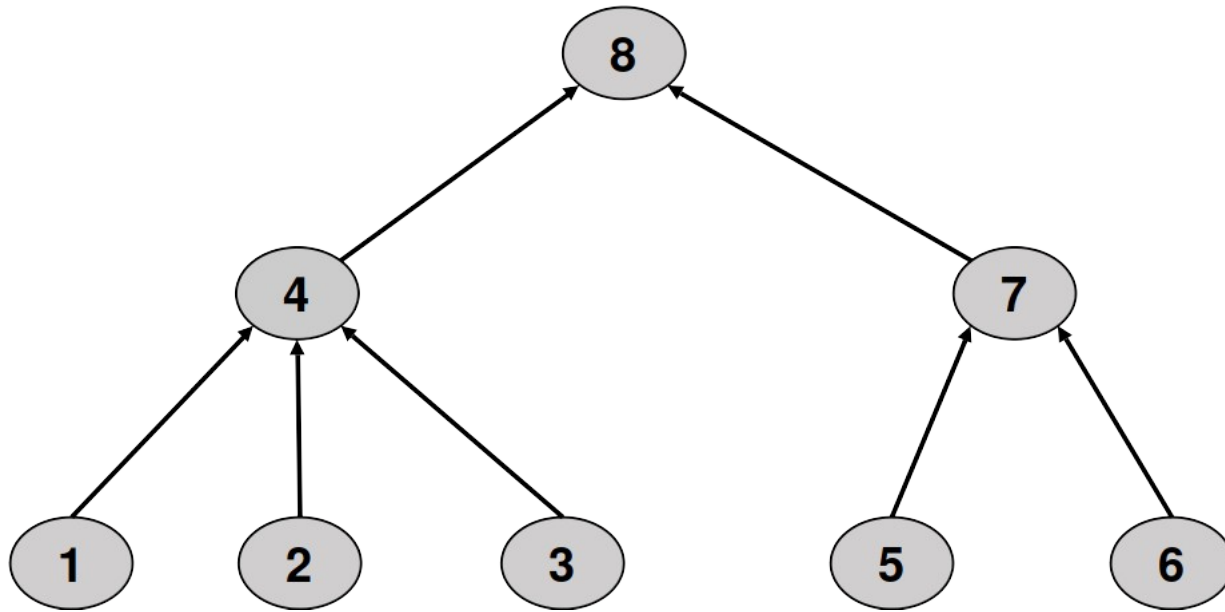
$S(*) = 1$

read → 4

4 = **LCA**({1, 2, 3})

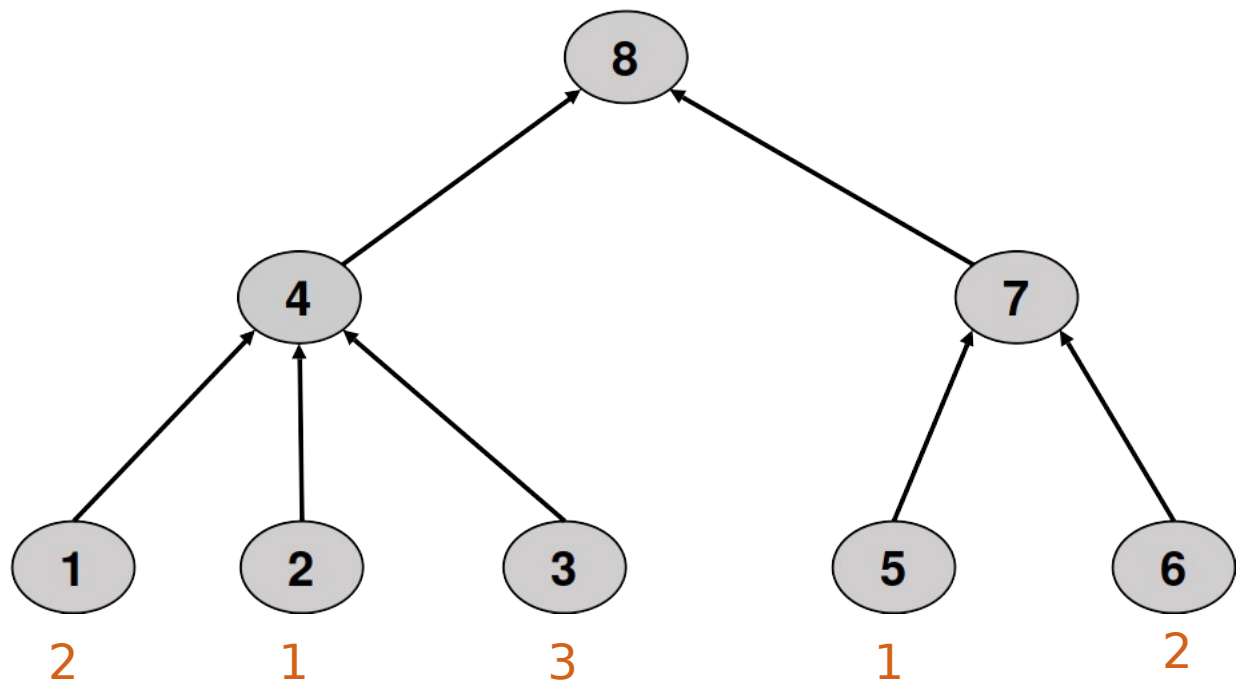


# What to do with multiple assignments?



R1  $\rightarrow \{1, 2\}$   
R2  $\rightarrow \{3, 7\} \rightarrow \{3, 5, 6\}$   
R3  $\rightarrow \{3\}$   
R4  $\rightarrow \{1, 3, 6\}$

# What to do with multiple assignments?



R1  $\rightarrow \{1, 2\}$   
R2  $\rightarrow \{3, 7\} \rightarrow \{3, 5, 6\}$   
R3  $\rightarrow \{3\}$   
R4  $\rightarrow \{1, 3, 6\}$

“Redistribute” reads

Consider all  $2^n - 1$  possible combinations

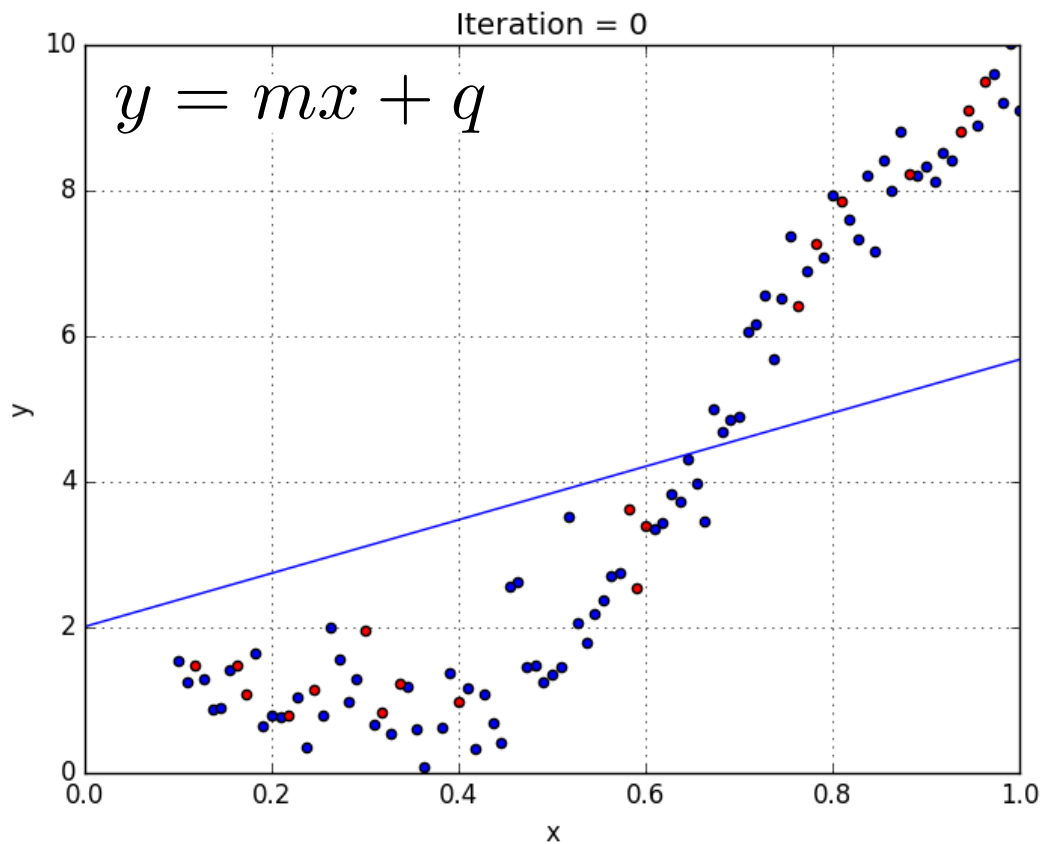
$\sim 10^{301}$  for  $n = 1000$

$\Sigma = 9$

but only 4 reads...



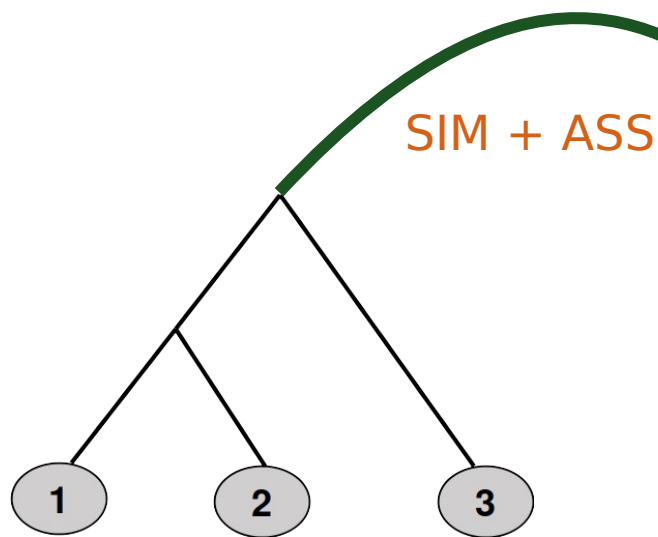
# Linear Regression!



# Acquiring statistics about mult. ass.

**G<sub>3</sub>:**

R1	→	{1, 2, 3}
R2	→	{2, 3}
R3	→	{3}
R[4-10]	→	{3} (x6)



<b>S</b>	G <sub>1</sub>	G <sub>2</sub>	G <sub>3</sub>
G <sub>1</sub>	<b>1</b>	0.1	<b>0.1</b>
G <sub>2</sub>	0.2	<b>1</b>	0.2
G <sub>3</sub>	0	0.1	<b>1</b>

→ “10% of reads simulated from genome 3 were assigned to a set containing G<sub>1</sub>”

# “Redistributing” multiple assignments

$$m = S \cdot r \quad \Rightarrow \quad \begin{cases} m_1 = r_1 + s_{12}r_2 + s_{13}r_3 \\ m_2 = s_{21}r_1 + r_2 + s_{23}r_3 \\ m_3 = s_{31}r_1 + s_{32}r_2 + r_3 \end{cases}$$

<b>S</b>	G <sub>1</sub>	G <sub>2</sub>	G <sub>3</sub>
G <sub>1</sub>	<b>1</b>	0.1	<b>0.1</b>
G <sub>2</sub>	0.2	<b>1</b>	0.2
G <sub>3</sub>	0	0.1	<b>1</b>

●

<b>r</b>
<b>?</b>
<b>?</b>
<b>?</b>

=

<b>m</b>
11
4
10

$$r = \arg \min_r \|m - S \cdot r\|^2$$



# “Redistributing” multiple assignments

$$m = S \cdot r \quad \Rightarrow \quad \begin{cases} m_1 = r_1 + s_{12}r_2 + s_{13}r_3 \\ m_2 = s_{21}r_1 + r_2 + s_{23}r_3 \\ m_3 = s_{31}r_1 + s_{32}r_2 + r_3 \end{cases}$$

<b>S</b>	G <sub>1</sub>	G <sub>2</sub>	G <sub>3</sub>				<b>r</b>				<b>m</b>
G <sub>1</sub>	<b>1</b>	0.1	<b>0.1</b>	●		=	10				11
G <sub>2</sub>	0.2	<b>1</b>	0.2				0				4
G <sub>3</sub>	0	0.1	<b>1</b>				10				10

Same thing in  
[**#ref.gen.**] dimensions

# Need regularization

- Many assignments are inaccurate
  - sequencing errors
  - something not in DB
  - $k$ -mers only heuristic
- $10^9$  reads +  $10^5$  ref. Genomes → anything could get few assignments
- Real scenario: 100s of known organisms in a sample → 1/100 of index
- Approximate the system while keeping results sparse
- Introduce penalties for “using too many variables”


## LASSO regressor



$$r = \arg \min_r (\|m - S \cdot r\|^2 + \lambda_1 \|r\|_1)$$

# Elastic Net

## (Zou & Hastie, 2005)

- $L_2$  regularization does not set regr. coeff. to 0
- Lasso ( $L_1$ ) → among highly correlated variables will choose only 1  
(at **random**)
-  Combine  $L_1$  and  $L_2$ !
- The quadratic part of the penalty:
  - encourages **grouping** effect
  - **stabilize** regularization path (selected variables)

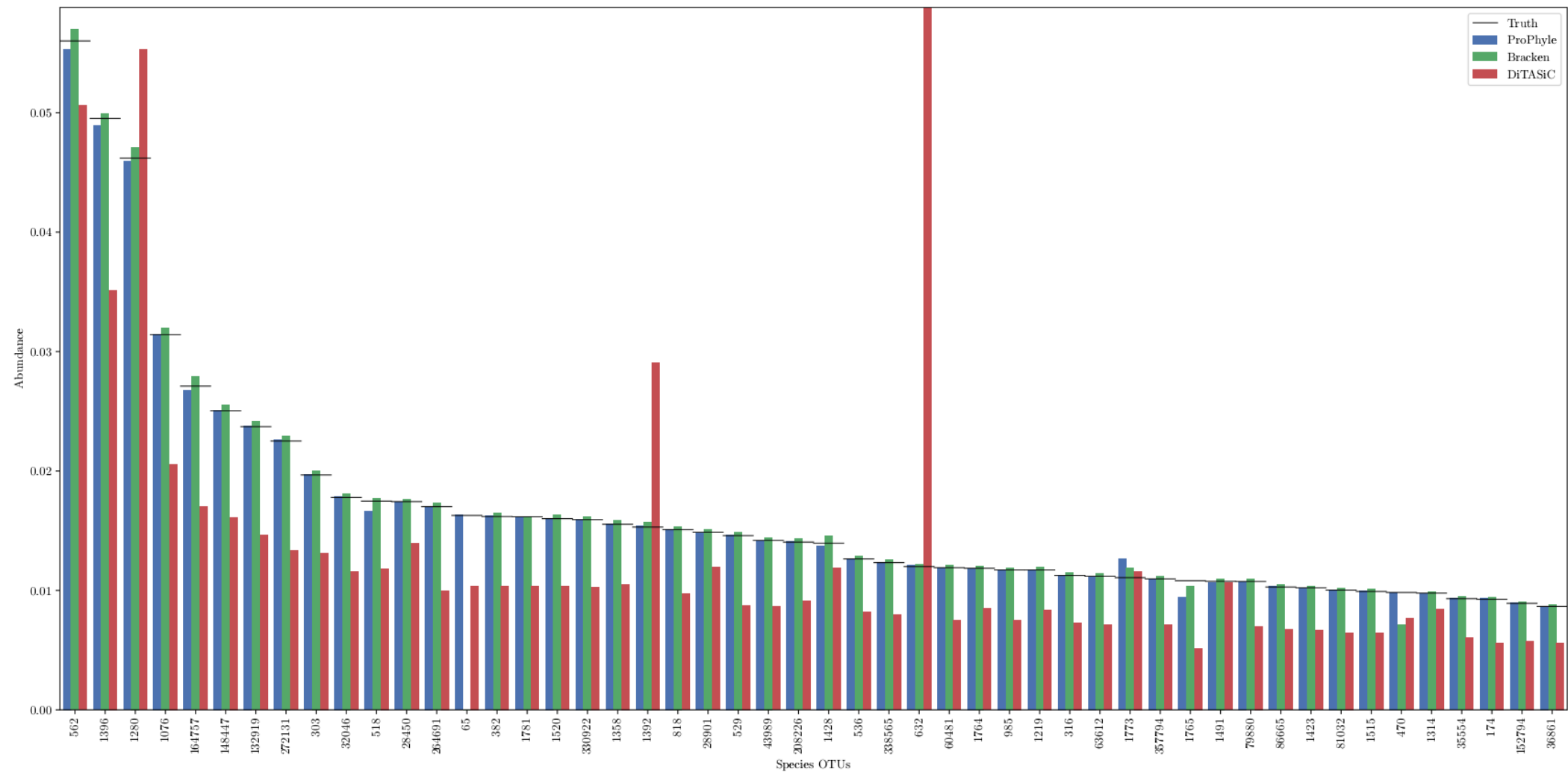
$$r = \arg \min_r (\|m - S \cdot r\|^2 + \lambda_1 \|r\|_1 + \lambda_2 \|r\|_2^2)$$

Can also simulate LASSO, which has desirable effects

# Simulation experiment (Mende et al, 2012)

- ~25M 75bp reads simulated from 193 genomes (85 species)
- includes multiple strains in the genera Bacillus and Mycobacterium
  - notoriously difficult to estimate ab. for (due to high similarity)
- Index for **ProPhyle**, **Bracken** and **DiTASiC** containing 1267 genomes
  - including simulated ones
  - $k = 31 \rightarrow 45$   $k$ -mers per read

- Residual Sum of Squares (RSS) error measure: 
$$\sum_{i=1}^n \varepsilon^2 = \sum_{i=1}^n (t_i - e_i)^2$$



# Results - Mende

		LASSO	E-Net	Bracken	DiTASiC
Genome	RSS	<b>7.52e-4</b>	9.32e-4	NA	2.26e-02
	FP# (ab)	<b>32 (5.60e-02)</b>	89 (8.57e-2)	NA	648 (3.22e-01)
	FN# (ab)	66 (1.06e-02)	<b>2 (2.02e-4)</b>	NA	32 (2.07e-02)
Species	RSS	3.84e-05	<b>6.68e-6</b>	2.80e-04	4.53e-02
	FP# (ab)	<b>1 (1.44e-04)</b>	26 (1.36e-3)	170 (9.73e-04)	381 (6.08e-02)
	FN# (ab)	<b>0 (0)</b>	0 (0)	1 (1.63e-02)	0 (0)
Genus	RSS	1.33e-05	<b>1.16e-06</b>	2.85e-04	5.45e-02
	FP# (ab)	<b>0 (0)</b>	12 (3.62e-5)	78 (3.23e-04)	194 (1.06e-02)
	FN# (ab)	<b>0 (0)</b>	0 (0)	1 (1.63e-02)	0 (0)

# “Stress Test”

- 10 samples with 500k read pairs of length 150bp
- Simulated from ref (250 genomes each, exponential distribution)
- Rare genomes ( $\sim 50$ ) have  $10^{-5}$  abundance  $\rightarrow$  5 reads
- Index for all (6171) “reference” and “representative” RefSeq genomes (Archaea, Bacteria, Fungi, Viruses)
- ProPhyle index:  $\sim 50\text{GB}$  **vs** Kraken DB:  $\sim 330\text{GB}$
- Kraken running for weeks on a powerful cluster

# Results – Stress

- OK scalability (~5 days to generate simulation matrix with 0.3 cov)
- As for Mende, can effectively use LASSO for FP  $\rightarrow 0$
- Grid-Search for optimal parameters in few minutes (needs ground truth)
- Many FN as intended (genomes with  $<10$  reads)

## **ISSUES:**

- Parameters choice extremely important
- Most samples fit perfectly (Pearson  $>0.95$ ), 2 of them **don't** ( $<0.4$ )
- Need more statistical info (e.g.  $P$ -values or other confidence est.)
- Cannot scale to index with  $\sim 100k$  ref



# Real – HMP pilot

- 6.5M reads of length 75bp
- Even mixture of DNA from 22 isolates
- Same index as *stress* (6171 RefSeq genomes)
- 6 out of 22 isolates only have relatives in the index
- 1 fungus with 18S gene only (assignments suggest low copy number)

# Results - HMP

		LASSO
Genome	RSS	0.15
	FP# (ab)	5 (0.01)
	FN# (ab)	3 (0.14)
Species	RSS	0.15
	FP# (ab)	3 (6.30e-03)
	FN# (ab)	1 (0.05)
Genus	RSS	0.15
	FP# (ab)	1 (5.66e-03)
	FN# (ab)	0 (0)

# Results - HMP

Illumina sequencing of HMP Mock Community even sample (SRR172902)

Metadata Analysis (alpha) Reads Download

**Warning:** experimental software

## Taxonomy Analysis

Unidentified reads: **28.13%**

Identified reads: **71.87%**

cellular organisms: **71.87%**

Bacteria: **71.37%**

Terrabacteria group: **41.44%**

Deinococcus-Thermus: **23.05%**

Deinococcus: **23.04%**

Deinococcus radiodurans

Deinococcus wulumuqi

Firmicutes: **13.51%**

Bacilli: **10.84%**

Bacillales: **5.3%**

Lactobacillales: **4.49%**

Clostridia: **2.59%**

Actinobacteria: **4.3%**

Proteobacteria: **17.69%**

Gammaproteobacteria: **10.0%**

delta/epsilon subdivisions: **1.78%**

Alphaproteobacteria: **2.07%**

Betaproteobacteria: **1.78%**

FCB group: **6.82%**

Bacteroidales: **6.81%**

Bacteroides: **6.4%**

Bacteroides vulgatus: **1.0%**

Archaea: **0.5%**

Eukaryota: **0.01%**

## Strong signals

SuperKingdom	Organism	Rank	%%	Kbp	weighted score
Bacteria	Deinococcus radiodurans	species	25.3	124,506	<b>38.2</b>
Bacteria	Acinetobacter baumannii	species	11.3	55,472	<b>13.9</b>
Bacteria	Bacteroides vulgatus	species	7.4	36,299	<b>7.1</b>
Bacteria	Propionibacterium	genus	3.6	17,558	17.6
Bacteria	Staphylococcus aureus	species	3.5	17,334	<b>6.1</b>
Bacteria	Streptococcus mutans	species	2.9	14,255	<b>7.3</b>
Bacteria	Clostridium beijerinckii	species	2.9	14,086	<b>2.4</b>
Bacteria	Helicobacter pylori	species	2.6	13,027	<b>8.0</b>
Bacteria	Rhodobacter sphaeroides	species	2.4	11,687	<b>2.5</b>
Bacteria	Streptococcus pneumoniae	species	2.1	10,342	<b>5.0</b>
Bacteria	Neisseria meningitidis	species	2.0	10,052	<b>4.7</b>
Bacteria	Listeria monocytogenes	species	1.8	9,057	<b>3.0</b>
Bacteria	Actinomyces odontolyticus ATCC 17982		1.1	5,646	5.6
Bacteria	Bacillus cereus group	species group	0.7	3,404	3.4
Bacteria	Pseudomonas	genus	0.6	3,158	3.2
Archaea	Methanobrevibacter smithii	species	0.5	2,472	<b>1.3</b>
Bacteria	Enterococcus	genus	0.4	2,015	2.0

		LASSO
Genome	RSS	0.15
	FP	5 (0.01)
	FN	3 (0.14)
Species	RSS	0.15
	FP	3 (6.30e-03)
	FN	1 (0.05)
Genus	RSS	0.15
	FP	1 (5.66e-03)
	FN	0 (0)

# Adjusted Results - HMP

		LASSO
Genome	RSS	<b>0.03</b>
	FP	5 (0.01)
	FN	3 ( <b>2.28e-6</b> )
Species	RSS	<b>0.03</b>
	FP	3 (6.30e-03)
	FN	1 ( <b>7.61e-07</b> )
Genus	RSS	<b>0.03</b>
	FP	1 (5.66e-03)
	FN	0 (0)

Suggests problem with:

- reference genomes (e.g. contamination)
- sample preparation (e.g. not so even mix)
- current computational approaches

# Conclusions and perspectives

- **ProPhyle** is a complete, **resource-frugal** and easy-to-use metagenomic classifier
- Lossless index, suitable for inaccurate phylogenetic trees
- Flexibility and feature richness  
(works with any tree, standard bioinformatics formats)

## Future directions:

- Assignment quality and read/ref.genome length heuristics for abundances
- ML framework to estimate optimal reg. parameters based on:
  - complexity of sample
  - # ref. genomes
  - sequencing technology



<http://github.com/karel-brinda/prophyle>



Read the Docs

<http://prophyle.rtdf.io>

**BIOCONDA**

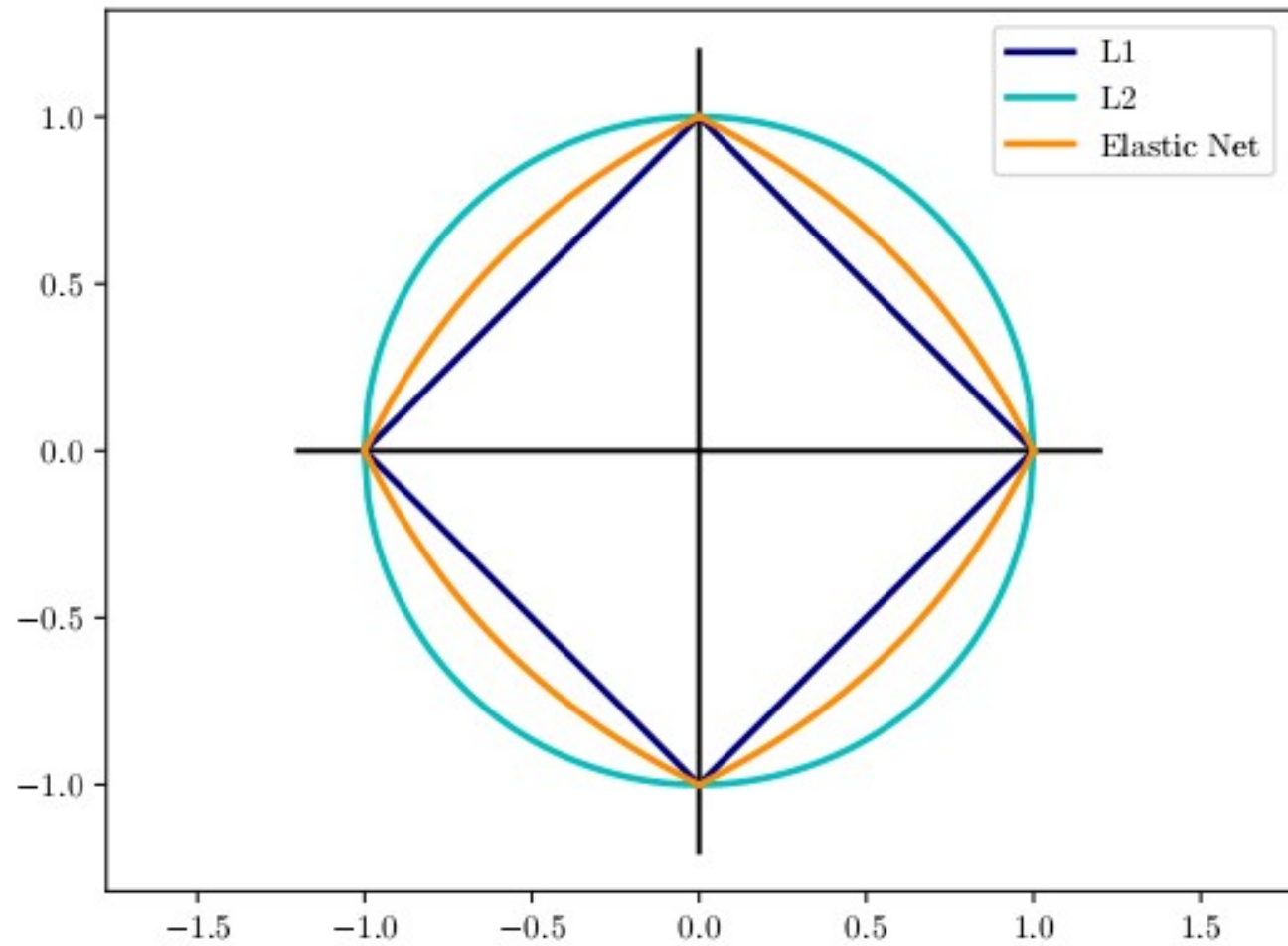
\$ conda install prophyle



\$ pip install prophyle



Thank you!



Lasso and Elastic-Net Paths

