

Mutual Information-based feature selection of phylo-k-mers

Nikolai Romashchenko, Benjamin Linard, Fabio Pardi, Eric Rivals

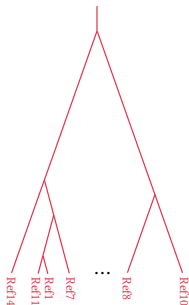
CNRS, University of Montpellier

17 November, 2022

PHYLOGENETIC PLACEMENT

```

Ref1: CTGTCCTATGAGAAGAT...
Ref2: CTTTTATGTGGGGATAAAGTG
Ref3: CTGTTTATACGGGGACGAAACGGC
Ref4: CTGTTGTTAGGGGAGCAAAAA
Ref5: CTTTTTATATGGGGAAGAT...
Ref6: CTATCAGAAAGGGGAAGAA...
Ref7: CTTTGGACAGGGACGAAAGA...
Ref8: CTTTTAGTAGGGGAAAGAA...
Ref9: CTTTTTATATGGGGAAGAT...
Ref10: CTTTTTCTCGTGAAAAAGGCA
Ref11:
Ref12: CTTTTGTTTTGGGGAATAATCGA
Ref13: CTGTGATTTCGGGGACGAAAGGC
Ref14: CTTTTTATATGGGGAAGAT...
Ref15: CTGTGATTTCGGGGACGAAAGAT
Ref16: CTTTAGGCGGGGAGTAAATGTG
Ref17: CTTTTTATATGGGGATAAAGTGT
Ref18: CTTTAGGGGGGAGTAAATGTG
Ref19: CTTTTGTATGGGGATTAAGTGC
Ref20: CTTTTATGCGGGGATAAAGGA
Ref21: CTGTTGTTGGTGAAGAAGGAC
Ref22: CTTTTGTATGGGGAAGAA...
Ref23: CTTTTATGCGGGGATAAAGGA
  
```



```

Q1: TGTTTTGGGAATAATCGA
Q2: CTTTTGTATGAAGAA
Q3: CTTTTATGCGGGG
  
```

The topology of the **reference** tree is given and remains fixed.

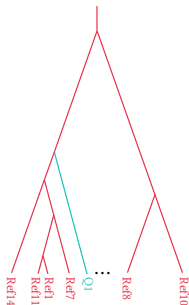
Queries are inserted one-by-one independently.

Possible with millions of queries.

PHYLOGENETIC PLACEMENT

```

Ref1: CTGTCCTATGAGGAAGAT...
Ref2: CTTTTATGTGGGGATAAAGTG
Ref3: CTGTTATACGGGACGAAACGGC
Ref4: CTGTTGTTAGGGAGCAAAAAA
Ref5: CTTTTATATGGGAAGAT...
Ref6: CTATCAGAAAGGGAAAGAA...
Ref7: CTTTGGACAGGGACGAAGA...
Ref8: CTTTTAGTAGGGAAAGAAAG...
Ref9: CTTTTATATGGGAAGAT...
Ref10: CTTTTTCTCGTGAAAAAGGCA
Ref11:
Ref12: CTTTTGTTTGGGAATAATCGA
Ref13: CTGTGATTTCGGGACGAAAGGC
Ref14: CTTTTATATGGGAAGAT...
Ref15: CTGTGATTTCGGGACGAAAGAT
Ref16: CTTTAGGCGGGGAGTAAATGTG
Ref17: CTTTTATATGGGGATAAAGTGG
Ref18: CTTTAGGGGGGAGTAAATGTG
Ref19: CTTTTGTATGGGATTAAGTGC
Ref20: CTTTTATGCGGGGATAAAGGA
Ref21: CTGTTGTTGGTGAAGAAGGAC
Ref22: CTTTTGTATGGGAAGAA...
Ref23: CTTTTATGCGGGGATAAAGGA
  
```



```

Q1: TGGTTGGGAATAATCGA
Q2: CTTTTGTATGAAGAA
Q3: CTTTTATGCGGGG
  
```

The topology of the **reference** tree is given and remains fixed.

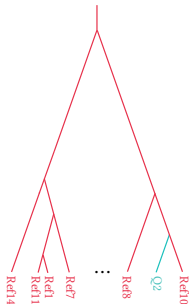
Queries are inserted one-by-one independently.

Possible with millions of queries.

PHYLOGENETIC PLACEMENT

```

Ref1: CTGTCCTATGAGGAAGAT...
Ref2: CTTTTATGTGGGGATAAAAGTG
Ref3: CTGTTATACGGGACGAAACGGC
Ref4: CTGTTGTTAGGGAAAGAAAAA
Ref5: CTTTTATATGGGAAGAT...
Ref6: CTATCAGAAAGGGAAAGAA...
Ref7: CTTTGGACAGGGACGAAGA...
Ref8: CTTTTAGTAGGGAAAGAAAG...
Ref9: CTTTTATATGGGAAGAT...
Ref10: CTTTTTCTCGTGAAAAAGGCA
Ref11:
Ref12: CTTTTGTTTTGGGAATAATCGA
Ref13: CTGTGATTTCGGGACGAAAGGC
Ref14: CTTTTATATGGGAAGAT...
Ref15: CTGTGATTTCGGGACGAAAAGAT
Ref16: CTTTAGGCAGGGAGTAAATGTG
Ref17: CTTTTATATGGGGATAAAAGTG
Ref18: CTTTAGGGGGGAGTAAATGTG
Ref19: CTTTTGTATGGGATTAAGTC
Ref20: CTTTTATGCGGGGATAAAGGA
Ref21: CTGTTGTTGGTGAAGAAGGAC
Ref22: CTTTTGTATGGGAAGAA...
Ref23: CTTTTATGCGGGGATAAAGGA
  
```



```

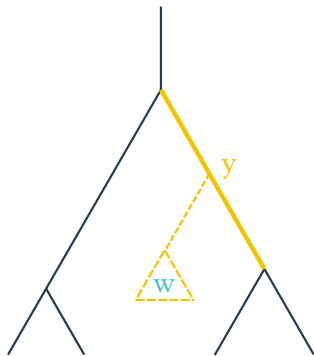
Q1: TGTTTTGGGAATAATCGA
Q2: CTTTTGTATGAAGAA
Q3: CTTTTATGCGGGG
  
```

The topology of the **reference** tree is given and remains fixed.

Queries are inserted one-by-one independently.

Possible with millions of queries.

PHYLO-K-MERS

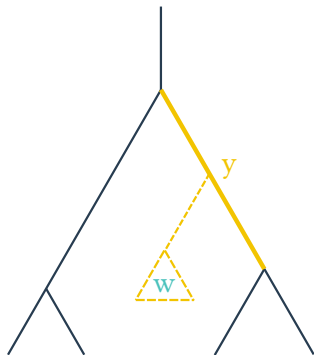


- ▶ w is k -mer (string of size k)
- ▶ y , a branch of the reference tree.

Phylo- k -mer is a triple
 $(w, y, S_y(w))$

$S_y(w)$ estimates the probability of observing w in sequences diverged from y .

PHYLO-K-MERS



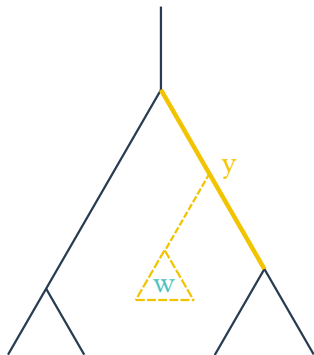
$S_y(w)$ estimates the probability of observing w in sequences diverged from y .

- ▶ w is k -mer (string of size k)
- ▶ y , a branch of the reference tree.

Phylo- k -mer is a triple
 $(w, y, S_y(w))$

k-mer	branch	score
AAA	y_1	0.9

PHYLO-K-MERS



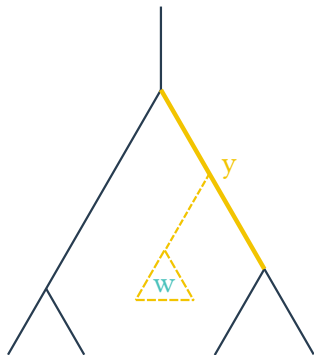
$S_y(w)$ estimates the probability of observing w in sequences diverged from y .

- ▶ w is k -mer (string of size k)
- ▶ y , a branch of the reference tree.

Phylo- k -mer is a triple
 $(w, y, S_y(w))$

k-mer	branch	score
AAA	y_1	0.9
	y_2	0.4

PHYLO-K-MERS



$S_y(w)$ estimates the probability of observing w in sequences diverged from y .

- ▶ w is k -mer (string of size k)
- ▶ y , a branch of the reference tree.

Phylo- k -mer is a triple
 $(w, y, S_y(w))$

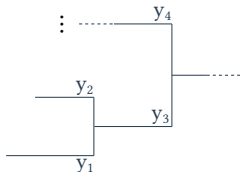
k-mer	branch	score
AAA	y_1	0.9
	y_2	0.4
ATC	y_2	0.3
	...	

RAPPAS2: PLACEMENT ALGORITHM

<i>database</i>		
AAA	Y_1	0.9
	...	
ATG	Y_2	0.4
	Y_3	0.2
	...	
TAT	Y_1	0.7
	Y_3	0.43

query
TATGAG..

Every k -mer of the query is searched in the database.



RAPPAS2: PLACEMENT ALGORITHM

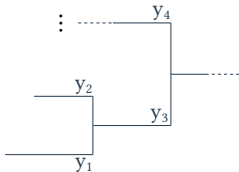
<i>database</i>		
AAA	y_1	0.9
	...	
ATG	y_2	0.4
	y_3	0.2
	...	
TAT	y_1	0.7
	y_3	0.43

query
TATGAG..

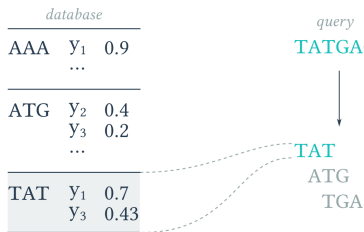
↓

TAT
ATG
TGA

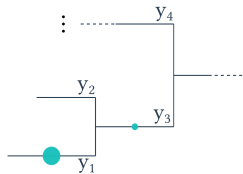
Every k -mer of the query is searched in the database.



RAPPAS2: PLACEMENT ALGORITHM



Every k -mer of the query is searched in the database.



RAPPAS2: PLACEMENT ALGORITHM



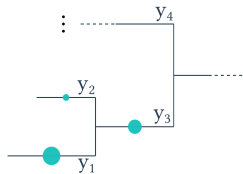
Every k -mer of the query is searched in the database.

Branch y is scored with

$$\prod_w S_y(w)$$

Equivalently:

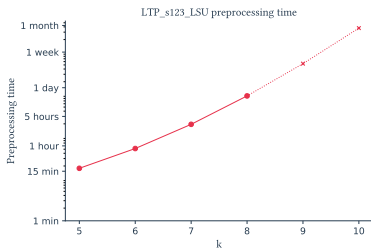
$$\sum_w \log S_y(w)$$



CHALLENGES OF THE METHOD

Archaea / Bacteria 23S rRNA

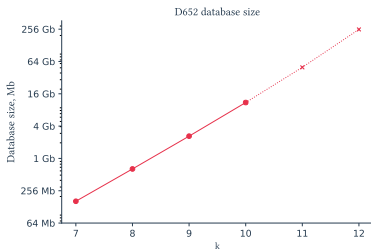
1.7Kbp, $\approx 14K$ taxa



Preprocessing is **long**

Bacteria 16S rRNA (bv)

1.7Kbp, 652 taxa



Databases are **large**

\implies can not be applied to long alignments and large trees.

INFORMATIVE K-MERS: INTUITION

ACAT	
branch	score
y ₁	ϵ
y ₂	ϵ
y ₃	0.87
y ₄	0.44
y ₅	ϵ

informative

Some k -mers are more *informative* for placement than others.

(better at distinguishing between branches)

INFORMATIVE K-MERS: INTUITION

ACAT		GATA	
branch	score	branch	score
y ₁	ϵ	y ₁	0.99
y ₂	ϵ	y ₂	0.99
y ₃	0.87	y ₃	0.99
y ₄	0.44	y ₄	0.99
y ₅	ϵ	y ₅	0.99

informative

not informative

Some k -mers are more *informative* for placement than others.

(better at distinguishing between branches)

Finding informative ones can be done with *feature selection*.

PLACEMENT AS CLASSIFICATION

Placement classifier:

$$f: X \rightarrow Y$$

X is a set of possible query sequences

Y are branches of the reference tree

Note that RAPPAS is a *classifier*, but **NOT** a *machine learning* algorithm

TEXT CLASSIFICATION (EXAMPLE)

Input: sequences of words

A classic example: spam detection.

Output: a label $y \in Y = \{spam, ham\}$

TEXT CLASSIFICATION (EXAMPLE)

Input: sequences of words

A classic example: spam detection.

Output: a label $y \in Y = \{spam, ham\}$

x	$\hat{y}(x)$
SALE -146%! Claim your prize...	spam
Hi guys! The meeting on reducing the number of meetings is...	ham

BERNOULLI NAÏVE BAYES TEXT CLASSIFICATION

query $B_{\mathbf{w}} = \{0, 1, 0 \dots, 0\}$ presence of words

parameters $P_y(\mathbf{w})$ $\mathbb{P}(\text{sale} \mid \text{spam}) = 0.7$

classifier $\prod_{\mathbf{w}: b_{\mathbf{w}}=1} P_y(\mathbf{w}) \prod_{\mathbf{w}: b_{\mathbf{w}}=0} (1 - P_y(\mathbf{w}))$

BERNOULLI NB AND RAPPAS

	Bernoulli NB	RAPPAS
goal	text classification	placement
query	$B_w = \{0, 1, 0, \dots, 0\}$ presence of words	$N_w = \{1, 0, 2, \dots, 0\}$ counts of k -mers
parameters	$P_y(w)$ $\mathbb{P}(\text{sale} \mid \text{spam}) = 0.7$	$S_y(w)$ $\mathbb{P}(\text{ACT} \mid y) = 0.42$
classifier	$\prod_{w: b_w=1} P_y(w) \prod_{w: b_w=0} (1 - P_y(w))$	$\prod_{w: n_w>0} S_y(w)^{n_w}$

BERNOULLI PLACEMENT

	Bernoulli placement	RAPPAS
goal	placement	placement
query	$B_w = \{0, 1, 0 \dots, 0\}$ presence of k -mers	$N_w = \{1, 0, 2, \dots, 0\}$ counts of k -mers
parameters	$P_y(w) := f \cdot S_y(w)$ corrected $S_y(w)$	$S_y(w)$ $\mathbb{P}(\text{ACT} y) = 0.42$
classifier	$\prod_{w: b_w=1} P_y(w) \prod_{w: b_w=0} (1 - P_y(w))$	$\prod_{w: n_w > 0} S_y(w)^{n_w}$

f is the relative query length

THE CONNECTION

$$f \rightarrow 0 \implies$$

Bernoulli placement = RAPPAS.

$f = \frac{|q|-k+1}{m-k+1}$, relative query length

$|q|$ is the query length

m is the alignment length

Interpretation: when placing short queries with long reference alignments.

MAXIMAL MUTUAL INFORMATION FILTER

Mutual Information (MI): how informative one random variable is about another one.

MI filter for RAPPAS maximizes

$$\mathbb{P}(B_w) \cdot \left(H(Y) - H(Y | B_w = 1) \right)$$

class variable Y , presence of k -mers B_w

MAXIMAL MUTUAL INFORMATION FILTER

Mutual Information (MI): how informative one random variable is about another one.

MI filter for RAPPAS maximizes

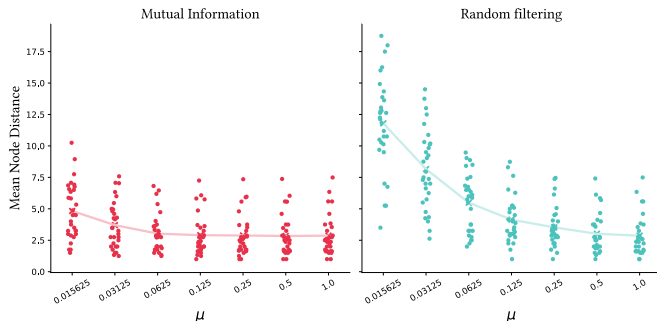
$$\mathbb{P}(B_w) \cdot \left(H(Y) - H(Y | B_w = 1) \right)$$

class variable Y , presence of k -mers B_w

Interpretation:

- ▶ how probable is the k -mer w
- ▶ how informative it is when observed

EXPERIMENTAL RESULTS

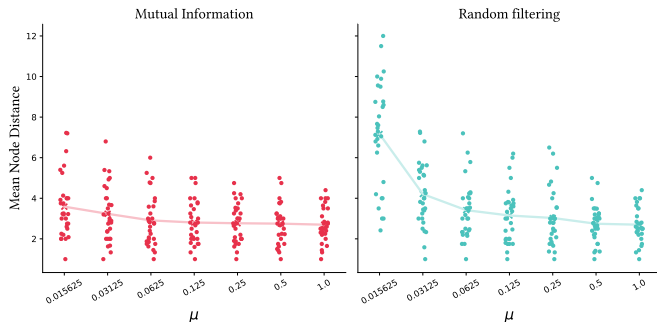


Data: D500 | Chloroplast *rbcL* | 1.4 Kbp | query length 300 bp | $k = 10$

μ is the fraction of phylo- k -mers kept

Interpretation: can keep only 6% of information with a negligible loss in accuracy.

EXPERIMENTAL RESULTS (CONTD.)



Data: D652 (bv) | Bacterial 16S rRNA | 1.7 Kbp | query length 300 bp | $k = 10$

μ is the fraction of phylo- k -mers kept.

Interpretation: can keep 6% of information a negligible loss in accuracy.

CONCLUSION

1. Mutual Information-based feature selection works (on plants, bacteria, viruses)
2. Phylo-*k*-mers work.
3. Filtering will help dealing with placing onto large phylogenies.

Thanks for your attention!