

Constructing phylo k-mers

Phylogenetically-informed k-mers
for phylogenetic placement and recombination detection

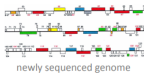
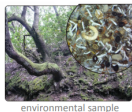
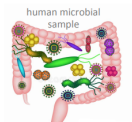
Nikolai Romashchenko



November 24, 2020

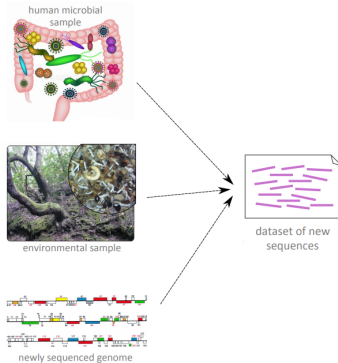
Phylo k-mers: Introduction

Many key bioinformatic tasks: how are new sequences related to those that we already know?



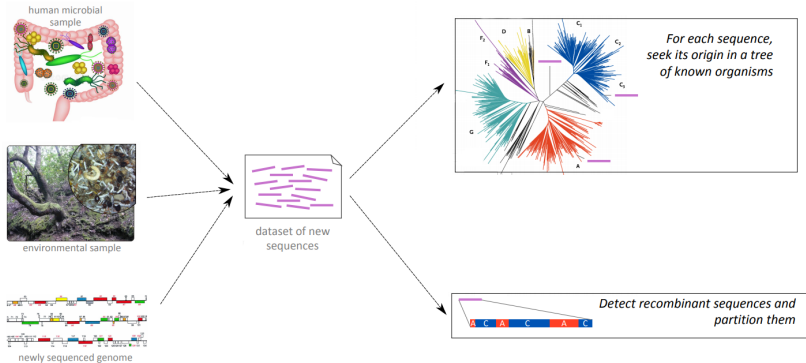
Phylo k-mers: Introduction

Many key bioinformatic tasks: how are new sequences related to those that we already know?



Phylo k-mers: Introduction

Many key bioinformatic tasks: how are new sequences related to those that we already know?

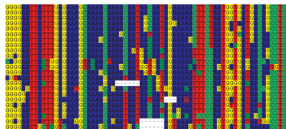


Phylo k-mers: Introduction (contd.)

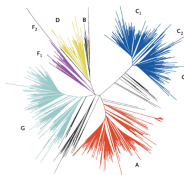
Many key bioinformatic tasks: how are new sequences related to those that we already know?

Alignment-based approaches:

known reference



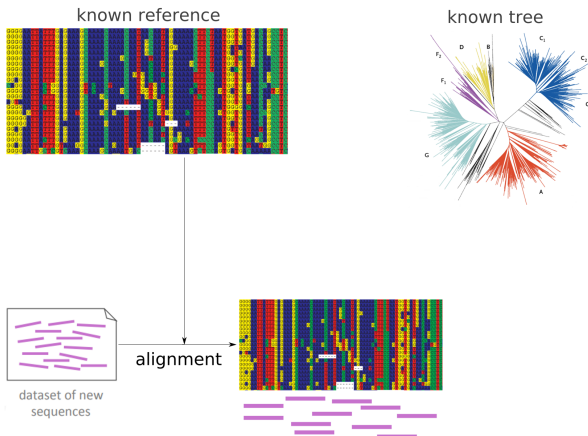
known tree



Phylo k-mers: Introduction (contd.)

Many key bioinformatic tasks: how are new sequences related to those that we already know?

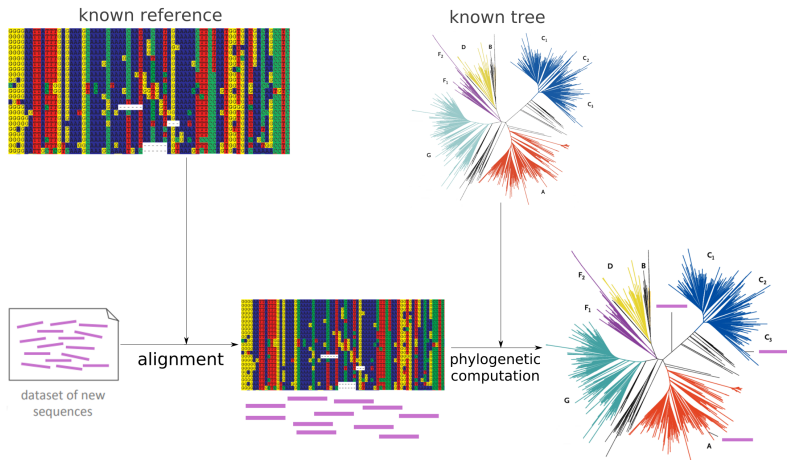
Alignment-based approaches:



Phylo k-mers: Introduction (contd.)

Many key bioinformatic tasks: how are new sequences related to those that we already know?

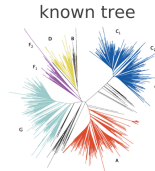
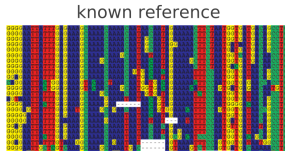
Alignment-based approaches:



Phylo k-mers: Our new approach

Many key bioinformatic tasks: how are new sequences related to those that we already know?

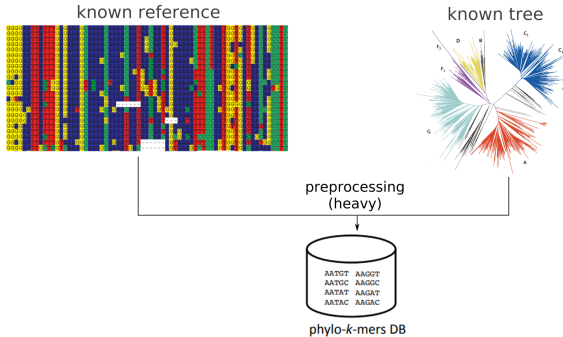
Our new approach:



Phylo k-mers: Our new approach

Many key bioinformatic tasks: how are new sequences related to those that we already know?

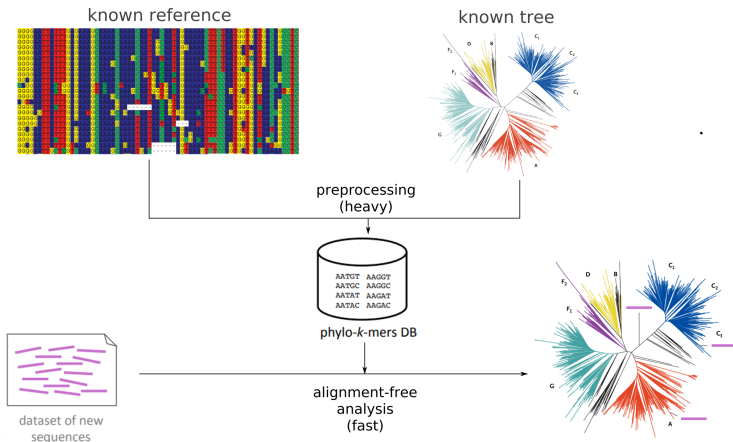
Our new approach:



Phylo k-mers: Our new approach

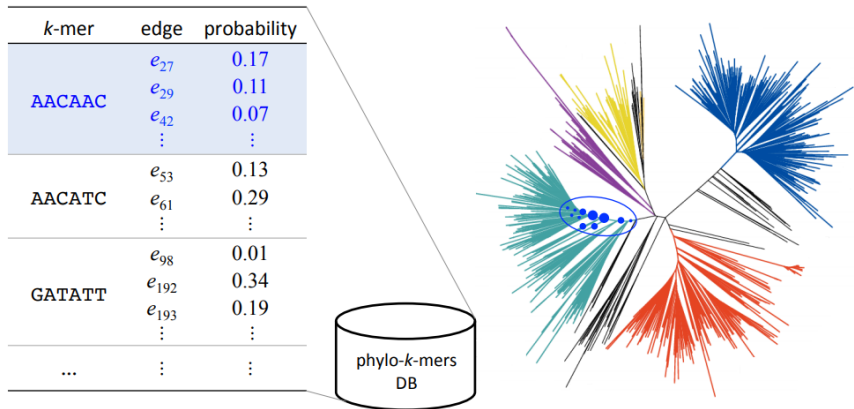
Many key bioinformatic tasks: how are new sequences related to those that we already know?

Our new approach:



Phylo k-mers: Content of databases

"Phylo k -mers" = sequences of length k that are phylogenetically informative: we know where they could have come from, and with what probability.

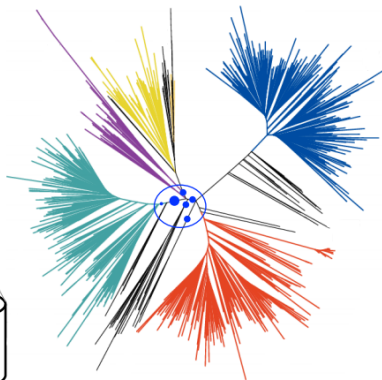
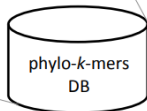


In practice we use longer k -mers ($k = 8, 10, 12 \dots$)

Phylo k-mers: Content of databases

"Phylo k -mers" = sequences of length k that are phylogenetically informative: we know where they could have come from, and with what probability.

k -mer	edge	probability
AACAAC	e_{27}	0.17
	e_{29}	0.11
	e_{42}	0.07
	\vdots	\vdots
AACATC	e_{53}	0.13
	e_{61}	0.29
	\vdots	\vdots
GATATT	e_{98}	0.01
	e_{192}	0.34
	e_{193}	0.19
	\vdots	\vdots
...	\vdots	\vdots

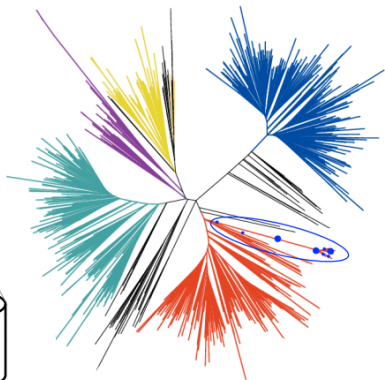
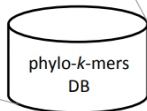


In practice we use longer k -mers ($k = 8, 10, 12 \dots$)

Phylo k -mers: Content of databases

"Phylo k -mers" = sequences of length k that are phylogenetically informative: we know where they could have come from, and with what probability.

k -mer	edge	probability
AACAAC	e_{27}	0.17
	e_{29}	0.11
	e_{42}	0.07
	\vdots	\vdots
AACATC	e_{53}	0.13
	e_{61}	0.29
	\vdots	\vdots
GATATT	e_{98}	0.01
	e_{192}	0.34
	e_{193}	0.19
	\vdots	\vdots
...	\vdots	\vdots



In practice we use longer k -mers ($k = 8, 10, 12 \dots$)

RAPPAS: phylogenetic placement

Rapid alignment-free phylogenetic identification of metagenomic sequences. B. Linard, K. Swenson, F. Pardi. Bioinformatics 2019:
doi.org/10.1093/bioinformatics/btz068

SHERPAS: viral recombination detection

[Accepted] Rapid screening and detection of inter-type viral recombinants using phylo-*k*-mers. G. Scholz, B. Linard, N. Romashchenko, E. Rivals, F. Pardi. Bioinformatics 2020.
Preprint: doi.org/10.1101/2020.06.22.161422

xpas is a library for phylo *k*-mer construction in SHERPAS and the new version of RAPPAS. Github: github.com/phylo42/xpas

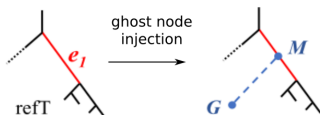
¹Yang, Z. (2007) PAML 4: Phylogenetic Analysis by Maximum Likelihood. Mol. Biol. Evol., 24, 1586–1591.

²Guindon, S. et al. (2010) New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. Systematic Biology 59 (3): 307–21.

³Kozlov, A. et al. (2019) RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. Bioinformatics, 35 (21). 4453–4455.

xpas is a library for phylo k -mer construction in SHERPAS and the new version of RAPPAS. Github: github.com/phylo42/xpas

1. Introducing ghost nodes (RAPPAS, xpas)



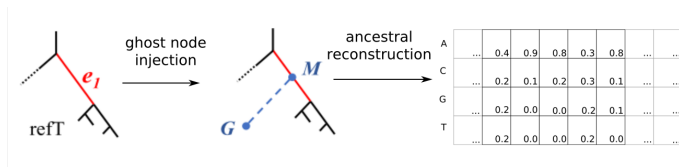
¹Yang, Z. (2007) PAML 4: Phylogenetic Analysis by Maximum Likelihood. Mol. Biol. Evol., 24, 1586–1591.

²Guindon, S. et al. (2010) New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. Systematic Biology 59 (3): 307–21.

³Kozlov, A. et al. (2019) RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. Bioinformatics, 35 (21). 4453–4455.

xpas is a library for phylo k -mer construction in SHERPAS and the new version of RAPPAS. Github: github.com/phylo42/xpas

1. Introducing ghost nodes (RAPPAS, xpas)
2. Ancestral reconstruction (PAML¹, PhyML², RAxML-NG³)



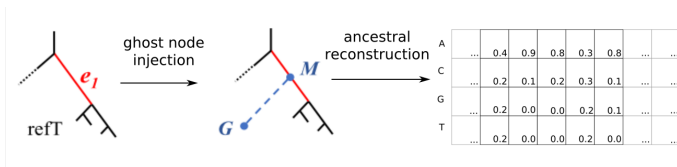
¹Yang, Z. (2007) PAML 4: Phylogenetic Analysis by Maximum Likelihood. Mol. Biol. Evol., 24, 1586–1591.

²Guindon, S. et al. (2010) New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. Systematic Biology 59 (3): 307–21.

³Kozlov, A. et al. (2019) RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. Bioinformatics, 35 (21). 4453–4455.

xpas is a library for phylo k -mer construction in SHERPAS and the new version of RAPPAS. Github: github.com/phylo42/xpas

1. Introducing ghost nodes (RAPPAS, xpas)
2. Ancestral reconstruction (PAML¹, PhyML², RAxML-NG³)
3. Calculating of phylo k -mers (RAPPAS, xpas)



¹Yang, Z. (2007) PAML 4: Phylogenetic Analysis by Maximum Likelihood. Mol. Biol. Evol., 24, 1586–1591.

²Guindon, S. et al. (2010) New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. Systematic Biology 59 (3): 307–21.

³Kozlov, A. et al. (2019) RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. Bioinformatics, 35 (21). 4453–4455.

Phylo k-mer construction: Intuition

We can straightforwardly compute scores of k -mers for a given position of the alignment, by just multiplying corresponding values.

		i			$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

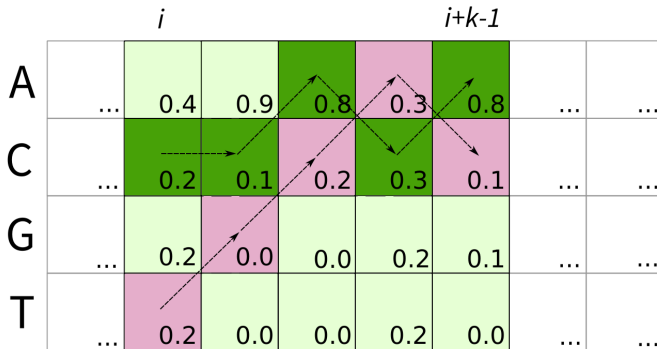
Phylo k-mer construction: Intuition

We can straightforwardly compute scores of k -mers for a given position of the alignment, by just multiplying corresponding values.

		i			$i+k-1$			
A C G T	...	0.4	0.9	0.8	0.3	0.8
	...	0.2	0.1	0.2	0.3	0.1
	...	0.2	0.0	0.0	0.2	0.1
	...	0.2	0.0	0.0	0.2	0.0

Phylo k-mer construction: Intuition

We can straightforwardly compute scores of k -mers for a given position of the alignment, by just multiplying corresponding values.



Phylo k-mer construction: Intuition

We can straightforwardly compute scores of k -mers for a given position of the alignment, by just multiplying corresponding values.

		i			$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

Problem

There are 4^k different combinations, and it must be done for every window of size k across the alignment.

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: A

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i		$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

Current l -mer: AA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i			$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i					$i+k-1$	
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAAAAA $\xrightarrow{\text{save}}$ DB

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: A

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i		$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

Current l -mer: AAAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i			$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAAAC $\xrightarrow{\text{save}}$ DB

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: A

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i		$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

Current l -mer: AA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i			$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

Current l -mer: AAAAG $\xrightarrow{\text{save}}$ DB

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: A

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i		$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

Current l -mer: AA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i			$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

Current l -mer: AAAA

Phylo k-mer construction: Naive algorithm

In the naive solution, we calculate scores of k -mers in the lexicographic order.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

Current l -mer: AAAAT $\xrightarrow{\text{save}}$ DB

That must be frustratingly long.

That must be frustratingly long.
Can we do any better?

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i			$i+k-1$		
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

1.	l -mer	score
	A	0.4

Current l -mer: A

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i		$i+k-1$				
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

- | | l -mer | score |
|----|----------|-------|
| 1. | A | 0.4 |
| 2. | AA | 0.36 |

Current l -mer: AA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i			$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288

Current l -mer: AAA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864

Current l -mer: AAAA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i		$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	...
C	...	0.2	0.1	0.2	0.3	0.1	...
G	...	0.2	0.0	0.0	0.2	0.1	...
T	...	0.2	0.0	0.0	0.2	0.0	...

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864
5.	AAAAA	0.06912

Current l -mer: AAAAAA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864

Current l -mer: AAAA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864
5.	AAAAC	0.00864

Current l -mer: AAAAC

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864

Current l -mer: AAAA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864
5.	AAAAG	0.00864

Current l -mer: AAAAG

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864

Current l -mer: AAAA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i				$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8	
C	...	0.2	0.1	0.2	0.3	0.1	
G	...	0.2	0.0	0.0	0.2	0.1	
T	...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAA	0.0864
5.	AAAAT	0.0

Current l -mer: AAAAT

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i			$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288

Current l -mer: AAA

Phylo k-mer construction: Dynamic algorithm

Let's save the current l -mer's score in a stack.

		i			$i+k-1$			
A	...	0.4	0.9	0.8	0.3	0.8
C	...	0.2	0.1	0.2	0.3	0.1
G	...	0.2	0.0	0.0	0.2	0.1
T	...	0.2	0.0	0.0	0.2	0.0

	l -mer	score
1.	A	0.4
2.	AA	0.36
3.	AAA	0.288
4.	AAAC	0.0864

Current l -mer: AAAC

In practice, we are interested in phylo k -mers with high final scores. The vast majority of phylo k -mers have low scores.

In practice, we are interested in phylo k -mers with high final scores. The vast majority of phylo k -mers have low scores.

Important optimization: do not store values less than a certain *threshold* value.

$$\text{score}(w) = \max(\text{score}(w), \text{threshold})$$

threshold becomes the default value of phylo k -mers in the database.

In practice, we are interested in phylo k -mers with high final scores. The vast majority of phylo k -mers have low scores.

Important optimization: do not store values less than a certain *threshold* value.

$$\text{score}(w) = \max(\text{score}(w), \text{threshold})$$

threshold becomes the default value of phylo k -mers in the database.

Benefits

- **Significantly** less values to store in the database
- Allows to give up with calculating whole sets of k -mers

Phylo k-mer construction: Branch-and-bound algorithm

Let's sort every column by score value.

	<i>i</i>		<i>i+k-1</i>				
	A	A	A	A	C		
...	0.4	0.9	0.8	0.3	0.8
	C	C	C	G	A		
...	0.2	0.1	0.2	0.3	0.1
	G	G	G	C	G		
...	0.2	0.0	0.0	0.2	0.1
	T	T	T	T	T		
...	0.2	0.0	0.0	0.2	0.0

l-mer score

Phylo k-mer construction: Branch-and-bound algorithm

Let's sort every column by score value.

The order of k -mers is now different, but we can cut out whole "branches" of k -mers if the score of a prefix is low enough.

	i		$i+k-1$					
	A	A	A	A	C			
...	0.4	0.9	0.8	0.3	0.8	
	C	C	C	G	A			
...	0.2	0.1	0.2	0.3	0.1	
	G	G	G	C	G			
...	0.2	0.0	0.0	0.2	0.1	
	T	T	T	T	T			
...	0.2	0.0	0.0	0.2	0.0	

	l -mer	score
1.	C	0.2
2.	CA	0.18
3.	CAC	0.036

Phylo k-mer construction: Branch-and-bound algorithm

Let's sort every column by score value.

The order of k -mers is now different, but we can cut out whole "branches" of k -mers if the score of a prefix is low enough.

	i		$i+k-1$				
	A	A	A	A	C		
...	0.4	0.9	0.8	0.3	0.8
	C	C	C	G	A		
...	0.2	0.1	0.2	0.3	0.1
	G	G	G	C	G		
...	0.2	0.0	0.0	0.2	0.1
	T	T	T	T	T		
...	0.2	0.0	0.0	0.2	0.0

	l -mer	score
1.	C	0.2
2.	CA	0.18
3.	CAC	0.036

$$s(\text{CAC}) = 0.036 \leq \text{thr}$$

Phylo k-mer construction: Branch-and-bound algorithm

Let's sort every column by score value.

The order of k -mers is now different, but we can cut out whole "branches" of k -mers if the score of a prefix is low enough.

	i		$i+k-1$				
	A	A	A	A	C		
...	0.4	0.9	0.8	0.3	0.8
	C	C	C	G	A		
...	0.2	0.1	0.2	0.3	0.1
	G	G	G	C	G		
...	0.2	0.0	0.0	0.2	0.1
	T	T	T	T	T		
...	0.2	0.0	0.0	0.2	0.0

	l -mer	score
1.	C	0.2
2.	CA	0.18
3.	CAC	0.036

$$s(\text{CAC}) = 0.036 \leq \text{thr}$$

We still recalculate the same suffixes.

Phylo k-mer construction: Divide-and-conquer algorithm

The branch-and-bound algorithm reuses prefix scores, but for two k -mers with the same suffix, it would calculate the suffix score twice.

To avoid this, we split the window into subwindows: $l = \left\lceil \frac{k}{2} \right\rceil, r = \left\lfloor \frac{k}{2} \right\rfloor$

	i		$i+l-1$		$i+k-1$			
	A	A	A	A	C			
...	0.4	0.9	0.8	0.3	0.8	
	C	C	C	G	A			
...	0.2	0.1	0.2	0.3	0.1	
	G	G	G	C	G			
...	0.2	0.0	0.0	0.2	0.1	
	T	T	T	T	T			
...	0.2	0.0	0.0	0.2	0.0	

Phylo k-mer construction: Divide-and-conquer algorithm

The branch-and-bound algorithm reuses prefix scores, but for two k -mers with the same suffix, it would calculate the suffix score twice.

To avoid this, we split the window into subwindows: $l = \left\lceil \frac{k}{2} \right\rceil, r = \left\lfloor \frac{k}{2} \right\rfloor$

	i			$i+l-1$		$i+k-1$			
	A	A	A	A	C				
...	0.4	0.9	0.8	0.3	0.8		
	C	C	C	G	A				
...	0.2	0.1	0.2	0.3	0.1		
	G	G	G	C	G				
...	0.2	0.0	0.0	0.2	0.1		
	T	T	T	T	T				
...	0.2	0.0	0.0	0.2	0.0		

Prefixes and suffixes are constructed in subwindows recursively, and then merged into phylo k -mers.

Windows of size 1 already list all possible phylo 1-mers.

Phylo k-mer construction: Window merge

l -window contains prefixes, r -window contains suffixes. We can take the cartesian product of scores, but it will produce a lot of k -mers with scores $\leq \text{threshold}$.

prefix	score		suffix	score		k-mer	score
AAA	0.288	×	AC	0.24	=	AAAAC	0.069
CAA	0.144		GC	0.24		AAAGC	0.069
...			...			CAAAC	0.034
TAA	0.144		TC	0.16		CAAGC	0.034
						...	
						TAATC	0.023

Instead, we sort suffixes, and calculate the *residual threshold* $\text{th}(\text{prefix})$ for every prefix. Suffix is taken to the product only if $\text{score}(\text{suffix}) \geq \text{th}(\text{prefix})$.

Phylo k-mer construction: Window merge

l -window contains prefixes, r -window contains suffixes. We can take the cartesian product of scores, but it will produce a lot of k -mers with scores $\leq \text{threshold}$.

prefix	score		suffix	score		k -mer	score
AAA	0.288		AC	0.24		AAAAC	0.069
CAA	0.144	\times	GC	0.24	$=$	AAAGC	0.069
...			...			CAAAC	0.034
TAA	0.144		TC	0.16		CAAGC	0.034
						...	
						TAATC	0.023

Instead, we sort suffixes, and calculate the *residual threshold* $\text{th}(\text{prefix})$ for every prefix. Suffix is taken to the product only if $\text{score}(\text{suffix}) \geq \text{th}(\text{prefix})$.

Phylo k-mer construction: Window merge

l -window contains prefixes, r -window contains suffixes. We can take the cartesian product of scores, but it will produce a lot of k -mers with scores $\leq \text{threshold}$.

prefix	score		suffix	score		k -mer	score
AAA	0.288	×	AC	0.24	=	AAAAC	0.069
CAA	0.144		GC	0.24		AAAGC	0.069
...			...			CAAAC	0.034
TAA	0.144		TC	0.16		CAAGC	0.034
						...	
						TAATC	0.023

Instead, we sort suffixes, and calculate the *residual threshold* $\text{th}(\text{prefix})$ for every prefix. Suffix is taken to the product only if $\text{score}(\text{suffix}) \geq \text{th}(\text{prefix})$.

- RAPPAS2, phylogenetic placement
- CLAPPAS, protein family classification
- Other applications?