

# Multiple Comparative Metagenomics using Multiset k-mer Counting

“Simka”



Pierre Peterlongo



SUMMER SCHOOL 2016 IN METAGENOMICS



# WARNING – This summer school

“  
Its content will focus on the taxonomic assignment and the functional analysis of metatranscriptomic and metagenomic data.

## Summer School 2016 in Metagenomics



**Metagenomics, the sequencing of DNA directly from a sample without first culturing and isolating the organisms, has become the principal tool of “meta-omic” analysis. It can be used to explore the diversity, function, and ecology of microbial communities.**

The aim of these 4 days workshop will be to give researchers and students an overview of the tools and bioinformatics techniques available for the analysis of next generation sequence data from microbial communities. Its content will focus on the taxonomic assignment and the functional analysis of metatranscriptomic and metagenomic data. The format will comprise a mixture of lectures and hands-on practical tutorials where students will process example data sets in real-time.

# WARNING – This summer school

“  
Its content will focus on the ~~taxonomic~~  
~~assignment~~ and the ~~functional analysis~~ of  
metatranscriptomic and metagenomic data.

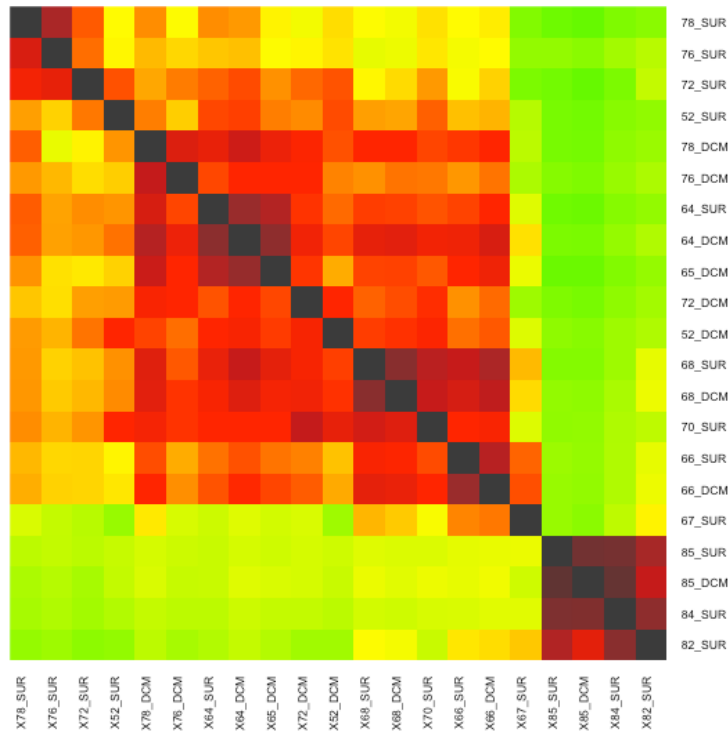
## Summer School 2016 in Metagenomics



**Metagenomics, the sequencing of DNA directly from a sample without first culturing and isolating the organisms, has become the principal tool of “meta-omic” analysis. It can be used to explore the diversity, function, and ecology of microbial communities.**

The aim of these 4 days workshop will be to give researchers and students an overview of the tools and bioinformatics techniques available for the analysis of next generation sequence data from microbial communities. Its content will focus on the taxonomic assignment and the functional analysis of metatranscriptomic and metagenomic data. The format will comprise a mixture of lectures and hands-on practical tutorials where students will process example data sets in real-time.

# “comparing metagenomic samples”



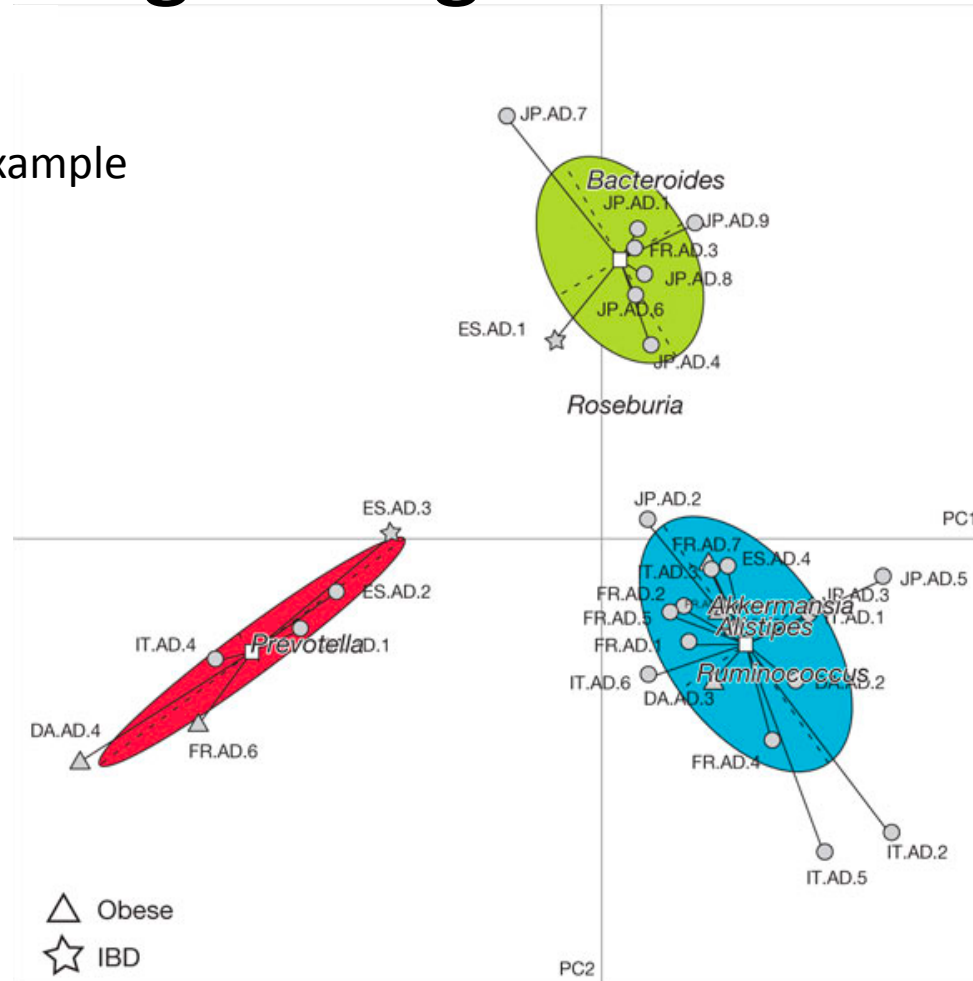
Similarity measure for each couple of datasets

# “comparing metagenomic samples”

Mettre ici des publis: Tara, HMP, New york metro, ...

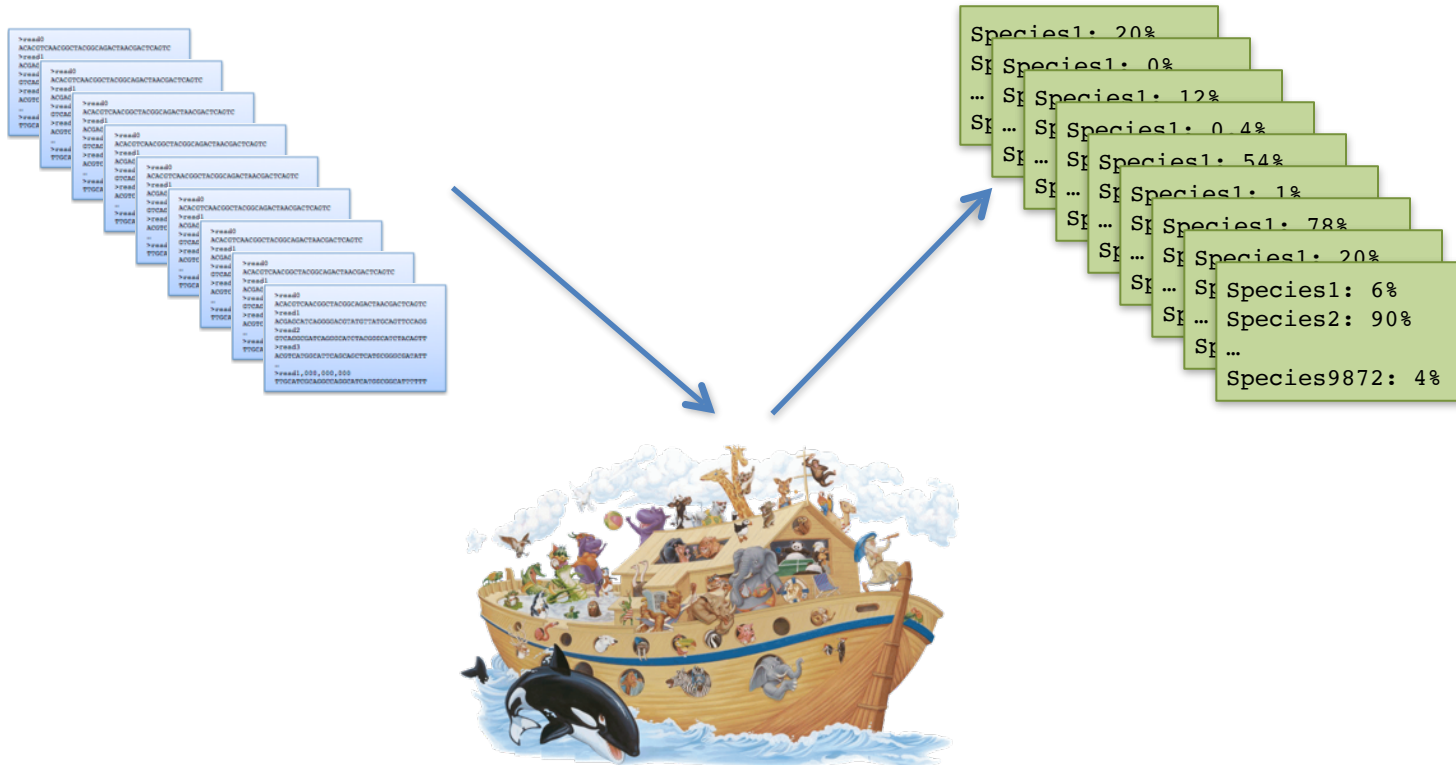
# “comparing metagenomic samples”

Human enterotypes example



# “reference-based comparison of metagenomic samples”

- From reads to taxonomic composition

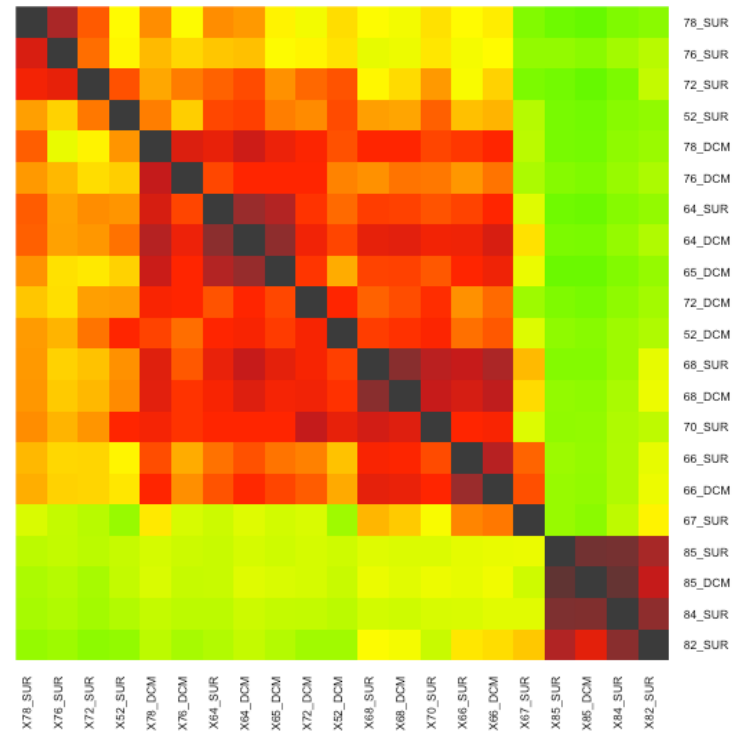
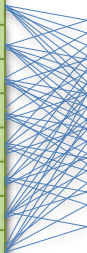


Credit: <http://niniejolie eklablog.com/>

# “reference-based comparison of metagenomic samples”

- Compare taxonomy composition

Species1: 20%
Species1: 0%
Species1: 12%
Species1: 0.4%
Species1: 54%
Species1: 1%
Species1: 78%
Species1: 6%
Species2: 90%
...
Species9872: 4%





# “reference-based comparison of metagenomic samples”

- Reference based limitations

- Databases not representative of diversity

“ we have only sequenced 10-22% of the total DNA on Earth

(Nature Review Microbiol. editorial, 2011)

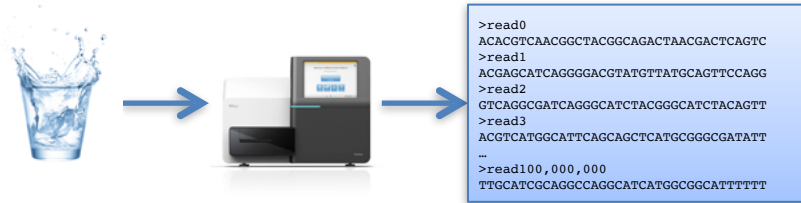


Credit: <http://niniejolie eklablog.com/>

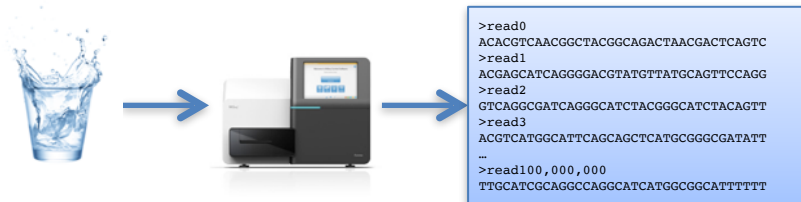
“Tara ocean example”



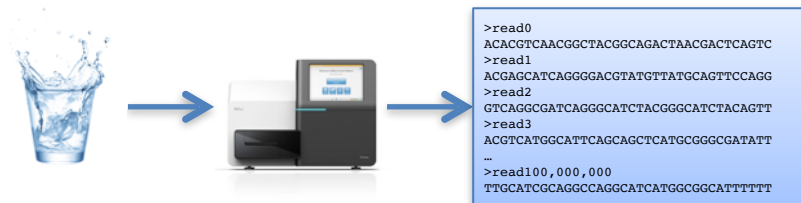
# Tara



- Hundreds of samples
- Billion of reads



...



# Tara



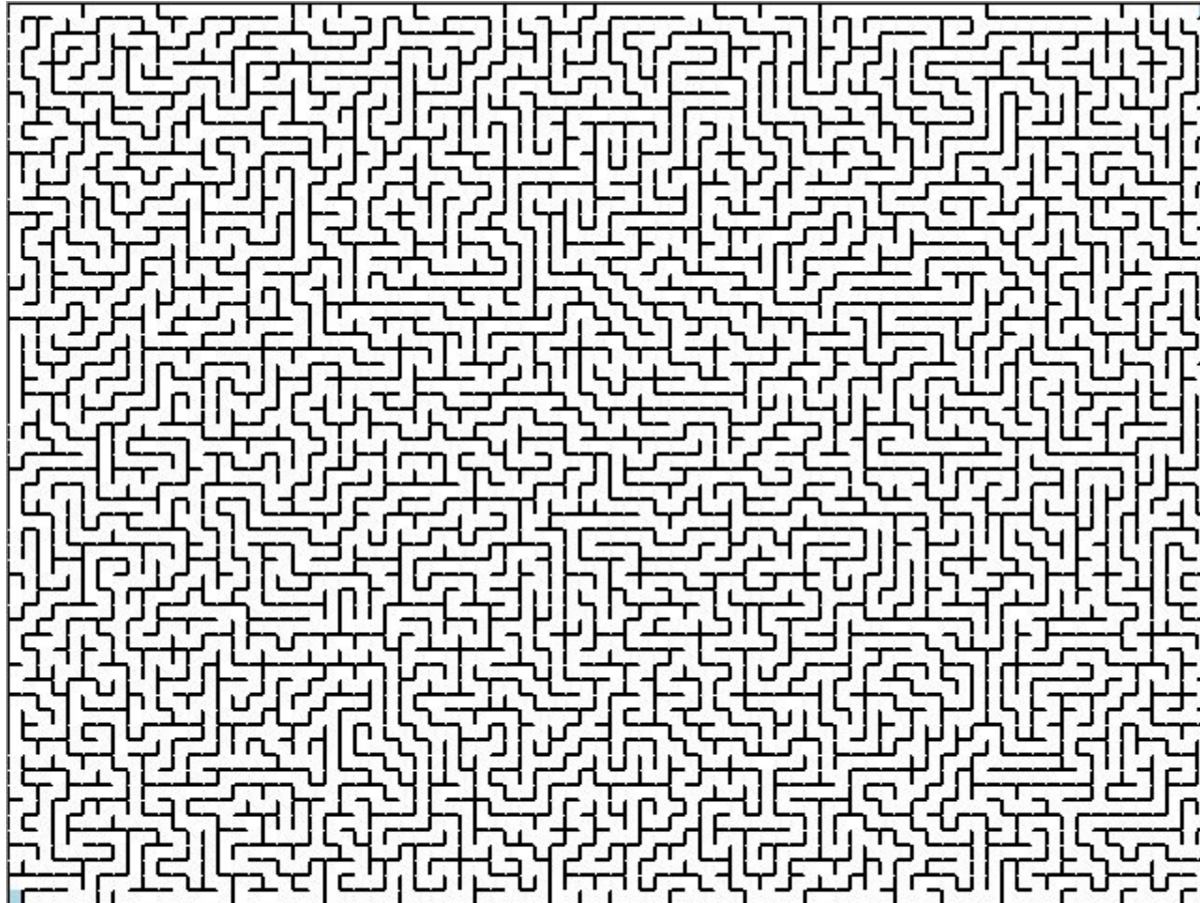
```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTAGCAGCTCATGGGGCGATATT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGCATTTTTT
```

# Assembly Mapping

**Sea water:**

< 5% assembled reads  
< 10% mapped reads

# “de-novo comparison: a complex problem”



Credit: <http://junesblog.org/>



# “de-novo comparison: a complex problem”

- Comparing two reads is **simple**
- Comparing
  - **100** read sets
  - each composed of **100 millions** reads is (terribly) **complex**



# “de-novo comparison: a complex problem”

- Comparing reads is **simple**
- Comparing
  - 100 read sets
  - each composed of 100 millions reads is (terribly) **complex**



$(100 \cdot 10^6)^2 \times 100^2$  Comparisons

```
TACGGGACTGAT-CAGACGTCAA
||||| ||||| || || |||||
ACGG--CTGATTCACTTCAAGG
```



# “de-novo comparison: a complex problem”

- Comparing reads is **simple**
- Comparing
  - 100 read sets
  - each composed of 100 millions reads is (terribly) **complex**



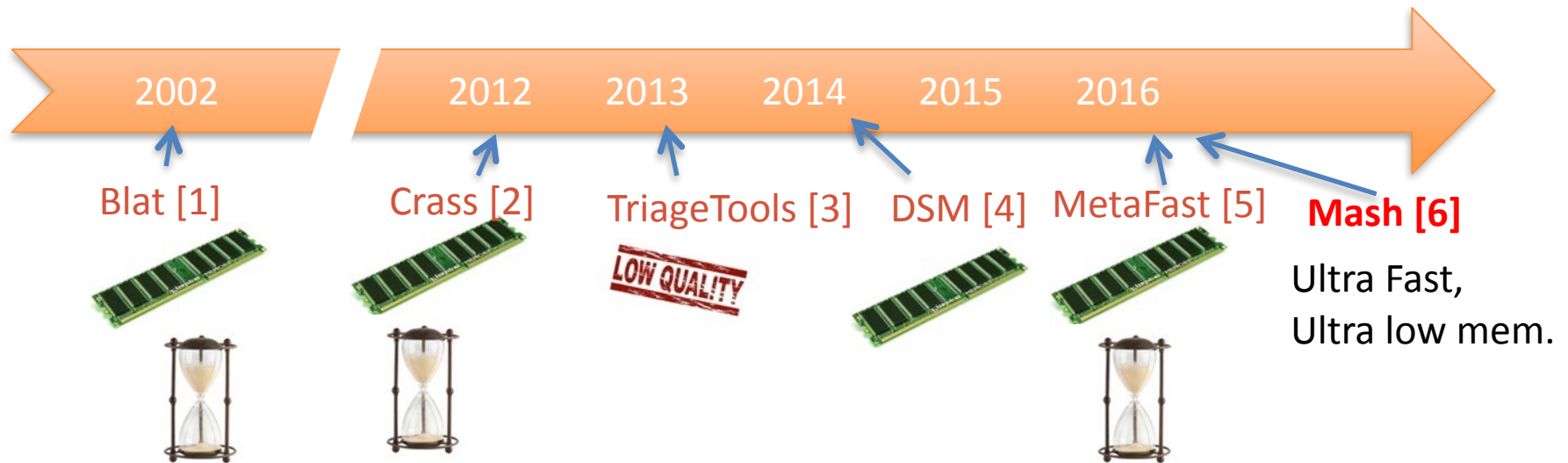
$$(100 \cdot 10^6)^2 \times 100^2 \times 1 \mu\text{sec}$$
$$= 10^{20} \mu\text{sec} = 3 \text{ billions centuries}$$



# What tool for *de novo* comparative metagenomics



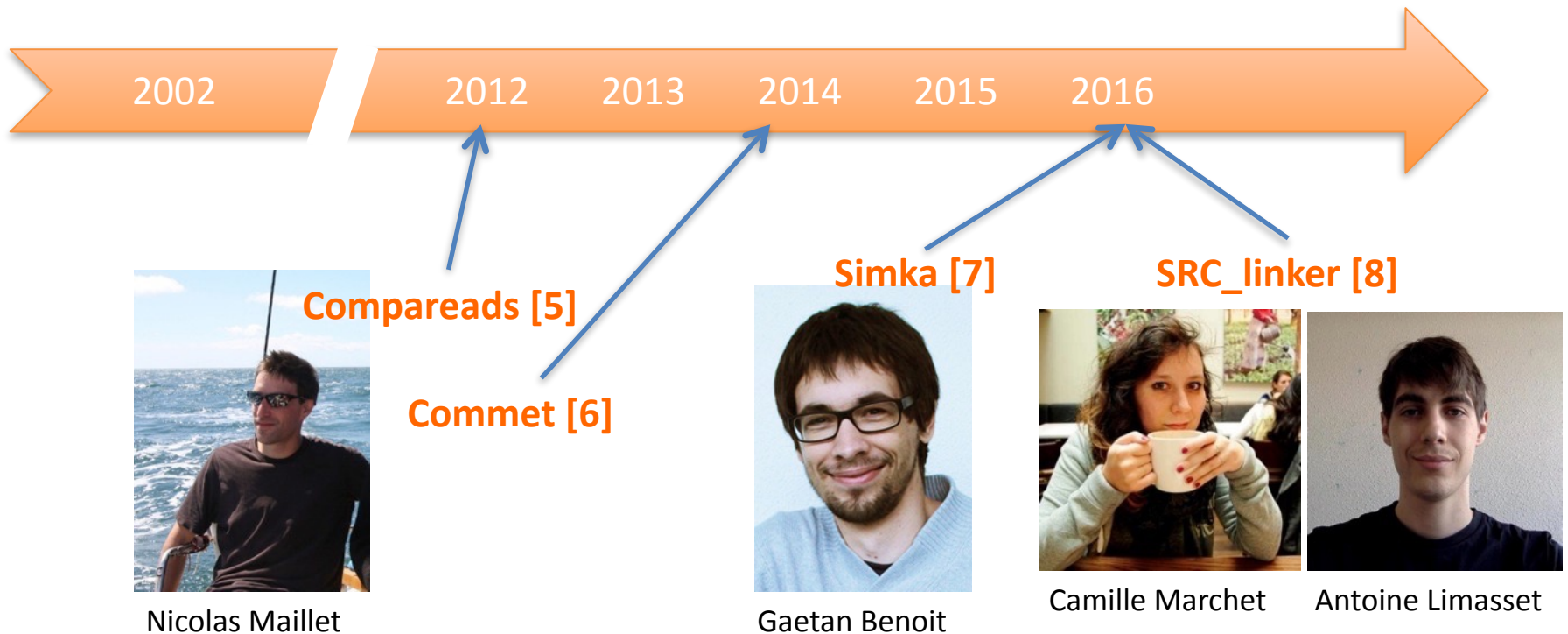
# *de novo* comparative metagenomics- *State of the art*



- [1] Kent 2002  
Computation of pairwise distances
- [2] Dutilh *et. al.* 2012  
Assembly based
- [3] Seth *et. al.* 2014  
Small kmer based

- [4] Fimereli *et al.* 2013  
Small k-mer based
- [5] Ulyantsev *et. al.* 2016  
Unitig & clustering based
- [6] Ondov *et. al* 2016  
Subsampling – 0/1  
Jaccard

# *de novo* comparative metagenomics- our proposals



[5] Compareads [Maillet *et. al.* 2012]

[6] Commet [Maillet *et. al.* 2014]

[7] Simka [Benoit *et al.* 2016]

[8] SRC\_linker [Marchet *et. al.* 2016]

# Main **algorithmic idea** of our tools



# Main idea

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTTCAGCAGCTCATGCGGGCGATATT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGCATTTTTT
```

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTTCAGCAGCTCATGCGGGCGATATT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGCATTTTTT
```

Read to read comparisons = way too long

```
TACGGGACTGAT-CAGACGTCAA
| | | | | | | | | | | | | |
ACGG--CTGATTCATACTTCAAGG
```

→ Similarity = 74%

# Main idea

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTACGAGCTCATGGCGGATATTT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGGATTTTTT
```

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTACGAGCTCATGGCGGATATTT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGGATTTTTT
```

All ideas based on *alignment-free* methods

*k-mer* based:

TACGGGACTGAT-CAGACGTCAA  
| | | | | | | | | | | | | | | |  
ACGG--CTGATTCATACTTCAAGG

→ Similarity = 4 shared *k*-mers:  
ACGG CTGA TGAT TCAA

56% of positions covered  
by a shared kmer

# Main idea

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTTCAGCAGCTCATGGCGGATATTT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGGCATTTTTT
```

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTTCAGCAGCTCATGGCGGATATTT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGGCATTTTTT
```

All ideas based on *alignment-free* methods

*k-mer* based:

TACGGGACTGATCAGACGTCAA  
ACGGCTGATTTCATACTTCAAGG



Similarity = 4 shared *k*-mers:  
ACGG CTGA TGAT TCAA

56% of positions covered  
by a shared kmer



# Main idea

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTTCAGCAGCTCATGGCGCATATTT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGCATTTTTT
```

```
>read0
ACACGTCAACGGCTACGGCAGACTAACGACTCAGTC
>read1
ACGAGCATCAGGGGACGTATGTTATGCAGTCCAGG
>read2
GTCAGGCGATCAGGGCATCTACGGGCATCTACAGTT
>read3
ACGTCATGGCATTTCAGCAGCTCATGGCGCATATTT
...
>read1,000,000,000
TTGCATCGCAGGCCAGGCATCATGGCGCATTTTTT
```

All ideas based on *alignment free* methods

```
TACGGGACTGAT-CAGACGTCAA
| | | | | | | | | | | | | |
ACGG--CTGATTCATACTTCAAGG
```





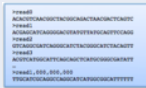

alignment

Similarity = 74%

k-mer based

56% of positions covered by a shared kmer

# Our (alignment free) proposals

		Commet	SRC_linker	Simka
Pro		✓	✗	✓
Scales up	Mem	✓	✗	✓
	Time	✗	✗	✓
		<hr/> VS 	<hr/> VS 	 VS 

# Simka focus



Gaetan Benoit



Sophie Schbath



Dominique  
Lavenier



Claire  
Lemaitre



Mahendra  
Mariadassou

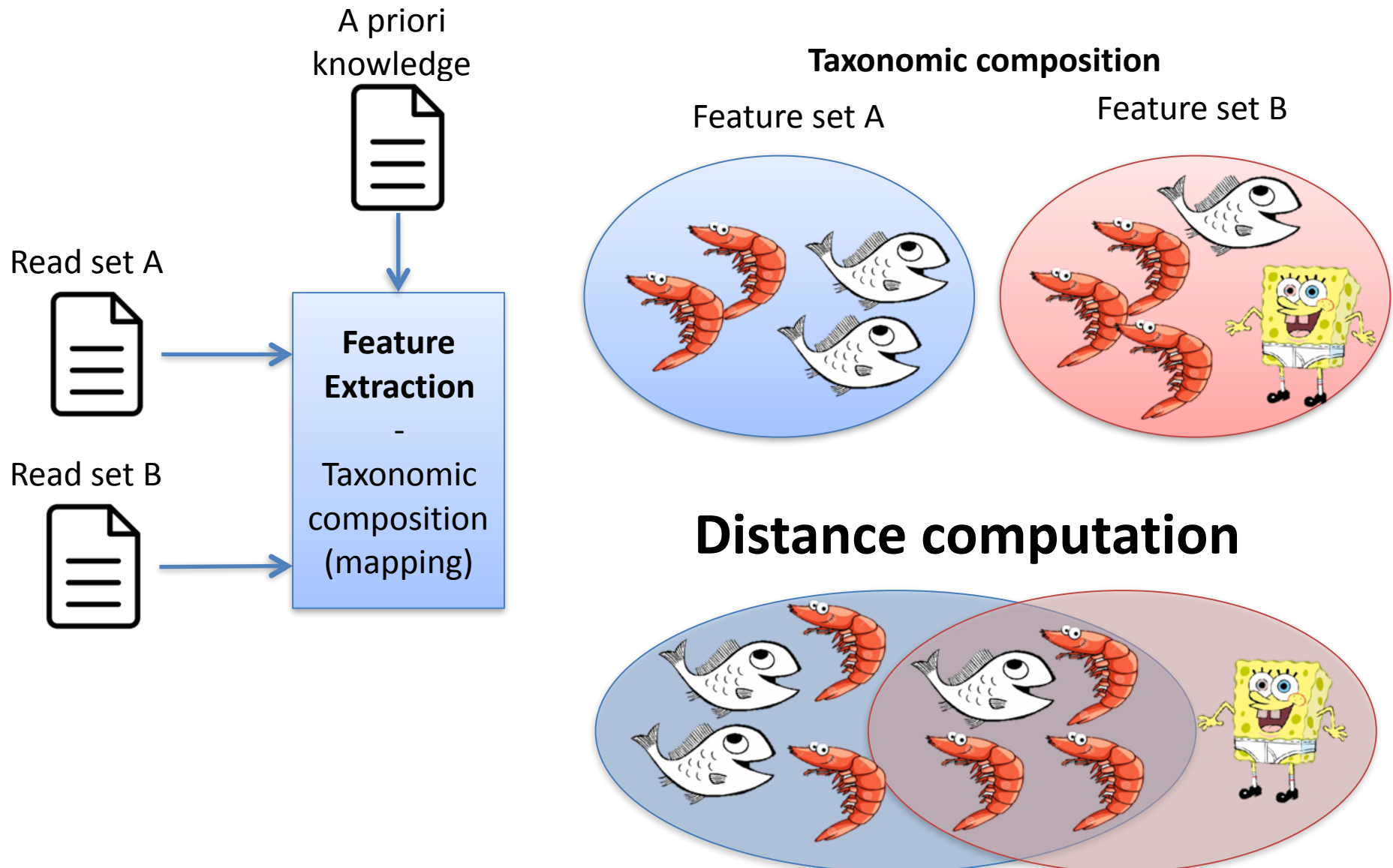


Erwan Drezen

## Multiple Comparative Metagenomics using Multiset $k$ -mer Counting

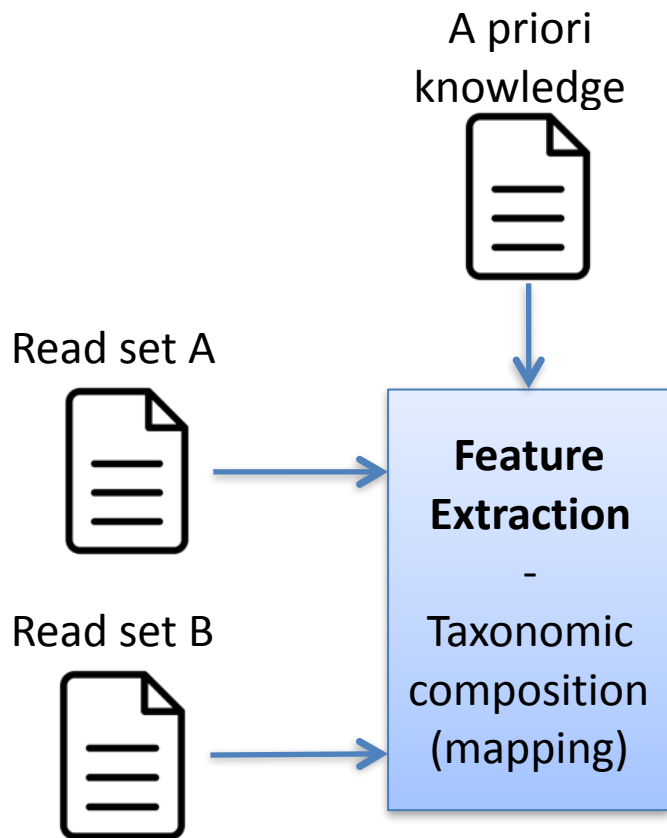
Gaëtan Benoit<sup>1,\*</sup>, Pierre Peterlongo<sup>1</sup>, Mahendra Mariadassou<sup>3</sup>, Erwan Drezen<sup>1,4</sup>, Sophie Schbath<sup>3</sup>,  
Dominique Lavenier<sup>1</sup>, Claire Lemaitre<sup>1</sup>

# Comparing 2 read sets – using a priori knowledge




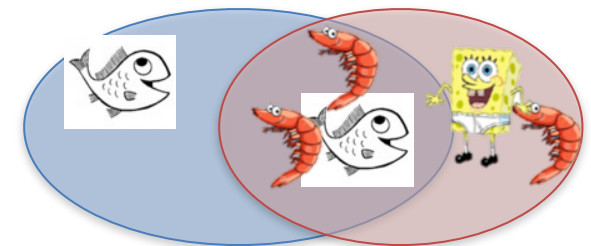
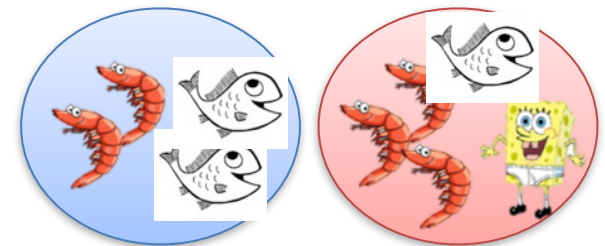
# Comparing 2 read sets – using a priori knowledge

$$Jaccard(A, B) = \underbrace{\sum_{s \in A \cap B} S_A + S_B}_{n \text{ (intersection)}} / \underbrace{\sum_{s \in A \cup B} S_A + S_B}_{U \text{ (union)}}$$



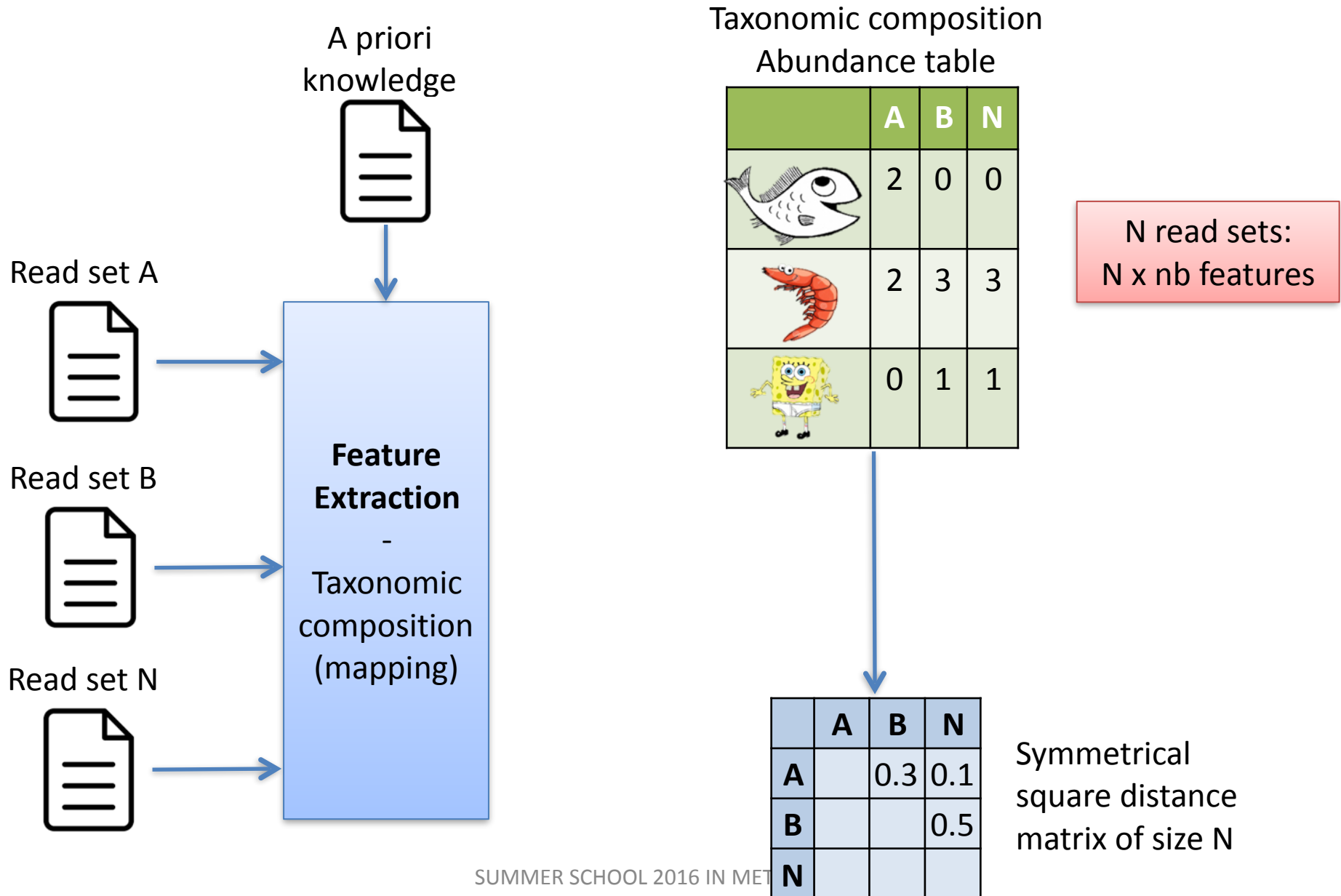
Taxonomic composition  
Abundance table

	A	B
	2	1
	2	3
	0	1

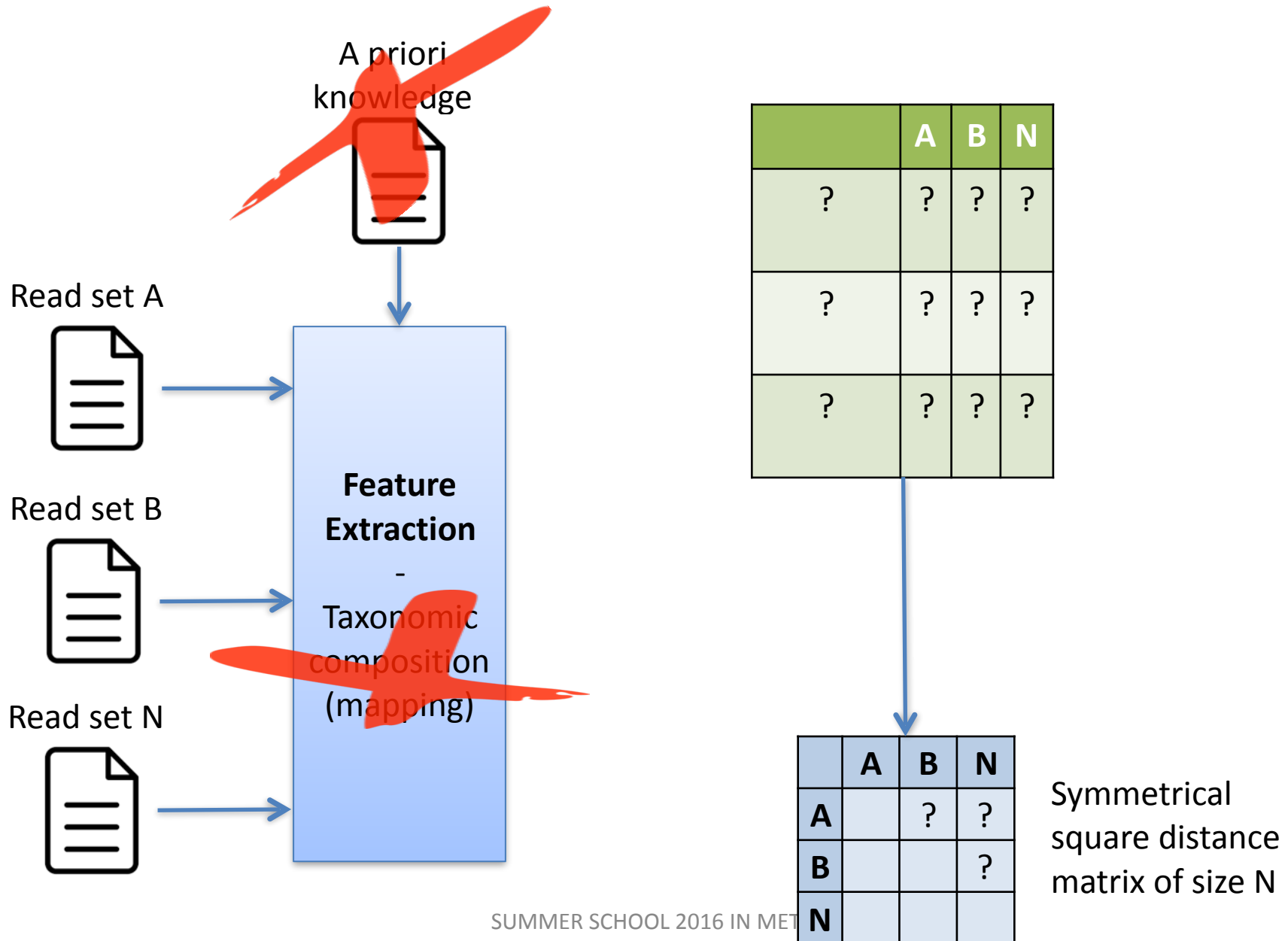


Jaccard(A,B) = 4 / 6

# Comparing N read sets – using a priori knowledge



# Comparing N read sets – *de novo*



# Comparing N read sets – *de novo*

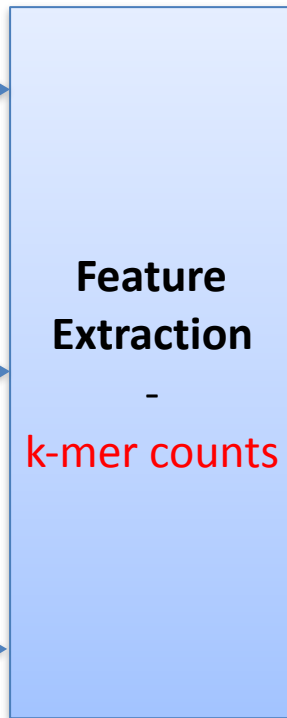
Read set A



Read set B



Read set N



	A	B	N
?	?	?	?
?	?	?	?
?	?	?	?



# Comparing N read sets – *de novo*

Kmer composition  
Abundance table

	A	B	N
ACGAG	2	4	0
CGAGC	2	1	9
GAGCT	0	0	5

Read set A



Read set B



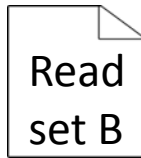
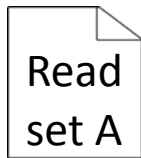
Read set N



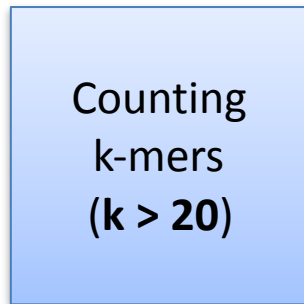
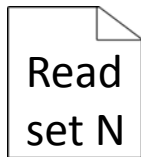
# Issues

Tara Oceans

HMP



...



billions

hundreds

	A	B	...	N
k1	12	11	...	0
k2	3	22	...	0
k3	0	3	...	4
k4	0	8	...	0
k5	0	0	...	2
k6	12	0	...	11
k7	0	1	...	3
k8	21	11	...	1
...	...	...	...	...

K-mer count matrix  $M$   
> 600 TB on HMP

# Issues

Tara Oceans  
HMP

Read  
set A

Read  
set B

Read  
set N



Counting  
k-mers  
( $k > 20$ )



billions

	A	B	...	...
k1	12	11	...	4
k2	3	22	...	0
k3	0	3	...	4
k4	0	8	...	0
k5	0	0	...	2
k6	12	0	...	11
k7	0	1	...	3
k8	21	11	...	1
...	...	...	...	...

K-mer count matrix  $M$   
> 600 TB on HMP

The abundance table is

too **huge**

hundreds

# Ecology distance computation

Most of the ecology distances are **additive over the lines** of the abundance table

$$Jaccard(A, B) = \frac{\sum_{s \in A \cap B} S_A + S_B}{\sum_{s \in A \cup B} S_A + S_B}$$

	A	B
k1	2	2
k2	1	0
k3	0	1
k4	1	1

Distance computation  
(n=intersection u=union)  
n=2 u=2

# Ecology distance computation

Most of the ecology distances are **additive over the lines** of the abundance table

$$Jaccard(A, B) = \frac{\sum_{s \in A \cap B} S_A + S_B}{\sum_{s \in A \cup B} S_A + S_B}$$

	<b>A</b>	<b>B</b>
<b>k1</b>	2	2
<b>k2</b>	1	0
<b>k3</b>	0	1
<b>k4</b>	1	1

Distance computation

(n=intersection u=union)

N=2      u=2

N=2      u=3

# Ecology distance computation

Most of the ecology distances are **additive over the lines** of the abundance table

$$Jaccard(A, B) = \frac{\sum_{s \in A \cap B} S_A + S_B}{\sum_{s \in A \cup B} S_A + S_B}$$

	<b>A</b>	<b>B</b>
k1	2	2
k2	1	0
<b>k3</b>	<b>0</b>	<b>1</b>
k4	1	1

Distance computation

(n=intersection u=union)

n=2      u=2

n=2      u=3

n=2      u=4

# Ecology distance computation

Most of the ecology distances are **additive over the lines** of the abundance table

$$Jaccard(A, B) = \frac{\sum_{s \in A \cap B} S_A + S_B}{\sum_{s \in A \cup B} S_A + S_B}$$

Distances can be computed **one line at a time**

	A	B
k1	2	2
k2	1	0
k3	0	1
k4	1	1

Distance computation  
(n=intersection u=union)

n=2      u=2

n=2      u=3

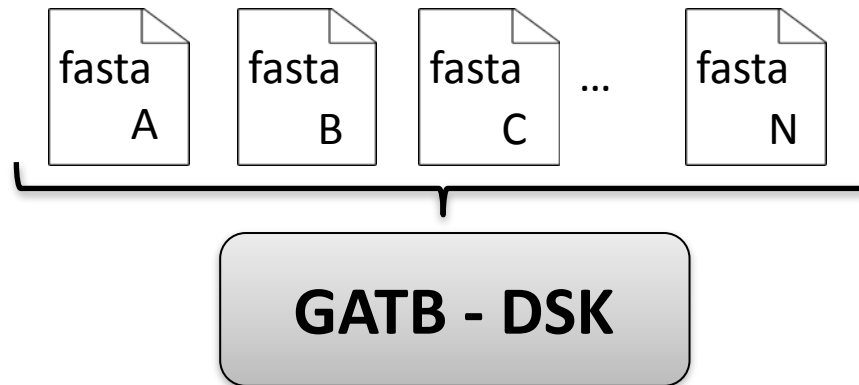
n=2      u=4

n=3      u=5

$$Jaccard(A, B) = 3 / 5$$

# Multiset kmer counting

- Count the kmers of N datasets **simultaneously**
  - Based on KMC2 algorithm (Deorowicz *et al.* Bioinformatics 2015)
  - Available in GATB library (Drezen *et al.* Bioinformatics 2014)



↓ Streaming for each distinct kmer

	A	B	C	...	N
ACGATC	0	4	52	...	0

**Its abundance in each dataset**



# From reads to counted kmers

GATB - DSK

basic idea

Can't be more simple:

1. Read kmers and write them in a file
2. Sort file
3. Identical kmers occur consecutively, count them

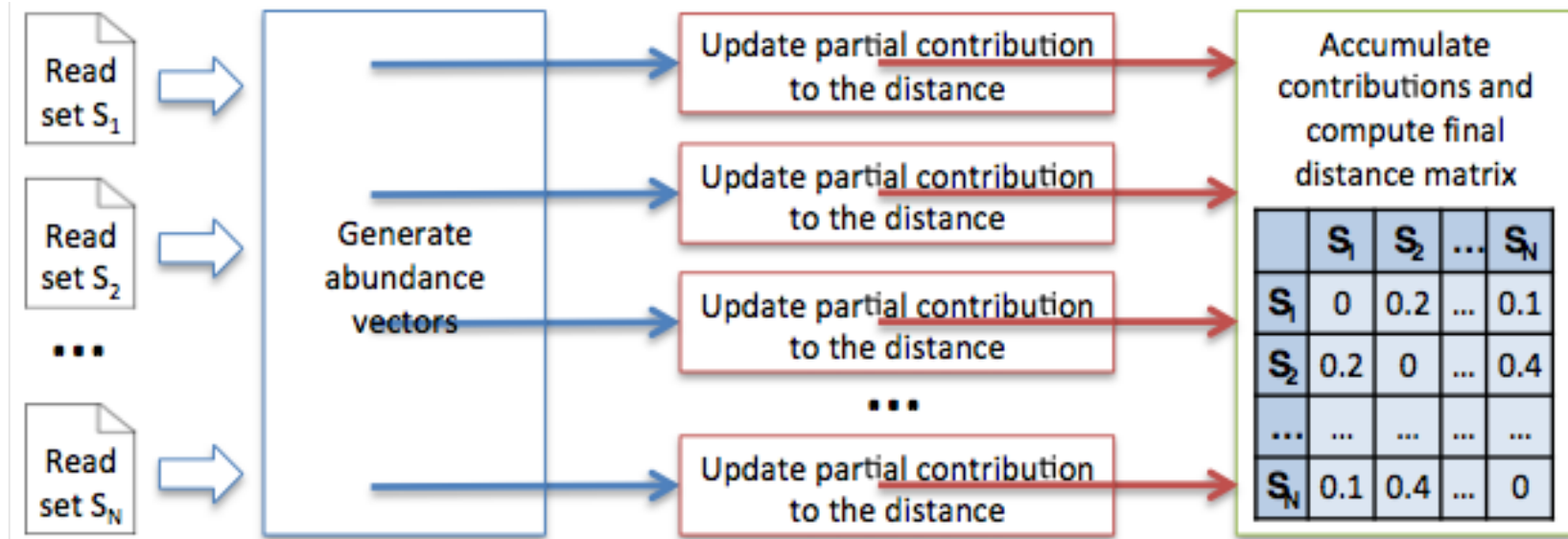
1. CAGG  
ACGG  
CAGG  
TTAC  
CAGG  
ACGG

2. ACGG  
ACGG  
CAGG  
CAGG  
CAGG  
TTAC

3. ACGG 2  
CAGG 3  
TTAC 1



# Dealing with a **streaming** of abundance vectors



# Simka performances



Full HMP project (690 samples, 32 billions reads)

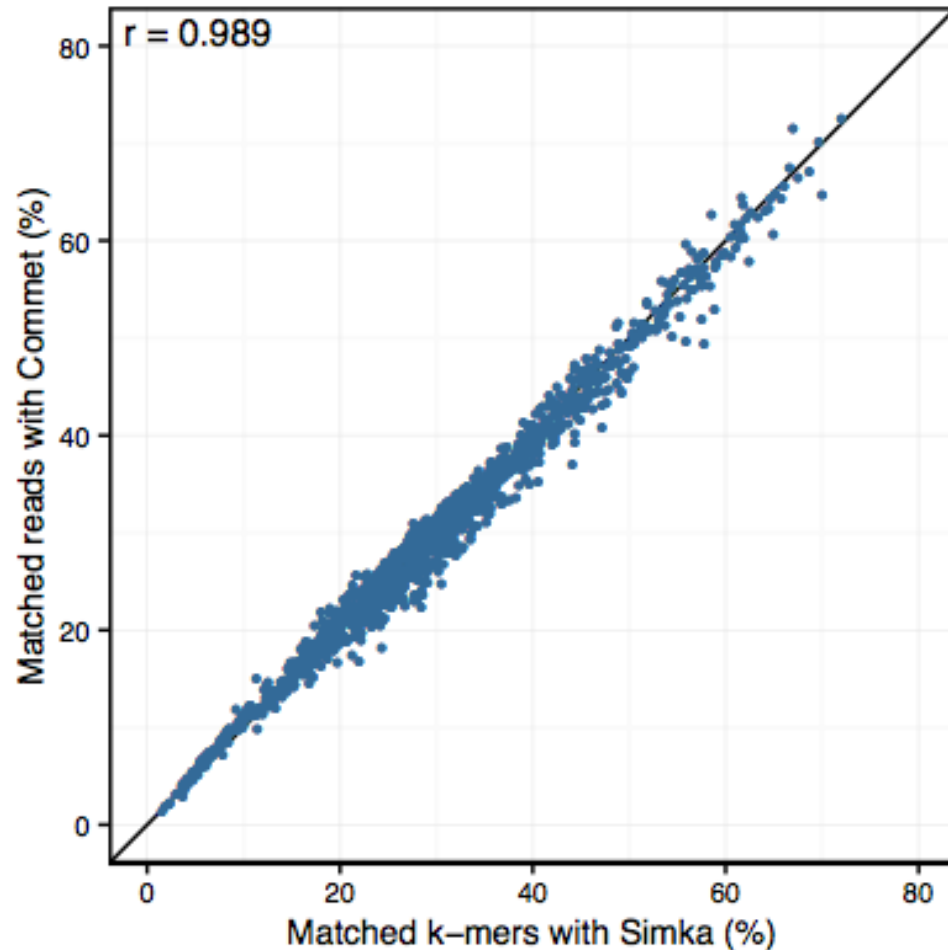
- Computation time: < 14h
- 64 GB memory
- ~1T disk (< half of input size)

# Simka computes several distances

Name	Definition	$C_{S_i}$	$f(x, y, X, Y)$	$g(x)$
<b>Quantitative distances</b>				
Chord	$\sqrt{2 - 2 \sum_w \frac{N_{S_i}(w)N_{S_j}(w)}{C_{S_i}C_{S_j}}}$	$\sqrt{\sum_w N_{S_i}(w)^2}$	$\frac{xy}{XY}$	$\sqrt{2 - 2x}$
Hellinger	$\sqrt{2 - 2 \sum_w \frac{\sqrt{N_{S_i}(w)N_{S_j}(w)}}{\sqrt{C_{S_i}C_{S_j}}}}$	$\sum_w N_{S_i}(w)$	$\frac{\sqrt{xy}}{\sqrt{XY}}$	$\sqrt{2 - 2x}$
Whittaker	$\frac{1}{2} \sum_w \frac{ N_{S_i}(w)C_{S_j} - N_{S_j}(w)C_{S_i} }{C_{S_i}C_{S_j}}$	$\sum_w N_{S_i}(w)$	$\frac{ xY - yX }{XY}$	$\frac{x}{2}$
Bray-Curtis	$\sum_w \frac{ N_{S_i}(w) - N_{S_j}(w) }{C_{S_i} + C_{S_j}}$	$\sum_w N_{S_i}(w)$	$\frac{ x - y }{X + Y}$	$x$
Kulczynski	$1 - \frac{1}{2} \sum_w \frac{(C_{S_i} + C_{S_j}) \min(N_{S_i}(w), N_{S_j}(w))}{C_{S_i}C_{S_j}}$	$\sum_w N_{S_i}(w)$	$\frac{(X + Y) \min(x, y)}{XY}$	$1 - \frac{x}{2}$
Jensen-Shannon	$\frac{1}{2} \sum_w \left[ \frac{N_{S_i}(w)}{C_{S_i}} \log \frac{2C_{S_j}N_{S_i}(w)}{C_{S_j}N_{S_i}(w) + C_{S_i}N_{S_j}(w)} + \frac{N_{S_j}(w)}{C_{S_j}} \log \frac{2C_{S_i}N_{S_j}(w)}{C_{S_j}N_{S_i}(w) + C_{S_i}N_{S_j}(w)} \right]$	$\sum_w N_{S_i}(w)$	$\frac{x}{X} \log \frac{2xY}{xY + yX} + \frac{y}{Y} \log \frac{2yX}{xY + yX}$	$\sqrt{\frac{x}{2}}$
Canberra	$\frac{1}{a + b + c} \sum_w \frac{ N_{S_i}(w) - N_{S_j}(w) }{N_{S_i}(w) + N_{S_j}(w)}$	-	$\frac{ x - y }{ x + y }$	$\frac{1}{a + b + c}x$
Jaccard	$1 - \sum_w \frac{(N_{S_i}(w) + N_{S_j}(w)) \mathbb{1}_{\{N_{S_i}(w)N_{S_j}(w) > 0\}}}{C_{S_i} + C_{S_j}}$	$\sum_w N_{S_i}(w)$	$\frac{(x + y) \mathbb{1}_{\{xy > 0\}}}{X + Y}$	$x$
<b>Qualitative distances</b>				
Chord/Hellinger	$\sqrt{2 \left( 1 - \frac{a}{\sqrt{(a+b)(a+c)}} \right)}$	-	-	-
Whittaker	$\frac{1}{2} \left( \frac{b}{a+b} + \frac{c}{a+c} + \left  \frac{a}{a+b} - \frac{a}{a+c} \right  \right)$	-	-	-
Bray-Curtis/Sorensen	$\frac{b+c}{2a+b+c}$	-	-	-
Kulczynski	$1 - \frac{1}{2} \left( \frac{a}{a+b} - \frac{a}{a+c} \right)$	-	-	-
Ochiai	$1 - \frac{a}{\sqrt{(a+b)(a+c)}}$	-	-	-
Jaccard	$\frac{b+c}{a+b+c}$	-	-	-
<b>Abundance-based (AB) variants of qualitative distances</b>				
AB-Jaccard	$1 - \frac{UV}{U+V-UV}$	-	-	-

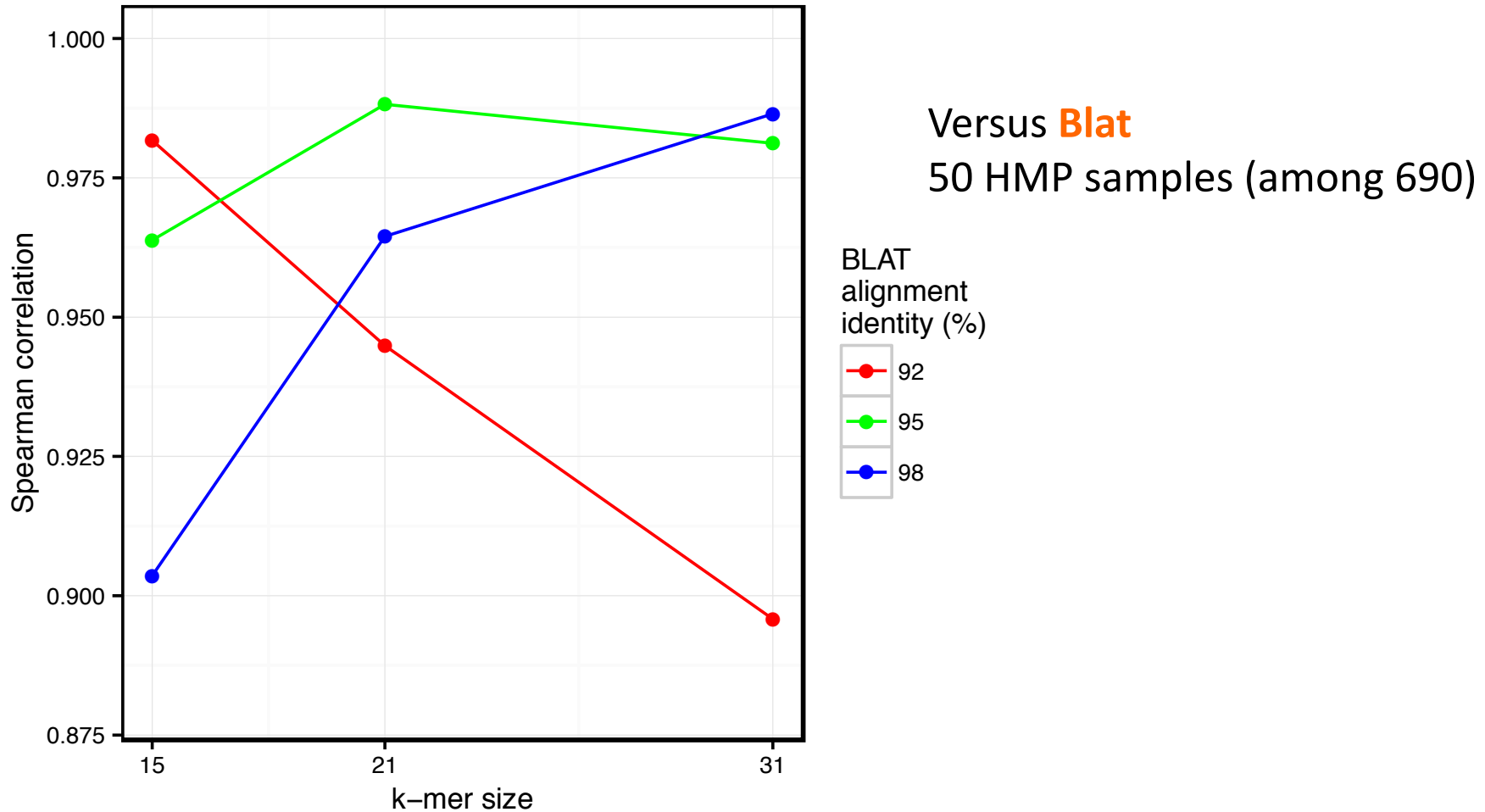


# Are kmer distances good enough?



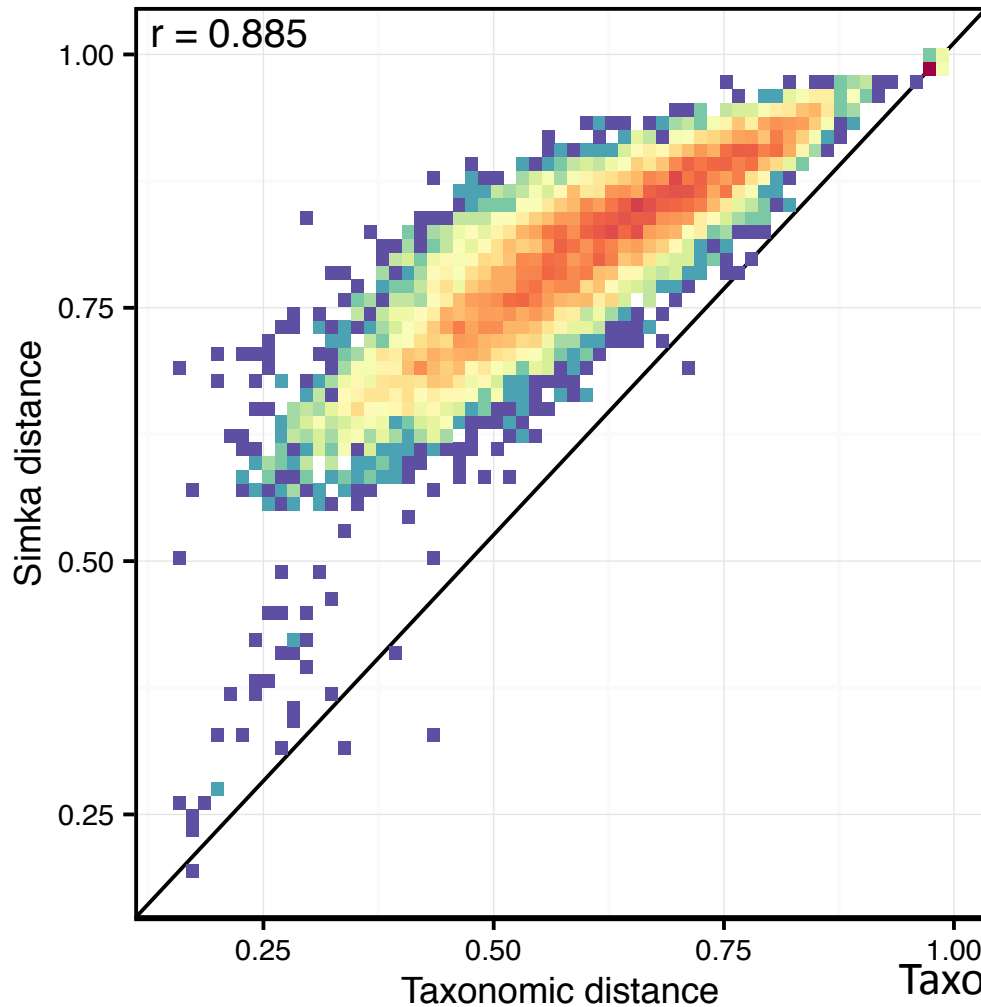
Versus **Commet**  
50 HMP samples (among 690)

# Are kmer distances good enough?

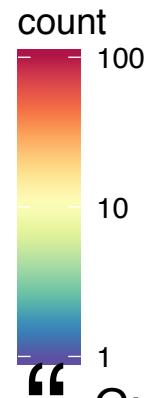




# Are kmer distances good enough?



Versus **Taxonomical dist.**  
HMP **GUT** samples



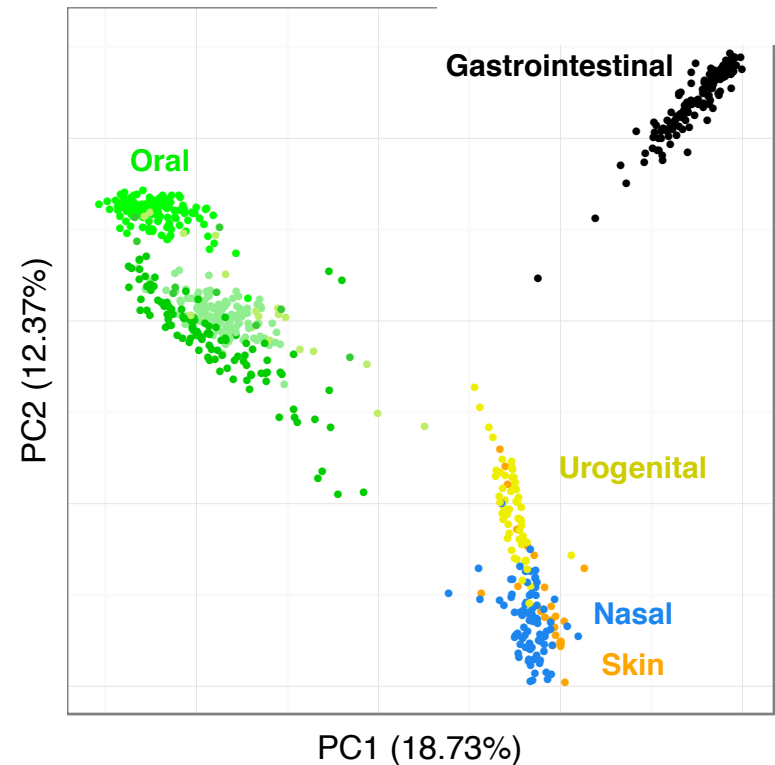
Mash Spearman = 0.51:

“ Gut samples differ more in terms of relative abundances of microbes than in terms of composition

Taxonomic distances are obtained from  
<http://www.hmpdacc.org/HMSCP/>

# Tests on full HMP project

Simka (14h, 62GB)

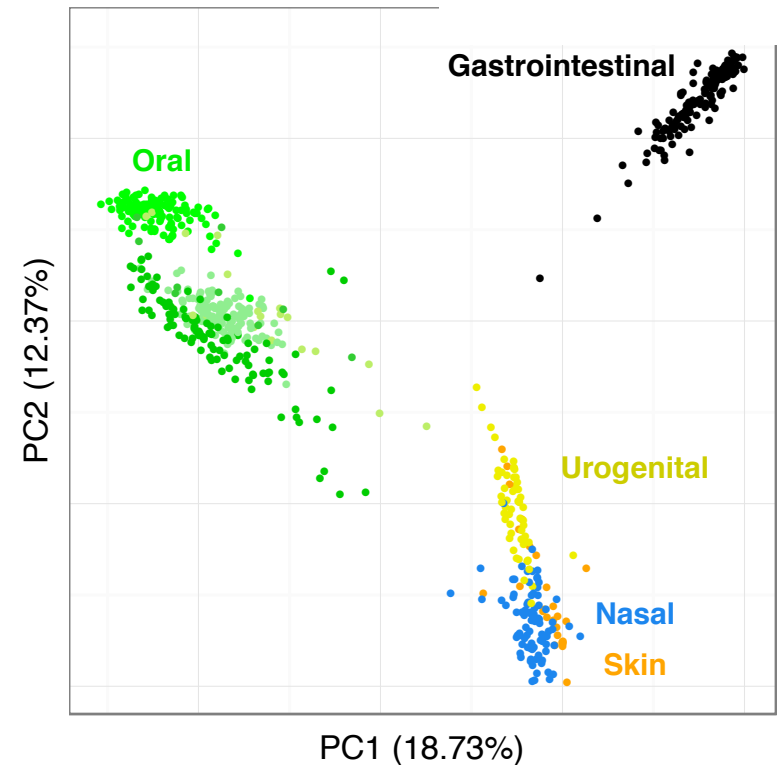
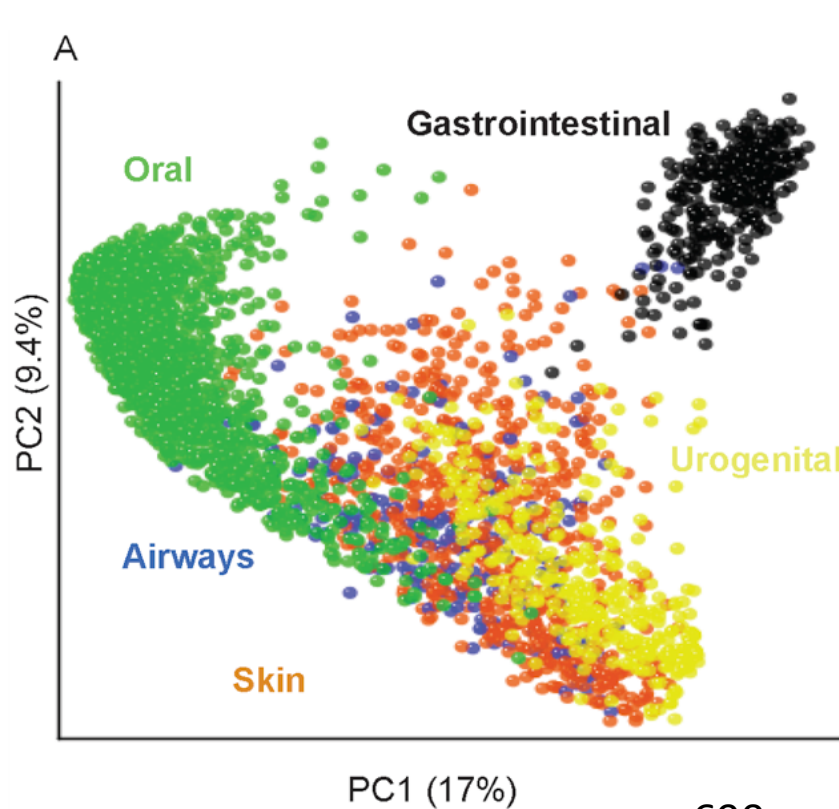


- 690 samples - 3727 GB
- 2x16 billions paired reads

# Tests on full HMP project

HMP – (OTU) [1]

Simka (14h, 62GB)



- 690 samples - 3727 GB  
- 2x16 billions paired reads

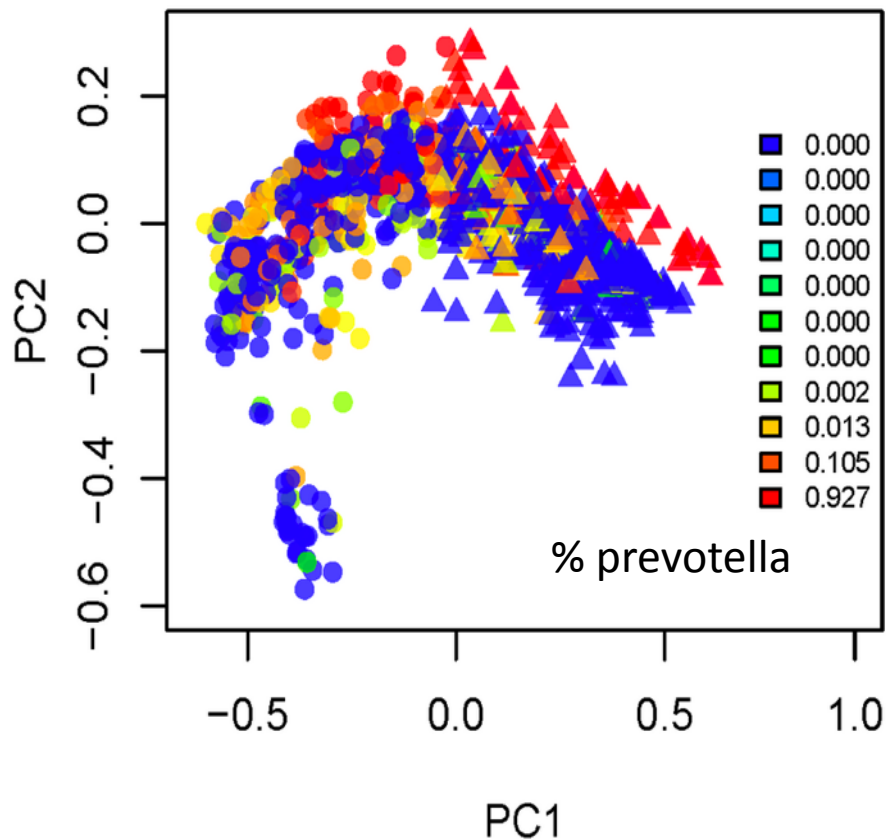
[1] Koren et al 2013

SUMMER SCHOOL 2016 IN METAGENOMICS

A Guide to Enterotypes across the Human Body [...]

# Tests on HMP project – GUT enterotypes

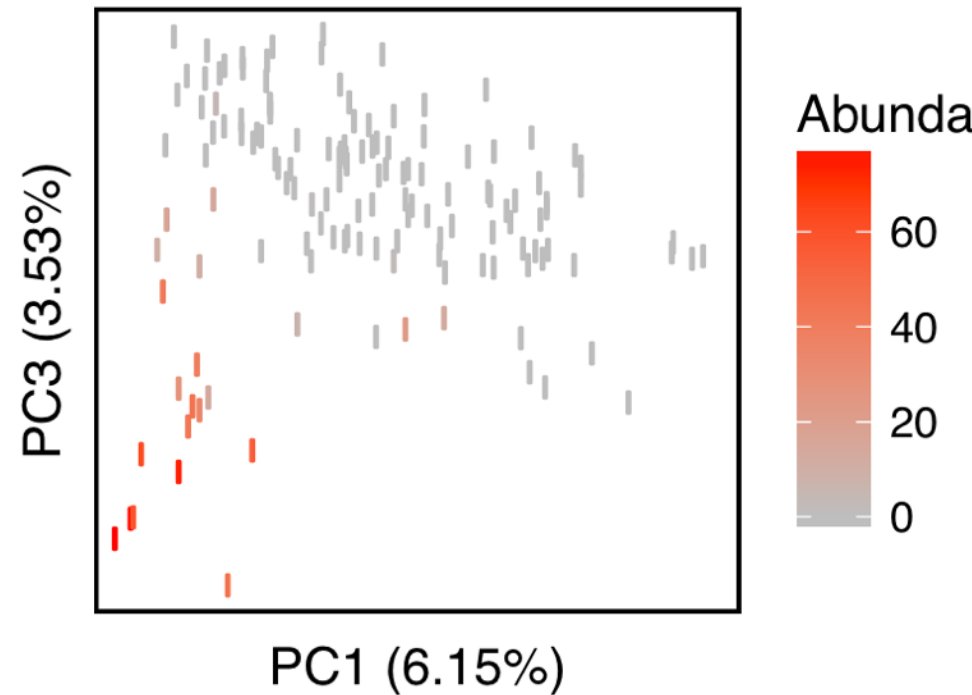
HMP (OTU) [1]



Simka

**B**

Prevotella

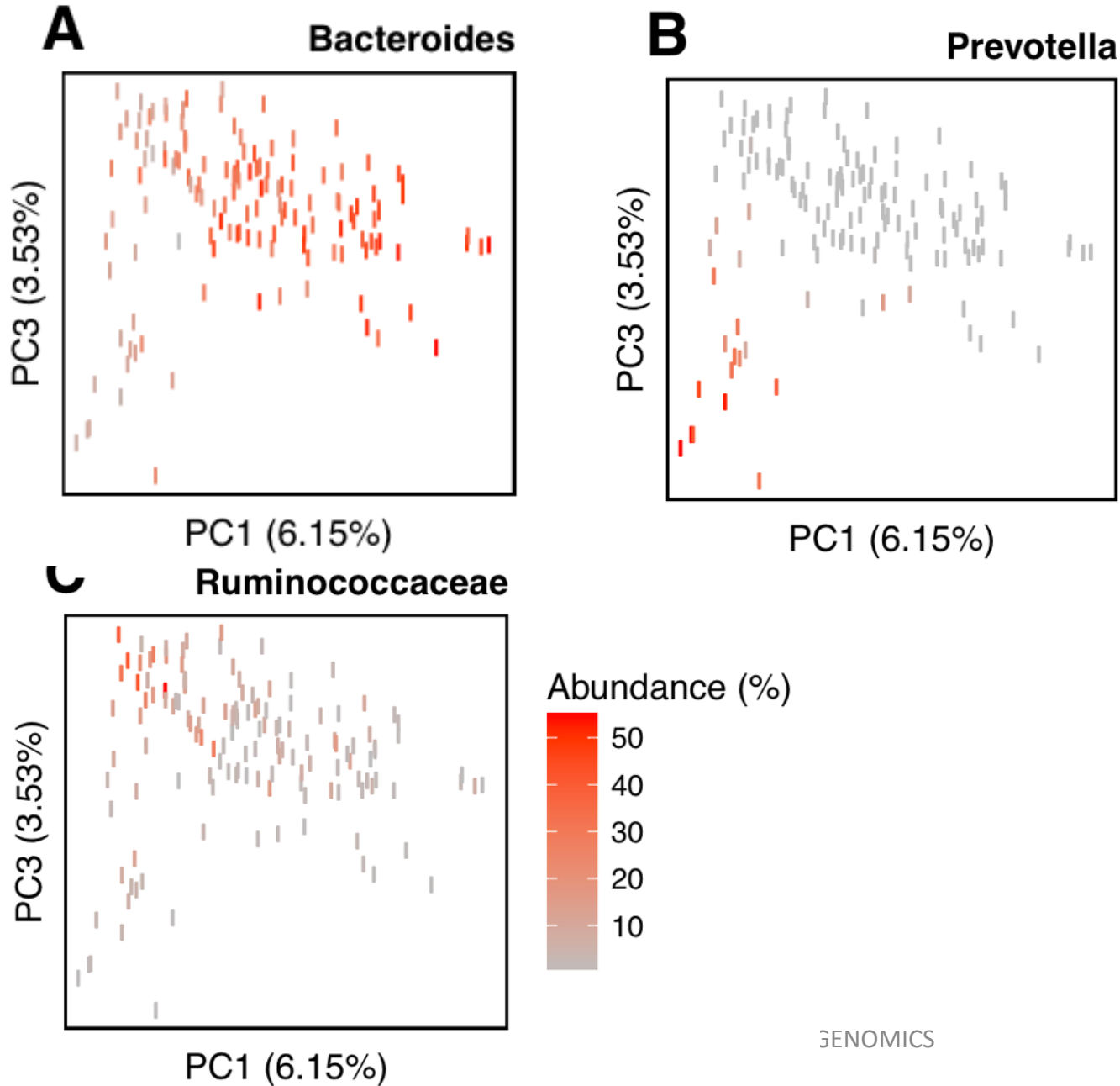


[1] Koren et al 2013

SUMMER SCHOOL 2016 IN METAGENOMICS

A Guide to Enterotypes across the Human Body [...]

# Tests on HMP project – GUT enterotypes



# Simka: Take Home Message



- Push button
- Many ecological distances
- Ultra Fast computation
- Acceptable memory footprint
- Does not provide links
- ***k-mers instead of species***
  - *under.: horizontal gene*
  - *over.: ≠ genome sizes*

Benoit, G., Peterlongo, P., Mariadassou, M., Drezen, E., Schbath, S., Lavenier, D., & Lemaitre, C.  
Multiple Comparative Metagenomics using Multiset k-mer Counting.  
arXiv id: 1604.02412 – Peer J. review process –




# Simka: Take Home Message



- Push button
  - Many ecological distances
  - Ultra Fast computation
  - Acceptable memory footprint
  - **Validates *k-mer* based distances**
- Does not provide links
  - ***k-mers instead of species***
    - *under.: horizontal gene*
    - *over.: ≠ genome sizes*

Benoit, G., Peterlongo, P., Mariadassou, M., Drezen, E., Schbath, S., Lavenier, D., & Lemaitre, C.  
Multiple Comparative Metagenomics using Multiset k-mer Counting.  
arXiv id: 1604.02412 – Peer J. review process –

# Future

	Commet	SRC_linke r	Simka	?
Proc 	✓	✗	✓	✓
Scale 	Mem ✓	✗	✓	✓
	Time ✗	✗	✓	✓
	VS	VS	VS	VS



# Future

## Algorithmic

- $k$ -mer subsampling
- Dynamic addition of new sets

## Applications

- Tara analyses
- ...

# Thanks!

## Tools

- **Commet**  
[github.com/pierrepeterlongo/commet](https://github.com/pierrepeterlongo/commet)
- **Simka**  
[github.com/GATB/simka](https://github.com/GATB/simka)
- **SRC\_linker**  
[github.com/GATB/rconnector](https://github.com/GATB/rconnector)

[pierre.peterlongo@inria.fr](mailto:pierre.peterlongo@inria.fr)



Olivier Jaillon



Antoine Limasset



Lucie Bittner



Mahendra  
Mariadassou



Thomas Vanier



Claire  
Lemaitre



Nicolas Maillet



Gaetan Benoit



Sophie Schbath



Dominique  
Lavenier



Guillaume  
Collet



Erwan Drezen



Camille Marchet









# FP & FN pregnancy test example

- False Negative:  
“You are not pregnant”
- False Positive:  
“You are pregnant”

