

Kaiju

Fast and sensitive taxonomic classification for metagenomics

Peter Menzel

 pmenzel@gmail.com  [ptr_menzel](https://twitter.com/ptr_menzel)

Metagenomics

Biological Question

What is the composition of the microbial community in a specific habitat?

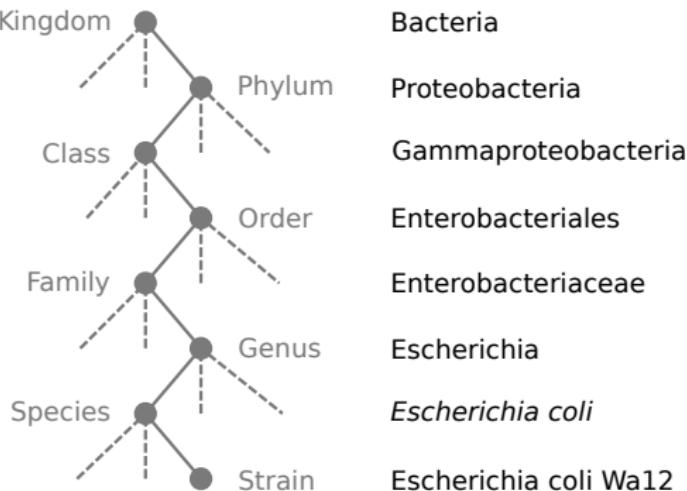
For example: drinking water, compost, human skin or gut, bioreactor, . . .



Metagenomics

Biological Question

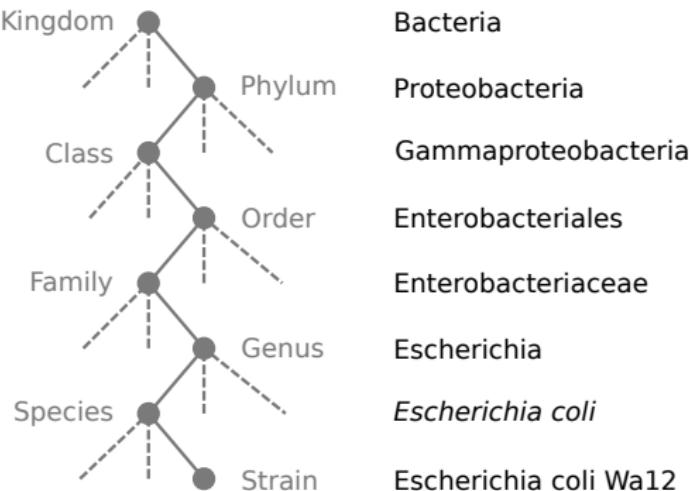
What is the composition of the microbial community in a specific habitat?
For example: drinking water, compost, human skin or gut, bioreactor, ...



Metagenomics

Biological Question

What is the composition of the microbial community in a specific habitat?
For example: drinking water, compost, human skin or gut, bioreactor, ...



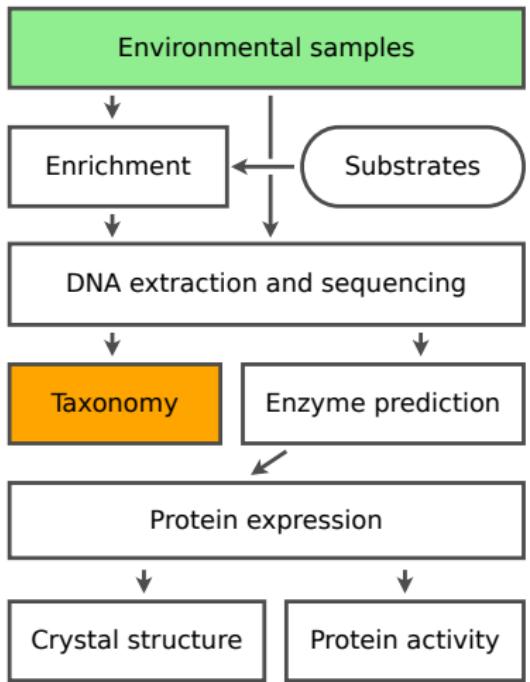
Metagenomic Sequencing

Sequencing of total genomic DNA.

Estimate relative abundances of taxa at a given taxonomic level.



Hotzyme: Discovery of novel hydrolases from hot springs



The Hotzyme Project - Metagenomics

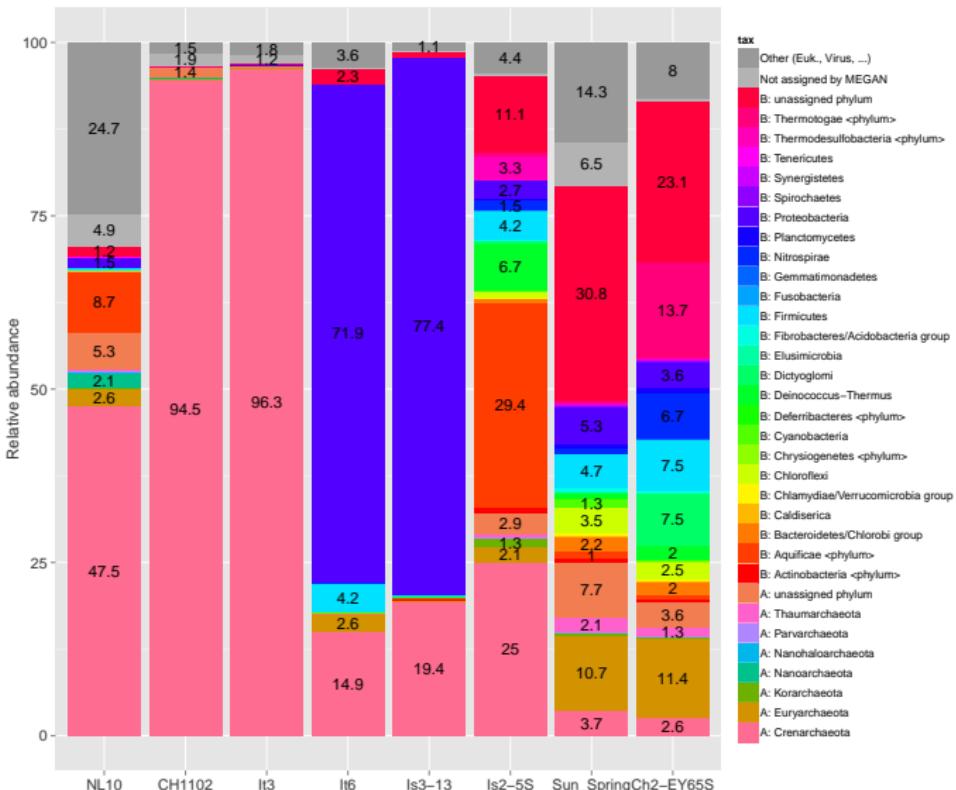
Whole genome sequencing of total DNA from samples

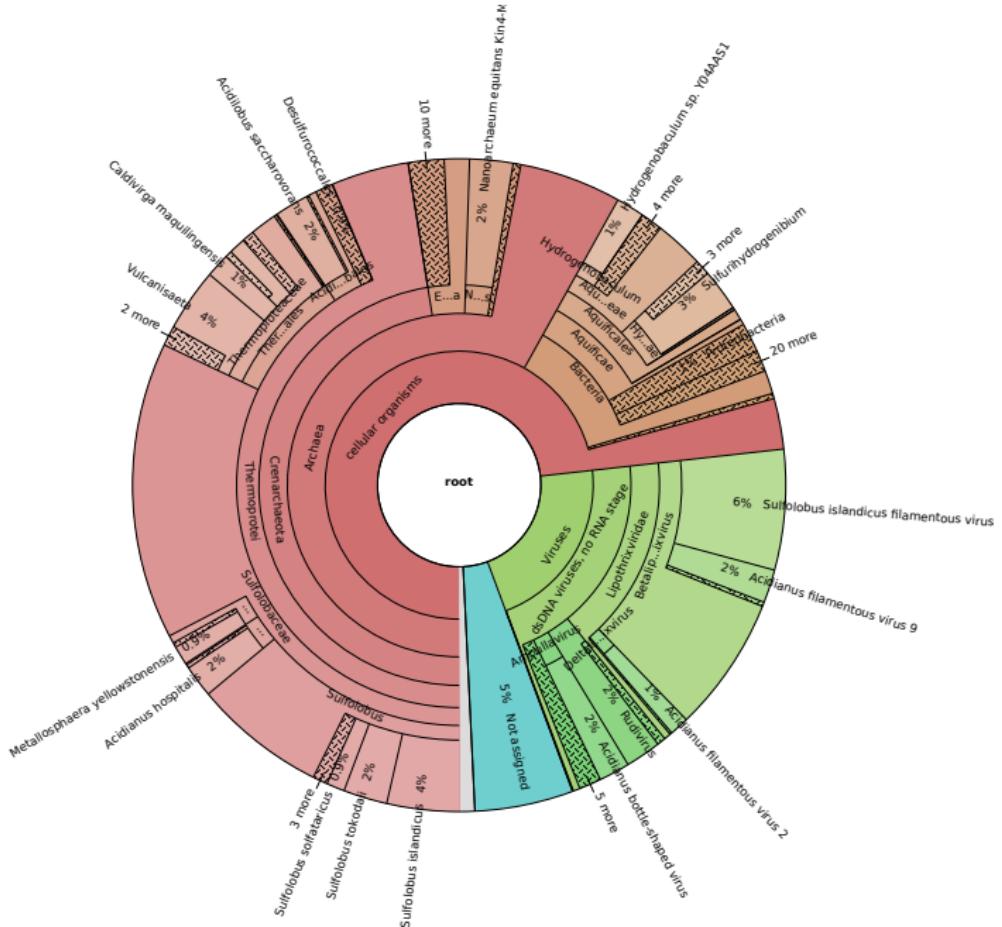
Sample Name	T in °C	pH	Site
Ch2-EY65S	65	7.0	Eryuan, China
Sun Spring	61-64	5.8-6.0	Uzon Caldera, Russia
It6	76	3.0	Pozzuoli, Italy
It3	86	5.5	Pisciarelli, Italy
Is2-5S	85-90	5.0	Grensdalur, Iceland
Is3-13	90	3.5-4.0	Krísuvík, Iceland
CH1102	79	1.8	YNP, USA
NL10	92	3.0-4.0	YNP, USA

Assembly of reads followed by aligning contigs to GenBank NR.



Diversity on phylum-level





Sequence classification

Computational problem

Assign individual HTS reads to taxa by comparison to a reference database.
Fast!



Sequence classification

Computational problem

Assign individual HTS reads to taxa by comparison to a reference database.
Fast!

Fast = Alignment-free

- Local alignment methods are too slow for dealing with millions of short reads (blastn, blastx, blat, RapSearch, ...)
- Comparison of genomic k -mers (exact matches of length k) between reads and database
- LMAT, Kraken, Clark, ...
- choice of k determines accuracy
- at least one k -mer needed per read
- $k = 31$ in Kraken and Clark



Sequence classification

Computational problem

Assign individual HTS reads to taxa by comparison to a reference database.
Fast!

Fast = Alignment-free

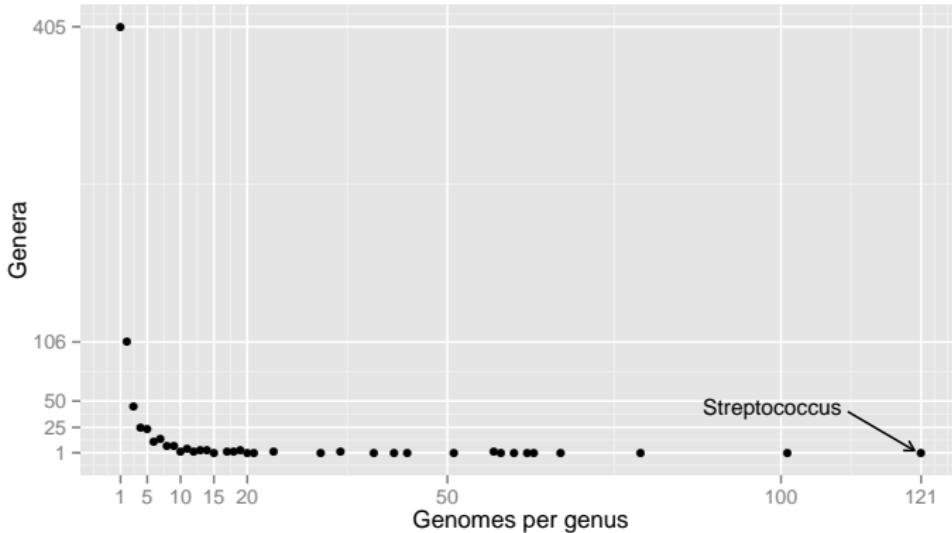
- Local alignment methods are too slow for dealing with millions of short reads (blastn, blastx, blat, RapSearch, ...)
- Comparison of genomic k -mers (exact matches of length k) between reads and database
- LMAT, Kraken, Clark, ...
- choice of k determines accuracy
- at least one k -mer needed per read
- $k = 31$ in Kraken and Clark

Hotzyme MGs

MG	Kraken
SunSpr	4.3 %
CH2	4.9 %
NL10	9.0 %
It6	14.0 %
CH11	14.7 %
Is3-13	33.5 %
Is2-5S	41.1 %
It3	88.1 %



Reference database bias



- over-representation of certain model organisms, pathogens, ...
- lack of species that are not possible to culture, extremophiles, ...
- > 1 500 genera without any sequenced genome
- evolutionary divergence between genomes in database and metagenomes
- requirement of long exact matches (e.g. $k=31$) is very stringent





- find maximum exact matches
- compare sequences on protein level
- focus on speed
- Burrows-Wheeler transform + FM-index



Kaiju's algorithm

Sequencing Read

1. Translation

Translate nucleotide sequence into amino acid sequences by the six possible reading frames and split into *fragments* at stop codons.



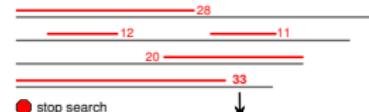
2. Sorting

MEM Sort fragments by length $> m$



3. Database search

Find MEMs with $\text{length} > m$

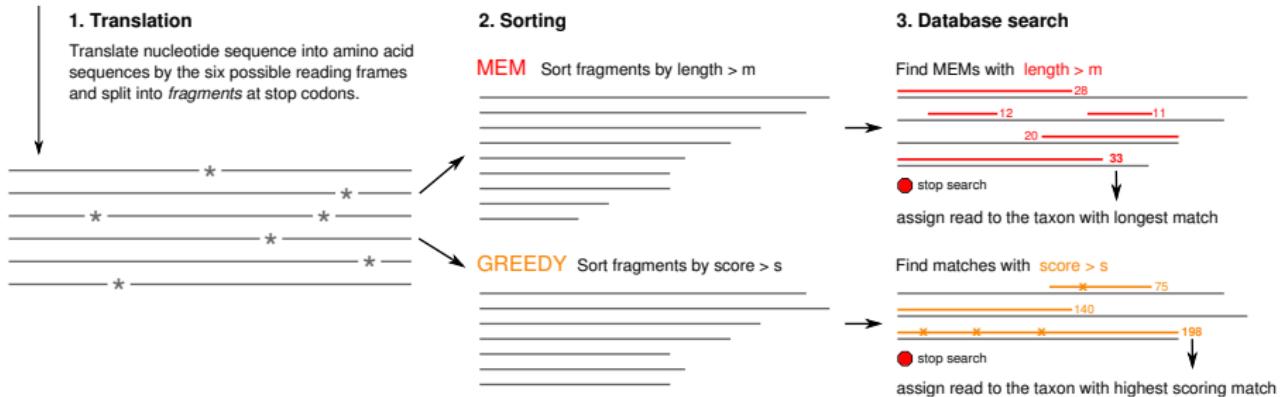


assign read to the taxon with longest match



Kaiju's algorithm

Sequencing Read



Benchmark

Simulate the problem of classifying a novel species / strain.



Benchmark

Simulate the problem of classifying a novel species / strain.

** Protocol **

- * 2,724 complete microbial genomes belonging to 692 genera
- * 882 of these genomes belong to a genus that has at least two and at most ten genomes

FOR EACH genome DO

1. make reference database **without** it
2. simulate Illumina and 454 reads
3. classify reads using the trimmed database

DONE



Benchmark

Simulate the problem of classifying a novel species / strain.

** Protocol **

- * 2,724 complete microbial genomes belonging to 692 genera
 - * 882 of these genomes belong to a genus that has at least two and at most ten genomes
- FOR EACH genome DO
1. make reference database **without** it
 2. simulate Illumina and 454 reads
 3. classify reads using the trimmed database
- DONE

Programs

Kaiju vs. Kraken vs. Clark

Sensitivity

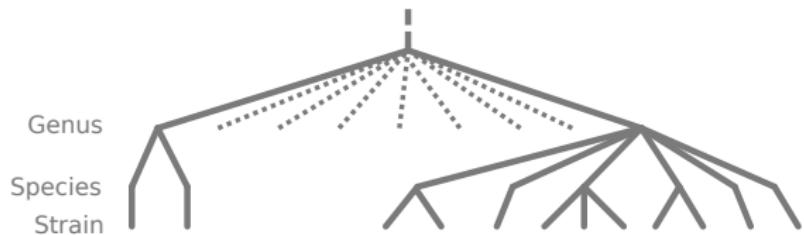
$$\frac{\text{Number of correct classified}}{\text{Total number of reads}} \times 100$$

Precision

$$\frac{\text{Number of correct classified}}{\text{Number of classified reads}} \times 100$$



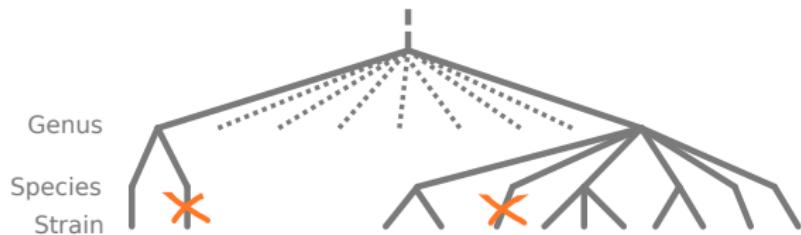
Genus categories



Genomes per genus (#Genomes in this category)	2 (212)	3 (135)	4 (100)	5 (115)	6 (66)	7 (98)	8 (64)	9 (72)	10 (20)
--	------------	------------	------------	------------	-----------	-----------	-----------	-----------	------------



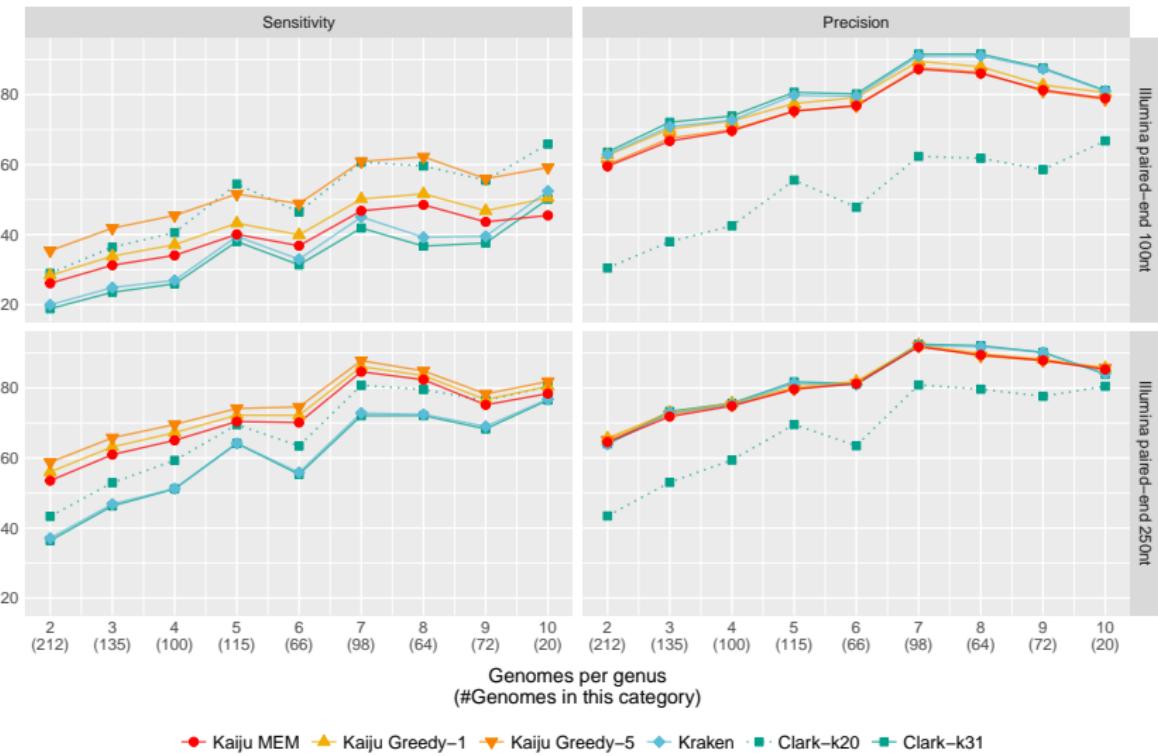
Genus categories



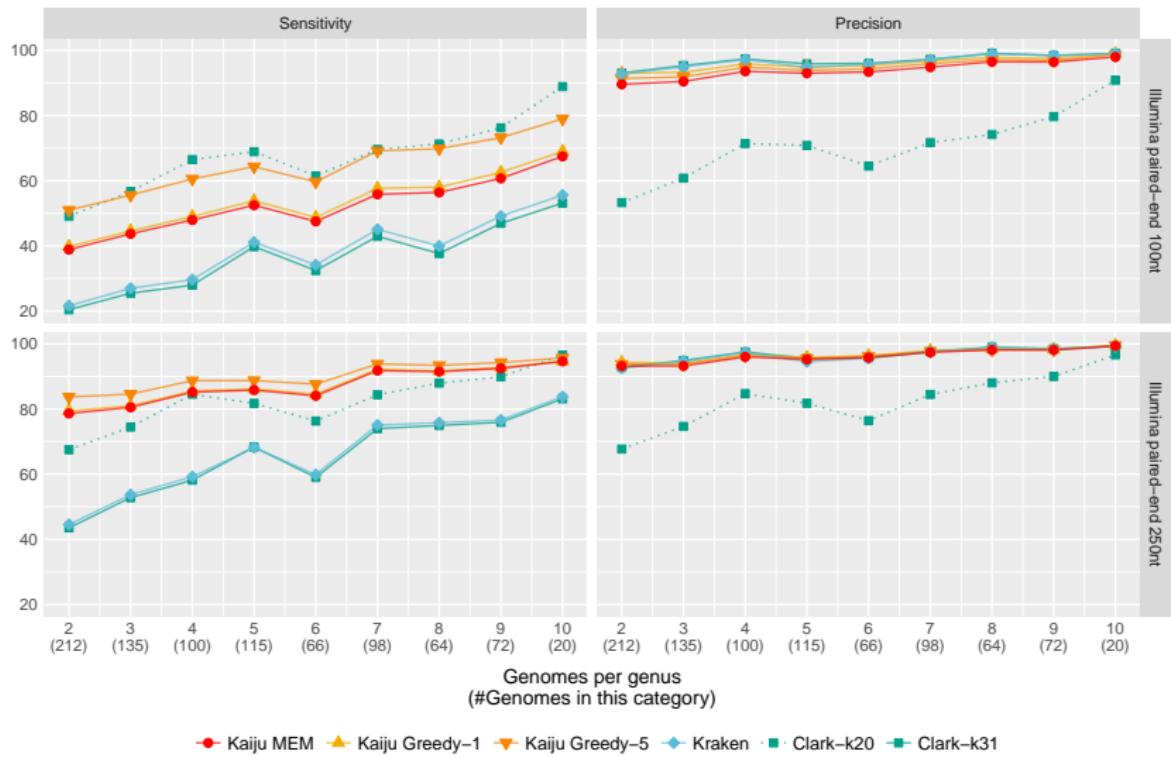
Genomes per genus (#Genomes in this category)	2 (212)	3 (135)	4 (100)	5 (115)	6 (66)	7 (98)	8 (64)	9 (72)	10 (20)
--	------------	------------	------------	------------	-----------	-----------	-----------	-----------	------------



Sensitivity vs precision on genus-level

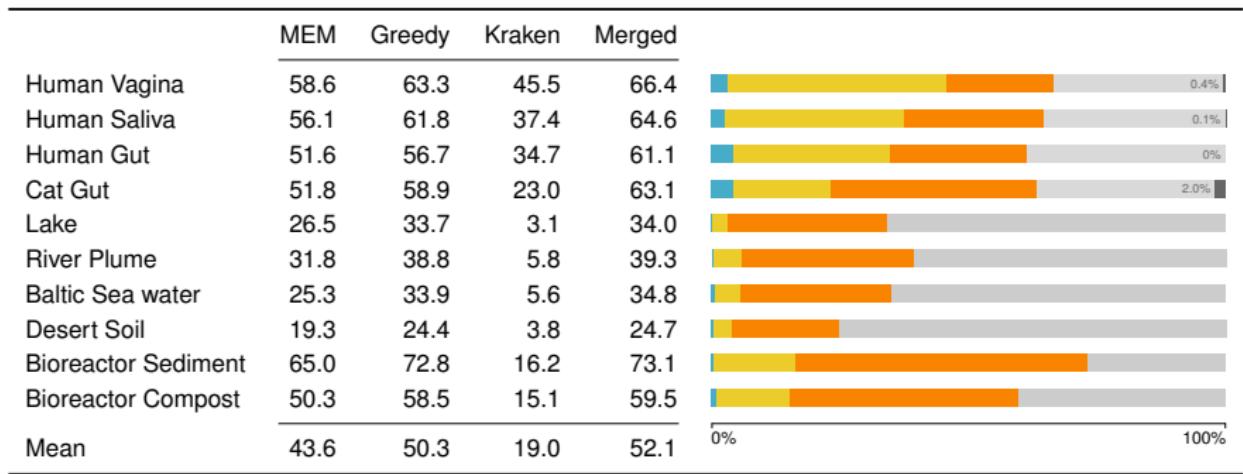


Sensitivity vs precision on phylum-level



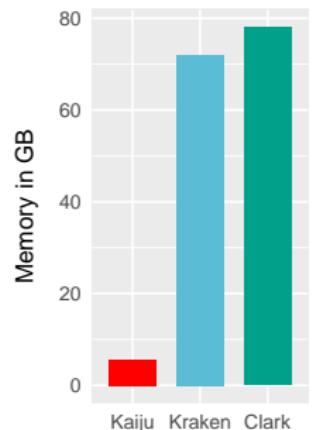
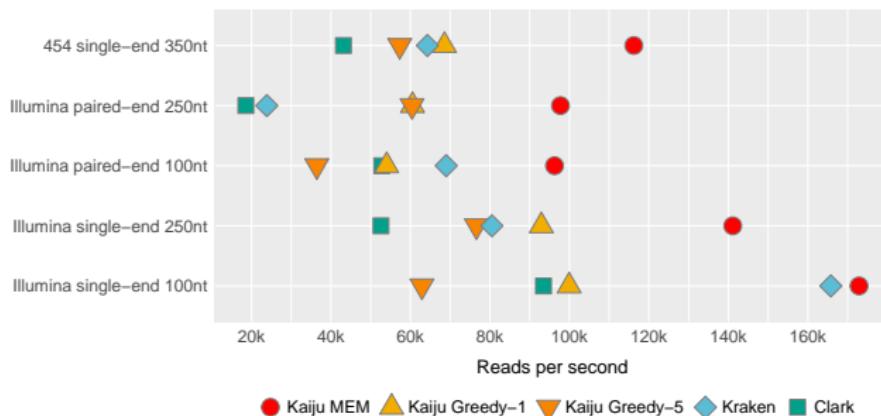
Real metagenomes

Percentage of classified reads from random metagenomic samples



Speed and memory requirement

using a reference database of 2724 complete genomes

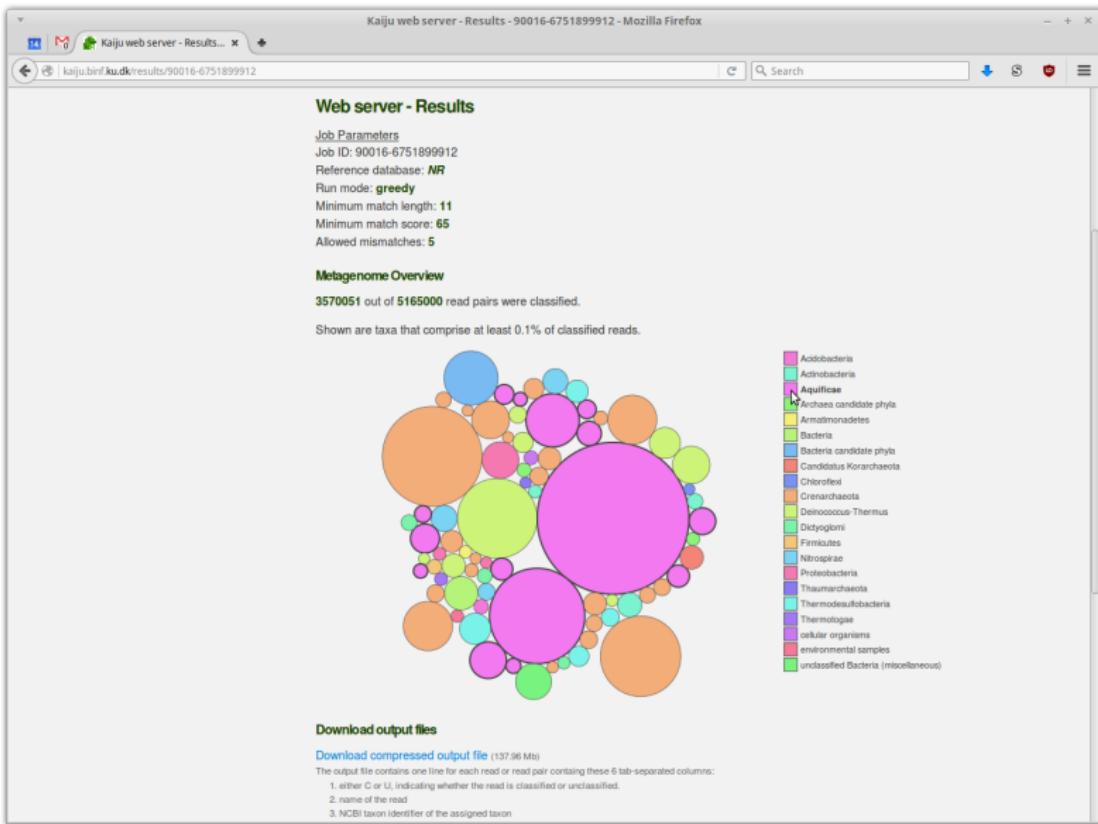


Web server - <http://kaiju.binf.ku.dk>

The screenshot shows a Mozilla Firefox browser window with the title "Kaiju web server - Submit Job - Mozilla Firefox". The address bar displays "kaiju.binf.ku.dk/server". The main content area features a green pixelated dino logo and the word "KAIJU". Below this, the text "Fast and sensitive taxonomic classification for metagenomics" is displayed. A sidebar on the left contains sections for "Reference databases" (Complete Genomes, NCBI BLAST NR) and a "Last updated" timestamp (2016-02-24). The central form is titled "Web server - Submit job" and includes fields for "Job Name" (with placeholder "Give a custom name to your submission."), "e-mail" (with placeholder "Receive a notification after the reads have been processed. [?]"), and "File with sequencing reads" (with placeholder "Supported formats are compressed FASTA and FASTQ [?]"). It also includes a "Select file" button, a file name input field containing "sun_spring.fasta.gz", a "Start upload" button, a progress bar, and a checkbox for "Upload a second file for paired-end sequencing".



Web server - Results page



Conclusion (last slide!)

Hotzyme MGs

MG	Kraken	Kaiju
SunSpr	4.3%	43.1%
CH2	4.9%	31.9%
NL10	9.0%	29.8%
It6	14.0%	49.7%
CH11	14.7%	65.5%
ls3-13	33.5%	63.4%
ls2-5S	41.1%	55.8%
It3	88.1%	92.1%



Conclusion (last slide!)

Hotzyme MGs

MG	Kraken	Kaiju
SunSpr	4.3%	43.1%
CH2	4.9%	31.9%
NL10	9.0%	29.8%
It6	14.0%	49.7%
CH11	14.7%	65.5%
ls3-13	33.5%	63.4%
ls2-5S	41.1%	55.8%
It3	88.1%	92.1%

Summary

- BWT + FM-index can be used with a database of all microbial/viral proteins
- very fast and small index size
- can also use full microbial BLAST *nr*
- Benchmark shows higher sensitivity with similar precision
- Up to ten times more reads classified in real metagenomes



Menzel P., Ng K.L., Krogh A. (2016) Fast and sensitive taxonomic classification for metagenomics with Kaiju. *Nature Communications*, 7:11257



Anders Krogh



Kim Ng

