



Recent progress

towards petabase-scale genomics

Rayan Chikhi, Institut Pasteur

Simons Institute, Very Large-scale 'Omics' 2022

Hello

- PI, Bioinformatics algorithms lab @ Institut Pasteur
- CV: PhD@ENS Rennes, Postdoc@PSU, CNRS

Research:

- *de novo* assembly
- k-mer methods
- metagenomics
- large-scale bioinfo



@RayanChikhi on Twitter 

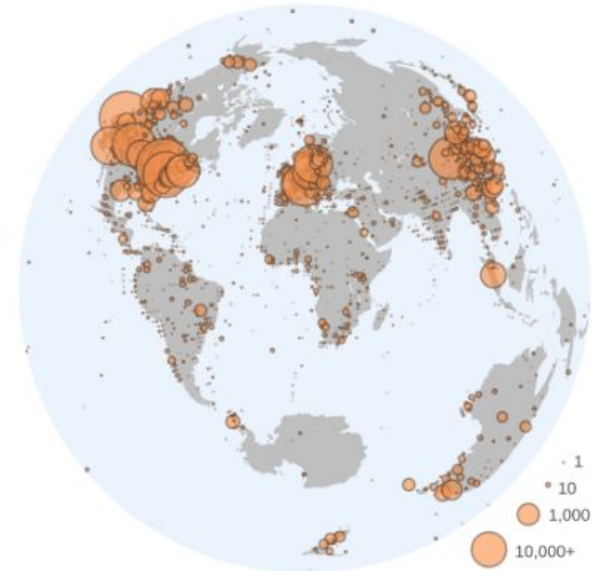
<http://rayan.chikhi.name>

Petabase-scale viral discovery

Rayan Chikhi, on behalf of the Serratus team

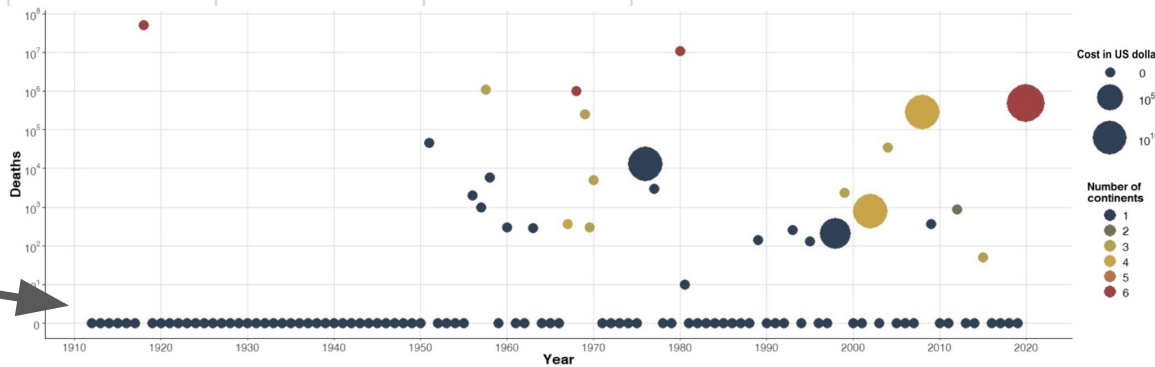
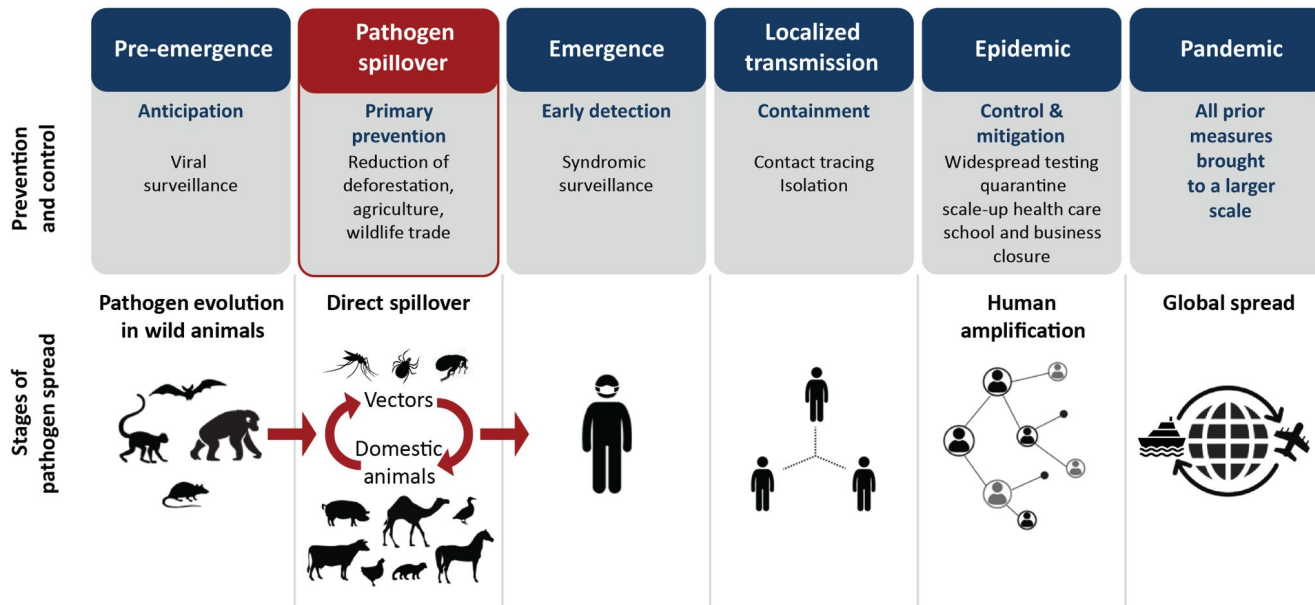
We analysed all available RNA sequencing data and discovered 10x more viruses species than previously known, including coronaviruses.

Nature, 2022



Viral surveillance in the age of pandemics

Source: <https://www.science.org/doi/10.1126/sciadv.abl4183>



Would be nice to contain viruses there

SARS-CoV-2 circulate(s|d) among animals

PETS & ANIMALS

Tiger at zoo in Knoxville tests positive for SARS-CoV-2, two others possibly infected

A veterinary team from the University of Tennessee College of Veterinary Medicine is taking care of the three tigers.

CNNWire By Joe Wenzel

Saturday, October 31, 2020

Ontario dog believed to be first in Canada to test positive for COVID-19

Officials said that the risk of infection and illness in most domestic animals is low

[KATYA SLEPIAN](#) / Oct. 26, 2020 1:45 p.m. / [CANADA & WORLD](#) / [NEWS](#)

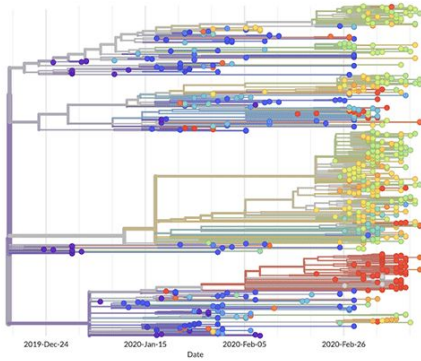


Denmark to cull mink herd over coronavirus mutation fears – here's what the science says

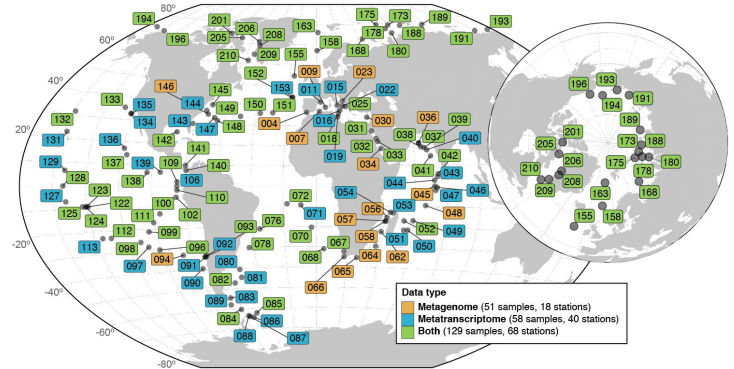
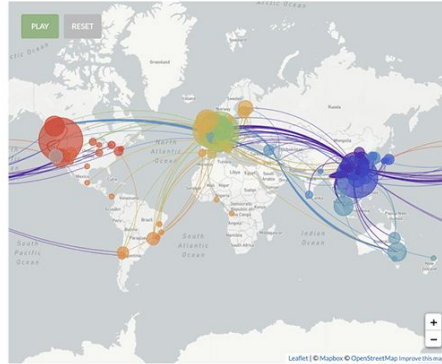
November 9, 2020 @ 4:56m EST



Enter sequencing efforts



Nextstrain



Tara Oceans, Salazar et al. (2019)





SRA

Sequence Read Archive (SRA) makes biological sequence data available to the research community to enhance reproducibility and allow for new discoveries by comparing data sets. The SRA stores raw sequencing data and alignment information from high-throughput sequencing platforms, including Roche 454 GS System®, Illumina Genome Analyzer®, Applied Biosystems SOLiD System®, Helicos Heliscope®, Complete Genomics®, and Pacific Biosciences SMRT®.

Search results

Items: 1 to 20 of 19964 NextSeq 500 paired end sequencing (ERR3407135)

Metadata Analysis (alpha) Reads Download

[NextSeq 500 paire](#)

1. 1 ILLUMINA (Illumina)
Accession: ERX34307

Filter: Find Filtered Download [What does it do?](#)

[What can the filter be applied to?](#)

[NextSeq 500 paire](#)

2. 1 ILLUMINA (Illumina)
Accession: ERX34307

[NextSeq 500 paire](#)

3. 1 ILLUMINA (Illumina)
Accession: ERX34307

< 1 1 346553 > View: biological reads technical reads

Reads (separated)

1. [ERR3407135.1](#) [ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:5421:
member: default

>gnl|SRA|ERR3407135.1.1 NB551234:144:HL523AFXY:1:11101:5421:1076 F (Biological)

ACCTGAGCGCGCAGCTCCAGTAAATCAAACGCGGCGGGAATTTGGGATGTTCCATCAGT
TTCCAGGCGCGTTTGCCCTGACGTCGCGACATGCGTAACTGAAGCTGCCAAATATCACGG
GTAAGCGTGGTAAGGCGTTTCGGGATCGCCA

2. [ERR3407135.2](#) [ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:2248:
member: default

>gnl|SRA|ERR3407135.1.2 NB551234:144:HL523AFXY:1:11101:5421:1076 R (Biological)

ATCAACAACAGCGGGAATACCACCTCTTCCAGCCGTTGTTCCAAACCAATACGCGTTAAT
TCACCGAAACCGCGACAGCGCAATGGAAACGCATCATTTGCCAGGTTTGCAGAATACGGA
AAACCGCATCCGAAACGAGATCGCGGTTAAT

3. [ERR3407135.3](#) [ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:2566:
member: default

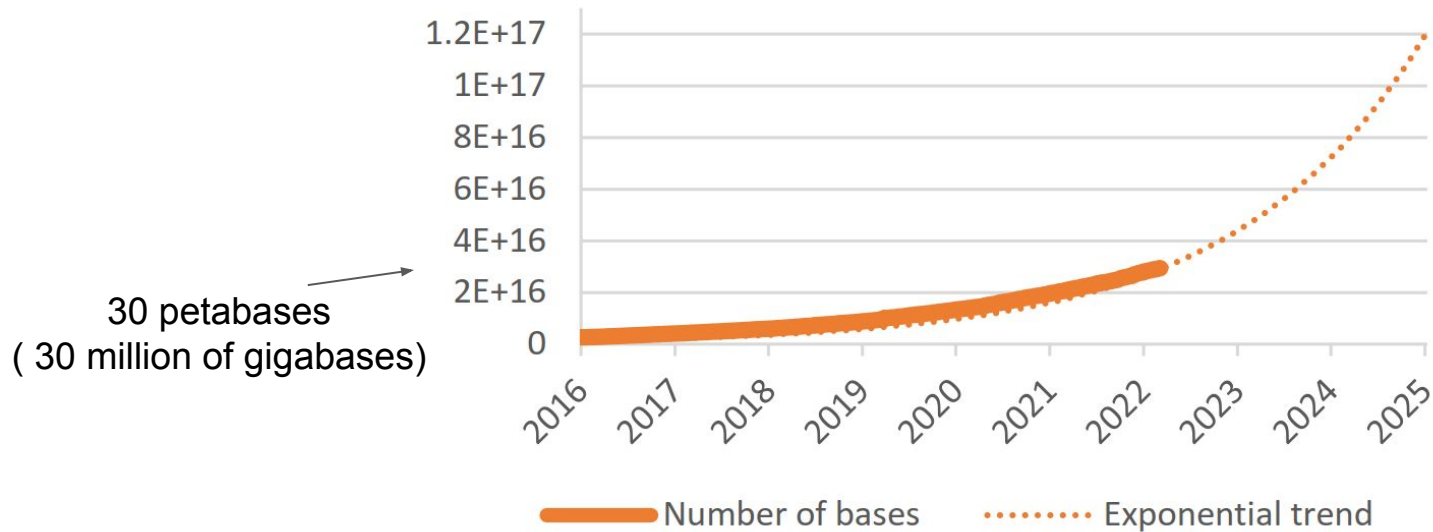
4. [ERR3407135.4](#) [ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:21195:
member: default

5. [ERR3407135.5](#) [ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:23504:
member: default

Growth of the Sequence Read Archive



NCBI SRA database : 30 PB



Data crypt

All the raw reads sleep
there, undisturbed



All RNA-seqs (2008-2020)
5 million samples, 10.2 Petabases

Downloading all
RNA-seq samples:

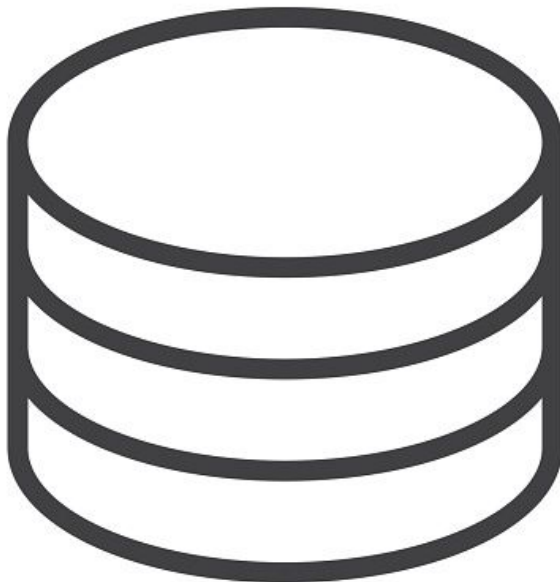


Guesstimate:

- How many years would it take to download 10 petabases (i.e. 10,000,000,000 MB) at 1 MB/sec?

Hint: ~30,000,000 seconds in a year

Downloading all
RNA-seq samples:



Google (10 petabytes divided by 1 megabyte) / (seconds per year) X

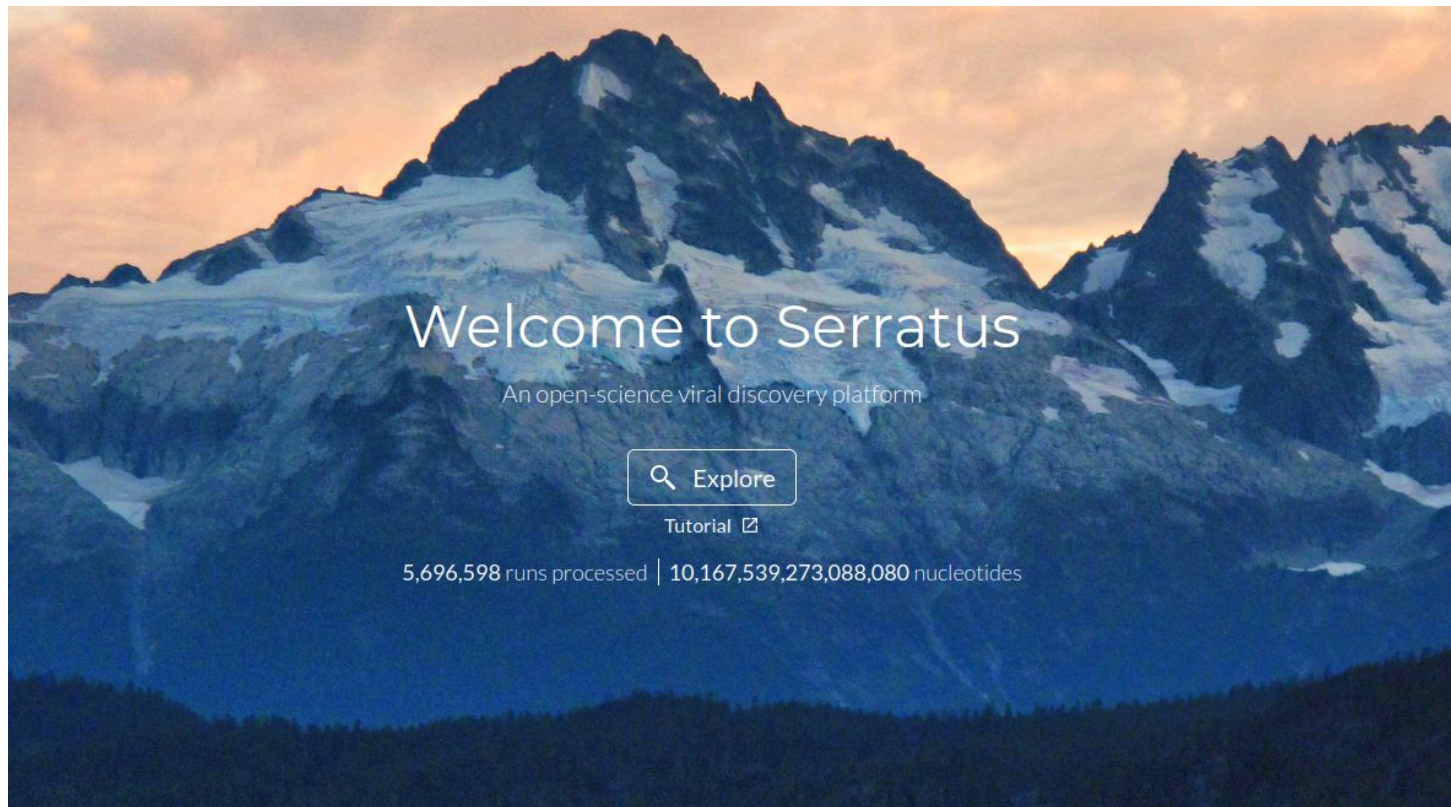
Tous Images Actualités Shopping Vidéos Plus Outils

Environ 291 000 résultats (0,57 secondes)

🕒 ((10 petabytes) / (1 megabyte)) / (seconds per year) =
316.887646408

years at 1 MB/s

Serratus: a cloud analysis of all RNA-seqs



Serratus: two analyses

1) Nucleotide alignments

all RNAseqs vs all RNA viral genomes

> Discovered new coronaviruses

2) Protein (translated) alignments

all RNAseqs vs a universal RNA virus gene

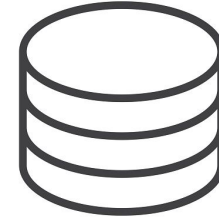
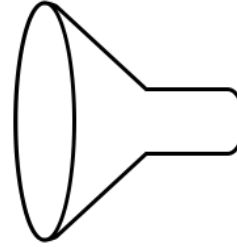
> Discovered 130,000 new RNA virus species

Analysis 1:



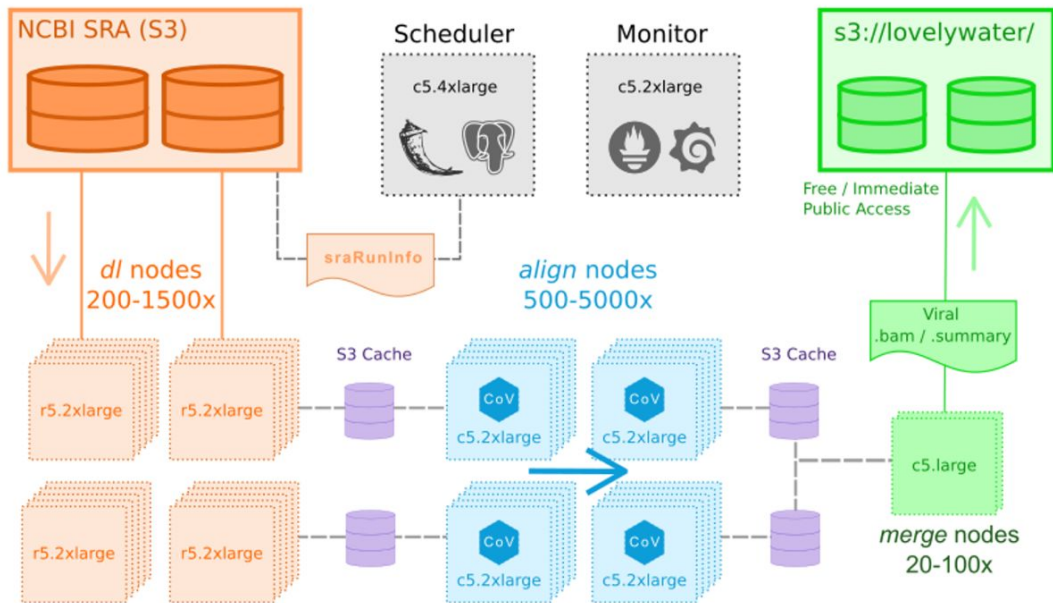
All RNA-seqs

**Serratus download &
align (bowtie2) to all
viral reference
genomes**



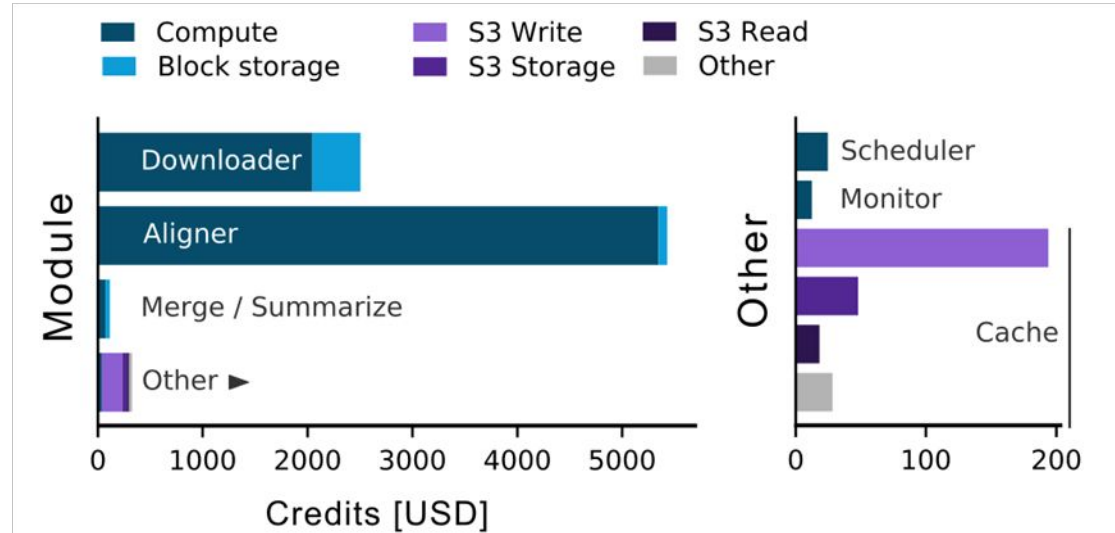
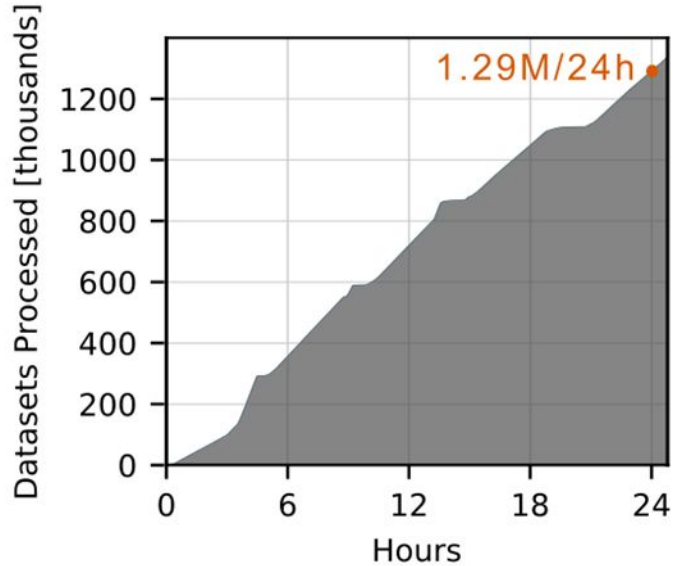
**55,715 CoV+
samples**

Serratus architecture



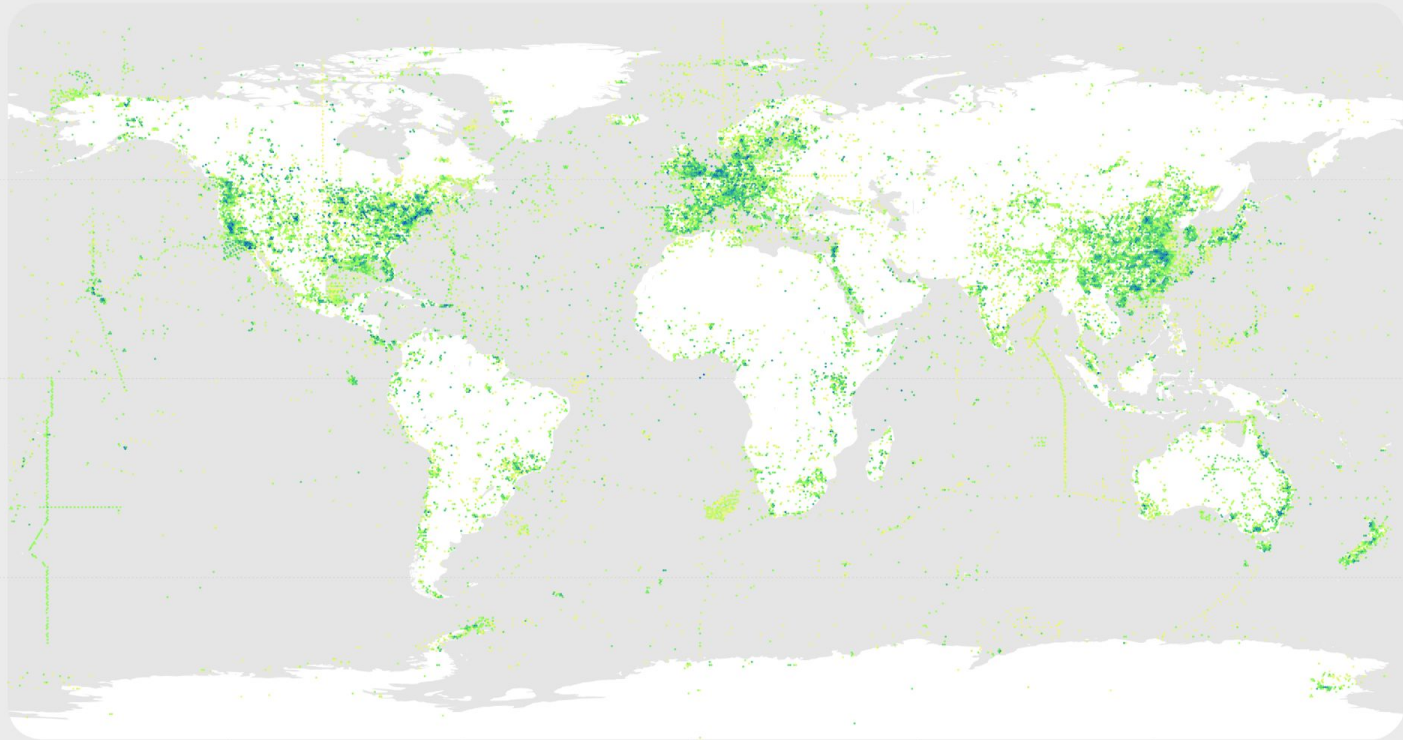
- Aggressively cost-optimized
- Native access to SRA on S3
- Dynamic scaling up to ~22,250s vCPU
- Open Source: GPLv3

Serratus performance & costs



1 million NGS libraries / day
\$0.005 / library

Geography of SRA samples



1 20 400 8000
Sequencing density (datasets)

Planetary DNA/RNA sequencing

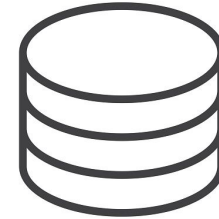
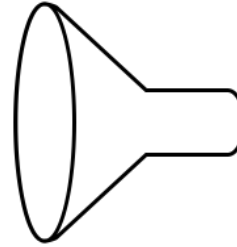


Analysis 2:



All RNA-seqs

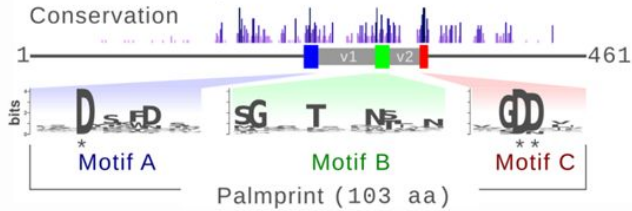
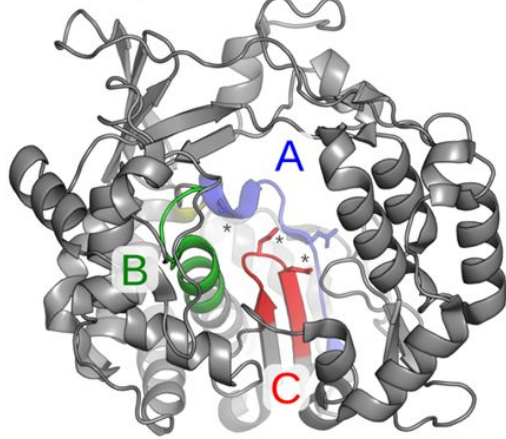
**Serratia download &
sensitive align
(DIAMOND2)
to all known versions of
RNA virus universal gene**



**aligned reads
(.bam files)**

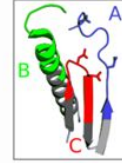
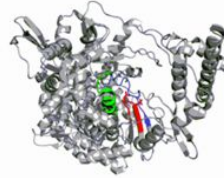
Analysis 2, search database: 15,060 known RNA viruses RdRP gene

Viral RdRP (Poliovirus)

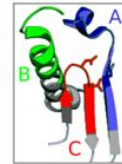
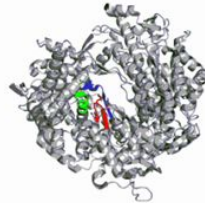


(Babaian & Edgar, 2021. bioRxiv)

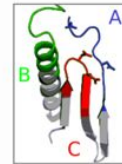
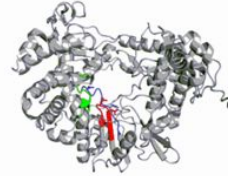
Coronaviridae



Reoviridae



Permutotetraviridae

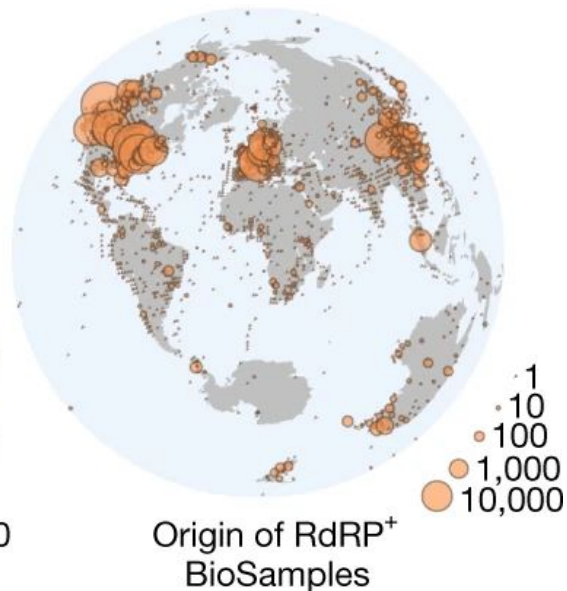
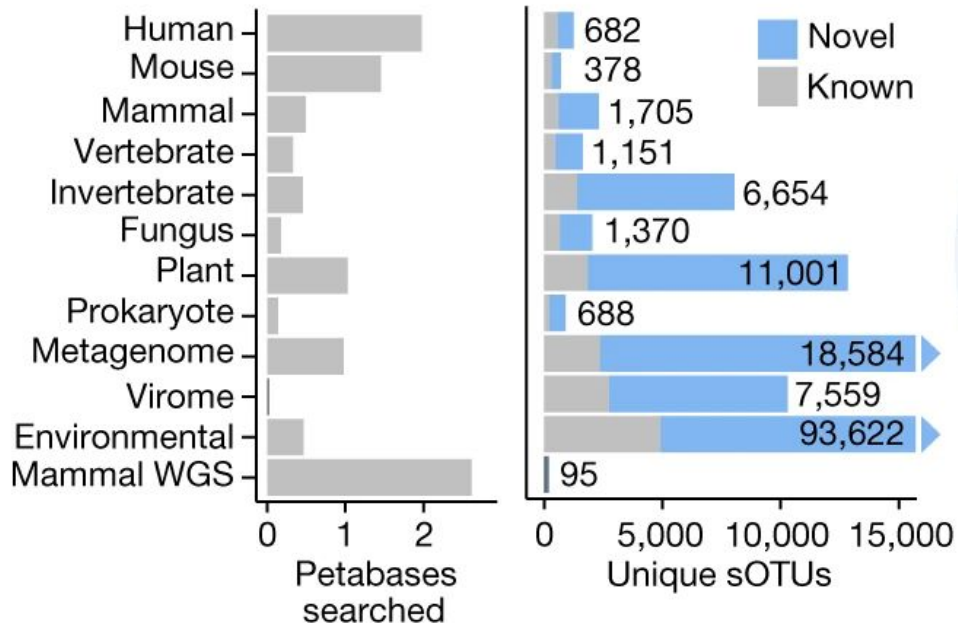


- RNA Virus “Palmprint”
- Species threshold: 90% amino-acid id

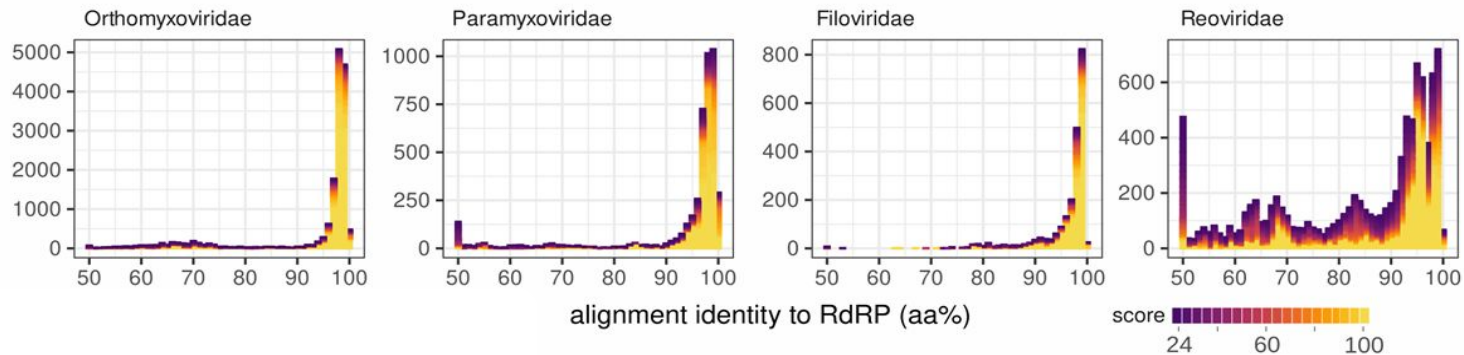
Assembly of all viral RdRPs (Analysis 2)

“Micro-assembly” of all RdRp-matching reads within each sample

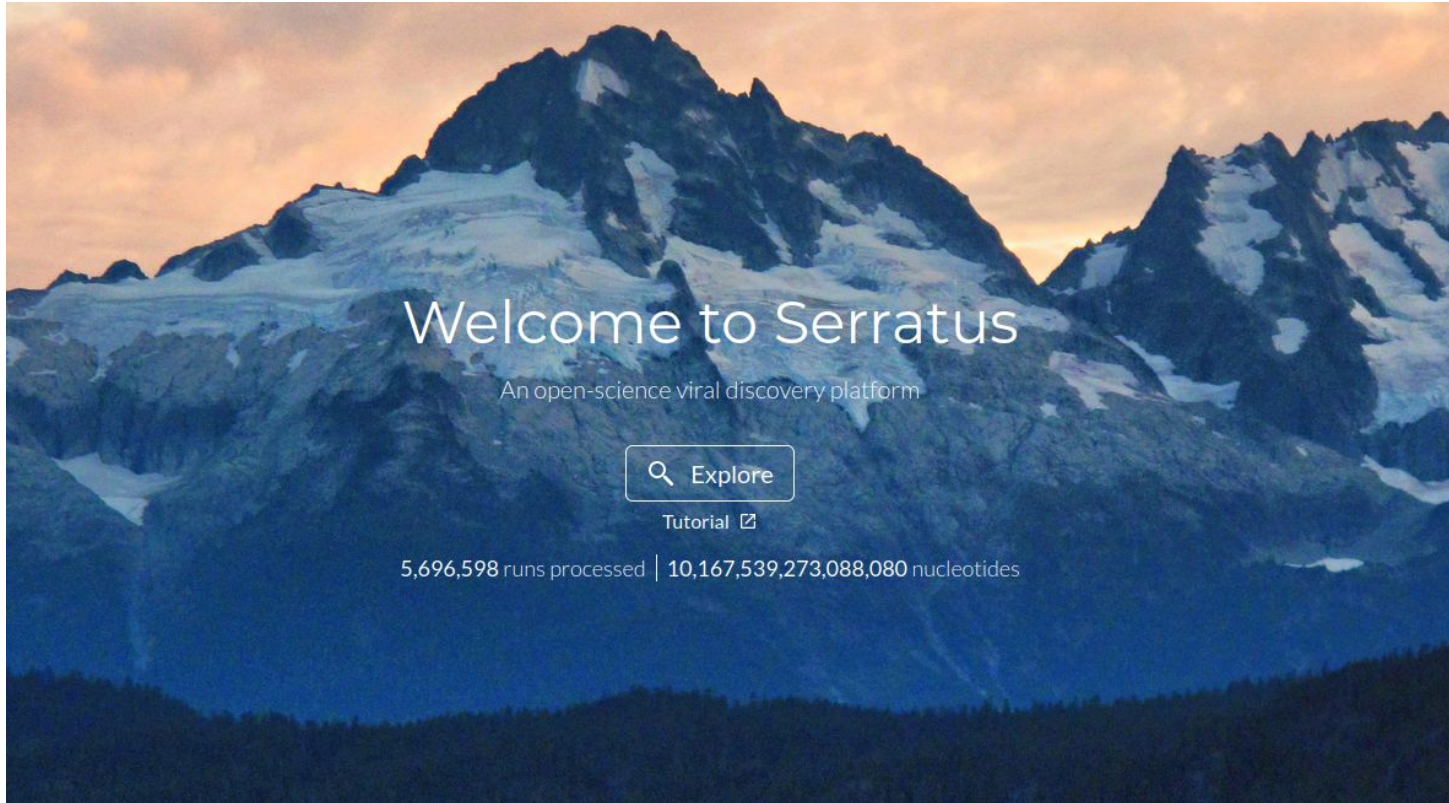
- SPAdes assembler & GNU parallel



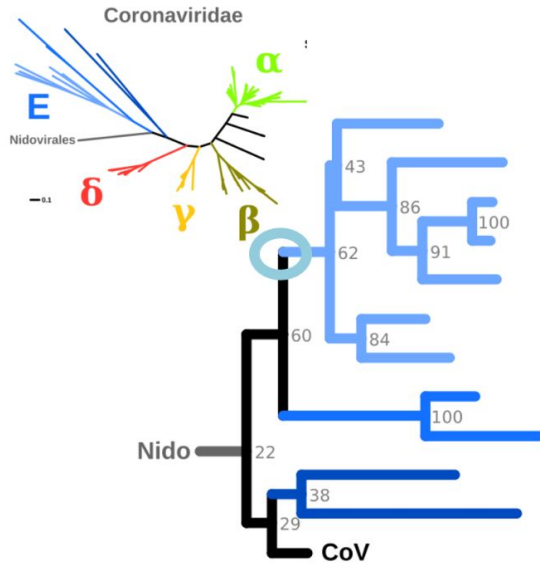
Number of samples



"petabase scale" on Google, `www.serratus.io`



Discovering new Coronaviruses

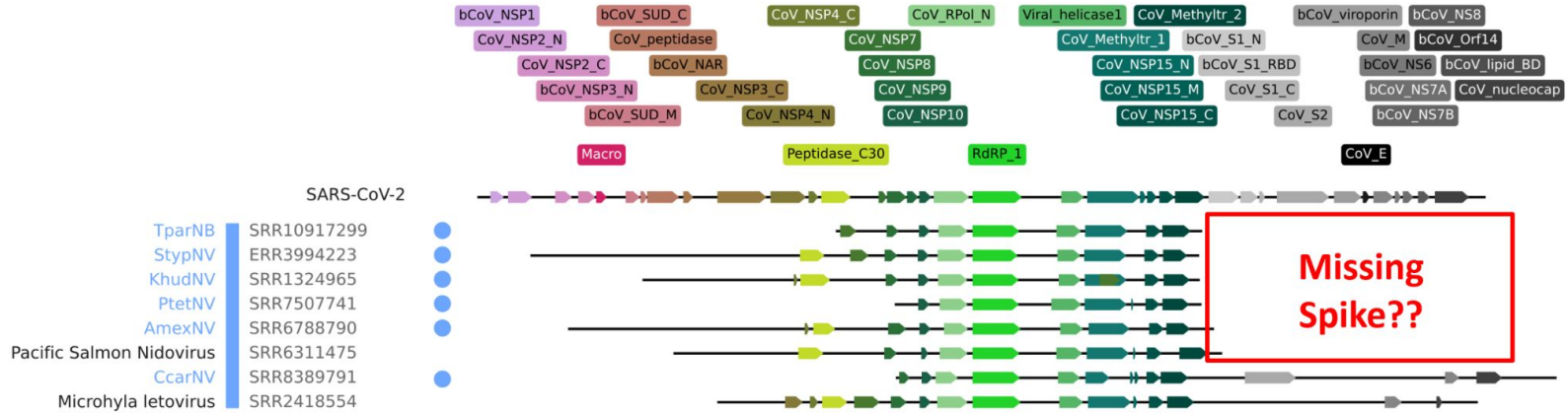


- AmexNV** SRR6788790
- PtetNV** SRR7507741
- HkudNV** SRR1324965
- StypNV** ERR3994223
- TparNV** SRR10917299
- Pacific Salmon Nidovirus
- AcaNV** SRR5997671
- SiINV** SRR12184956
- MalbNV** SRR10402291
- HtraNV** SRR8389791
- Microhylla Letovirus

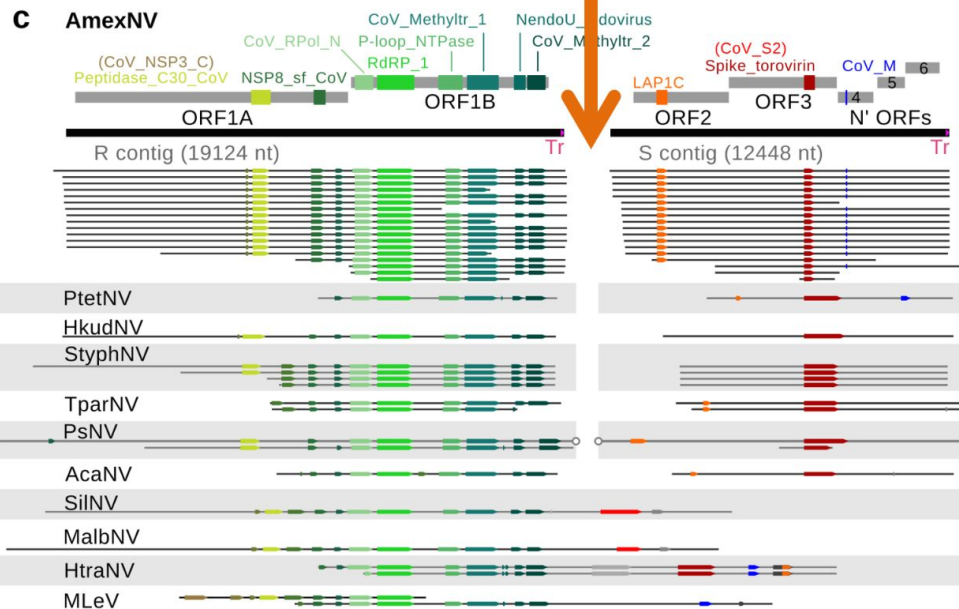
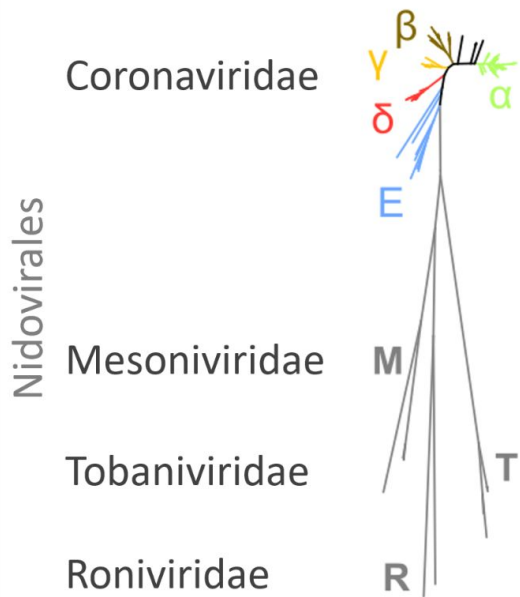


Slides credit: A. Babaian

Discovering new Coronaviruses



Segmented Coronaviruses?



Re-writing the textbook definition of a Coronavirus

Metagenome & metavirome assembly

Usually:

Reconstruct *all* the genomes

Computationally
intensive

Analysis 1:

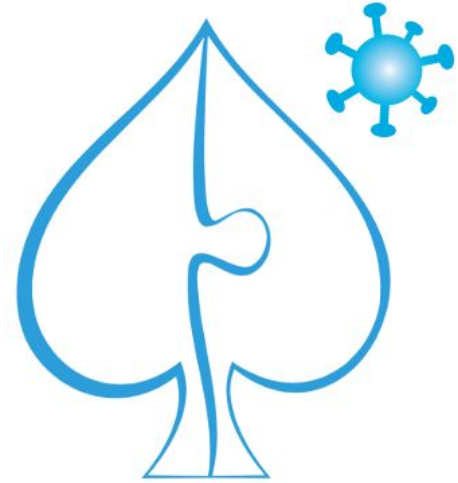
Reconstruct only CoV genome(s)

Computationally
intensive

Analysis 2:

Reconstruct only RdRP genes

Computationally
easy(/ier) :)



SPAdes assembler

rnaSPAdes

coronaSPAdes

(screenshot: P. Barbera)

<input type="checkbox"/>	Name	Instance ID	Instance Type	Availability Zone	Instance State
<input type="checkbox"/>	Compute	i-004fc86f836336d17	c5.9xlarge	us-east-2a	● running
<input type="checkbox"/>	Compute	i-01af64dd577f162b5	c5.9xlarge	us-east-2a	● running
<input type="checkbox"/>	Compute	i-0879ad68f76a4a54e	c5.9xlarge	us-east-2a	● running
<input type="checkbox"/>	Compute	i-094ddc9b931fde962	c5.9xlarge	us-east-2a	● running
<input type="checkbox"/>	Compute	i-0c8f6d93593531c32	c5.9xlarge	us-east-2a	● running
<input type="checkbox"/>	Compute	i-0e08ab6c5a3d0ce3f	c5.9xlarge	us-east-2a	● running
<input type="checkbox"/>	Compute	i-0ea10648adeeabf68	c5.9xlarge	us-east-2a	● running

AWS Batch framework for large-scale assembly

Peak:
~28,000 vCPUs

AWS Batch > Dashboard Last updated: 07:11:08 PM. Auto-refreshes every 60 seconds

Dashboard

Jobs overview

<p>RUNNABLE</p> <h1>450</h1>	<p>RUNNING</p> <h1>173</h1>	<p>SUCCEEDED</p> <h1>48</h1>	<p>FAILED</p> <h1>817</h1>
------------------------------	-----------------------------	------------------------------	----------------------------

Job queue overview

Job queue	SUBMITTED	PENDING	RUNNABLE	STARTING	RUNNING	SUCCEEDED	FAILED
RayanUnitigsBatchProcessingJobQueue	0	0	0	0	0	● 0	⊗ 0
RayanSerratusDLBatchProcessingJobQueue	0	0	0	0	0	● 0	⊗ 0
RayanSerratusAssemblyBatchJobQueue	0	0	450	7	173	● 48	⊗ 817

But, for all-RdRp assembly (Analysis 2)..

With a single “bigger” instance (c6a.48xlarge, 192 cores)

10^5 viral species known, 10^8 left to discover

What's next?

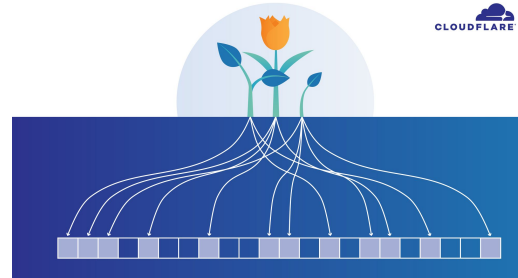
- DNA viruses
- Lower homology detection with known RdRPs
 - Replacing Bowtie 2 / Diamond by ...?
- A global **index of the SRA**
 - nearly feasible with k-mers already
 - would only support exact search

Deep embedding and alignment of protein sequences

Felipe Llinares-López, Quentin Berthet, Mathieu Blondel,
Olivier Teboul and Jean-Philippe Vert*

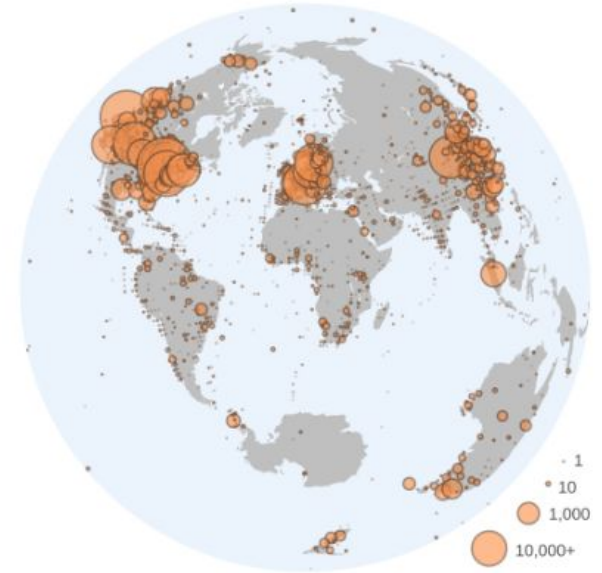
Google Research, Brain team, Paris, France

November 15, 2021



Summary:

- **132,260 novel RNA virus species**
- **1 new group of CoV-like segmented virus**
- **hyper-compressed (300-500 nt) Zetaviruses**
53 novel deltaviruses (cancer),
252 huge phages, ..



All our data is accessible:

<https://github.com/ababaian/serratus/wiki/Access-Data-Release>

7 TB of alignments and assemblies

More details:

<https://www.nature.com/articles/s41586-021-04332-2>

<https://github.com/ababaian/serratus/>

Chat with us on Slack:

https://join.slack.com/t/hackseq-rna/shared_invite/zt-ewlzh9qf-SiNkxvVTJflcutFN0h5jIQ

Petabase-scale sequence alignment catalyses viral discovery

[Robert C. Edgar](#), [Jeff Taylor](#), [Victor Lin](#), [Tomer Altman](#), [Pierre Barbera](#), [Dmitry Meleshko](#), [Dan Lohr](#), [Gherman Novakovsky](#), [Benjamin Buchfink](#), [Basem Al-Shayeb](#), [Jillian F. Banfield](#), [Marcos de la Peña](#), [Anton Korobeynikov](#), [Rayan Chikhi](#) & [Artem Babaian](#) 

Nature **602**, 142–147 (2022) | [Cite this article](#)

32k Accesses | **1024** Altmetric | [Metrics](#)

Abstract

Public databases contain a planetary collection of nucleic acid sequences, but their systematic exploration has been inhibited by a lack of efficient methods for searching this corpus, which (at the time of writing) exceeds 20 petabases and is growing exponentially¹. Here we developed a cloud computing infrastructure, Serratus, to enable ultra-high-throughput sequence alignment at the petabase scale. We searched 5.7 million biologically diverse samples (10.2 petabases) for the hallmark gene RNA-dependent RNA polymerase and identified well over 10⁵ novel RNA viruses, thereby expanding the number of known species by roughly an order of magnitude. We characterized novel viruses related to coronaviruses, hepatitis delta virus and huge phages, respectively, and analysed their environmental reservoirs. To catalyse the ongoing revolution of viral discovery, we established a free and comprehensive database of these data and tools. Expanding the known sequence diversity of viruses can reveal the evolutionary origins of emerging pathogens and improve pathogen surveillance for the anticipation and mitigation of future pandemics.

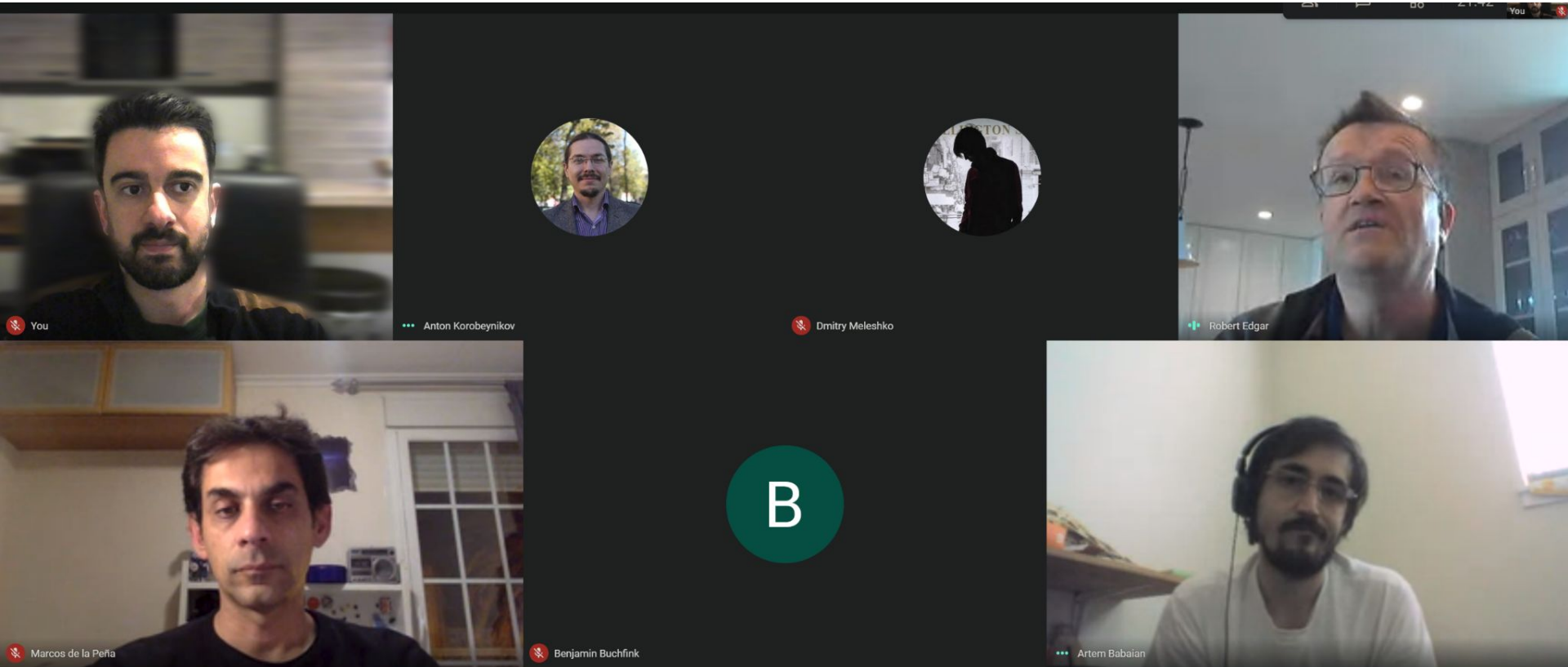


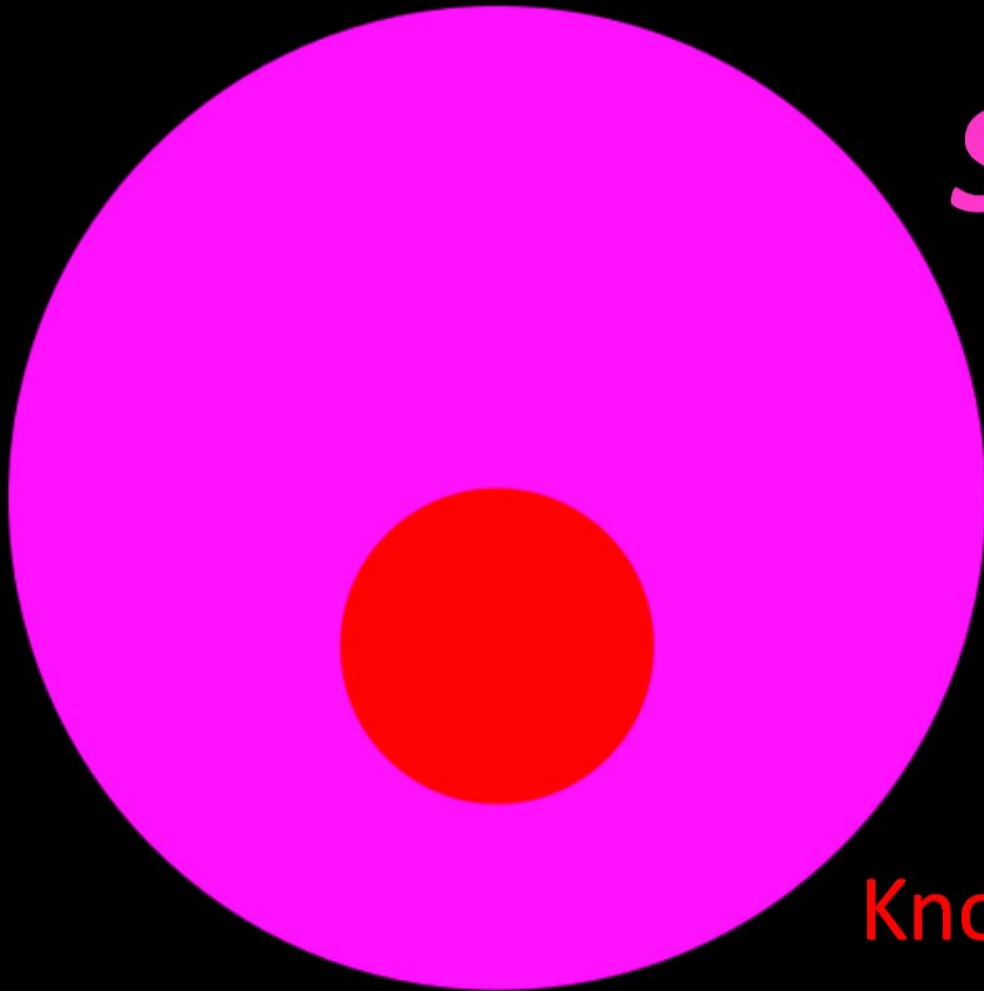
Digital Collaboration

- Anton Korobeynikov (St. Petersburg)
- Artem Babaian (Vancouver)
- Basem Al-Shayeb (Berkeley)
- Benjamin Buchfink (Tubingen)
- Dan Lohr (Boulder)
- Dmitry Meleshko (Ithaca)
- Gherman Novakovsky (Vancouver)
- Jeff Taylor (Vancouver)
- Jillian F. Banfield (Berkeley)
- Marcos de la Pena (Valencia)
- Pierre Barbera (Heidelberg)
- Rayan Chikhi (Paris)
- Robert C. Edgar (Sonoma)
- Tomer Altman (San Francisco)
- Victor Lin (Gainsville)

All equal contributions

We never met IRL





Serratus

Known RNA Virome

Earth's Virome

We are here



Summary

- Lots of genomics data
- Many great analyses could be made
- Cloud helps at the largest scale

Credits

Some of the people who initiate these “small-group but large-scale” analyses:

C. Titus Brown, Ben Langmead, Artem Babaian, Rob Finn, Adam Phillippy, Andre Kahles, Zamin Iqbal, Carl Kingsford, Rob Patro, Christina Boucher, Pierre Peterlongo, Olivier Jaillon, Dominique Lavenier, Antoine Limasset, Camille Marchet, Daniel Gautheret, Thérèse Commes, and many others I forget to mention

Additional credits:

k-mer people

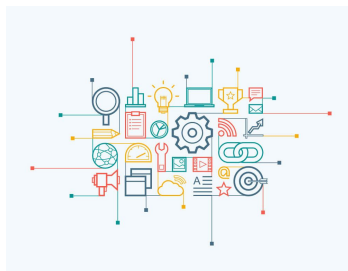
Slides advise: Michel Attafeu, Sophie Shaw, Na, Cami, Karin, M, Malfoy

Sequence Bioinformatics

@ Institut Pasteur



Genomes &
metagenomes
assembly



Algorithms and
data structures
on k-mers



Sequence
search in very
large datasets



Pangenomics

Vielen Dank für ihre
Aufmerksamkeit!



Bonus slides

Part 1.5:
*“Spill the beans!
Where is this
magical bigger data
you speak of?”*



GenBank



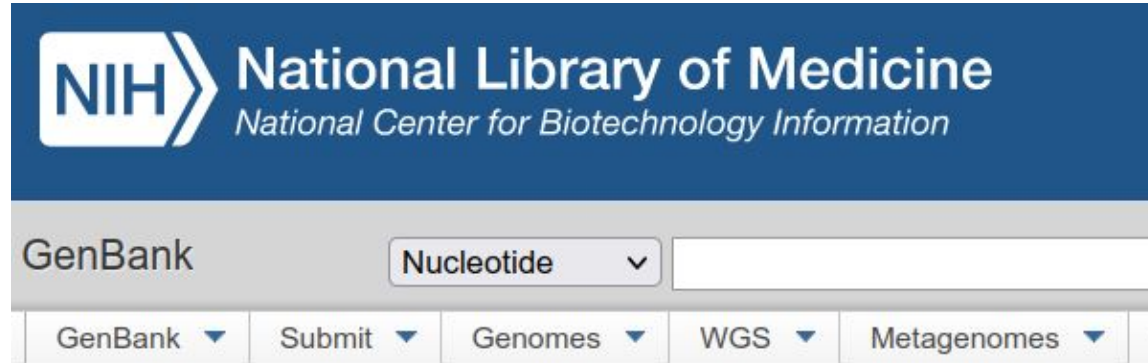
Type: assemblies

Size: 1.2 TB ([April 2022](#))

Diversity: high

Particularity: all sequences are *annotated*

NCBI WGS



Whole Genome Shotgun Submissions

What is Whole Genome Shotgun (WGS)?

Whole Genome Shotgun (WGS) projects are genome assemblies of incomplete genomes of eukaryotes that are generally being sequenced by a whole genome shotgun strategy.

Type: assemblies


Size: 16 TB ([April 2022](#))

Diversity: high

Difference with GenBank: sequences are not necessarily annotated

NCBI SRA

SRA [Advanced](#) [Help](#)



SRA

Sequence Read Archive (SRA) makes biological sequence data available to the research community to enhance reproducibility and allow for new discoveries by comparing data sets. The SRA stores raw sequencing data and alignment information from high-throughput sequencing platforms, including Roche 454 GS System®, Illumina Genome Analyzer®, Applied Biosystems SOLiD System®, Helicos Heliscope®, Complete Genomics®, and Pacific Biosciences SMRT®.

Search results

Items: 1 to 20 of 19964 **NextSeq 500 paired end sequencing (ERR3407135)**

[Metadata](#) [Analysis \(alpha\)](#) [Reads](#) [Download](#)

[NextSeq 500 paire](#)

1. **1 ILLUMINA (Illumina)**
Accession: ERX34307

Filter: [What does it do?](#)

[What can the filter be applied to?](#)

[NextSeq 500 paire](#)

2. **1 ILLUMINA (Illumina)**
Accession: ERX34307

< 1 1 346553 >

View: biological reads technical reads

Reads (separated)

1. [ERR3407135.1 ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:5421:
member: default

>gnl|SRA|ERR3407135.1.1 NB551234:144:HL523AFXY:1:11101:5421:1076 F (Biological)

ACCTGAGCGCGCAGCTCCAGTAAATCAAACGCGCGCGGAAATTTGGGATGTTCCATCAGT

TTCAGGGCGCTTTGCCCTGACGTGCGGACATGCGTAACTGAAGCTGCCAAATATCACGG

GTAAGCGTGGTAAGCGCTTTCGGATCGCCA

< >

2. [ERR3407135.2 ERS3549882](#)
name: NB551234:144:HL523AFXY:1:11101:2248:
member: default

>gnl|SRA|ERR3407135.1.2 NB551234:144:HL523AFXY:1:11101:5421:1076 R (Biological)

ATCAACAACAGCGGGAATACCACCTCTTCCAGCCGTTGTTTCCAAACAAATACCGGTTAAT

TCACCGAAACCGGACAGCGCAATGGAACGCATCATTTGCCGAGGTGTTGCAGAAATACGGA

AAACCGCATCCGAAACGAGATGCGCGTTAAT

[NextSeq 500 paire](#)

3. **1 ILLUMINA (Illumina)**
Accession: ERX34307

3. [ERR3407135.3 ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:2566:
member: default

4. [ERR3407135.4 ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:21199:
member: default

5. [ERR3407135.5 ERS3549882](#)

name: NB551234:144:HL523AFXY:1:11101:23504:
member: default

What is STAT good for?


- Say you have a model organism
 - Search for all sequencing data containing that organism
 - Find host-associations
 - Find co-occurrences with other species
- Say you have a set of samples
 - Determine set of species in them
 - Find other similar samples
- etc..

NCBI STAT

A taxonomic index of all sequencing data

Method | [Open Access](#) | [Published: 20 September 2021](#)

STAT: a fast, scalable, MinHash-based k -mer tool to assess Sequence Read Archive next-generation sequence submissions

[Kenneth S. Katz](#) , [Oleg Shutov](#), [Richard Lapoint](#), [Michael Kimelman](#), [J. Rodney Brister](#) & [Christopher O'Sullivan](#)

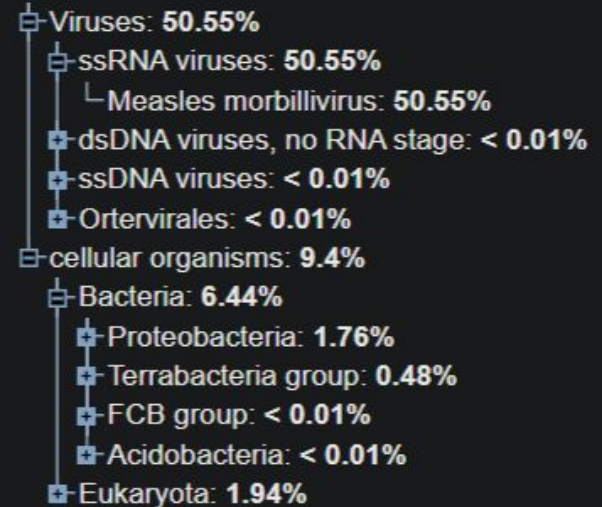
[Genome Biology](#) **22**, Article number: 270 (2021) | [Cite this article](#)

"we have processed more than 27.9 Peta base pairs from runs"

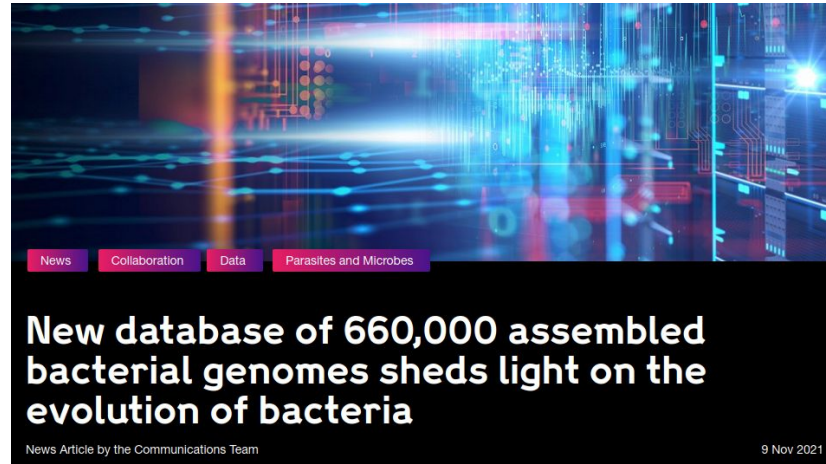
Taxonomy Analysis

Unidentified reads: 40.04%

Identified reads: 59.96%



Blackwell, ..., Iqbal's 661k bacterial genomes collection



Type: assemblies

Size: 2.5 TB

Diversity: medium

dBG? yes

Results: Pangenome graph of 661,405 bacterial genomes

Data from Blackwell et al, 2021:

2.9T 661k_assemblies.fa

1.6T 661k_assemblies.fa.lz4

```
rust-mdbg -k 10 -l 12 --density 0.001 --minabund 1 661k_assemblies.fa.lz4
```

Largest 5
connected
components:



Taxons in component

18

22

4

22

10

Dominant species

*Mycobacterium
tuberculosis*

*Salmonella
enterica*

*Burkholderia
gladioli*

*Pseudomonas
protegens*

*Cupriavidus
alkaliphilus*

Many others (often metagenomic)



Rayan Chikhi
@RayanChikhi

In this thread we are releasing a concatenated FASTA file of all assemblies produced by Serratus: 59,256 SRA accessions, 5.9 terabases total.



Uros @uki156 · Mar 22

Replying to @RayanChikhi

When you said "in this thread we are releasing", I was hoping you were actually going to tweet out the entire thing

Resource | [Open Access](#) | [Published: 20 July 2020](#)

A unified catalog of 204,938 reference genomes from the human gut microbiome

[Alexandre Almeida](#) ✉, [Stephen Nayfach](#), [Miguel Boland](#), [Francesco Strozzi](#), [Martin Beracochea](#), [Zhou Jason Shi](#), [Katherine S. Pollard](#), [Ekaterina Sakharova](#), [Donovan H. Parks](#), [Philip Hugenholtz](#), [Nicola Segata](#), [Nikos C. Kyrpides](#) & [Robert D. Finn](#) ✉

MGNify: a database of assemblies of metagenome studies from ENA searchable by metadata

EMBL-EBI | MGNify

MGNify

Submit, analyse, discover and compare microbiome data

Search MGNify

Example searches: [Tara oceans](#), [MGYS00000410](#), [Human Gut](#)

Anton Korobeynikov
23:42 Hier ✓

Lots of stuff in MGNify: <https://ebi-metagenomics.github.io/blog/>

[Overview](#) [Submit data](#) [Text search](#) [Sequence search](#) [Browse data](#)

Search by

[Text search](#) →

Name, biome, or keyword

[Sequence search](#) →

Sequence search

Or by data type

xxx Analysis types

356039 amplicon

28873 assemblies

2039 metabarcoding

33827 metagenomes

2205 metatranscriptomics

Public data

8696 studies

661121 samples

444172 analyses

9421 genomes in 4 MAG catalogues