seqera

# Data in life sciences

Views and perspectives on
challenges in the field

**seqera**

# Data in life sciences

- Background and introduction
- Big data keeps getting bigger
- Artificial intelligence
- Understanding and trust
- On the importance of being open
- Future challenges
- Conclusion

Background and introduction

# Background and introduction

**Phil Ewels, PhD**

Product Manager for Open Source

phil.ewels@seqera.io

seqera

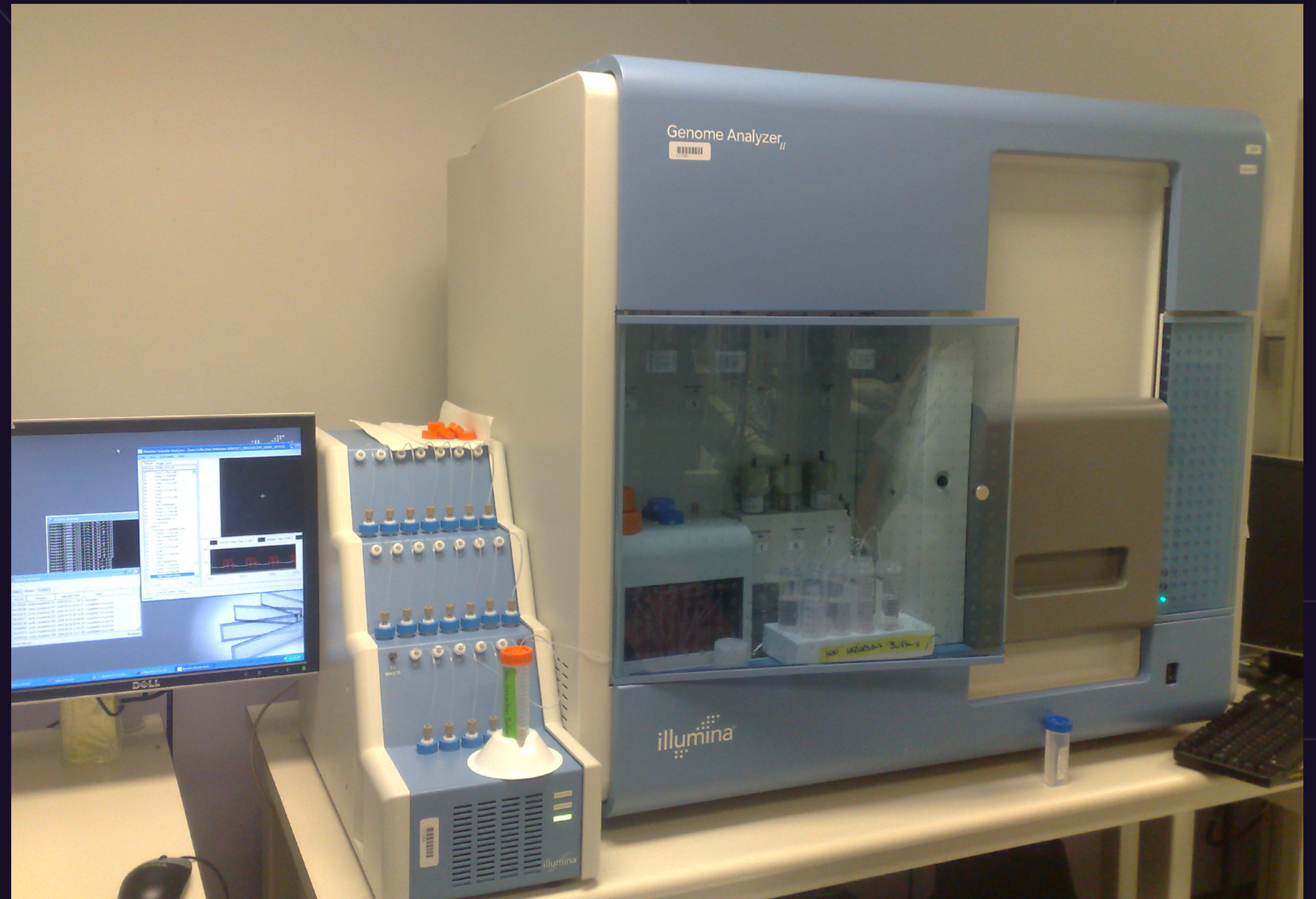# Background and introduction

PhD in epigenetics 2008-12
University of Cambridge

Dawn of NGS:
Big data came to biology

1 × 35bp
22M reads / run
Introduction to bioinformatics



Illumina Genome Analyzer IIx

# Background and introduction

Postdoc in bioinformatics

Babraham Institute, Cambridge

Got the bug for building
open-source software

Wrote my own workflow tool

Started writing data-vis scripts

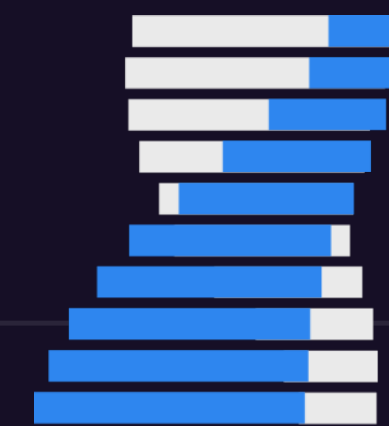# Background and introduction

Moved to Sweden in 2014
Joined NGI at SciLifeLab

Started building software to
handle the scale of data

SciLifeLab

NATIONAL
GENOMICS
INFRASTRUCTURE

# Background and introduction

Started building software to handle the scale of data

Wrote and released MultiQC
Adopted Nextflow, started nf-core

SciLifeLab

NATIONAL GENOMICS INFRASTRUCTURE

nextflow

multiqc

nf-core

# Background and introduction

Joined Seqera in 2022 as
employee #21

Set up the community team,
now product manager for OSS*

**seqera**
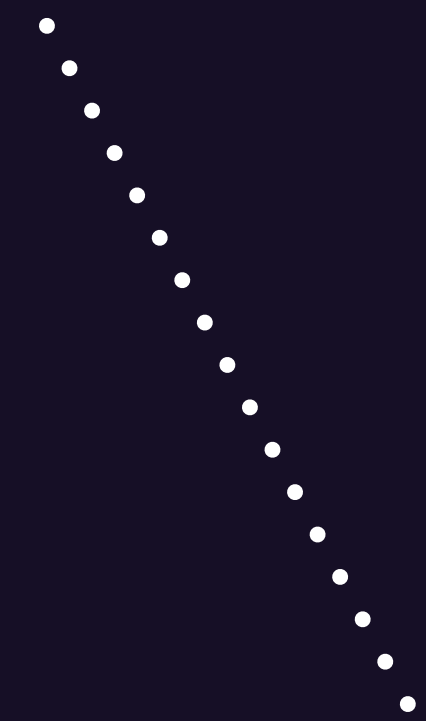the modern biotech stack
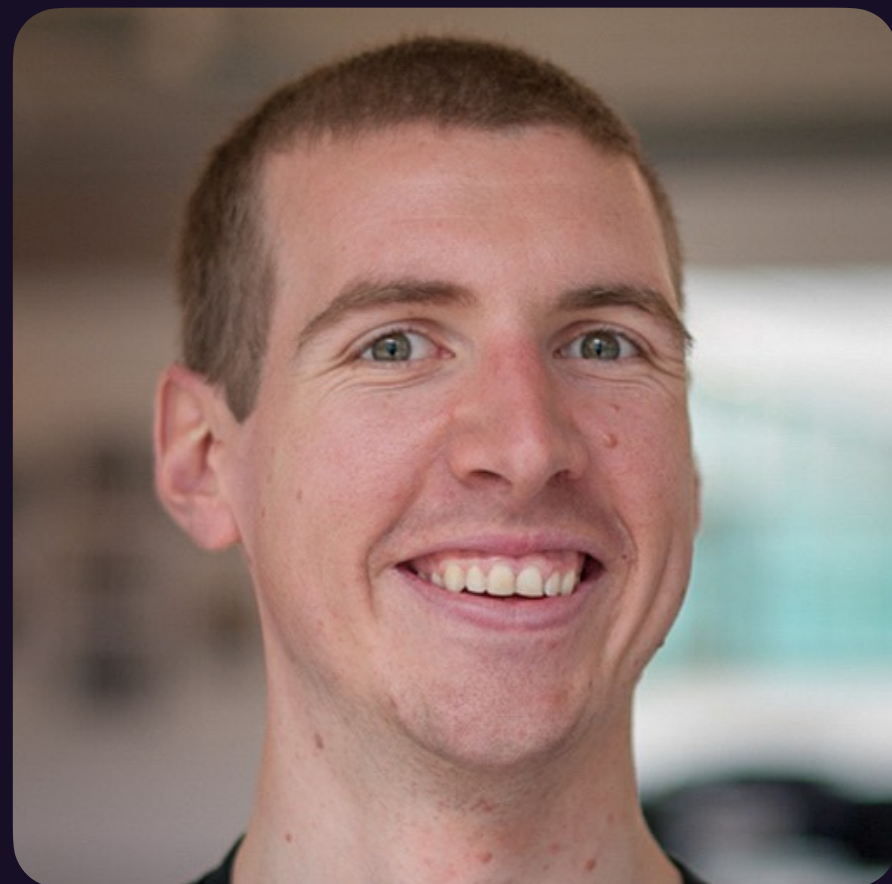
**nextflow**
workflow management

**multiqc**
reporting and analytics

**fusion**
cloud native file-system

**wave**
container provisioning

**nf-core**

# Background and introduction

seqera

# Big data keeps on getting bigger

# Big data keeps getting bigger

NCBI Sequence Read Archive
Storage Footprint

31.4 PB

Bytes stored

4×10^+16

3×10^+16

2×10^+16

1×10^+16

06/05/2007    10/02/2011    07/05/2014    03/31/2017    12/29/2019    09/24/2022

# Big data keeps getting bigger

GenBank
Database Size

Bases

3×10^+13

2.25×10^+13

1.5×10^+13

7.5×10^+12

0×10^+00

GenBank          WGS

Dec 1982     Jun 1993     Apr 2000     Dec 2006     Aug 2013     Apr 2020

# Big data keeps getting bigger



PDB
Database Size

Entries

250,000

187,500

125,000

62,500

0

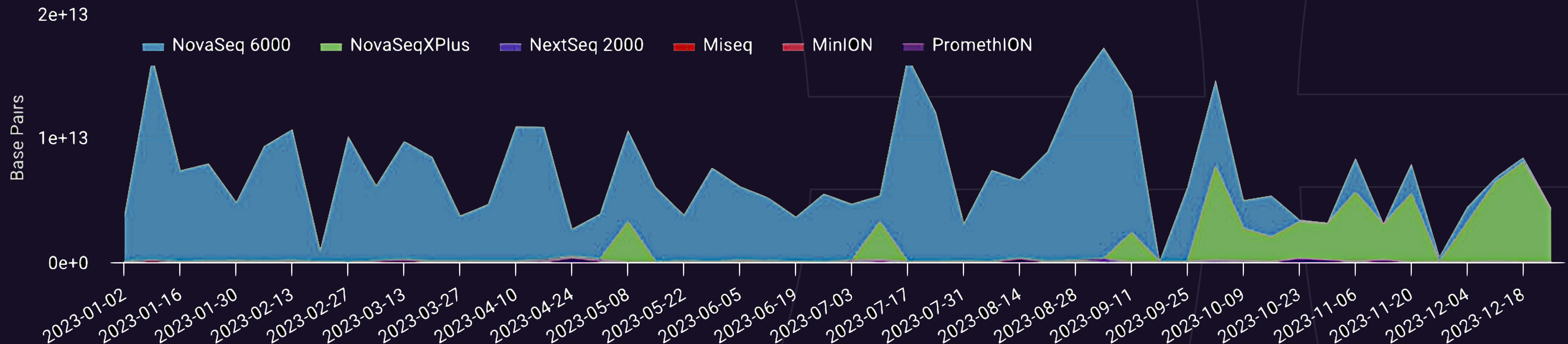1977/11    1985/04    1991/07    1998/01    2004/07    2010/10    2017/01    2023/04

# Big data keeps getting bigger

Reproducible analysis of
genomics data at scale

~1 Tbp sequencing per day in 2023   *(from the Stockholm site only)*

SciLifeLab

NATIONAL
GENOMICS
INFRASTRUCTURE

nextflow

# Big data keeps getting bigger

Reproducible analysis of
any kind of data at scale

nextflow

## Compute platforms

slurm | grid | kubernetes | Google GKE | Amazon EKS

AWS BATCH | Azure Batch | Google Life Sciences | Altair | PBS Works™ | IBM Spectrum LSF

# Big data keeps getting bigger

Reproducible analysis of
any kind of data at scale

**nextflow**

Scalable

Portable

Reproducible

## Compute platforms

slurm | grid | kubernetes | Google GKE | Amazon EKS

AWS BATCH | Azure Batch | Google Life Sciences | Altair | PBS Works™ | IBM Spectrum LSF

## Storage and data

NFS | aws | S3 | Azure Cloud Storage | Google Cloud | SQL

## Container technologies / SCM

docker | CONDA | Singularity | GitHub | Bitbucket

# Artificial Intelligence

# Artificial Intelligence



ML + AI arXiv
papers per month

# Artificial Intelligence

**Number of notable machine learning models by sector, 2003–23**
Source: Epoch, 2023 | Chart: 2024 AI Index report



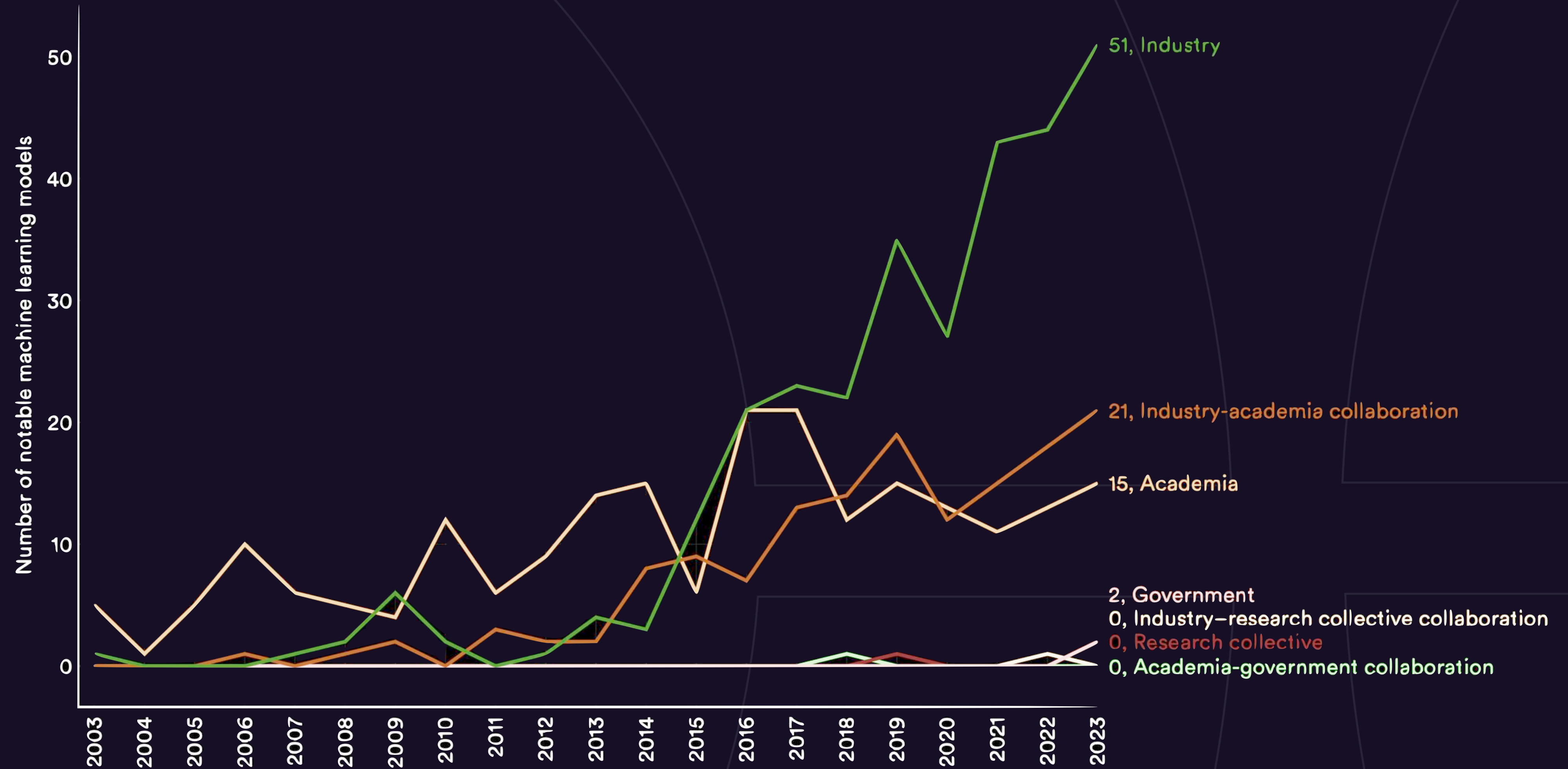- 51, Industry
- 21, Industry–academia collaboration
- 15, Academia
- 2, Government
- 0, Industry–research collective collaboration
- 0, Research collective
- 0, Academia–government collaboration

# Artificial Intelligence

## nextflow

nextflow-io/nf-gpt

```groovy
def query = '''
        Extract information about a person from In 1968, amidst the fading echoes of Independenc
        a child named John arrived under the calm evening sky. This newborn, bearing the surname
        marked the start of a new journey.
        '''


def response = gptPromptForData(query, schema: [firstName: 'string', lastName: 'string', birthDa

println "First name: ${response[0].firstName}"
println "Last name: ${response[0].lastName}"
println "Birth date: ${response[0].birthDate}"
```

# Artificial Intelligence



Future: Explanation of errors?

⚠ The error message `command not found` with an exit status of **127** typically indicates that the command (in this case, `fastqc`) is not found in the system's **PATH.** This means the system does not know where to locate the `fastqc` executable to run it.

ⓘ This is likely because you have not specified a container or environment profile with the `-profile` option when running the pipeline. This is essential as it ensures that all the required tools, such as `fastqc`, are available and correctly configured in a reproducible manner.

💡 The best way to fix this is by specifying a profile that suits your system's capabilities. Commonly used profiles include `docker`, `singularity`, `podman`, or `conda`. Using one of these profiles ensures that the necessary software environments are set up automatically, avoiding the need to manually manage dependencies and **PATH** configurations.

▢ For example, if you have Docker installed, you can modify your command like this:

```
nextflow run nf-core/rnaseq --input mysamples.csv --outdir ./results -profile docker
```

seqera

Big data + AI results are only
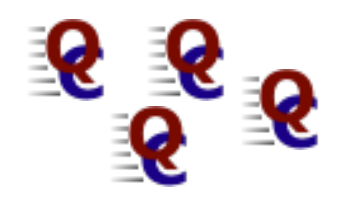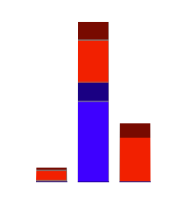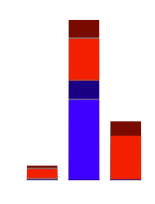useful if they can be
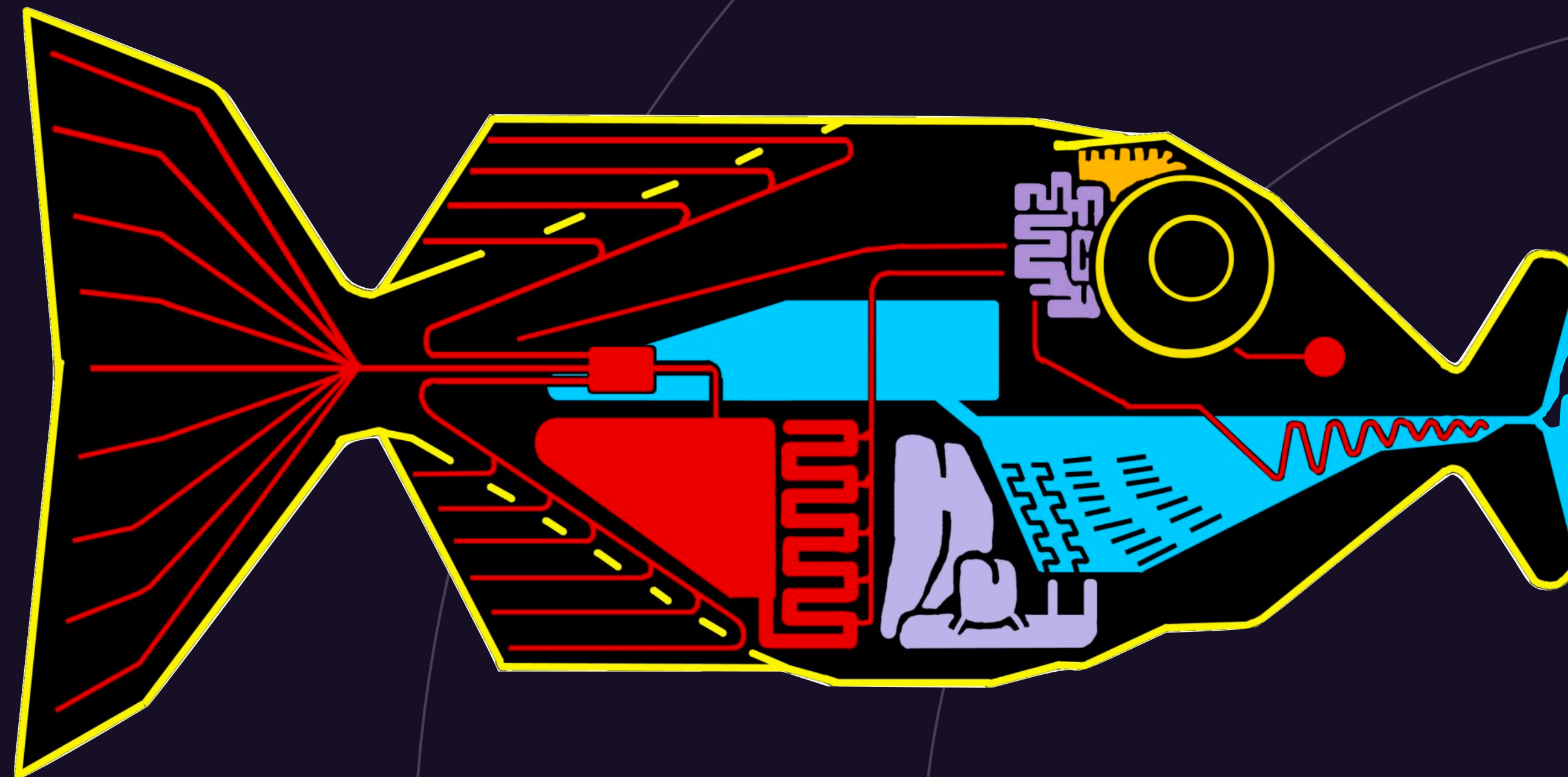understood and trusted

DNA-club 2015

# Understanding and trust

Bioinformatics log outputs →

Human-readable report

multiqc

# Understanding and trust



multiqc

Visualises metrics across many tools and many samples

Collects software versions automatically

## Software Versions

Software Versions lists versions of software tools extracted from file contents.

📋 Copy table

| Group | Software | Version |
|---|---|---|
| FASTQC | fastqc | 0.11.9 |
| STAR_ALIGN | star | 2.6.1d |
|  | samtools | 1.10 |
|  | gawk | 5.1.0 |
| SALMON_QUANT | salmon | 1.10.1 |

# Understanding and trust

SWEDAC

FDA U.S. FOOD & DRUG ADMINISTRATION

Five Safes

RO-Crate

BioCompute Objects

# Understanding and trust

**nextflow**

nextflow-io/nf-prov

Automatically generate standards-compliant provenance reports

**RO-Crate**

**nf-core**

Automatically find contributors and link to GitHub and ORCiD
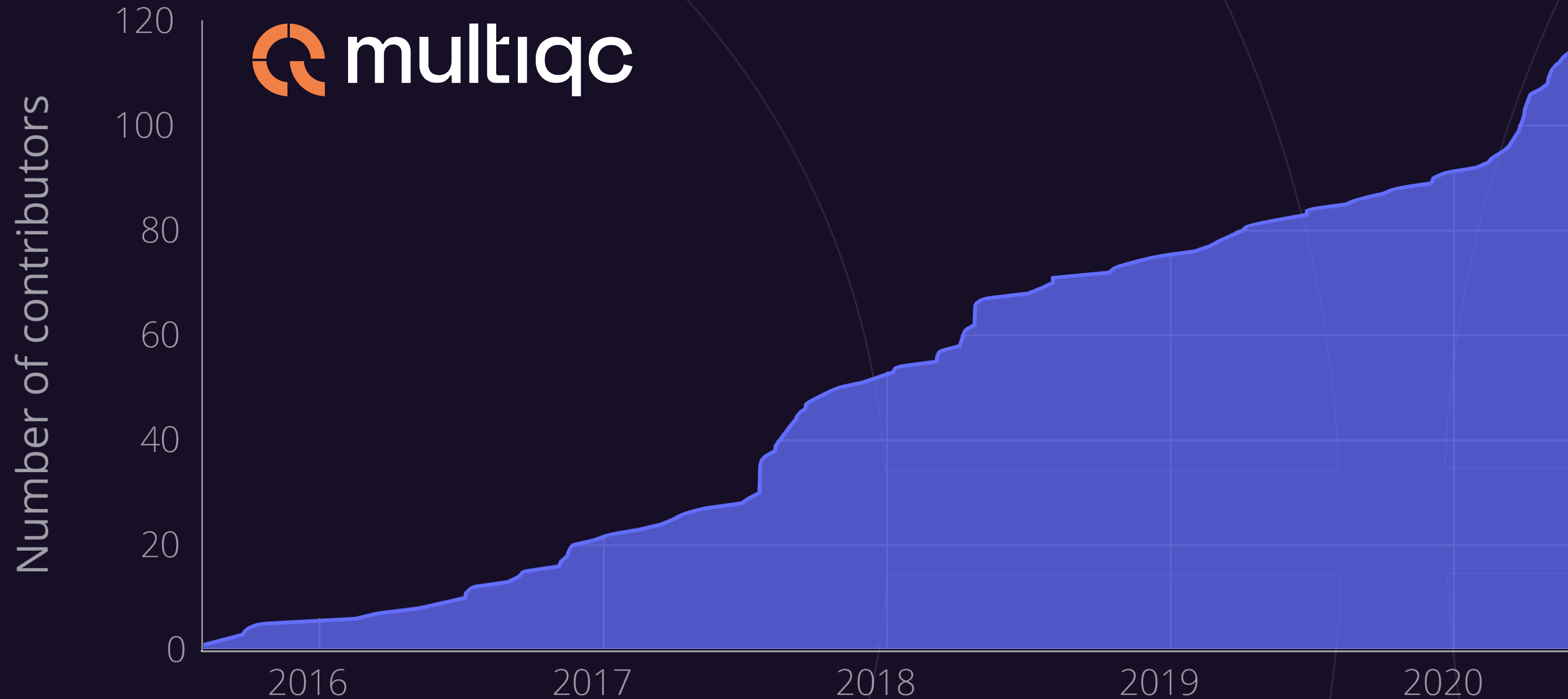
**BioCompute Objects**

seqera

# On the importance of being open

# On the importance of being open

# On the importance of being open

multiqc

**0.45** ⊙ Open issue
*per day*

**0.23** ⑂ Open pull request
*per day*

# On the importance of being open

Reproducible analysis of
genomics data at scale

Founded the nf-core community -
by removing institutional branding

Now 9000 community members on
Slack and still growing fast

SciLifeLab

NATIONAL
GENOMICS
INFRASTRUCTURE

nextflow

nf-core

# On the importance of being open



**nf-core**

**Plus ~1000 pipelines**
built using the nf-core template on GitHub

Legend: Released, In development

# On the importance of being open

🤗 Hugging Face

SciLifeLab
Serve

AlphaFold

# On the importance of being open

## nf-core

Domain experts are more important than ever

New for nf-core:
Special interest groups

Inter-disciplinary collaboration
can yield incredible results

**Futurism**

SELF KRYPTONITE  |  7.12.23, 3:56 PM EDT *by* MAGGIE HARRISON DUPRÉ

## AI Loses Its Mind After Being Trained on AI-Generated Data

"As the use of generative models continues to grow rapidly, this situation will only accelerate."

/ Artificial Intelligence  / Ai  / Ai Chatbots  / Ai Training

Image by Getty Images

seqera

# Future challenges

# Future challenges

## Accessible to anyone

Intuitive and easy to use for people from any background.

## Suitable for any use case

Generalist interface that can be used with any data type.

## Powerful at any scale

Can be run on your laptop or scaled to a production cluster with millions of samples.

# (choose two?)

# Future challenges

Accessible to anyone

Suitable for any use case

Powerful at any scale

# Future challenges

## Generally useful

Useful for the majority of people running
with this type of data / analysis.

## Specific

Analysis that can be applied to a
specific research question.

## Maintainable

Clean code base without excessive logic
or parameter space.

# The "final mile" of analysis

# Future challenges

Generally useful

nf-core

Specific

Modularity of components

Chaining of workflows

Importing and extending workflows

Maintainable

# Future challenges

Generally useful

Specific

Maintainable

**nf-core**

**multiqc**

"Custom content"

Use as library / use within notebooks

# Future challenges

Generally useful

Specific

Maintainable


nf-core


multiqc


seqera

Data Studios - eg. Notebooks, but any container:
interactive environments for downstream analysis

# Future challenges

## Lossy storage

What can we afford to throw away?

## Green computing

Justifying the cost of your data centre

## Heterogenous data

Mixing and matching data types for new science

## How much is enough?

# Future challenges

Loïc Lannelongue, Sabrina Krakau
green-algorithms.org
nextflow-io / nf-co2footprint

Lossy storage

Green computing

Heterogenous data



## Nextflow CO$_2$e footprint report

[special_davinci] *(resumed run)*

Workflow execution completed successfully!

**Run times**
12-Oct-2023 08:50:47 - 12-Oct-2023 08:55:44 (duration: **4m 57s**)

**Nextflow command**

```
/home-link/qeakr01/development/nextflow/launch.sh run nf-core/mag -r 2.3.0 -profile cfc_dev -c ../co2_stuff.cpus_8.config
co2footprint --input 's3://ngi-igenomes/test-data/mag/samplesheets/samplesheet.full.csv' --skip_binning --centrifuge_db
kraken2_db false --skip_prokka --outdir results_cpus8 -resume
```

**Nextflow version**
version 23.07.0-edge, build 5870 (22-07-2023 15:44 UTC)

## Total CO$_2$e footprint measures

| 1.23 Kg | 3.63 KWh | 1.34 | 7.02 |
|---------|----------|------|------|
| ☁ | 🔌 | 🌲 | 🚗 |
| CO$_2$e emissions | Energy consumption | Tree months | km by car | Flight

# Conclusion

# Conclusion

Join at slido.com

## #3201 014

# Conclusion

## Work in the open
Find your next collaborator online!

## Join initiatives
Don't assume you need to DIY

## Build for the future
A tool is for life, not just for Christmas

### Phil Ewels, PhD

Product Manager for Open Source

phil.ewels@seqera.io

seqera