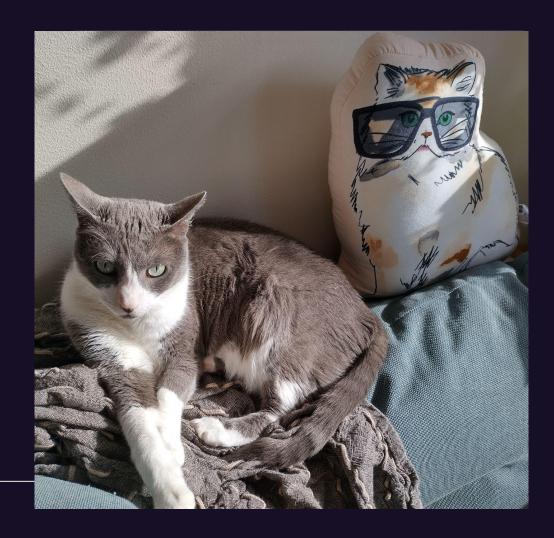
## References data

Maxime U Garcia SciLifeLab | National Genomics Infrastructure nf-core

ANF - workflow et reproductibilité



### The Challenge of Reproducibility

Hidden reproducibility issues are like an iceberg

"

First, we tried to re-run the analysis with the code and data provided by the authors.

Second, we reimplemented the whole method in a Python package...

Experimenting with reproducibility: a case study of robustness in bioinformatics <a href="https://doi.org/10.1093/gigascience/giy077">https://doi.org/10.1093/gigascience/giy077</a>

The five pillars of computational reproducibility: bioinformatics and beyond https://doi.org/10.1093/bib/bbad375





#### Vocabulary of Reproducibility

Permalink: http://nojhan.net/tfd/vocabulary-of-reproducibility.html 2025-03-31 © nojhan — CC-8Y-SA v4
Reference: http://arxiv.org/abs/2102.03380v2

### Random Factors e.g. seed, batch

Artifacts e.g. input data

Fixed Factors e.g. source code



Exactly repeat the original experiment, generating precisely the same results.









anomaly.









Test whether one can independently reach the same conclusion.







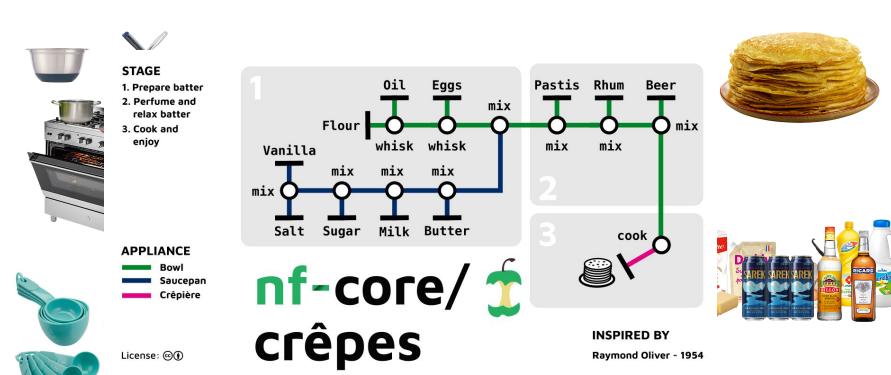








## So let's make crêpes, but France 1954 style



# Input data for pipelines

- Data to be analyzed
- Reference data
- Metadata



# nf-core T



### nf-core - what do we do?

All the resources we use

- Public data
- AWS iGenomes
- annotation-cache
- test-data
- Bundled with pipelines
- User provided



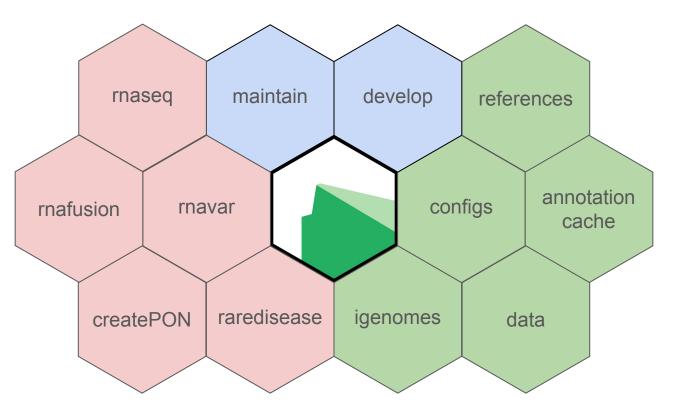
### Sarek

The park



Sarek is a National Park in Northern Sweden. Pretty inaccessible. With beautiful landscape. The name is an allegory for the genomic landscape to explore with the pipeline.

# Sarek and its ecosystem



### Annotation-cache

Many years ago, annotation only done using only local cache on the institutional cluster Difficulties for updating cache, I once waited for more than a year on my request BAD IDEA: Building containers with pre-downloaded cache

Gained admin access to a new cluster, so made full usage of local cache

Hassle to maintain, update, add data and not optimal on cloud

New option to download cache within Sarek

Discussion with AWS to make an open data resource (inspired by AWS iGenomes) Annotation Cache creation

