

# nf-core/sarek: an open-source pipeline for germline, tumor-only, and somatic analysis of NGS data

Friederike Hanssen<sup>1</sup>, Maxime Garcia<sup>2</sup>, Lasse Folkersen<sup>3</sup>, Susanne Jodoin<sup>1</sup>, Oskar Wacker<sup>1</sup>, Anders Sune Pedersen<sup>4</sup>, Edmund Miller<sup>5</sup>, Francesco Lescai<sup>6</sup>, Julius Joos<sup>7</sup>, nf-core community, Gisela Gabernet<sup>1</sup>, Sven Nahnsen<sup>1</sup>

<sup>1</sup>Quantitative Biology Center (QBiC), University of Tübingen, Tübingen <sup>2</sup>SciLifeLab, Karolinska Institutet, Stockholm <sup>3</sup>Nucleus Genomics Ltd., New York <sup>4</sup>Danish National Genome Center, Copenhagen <sup>5</sup>University of Texas, Dallas <sup>6</sup>Department of Biology and Biotechnology, University of Pavia, Pavia <sup>7</sup>Internal Medicine, University Hospital, Tübingen

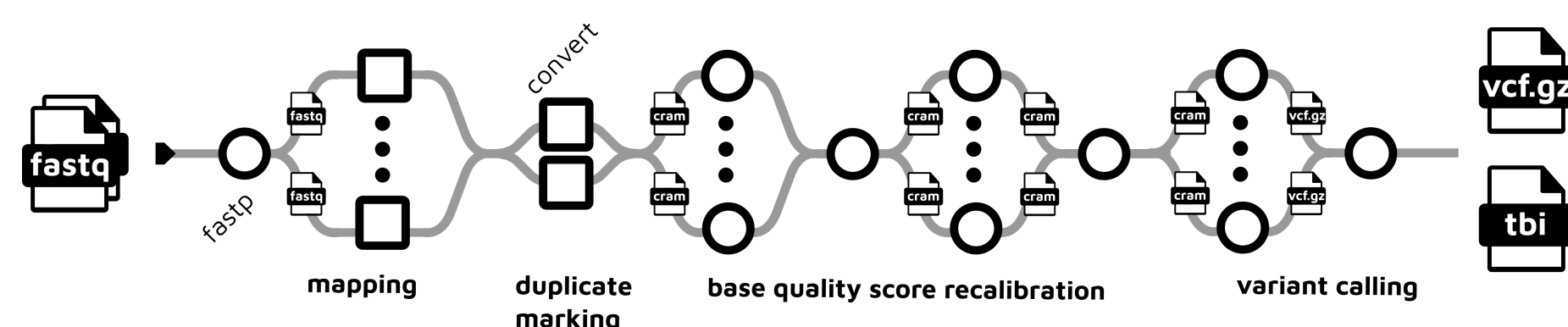
## 1. Introduction

Somatic variant calling studies often include many patients with dataset sizes varying widely between oncopanel, whole-exome, and whole-genome sequencing data. nf-core/sarek<sup>1</sup> is an established pipeline for exploring single-nucleotide variants, structural variation, microsatellite instability, and copy-number alterations of germline, tumor-only, and paired tumor-normal short-reads.

nf-core/sarek is part of nf-core<sup>2</sup>, a community project which provides an infrastructure to create reproducible, scalable, and portable open-source Nextflow<sup>3</sup>-based pipelines.

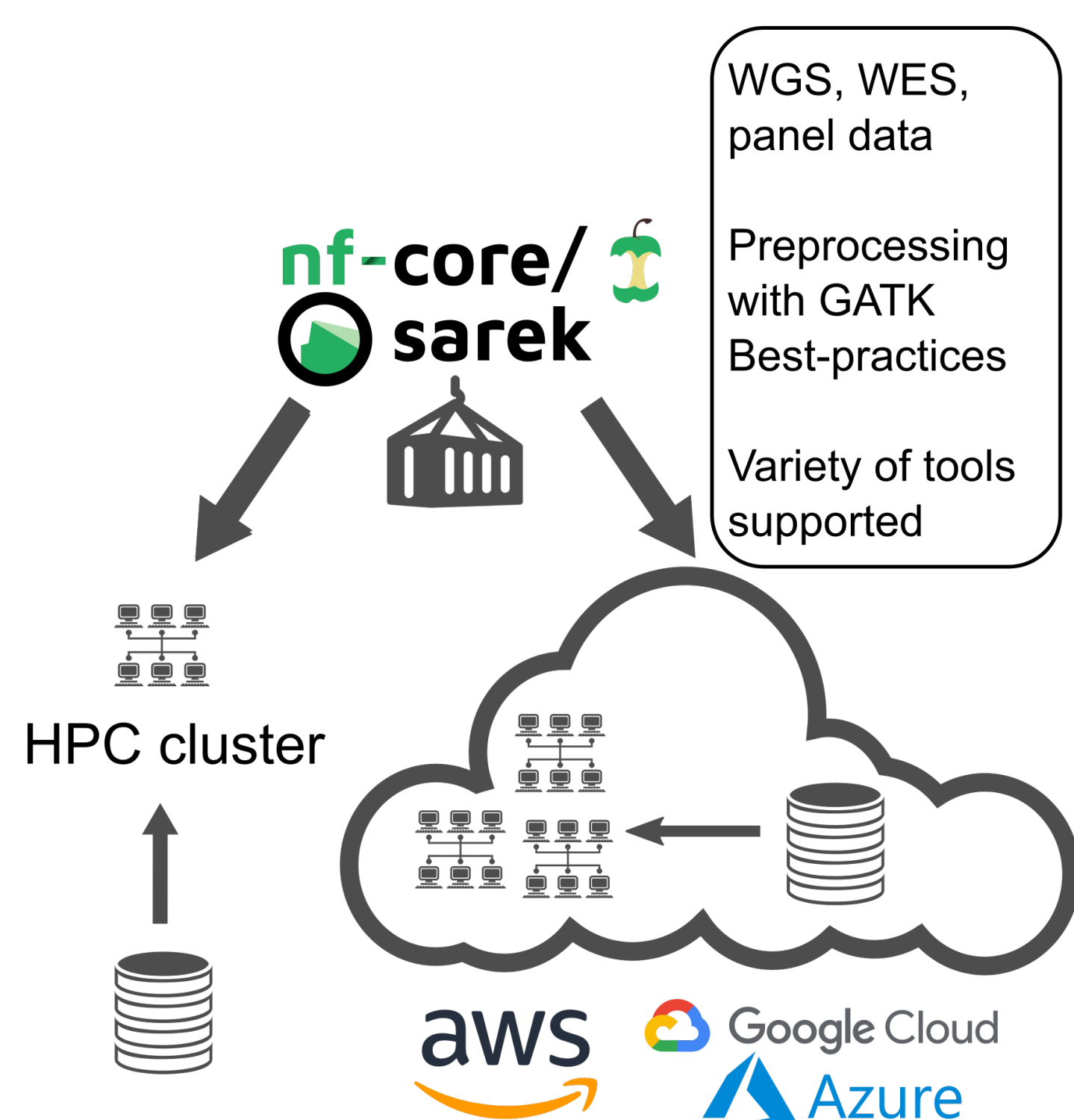
Here, we show the latest updates including improvements to the data flow and tool selection reducing time and compute resources, and modularization improving code maintainability.

## 2. Overview

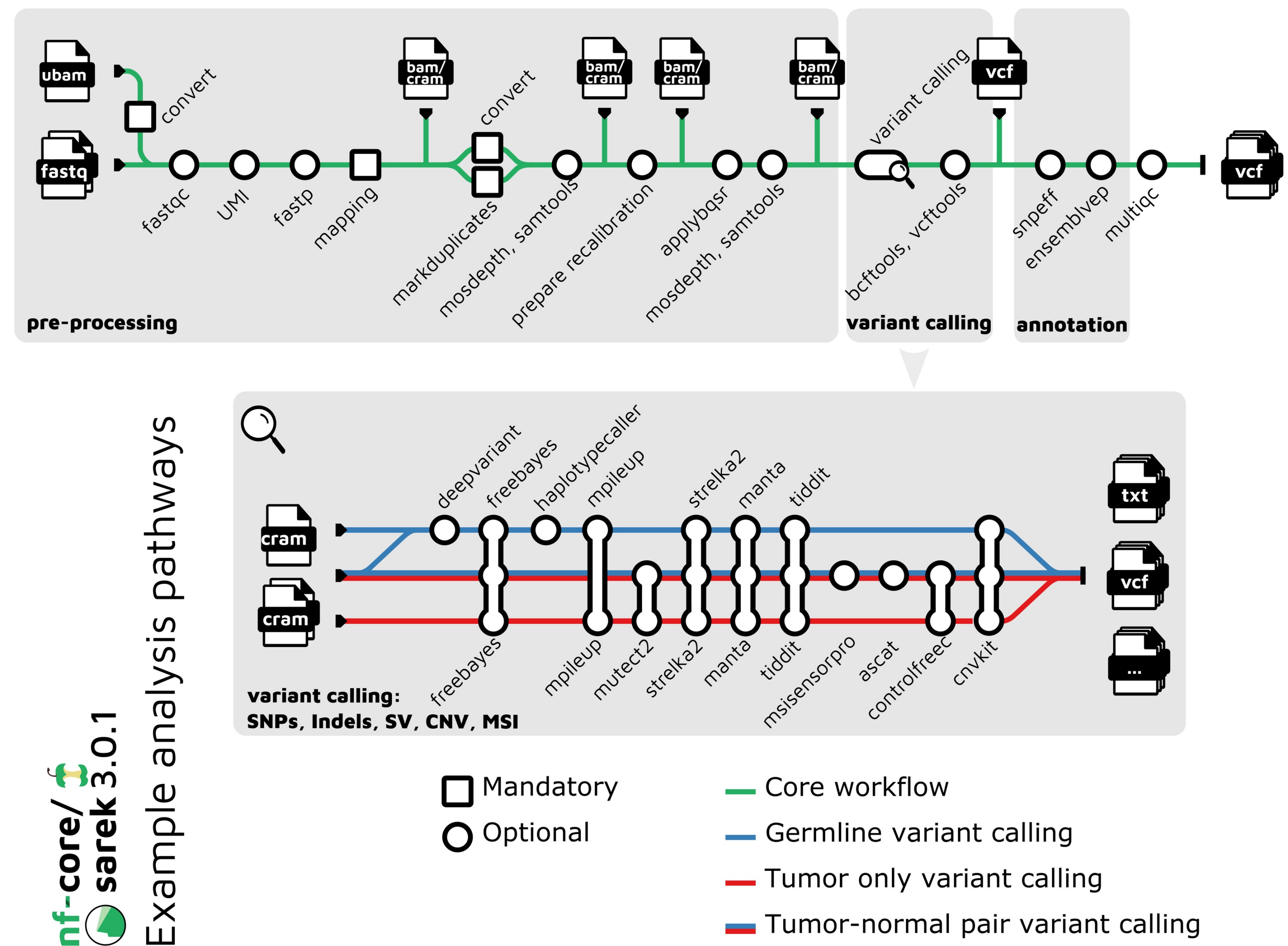


- FASTQ or BAM inputs are split into files of equal size before alignment to speed up computation
- Resulting BAM files are then merged and duplicate marked in one step before they are converted into CRAM format.
- Subsequent steps are run on multiple genomic regions in parallel. By default an interval file with chromosomes cut at their centromeres is used for WGS, and a user-supplied target bed file is used for WES or panel data.
- For all data types, small regions are grouped resulting in approximately equal sizes being processed together.

## 3. Nextflow pipelines are portable



nf-core pipelines are containerized and versioned to ensure reproducible analyses. They are portable to different compute infrastructures mitigating time-consuming data transfer.

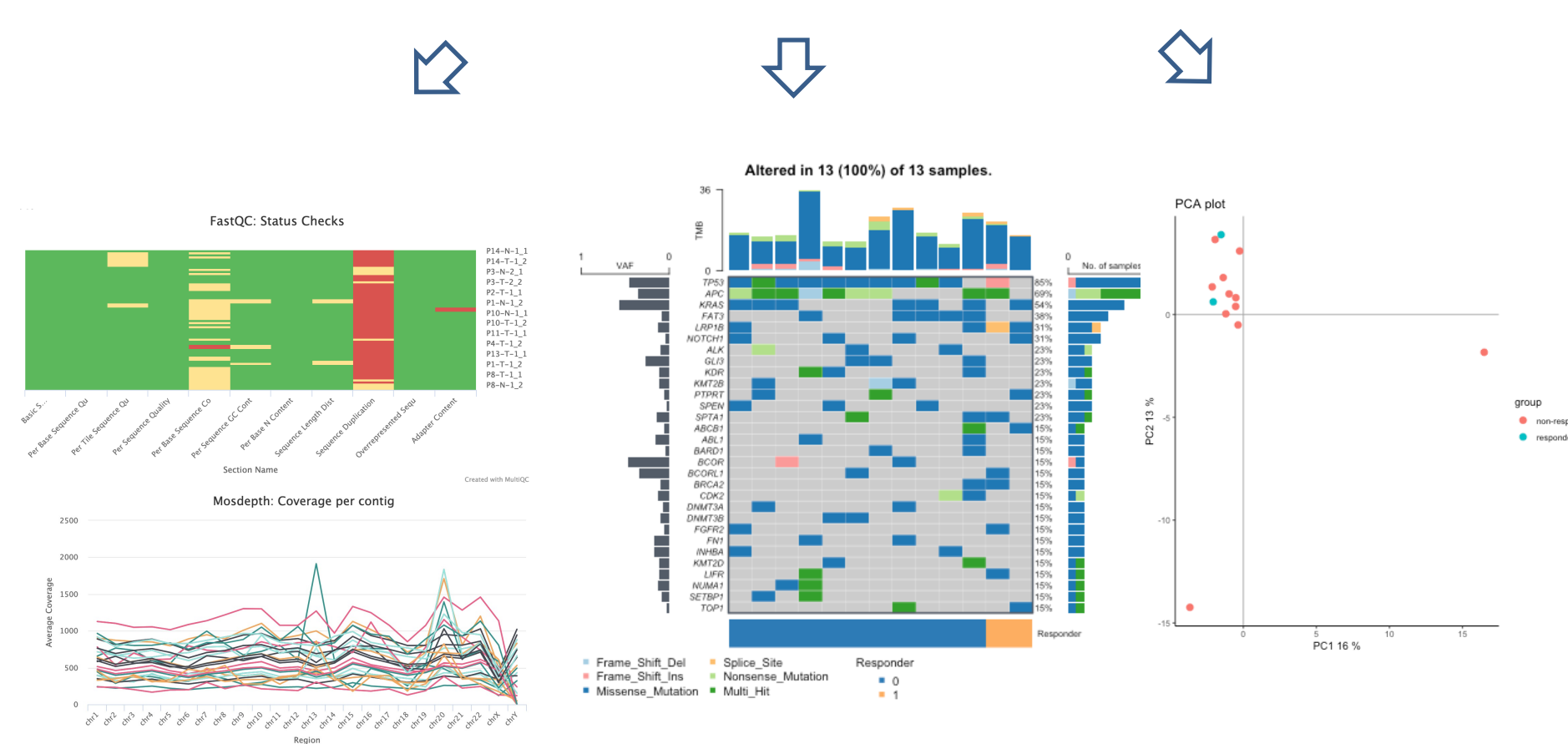


Adapted from: Fellows Yates, James A., et al. PeerJ 9 (2021).

Pipeline metromap showing a high-level view of the different analysis steps. The pipeline can be started from six different entry points and run through all subsequent tasks. All optional tools can be selected in any combination. This allows to recompute and extend the results throughout a project's duration.

## 4. Application

Showcase of somatic SNP/Indel calling results (Strelka2) on 13 onco-panels datasets comparing treatment responses



Extensive QC report using the tool MultiQC

Oncoplot generated with the R package maftools to compare SNP variants of both groups

PCA plot to evaluate clustering of the mutational signatures

Tools	SNP recall	SNP precision	F1
Haplotypcaller	0.990652	0.992179	0.9915
Deepvariant	<b>0.992866</b>	0.99787	<b>0.99536</b>
Freebayes	0.992847	0.954791	0.97345
Strelka2	0.983946	<b>0.998364</b>	0.9911

Tools	INDEL recall	INDEL precision	F1
Haplotypcaller	0.971052	0.985107	0.97803
Deepvariant	<b>0.983122</b>	0.988865	<b>0.98599</b>
Freebayes	0.945148	0.94993	0.94753
Strelka2	0.959560	<b>0.990090</b>	0.97458

Benchmark of germline variant calling on the Genome in a Bottle sample HG0002 (26X coverage, Illumina NovaSeq, concatenated fastq files) using the pipeline's default settings. Variants in the high-confidence regions were evaluated.

## Conclusion

- nf-core/sarek is a high-throughput reproducible pipeline ready to be used in high throughput variant calling projects.
- As a showcase project, 161 WGS germline samples were already analyzed with SNP, SV and CNV calling on a local HPC
- Cost and time evaluation on AWS cloud is currently under way.
- Continuous optimization & addition of community-requested tools
- Possible application: Reanalysis of ICGC /TCGA cohorts for comparative analyses with local cohorts

## Join us



<https://nf-co.re/sarek>

## References

- Garcia et al. (2020), F1000Research 9:63
- Ewels et al. (2020), Nature Biotechnology 38, 276–278
- Di Tommaso et al. (2017), Nature Biotechnology, 35(4), 316–319

## Acknowledgements

We would like to acknowledge funding from the Excellence cluster iFIT, and the SFB 209 & Amazon Web Services for cloud computing.

We are grateful to the nf-core and nextflow community for their support during the development.