🌐 http://sarek.scilifelab.se/

- Analysis germline and somatic workflow

🌐 http://sarek.scilifelab.se/

- Analysis germline and somatic workflow
- Whole genome or targeted sequencing

- Analysis germline and somatic workflow
- Whole genome or targeted sequencing
- Developed with NGI and NBIS

🌐 http://sarek.scilifelab.se/

- Analysis germline and somatic workflow
- Whole genome or targeted sequencing
- Developed with NGI and NBIS
- Support from The Swedish Childhood Tumor Biobank

🌐 https://www.nextflow.io/



🌐 https://www.sylabs.io/singularity/

# nextflow

🌐 https://www.nextflow.io/

Data-driven workflow language



🌐 https://www.sylabs.io/singularity/

HPC specific container engine

# AWS iGenomes

🌐 https://ewels.github.io/AWS-iGenomes/

- Human `GRCh37`
- Human `GRCh38`

# AWS iGenomes

🌐 https://ewels.github.io/AWS-iGenomes/

- Human `GRCh37`
- Human `GRCh38`
- Dog `CanFam3.1` 🔧
- Mouse `GRCm38` 🔧

Based on GATK Best Practices (GATK 4.0)

🌐 https://software.broadinstitute.org/gatk/best-practices/

Based on GATK Best Practices (GATK 4.0)

- Reads mapped to reference genome with `bwa mem`

🌐 https://software.broadinstitute.org/gatk/best-practices/

Based on GATK Best Practices (GATK 4.0)

- Reads mapped to reference genome with `bwa mem`
  - FASTQs or BAMs 🔧

🌐 https://software.broadinstitute.org/gatk/best-practices/

Based on GATK Best Practices (GATK 4.0)

- Reads mapped to reference genome with `bwa mem`
  - FASTQs or BAMs 🔧
- Duplicates marked with `picard MarkDuplicates`

🌐 https://software.broadinstitute.org/gatk/best-practices/

Based on GATK Best Practices (GATK 4.0)

- Reads mapped to reference genome with `bwa mem`
  - FASTQs or BAMs 🔧
- Duplicates marked with `picard MarkDuplicates`
- Recalibrate with `GATK BaseRecalibrator`

- SNVs and small indels:

- SNVs and small indels:
    - Freebayes
    - HaplotypeCaller
    - MuTect2
    - Strelka2

- SNVs and small indels:
  - Freebayes 🔴
  - HaplotypeCaller 🔵
  - MuTect2 🔴
  - Strelka2 🔵/🔴
- Structural variants:

## Variant Calling

- SNVs and small indels:
  - Freebayes 🔴
  - HaplotypeCaller 🔵
  - MuTect2 🔴
  - Strelka2 🔵/🔴
- Structural variants:
  - Manta 🔵/🔴
- Sample heterogeneity, ploidy and CNVs:

## Variant Calling

- SNVs and small indels:
  - Freebayes 🔴
  - HaplotypeCaller 🔵
  - MuTect2 🔴
  - Strelka2 🔵/🔴
- Structural variants:
  - Manta 🔵/🔴
- Sample heterogeneity, ploidy and CNVs:
  - ASCAT 🔴
  - Control-FREEC 🔧 🔴

## Annotation

- VEP and SnpEff
- 🗄 ClinVar, COSMIC, dbSNP, GENCODE, gnomAD, polyphen, sift, etc.

## Annotation

- VEP and SnpEff
- 🗄 ClinVar, COSMIC, dbSNP, GENCODE, gnomAD, polyphen, sift, etc.
- Possibility to use cache directories 🔧

- VEP and SnpEff
- 📚 ClinVar, COSMIC, dbSNP, GENCODE, gnomAD, polyphen, sift, etc.
- Possibility to use cache directories 🔧
- Prioritization 🔧
    - Rank scores are computed for all variants, and can be explored

🌐 https://www.sylabs.io/singularity/

🌐 https://www.sylabs.io/singularity/

- Available on `rackham` and/or `bianca`
- `/sw/data/uppnex/ToolBox/sarek`

🌐 https://www.sylabs.io/singularity/

- Available on `rackham` and/or `bianca`
- `/sw/data/uppnex/ToolBox/sarek`
  - Updated by myself at each new Sarek release

🌐 [https://www.sylabs.io/singularity/](https://www.sylabs.io/singularity/)

- Available on `rackham` and/or `bianca`
- `/sw/data/uppnex/ToolBox/sarek`
  - Updated by myself at each new Sarek release
- Next step `Sarek` module

🌐 https://bioconda.github.io/

🌐 https://bioconda.github.io/

- Execute Sarek within a conda environment

🌐 https://aws.amazon.com/

- Improving AWS usage

# Acknowledgments

**Barntumörbanken**
Elisa Basmaci
Szilveszter Juhos
Gustaf Ljungman
Monica Nistèr
Gabriela Prochazka
Johanna Sandgren
Teresita Díaz De Ståhl
Katarzyna Zielinska-Chomej

**Grupp Nistèr**
Saad Alqahtani
Min Guo
Daniel Hägerstrand
Anna Hedrén
Martin Proks
Rong Yu
Jian Zhao

**NGI**
Johannes Alneberg
Anandashankar Anil
Franziska Bonath
Orlando Contreras-López
Phil Ewels
Sofia Haglund
Max Käller
Anna Konrad
Pär Lundin
Remi-Andre Olsen
Senthilkumar Panneerselvam
Fanny Taborsak
Chuan Wang

**Clinical Genetics**
Jesper Eisfeldt

**NBIS**
Sebastian DiLorenzo
Malin Larsson
Marcel Martin
Markus Mayrhofer
Björn Nystedt
Markus Ringnér
Pall I Olason
Jonas Söderberg

**Clinical Genomics**
Kenny Billiau
Hassan Foroughi Asl
Valtteri Wirta

**Nextflow folks**
Paolo Di Tommaso
Sven Fillinger
Alexander Peltzer

## Any questions?

🌐 https://maxulysse.github.io/dnaclub2019

🐙 https://github.com/SciLifeLab/Sarek

💬 https://gitter.im/SciLifeLab/Sarek

🌐 http://sarek.scilifelab.se/

✳️ #sarek-pipeline