

Nextflow Data Processing & Configuration

⚠ Our pipeline configuration is still in-development, and the contents of this document are subject to change.

<https://nf-co.re/sarek/usage v2.7.1>

Data Type	Method	Output	Tried yet?
WES or WGS	DeepVariant	Germline SNV, INDEL	Yes
WES or WGS	Strelka, Mutect2, Freebayes ?	Somatic SNV, INDEL	
WES or WGS	TBD	Germline and Somatic Structural Variants	
WES or WGS	TBD	Germline and Somatic CNV	
WES or WGS	TBD	Tumor MSI	
SNV, INDEL variants	TBD	Annotated Variants	

<https://nf-co.re/maseq v3.5>

Data Type	Method	Output	Tried yet?
RNA-Seq	Salmon	Gene expression counts	Yes

bam/cram to fastq conversion

When fastq files are not available, cram/bam files are converted to fastq using this pipeline: [GitHub - qbic-pipelines/bamtofastq](#) (v1.2.0).

WES and WGS Variant Calling (SNV & INDEL)

Germline SNV + INDEL

This involves transformation of WES `fastq` or `cram` files to variant call files in VCF format (`.vcf` files).

As of Jan 2022, the `reference genome` used is `GRCh38`.

The processing steps include the following:

- Raw fastq files uploaded to Synapse by researcher in a folder with name format `experiment_name_rnaseq_fastq_date`. No white space should be present in the filenames (all filenames should have `_` for whitespaces).
- All experiment and sample related annotations need to be added on Synapse before processing can start. This is a required step so that a sample sheet can be generated to trigger the processing workflow
- The sample sheet should contain the following information in the following format (saved as a `.txt` file) :

sample	subject	status	sex	file_1	file_2	lane	parentID
Synapse specimenID	Synapse individualID	1 (Tumor = 1, Normal 0)	XX or XY	synID	synID	Lane information	SynapseID of parent folder

- The files are pulled into NextFlow workflow setup and processed using the following versions of software:

```
1 nf-core/sarek v2.7.1
2 Nextflow v21.10.5
3 BWA 0.7.17
```

```
4 GATK v4.1.7.0
5 FreeBayes v1.3.2
6 samtools v1.9
7 Strelka v2.9.10
8 Manta v1.6.0
9 TIDDIT v2.7.1
10 AlleleCount v4.0.2
11 ASCAT v2.5.2
12 Control-FREEC vv11.6
13 msisensor v0.5
14 SnpEff v4.3t
15 VEP v99.2
16 MultiQC v1.8
17 FastQC v0.11.9
18 bcftools v1.9
19 CNVkit v0.9.6
20 htlib v1.9
21 QualiMap v2.2.2-dev
22 Trim Galore v0.6.4_dev
23 vcftools v0.1.16
24 R v4.0.2
```

Commands used for running JHU samples on DeepVariant:

All files and sample sheets are first staged in S3 buckets linked to NFTower. then the following command are used to launch the processing pipeline.

Params:

```
1 input: s3://jhu-biobank-nf-project-tower-bucket/jobs/02-sage-sarek-2.7.1-deepvariant/inputs/sample-sheet.tsv
2 outdir: s3://jhu-biobank-nf-project-tower-bucket/jobs/02-sage-sarek-2.7.1-deepvariant/outputs/
3 genome: GRCh38
4 igenomes_base: s3://sage-igenomes/igenomes
5 model_type: WES
6 tools: "deepvariant"
```

Config:

```
1 process {
2   errorStrategy = 'retry'
3   maxRetries = 3
4
5   withLabel:deepvariant {
6     container = "google/deepvariant:1.1.0"
7     cpus = 24
8   }
9 }
```

Pre-run Script:

```
1 export NXF_VER=21.10.5
```

Profiles:

```
1 aws_tower
```

Somatic SNV + INDEL

TBD

Annotated Variants

Currently, germline variant calls in VCF format are being processed manually using VEP and `vcf2maf`

RNA SEQUENCING DATA QUANTIFICATION

Processing RNA-seq files involve transformation of raw data (`fastq` files) to transcript counts (`quants.sf` files).

The quantification software of choice is `Salmon`.

As of Jan 2022, the `reference genome` used is `GRCh38`.

Processing involves the following steps:

- Raw fastq files uploaded to Synapse by researcher in a folder with name format `experiment_name_rnaseq_fastq_date`. No white space should be present in the filenames (all filenames should have `_` for whitespaces).
- All experiment and sample related annotations need to be added on Synapse before processing can start. This is a required step so that a sample sheet can be generated to trigger the processing workflow
- The sample sheet should contain the following information in the following format (saved as a `.csv` file):

sample	single_end	fastq_1	fastq_2	strandedness
Synapse specimenID	0 (1 if paired-end)	synID	synID	reverse or forward

- The files are pulled into NextFlow workflow setup and processed using the following versions of software:

```
1 BEDTOOLS_GENOMECOV:
2 bedtools: 2.30.0
3 CAT_FASTQ:
4 cat: 8.3
5 CUSTOM_DUMPSOFTWAREVERSIONS:
6 python: 3.9.5
7 yaml: 5.4.1
8 DESEQ2_QC_STAR_SALMON:
9 bioconductor-deseq2: 1.28.0
10 r-base: 4.0.3
11 DUPRADAR:
12 bioconductor-dupradar: 1.18.0
13 r-base: 4.0.2
14 FASTQC:
15 fastqc: 0.11.9
16 GET_CHROM_SIZES:
17 samtools: 1.1
18 GTF_GENE_FILTER:
19 python: 3.8.3
20 PICARD_MARKDUPLICATES:
21 picard: 2.25.7
22 PRESEQ_LCEXTRAP:
23 preseq: 3.1.1
```

```
24 QUALIMAP_RNASEQ:
25 qualimap: 2.2.2-dev
26 RSEM_PREPAREREFERENCE_TRANSCRIPTS:
27 rsem: 1.3.1
28 star: 2.7.6a
29 RSEQC_BAMSTAT:
30 rseqc: 3.0.1
31 RSEQC_INFERENCE:
32 rseqc: 3.0.1
33 RSEQC_INNERDISTANCE:
34 rseqc: 3.0.1
35 RSEQC_JUNCTIONANNOTATION:
36 rseqc: 3.0.1
37 RSEQC_JUNCTIONSATURATION:
38 rseqc: 3.0.1
39 RSEQC_READDISTRIBUTION:
40 rseqc: 3.0.1
41 RSEQC_READDUPLICATION:
42 rseqc: 3.0.1
43 SALMON_QUANT:
44 salmon: 1.5.2
45 SALMON_SE_GENE:
46 bioconductor-summarizedexperiment: 1.20.0
47 r-base: 4.0.3
48 SALMON_TX2GENE:
49 python: 3.8.3
50 SALMON_TXIMPORT:
51 bioconductor-tximeta: 1.8.0
52 r-base: 4.0.3
53 SAMPLESHEET_CHECK:
54 python: 3.8.3
55 SAMTOOLS_FLAGSTAT:
56 samtools: 1.13
57 SAMTOOLS_IDXSTATS:
58 samtools: 1.13
59 SAMTOOLS_INDEX:
60 samtools: 1.13
61 SAMTOOLS_SORT:
62 samtools: 1.13
63 SAMTOOLS_STATS:
64 samtools: 1.13
65 STAR_ALIGN:
66 star: 2.6.1d
67 STRINGTIE:
68 stringtie: 2.1.7
69 TRIMGALORE:
70 cutadapt: 3.4
71 trimgalore: 0.6.7
72 UCSC_BEDCLIP:
73 ucsc: 377
74 UCSC_BEDGRAPHTOBIGWIG:
75 ucsc: 377
76 Workflow:
77 Nextflow: 21.10.5
78 nf-core/rnaseq: '3.4'
```

Command used to process JHU Biobank samples:

Params:

```
1 input: s3://jhu-biobank-nf-project-tower-bucket/jobs/01-nfcore-rnaseq-3.4/inputs/sample-sheet.csv
2 outdir: s3://jhu-biobank-nf-project-tower-bucket/jobs/01-nfcore-rnaseq-3.4/outputs/
3 genome: GRCh38
4 igenomes_base: s3://sage-igenomes/igenomes
```

Config:

```
1 process {
2   errorStrategy = 'retry'
3   maxRetries = 3
4 }
```

Pre-run script:

```
1 export NXF_VER=21.10.5
```

Profile:

```
1 aws_tower
```