
Implementing Scalable Bioinformatic Workflows in Snakemake and Nextflow

Nathan S. Watson-Haigh
Radosław Suchecki

Licensing

This work is licensed under a Creative Commons Attribution 3.0 Unported License and the below text is a summary of the main terms of the full Legal Code (the full licence) available at <http://creativecommons.org/licenses/by/3.0/legalcode>.

You are free:

- to copy, distribute, display, and perform the work
- to make derivative works
- to make commercial use of the work

Under the following conditions:

Attribution - You must give the original author credit.

With the understanding that:

Waiver - Any of the above conditions can be waived if you get permission from the copyright holder.

Public Domain - Where the work or any of its elements is in the public domain under applicable law, that status is in no way affected by the license.

Other Rights - In no way are any of the following rights affected by the license:

- Your fair dealing or fair use rights, or other applicable copyright exceptions and limitations;
- The author's moral rights;
- Rights other persons may have either in the work itself or in how the work is used, such as publicity or privacy rights.

Notice - For any reuse or distribution, you must make clear to others the licence terms of this work.



Contents

Licensing	3
Contents	4
Workshop Information	5
The Trainers	6
Providing Feedback	7
Document Structure	7
Introduction to Snakemake	9
Key Learning Outcomes	10
Resources Required	10
Useful Links	10
Setting Up Your Environment	11
Your First Minimal Snakefile	13
Generalising Rules with Wildcards	14
Submitting Jobs to Slurm	15
A Bash “Pipeline”	16
Reimplementing A Workflow in Snakemake	17
Snakemake Troubleshooting	30
Introduction to Nextflow	33
Key Learning Outcomes	34
Resources Required	34
Useful Links	35
Introduction	36
Setting Up Your Environment	37
Nextflow basics	39
Example workflow	43
Modify/extend the workflow	52
Your own workflow (TODO: replace with variant calling?)	53
Troubleshooting	54
Space for Personal Notes or Feedback	55

Workshop Information

The Trainers

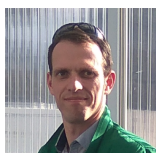


Dr. Nathan S. Watson-Haigh

Senior Bioinformatician

Bioinformatics Hub, University of Adelaide

nathan.watson-haigh@adelaide.edu.au



Dr. Radosław Suhecki

Research Scientist

Crop Bioinformatics and Data Science, CSIRO

rad.suhecki@csiro.au

Providing Feedback

While we endeavour to deliver a workshop with quality content and documentation in a venue conducive to an exciting, well run hands-on workshop with a bunch of knowledgeable and likable trainers, we know there are things we could do better.

Whilst we want to know what didn't quite hit the mark for you, what would be most helpful and least depressing, would be for you to provide ways to improve the workshop. i.e. constructive feedback. After all, if we knew something wasn't going to work, we wouldn't have done it or put it into the workshop in the first place!

Clearly, we also want to know what we did well! This gives us that "feel good" factor which will see us through those long days and nights in the lead up to such hands-on workshops!

With that in mind, we'll provide a some high tech mechanism through which you can provide anonymous feedback during the workshop:

1. Some empty ruled pages at the back of this handout. Use them for your own personal notes or for writing specific comments/feedback about the workshop as it progresses.

Document Structure

We have provided you with an electronic copy of the workshop's hands-on tutorial documents. We have done this for two reasons: 1) you will have something to take away with you at the end of the workshop, and 2) you can save time (mis)typing commands on the command line by using copy-and-paste.

We advise you to use Acrobat Reader to view the PDF. This is because it properly supports some features we have implemented to ensure that copy-and-paste of commands works as expected. This includes the appropriate copy-and-paste of special characters like tilde and hyphens as well as skipping line numbers for easy copy-and-paste of whole code blocks.



While you could fly through the hands-on sessions doing copy-and-paste, you will learn more if you use the time saved from not having to type all those commands, to understand what each command is doing!

The commands to enter at a terminal look something like this:

```
1 tophat --solexa-quals -g 2 --library-type fr-unstranded -j \  
  annotation/Danio_rerio.Zv9.66.spliceSites -o tophat/ZV9_2cells \  
  genome/ZV9 data/2cells_1.fastq data/2cells_2.fastq
```

The following styled code is not to be entered at a terminal, it is simply to show you the syntax of the command. You must use your own judgement to substitute in the correct arguments, options, filenames etc

```
tophat [options]* <index_base> <reads_1> <reads_2>
```

The following is an example of how R commands are styled:

```
1 R --no-save
2 library(plotrix)
3 data <- read.table("run_25/stats.txt", header=TRUE)
4 weighted.hist(data$short1_cov+data$short2_cov, data$lgth, breaks=0:70)
5 q()
```

The following icons are used in the margin, throughout the documentation to help you navigate around the document more easily:



Important



For reference



Follow these steps



Questions to answer



Warning - STOP and read



Bonus exercise for fast learners



Advanced exercise for super-fast learners

Introduction to Snakemake

Primary Author(s):
Nathan S. Watson-Haigh nathan.watson-haigh@adelaide.edu.au

Contributor(s):

Key Learning Outcomes

After completing this module the trainee should be able to:

- Install Snakemake in a conda environment
- Execute a Snakemake workflow
- Use the provided “profile” to execute jobs on a compute cluster
- Write simple Snakemake rules capable of generating some output(s) by executing some code which operates on some input(s)

Resources Required

For the purpose of this training you need access to:

- A compute cluster with the `module` command available to you for loading software
- Singularity (<https://sylabs.io/singularity/>) - available as a module on the above cluster
- Conda(<https://www.anaconda.com/distribution/>) - available as a module on the above cluster

Tools Used

Snakemake

<https://snakemake.readthedocs.io>

Graphviz

<https://www.graphviz.org>

Useful Links

Slurm Documentation

<https://slurm.schedmd.com/documentation.html>

Setting Up Your Environment

For the purpose of the workshop we will be working on the head node of an HPC cluster running slurm (<https://slurm.schedmd.com/documentation.html>). This is the most likely infrastructure that fellow bioinformaticians already find themselves using on a regular basis. We also assume that the cluster provides the `module` command for you to load software and the modules **Anaconda3** and **Singularity** are available to use.

The execution of the Snakemake workflow will actually take place on the cluster head node with jobs being submitted to Slurm for queuing and processing. From the head node, Snakemake will monitor the submitted jobs for their completion status and submit new jobs as dependent jobs complete successfully.

Connect to the Cluster Head Node



First up, let's connect to the head node of the HPC cluster using `ssh`.

See your local facilitator for connection details. You will have one user account per person.

Monitoring Slurm Jobs



You can monitor all jobs in the slurm queue, or just your own job(s) using the slurm command `squeue`:

```
1 # All jobs in the queue
2 squeue
3
4 # Just your own jobs
5 squeue --user ${USER}
```

For convenience we have provided you with the `sq` function which produces nicer output than the default `squeue` and only shows your own jobs:

```
1 # Your own jobs
2 sq
3
4 # Someone else's jobs
5 sq --user ${SOMEONE_ELSE}
```

Install Snakemake

The recommended installation route for Snakemake is through a conda environment (https://snakemake.readthedocs.io/en/stable/getting_started/installation.html). As such, you need Anaconda3, usually available to you on your cluster via the module system.



```

1 # We use a specific version for reproducibility reasons
2 # Find the latest version: https://anaconda.org/search?q=snakemake
3 SNAKEMAKE_VERSION="5.5.4"
4
5 # Load miniconda
6 module load \
7     miniconda3-4.6.14-gcc-5.4.0-kkzv7zk
8
9 #####
10 # One-time commands
11 #####
12 # Integrate conda into bash
13 conda init bash
14 . ~/.bashrc
15
16 # Change the default location into which conda saves packages
17 # and environments
18 conda config --prepend pkgs_dirs /shared/${USER}/.conda/pkgs
19 conda config --prepend envs_dirs /shared/${USER}/.conda/envs
20
21 # Change the default channels used for finding software and
22 # resolving dependencies
23 conda config --add channels defaults
24 conda config --add channels bioconda
25 conda config --add channels conda-forge
26 #####

```



Do NOT run the following command! This is provided for future reference so you know how to Install Snakemake on another system. Rather than creating the conda environment from scratch, we'll simply copy a pre-existing directory so we save time, and possible headaches.

```

1 # Install snakemake using conda
2 # This might take 5-10mins
3 conda create \
4     --name snakemake \
5     --yes \
6     snakemake=${SNAKEMAKE_VERSION:-5.5.4}

```

Snakemake installation is now complete.



For the purposes of this workshop, simply copy the following `.conda` directory and you will have Snakemake setup and ready to go:

```
1 mkdir --parents /shared/${USER}
2 cp --recursive \
3   /shared/ubuntu/.conda \
4   /shared/${USER}/
```

All that is left to do is to activate the environment which will make `snakemake` available on the command line:

```
1 # Activate the newly created conda environment
2 conda activate snakemake
```

Integrate Snakemake autocompletion into bash:

```
1 complete -o bashdefault -C snakemake-bash-completion snakemake
```

Test if Snakemake is actually working:

```
1 snakemake --version
```

If you experience problems with the installation, head to the [Troubleshooting](#) section for help.



While waiting for others to catch up, why not have a look into how you would go about updating Snakemake within this conda environment if there is a new version available.

Your First Minimal Snakefile

To get started with Snakemake, all you need to do is create a **Snakefile** (note the capitalisation) containing a rule which specifies how to create an output file.

Setup a working directory for this task:

```
1 mkdir --parents /shared/${USER}/snakemake/minimal
2 cd /shared/${USER}/snakemake/minimal
```

Create a file called **Snakefile** and add the following content:

```
1 rule hello_world:
2     output:
3         "Hello/World.txt",
4     shell:
5         """
6         echo "Hello, World!" > {output}
```



```
7 """
```

You can now run this workflow in one of 3 ways:

```
1 # Request Snakemake to generate the specific output file
2 # "Hello/World.txt"
3 snakemake Hello/World.txt
4
5 # Request Snakemake to execute the specific rule "hello_world"
6 snakemake hello_world
7
8 # Request Snakemake to execute the first rule in the Snakefile
9 snakemake
```



What happens if you run one of the above commands two or more times? Why?

Generalising Rules with Wildcards

The original `hello_world` rule wasn't very flexible. We couldn't say "Hello, World!" in Spanish, Polish or French. However, we can generalise the rule using "wildcards":

```
1 rule hello_world:
2     output:
3         "{cheer}/{world}.txt",
4     shell:
5         """
6         echo "{wildcards.cheer}, {wildcards.world}!" > {output}
7         """
```

Now we can use whatever language we want:

```
1 # In English
2 snakemake Hello/World.txt
```

```
3
4 # In Polish
5 snakemake Czesc/Swiat.txt
6
7 # In French and Spanish at the same time
8 snakemake Monde/Monde.txt Ciao/Mondo.txt
```

Take a look at the files created:

```
1 tree ./
```

Submitting Jobs to Slurm

By default, Snakemake executes jobs on the same computer on which it is running. For Snakemake to be able to submit jobs to a cluster resource management/queuing system, such as Slurm, we can use a “profile” which conveniently contains scripts for job submission and monitoring as well as setting some additional Snakemake command line arguments so it can “talk” to a cluster backend.

To avoid having to delve into implementing our own “profile” for use with our Slurm cluster, we have created a Slurm profile ready for you to use. So lets grab it:

```
1 # Ensure a working directory exists and move into it
2 mkdir --parents /shared/${USER}/snakemake/tutorial
3 cd /shared/${USER}/snakemake/tutorial
4
5 # Clone the Snakemake template repository from GitHub
6 git clone https://github.com/UofABioinformaticsHub/snakemake_template ./
7
8 # Checkout the "simple" branch
9 git checkout simple
```

The Snakefile in this branch of the repository is the same “Hello, World!” example you created above, with wildcards. Lets see how we use the provided “profile” to get Snakemake to submit jobs to Slurm:

```
1 snakemake \
2   --profile profiles/slurm \
3   Hello/World.txt Czesc/Swiat.txt Monde/Monde.txt Ciao/Mondo.txt
```

If the `STDOUT` and `STDERR` of the command(s) in a rule are not explicitly sent to a file, then they will end up in Slurm’s log file for a particular job which is normally something like `slurm-<job_id>.out`. This isn’t that helpful for debugging purposes, so the provided profile changes this to `logs/<rule_name>/<wildcards>.out` e.g. `logs/hello_world/cheer=Ciao,world=Mondo.out`.

Cleanup after yourself!

```
1 | snakemake \  
2 |   --delete-all-output \  
3 |   Hello/World.txt Czesc/Swiat.txt Monde/Monde.txt Ciao/Mondo.txt
```

A Bash “Pipeline”

A bioinformatic “pipeline” is commonly a single, monolithic bash script which performs all the tasks which need to be performed. For example, someone might have written a script for performing the following tasks:

- Run FastQC across all the raw read files
- Adapter, quality, and read length filtering using Trimmomatic
- Aggregating FastQC reports from the raw reads using MultiQC
- Index the reference FASTA file
- Perform a `bwa-mem` read alignment

We have this script available for you on the `tutorial` branch, switch to it and have a look:

```
1 | git checkout tutorial  
2 |  
3 | less analysis.sh
```

While the author of such a script should be commended for their efforts in documenting their analysis using a script, it has several significant limitations:

Not parallelised

loops over input files, executing independant commands in sequential order

Resources over-specified

the compute resources needed by the script are dictated by the command(s) with the largest requirement(s)

Not idempotent

significant programming logic is needed to wrap around commands to detect failures and only execute parts of the analysis which failed in earlier attempts



How might *you* modify the above script to:

- Add new samples
- Rerun the script if you find one of the files generated is corrupt
- Include readgroup information at the `bwa-mem` step (`-b argumanet`)

How would you avoid rerunning commands which take a long time and already completed sucessfully on a previous run e.g. the reference index, `bwa-mem` etc?

Reimplementing A Workflow in Snakemake

We will walk you through the steps of reimplementing the first few steps of the above script into a Snakemake workflow. Along the way, we will introduce the core concepts of Snakemake and then ask you to reimplement the `bwa-mem` step yourselves. For those working quickly, you will have the opportunity to reimplement the `multiqc` step. This will provide you with a foundation for you to be able to convert your own workflows into Snakemake rules and begin reaping the rewards of being able to run your analyses in Snakemake.

Getting the Data

We've provided you with some real data whole genome sequencing (WGS) data from wheat together with a small chunk of the wheat genome. The data set is small enough so each step in the analysis will take less than a couple of minutes to run. We have a copy of this data available locally to save on bandwidth, time and the possibility we are detected as a DDoS attack on some poor remote server!

```
1 # Get a copy of the data
2 cp --recursive \
3   /shared/data/{raw_reads,references,misc} \
4   ./
5
6 # Have a look at what files we'd provided
7 tree raw_reads references misc
```

Implementing BWA Indexing

We have provided an out-of-the-box **Snakefile** capable of indexing the provided reference sequence, together with comments. Lets have a bit of a play before we get around to actually running the workflow.

```
1 less Snakefile
2
3 # All these commands have the same effect
4 snakemake --dryrun bwa_index
5 snakemake --dryrun all
6 snakemake --dryrun
```

The effect `--dryrun` is to simply show you what “would” be run, without actually running it. It's useful to ensure you're going to get what you though, especially as your workflows get larger and more interconnected.

Another useful feature is to generate a directed acyclic graph (DAG) of the jobs which comprise the workflow and how they are linked together. Although for this workflow is not yet that impressive, but we'll have a look at how we generate the DAG:

```
1 snakemake \
2   --dag \
3   | dot -Tpdf \
4   > dag1.pdf
```

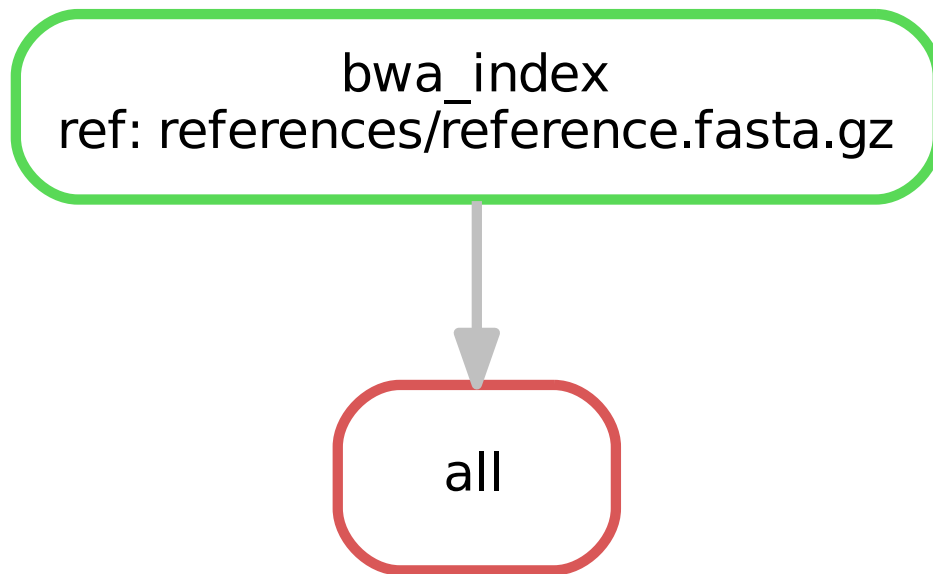


Figure 1: DAG of jobs showing `bwa_index` job dependant on the `all` pseudo-rule.

Implementing FastQC

Lets add a rule for performing FastQC on our input files. Looking at the `analysis.sh` file we see the following command is executed for each `SAMPLE` while iterating over the `SAMPLES` list:

```
1 fastqc --threads 1 \  
2   raw_reads/${SAMPLE}_R1.fastq.gz \  
3   raw_reads/${SAMPLE}_R2.fastq.gz
```

This command can be converted into a Snakemake rule by adding the following rule to the Snakefile:

```
1 rule fastqc:  
2   input:  
3     r1 = "raw_reads/{SAMPLE}_R1.fastq.gz",  
4     r2 = "raw_reads/{SAMPLE}_R2.fastq.gz",  
5   output:  
6     zip = [ "raw_reads/{SAMPLE}_R1_fastqc.zip", \  
7            "raw_reads/{SAMPLE}_R2_fastqc.zip" ],  
8     html = [ "raw_reads/{SAMPLE}_R1_fastqc.html", \  
9              "raw_reads/{SAMPLE}_R2_fastqc.html" ],  
10  shell:  
11    """  
    fastqc --threads 1 {input.r1} {input.r2}  
    """
```

Now we can run Snakemake and request a “target” file which matches an output files defined by the above `fastq` rule:

```
1 | snakemake --dryrun raw_reads/ACBarrie_R1_fastqc.html \  
   | raw_reads/ACBarrie_R2_fastqc.html
```

There are a few improvements we can make to this rule:

- We don't need to process both the R1 and R2 read files with the same FastQC job. We can operate on one read file at a time. By doing this, Snakemake will be able to execute the FastQC job for each file in parallel.
- We want a convenient way of generating FastQC outputs for ALL samples without typing them all at the command line.

Improving FastQC Parallelisation

Add the following new rule to your Snakefile:

```
1 | rule fastqc_single_input:  
2 |     input:  
3 |         "raw_reads/{prefix}.fastq.gz",  
4 |     output:  
5 |         zip = "raw_reads/{prefix}_fastqc.zip",  
6 |         html = "raw_reads/{prefix}_fastqc.html",  
7 |     shell:  
8 |         ""  
9 |         fastqc --threads 1 {input}  
10 |         ""
```

Now run the same Snakemake dryrun command as before:

```
1 | snakemake --dryrun raw_reads/ACBarrie_R1_fastqc.html \  
   | raw_reads/ACBarrie_R2_fastqc.html
```



Why did Snakemake complain about an `AmbiguousRuleException`?

How could we fix it? Hint: <https://snakemake.readthedocs.io/en/stable/snakefiles/rules.html#handling-ambiguous-rules>

Go ahead and delete the `fastqc` rule in favour of the `fastqc_single_input` rule. Do the same dryrun and see how many jobs Snakemake would run in order to create those files:

```
1 | snakemake --dryrun raw_reads/ACBarrie_R1_fastqc.html \  
   | raw_reads/ACBarrie_R2_fastqc.html
```

Pseudo-Rules

We can use “pseudo-rules” to define a list of target filenames for creation when we use the rule name as a “target”. Pseudo-rules consist of just an `input` directive:

```
1 | rule all:  
2 |     input:  
3 |         "raw_reads/ACBarrie_R1_fastqc.html",  
4 |         "raw_reads/ACBarrie_R2_fastqc.html",
```

By convention, the first pseudo-rule in the `Snakefile` is called `all` and specifies all the output filenames of the workflow. This now means we can execute a workflow in any of the following ways:

```
1 | # Not specifying a target will result in Snakemake executing the  
2 | # first rule in the Snakefile ("all" in this case)  
3 | snakemake --dryrun  
4 |  
5 | # Explicitly request the "all" rule  
6 | snakemake --dryrun all
```

When workflows get larger and the lists of filenames get bigger, specifying long lists of filenames in pseudo-rules can start to feel cumbersome. Since Snakemake syntax is an extension of Python, we can start to use some Python data structures and functions to help.

Add the following Python list of sample names (with most commented out for now) at the top of the file:

```
1 SAMPLES = [  
2     "ACBarrie",  
3     "Alsen",  
4     # "Baxter",  
5     # "Chara",  
6     # "Drysedale",  
7     # "Excalibur",  
8     # "Gladius",  
9     # "H45",  
10    # "Kukri",  
11    # "Pastor",  
12    # "RAC875",  
13    # "Volcanii",  
14    # "Westonia",  
15    # "Wyalkatchem",  
16    # "Xiaoyan",  
17    # "Yitpi",  
18    ]
```

Add FastQC output files for all samples in the **SAMPLES** list, as well as both read files, as new targets to the existing **all** rule. We'll make use of the **expand()** function to simplify things somewhat. The resulting **all** pseudo-rule should look like this:

```
1 rule all:  
2     input:  
3         expand("references/reference.fasta.gz.{ext}", ext=['amb', 'ann', \  
4             'bwt', 'pac', 'sa']),  
5         expand("raw_reads/{SAMPLE}_{read}_fastqc.html", SAMPLE=SAMPLES, \  
6             read=['R1', 'R2']),
```

Lets take a look at what jobs would be run if we run the whole workflow. Remember, the following commands are equivalent:

```
1 # Explicitly run the "all" pseudo-rule  
2 snakemake --dryrun all  
3  
4 # Run the first rule in the Snakefile. This should be the "all" rule by \  
5     convention  
6 snakemake --dryrun
```

Lets look at the DAG for the workflow:

```

1 | snakemake \
2 |   --dag \
3 |   | dot -Tpdf \
4 |   > dag2.pdf

```

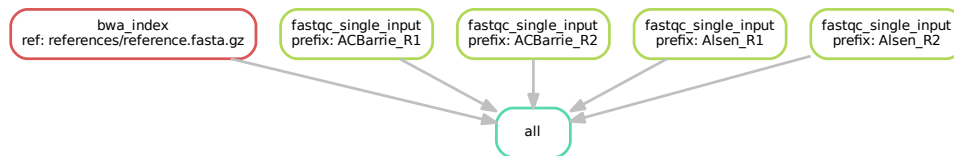


Figure 2: DAG of jobs showing `bwa_index` and several `fastqc_single_input` job dependant on the `all` pseudo-rule.

Executing the Workflow on Slurm

Up until now, we've just been playing around with `--dryrun`, so lets move on and start executing the workflow on the Slurm cluster!

Remember, we need to use the Slurm profile we've provided you with so Snakemake knows how to communicate with Slurm. In addition, we're also going to execute the jobs within a singularity container which has the tools we need already installed inside it.

```

1 | # Make sure Singularity is available
2 | module load \
3 |   singularity-3.2.1-gcc-5.4.0-tn5ndnb
4 |
5 | # Execute the workflow
6 | snakemake \
7 |   --profile profiles/slurm \
8 |   --use-singularity

```

Depending on how quickly everyone else is in executing their workflows, you might get to see your jobs in the Slurm queue by executing this in another window:

```

1 | sq

```

Implementing Trimmomatic

We've gone through implementing the FastQC command as a Snakemake rule and demonstrated the core concepts of Snakemake along the way. We'll go through implementing one more command as a Snakemake rule before you go off and try one on your own!

If you compare the `trimmomatic` command in `analysis.sh` to the rule provided below, you will see that we have simply pulled out all references to input or output files into the `input` or `output` directives. Where we had used the bash variable `${SAMPLE}` in the filenames we are now using Snakemake "wildcards". It is almost the same syntax - just notice the absence of the `$` but the curly braces are retained. The biggest changes seen are in the `shell` directive, where we now have to refer to the input and output files via `{input.r1}`, `{output.r1_unpaired}` etc.

```

1 rule trimmomatic:
2     input:
3         r1 = "raw_reads/{SAMPLE}_R1.fastq.gz",
4         r2 = "raw_reads/{SAMPLE}_R2.fastq.gz",
5         adapters = "misc/trimmomatic_adapters/TruSeq3-PE.fa"
6     output:
7         r1 = "qc_reads/{SAMPLE}_R1.fastq.gz",
8         r2 = "qc_reads/{SAMPLE}_R2.fastq.gz",
9         r1_unpaired = "qc_reads/{SAMPLE}_R1.unpaired.fastq.gz",
10        r2_unpaired = "qc_reads/{SAMPLE}_R2.unpaired.fastq.gz",
11    shell:
12        """
13        trimmomatic PE \
14            -threads 1 \
15            {input.r1} {input.r2} \
16            {output.r1} {output.r1_unpaired} \
17            {output.r2} {output.r2_unpaired} \
18            ILLUMINACLIP:{input.adapters}:2:30:10:3:true \
19            LEADING:2 \
20            TRAILING:2 \
21            SLIDINGWINDOW:4:15 \
22            MINLEN:36
23        """

```

Next, we need to add the trimmomatic output files to our `all` pseudo-rule to make it convenient to create them. Your `all` pseudo-rule should look like this:

```

1 rule all:
2     input:
3         expand("references/reference.fasta.gz.{ext}", ext=['amb', 'ann', \
4             'bwt', 'pac', 'sa']),
5         expand("raw_reads/{SAMPLE}_{read}_fastqc.html", SAMPLE=SAMPLES, \
6             read=['R1', 'R2']),
7         expand("qc_reads/{SAMPLE}_{read}.fastq.gz", SAMPLE=SAMPLES, \
8             read=['R1', 'R2']),

```


Implementing BWA-MEM

Now is your opportunity to put into practice what you have learnt from the above walk-thoughts of implementing FastQC and Trimmomatic commands. Your task is to implement the `bwa mem` command into a Snakemake rule.

Here are some questions to get you thinking as you try to implement this rule:



What input read files are required for the command/rule?

Does the command/rule need the FASTA reference file or the index files as input?

The BWA-MEM command uses a “prefix” to the FASTA index files, not the index filenames themselves. How will you specify this in the `shell` directive?

If you hard-coded it, what would the rule look like?



Hard-coding the path is simple, but not ideal. What if you changed the name of the reference file or wanted to use the rule with a different project? You would have to modify the paths in multiple places, once in the `input` directive and once in the `shell` directive.

Moving the hard-coded path out of the `shell` directive into the `params` directive (see: <https://snakemake.readthedocs.io/en/stable/snakefiles/rules.html#non-file-parameters-for-rules>).



Now you are making use of the `params` directive, this opens up the possibility of using some Python to do some string manipulations on the paths defined in the `input` directive. In particular, we can use a Python lambda function in the `params` directive (see: <https://snakemake.readthedocs.io/en/stable/snakefiles/rules.html#non-file-parameters-for-rules>).

How might you change a hard-coded path in the `params` directive to use a lambda function which manipulates the index file path(s) set in the `input` directive to define the prefix? Hint: take the path of one of the index files and remove the last few characters corresponding to the last file extension. You will probably need to do some reading of the Snakemake and/or Python documentation.

Now add the BAM files corresponding to the all the samples in the `SAMPLES` list, that you want the BWA-MEM rule to produce, to the `all` pseudo-rule.

Don't worry if you didn't complete the above implementation of the BWA-MEM command, we have a git repository branch with the rules we developed. Get it, execute the workflow and generate the DAG:

```
1 # Checkout the branch with our implemetation of these rules
2 git checkout final
3
4 # Execute the workflow
5 snakemake \
6     --profile profiles/slurm \
7     --use-singularity
8
9 # Generate a DAG
10 snakemake \
```

```

11 --dag \
12 | dot -Tpdf \
13 > dag3.pdf

```

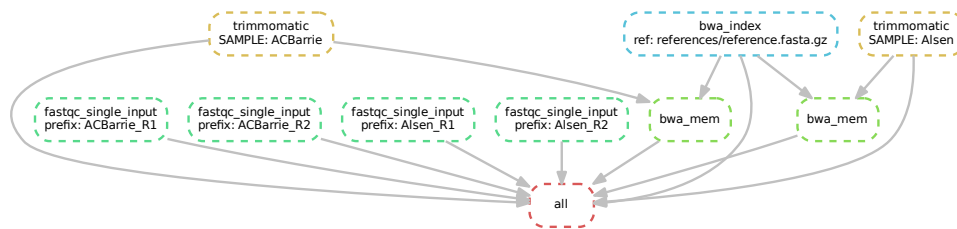


Figure 3: DAG of jobs showing the dependencies which exist in our final implementation of the `analysis.sh` workflow.

Adding New Samples

Our `SAMPLES` list contains a lot of samples which are currently commented out. Lets uncomment them and have a look at some other features of Snakemake:

```

1 # Manually uncomment the samples or use this sed command
2 sed -i 's/^# "/" "/" Snakefile

```

With so many more samples, the DAG becomes next to useless:

```

1 # Generate a DAG
2 snakemake \
3 --dag \
4 | dot -Tpdf \
5 > dag4.pdf

```



Figure 4: DAG of jobs for the whole workflow consisting of 16 samples.

Instead, the “rulegraph” might provide a better view of the workflow. Unlike the DAG, that shows the individual jobs and their dependencies, the rulegraph shows only the rules and their dependencies so provides a simplified view of the workflow:

```

1 # Generate a rulegraph
2 snakemake \
3 --rulegraph \
4 | dot -Tpdf \
5 > rulegraph.pdf

```

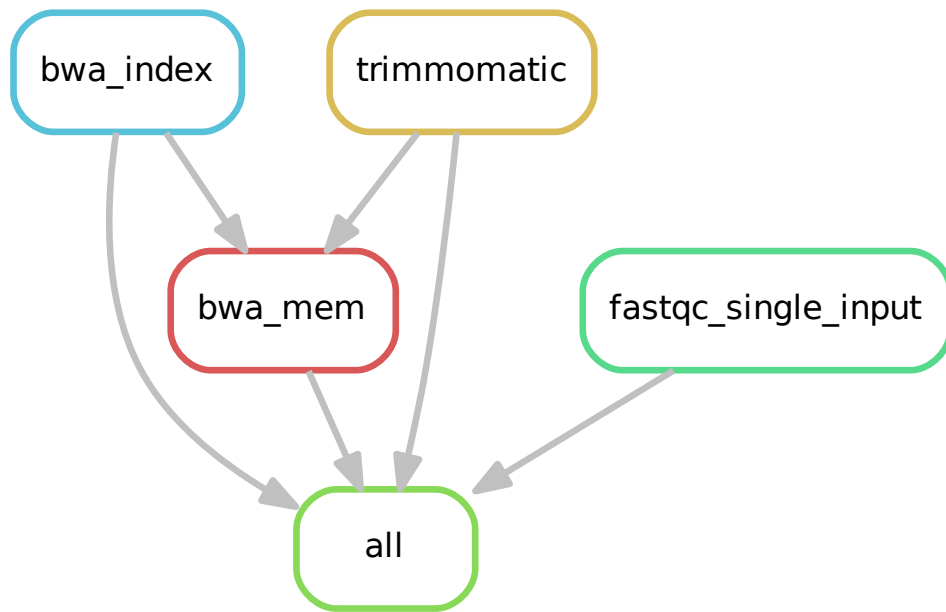


Figure 5: Rulegraph for the whole workflow.

Execute the rest of the workflow:

```
1 # Execute the workflow
2 snakemake \
3   --profile profiles/slurm \
4   --use-singularity
```



How many `fastqc_single_input` and total number of jobs are run as part of the whole workflow? Hint: try using `--forceall` in combination with `--dryrun`.

Using the Snakemake help, which command line argument can be used to get Snakemake to print the shell commands associated with each job during a `dryrun`?

Using the Snakemake help, which command line argument can be used to delete all the outputs associated with a given “target”?

Snakemake Troubleshooting

Snakemake Install

If you have a broken or incomplete snakemake installation, try the following steps to fix things:

```
1 # deactivate the snakemake conda environment if it is already active
2 conda deactivate
3
4 # Delete the snakemake conda environment
5 conda env remove --name snakemake
```

Now try reinstalling snakemake.

Conda Software Environment Setup

If your job failed or timed out, you will need to re-run conda software environment setup job again. However, you may first need to release the Snakemake lock which protects you from running multiple instances of the same workflow at the same time:

```
1 snakemake \  
2 --unlock
```

To ensure Snakemake starts with a clean slate, delete the “hidden” `.snakemake` directory:

```
1 rm -rf .snakemake
```

Getting Going After a Disconnect

If you find that your connection to the server has been dropped, you can get yourself going again using this convenient block of commands:

```
1 # Load the required software modules  
2 module load \  
3 miniconda3-4.6.14-gcc-5.4.0-kkzv7zk \  
4 singularity-3.2.1-gcc-5.4.0-tn5ndnb  
5  
6 # Activate the snakemake conda environment and integrate shell \  
7   autocompletion into bash  
8 conda activate snakemake  
9 complete -o bashdefault -C snakemake-bash-completion snakemake  
10  
11 # Move to the correct directory location  
12 cd /shared/${USER}/snakemake-tutorial
```

Introduction to Nextflow

Primary Author(s):

Radosław Suchecki rad.suchecki@csiro.au

Contributor(s):

Nathan S. Watson-Haigh nathan.watson-haigh@adelaide.edu.au

Key Learning Outcomes

After completing this module the trainee should be able to:

- Install Nextflow and execute an existing Nextflow workflow locally
- Modify the workflow to allow its execution on a compute cluster
- Write simple Nextflow process definitions and connect them with channels
- Apply operators to transform items emitted by a channel
- Leverage Nextflow's implicit parallelisation to process multiple data chunks independently

Resources Required

For the purpose of this training you need access to:

- A compute cluster with the `module` command available to you for loading software
- <https://sylabs.io/singularity/> Singularity - available as a module on the above cluster
- <https://www.anaconda.com/distribution/> conda - available as a module on the above cluster

Tools Used

Nextflow

<https://nextflow.io>

Graphviz

<https://www.graphviz.org>

Useful Links

Nextflow Documentation

<https://www.nextflow.io/docs/latest/index.html>

Nextflow Patterns

<http://nextflow-io.github.io/patterns/>

Slurm Documentation

<https://slurm.schedmd.com/documentation.html>

Introduction

Setting Up Your Environment

For the purpose of the workshop we will be working on the head node of an HPC cluster running [Slurm](#). This is the most likely infrastructure that fellow bioinformaticians already find themselves using on a regular basis. We also assume that the cluster provides the `module` command for you to load software and the modules Java and Singularity are available to use.

The execution of the Nextflow workflow will take place on the cluster head node with jobs being submitted to Slurm for queuing and processing. From the head node, Nextflow will monitor the submitted jobs for their completion status and submit new jobs as dependent jobs complete successfully.

Connect to the Cluster Head Node



First up, lets connect to the head node of the HPC cluster using `ssh`.

See your local facilitator for connection details. You should have one user account per person.

Install nextflow



```
1 # Load the Java module on your cluster
2 # If it's unavailable contact the cluster sysadmin
3 module load openjdk-1.8.0_202-b08-gcc-5.4.0-sypwasp
4
5 # Download and install nextflow executable
6 curl -s https://get.nextflow.io | bash
7
8 # You should now be able to run it
9 ./nextflow help
```

The installation should have placed the executable in your working directory. It is preferable to move the executable to a directory accessible via `$PATH`, to be able to run `nextflow` rather than having to remember to type the full `/path/to/nextflow` each time you want to run it.

Depending on the system this may suffice:



```
1 mkdir -p $HOME/bin
2 mv ./nextflow $HOME/bin
```

You should now be able to run `nextflow` without specifying the location of the binary. Let's see if it works by running a script which is nextflow's take on 'hello world'.

Hello (nextflow) world!



```
1 nextflow run rsuhecki/hello
```

Nextflow will pull the `rsuhecki/hello` GitHub repository and run its main script.



We are relying on nextflow's integration with git and git registries. The **alternative** would be to

```
git clone https://github.com/rsuhecki/hello.git
nextflow run hello/main.nf
```

In which case the location of the cloned repository will be different to the one used by nextflow. You will also not have access to nextflow-git integration functionality.



Where do we find the local copy of `hello`? Hint: try `nextflow` commands related to pipeline sharing, such as `list` and `info`.

For now, we are mostly interested in the local path to the repository, the file name of the main script and its contents, which we will discuss next.



While waiting for others to catch up, why not have a look into how you would go about pulling and removing local clones of remote repositories using nextflow.



What revisions (git branches or tags) are available for `nextflow-io/hello`? How would you run a specific revision?

Nextflow basics

Processes and channels

- *process* – a wrapper for a language-agnostic script which ensures isolation of the executed code.
- *channel* – an asynchronous¹ FIFO queue which facilitates data flow to/from/between processes by linking their outputs/inputs.

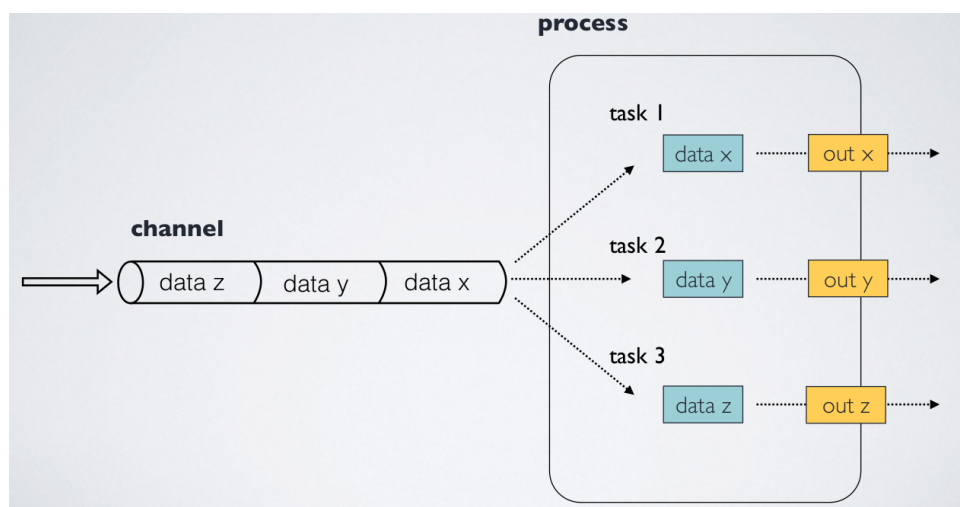


Figure 6: Nextflow building blocks: a *channel* “feeding” a processes. A *task* is an instance of a process. An isolated task is created for each emission (data chunk) from the input channel. Credit: Evan Floden

The main script

A nextflow script file name can be anything but in most cases it is best to stick to the default `main.nf`. The main script for the ‘hello’ example is as follows:

```

1  #!/usr/bin/env nextflow
2  echo true
3
4  cheers = Channel.from 'Bonjour', 'Ciao', 'Hello', 'Hola'
5
6  //setting default value, to be modified at runtime
7  params.world = 'world'
8
9  process sayHello {

```

¹send operation completes immediately, receiving stops the receiving process until the message has arrived

```

10  input:
11      val x from cheers
12  script:
13      """
14      echo '$x $params.world!'
15      """
16  }

```

A channel called `cheers` is created and emits each of the listed strings separately. A separate instance of the process `sayHello` is executed for each emission.



The content of the above script can be broken down as follows:

- The shebang line (line 1) is optional.
- Setting `echo true` will output `stdout` of (every) process to the terminal - not advised for real world applications.
- `Channel.from(some_list)` creates a channel emitting the list elements one by one.
- [Process](#) definition (lines 6-13)
 - Input block (lines 7-8)
 - Script block (lines 9-12)
- The `$x` in the script block is a nextflow variable local to the process, not a bash variable.
- Indentation is inconsequential.

In addition [process directives](#) could be inserted above the input block.

Hello HPC!

The nextflow hello example shown us how the `sayHello` process was executed separately for each input string as a separate *task*, but all the tasks were executed locally on our cluster's head node. We would now like each task to be submitted as a batch job for execution on one of the compute nodes.



```

1  nextflow run rsuckecki/hello -revision slurm

```

This is the modified version of the `main.nf` script. Submission to Slurm was achieved by adding `executor 'slurm'` directive to the process definition.

```

1  #!/usr/bin/env nextflow
2  echo true
3

```



```
4 | cheers = Channel.from 'Bonjour', 'Ciao', 'Hello', 'Hola'
5 |
6 | //setting default value, to be modified at runtime
7 | params.world = 'world'
8 |
9 | process sayHello {
10 |   executor 'slurm'
11 |
12 |   input:
13 |     val x from cheers
14 |   script:
15 |     """
16 |     echo "$x $params.world from \${HOSTNAME} on Slurm!"
17 |     """
18 | }
```

You might also have noticed that we have modified the script block so that the messages printed to the terminal include the name of the compute node on which a given task is executed.



Note the difference between how nextflow variables (`$x`, `$params.world`) and bash variables (`$HOSTNAME`) are included in the script block. There are alternative ways of including variables in scripts for execution by nextflow processes which may be more convenient if your script contains multiple special characters.

Hello task caching!

When the pipeline is launched with the `-resume` option, any attempt to execute already executed process with the same inputs, will cause the process execution to be skipped, producing the stored data as the output.

In this toy example we do not specify any outputs but the ‘hello’ messages printed to the terminal reflect this behaviour.



```
1 | nextflow run rsucecki/hello -revision slurm -resume
```

To avoid unintentionally re-computing long running tasks you may consider always running your pipelines with `-resume` and only omitting it on rare occasions when you want to re-compute the results even though inputs have not changed.

<https://www.nextflow.io/docs/latest/process.html#cache>

Hello command line options

Single-dashed options are reserved for nextflow engine (`-resume`, `-revision`, `-ansi-log false` etc). The double-dashed options are all yours and you are free to use them for your workflow. When you `nextflow run some_script.nf --foo bar`, the value of the parameter ('bar') will be accessible in `main.nf` as `params.foo` and within a script block as `$params.foo`.



In the 'hello' example we use `params.world` which by default is set to 'word', so lets try to use an alternative string.

```
1 | nextflow run rsuckecki/hello -revision slurm --world Mundo
```

Goodbye Hello

Nextflow facilitates but does not enforce separation of workflow logic from the configuration of compute and software environments as well as from other properties of the workflow. As such, you *could* get by developing nextflow workflows without worrying about that aspect – but you would be missing a lot in terms of flexibility, extensibility, portability and more

Nextflow looks for workflow configuration primarily in `nextflow.config` file, and additional config files can be included. Unsurprisingly the 'hello' example does not require much configuration, we would also like to crunch some real, albeit small, data.



This is mostly symbolic

```
1 | nextflow drop rsuckecki/hello
```

Let's have a play with a slightly more practical workflow.

Example workflow

We are going to work with an example Nextflow workflow to demonstrate how they are run, improve your understanding of *processes* and *channels* and finally introduce *operators*, which are applied to channels to shape and direct flowing data.

This example workflow consists of the following steps:

- Running FastQC across the raw reads
- Aggregating the raw read FastQC reports using MultiQC
- Performing adapter, quality, and read length filtering using Trimmomatic
- Running FastQC across the QC'd reads
- Aggregating the QC read FastQC reports using MultiQC
- Indexing the reference FASTA file
- Performing a `bwa-mem` read alignment



Although not necessary for simply running the pipeline, in the training context it makes sense to start by cloning the workflow repository and moving to the directory.

```
1 mkdir -p /shared/${USER}/nextflow-tutorial
2 cd /shared/${USER}/nextflow-tutorial
3 git clone \
    https://github.com/csiro-crop-informatics/nextflow-embl-abr-webinar.git \
    example_workflow
4 cd example_workflow
5 git checkout noslurm
6 git branch
```

This time, in addition to `main.nf` we have a separate script which downloads the required data sets, which include a small reference FASTA file and 16 pairs of FASTQ files, each for a different bread wheat accession.

```
1 nextflow run setup_data.nf
```

If successful, we could now try to run the workflow...

```
1 nextflow run main.nf
```



This is expected to fail.

Unless all the software required by the pipeline is available on the `$PATH`, which we don't expect, the pipeline should terminate with an error. The output information may help you identify the cause. Try to relate the error message to the relevant section of the main script (`main.nf`).



Which process has failed? What was the underlying cause?

There are two main issues with executing this workflow as is,

1. Third-party software tools have not been made available to the workflow.
2. We are trying to run the entire workflow on the cluster's head node.

There are different ways in which these issues could be addressed, for example using process *directives* at the top of each process definition. Depending on your cluster configuration this could be for example:

```
1 process foo {  
2   executor 'slurm'  
3   module 'samtools/1.9'  
4   //further code omitted
```

This is a perfectly valid syntax, which can be convenient, particularly during pipeline development, but for more portable workflows it is preferable to keep compute and software environment configuration separate from pipeline logic – in simple terms not in the workflow script (`main.nf`).

The config file(s) and profiles

Workflow configuration belongs in `nextflow.config` file. Transferring the above mention *directives* from process definitions in `main.nf` to `nextflow.config` would make things slightly better, e.g.

```
1 #nextflow.config
2
3 process.executor = 'slurm'
4 process.module = 'samtools/1.9'
```

or using the preferred syntax

```
1 process {
2     executor = slurm
3     module = 'samtools/1.9'
4 }
```

This is however still a bit rigid.

- You may be developing your pipeline on a local machine or a server where software modules are not available.
- If developing directly in the cluster environment, you may prefer your quick test runs to happen either on the head node or in an interactive session you are using, rather than always having jobs submitted to sit in the always-busy cluster queue.

Nextflow enables the definition of *profiles* which make it easy to run a workflow with different configuration settings, including, but not limited to executors and software environment.

For our pipeline we have defined several *profiles*, which allow us to execute the logic in `main.nf` while providing the required software either by creating a `conda` environment or by using Docker or Singularity containers where the `conda` environment has already been captured.

Relevant profiles

Identify the profile definitions in `nextflow config`. The ones most immediately relevant are:

```
1 profiles {
2     //SOFTWARE
3     conda {
4         process {
5             conda = "$baseDir/conf/conda.yaml"
6         }
7     }
8     singularity {
```



```
9     process {
10         container = \
            'shub://csiro-crop-informatics/nextflow-embl-abr-webinar'
11     }
12     singularity {
13         enabled = true
14         autoMounts = true
15         cacheDir = "singularity-images"
16     }
17 }
18 }
```

As you can see, Nextflow makes it really easy to define software environment via Singularity or Conda².

Given that Singularity is available on our cluster, let's start by using that profile, as the most robust way of setting up the software environment.

We will need Singularity for nextflow to be able to pull the container image from Singularity Hub and run the containerised software. By default the pipeline will process reads for a single accession – our head node should be able to handle this.



```
1 # Load the Singularity module
2 # If it is unavailable contact the cluster sysadmin
3
4 module load singularity-3.2.1-gcc-5.4.0-tn5ndnb
5
6 # Run the workflow
7
8 nextflow run main.nf -profile singularity
```

This is sufficient when running a workflow locally, in an interactive session or on a standalone server. The next step is to get nextflow to make use of the HPC batch submission system, to be able to run the full workflow without unleashing your sysadmins wrath.



Edit `nextflow.config`. Your task is to add a `slurm` profile which will set the appropriate executor.

²We also have a docker profile which you may find useful if you decide to run the workflow on your machine



What is your `slurm` profile configuration and where do you place it in `nextflow.config`?

There are of course many settings that can and in some cases must be set – refer to [executors section of Nextflow documentation](#)³. For running real-life pipelines in a cluster environment you will also use [directives](#)⁴ controlling the resources (`cpus`, `memory`, `time`) requested for each job. Other possibly relevant directives include `queue` and `scratch`.

Cluster run

To avoid running the workflow on our head node or in an interactive session, we will use the `slurm` profile you have defined⁵. As before, the software environment will be handled via the `singularity` profile. For that, we will need Singularity on the head node for nextflow to be able to pull the container image from Singularity Hub (we could also use a locally stored image). Singularity will also be required on the compute nodes which will run the individual tasks, but this should happen seamlessly if an appropriate module is loaded on the head node, otherwise the required module would also have to be specified in the workflow configuration files.



By default a single accession will be processed. You may use the `-resume` flag to avoid re-computing already existing results.

```
1 # Load the Singularity module on your cluster
2 # If it is unavailable contact the cluster sysadmin
3
4 module load singularity-3.2.1-gcc-5.4.0-tn5ndnb
5
6 # Run the workflow
7
8 nextflow run main.nf -profile slurm,singularity -resume
```

³<https://www.nextflow.io/docs/latest/executor.html>

⁴<https://www.nextflow.io/docs/latest/process.html#directives>

⁵If you are struggling and can't get help, try: `git stash && git checkout workshop`

Under the hood

If you think you are ready to look under the hood and try to work out how nextflow stages process inputs, wraps process script blocks and submits them to the cluster, here is a start.



```
1 # Remove the work directory to limit the number of task directories \
  to look at
2 rm -r work
3 # Re-run for a single sample
4 nextflow run main.nf -profile slurm,singularity
5 # Take a peak
6 ls -la work/ | less
7 # or
8 tree -ah work/ | less
```

Each task is executed in a separate directory and every abbreviated hash displayed in the terminal can be related to a specific sub-directory of `./work`, such as `work/d2/c4517b0a81f61ceca29ec355ddeaa6/` in which you may find

```
1 # NF generated files
2 .command.begin
3 .command.err
4 .command.log
5 .command.out
6 .command.run
7 .command.sh
8 .command.trace
9 .exitcode
10
11 # Output file
12 H45.bam
13
14 # Symlinks to input files
15 H45_R1.paired.fastq.gz
16 H45_R2.paired.fastq.gz
17 reference.fasta.gz.amb
18 reference.fasta.gz.ann
19 reference.fasta.gz.bwt
20 reference.fasta.gz.pac
21 reference.fasta.gz.sa
```

Identify and investigate hidden file (starting with dot) containing the executed script and the one containing cluster and container handling.

Cluster run - all accessions

We have successfully submitted workflow to the cluster.

To be sure, feel free to re-run it again (and again, and again...) with `-resume` to avoid wasting CPU cycles.



```
1 | nextflow run main.nf -profile slurm,singularity -resume
```

If all went well, the workflow successfully processed a single accession, let's have a closer look at the script to better understand how it handles the inputs before we proceed to run it on all the accessions.



In `main.nf` we create a channel which reads pairs of FASTQ files from a sub-directory of the `./data`. We then apply some operators.

```
Channel.fromFilePairs("data/${region}/*_R{1,2}.fastq.gz")
    .take ( params.take == "all" ? -1 : params.take )
    .into { readPairsChannelA; readPairsChannelB }
```

1. Identify the two operators, refer to [nextflow documentation^a](https://www.nextflow.io/docs/latest/operator.html) as required and explain the purpose of each of the two operators.

2. How can you run the workflow for more than one accession? How about all of them? Recall that workflow parameters use double-dash syntax. Run the relevant commands.

^a<https://www.nextflow.io/docs/latest/operator.html>

Monitoring your jobs on our cluster



You can monitor your job(s) in the slurm queue using the slurm command `squeue`:

```
1 | squeue --user ${USER}
```

For convenience you are also provided with the `sq` function which produces nicer output and by default only shows your own jobs:

```
1 | sq
2
```

```

3 # Someone elses jobs
4 sq --user ${SOMEONE_ELSE}

```

If you want to see all jobs in the queue:

```

1 squeue

```



For an optional exercise you may try to re-run the workflow with `conda`. For that, you'll need to find and load a `conda` module before re-running the workflow with appropriate profile. Don't forget to use the `-resume` flag.

If you remembered to use `-resume`, why do you think it appeared to not make a difference?

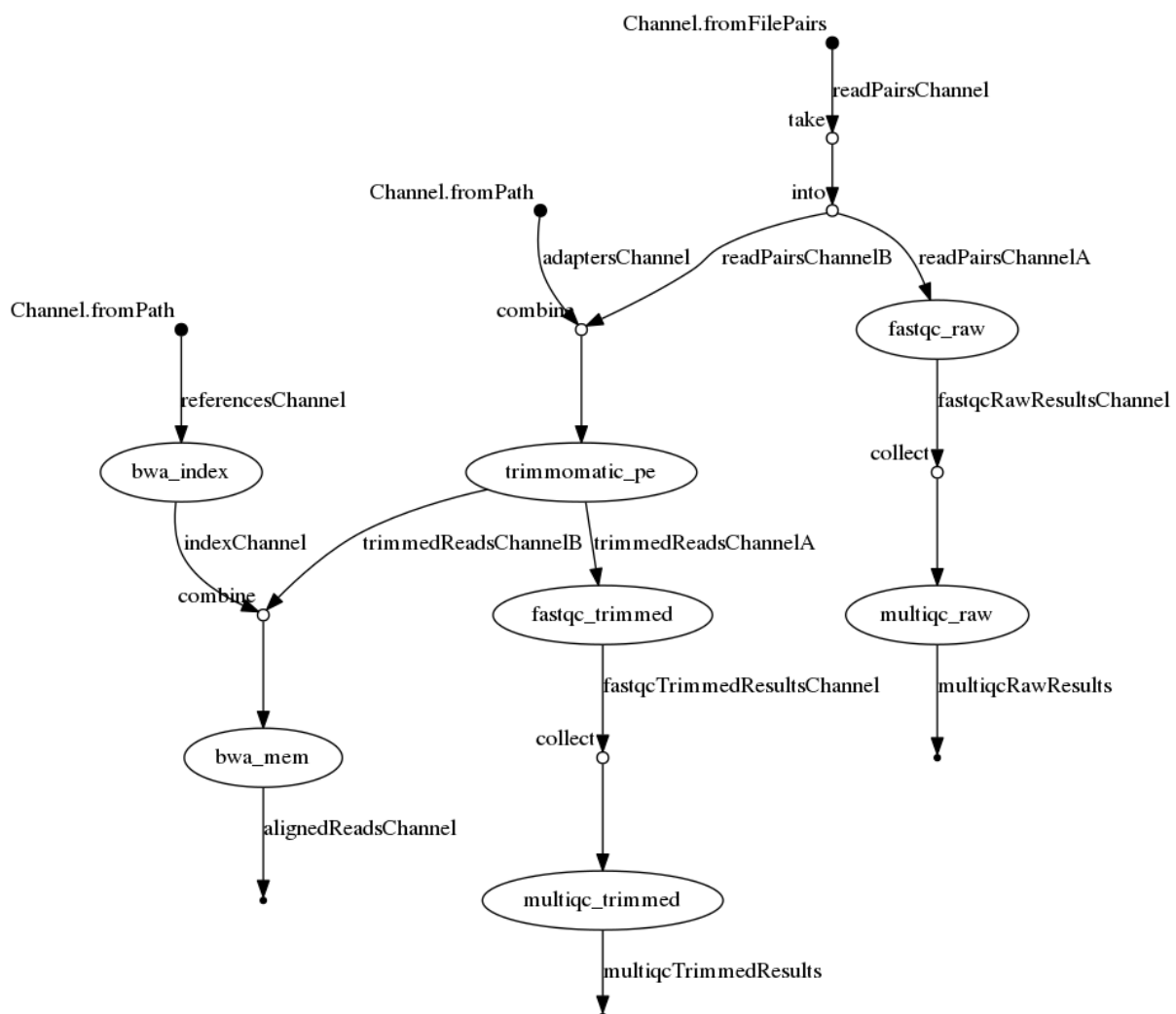


Figure 7: The example workflow



Investigate `main.nf` alongside Figure 7.

Which [nextflow operators](#)^a, in addition to the previously discussed, are used and for what purposes?

^a<https://www.nextflow.io/docs/latest/operator.html>

Workflow outputs

We now have each task nicely isolated in a separate sub-directory under `work`, but how do I find my results? Was it `work/a7:fc9339a827fb4b34d2408e1c3ee29c` or maybe `work/3c:8fdf958e96b448ecb83bd7806af382`? This should be handled by applying the [publishDir directive](#)⁶ to selected processes. As with other directives, this can be included at the top of the process block or in a configuration file using [process selectors](#)⁷ to apply the directive to one or more relevant process. To keep things tidy-ish, we define the publishing of the outputs in a separate file which we `includeConfig 'conf/publish.config'` in `nextflow.config`.

In `conf/publish.config` we only really use the `withName` selectors. The alternative `withLabel` selectors are convenient e.g. when outputs of multiple processes are to be gathered in one location, in which case we attach the same `label` to each of those processes.

Modify/extend the workflow



Edit `main.nf`. Your task is to add a process which will merge the bam files produced by the `bwa.mem` process.

⁶<https://www.nextflow.io/docs/latest/process.html#publishdir>

⁷<https://www.nextflow.io/docs/latest/config.html#process-selectors>



How do you ensure that **all** BAM files end up in the same instance of your process? Demonstrate your process definition to your facilitator.

Where can we find the merged BAM file? Can you publish it to a human-readable location? Hint: only declared outputs can be published.



Modify your merge process to allow samtools to use 2 cpus with `--threads 2`, don't forget to modify your process configuration to request 2 cpus per task.

Your own workflow (TODO: replace with variant calling?)

It is time to have a go at your own pipeline. Since we have some inputs and configuration files at hand, you can start a `own.nf` script file in the current directory and read the input files from `./data`.



The simple pipeline should include the following:

- Code for reading FASTQ read files from `./data` individually (i.e. not as pairs) into a channel.
- A process which will take a read file, count the reads and output the file name alongside the read count.
- A way of aggregating the individual count files into a single csv file. This could be done in another process or using an operator.

Troubleshooting

Disconnected from the cluster?

Missing modules - new shell session?



Make sure all the required modules are loaded.

```
1 # Java - essential for nextflow
2 module load openjdk-1.8.0_202-b08-gcc-5.4.0-sypwasp
3
4 # Singularity - our go to system for providing software for the example \
   workflow
5 module load singularity-3.2.1-gcc-5.4.0-tn5ndnb
6
7
8 # If using conda
9 module load miniconda3-4.6.14-gcc-5.4.0-kkzv7zk
```

Space for Personal Notes or Feedback

[illegible]

[illegible]

[illegible]

[illegible]