

TRAINER'S MANUAL

Implementing Scalable Bioinformatic Workflows in Snakemake and Nextflow

Nathan S. Watson-Haigh
Radosław Suchecki

Across Australia

Aug/Sep 2019

TRAINER'S MANUAL

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	7
The Trainers	8
Providing Feedback	9
Document Structure	9
Provided Compute Infrastructure	10
Connecting to the Cluster	11
Monitoring Slurm Jobs	11
Introduction to Snakemake	13
Key Learning Outcomes	14
Resources Required	14
Tools Used	14
Useful Links	14
Snakemake is Like Making Sunday Dinner	15
Setting Up Your Environment	15
Installing Snakemake	15
Your First Snakemake Workflow	17
Generalising Rules with Wildcards	18
Submitting Jobs to Slurm	19
Reimplementing a Workflow in Snakemake	20
Getting the Code	20
Getting the Data	21
Implementation of BWA Indexing	21
Implementing FastQC	24
Implementing Trimmomatic	27
Implementing BWA-MEM	29
Adding New Samples	32
Large Workflows	34
Troubleshooting	36
Getting Going After a Disconnect	36
Introduction to Nextflow	37
Key Learning Outcomes	38
Resources Required	38
Tools Used	38

Useful Links	39
Introduction	40
Setting Up Your Environment	41
Connect to the Cluster Head Node	41
Install nextflow	42
Standard installation	42
Conda installation	42
Running Nextflow	45
Hello (nextflow) world!	45
Hello command line options	46
Nextflow basics	47
Processes and channels	47
The script	47
Hello HPC!	48
Reimplementing a workflow in Nextflow	51
Getting the code	51
Getting the Data	51
The software environment	52
The config file(s) and profiles	54
Reimplementing a workflow in Nextflow - a walk-through	56
Implementing FastQC	57
Implementing BWA Indexing	60
Implementing Trimmomatic	62
Implementing BWA-MEM	65
Under the hood	67
Workflow outputs	68
Bonus tasks	68
Troubleshooting	70
Disconnected from the cluster?	70
New shell session?	70
Space for Personal Notes or Feedback	73

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

Provided Compute Infrastructure

For the purposes of this training, we are providing you with access to a virtual cluster running on top of AWS. The specification of this cluster is as follows:

Head Node

m5d.4xlarge (16 vCPU, 64G RAM and 2x300 SSD)

Compute Nodes

t2.medium, min: 30 max: 50

Shared Storage

1000G

Connecting to the Cluster



First up, lets 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.

Upon connecting, feel free to use **screen** or **tmux** if you are familiar with either of those tools.

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**:

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

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>

Snakemake is Like Making Sunday Dinner

The way Snakemake approaches running workflows is a bit like the way you prepare for dinner/tea. First, you think about what you want for dinner/tea, say a Sunday roast. To create the Sunday roast, you need meat and vegetables, all of which need preparing (e.g. peeling and seasoning). If you haven't got an ingredient, you go out to the shop and buy it.

With Snakemake, you decide what output files you want to create, say some BAM files. To create the BAM files, you need FASTQ files and a reference genome, all of which need preparing (e.g. quality/adaptor trimming and indexing). If you haven't got those files, you need to create them.

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.

Installing 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 . ${HOME}/.bashrc
15
16 # Change the default location into which conda saves packages
17 # and environments
18 conda config --prepend pkgs_dirs /shared/${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 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.

```
1 | conda update \
2 | snakemake
```

Your First Snakemake Workflow

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/${USER}/snakemake/hello
2 | cd /shared/${USER}/snakemake/hello
```

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?

Snakemake found the requested output file was already there so decided not to regenerate it.

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" (<https://snakemake.readthedocs.io/en/stable/snakefiles/rules.html#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 - nothing should be done since the file
2 # already exists
3 snakemake Hello/World.txt
4
5 # In Polish
6 snakemake Czesc/Swiat.txt
7
8 # In French and Spanish at the same time
9 snakemake Monde/Monde.txt Ciao/Mondo.txt
```

Take a look at the files created and their contents:

```
1 tree ./
2 cat */*.txt
```

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 let's 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 "hello" branch
9 git checkout hello
```

The Snakefile in this branch of the repository is the same “Hello, World!” example you created above, with wildcards. Let's 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
4
5 # See what files we have
6 tree
```

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`. See:

```
1 tree logs/
```

We've finished with this simple “Hello, World!” example, so cleanup after yourself:

```
1 snakemake \
2   --delete-all-output \
3   Hello/World.txt Czesc/Swiat.txt Monde/Monde.txt Ciao/Mondo.txt
4
5 # See what files we have
6 tree
```

Reimplementing a Workflow in Snakemake

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 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 Code

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

```
1 git checkout walkthrough
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` argument)

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

With difficulty! Enter - workflow management systems like Snakemake or Nextflow.

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
```

Implementation of 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 # These commands have the same effect
4 snakemake --dryrun all
5 snakemake --dryrun
```



Why doesn't this work:

```
1 | snakemake --dryrun bwa_index
```

Because the rule makes use of wildcards. Effectively, the wildcards are NULL values when we specify the rule to run. We need to specify the target filenames and they will get matched against the wildcards present in the defined **output**.

Do we have to specify all 5 of the BWA index files in the **all** pseudo-rule?

No. Snakemake determines which rule needs to be run in order to create each of the 5 index files specified. Since all these files are created by the **bwa_index** rule, snakemake knows it only needs to run that rule once in order to create all 5 files. We could easily have specified only one of the index files and this would have resulted in the same jobs being run.

Why do we have to specify all 5 index files as **output** of the **bwa_index** rule?

Because we want Snakemake to keep track of which files get generated as part of the workflow. This is how we are able to define dependencies between rules.

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 thought, 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
```

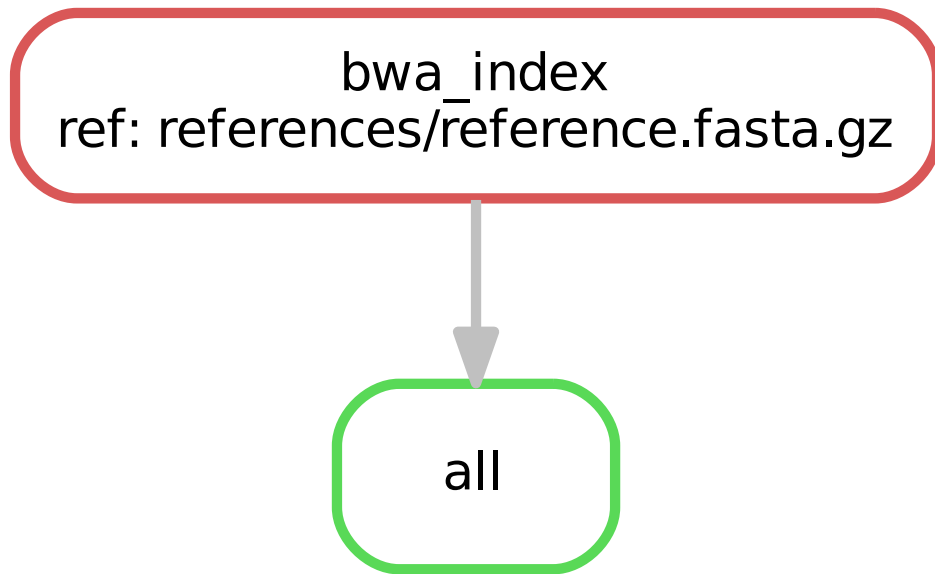


Figure 1: DAG of jobs showing `bwa_index` job dependant on the `all` pseudo-rule.

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 = [  
7             "raw_reads/{SAMPLE}_R1_fastqc.zip",  
8             "raw_reads/{SAMPLE}_R2_fastqc.zip",  
9         ],  
10        html = [  
11            "raw_reads/{SAMPLE}_R1_fastqc.html",  
12            "raw_reads/{SAMPLE}_R2_fastqc.html",  
13        ],  
14    shell:  
15        ""  
16        fastqc --threads 1 {input.r1} {input.r2}  
17        ""
```

Improving FastQC Parallelisation

There are a few improvements we can make to the `fastqc` 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.

Change the `fastqc` rule to the following:

```
1 rule fastqc:  
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 {threads} {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
```

Rule-Specific Resource Specification

The profile provided in this repository specifies default values for Slurm resources. However, it is possible to provide rule-specific overrides so that jobs can make use of more time, memory, cores etc. The ‘cluster-configs/default.yaml’ file is where these settings can be modified.

```
1  less cluster-configs/default.yaml
```

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}",  
4             ext=['amb', 'ann', 'bwt', 'pac', 'sa']  
5         ),  
6         expand("raw_reads/{SAMPLE}_{read}_fastqc.html",  
7             SAMPLE=SAMPLES,  
8             read=['R1', 'R2']  
9         ),
```

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  
5 # "all" rule by convention  
6 snakemake --dryrun
```

Lets look at the DAG for the workflow:

```
1 snakemake \  
2 --dag \  

```

```

3 | dot -Tpdf \
4 | > dag2.pdf

```

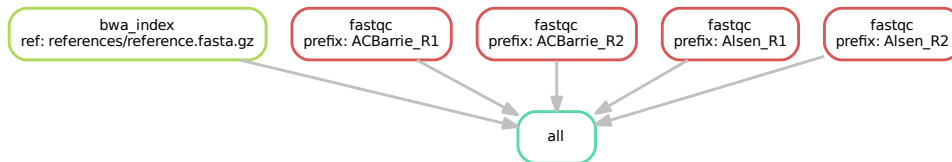


Figure 2: DAG of jobs showing `bwa_index` and several `fastqc` job dependant on the `all` pseudo-rule.

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 {threads} \
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}",
4             ext=['amb', 'ann', 'bwt', 'pac', 'sa']
5         ),
6         expand("raw_reads/{SAMPLE}_{read}_fastqc.html",
7             SAMPLE=SAMPLES,
8             read=['R1', 'R2']
9         ),
10        expand("qc_reads/{SAMPLE}_{read}.fastq.gz",
11            SAMPLE=SAMPLES,
12            read=['R1', 'R2']
13        ),

```

We won't run the workflow in a piecemeal fashion, we'll save all our jobs for running a bit later. So for now, let's look at the dryrun again:

```

1 # Run the first rule in the Snakefile. This should be the
2 # "all" rule by convention

```

```
3 | snakemake --dryrun
```

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?

QC'd R1 reads for a sample: `qc_reads/{SAMPLE}_R1.fastq.gz` QC'd R2 reads for a sample: `qc_reads/{SAMPLE}_R2.fastq.gz`

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

The rule doesn't need the FASTA, it needs the index files:

- `references/reference.fasta.gz.amb`
- `references/reference.fasta.gz.ann`
- `references/reference.fasta.gz.bwt`
- `references/reference.fasta.gz.pac`
- `references/reference.fasta.gz.sa`

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?

```
1 | shell:
2 |     """
3 |     bwa mem -t {threads} \
4 |         references/reference.fasta.gz \
5 |         {input.r1} {input.r2} \
6 |         | samtools view -b \
7 |         > {output}
8 |     """
```



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>).

```
1  params:
2      prefix = "references/reference.fasta.gz",
3  shell:
4      """
5      bwa mem -t {threads} \
6          {params.prefix} \
7          {input.r1} {input.r2} \
8          | samtools view -b \
9          > {output}
10     """
```



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.

```

1 input:
2 reference = expand("references/reference.fasta.gz.{ext}",
3 ext=["amb", "ann", "bwt", "pac", "sa"]
4 ),
5 ...
6 params:
7 prefix = lambda wildcards, input: input["reference"][0][: -4],
8 shell:
9 """
10 bwa mem -t {threads} \
11     {params.prefix} \
12     {input.r1} {input.r2} \
13     | samtools view -b \
14     > {output}
15 """

```

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.

```

1 input:
2 reference = expand("references/reference.fasta.gz.{ext}",
3 ext=["amb", "ann", "bwt", "pac", "sa"]
4 ),
5 ...
6 params:
7 prefix = lambda wildcards, input: input["reference"][0][: -4],
8 shell:
9 """
10 bwa mem -t {threads} \
11     {params.prefix} \
12     {input.r1} {input.r2} \
13     | samtools view -b \
14     > {output}
15 """

```

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

```

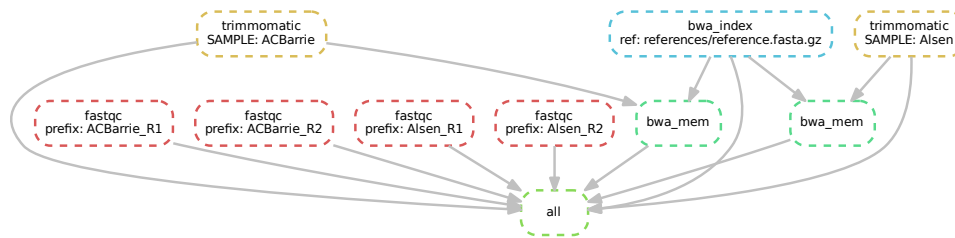


Figure 3: DAG of jobs showing the dependencies which exist in our final implementation of the `analysis.sh` workflow.

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/^# \{2\}\"/ "/" Snakefile

```

With so many more samples, the DAG becomes next to useless with `dot`'s default layout algorithm:

```

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
```

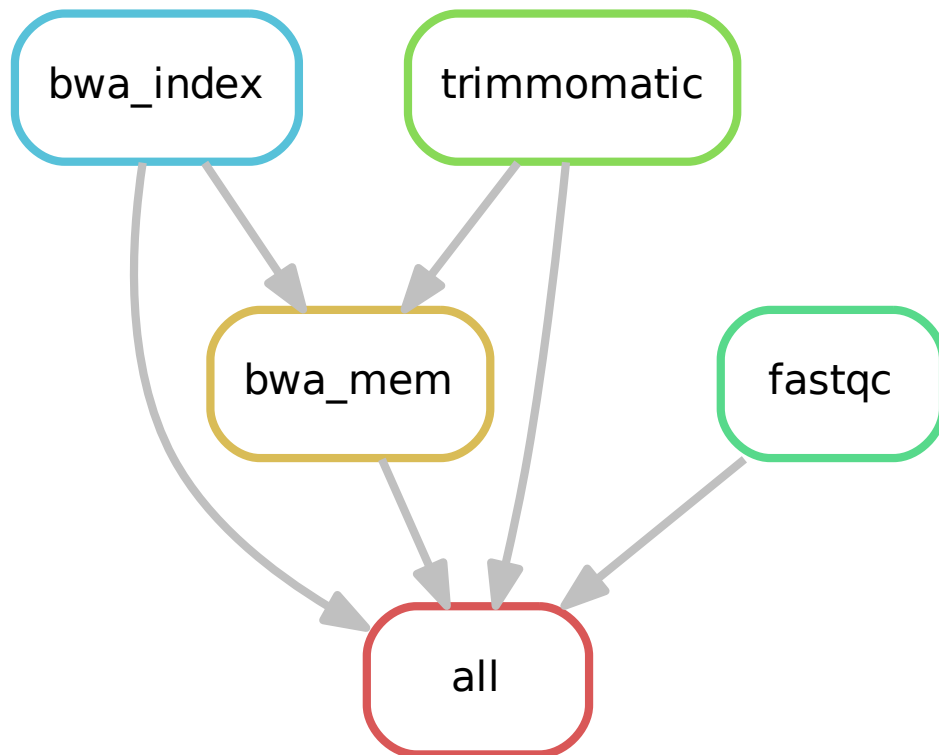


Figure 5: Rulegraph for the whole workflow.

Execute the rest of the workflow:

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

In a different terminal, have a look at your jobs in the queue:

```
1 sq
```



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

fastqc: 32

Total: 66

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**?

```
1 snakemake \  
2 --dryrun \  
3 --printshellcmds \  
4 --forceall
```

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

```
1 snakemake \  
2 --delete-all-outputs
```

Large Workflows

A workflow I have developed is capable of mapping RNA-Seq, WGS, exome capture and Iso-Seq data to the wheat genome. In addition, downstream variant analysis and creation of various tracks for JBrowse (e.g. BigWig) are possible. Much of the workflow involves broking the tasks down so each of the chromosomes is proccessed independantly, and in parallel by the workflow. All up, there are more than 10,000 individual jobs representing 14 rules. DAG's of this size are best opened and layed out by Gephi using the **ForceAtlas 2** algorithm:

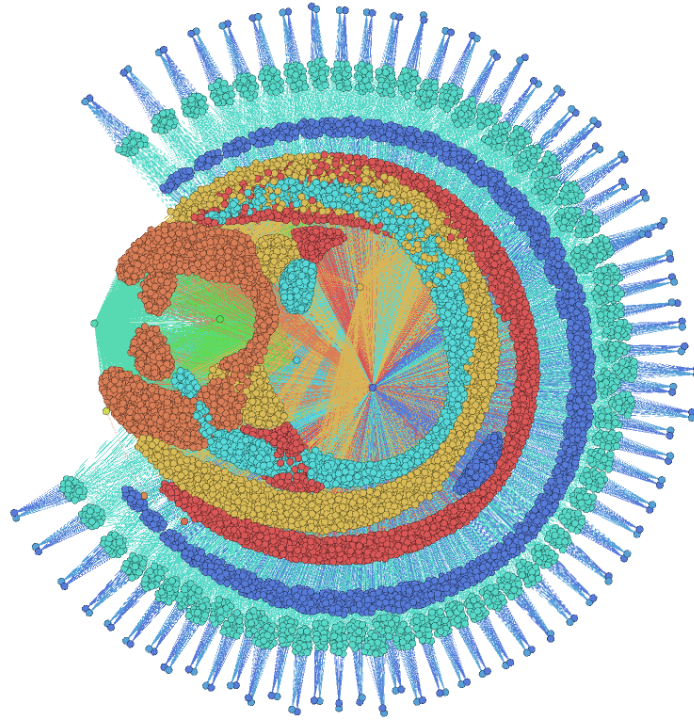


Figure 6: Large DAG, from the DAWN workflow, consisting of 10,587 nodes/jobs and 32,353 edges/dependencies. Layed out using ForceAtlas 2 algorithm of Gephi

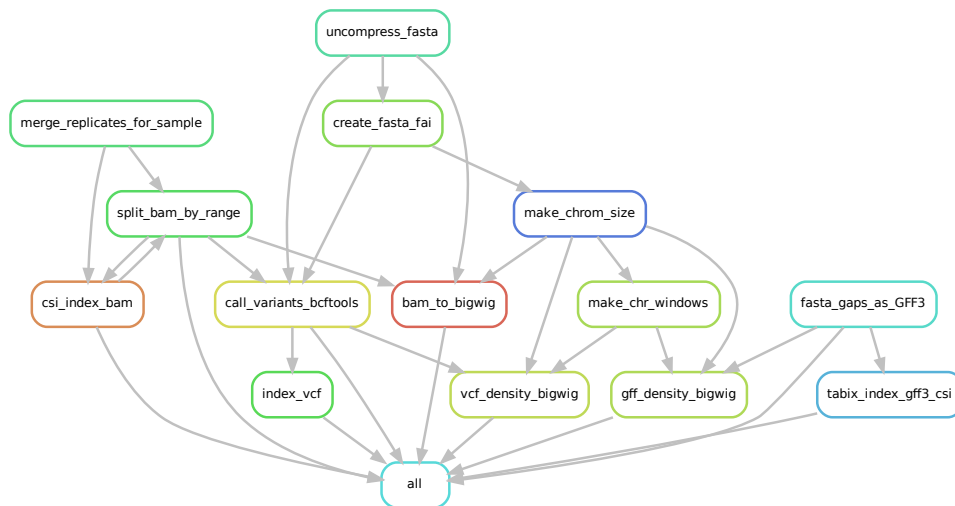


Figure 7: Rulegraph, from the DAWN workflow.

Troubleshooting

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
7 # integrate shell autocompletion into bash
8 conda activate snakemake
9 complete -o bashdefault -C snakemake-bash-completion snakemake
10
11 # Move to the correct directory location, most
12 # likely here:
13 cd /shared/${USER}/snakemake/tutorial
```

Introduction to Nextflow

Primary Author(s):

Radosław Suhecki rad.suhecki@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

Standard installation



Do NOT run the following commands! This is provided for future reference so you know how to install Nextflow on another system.

```
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 -version
```

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.

Conda installation



```
1 # Load miniconda
2 module load \
3     miniconda3-4.6.14-gcc-5.4.0-kkzv7zk
4
5 #####
6 # One-time commands
7 #####
8 # Integrate conda into bash
9 conda init bash
10 . ${HOME}/.bashrc
11
12 # Change the default location into which conda saves packages
13 # and environments
14 conda config --prepend pkgs_dirs /shared/${USER}/.conda/pkgs
```

```
15 conda config --prepend envs_dirs /shared/${USER}/.conda/envs
16
17 # Change the default channels used for finding software and
18 # resolving dependencies
19 conda config --add channels defaults
20 conda config --add channels bioconda
21 conda config --add channels conda-forge
22 #####
```



Do NOT run the following command! This is provided for future reference so you know how to install Nextflow on another system using conda. 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 nextflow using conda
2 conda create \
3   --name nextflow \
4   --yes \
5   nextflow=${NEXTFLOW_VERSION:-19.04.0}
```

Nextflow installation is now complete.

Now simply copy the following `.conda` directory and you will have Nextflow 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 `nextflow` available on the command line:



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

For the duration of the workshop we may fix the version of nextflow we will use, and also opt for simpler but more verbose terminal logging which may be more useful when you get started with Nextflow or when you start developing a new pipeline.



```
1 export NXF_VER=19.04.0
2 export NXF_ANSI_LOG=false
3
```

```
4 | # now try
5 | nextflow -version
```

To revert to the default logging you can either change the value of the environmental variable to `false` or use `-ansi-log false` at run time.

Running Nextflow

We start with Nextflow take on ‘Hello world’

Hello (nextflow) world!

Set-up a working directory for this task:



```
1 | mkdir -p /shared/${USER}/nextflow
2 | cd /shared/${USER}/nextflow
```

Clone the hello-world example and go to its directory



```
1 | git clone https://github.com/rsucheki/hello.git
2 | cd hello
```

Feel free to investigate the content of the `main.nf` script, but we will have a closer look further on.

```
1 | cat main.nf
```

Run the `main.nf` script and observe the terminal output

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

Among the terminal output you should see that the local executor is being used and that a number of `sayHello` processes (tasks) have been submitted. You should also see the hellos being printed to the terminal.

Run the hello script again a few times and observe the effect of using the `-resume` flag

```
1 | nextflow run main.nf
2 | nextflow run main.nf -resume
```

You may run the hello script a few more times with and without `-resume`.



How does the output differ depending on whether the `-resume` flag is used? **With `-resume`, Nextflow will not re-run the tasks and report them as cached.**

Hello command line options

Single-dashed options are reserved for Nextflow engine (`-resume`, `-ansi-log false` etc). The double-dashed options are all yours and you are free to use them for your workflow.

In the ‘hello’ example we use the variable `params.universe` which is by default set to ‘World’. To set the variable to a different value, say ‘Mundo’, we use the variable name prefixed by a double-dash as runtime option and the new string as its parameter, e.g. `--universe Mundo`.

The value of the parameter (‘Mundo’ in this case) will be accessible in `main.nf` as

- `params.universe` in “plain” code
- `$params.universe` in GStrings (interpolated strings) – including the script block



Lets re-run the script.

```
1 | nextflow run main.nf -resume --universe Mundo
```



What effect did the `-resume` flag had this time? The use of `-resume` was inconsequential as the process `sayHello` was never previously executed with input that includes the value of `params.universe` set to ‘Mundo’ rather than ‘World’



When the pipeline is launched with the `-resume` option, any attempt to execute previously executed process with the same inputs will cause the process execution to be omitted, producing the previously stored outputs.

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

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.

For more on task caching see <https://www.nextflow.io/docs/latest/process.html#cache>

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.

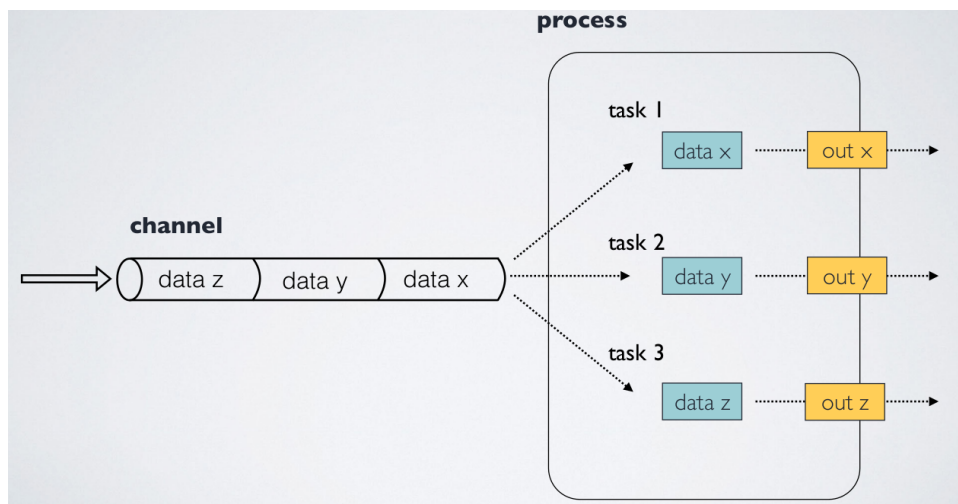


Figure 8: Nextflow building blocks: a *channel* “feeding” a processes. A *task* is an instance of a process. An isolated task is created for each emission (data chunk) from the input channel. *Credit: Evan Floden*

The 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 |
3 | cheers = Channel.from 'Bonjour', 'Ciao', 'Hello', 'Hola', 'Czesc'
4 |
5 | params.universe = 'World'
6 |
7 | process sayHello {
8 |     echo true
9 | }

```

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

```
10 input:
11     val cheer from cheers
12
13 script:
14     """
15     echo "$cheer $params.universe!"
16     """
17 }
```

A channel called `cheers` is created and emits each of the listed strings separately. A separate instance of the process `sayHello`, i.e. a *task*, is executed for each emission.

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

- The shebang line (line 1) is optional.
- `Channel.from(some_list)` creates a channel emitting the list elements one by one.
- `params.universe = 'World'` sets the default value of the command-line accessible option (e.g. `--world Mundo`)
- [Process](#) definition (lines 7-16)
 - Directives are placed at the top, directive `echo` controls whether `stdout` of the process is redirected to your terminal
 - Input block (lines 10-11)
 - Script block (lines 13-16)
 - * It is assumed to be a bash script unless an alternative shebang line (e.g. `/usr/bin/env python3`) is specified at the top of the block
 - * The `$cheer` in the script block is a nextflow variable local to the process, **not** a bash variable.
- Indentation is inconsequential.

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.

Thanks to Nextflow's out-of-the-box support for different compute environments, this only requires changing the process executor from the default `local` to `slurm`. This could be done via configuration files (more about that later) or using the directive `executor 'slurm'` in the process definition. However, for our simple use-case it should suffice to use the `NXF_EXECUTOR` environmental variable. Try that without `-resume` to avoid the cached results being reported.



```
1 NXF_EXECUTOR=slurm nextflow run main.nf
```



Play spot-the-difference – how can you tell if the tasks were submitted on the cluster? The terminal output should include the line `executor > slurm (5)` (or similar), but may differ depending on the version of Nextflow and whether ANSI logging is switched on or off.



Let's modify our script slightly to make it easier to see how processes are executed. Replace the single line within the script block

```
echo "$scheer $params.universe!"
```

with the following

```
echo "$scheer ${params.universe}! -- from ${task.executor} \${HOSTNAME}"
```



Note the difference between how nextflow variables (`$x`, `$params.universe`) 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. See for example [Nextflow documentation of the alternative shell block^a](https://www.nextflow.io/docs/latest/process.html#process-shell) or consider using alternative string delimiters, particularly the [dollar slashy^b](http://groovy-lang.org/syntax.html#_dollar_slashy_string). More broadly, remember that

- single- and triple-single-quote delimited strings are literal
- double- and triple-double-quote delimited strings are GStrings^c (interpolated strings).

^a<https://www.nextflow.io/docs/latest/process.html#process-shell>

^bhttp://groovy-lang.org/syntax.html#_dollar_slashy_string

^chttps://groovy-lang.org/syntax.html#_string_interpolation



Now re-run...

```
1 NXF_EXECUTOR=slurm nextflow run main.nf
```

The messages printed to the terminal by each task should now include the name of the executor being used and the name of the compute node on which the task was executed.



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* develop 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.

Let’s have a play with a slightly more practical workflow which will help us with exploring other features of Nextflow.

Reimplementing a workflow in Nextflow

In the Snakemake section of our tutorial, we include an example bash “pipeline”, which includes a number of common tasks wrapped in a single [bash script](#)². The workflow includes the following steps

- Run FastQC across the raw read (FASTQ) files
- Adapter, quality, and read length filtering using Trimmomatic
- Index the reference FASTA file
- Perform a `bwa-mem` read alignment

Before we start developing the Nextflow script, we need to set-up

- shared code base to be able to keep our efforts in sync when required
- input data for our workflow
- software environment including third-party tools

Getting the code



Start by cloning the workflow repository and moving to the directory.

```
1 mkdir -p /shared/${USER}/nextflow
2 cd /shared/${USER}/nextflow
3 git clone https://github.com/rsucheki/nextflow-walkthrough.git
4 cd nextflow-walkthrough
```

Getting the Data

We have provided you with some real whole genome sequencing (WGS) data from bread wheat together with a small chunk of the wheat genome. The data set is small enough for each step in the analysis to 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/ \
4   ./
5
```

²https://github.com/UofABioinformaticsHub/snakemake_template/blob/tutorial/analysis.sh

```
6 | # Have a look at what files we'd provided
7 | tree data
```



To be able to run this tutorial outside the HPC cluster provided for the workshop you may use the provided script to download the required input files

```
1 | #nextflow run setup_data.nf
```

The software environment

There are multiple ways in which you can provision a software environment for your workflow. You may have all the required tools installed on your system or available as a collection of pre-compiled binaries. Perhaps a friendly sysadmin installs convenient modules on your HPC cluster? You may also be using conda or perhaps docker or singularity containers? All those options are compatible with Nextflow.

We have provided a rudimentary script `versions.nf` which tries to run each of the software tools, and report their version numbers.

```
1 | #!/usr/bin/env nextflow
2 |
3 | process get_versions {
4 |     echo true
5 |
6 |     script:
7 |     """
8 |     fastqc --version
9 |     multiqc --version
10 |    echo -n "bwa " && bwa 2>&1 | grep 'Version'
11 |    echo -n "samtools " && samtools 2>&1 | grep 'Version'
12 |    """
13 | }
```



Try to run

```
1 | nextflow run versions.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`).



Identify the exact **error**. The cause was likely “.command.sh: line 2: bwa: command not found”. See line 23 in the following example:

```

1 N E X T F L O W ~version 19.04.0
2 Launching `versions.nf` [elegant_morse] - revision: 525c4d8fea
3 [warm up] executor > local
4 [0a/26fff6] Submitted process > get_versions
5 ERROR ~Error executing process > 'get_versions'
6
7 Caused by:
8   Process `get_versions` terminated with an error exit status (127)
9
10 Command executed:
11
12   fastqc --version
13   multiqc --version
14   echo -n "bwa " && bwa 2>&1 | grep 'Version'
15   echo -n "samtools " && samtools 2>&1 | grep 'Version'
16
17 Command exit status:
18   127
19
20 Command output:
21   (empty)
22
23 Command error:
24   .command.sh: line 2: fastqc: command not found
25
26 Work dir:
27   /shared/training001/nextflow/example_workflow/work/0a/26fff6b11f473b0c3f0b3961608462
28
29 Tip: view the complete command output by changing to the process \
   work dir and entering the command `cat .command.out`
30
31 -- Check '.nextflow.log' file for details

```

There are different ways in which we could ensure that the relevant software is available, for example using process *directives* at the top of each process definition. We could also use the opportunity to change the executor from local to Slurm, for example:

```

1 process get_versions {
2   module 'fastqc-0.11.7-gcc-5.4.0-z5pkqvq'
3   executor 'slurm'
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 (primarily) in the `nextflow.config` file. Transferring the above mention *directives* from process definitions in `main.nf` to `nextflow.config` would make things slightly better. The directives must be applied in the *process config scope*³ in `nextflow.config`.

```
1 process.executor = 'slurm'
2 process.module = 'fastqc-0.11.7-gcc-5.4.0-z5pkqvq'
```

or using the preferred syntax

```
1 process {
2     executor = slurm
3     module = 'fastqc-0.11.7-gcc-5.4.0-z5pkqvq'
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 from `main.nf` while providing the required software either by creating a `conda` environment or by using `Docker` or `Singularity` containers in which the `conda` environment has already been captured.



A separate point is that you may want to define different software or compute environments for different processes. Profiles can go hand-in-hand with *process selectors*^a which make it easy to apply configuration, (including software and compute config) to individual processes or subsets of processes.

```
1 process {
2     withName: get_versions {
3         executor = 'slurm'
4         cpus = 1
5         time = 1.m
```

³<https://www.nextflow.io/docs/latest/config.html#config-scopes>

```
6     memory = 256.MB
7   }
8   withLabel: snail {
9     time = 48.h
10  }
11 }
```

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

We can select processes by name or by a label assigned to one or more processes via the [label directive](#)⁴. Process selectors can also be really useful outside profile definitions – we will revisit them later.

Relevant profiles



Have a look inside `nextflow.config`, and locate the process definitions

```
1 less nextflow.config
```

The ones most immediately relevant are:

```
1 profiles {
2   //EXECUTORS
3   slurm {
4     process {
5       executor = 'slurm'
6     }
7   }
8   //SOFTWARE
9   conda {
10    process {
11      conda = "$baseDir/conf/conda.yaml"
12    }
13  }
14  singularity {
15    process {
16      container = '/shared/.singularity/nextflow-embl-abr-webinar.simg'
17    }
18    singularity {
19      enabled = true
20      autoMounts = true
21    }
22  }
23 }
```

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

As you can see, Nextflow makes it really easy to define software environment via Singularity or Conda⁵.

Given that Singularity is available on our cluster, we will use the `singularity` profile for software environment using `-profile singularity`. Combinations of profiles can be selected at runtime so we will also use the `slurm` profile which will ensure the tasks are not executed on the head node but submitted to the cluster.

In addition to this basic configuration, there are many settings that can and in some cases must be set – refer to [executors section of Nextflow documentation](#)⁶. For running real-life pipelines in a cluster environment you will also use [directives](#)⁷ controlling the resources (`cpus`, `memory`, `time`) requested for each job. Other directives relevant in HPC context might include `queue` and `scratch`.



```
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 versions.nf -profile slurm,singularity
```



Setting up profiles and keeping configuration separate from the main script is an upfront investment which will benefit you many times over. Many aspects of the configuration are consistent across quite different workflows so you get to re-use and refine them with each new pipeline you work on. This feeds back to your workflows being more portable and robust.

Reimplementing a workflow in Nextflow - a walk-through

Recall the steps of the original bash workflow

- Run FastQC across the raw read (FASTQ) files
- Adapter, quality, and read length filtering using Trimmomatic
- Aggregating FastQC reports from the raw reads using MultiQC
- Index the reference FASTA file

⁵We also have a docker profile which you may find useful if you decide to run the workflow on your machine

⁶<https://www.nextflow.io/docs/latest/executor.html>

⁷<https://www.nextflow.io/docs/latest/process.html#directives>

- Perform a `bwa-mem` read alignment

We will walk you through reimplementing the first few steps into a Nextflow workflow. Along the way, we will introduce the core concepts of Nextflow and then ask you to reimplement the `bwa-mem` step yourselves. If you finish that with some time to spare, you will have the opportunity to reimplement the `multiqc` step and tackle other bonus exercises. This will provide you with a foundation for you to be able to port your own workflows into Nextflow.

Implementing FastQC

In this section you will learn how to

- Read input files into a *channel*
- View items emitted by a *channel*
- Restrict how many items are emitted through a *channel* using *operators*
- Effortlessly parallelize the task by plugging a *channel* into a *process*
- Reference emitted items within processes' script block



Paste and save the following into a new file `main.nf`

```
1 Channel.fromPath("data/raw_reads/*.fastq.gz")
2   .view()
```

We create a (yet unnamed) channel and use the `view` operator to print the names of emitted files to the terminal. If you now run

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

You will see a long list of files being emitted by the channel.

While developing a workflow we tend to prefer to only use a small subset of data. We could restrict the number of FASTQ files emitted by the channel e.g. by pointing to a specific file or a subset of files, in this case the two read files available for Baxter

```
1 Channel.fromPath("data/raw_reads/Baxter_R{1,2}.fastq.gz")
2   .view()
```

We can do better! Let's use one of NF operators capable of limiting the number of items emitted by the channel. Either `one` of `first()`, `last()` or `take(n)` will suffice.

```
1 Channel.fromPath("data/raw_reads/*.fastq.gz")
2   .take(1)
3   .first()
```

```
4 .last()
5 .view()
```

Having all three operators is redundant, we will stick with *take*, due to extra flexibility it gives us. The last thing to do is to assign our channel to a variable whose name we can use elsewhere in the script. The simplest form would be

```
1 readsForQcChannel = Channel.fromPath("data/raw_reads/*.fastq.gz")
2   .take(1)
3   .view()
```

But since we often chain multiple operations, the following syntax is preferred due to consistent left-to-right flow.

```
1 Channel.fromPath("data/raw_reads/*.fastq.gz")
2   .take(1)
3   .view()
4   .set { readsForQcChannel }
```

For consistency, we will stick to this syntax throughout this workshop

Update `main.nf` and run it again

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

You should see that only a single item is emitted by the channel.

Now that we have our input channel ready and emitting, we can remove the view operator and focus on defining a process to implement the original bash snippet.

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

In this case, there is no benefit from processing files as pairs. In fact, separating the pairs should in principle allow more efficient allocation of resources as we will end up with twice the number of tasks, but each being half the size of the original.

In `main.nf` we start our process definition by giving it a name and defining the input block

```
1 process fastqc {
2
3   input:
4     file(reads) from readsForQcChannel
5
6 }
```

We defined process input to be file taken from the `readsForQcChannel` which we assign to variable `reads`. We can now proceed to defining our script block.

```

1 process fastqc {
2
3   input:
4     file(reads) from readsForQcChannel
5
6   script:
7     """
8     fastqc --threads 1 ${reads}
9     """
10  }

```

In the script block we address the input file via its variable name as `$reads` or `${reads}`.



Rather than fixing the number of threads inside the script block to 1, we could use the `task.cpus` variable.

How can you set the value of the `task.cpus` variable? Hint: look into *process directives* documentation.

```

1 process fastqc {
2   cpus 2
3
4   input:
5     file(reads) from readsForQcChannel
6
7   script:
8     """
9     fastqc --threads ${task.cpus} ${reads}
10    """
11  }

```



We should now be able to run fastqc on a single FASTQ file

```

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

```

How about running on a handful of samples? We could simply increase the number in our `.take(1)` call, but rather than hard-coding that, we can use command line arguments. It is a good practice to set a default value either in `nextflow.config` or in `main.nf`.

```

1 params.n = 1 //set default value
2
3 Channel.fromPath("data/raw_reads/*.fastq.gz")
4   .take( params.n )
5   .set { readsForQcChannel }

```



Now run the workflow for $n = 5$ FASTQ files

```
1 | nextflow run main.nf -profile slurm,singularity -resume --n 5
```



How would can we run it for all samples? The *take* operator returns all emissions when set to -1, so `--n -1` should do it. If we'd rather have `--n all`, then we modify our code:

```
1 | .take ( params.n == all ? -1 : params.n )
```



If the above tasks caused you some un-recoverable issues you can rename or delete your `main.nf` and check-out a revision where the above steps have been captured.

```
1 | mv main.nf myfastqc.nf
2 | git checkout fastqc
```

You should now be able to

- Read input files into a *channel*
- View items emitted by a *channel*
- Restrict how many items are emitted through a *channel* using *operators*
- Effortlessly parallelize the task by plugging a *channel* into a *process*
- Reference emitted items within processes' script block

Implementing BWA Indexing

In this section you will learn how to

- group *process* output files and values in a **set** to be emitted through an output channel

In the original bash script this task consists of a simple call to `bwa` with the input file name hard-coded.

```
1 | bwa index -a bwtsw references/reference.fasta.gz
```

We start by creating a channel which pass the reference file to our indexing process.

```
1 Channel.fromPath('data/references/reference.fasta.gz').set { \
    referencesChannel }
```

We can now move on to defining the process which will *consume* the `referencesChannel`. The bare-minimum would probably be a process definition with an `input` block and a `script` block.

```
1 process bwa_index {
2   input:
3     file(ref) from referencesChannel
4
5   script:
6     """
7     bwa index -a bwtsw ${ref}
8     """
9 }
```

That should be sufficient to run the indexing, but given that we want to use the generated index, in another process, we also need to define the outputs for the current one. The expected output files are:

```
1 reference.fasta.gz.amb
2 reference.fasta.gz.ann
3 reference.fasta.gz.bwt
4 reference.fasta.gz.pac
5 reference.fasta.gz.sa
```

To capture these as outputs we could simply declare that all (non-input) files constitute the desired output.

```
1 output:
2   file("*") into indexChannel
```

We could also be slightly more explicit and declare that outputs are all (non-input) files sharing a prefix

```
1 output:
2   file("${ref}.*") into indexChannel
```

or even painstakingly list all the expected files... not recommended. However, one of the great things about Nextflow is that thanks to process isolation, you typically don't have to pay much attention to input/output file names unless a tool you use have some specific requirements in this regard. For example, to run `bwa` alignment later on, we will need to pass it the common prefix of the names of the index files.

In this case rather than simply outputting the index files, and then parsing the prefix out of the file name(s), we opt to output our reference file name alongside the generated index file. To do that in a single emission we wrap these in a *set* (in other words a tuple - an ordered list of elements).

```

1  output:
2  set val("${ref}"), file("*") into indexChannel

```

The emitted set will consist of two elements

- The reference file name (String)
- The list of index files (not just file names – objects of class `java.nio.file.Path`).

This is our complete indexing process definition

```

1  process bwa_index {
2
3  input:
4  file(ref) from referencesChannel
5
6  output:
7  set val("${ref}"), file("*") into indexChannel
8
9  script:
10  """
11  bwa index -a bwtsw ${ref}
12  """
13  }

```

You could temporarily add `indexChannel.view()` to see how the output looks.

Check if nothing is broken



```

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

```



If the above tasks caused you some un-recoverable issues you can rename or delete your `main.nf` and check-out a revision where the above steps have been captured.

```

1  mv main.nf myindexing.nf
2  git checkout index

```

Implementing Trimmomatic

We will guide you through the implementation of one more process before you take over! The task is to adapter/quality trim raw reads as pairs and output them into a channel to be consumed by an alignment process.

The original bash implementation was as follows

```

1  #####
2  mkdir -p qc_reads
3  trimmomatic PE \
4  -threads 1 \

```

```

5 raw_reads/${SAMPLE}_R1.fastq.gz raw_reads/${SAMPLE}_R2.fastq.gz \
6 qc_reads/${SAMPLE}_R1.fastq.gz \
   qc_reads/${SAMPLE}_R1.unpaired.fastq.gz \
7 qc_reads/${SAMPLE}_R2.fastq.gz \
   qc_reads/${SAMPLE}_R2.unpaired.fastq.gz \
8 ILLUMINACLIP:misc/trimmomatic_adapters/TruSeq3-PE.fa:2:30:10:3:true \
   \
9 LEADING:2 \
10 TRAILING:2 \
11 SLIDINGWINDOW:4:15 \
12 MINLEN:36

```

This time we have to process the reads **as pairs** so using `Channel.fromPath()` will not suffice. Conveniently, NF has been developed by bioinformaticians so it smoothly handles paired files.

```

1 Channel.fromFilePairs("data/raw_reads/*_R{1,2}.fastq.gz")
2   .take (1)
3   .view ()
4   .into { readPairsForTrimmingChannel }

```



As before, during development we

- use the `.take(n)` operator to limit the number of emissions coming through the channel
- use `.view()` to sneak-peak at what exactly is emitted

The output from the `.view()` operator in this case, after some re-formatting would be something like

```

1 [
2   Gladius,
3   [
4     /path/to/data/raw_reads/Gladius_R1.fastq.gz,
5     /path/to/data/raw_reads/Gladius_R2.fastq.gz
6   ]
7 ]

```

The single emission consists of a set of two elements.

1. The sample name (accession) i.e. the shared part of the file names captured by the `'*'` glob.
2. The list of two FASTQ files

One more thing we need, is a file containing the adapter sequences for trimmomatic. This could be done in a number of ways, but it is best to stick to the dataflow paradigm and set-up another channel for that.

```
1 Channel.fromPath('data/misc/trimmomatic_adapters/TruSeq3-PE.fa')
2 .set{ adaptersChannel }
```

We now have two input channels, but how do we handle multiple input channels?

If we just declare each separately

```
1 process trimmomatic {
2   input:
3     set val(sample), file(reads) from readPairsForTrimmingChannel
4     file(adapters) from adaptersChannel
5   ...
```

This would work perfectly fine... but only for the first pair of reads. There is only a single adapters file emitted by the `adaptersChannel`, so once this channel no longer emits outputs, the process will not run again and will never read the next pair of FASTQ files from `readPairsForTrimmingChannel`.

What we really want is for the adapters file to be used in combination with each pair of FASTQ files. Enter the `combine()` operator.

```
1 process trimmomatic_pe {
2
3   input:
4     set file(adapters), val(sample), file(reads)
5     from adaptersChannel.combine(readPairsForTrimmingChannel)
6   ...
```

This covers our inputs and we can now focus on the outputs and the script block. As mentioned earlier we don't have to be too specific about the output file names but just have to make it clear to NF which output files are to be sent to an output channel. In this case these are the ones with names ending with `paired.fastq.gz`. We also output the `sample` variable to keep track of which sample is being processed.

```
1 process trimmomatic_pe {
2
3   input:
4     set file(adapters), val(sample), file(reads)
5     from adaptersChannel.combine(readPairsForTrimmingChannel)
6   output:
7     set val(sample), file('*.paired.fastq.gz') into trimmedReadsChannel
8   ...
```

The translation of the trimmomatic call from plain bash to bash embedded in the script block is straightforward.

1. Explicit paths to the input files are replaced with `${reads}` variable. The read files could also be addressed individually as `${reads[0]}` and `${reads[1]}`.
2. Explicit output file paths for trimmomatic are replaced with convenient short names

which need not be globally unique.

```

1 process trimmomatic_pe {
2   input:
3     set file(adapters), val(sample), file(reads) from \
        adaptersChannel.combine(readPairsForTrimmingChannel)
4
5   output:
6     set val(sample), file('*.paired.fastq.gz') into trimmedReadsChannel
7
8   script:
9     """
10    trimmomatic PE \
11    ${reads} \
12    R1.paired.fastq.gz \
13    R1.unpaired.fastq.gz \
14    R2.paired.fastq.gz \
15    R2.unpaired.fastq.gz \
16    ILLUMINACLIP:${adapters}:2:30:10:3:true \
17    LEADING:2 \
18    TRAILING:2 \
19    SLIDINGWINDOW:4:15 \
20    MINLEN:36 \
21    -Xms256m \
22    -Xmx256m
23    """
24 }
```



If the above tasks caused you some un-recoverable issues you can rename or delete your main.nf and check-out a revision where the above steps have been captured.

```

1 mv main.nf mytrim.nf
2 git checkout trim
```

Implementing BWA-MEM

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

The relevant snippet from the bash pipeline is as follows

```

1 bwa mem -t 1 \
2   references/reference.fasta.gz \
3   qc_reads/${SAMPLE}_R1.fastq.gz qc_reads/${SAMPLE}_R2.fastq.gz \
4   | samtools view -b \
5   > mapped/${SAMPLE}.bam
```

Here are some questions to get you thinking



The mapping process requires an index comprised of several files but takes their common prefix (basename) as argument. How can you pass that prefix to the command?

You *could* hard-code it (not a good practice) or parse it out from the index file names.

Can you do it without hard-coding it or having to parse it out from the input file name(s)?

Our indexing process, in addition to the index files also outputs the prefix used when building the index.

How do you make sure that each instance of the process - and not just the first one - gets the index files in (*cough*) combination with a pair of FASTQ files?

Use the *combine* operator to combine the `indexChannel` with the `trimmedReadsChannel`.

```
1 process bwa_mem {
2   input:
3     set val(prefix), file(index), val(sample), file(reads) from \
4       indexChannel.combine(trimmedReadsChannel)
5
6   output:
7     file '*.bam'
8
9   script:
10    """
11    bwa mem -t ${task.cpus} ${ref} ${reads} \
12    | samtools view -b > ${sample}.bam
13    """
14 }
```

If you would like to compare your answer with the “official” one, ask your facilitator for instructions.

```
1 mv main.nf mymap.nf
2 git checkout map
```

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.

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` directive⁸ 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`⁹ to apply the directive to one or more relevant process.

Bonus tasks

If you exhaust the core contents of the workshop you could proceed with solving the following bonus exercises.



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.

```
1 # Find the appropriate module name
2 module av -l 2>&1 | grep conda
3
4 # Load the module
5 module load miniconda3-4.6.14-gcc-5.4.0-kkzv7zk
6
7 # Run with conda
8 nextflow run main.nf -profile conda,slurm --take -1 -resume
```

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

We have switched from singularity to conda so the software environment has changed.



Edit `main.nf`. Your task is to add `merge_bams` process which will merge the bam files produced by the `bwa_mem` process. The relevant `samtools` command should follow this pattern.

```
1 samtools merge outfile.bam infile1.bam infile2.bam infile3.bam
```

⁸<https://www.nextflow.io/docs/latest/process.html#publishdir>

⁹<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? Using the `.collect()` operator. Demonstrate your process definition to your facilitator.

```
1 process merge_bams {
2     input:
3         file('*.bam') from alignedReadsChannel.collect()
4
5     script:
6         """
7     samtools merge ${params.n}_samples_megred.bam *.bam
8     """
9 }
```

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

```
1     publishDir 'results/aligned', mode: 'copy'
2
3     output:
4         file '*.bam' into alignedReadsChannel
```

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.

```
1 process merge_bams {
2     cpus 2
3
4     input:
5         file('*.bam') from alignedReadsChannel.collect()
6
7     script:
8         """
9     samtools merge --threads ${task.cpus} \
10         ${params.n}_samples_megred.bam *.bam
11     """
12 }
```



Add `multiqc` process to summarize all FastQC outputs in a single report. Make sure the outputs are easily findable e.g. in `./results/multiqc`.

```
1 process multiqc {
2     publishDir 'results/multiqc', mode: 'copy'
3 }
```

```
4   input:
5     file '*' from fastQCdChannel.collect()
6
7   output:
8     file '*'
9
10  script:
11  """
12  multiqc .
13  """
14 }
```

Troubleshooting

Disconnected from the cluster?

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
```

Fix Nextflow version to the one used throughout the workshop and switch ansi logging off.



```
1 export NXF_VER=19.04.0
2 export NXF_ANSI_LOG=false
```

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
```

Space for Personal Notes or Feedback

[illegible]

[illegible]

[illegible]

[illegible]