

TRAINER'S MANUAL

Implementing Scalable Bioinformatic Workflows in Snakemake

Nathan S. Watson-Haigh

Sydney, Australia

12 December 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	5
The Trainers	6
Providing Feedback	7
Document Structure	7
Provided Compute Infrastructure	8
Connecting to the Cluster	9
Monitoring Slurm Jobs	9
Introduction to Snakemake	11
Key Learning Outcomes	12
Resources Required	12
Tools Used	12
Useful Links	12
Snakemake is Like Making Sunday Dinner	13
Setting Up Your Environment	13
Installing Snakemake	13
Your First Snakemake Workflow	15
Generalising Rules with Wildcards	16
Submitting Jobs to Slurm	17
Reimplementing a Workflow in Snakemake	18
Getting the Code	18
Getting the Data	19
Implementation of BWA Indexing	19
Implementing FastQC	22
Implementing Trimmomatic	26
Implementing BWA-MEM	27
Adding New Samples	32
Large Workflows	34
Troubleshooting	36
Getting Going After a Disconnect	36
Space for Personal Notes or Feedback	37

Workshop Information

The Trainers



Dr. Nathan S. Watson-Haigh

Senior Bioinformatician

Bioinformatics Hub, University of Adelaide

nathan.watson-haigh@adelaide.edu.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="5.8.1"
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. We have already run this for you and doing so again will cause delays.

```
1 ## Install snakemake using conda
2 ## This might take 5-10mins
3 #SNAKEMAKE_VERSION="5.8.1"
4 #
5 #conda create \
6 # --name snakemake \
7 # snakemake=${SNAKEMAKE_VERSION}
```

Snakemake installation is now complete.



All you need to do is to activate the environment which will make **snakemake** available to you 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
```



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:

Like Python, indentation in a **Snakefile** matters! For consistency, I'd recommend tab indentations.

```
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
```



There are examples of community developed profiles, available for other compute back-ends, see: <https://github.com/Snakemake-Profiles/doc>.

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
- 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 argumanet`)

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. Let's 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
```



If you want to view the **dag1.pdf** file yourself, you will need to copy the file to your local computer using **scp**, Filezilla (<https://filezilla-project.org>), WinSCP (<https://winscp.net>) or similar. The DAG generated should look like that shown in Figure 1.

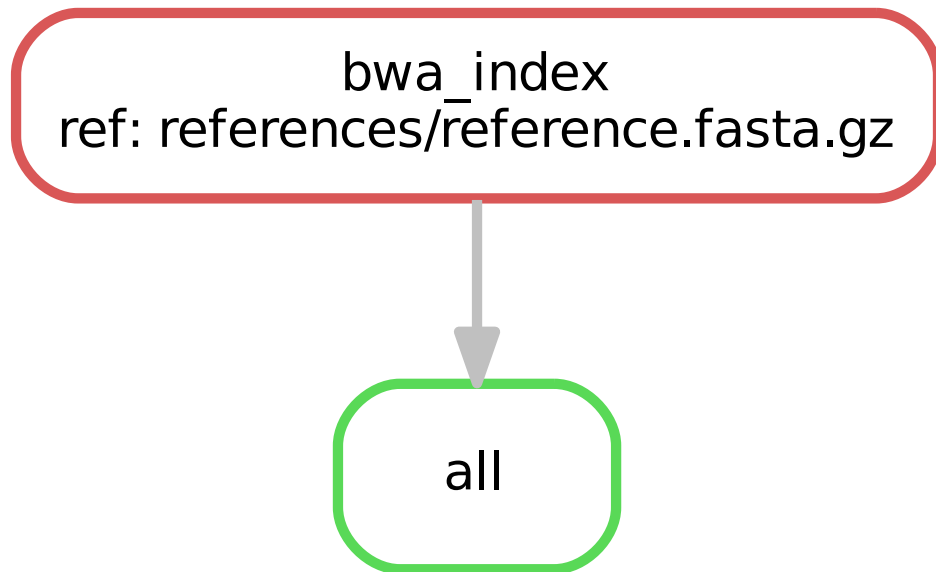


Figure 1: DAG of jobs showing `bwa_index` job dependant on the `all` pseudo-rule.

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
```



If you want to view the `dag2.pdf` file yourself, you will need to copy the file to your local computer using `scp`, Filezilla (<https://filezilla-project.org>), WinSCP (<https://winscp.net>) or similar. The DAG generated should look like that shown in Figure 2.

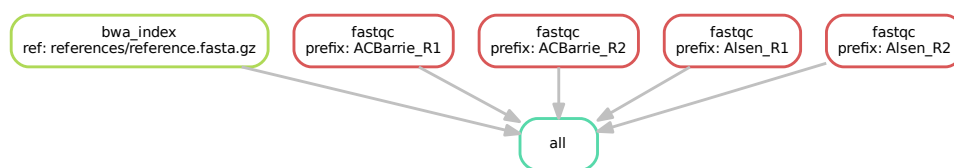


Figure 2: DAG of jobs showing `bwa_index` and several `fastqc` job dependant on the `all` pseudo-rule.

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.

Move 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 all the samples in the `SAMPLES` list, that you want the BWA-MEM rule to produce, to the `all` pseudo-rule.

```

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  ),
14  expand("mapped/{SAMPLE}.bam",
15  SAMPLE=SAMPLES,
16  ),
17  input:
18    reference = expand("references/reference.fasta.gz.{ext}",
19    ext=["amb", "ann", "bwt", "pac", "sa"]
20    ),
21    ...
22  params:
23    prefix = lambda wildcards, input: input["reference"][0][:4],
24    shell:
25      """
26      bwa mem -t {threads} \
27        {params.prefix} \
28        {input.r1} {input.r2} \
29        | samtools view -b \
30        > {output}
31      """

```

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

```



If you want to view the `dag3.pdf` file yourself, you will need to copy the file to your local computer using `scp`, Filezilla (<https://filezilla-project.org>), WinSCP (<https://winscp.net>) or similar. The DAG generated should look like that shown in Figure 3.

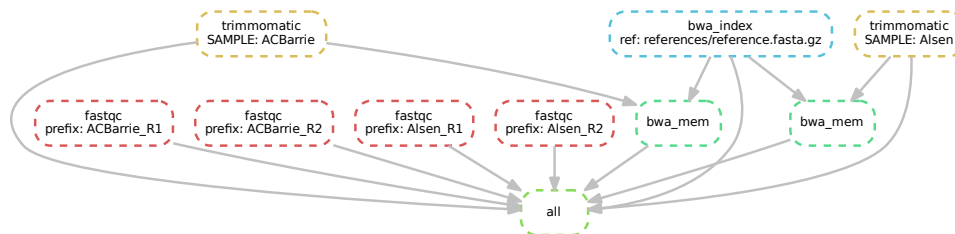


Figure 3: DAG of jobs showing the dependencies which exist in our final implementation of the `analysis.sh` workflow.

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

```



If you want to view the `dag4.pdf` file yourself, you will need to copy the file to your local computer using `scp`, Filezilla (<https://filezilla-project.org>), WinSCP (<https://winscp.net>) or similar. The DAG generated should look like that shown in Figure 4.



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
```



If you want to view the `rulegraph.pdf` file yourself, you will need to copy the file to your local computer using `scp`, Filezilla (<https://filezilla-project.org>), WinSCP (<https://winscp.net>) or similar. The rulegraph generated should look like that shown in Figure 5.

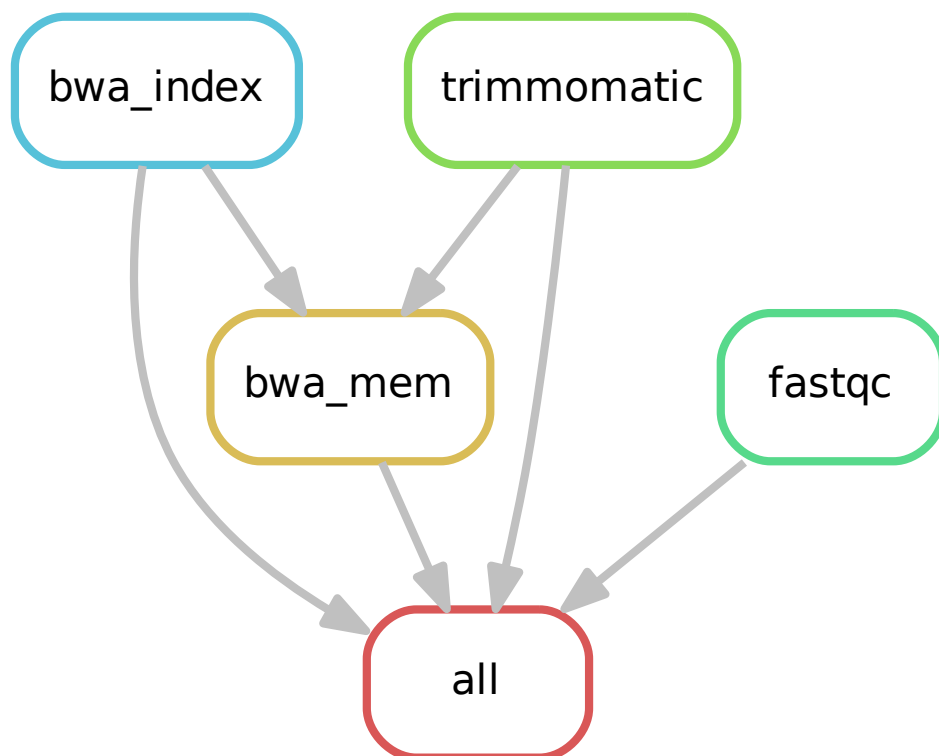


Figure 5: Rulegraph for the whole workflow.

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 processed 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:

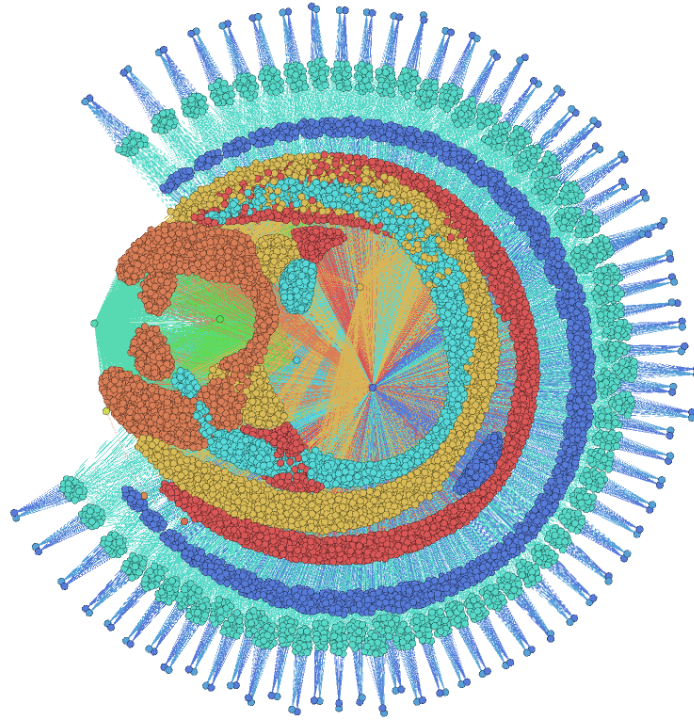


Figure 6: Large DAG, from the DAWN workflow, consisting of 10,587 nodes/jobs and 32,353 edges/dependencies. Layed out using ForceAtlas 2 algorithm of Gephi

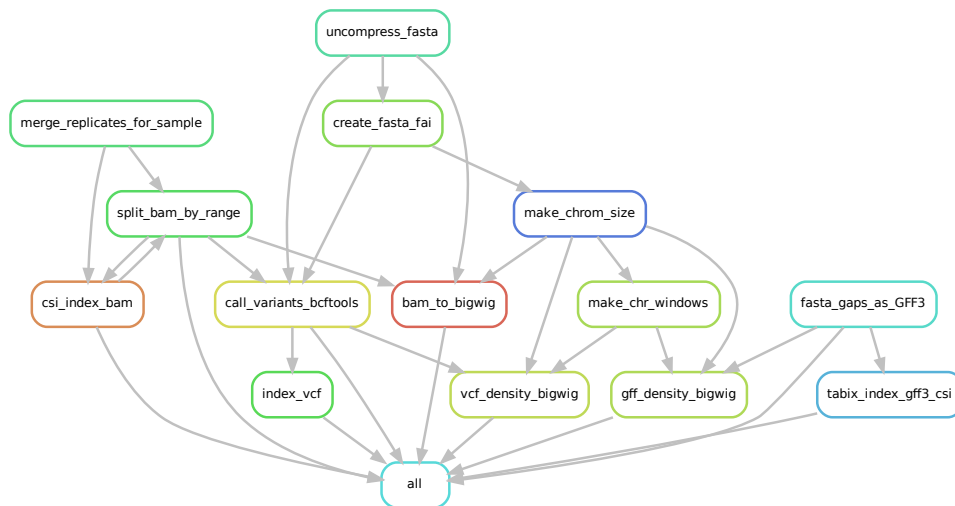


Figure 7: Rulegraph, from the DAWN workflow.

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
```

Space for Personal Notes or Feedback

[illegible]

[illegible]

[illegible]

[illegible]