

RNAseq reads to differential genes and pathways

27 & 28 September 2022

Archive of questions from the workshop Slack Channel

(In the order they were asked during the workshop)

Questions about the pipeline and RNAseq

Question/comment	Answer
Can you explain allele specific expression more?	This is about differences in expression between alleles - e.g. expression of a gene inherited from one parent compared to the other. This analysis pipeline is different to the one we will go through today
Are these column headings for the sample sheet are must-have?	Yes, there is a naming convention that the programs used for RNAseq data follow.
So is it convention that fastq_1 always the forward?	Yes. They can also be L1 and R1
If you have your own reference genome and annotate it using say example Augustus, then how do you use it in the nextflow?	If you have your own reference genome and annotate it using say example Augustus, then how do you use it in the nextflow. in the nextflow command, you can change the parameters to use your own reference and annotation files with the --fasta and --gtf options
Does Nextflow have pipelines for RNA-seq analysis of prokaryotes?	you can browse nf-core pipelines, which can be found here: https://nf-co.re/ There are currently 69 pipelines available here
Is there a list of reference genomes that are available onNextflow?	It depends on the specific workflow you're using. If you go to https://nf-co.re/ you can see all the available pipelines and read more about each pipeline Also, nf-core give you access to Illumina's

	<p>iGenomes database https://nf-co.re/usage/reference_genomes which has pre-stored and indexed reference genomes for a few species.</p>
<p>in regards to the <code>-profile singularity</code> option, is singularity a tool we also need to install?</p>	<p>Singularity is a software that runs containers. Nextflow handles the installation for you</p>
<p>My model organism is a cow and i dont see a cow reference genome in the reference genome database. So all i have to do is download the genome sequence and the annotation files from RefSeq in the working directory?</p>	<p>Yes, that's right.</p>
<p>In terms of customizing the workflow, if I like to adjust certain flags for some of the steps, for example for STAR, adjusting the read length taken into the pipeline or adjust the strand specific flags, how can I adjust specific parameters please?</p>	<p>The link to nf-co.re/rnaseq parameters: https://nf-co.re/rnaseq/3.8.1/parameters</p>
<p>Day 2</p>	
<p>Can we do parallelization (for example alignment of all samples at the same time) inside nextflow?</p>	<p>Yes, if you have sufficient compute resources nextflow can run many jobs at once</p> <p>We did that yesterday. Download the <code>timeline.html</code> report from your run yesterday to see which processes were able to run in parallel. Each process run by nf-core has specific cpu and memory requirements and your VM has a set number of CPUs and memory. Nextflow will run as many processes as it can in parallel given the resources available to it on the machine.</p>
<p>Can we use normalized/standardized gene expression data for DE analysis?</p>	<p>Depends on the tool you use. Today we will use DESeq2 which requires raw counts. This is because DESeq2 does normalization internally (and has it's own unique method to do so).</p>
<p>In the nextflow of nf-core, what is the difference between the fasta and Star-index file? My understanding is star-index are used for STAR alignment and quantification, while fasta is used for Salmon to generate</p>	<p>Index files are used by tools (like STAR, salmon) to perform look-up functions faster (i.e. when it is aligning to a reference genome, it can do this faster). The index method/files are tool specific</p>



<p>salmon index file for counts and possible for bam etc.? Is that right?</p>	<p>Sure, but the index file are normally come from the reference genome, right? I guess from metagenomics, we use the assembly for both index and fasta, maybe different in human genome, as the index can be very specific sets of genes only.</p> <p>That's right - the index file is specific to the reference genome and tool (and even tool version). You should generate index files with your FASTA reference file & tools that you use</p>
<p>Is there a list of R studio containers that are ready to be used through nimbus please?</p>	<p>RStudio is installed already on Nimbus if you set up your instance using the bio image. See this documentation for more info on how to open it: https://support.pawsey.org.au/documentation/display/US/Nimbus+for+Bioinformatics#NimbusforBioinformatics-RStudio</p> <p>See also: https://support.pawsey.org.au/documentation/display/US/Run+RStudio+Interactively</p> <p>We're working with a container we developed specifically for today's training, you can access it at the same cvmfs path if you'd like to use it for yourself and we plan to release it publicly soon.</p>
<p>Why do we use Rstudio to do DE rather than a nextflow pipeline?</p>	<p>DE analysis (rather than data processing as we did on day 1) involves a lot more interactive analysis, looking at plots, etc. So, it is a nicer way to perform this analysis.</p> <p>I would add, yesterday, we were processing our data from raw sequence reads to generate a dataset that could be analysed. Its ok to automate processing steps which are generally more standardised across experiments compared with analysis that changes a lot depending on your research questions.</p>
<p>Why does DESeq2 uses a negative binomial distribution model?</p>	<p>Take a look at their documentation, go down to Theory behind DESeq2 in http://bioconductor.org/packages/release/bioc/vignettes/DESeq2/inst/doc/DESeq2.html</p>



Australian
BioCommons

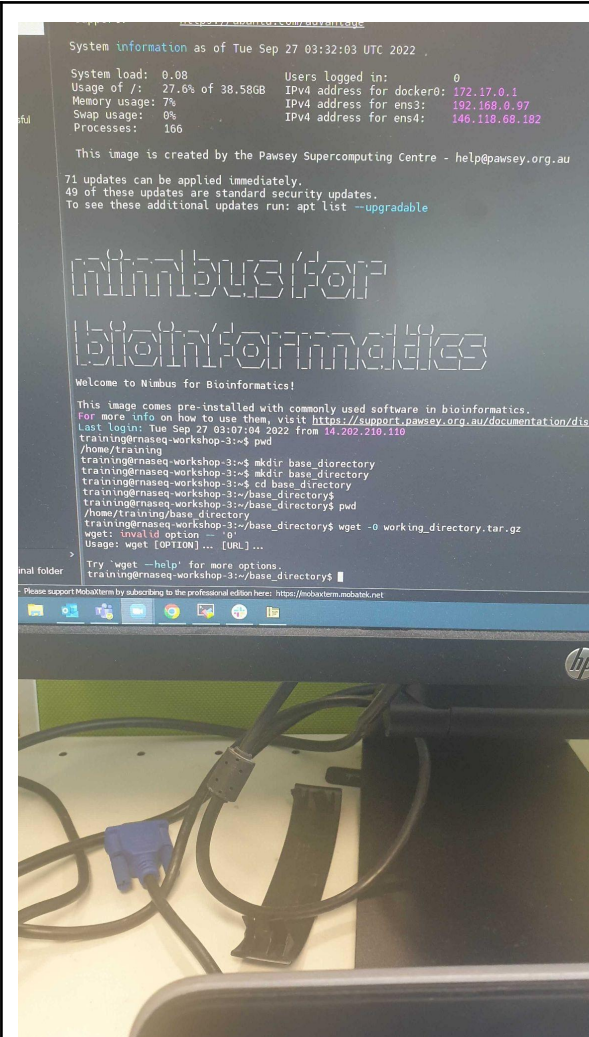


THE UNIVERSITY OF
SYDNEY

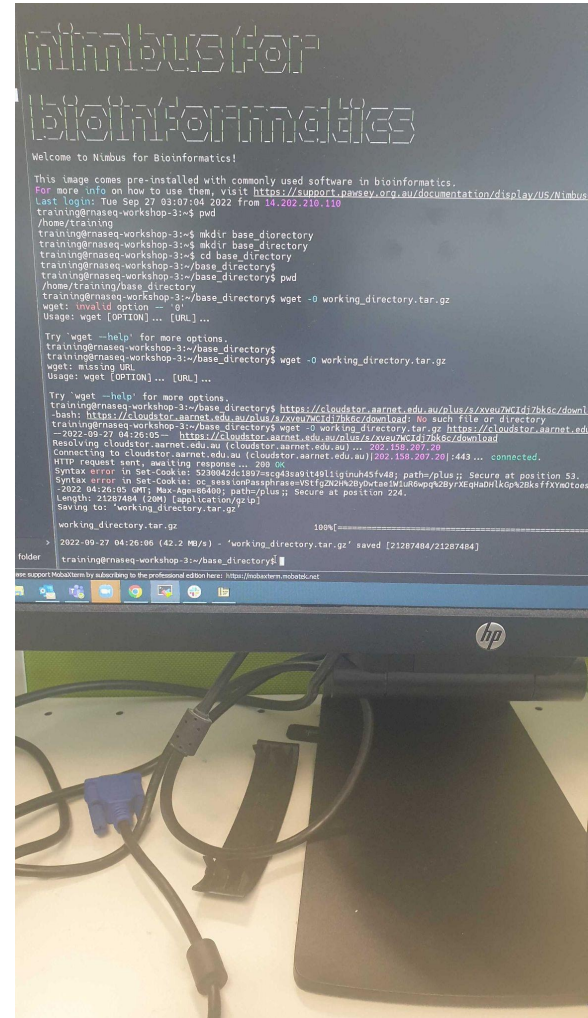
	<p>Here is an example of a weird dispersion plot : https://support.bioconductor.org/p/107937/ And a discussion as to why. Michael Love is one of the DESeq2 authors and is actively on the bioconductor forum if you want to ask the authors a question. Otherwise do ask a local statistician/bioinformatician</p>
Isn't p.adjust the same as qvalue?	<p>Yep, the qvalue package doing slightly different</p> <p>clusterProfiler is just offering you a bit of choice in the types of padj methods you can apply, whereas the qvalue applied here is just the false discovery rate</p>
Why do we conduct DE for up- and down-regulatory genes separately?	<p>It's much easier to interpret this way. We can identify which group of genes/enrichments are upregulated together in the knockout as compared to the wild, and vice versa</p>

Error messages

Question/comment	Answer
When connecting to Nimbus my connection timed out	Try connecting to a different internet network (e.g. Eduroam). Some institutional internet connections block access to Nimbus.



Ah ok, its meant to be -O (the letter O) not the number 0



Looks like the file downloaded successfully. You can check by typing the command `ls`

After `ls` do `tar -zxvf working_directory.tar.gz`

then `cd working_directory`

those commands 1. unzip the file you just downloaded 2. change your directory to `working_directory`



```

drwxr-xr-x  2 othman  staff   64B 27 Sep 14:12 base_directory
drwxr-xr-x 10 othman  staff  320B 31 Mar 10:20 lesson2
drwxr-xr-x  3 othman  staff   96B 18 Feb 2022 opt
(base) othman@i184-12-208 ~ % cd cd base_directory
cd: string not in pwd: cd
(base) othman@i184-12-208 ~ % cd base_directory
(base) othman@i184-12-208 base_directory % pwd
/Users/othman/base_directory
(base) othman@i184-12-208 base_directory % wget -O working_directory.tar.gz https://cloudstor.aarnet.edu.au/plus/s/xveu7WCIdj7bk6c/download
zsh: command not found: wget
(base) othman@i184-12-208 base_directory % wget -O working_directory.tar.gz https://cloudstor.aarnet.edu.au/plus/s/xveu7WCIdj7bk6c/download
zsh: command not found: wget
(base) othman@i184-12-208 base_directory % wget -O working_directory.tar.gz https://cloudstor.aarnet.edu.au/plus/s/xveu7WCIdj7bk6c/download
zsh: command not found: wget
(base) othman@i184-12-208 base_directory % wget -O working_directory.tar.gz https://cloudstor.aarnet.edu.au/plus/s/xveu7WCIdj7bk6c/download
zsh: command not found: wget
(base) othman@i184-12-208 base_directory % wget -O working_directory.tar.gz https://cloudstor.aarnet.edu.au/plus/s/xveu7WCIdj7bk6c/download
zsh: command not found: wget
(base) othman@i184-12-208 base_directory % wget -O working_directory.tar.gz https://cloudstor.aarnet.edu.au/plus/s/xveu7WCIdj7bk6c/download
zsh: command not found: wget
(base) othman@i184-12-208 base_directory % wget -O working_directory.tar.gz https://cloudstor.aarnet.edu.au/plus/s/xveu7WCIdj7bk6c/download
zsh: command not found: wget
(base) othman@i184-12-208 base_directory %
/Users/othman/base_directory
(base) othman@i184-12-208 base_directory %

```

You aren't logged into your nimbus vm. it appears you're on your local computer.

Follow directions here to login
<https://sydney-informatics-hub.github.io/rna-seq-pt1-quarto/setup.html#connect-to-nimbus>

I get command not found

```

/home/mobaxterm/base_directory/working_directory
27/09/2022 14:32:39 nextflow run $cvmfs_path/nfcore_pipeline/rnaseq/ \
--input samplesheet.csv \
--profile singularity \
--fasta $cvmfs_path/Mouse_chr18_reference/chr18.fa \
--gtf $cvmfs_path/Mouse_chr18_reference/chr18_startOfLine.gtf \
--star_index $cvmfs_path/Mouse_chr18_reference/chr18_STAR_singularity_index/ \
--max_memory '6 GB' --max_cpus 2 \
--outdir results \
--with-report execution_report.html \
--with-timeline timeline_report.html \
--with-dag flowchart.png
nextflow: command not found
/home/mobaxterm/base_directory/working_directory
27/09/2022 14:32:45

```

Looks like you are not logged onto Nimbus.

follow directions here to login
<https://sydney-informatics-hub.github.io/rna-seq-pt1-quarto/setup.html#connect-to-nimbus>

my command failed: [bd/239935] NOTE: Process `NFCORE_RNASEQ:RNASEQ:ALIGN_STAR:STAR_ALIGN (SRR3473985)` terminated with an error exit status (137) -- Execution is retried (1)

terminal@training: ~

```

* Documentation: https://help.ubuntu.com
* Management: https://landscape.canonical.com
* Support: https://ubuntu.com/advantage

System information as of Tue Sep 27 04:38:02 UTC 2022
System load: 0.07          Users logged in: 0
Usage of /: 27.0% of 30.86GB  IPv4 address for dockero: 172.17.0.1
Memory usage: 7%          IPv4 address for ens3: 192.168.0.167
Swap usage: 0%           IPv4 address for ens4: 186.138.07.192
Processes: 163

This image is created by the Pausey Supercomputing Centre - help@pausey.org.au

71 updates can be applied immediately.
49 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable

nimbus for bioinformatics
Welcome to Nimbus for Bioinformatics!

This image comes pre-installed with commonly used software in bioinformatics.
For more info on how to use them, visit https://support.pausey.org.au/documentation/playing-with-nimbus-for-bioinformatics
Last login: Tue Sep 27 04:35:44 2022 from 131.101.69.4
training@rnaseq-workshop-22:~$

```

You are logged in now, just make sure to keep working in this terminal window.

you can try to rerun your pipeline by adding -resume to the end of the command. Like this:

```

nextflow run
$cvmfs_path/nfcore_pipeline/rnaseq/ \
--input samplesheet.csv \
--profile singularity \
--fasta

```




It seems like the command from the previous step didn't run for me. I'm stuck with this screen

```
training@rnaseq-workshop-45:~/base_directory/working_directory$ ls -lh results
ls: cannot access 'results': No such file or directory
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ cd results
-bash: cd: results: No such file or directory
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ ^C
training@rnaseq-workshop-45:~/base_directory/working_directory/results$
```

try `ctrl + c` and then try to rerun the command

```
training@rnaseq-workshop-45:~/base_directory/working_directory$ cd results
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ ls -lh results
ls: cannot access 'results': No such file or directory
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ cd results
-bash: cd: results: No such file or directory
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ ^C
training@rnaseq-workshop-45:~/base_directory/working_directory/results$
```

you're already in the results directory

how do i get ls-lh to work?

it's `ls -lh` with a space in the middle

Has anyone have permission issue to transfer data?

```
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ cd results
```

you've only got `train@` you need to have `training@`

I would recommend only downloading the fastqc and trimalore folders, There are some big alignment files in there

```
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ cd results
-bash: cd: results: No such file or directory
training@rnaseq-workshop-45:~/base_directory/working_directory/results$ ^C
training@rnaseq-workshop-45:~/base_directory/working_directory/results$
```

You are missing a space just before `./RNASeq_workshop/`

there is no "." in fastqc



<p>Day 2</p>	
<p>I don't have the warning about the package on my RStudio. Is this something to worry about?</p>	
<p>I'm not getting a new "ready" line after running the password command</p>	<p>that means it's working! now go to chrome browser</p> <p>then Type 146.118.XX.XX:8787 in your browser where the XX.XX will be replaced by your IP specific digits:</p>
<p>I don't get the warning when I run the 'DESeqDataSetFromMatrix' chunk, will assume this is fine but unsure why it's different from Nandan's</p>	<p>Neither do I, but I can run all the code fine so we can probably disregard it.</p>
<p>Just getting an error on the volcano plot</p> <pre> 452 # Add significance lines at log2FoldChange -1, 1 and pvalue 0.05 453 p2 <- p + geom_vline(xintercept=c(-1, 1), col = "red") + 454 geom_hline(yintercept=-log10(0.05), col = "red") 455 456 # Print the plot 457 p2 458 - </pre> 	<p>I think there's a typo somewhere, can you compare/use this code:</p> <pre> # Create a basic volcano plot (scatter plot) with x axis = LogFC, y axis = # -log10(pvalue) resdata <- as.data.frame(res) # Define whether genes are significantly DE or not and store this # in a new column called DE resdata\$Significant <- "No" resdata\$Significant[resdata\$log2FoldChange > 1 & resdata\$pvalue < 0.05] <- "Upregulated" resdata\$Significant[resdata\$log2FoldChange < -1 & resdata\$pvalue < 0.05] <- "Downregulated" # Create the volcano plot p <- ggplot(data=resdata, aes(x=log2FoldChange, y=-log10(pvalue), col=Significant)) + geom_point() # Add significance lines at # log2FoldChange -1, 1 and pvalue 0.05 p2 <- p + geom_vline(xintercept=c(-1, 1), col = "red") + geom_hline(yintercept=-log10(0.05), col = "red") # Print the plot p2 </pre>
<p>My knitting error is as follows: No protocol specified</p>	<p>looks like you have a View() command in one of your code blocks.</p>



Australian
BioCommons



THE UNIVERSITY OF
SYDNEY

```
Quitting from lines 568-603
(rnaseq_DE_analysis_Day2.Rmd)
Error in .External2(C_dataviewer, x,
title) : unable to start data viewer
Calls: <Anonymous> ...
withCallingHandlers -> withVisible ->
eval -> eval -> View
In addition: Warning messages:
1: In DESeqDataSet(se, design = design,
ignoreRank) :
  some variables in design formula are
characters, converting to factors
2: Removed 8 rows containing missing
values (geom_point).
3: In View(sig.up) : unable to open
display
Execution halted
```

Try hash it out or remove it and then knit again.