

Bioinformatics for transcriptional and genome variation

Kevin Murray

@kdmurray91

kevin@kdmurray.id.au

Borevitz Lab, ANU

RNAseq bioinformatics session, 2015-04-29

- ▶ *Disclaimer: This is a whirlwind tour!*

- ▶ *Disclaimer: This is a whirlwind tour!*
- ▶ How we do RNAseq
- ▶ Experimental design: What we've done, how important it is
- ▶ RNAseq analysis pipelines
- ▶ DNA analysis pipelines

RNAseq in the capital

- ▶ We've done a fair bit (≈ 15 lanes, mostly Pogson)
- ▶ Largely small-scale, < 50 samples per experiment

RNAseq in the capital

- ▶ We've done a fair bit (≈ 15 lanes, mostly Pogson)
- ▶ Largely small-scale, < 50 samples per experiment
- ▶ Developed our own sequence analysis pipelines

RNAseq in the capital

- ▶ We've done a fair bit (≈ 15 lanes, mostly Pogson)
- ▶ Largely small-scale, < 50 samples per experiment
- ▶ Developed our own sequence analysis pipelines
- ▶ Wet lab: Mostly NEB/TruSeq kits; few attempts at custom library prep



Starting at the start. . .

► **Experimental design is key**

- Design of sampling: randomisation at each step
- Replication vs coverage trade-off
- Speed is essential during collection: responses can be fast

► **Experimental design is key**

- Design of sampling: randomisation at each step
- Replication vs coverage trade-off
- Speed is essential during collection: responses can be fast

► **Example experiment:**

- Col-0
- Three growth conditions
- Before and after 1000 μ E light treatment
- 3 biological reps each group
- Tissue harvested within 60 seconds of end of stress
- TruSeq RNASeq kits, 12 samples/lane
- 24 samples

Deeper not always better

- ▶ Given the same amount of sequencing, go deep or wide?



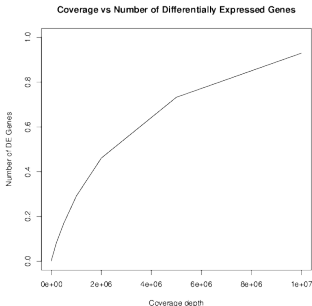
Deeper not always better

- ▶ Given the same amount of sequencing, go deep or wide?
- ▶ For our experiments, **wide, more reps, almost always better**



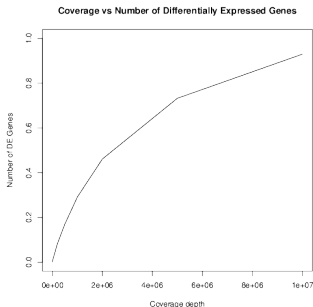
Deeper not always better

- ▶ Given the same amount of sequencing, go deep or wide?
- ▶ For our experiments, **wide, more reps, almost always better**



Deeper not always better

- ▶ Given the same amount of sequencing, go deep or wide?
- ▶ For our experiments, **wide, more reps, almost always better**



- ▶ **Your mileage may vary!**
- ▶ See Kliebenstein, (2012) FIPS: Exploring the Shallow End; Estimating Information Content in Transcriptomics Studies.



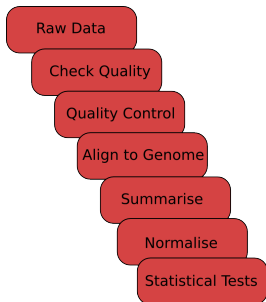
Sequence analysis

- ▶ Data is rawer now than with microarrays
 - ▶ Needs significant computational resources
 - ▶ At large scale, can be a bottleneck

Sequence analysis

- ▶ Data is rawer now than with microarrays
 - ▶ Needs significant computational resources
 - ▶ At large scale, can be a bottleneck
- ▶ Have developed pipelines to do this efficiently
- ▶ `https://github.com/kdmurray91/RNAseqPipeline`
- ▶ `https://github.com/pedrocrisp/NGS-pipelines`

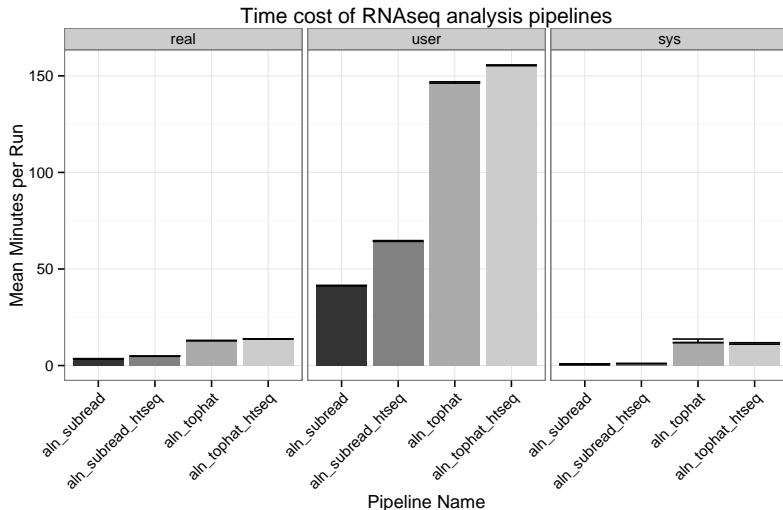
Sequence analysis pipeline



- ▶ fastqc
- ▶ scythe
- ▶ sickle
- ▶ subread/subjunc
- ▶ featurecounts
- ▶ edgeR
 - ▶ TMM normalisation
 - ▶ exactTest or glmFit
 - ▶ Also using limma's voom
- ▶ R scripts for post-analysis
 - ▶ G0seq
- ▶ Diagnostic plots **highly recommended!**

Pipeline performance

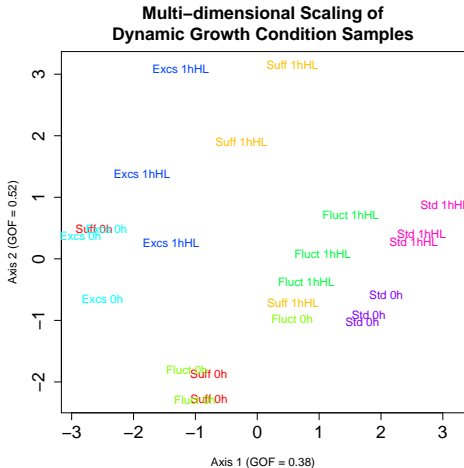
- Outperforms others by $> 2 - 3\times$





MDS Plots save time!

- If your reps don't cluster, time to cry into beer.





Existing DNA variation pipelines

- ▶ Genotyping-by-sequencing
 - ▶ Reference & *de-novo* analysis
 - ▶ Porting to NCI NF
 - ▶ Currently manual, working to automate
 - ▶ Processed > 5000 samples
- ▶ Reference-based genotype calling
 - ▶ Pipelines exist
 - ▶ Not used a lot, requires deeper coverage
 - ▶ See Norman's talk yesterday

Novel algorithms for DNA variation

- ▶ My PhD topic

Novel algorithms for DNA variaiton

- ▶ My PhD topic
- ▶ Our wish-list:
 - ▶ Reference & alignment free
 - ▶ Tolerates any sequencing kind
 - ▶ Works with “wide and shallow” experiments: e.g. 1000 samples at 1x

Novel algorithms for DNA variaiton

- ▶ My PhD topic
- ▶ Our wish-list:
 - ▶ Reference & alignment free
 - ▶ Tolerates any sequencing kind
 - ▶ Works with “wide and shallow” experiments: e.g. 1000 samples at 1x
- ▶ e.g: gearing up to sequence 7500+ Eucalyptus, generating over 10 TB **raw** sequence data.

- ▶ My PhD topic
- ▶ Our wish-list:
 - ▶ Reference & alignment free
 - ▶ Tolerates any sequencing kind
 - ▶ Works with “wide and shallow” experiments: e.g. 1000 samples at 1x
- ▶ e.g: gearing up to sequence 7500+ Eucalyptus, generating over 10 TB **raw** sequence data.
- ▶ *k*-mer analysis: analyse *k*-length words of sequence
 - ▶ Fast
 - ▶ Constant-memory (with `khmer`)
 - ▶ Scalable (linear time w/ number of samples)
 - ▶ Parallelisable (within & across nodes)

Thanks

- ▶ Borevitz lab (Norman, Justin, Megan, Steve)
- ▶ Pogson lab (Pete Crisp)
- ▶ Genome Discovery Unit
- ▶ Slides at git.io/vfAof

Grab-bag of capabilities

- ▶ Confirm genotype using RNAseq reads
- ▶ Check technical reps are true