

# NovaSeq 6000 vs HiSeq 2500

## Evaluation by mRNA-seq

Joel Parker

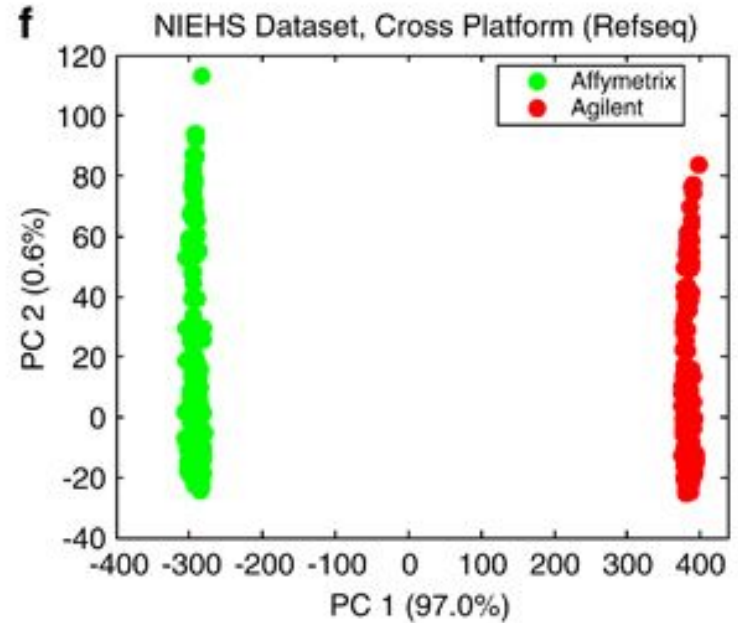
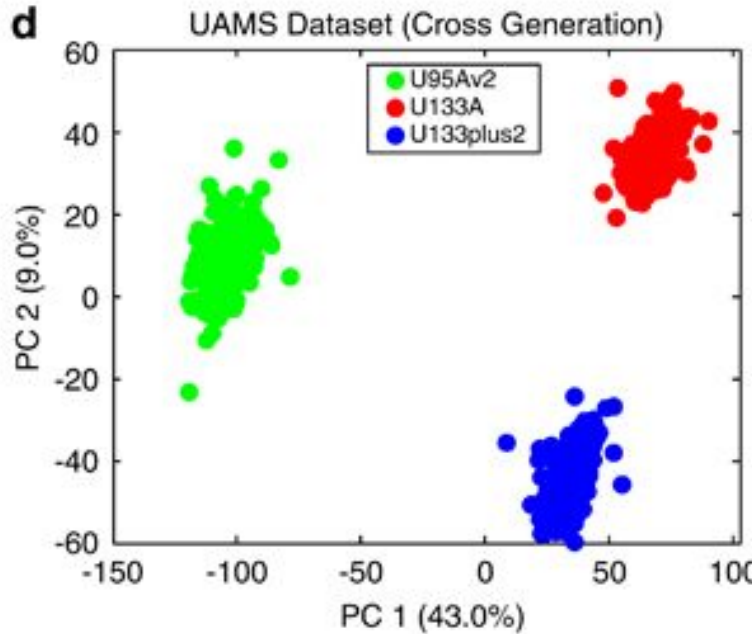
Lineberger Bioinformatics Core

<https://lbc.unc.edu/>



UNC  
LINEBERGER

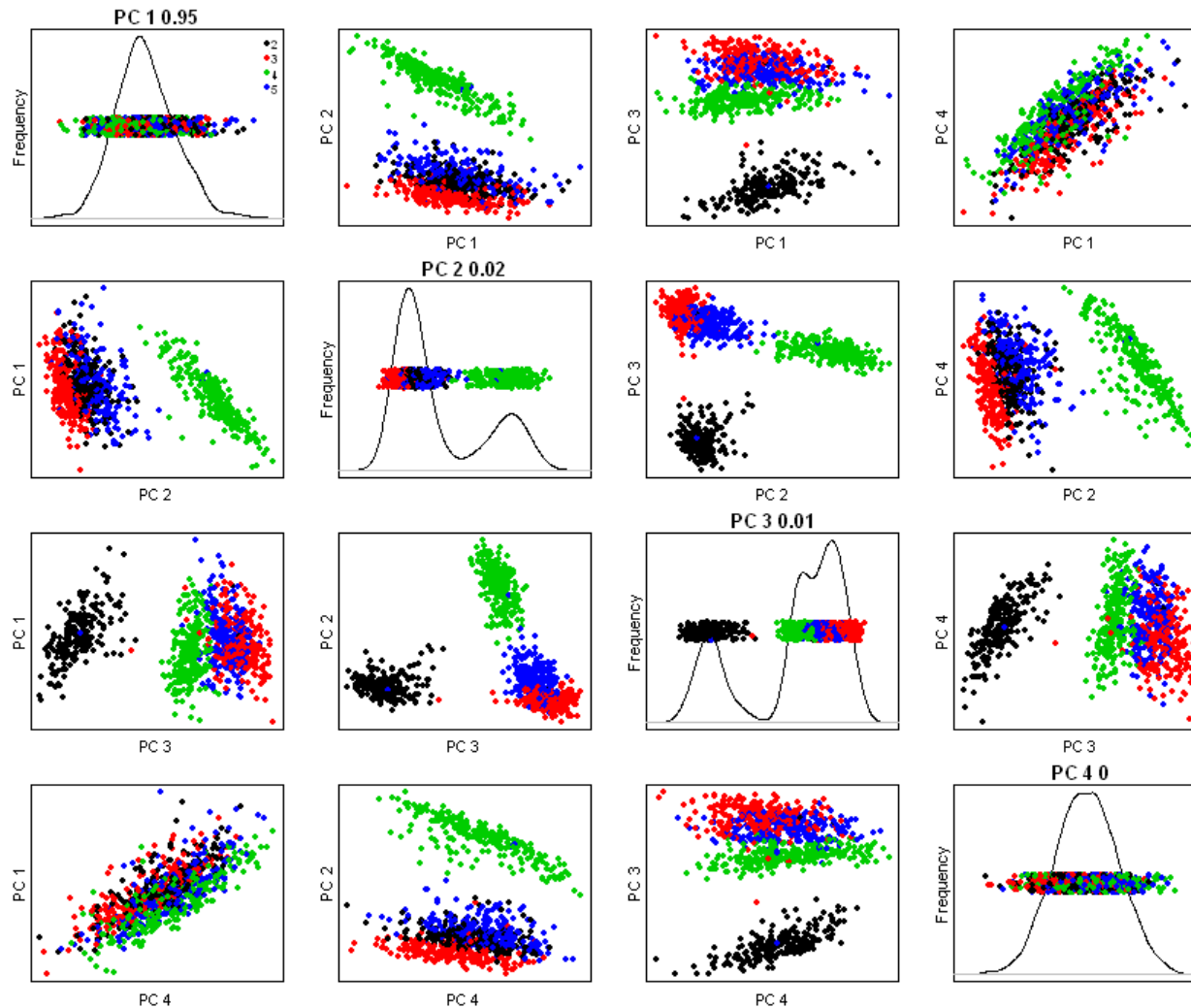
# Instrument Bias



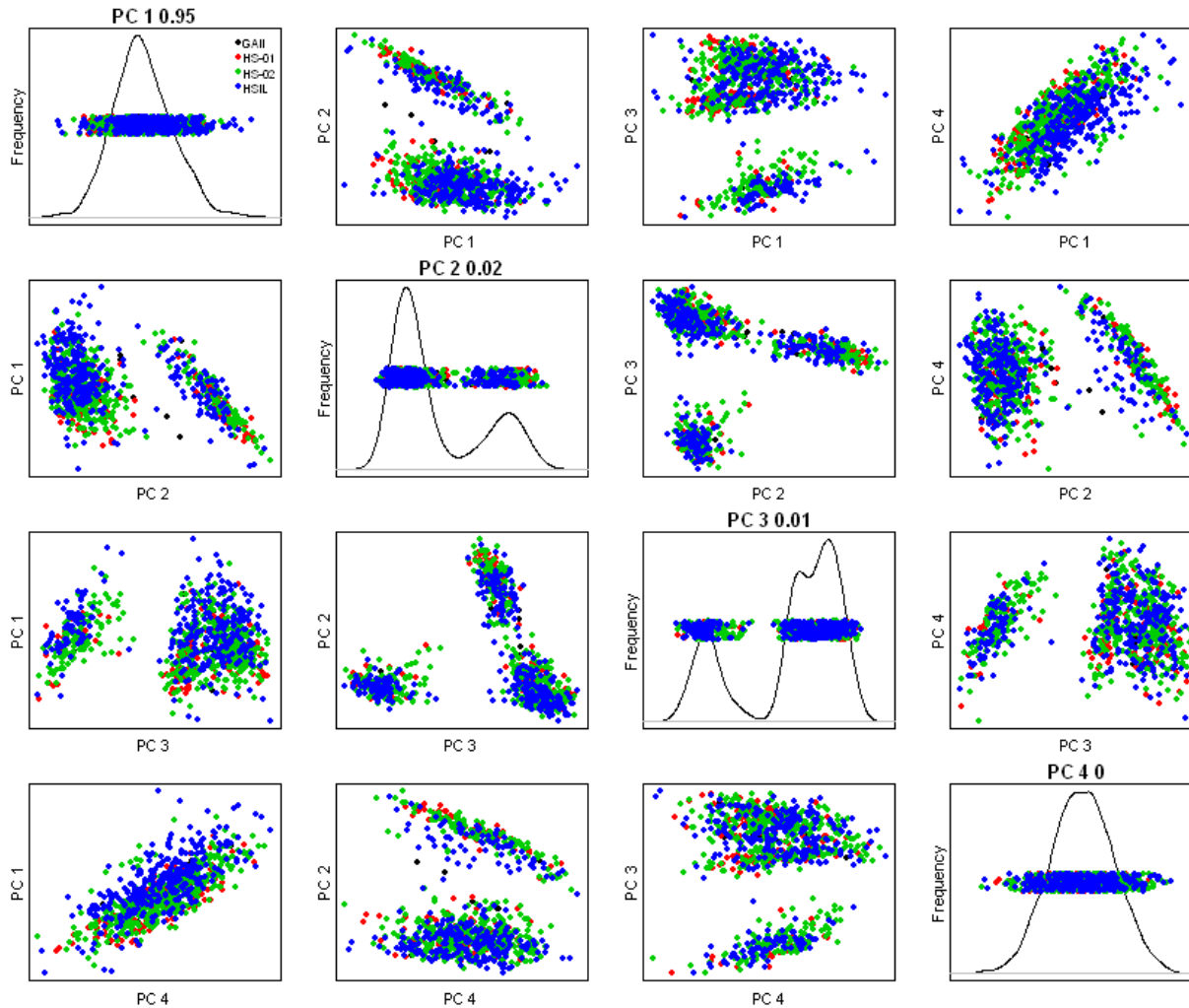
Platform	Samples in GEO
U133plus2	146,142
U133A	22,283
U95	6,446
Agilent 4x44k	14,367

Luo et al., MAQC-II 2010

# Instrument bias in mRNA-seq



# Instrument bias in mRNA-seq

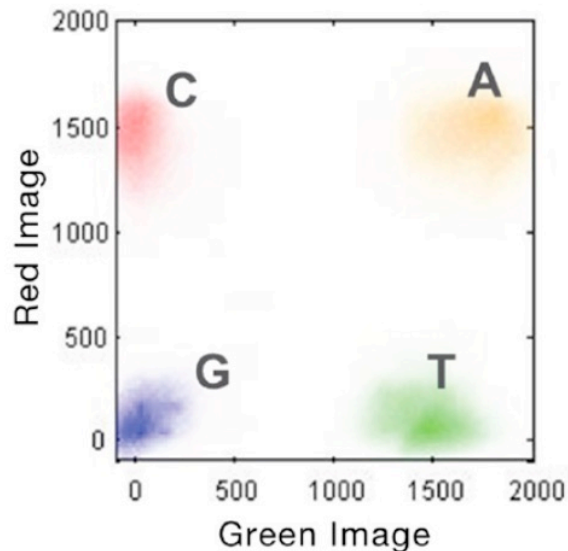


4 machines  
2 sites  
6 batches – 2 mos

**GAIIX**  
**HS-01**  
**HS-02**  
**HS-IL**

# Hiseq vs Novaseq

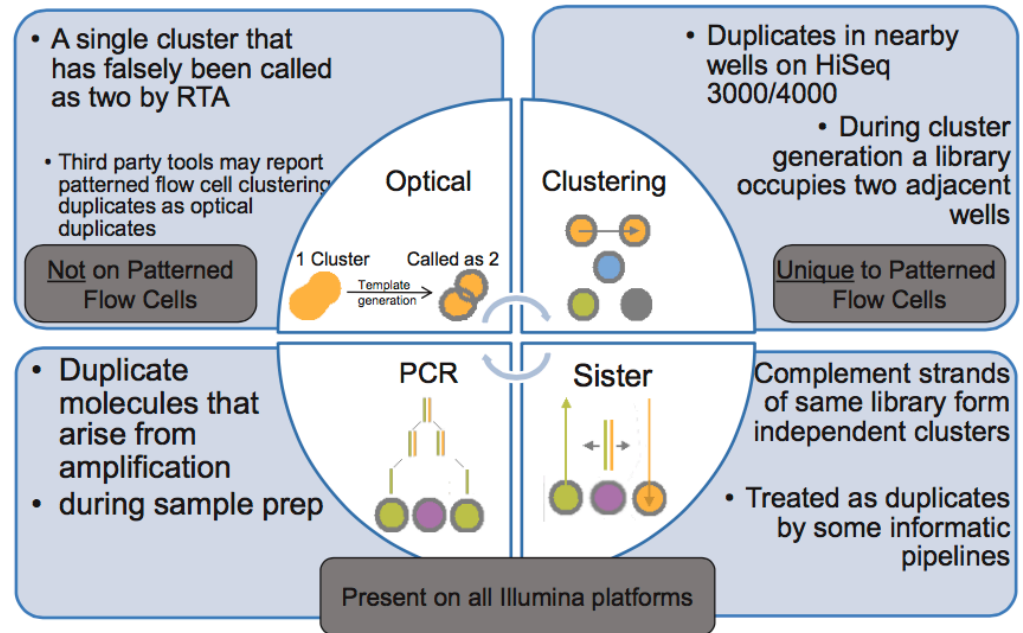
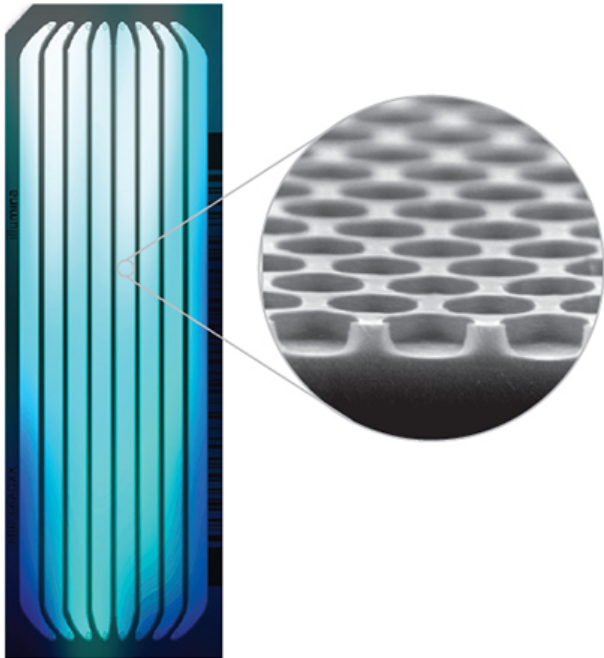
- Hiseq2500 uses 4 channels while Novaseq only uses 2
  - G is represented by the lack of signal, previously called N
  - Poor quality reads may show up as polyG



- Increased sensitivity to imbalanced indexes during calibration
- cutadapt permits removing trailing 'G's

# HiSeq vs Novaseq

- Ordered flowcells
  - Reports of barcoded reads spilling over into adjacent wells

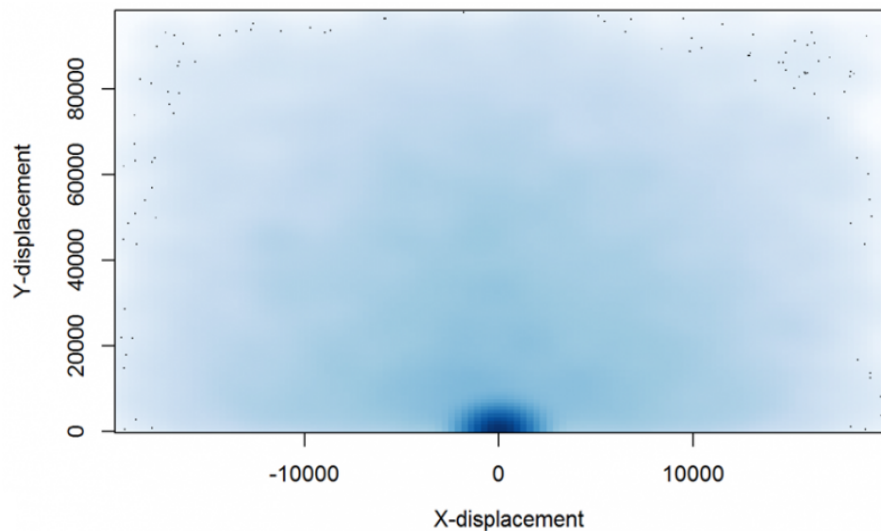


# Hiseq vs Novaseq

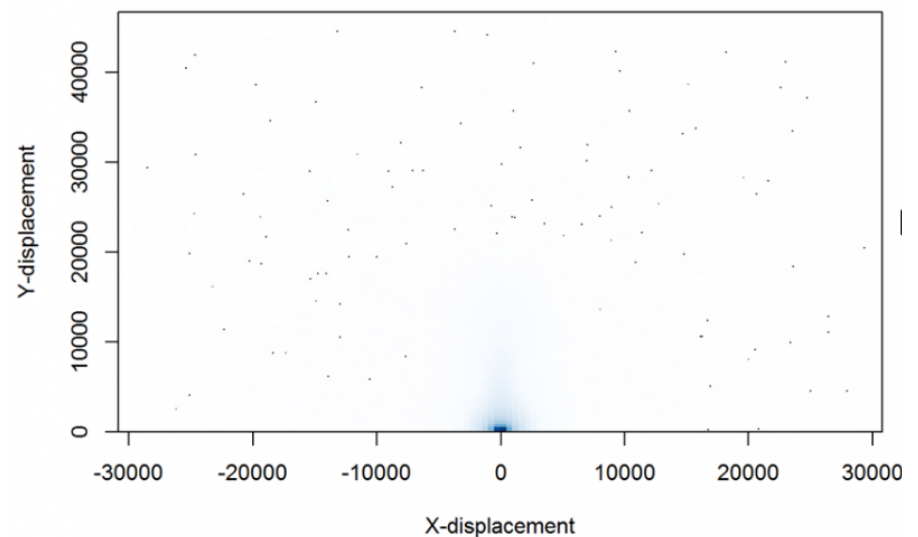
- Ordered flowcells
  - Reports of barcoded reads spilling over into adjacent wells

## Relative position of duplicates

Hiseq 2500

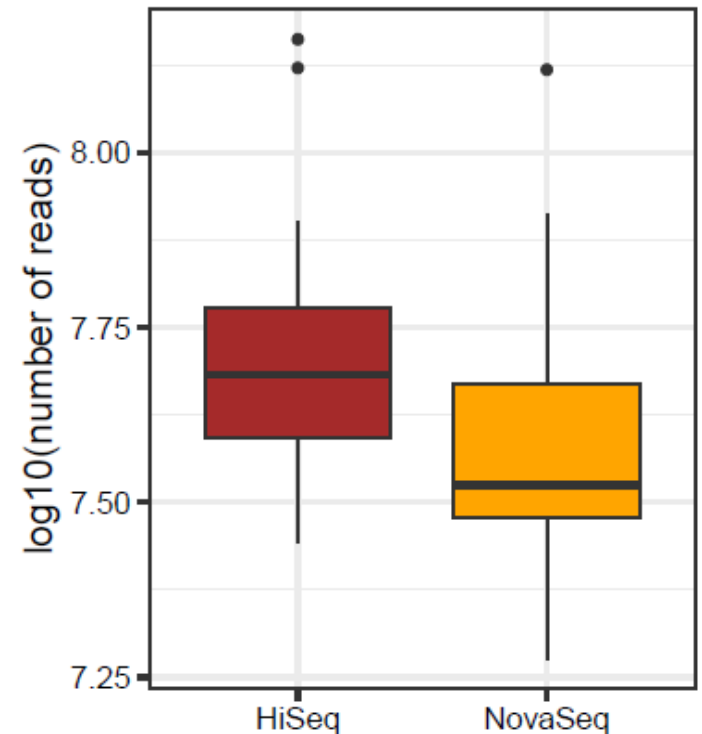


Novaseq



# Evaluation

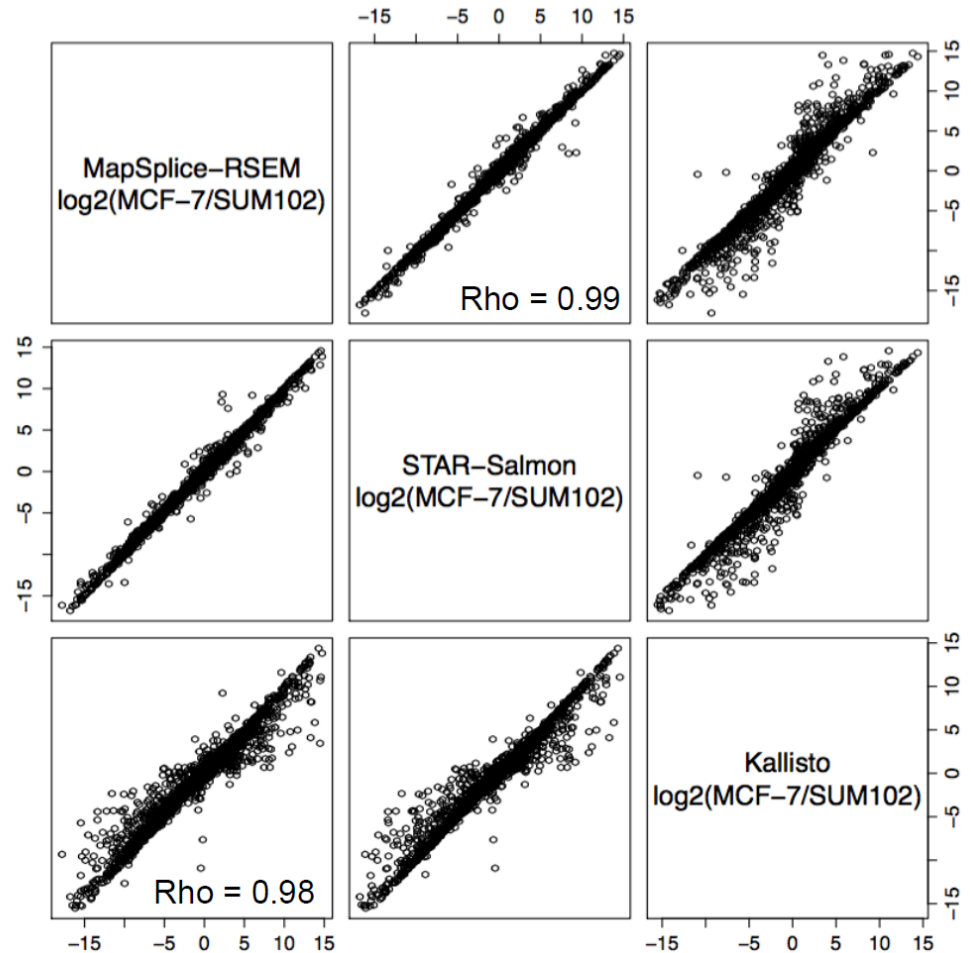
- RNA isolated from 48 GEMM derived tumors
- TruSeq stranded mRNA
- Single barcode
- 2 NovaSeq lanes (S1) with 24 samples / lane
- 12 HiSeq lanes with 4 samples / lane
- Identical library on both machines



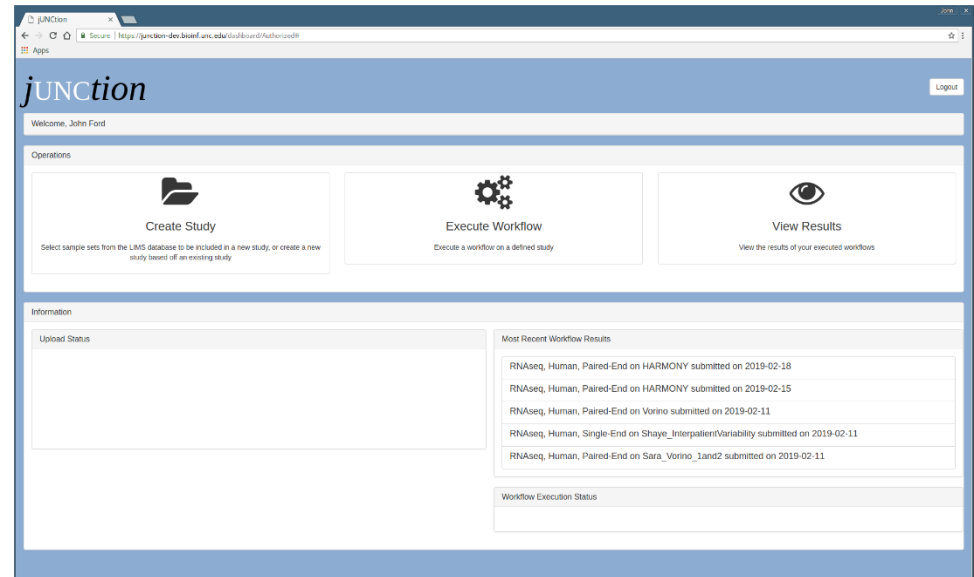
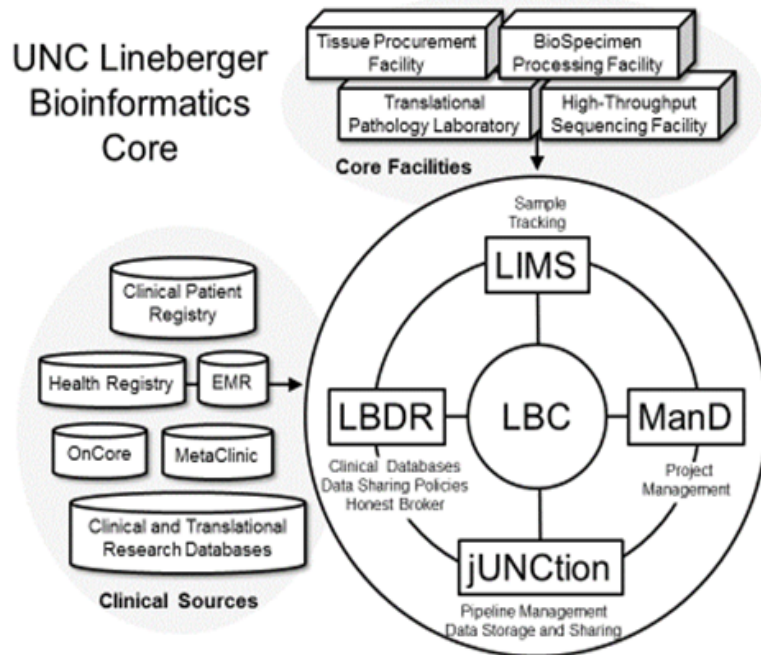


# RNA-seq Workflow

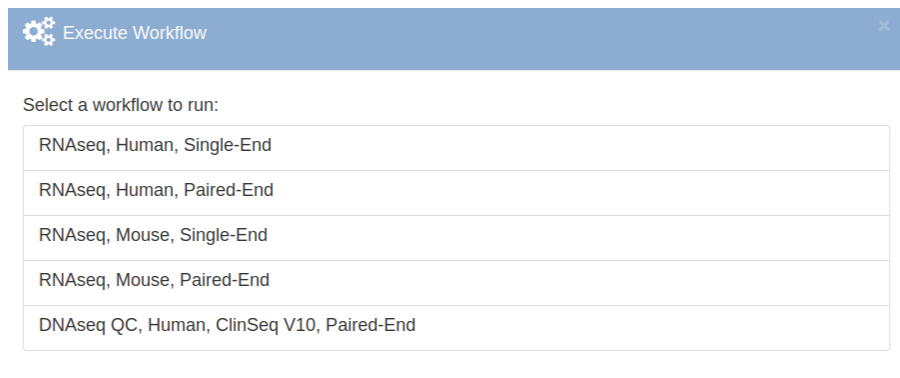
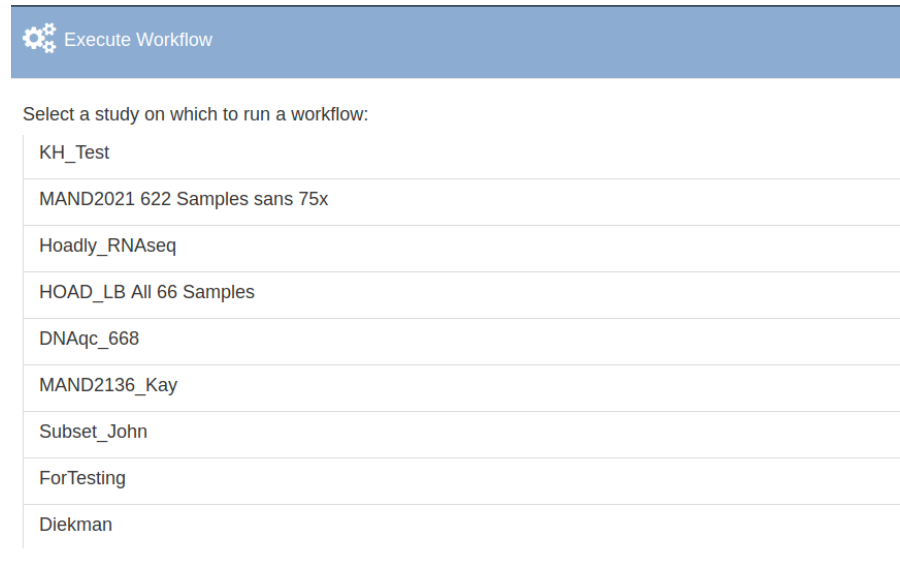
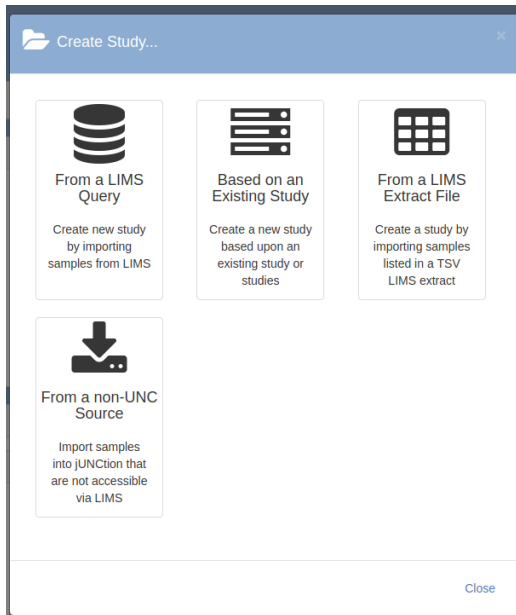
	Legacy	Current
Alignment	MapSplice	STAR
Expression	RSEM	Salmon
Normalization	Upper quartile	Upper quartile
Genome	GRCh37	GRCh38
Transcriptome	UCSC knownGene	GenCode v22
Runtime (min)	~800	80
GB RAM	16	4



# RNA-seq Workflow



# RNA-seq Workflow

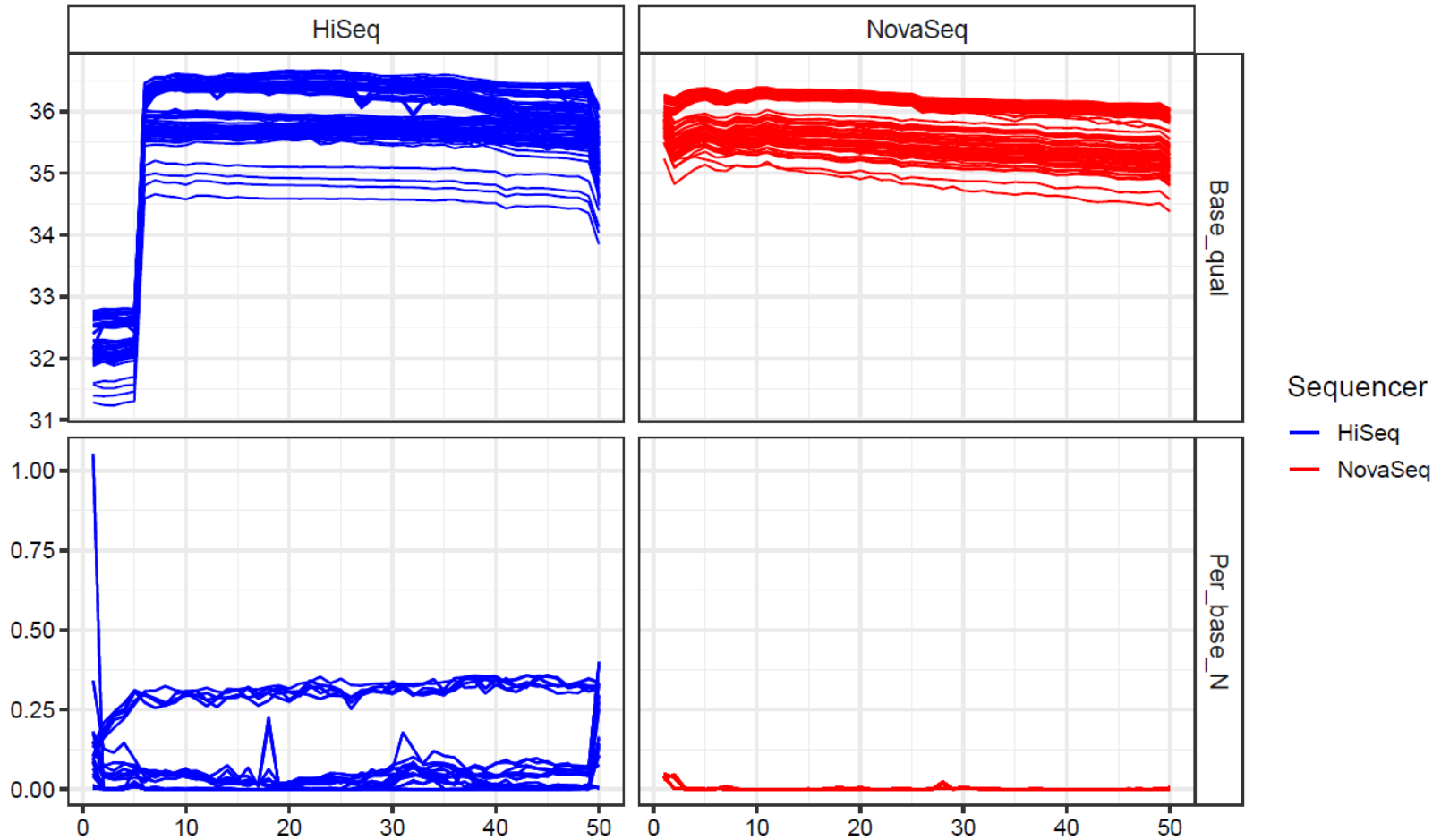


User & Lab (group) level security policies based on onyen

Provenance

User friendly

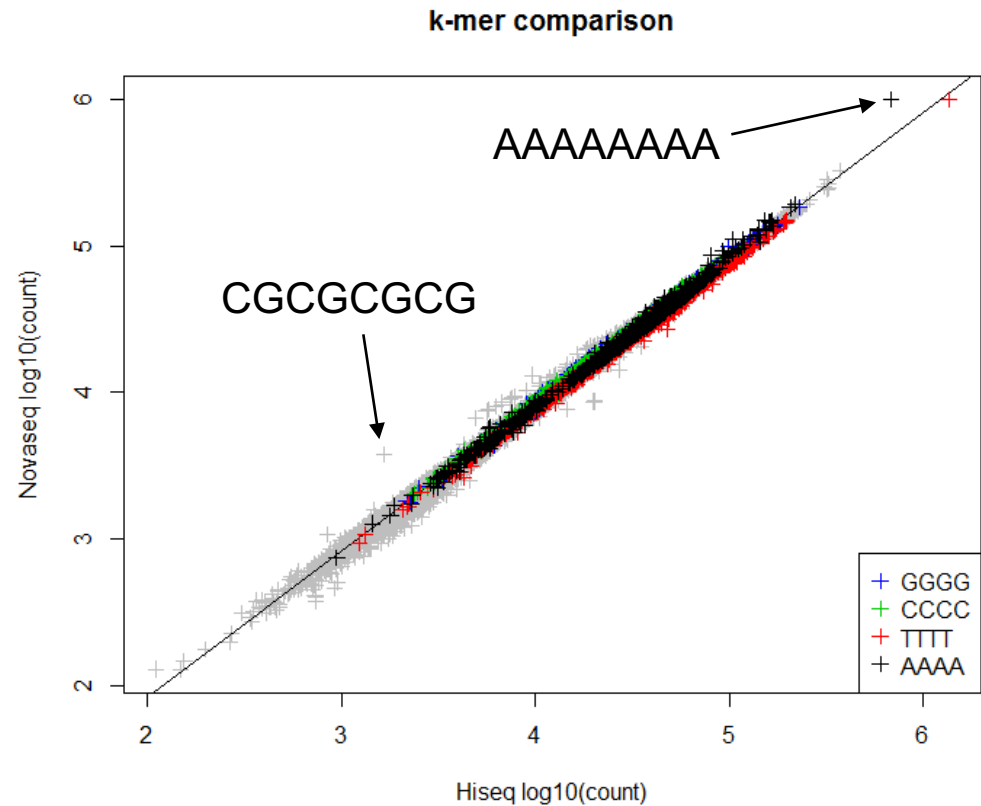
# Base Quality



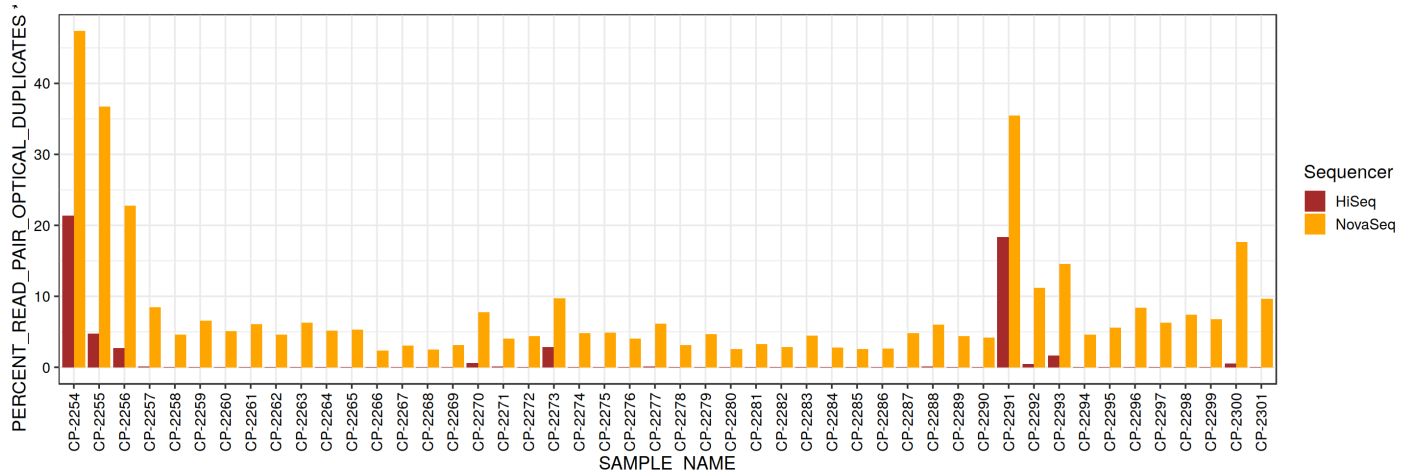
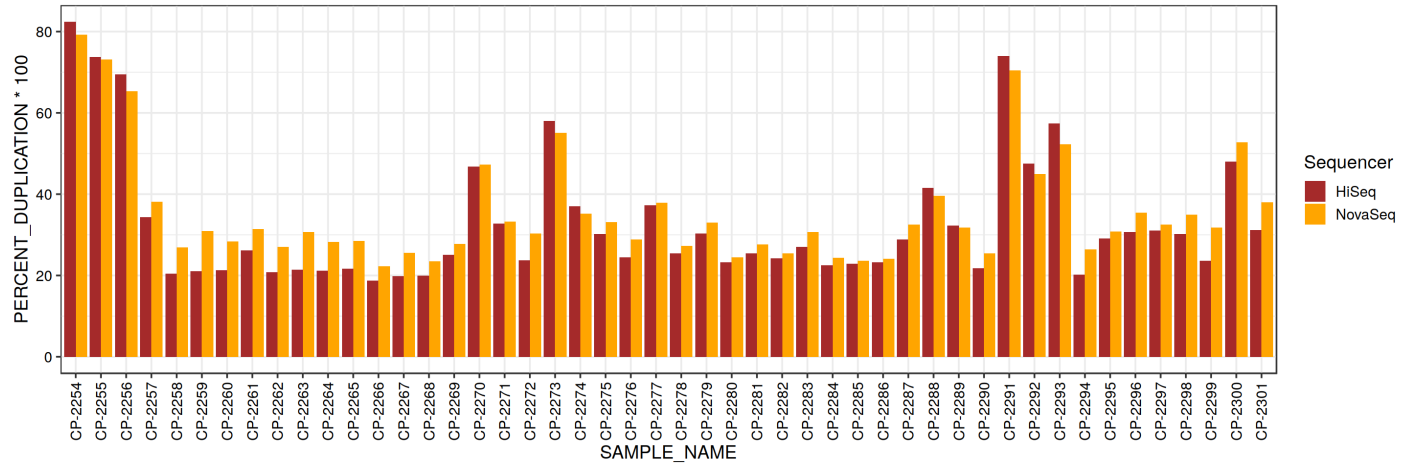
Novaseq base qualities are binned to  
2, 12, 23, 37

# Base composition

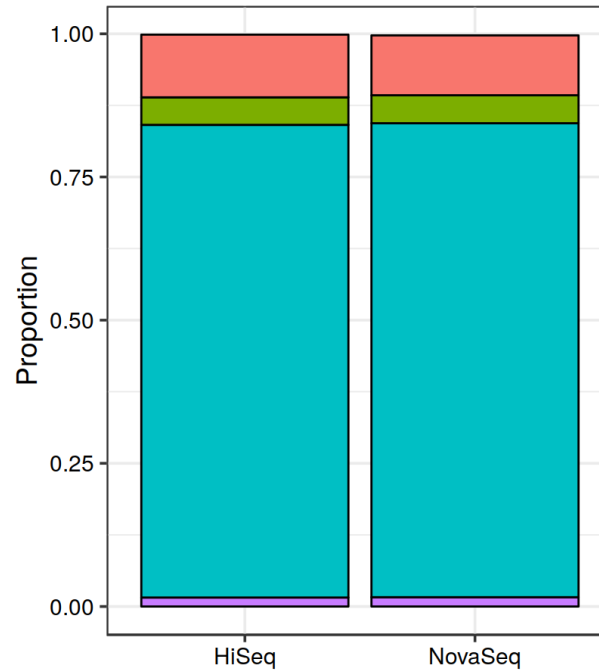
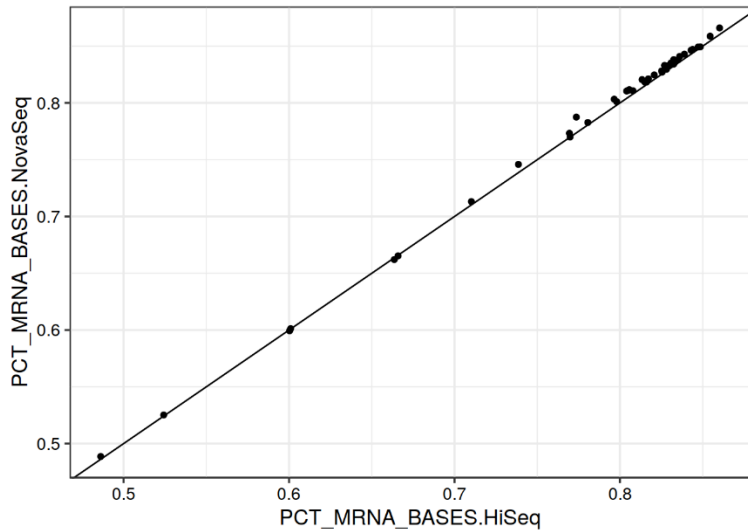
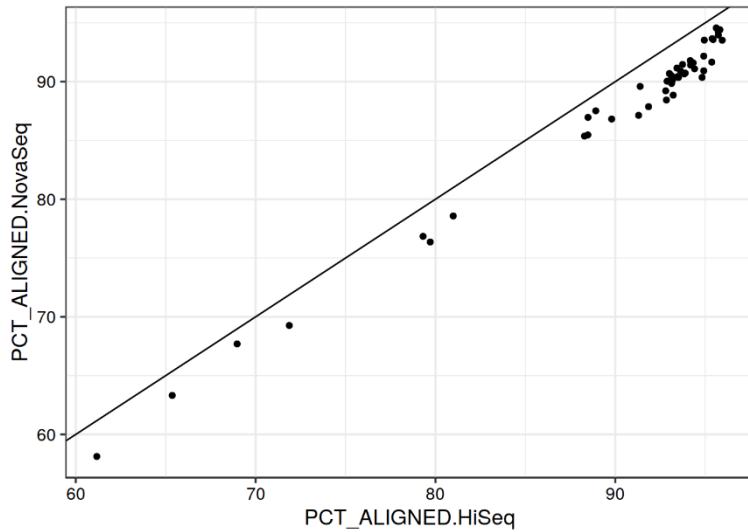
- K-mer counts (k=8 shown) tabulated from fastqs
- Expected bias in G rich sequences is not observed
- polyA and CG repeating sequences are mildly enriched in the Novaseq run



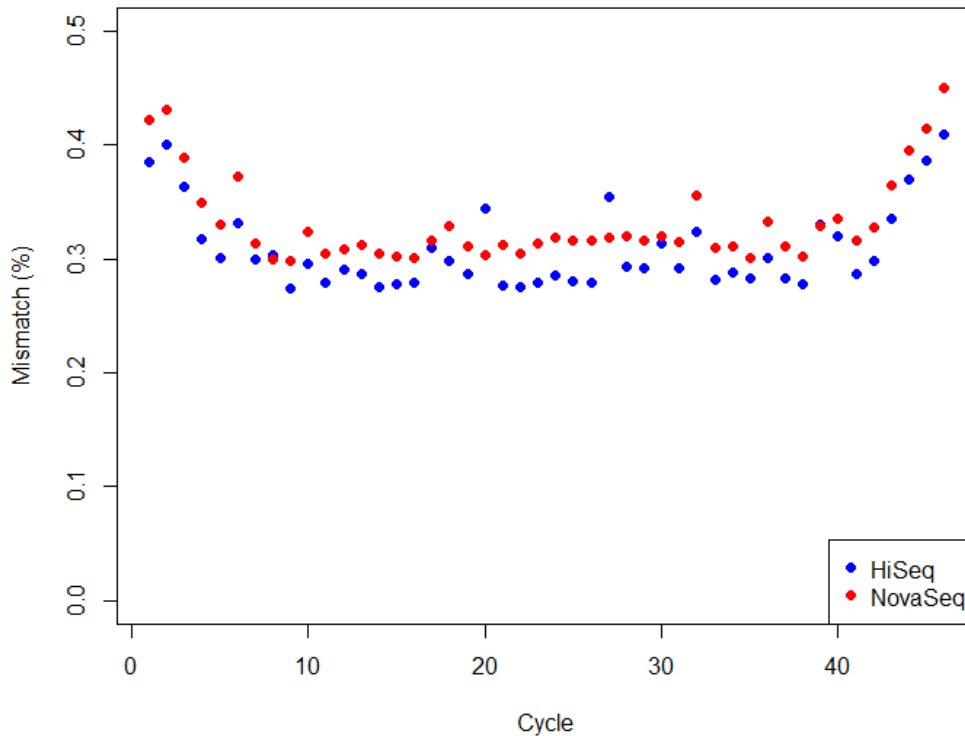
# Duplication



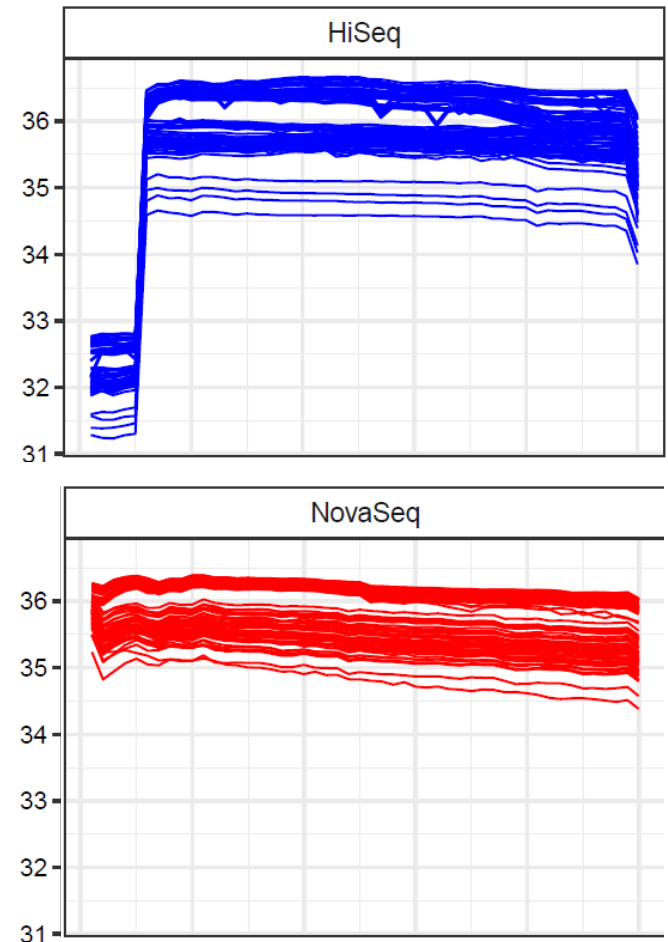
# Alignment



# Alignment



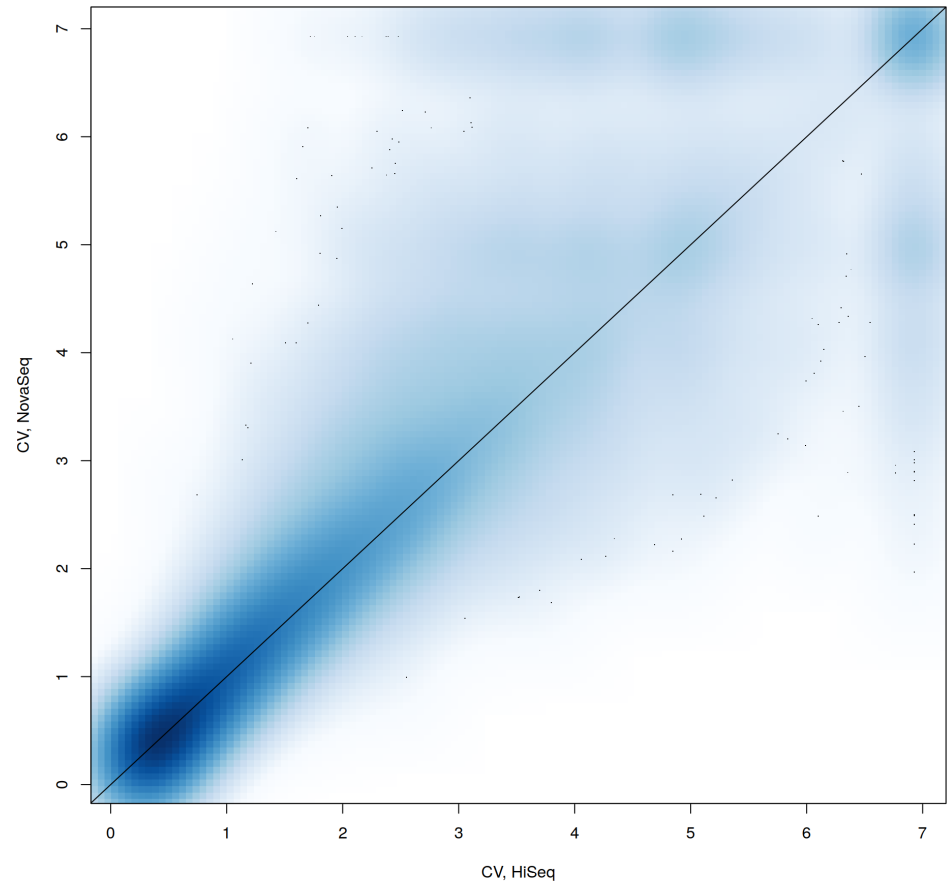
- Modest increase in mismatch rate
- Estimated Q scores do not appear accurate in the initial cycles



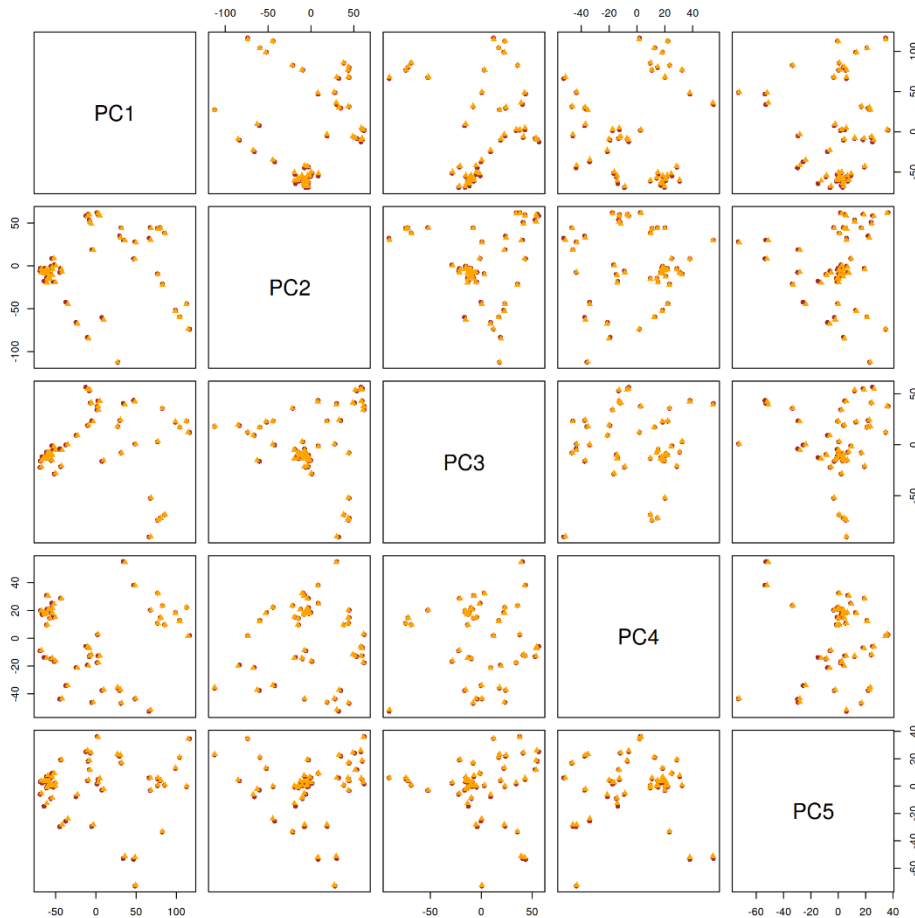


# Repeatability

- Gene coefficient of variations (CV) estimated for each cohort
- CV estimates indicate high concordance in expression variation

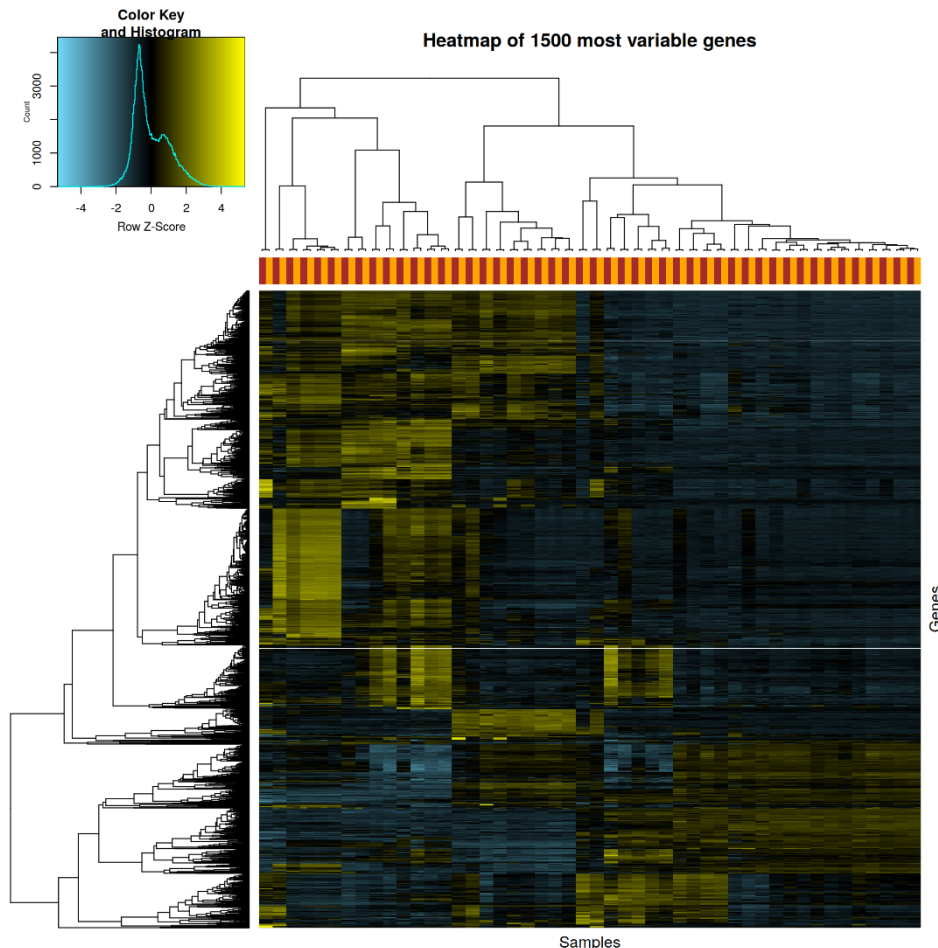


# Unsupervised Comparison



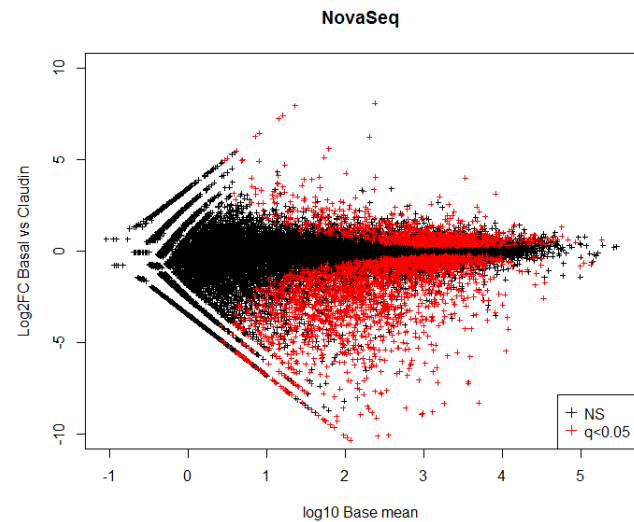
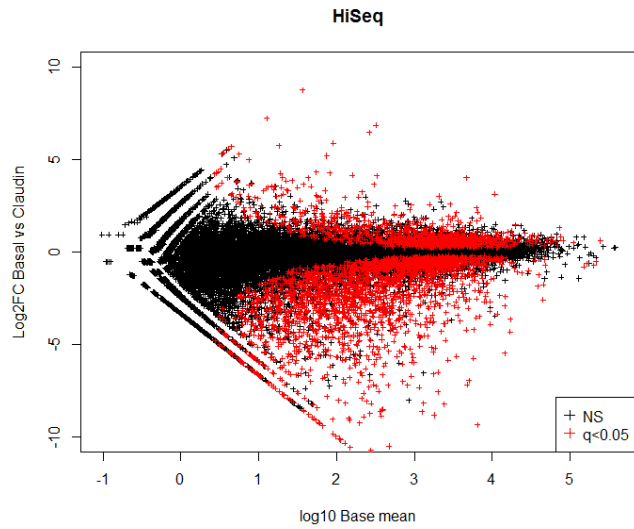
- PCA performed on log transformed, upper quartile normalized count estimates
- Nuisance variation due to instrument bias is not observed in any of the top components of variation
- PC35 (0.62% of expression variation) is the most highly ranked component associated with instrument (WRS  $p < 0.001$ )

# Unsupervised Comparison

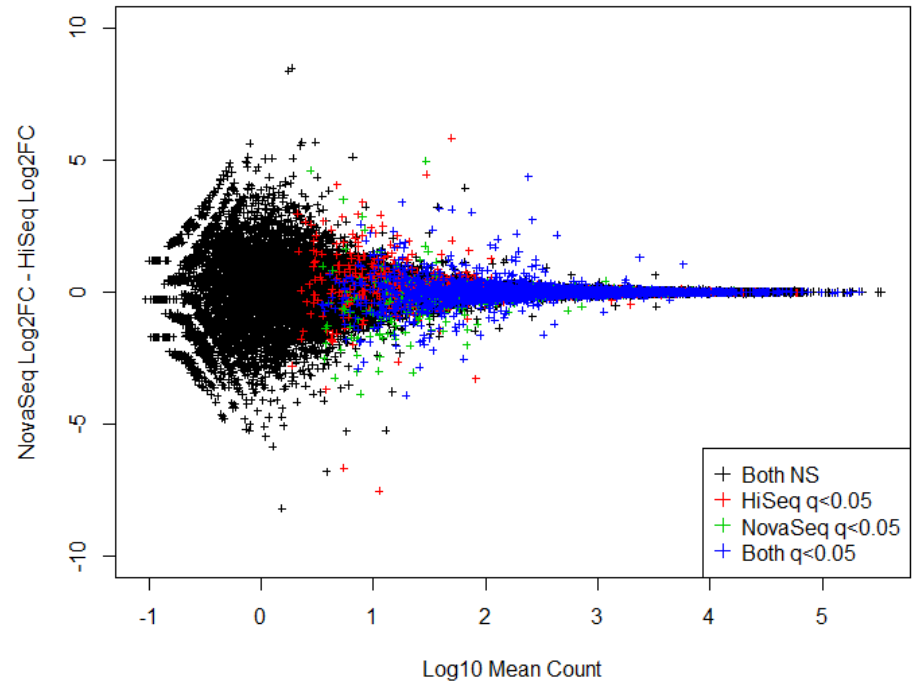


- Unsupervised hierarchical cluster analysis using the 1500 the most variable genes
- All paired samples are more similar to one another than other samples from the same instrument
- Magnitude of variation is also preserved (as before in CV plot)

# Supervised Comparison

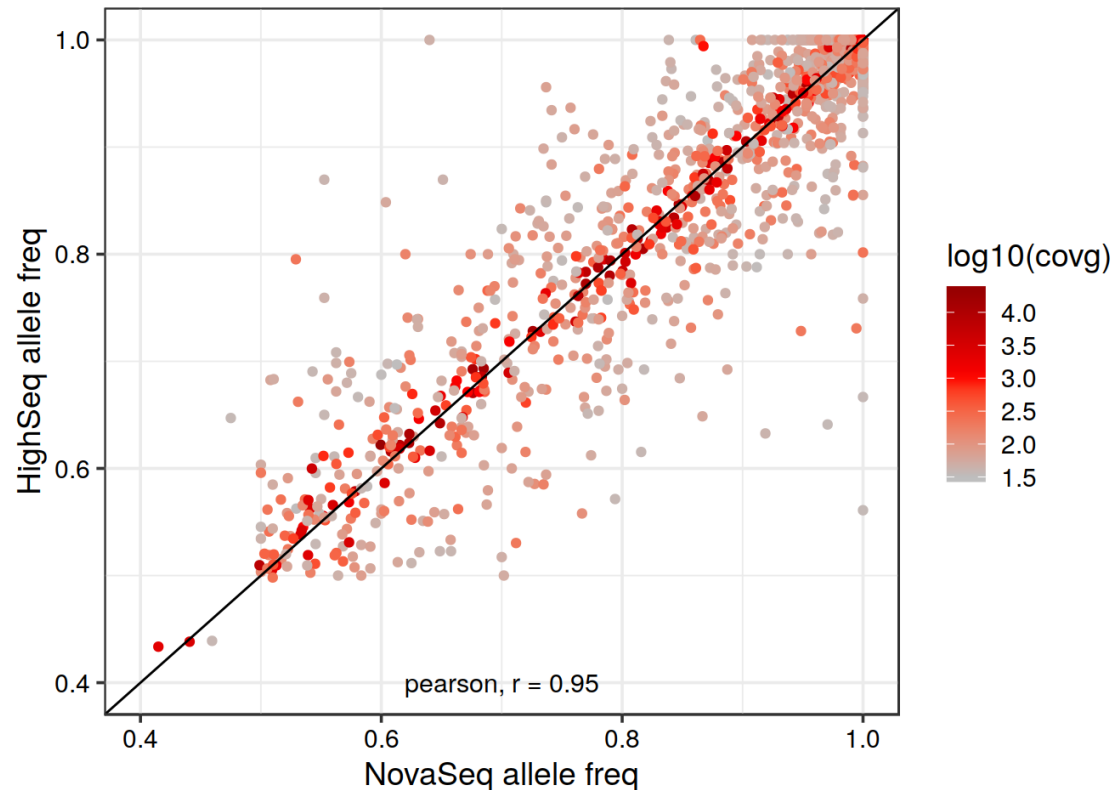


DESeq2  
Claudin-low (n=4) vs Basal-like (n=4)



# VAF Comparison by RNA-seq

- Exome sequencing available for one sample (HiSeq)
- All coding variants identified in DNA were quantified in RNA
- Allele frequencies by Novaseq are concordant with those of HiSeq
- High expected agreement for sequencing applications



# Summary

- Systematic bias is expected when changing protocols
- Expected sources of bias – sequencing chemistry and patterned flow cells – did not broadly affect experimental results
- The magnitude of instrument bias is negligible relative observed biological variation

# Acknowledgements

- Lineberger Bioinformatics

Sara Selitsky



David Marron



John Ford



Lisle Mose



- Perou lab

Xiaping He

Lynn Challot

- HTSF

Amy Perou

