



# Whole genome Next Generation DNA Sequencing (NGS) and Third Generation DNA Sequencing (TGS) data analysis

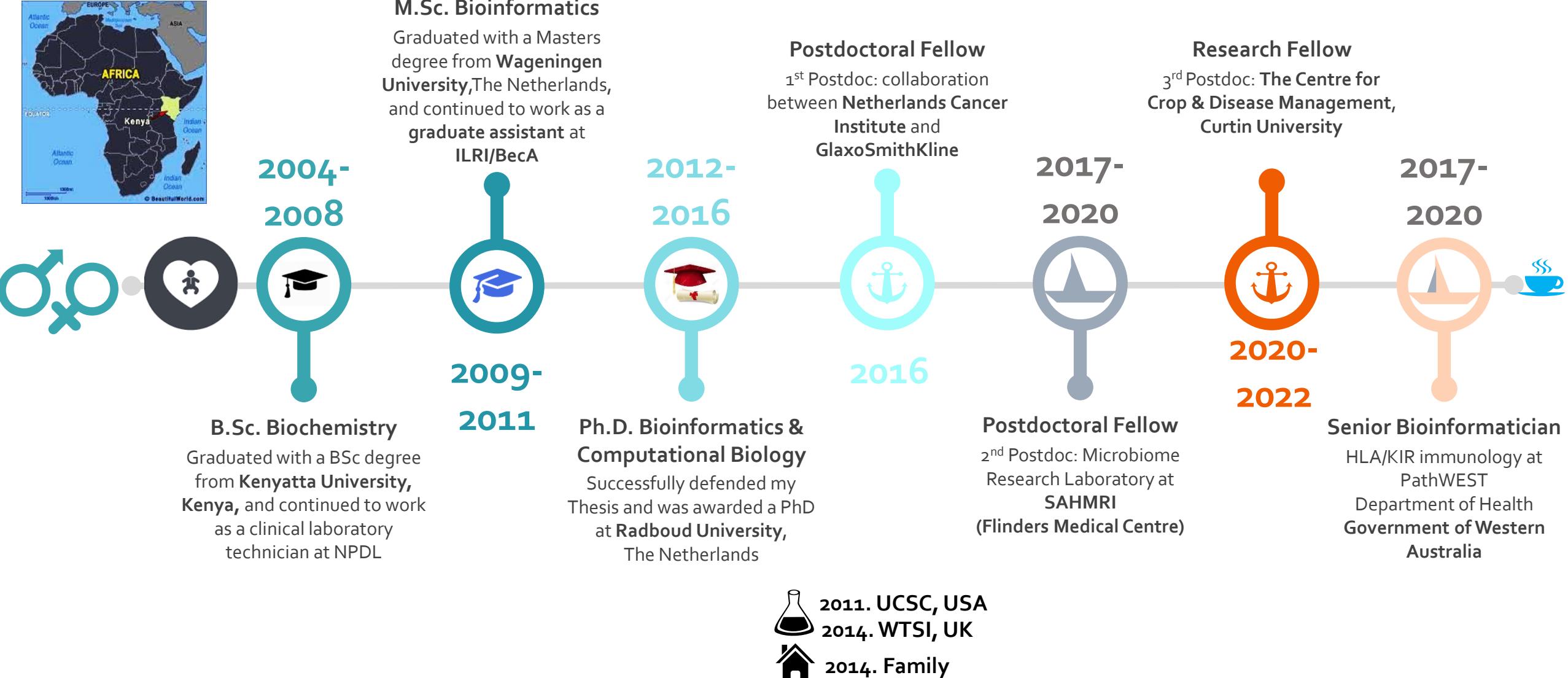
---

Fredrick M. Mobegi, PhD  
BHKi hybrid seminar series



[linkedin.com/fmobegi](https://linkedin.com/fmobegi)  
[twitter.com/mobeginomics](https://twitter.com/mobeginomics)  
[github.com/fmobegi](https://github.com/fmobegi)  
[rpubs.com/fmobegi](https://rpubs.com/fmobegi)

# Once upon a time in a small village...



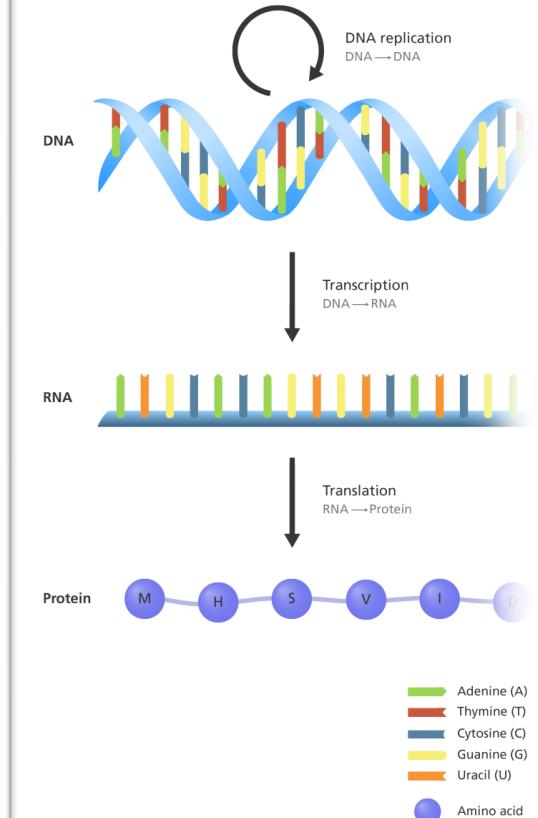
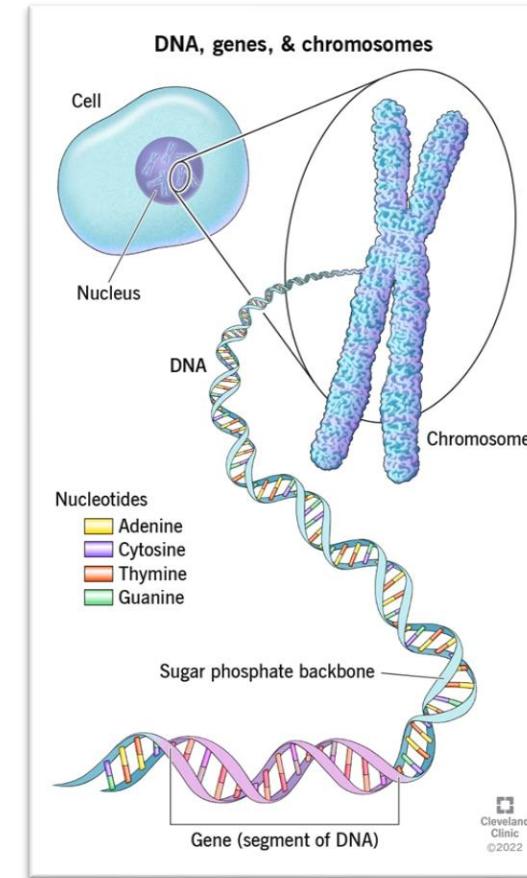
# Let's get some terminology out of the way...

Genetics is the study of how genes and traits are passed down from one generation to the next: **Heredity**

Genomics is the study of the total or part of the genetic or epigenetic sequence information of organisms, and attempts to characterize the structure and function of these sequences and of downstream biological products: **Architecture**

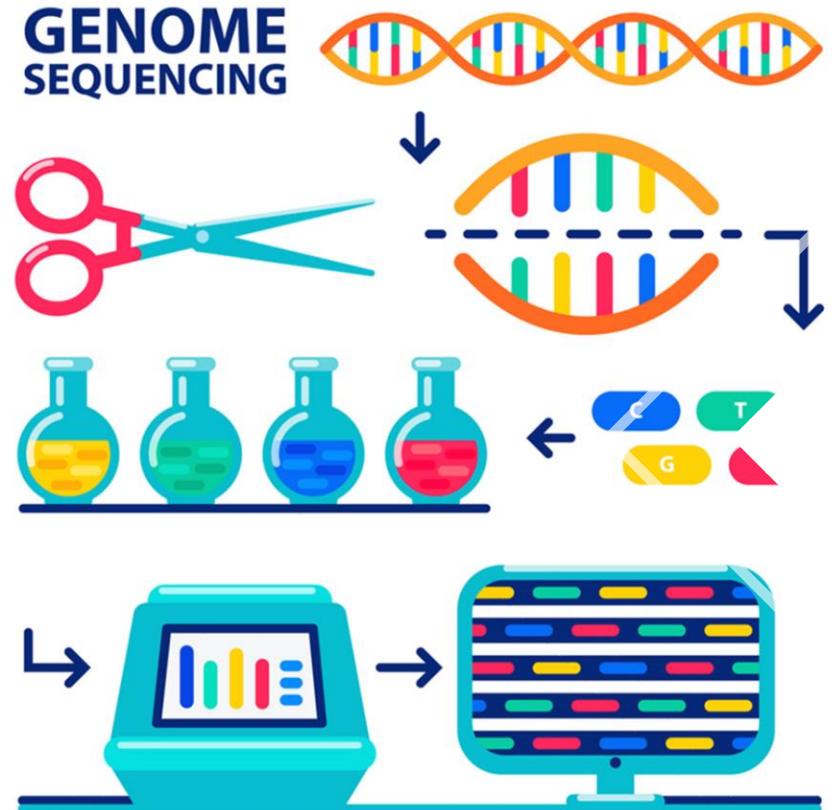
A genome is an organism's complete set of deoxyribonucleic acid (DNA), a chemical compound that contains the genetic instructions needed to develop and direct the activities of every organism.

A gene is the basic physical and functional unit of heredity.

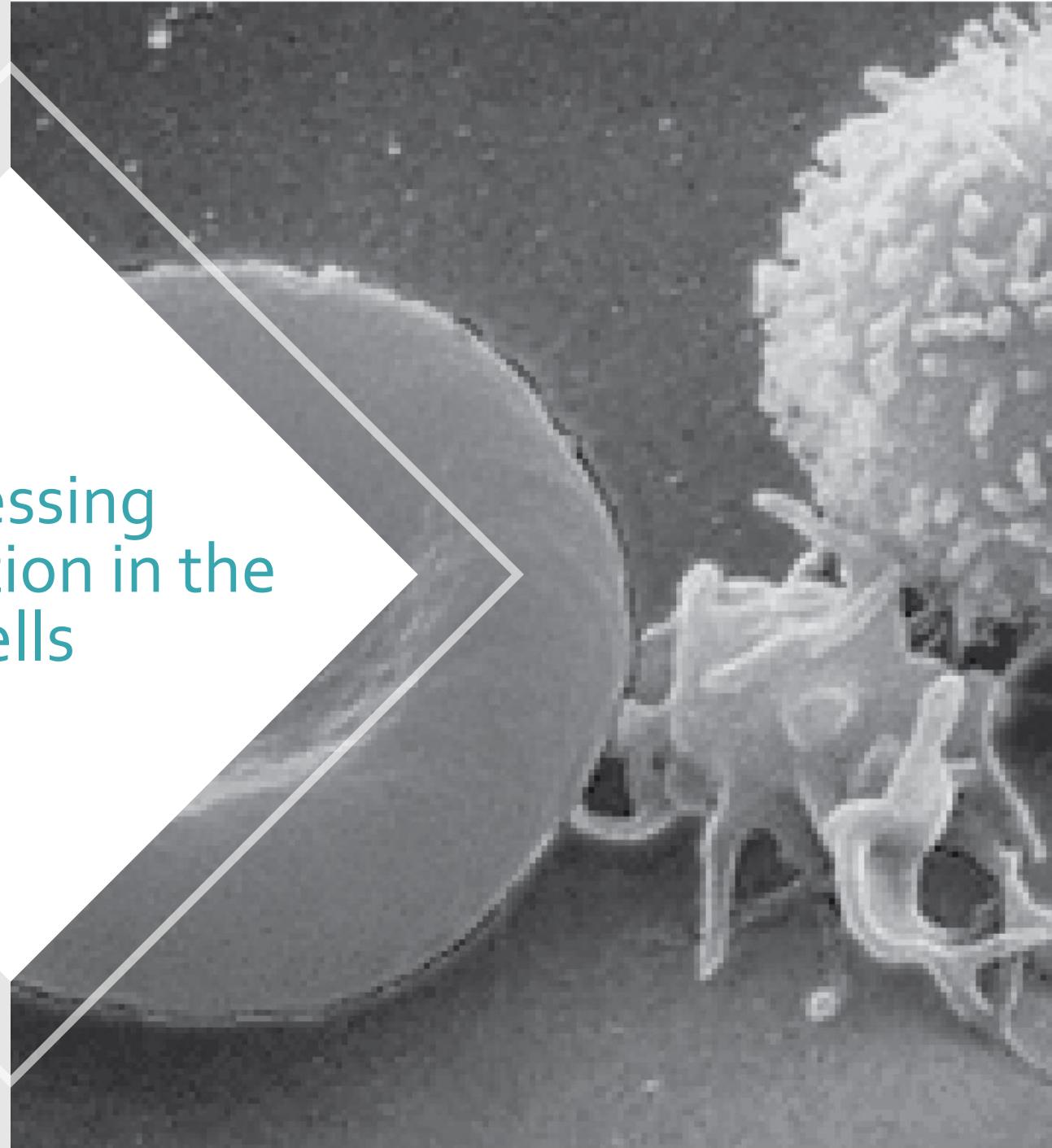


- <https://www.who.int/news-room/questions-and-answers/item/genomics>
- <https://www.amnh.org/explore/ology/genetics/what-is-genetics>

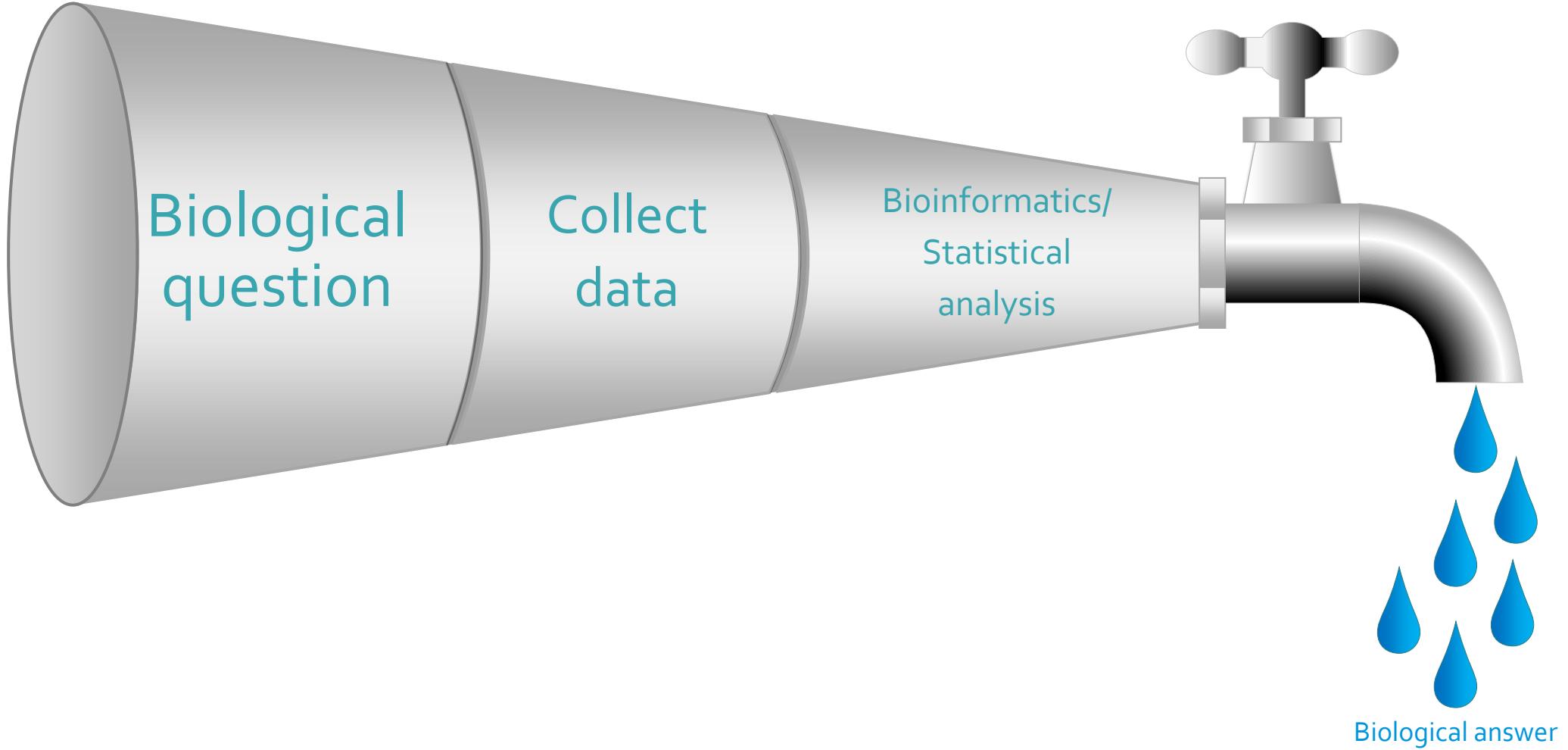
## GENOME SEQUENCING



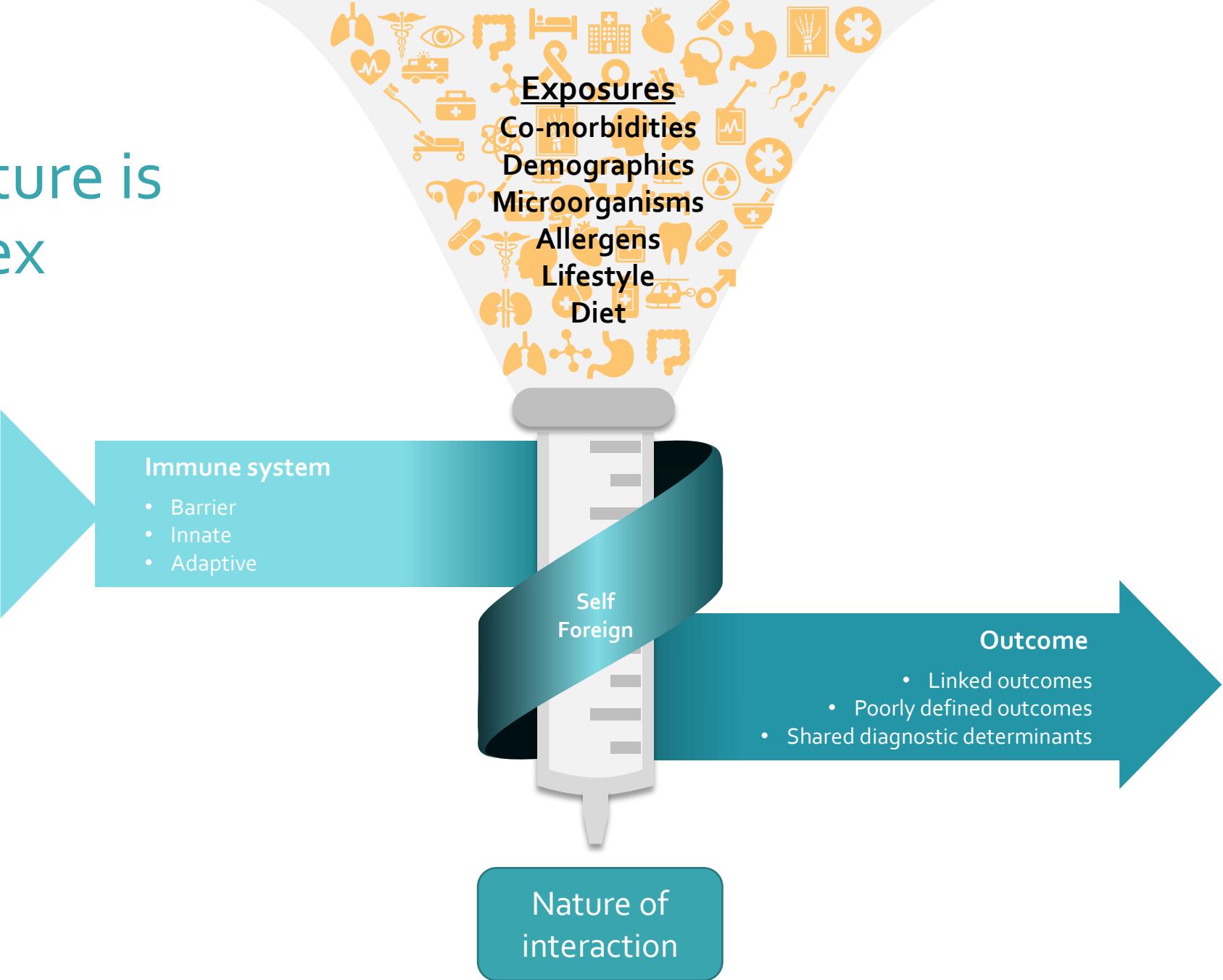
Accessing  
information in the  
cells



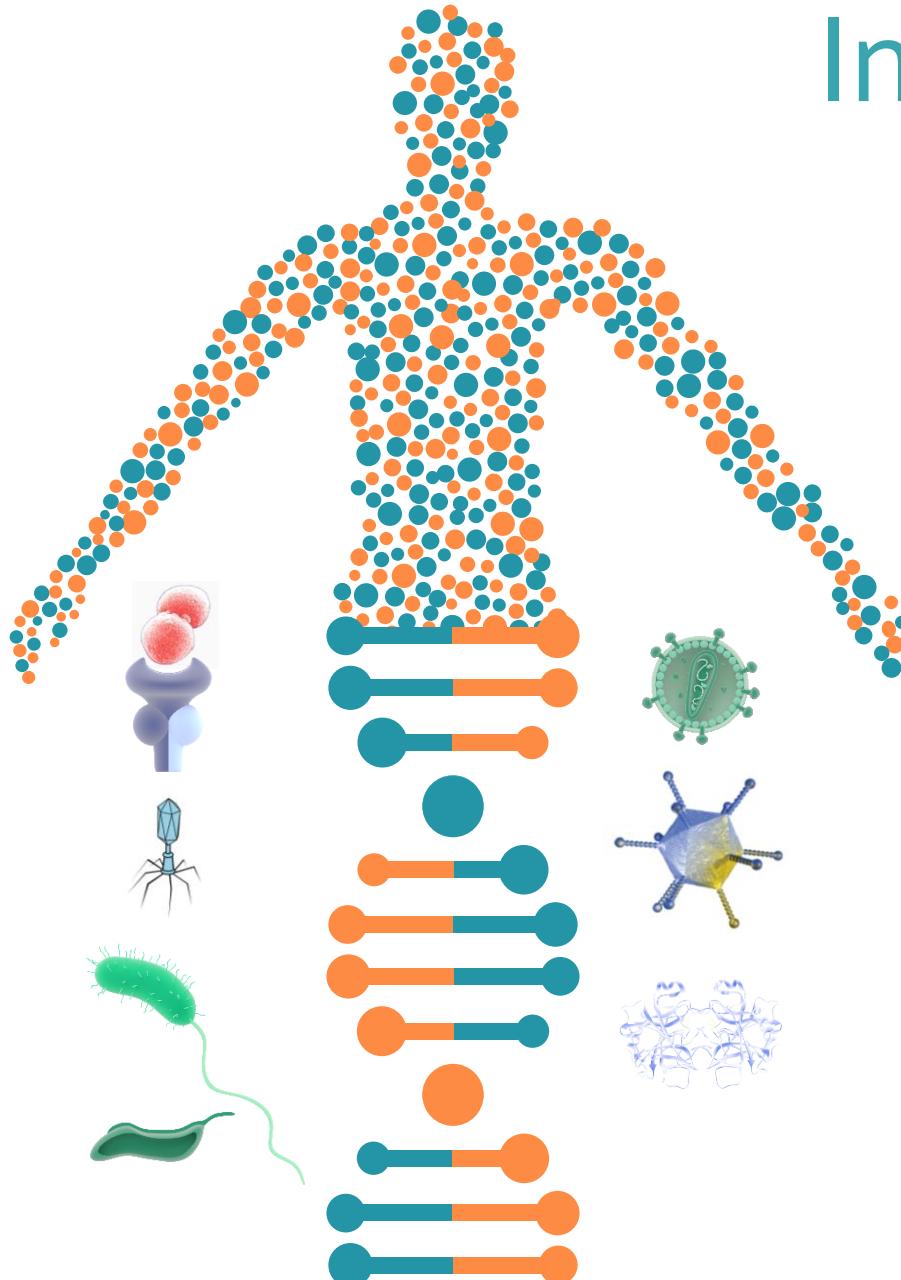
In an ideal situation, we like our science simple...



# But nature is complex



# Integrated genomics concept



- The immune system's ability to distinguish self from non-self is a critical process for our health.
- Breakdown of checkpoints for the "self-nonself" discrimination results in the hundred different autoimmune diseases
- Technological advances in computing and genomics, and the reducing cost of sequencing have opened-up possibilities to "cut through" the immunological complexity of clinical disease to identify causal genes, molecules, and cellular mechanisms.

Malfunctioning immune system has been associated with the development or progression of several human diseases.

## Cardiometabolic disease

- Obesity
- Diabetes
- CVD /CKD
- Asthma
- coronary artery disease
- Heart disease
- stroke

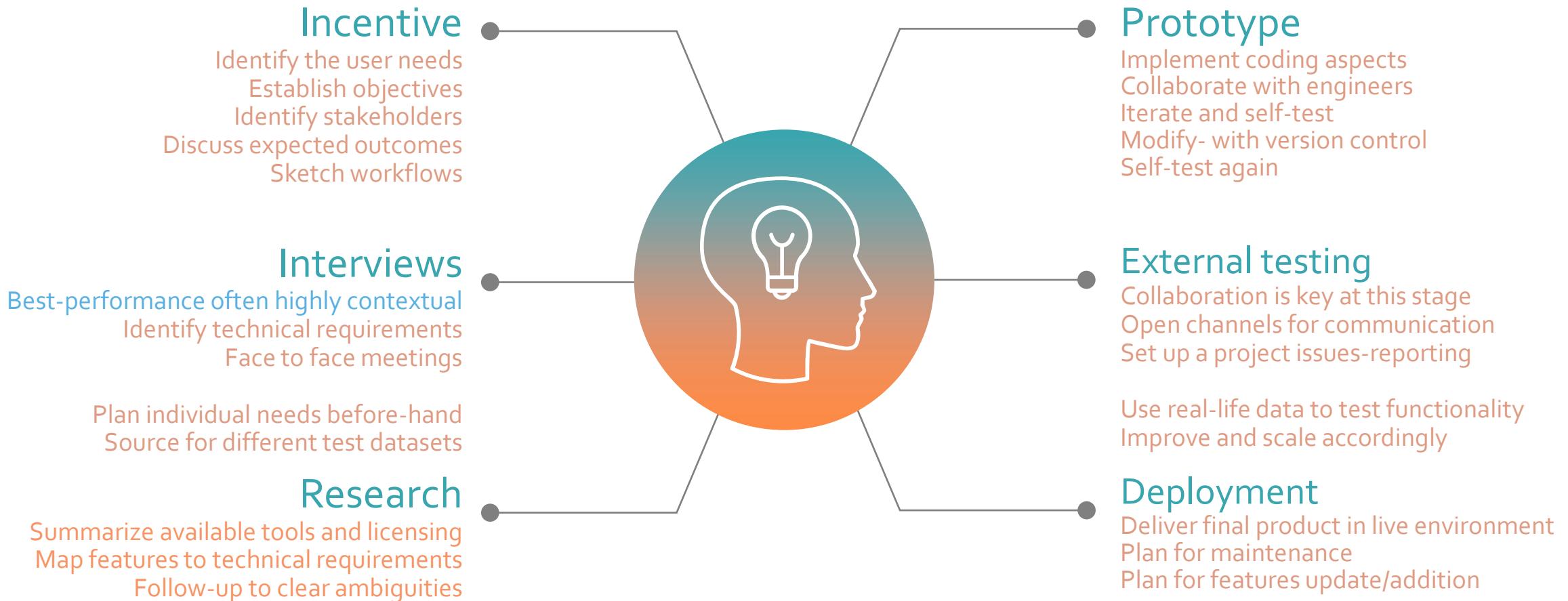
## Inflammation

- IBD
- Ulcerative colitis
- Crohn's disease
- Multiple sclerosis
- Systemic lupus erythematosus
- Arthritis

## Others

- Cancer
- Alzheimer's
- Autoimmune diseases
- Parkinson's
- Allergies
- Graves' disease
- Coeliac disease

# Genomics workflow projects: planning



# Genomics workflow projects: IT infrastructure

**Developer**

Determine resources needed  
*Profiling and Timing Code  
Tracing & visualisation*

Scalable strategies  
*Shared-Memory Multicore Architecture  
Special hardware*

Portability and platform scalability  
*Multi-Node HPC scalability  
Cloud scalability  
Container Scalability*



**Support**  
Kind of support  
*Full-time support  
occasional support*

**Efficiency**  
*Serverless web services  
HPC/Cloud services*

**Level of support**  
*MPI tools require seasoned engineers to handle memory control, data communication, and task synchronization*

## Parallelization

How much compute power is needed?

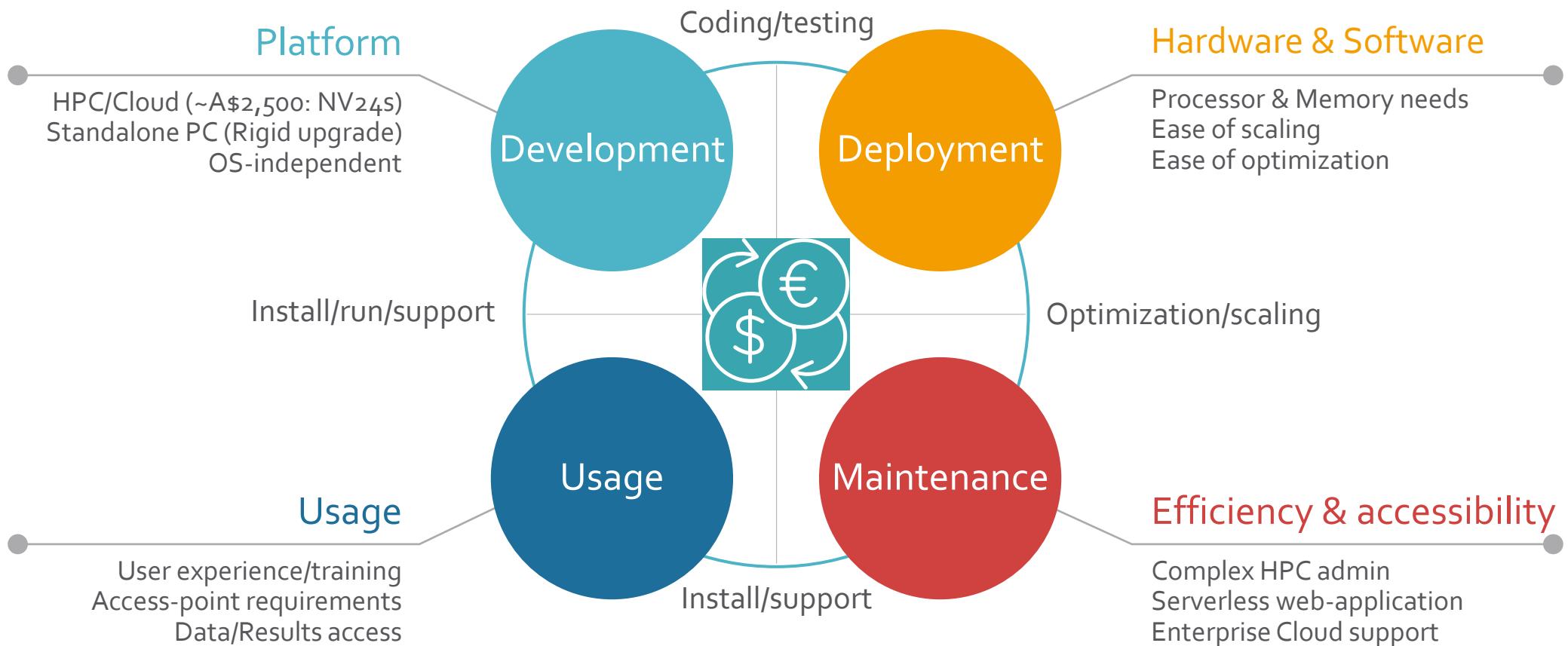
## Scalability and optimization

Can we extend beyond primary state?

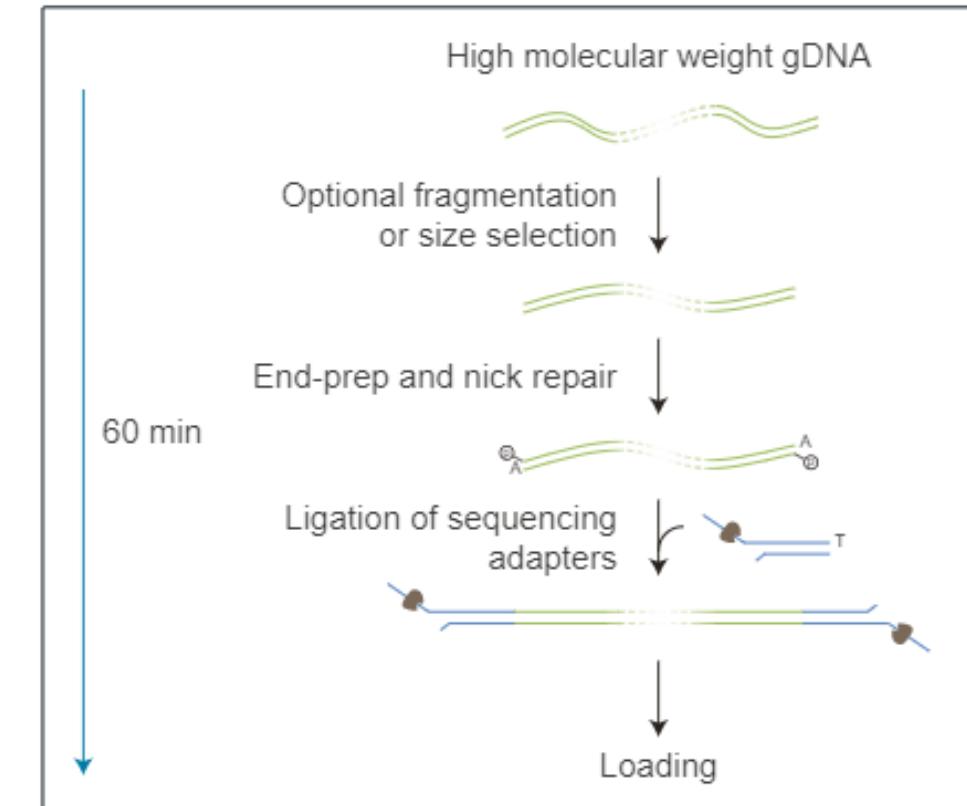
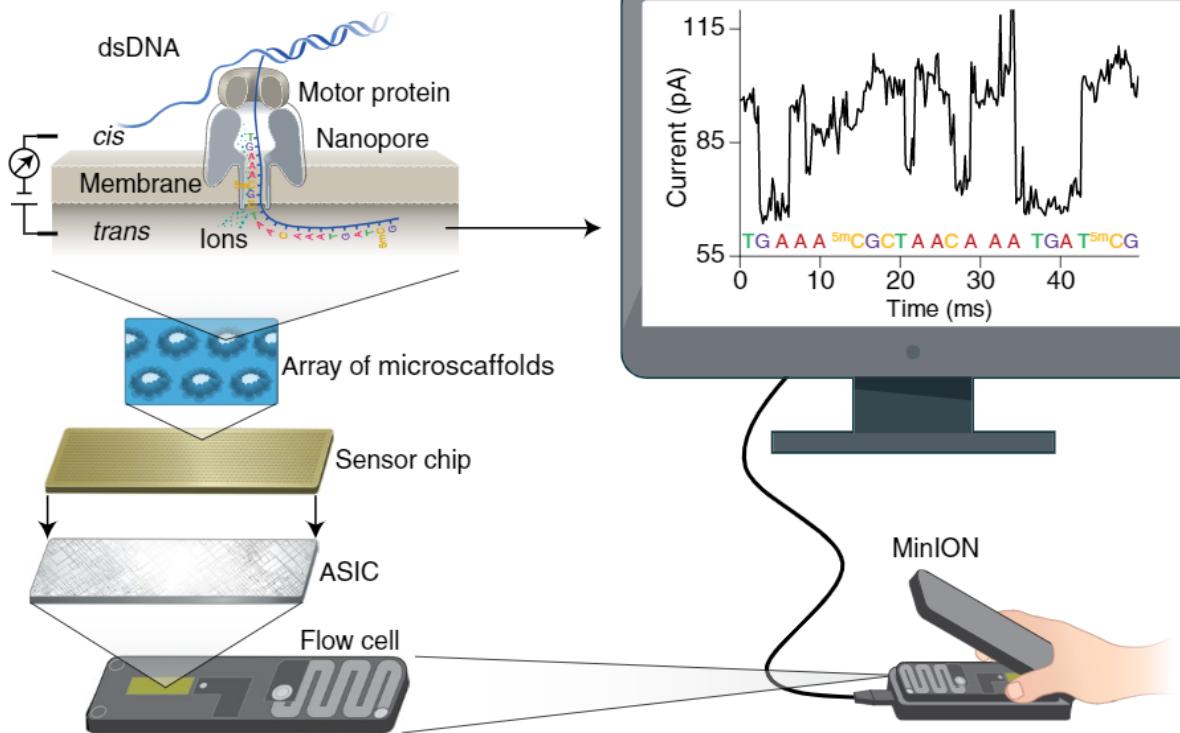
## Integration & portability

Can we use existing resources?

# Genomics workflow projects: Budget



# Genome sequencing: Oxford Nanopore Technology



# Genome sequencing: Oxford Nanopore Technology



SINGLE MOLECULE  
SEQUENCING



ENABLES REAL-TIME  
ANALYSIS



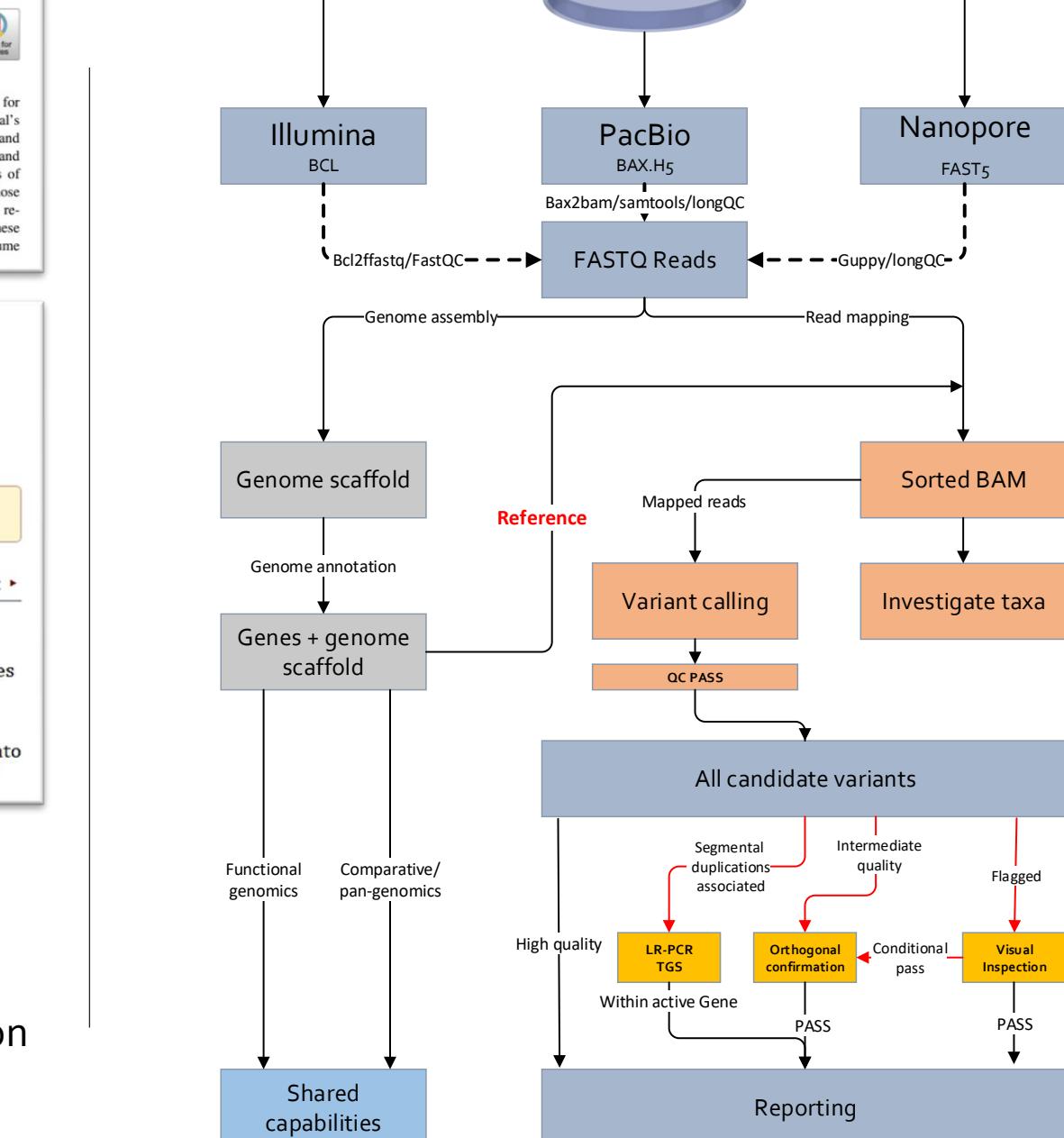
DIRECT RNA  
SEQUENCING

## IuAC: Implementing quality control best practices for genome sequencing and exome sequencing data

Manavalan Gajapathy<sup>1</sup>, Brandon Wilk<sup>1</sup>, Donna Brown<sup>1</sup>, Elizabeth Worthey<sup>1</sup>

The University of Alabama at Birmingham

**Introduction:** Quality Control (QC) of genome sequencing (GS) and exome sequencing (ES) data is necessary to ensure that data are of sufficient quality for downstream analyses. This would ensure that sequenced reads pass expected measures of quality, inferences from sequenced data match sequenced individual's expected metadata (sex, ancestry, relatedness) to identify sample swaps, samples are free of contamination from other human samples or other species, and unexpected batch effects are caught. While several QC tools are available to perform QC at various levels post sequencing, output needs to be reviewed and interpreted in a very manual process. Such manual review can be a challenge, especially for large projects involving hundreds of samples, in terms of standardization and consistency, as it can be subjective based on the reviewer. Further, logging the results of QC review and disseminating them with those involved in downstream analysis in an understandable format can be challenging, and could result in downstream users not utilizing the QC review or re-reviewing them on their own and thereby wasting time and effort. Here we present the best practices that we have implemented to remove and/or alleviate these challenges. We have developed a pipeline to run QC tools at various stages of secondary analysis, visualization and sharing of results in an easy-to-consume



## Sequencing quality assessment tools to enable data-driven informatics for high throughput genomics

Richard M. Leggett,\* Ricardo H. Ramirez-Gonzalez, Bernardo J. Clavijo, Darren Waite, and Robert P. Davey

► Author information ► Article notes ► Copyright and License information ► Disclaimer

This article has been cited by other articles in PMC.

### Abstract

Go to: ►

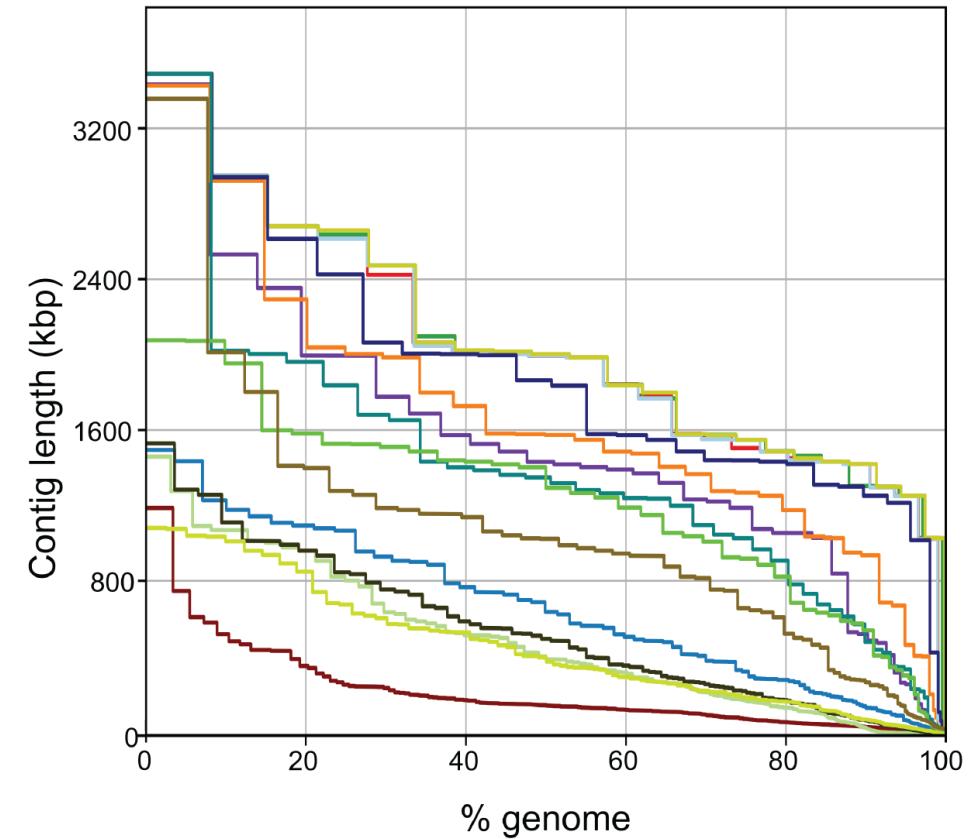
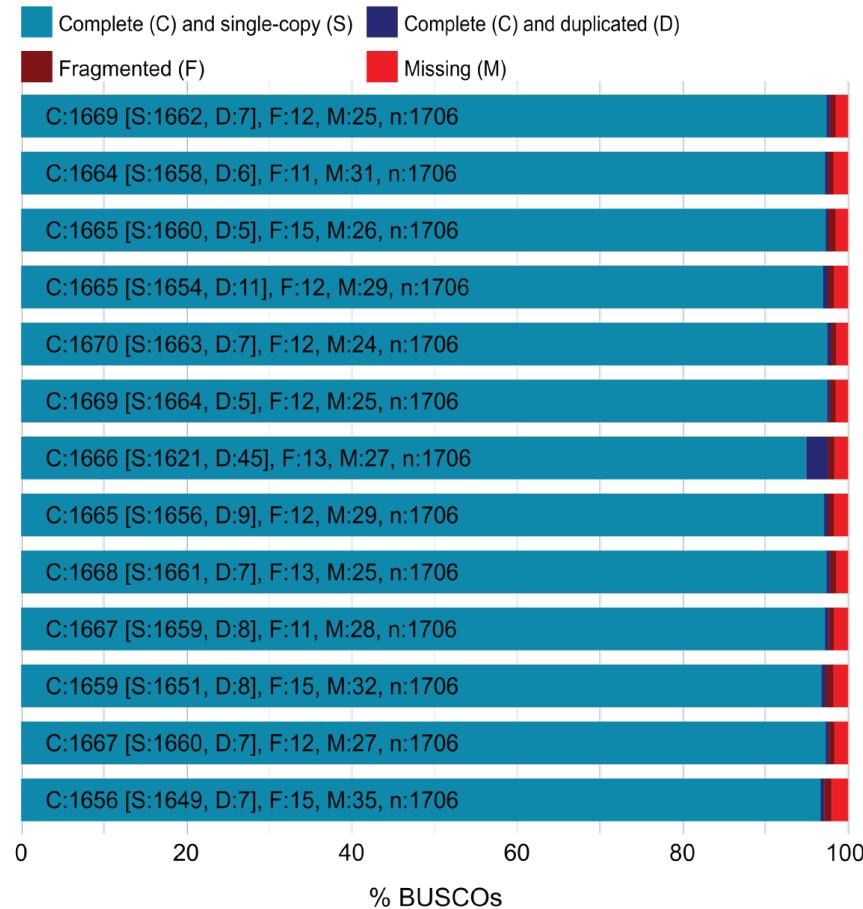
The processes of quality assessment and control are an active area of research at The Genome Analysis Centre (TGAC). Unlike other sequencing centers that often concentrate on a certain species or technology, TGAC applies expertise in genomics and bioinformatics to a wide range of projects, often requiring bespoke wet lab and *in silico* workflows. TGAC is fortunate to have access to a diverse range of sequencing and analysis platforms, and we are at the forefront of investigations into library quality and sequence data assessment. We have developed and implemented a number of

- Virtually all downstream analysis and interpretation processes rely on accurate data QC and validation.
- Most genomics workflows are interconnected on the critical reads pre-processing steps.

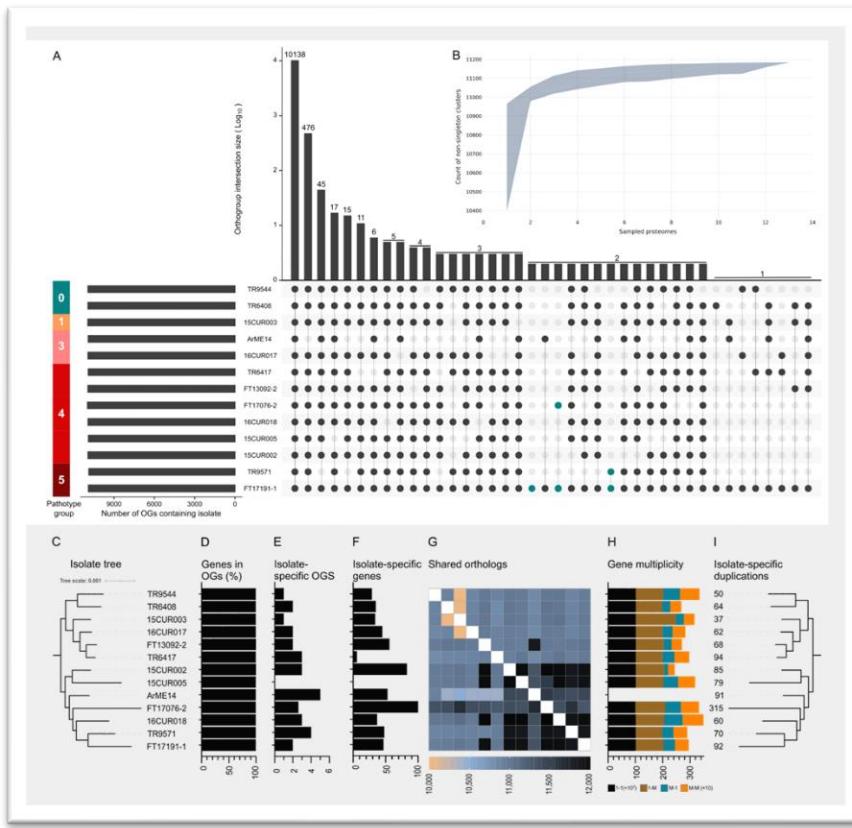
# Genome assembly



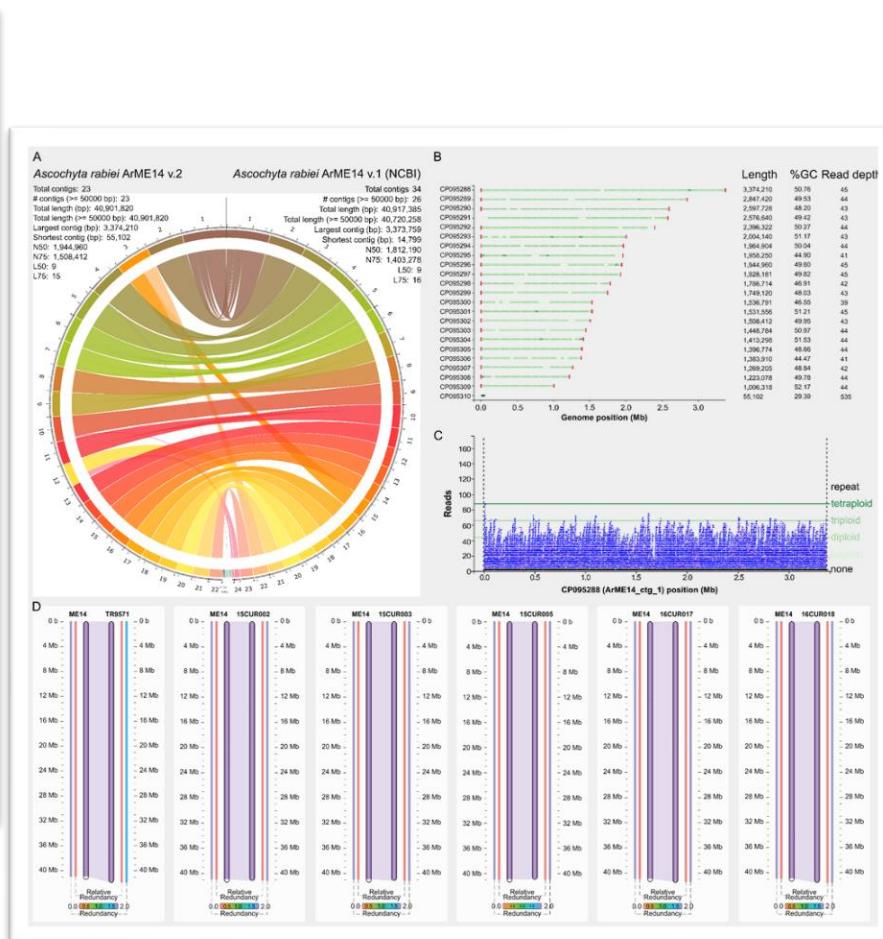
# Genome assembly: quality assessment



# Comparative analysis

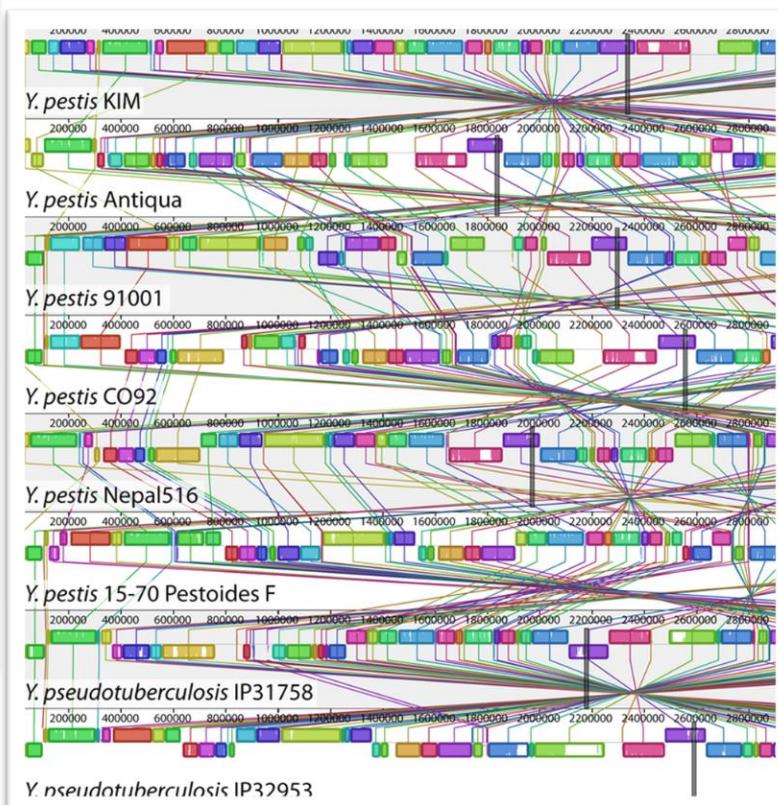


Orthologs

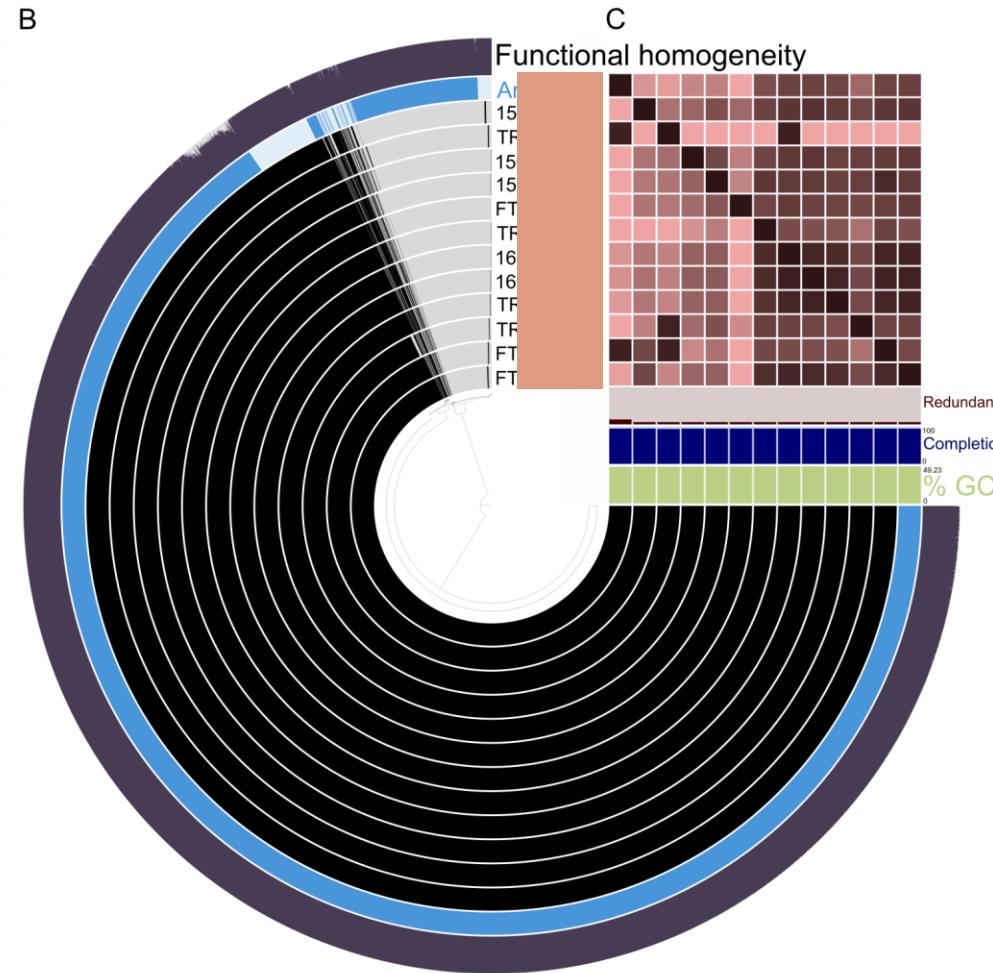
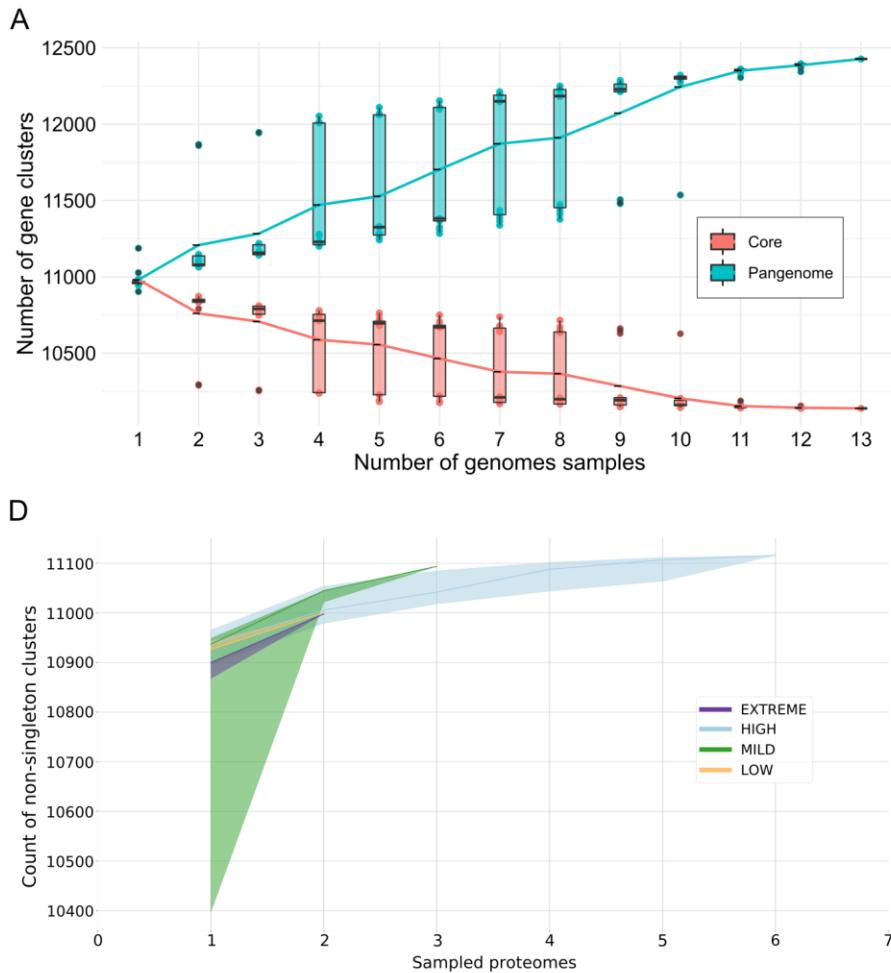


Genome structure

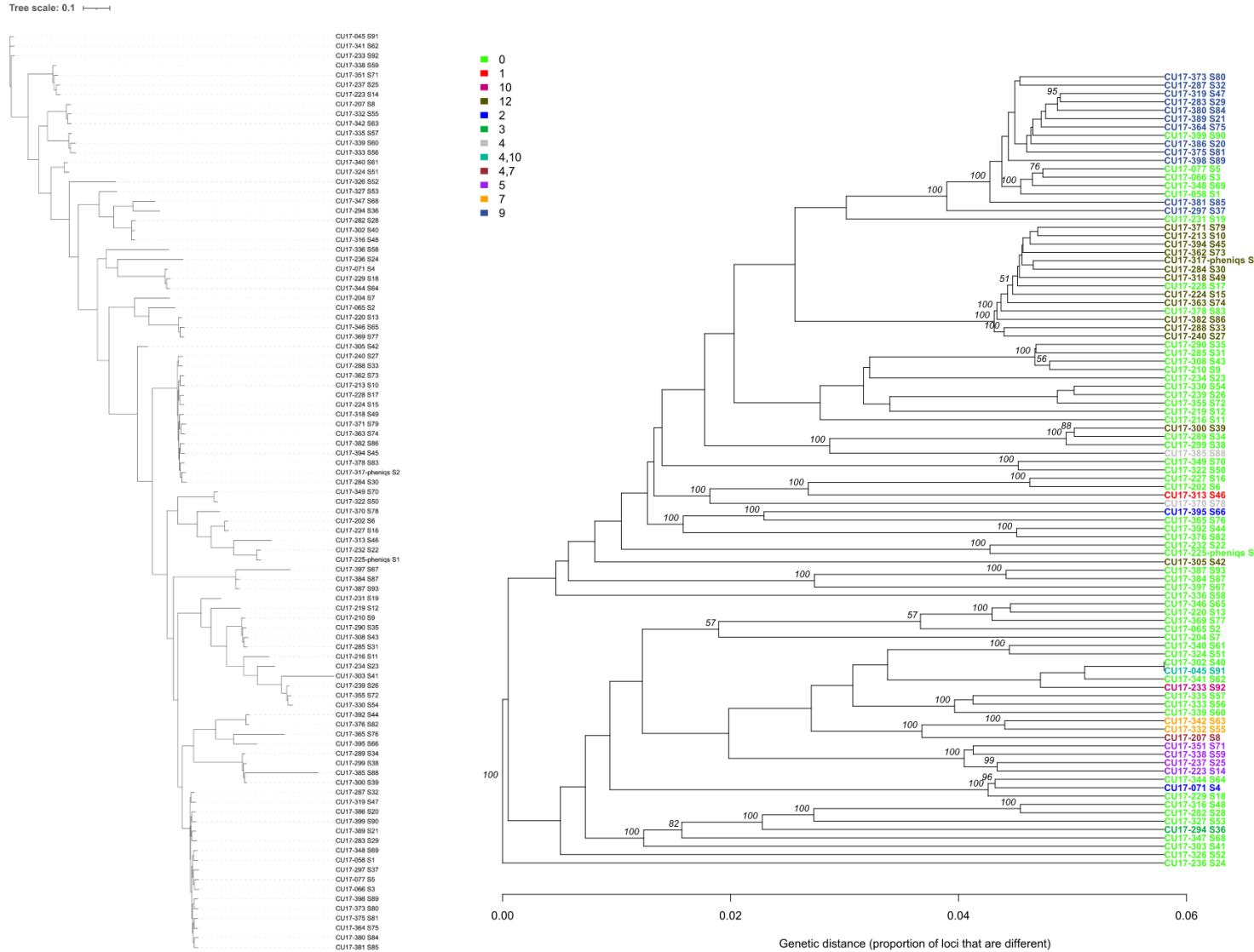
Whole genome alignment



# Pan-genomics



# Phylogenetic analysis



# Variant calling

## Quality trim and filter

### Assess quality

- FastQC (NGS)
- LongQC/Falco (TGS)

### Trimming and filtering

- Trimmomatic
- Canu2

### Assess quality

- FastQC (NGS)
- LongQC/Falco (TGS)

FASTQ



BAM



Sorted-BAM



VCF



VCF;HTML;SVG;TXT



## Align to reference

### Index reference

- Aligner-specific

### Align QC reads

- BWA (splice-unaware)
- STAR; Hisat2

### Mark duplicates

- Picard; Sambamba
- GATK

### Compress SAM output

- Samtools

## Alignment cleanup

### Sort by coordinates

- Samtools sort

### Index BAM (accessible)

- Samtools index

### Assess BAM Quality

Coverage depth; %mapped;  
%unmapped; MapQ score

- Samtools flagstat
- QualiMap; BEDtools

### Sample quality

- KING; VerifyBamID

## Variant calling

### Generate VCF

- GATK; FreeBayes; BCFtools (SNV/indels)
- VarScan2; DeepSNV (Somatic)
- CoNVEX; ExomeCNV (CNV)
- InDelible; Sniffles; SVMerge (SV)

### Filter VCF & annotate

#PASS; #QUAL; #DP; SNV; INDELS  
CDS; Non-CDS(Intron; Silent, UTR)

- GATK; BCFtools
- VCFAnnotate; AWK; ALFA

## Variant reporting

### Functional implications

- AnnoVAR; ENSEMBL VEP; SnpEff
- Clinical data, family history, signaling pathways

### Disease/drug-associated

- ClinVar; ClinPred; Clinotator >> JIGV report

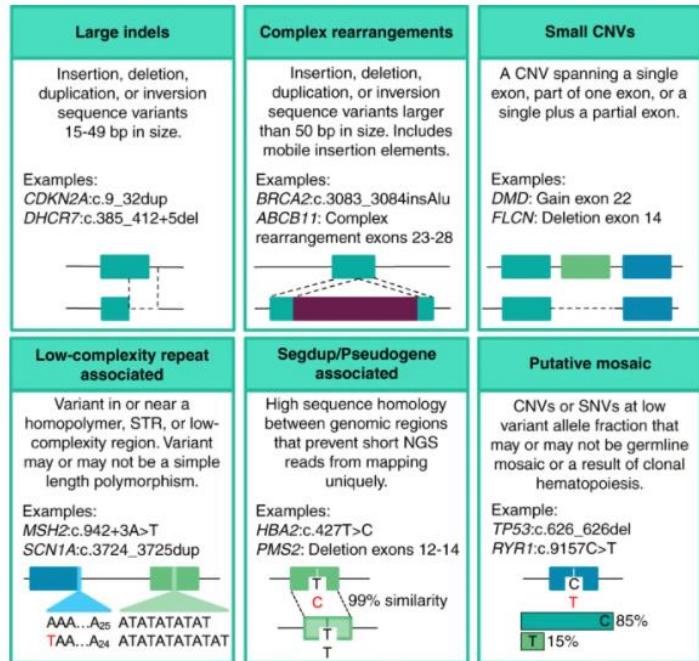
- Variant Annotation Integrator - UCSC Genome Browser
- Overall pipeline run summaries: MultiQC

## Pipelining in NetFlow:

- Fast to prototype and has scalability without extra code.
- Resumable tasks with continuous checkpoints.
- Portable (AZURE/AWS/HPC), reproducible, multi-language.
- Supports software containers and use of existing code "as is".
- Internal parallelization eliminates resources wastage.
- Focused resources-allocation to compute-intensive processes

# Expecting the expected...

Fig. 1: Technically challenging variant types.



Variants were categorized as being technically challenging, or not, based on these six criteria. Note that some variants could be considered challenging for multiple reasons (e.g., a single-exon deletion within a segmentally duplicated region). Examples provided are variants observed in the prevalence analysis. Detailed criteria are provided in the Supplemental Methods. CNV copy-number variant, indels, insertions or deletions, NGS next-generation sequencing, Segdup, segmental duplication, SNVs single-nucleotide variants, STR short tandem repeat.

Despite advances in technology, some variants are still complex to predict reliably

The Journal of Molecular Diagnostics  
Volume 21, Issue 2, March 2019, Pages 318-329  
  
Regular article  
A Rigorous Interlaboratory Examination of the Need to Confirm Next-Generation Sequencing-Detected Variants with an Orthogonal Method in Clinical Genetic Testing  
Stephen E. Lincoln \* & Rebecca Truty \*, Chiao-Feng Lin †, ‡, Justin M. Zook §, Joshua Paul \*, Vincent H. Ramey \*, Marc Salit §, ¶, Heidi L. Rehm †, ‡, ||, Robert L. Nussbaum \*||, Matthew S. Lebo †, ‡, \*\*, ||  
[Show more](#) [Add to Mendeley](#)

Standardized confirmation protocols are still in their infancy stages



Review | Full Access

## Integrating germline variant assessment into routine clinical practice for myelodysplastic syndrome and acute myeloid leukaemia: current strategies and challenges

Kiran Tawana, Anna L. Brown, Jane E. Churpek

First published: 18 October 2021 | <https://doi.org/10.1111/bjh.17855> | Citations: 1

SECTIONS

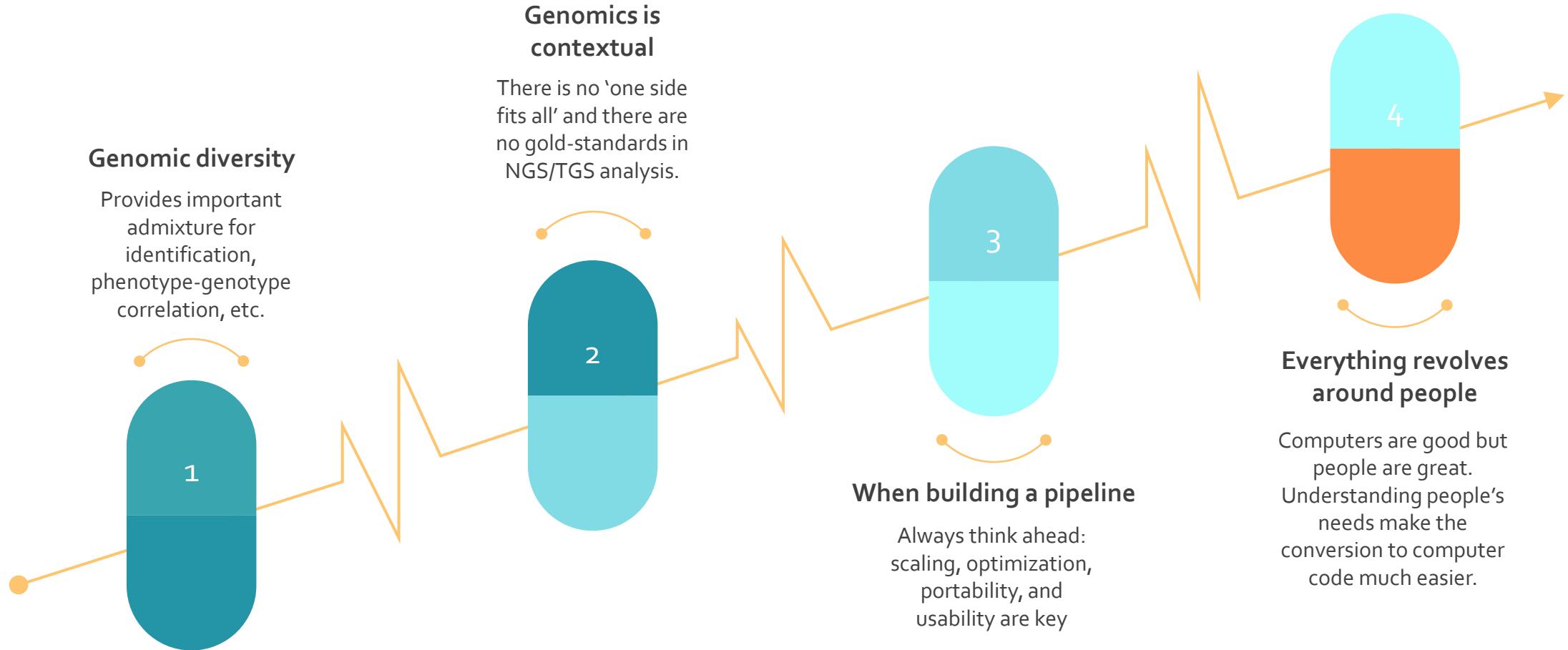
PDF TOOLS SHARE

## Summary

Over the last decade, the field of hereditary haematological malignancy syndromes (HHMSs) has gained increasing recognition among clinicians and scientists worldwide. Germline mutations now account for almost 10% of adult and paediatric myelodysplasia/acute myeloid leukaemia (MDS/AML). As our ability to diagnose HHMSs has improved, we are now faced with the challenges of integrating these advances into routine clinical practice for patients with MDS/AML and how to optimise management and surveillance of patients and asymptomatic carriers. Discoveries of novel syndromes combined with clinical, genetic and epigenetic profiling of tumour samples, have highlighted unique patterns of disease evolution across HHMSs. Despite these advances, causative lesions are detected in less than half of familial cases and evidence-based guidelines are often lacking, suggesting there is much still to learn. Future research efforts are needed to sustain current momentum within the field, led not only by advancing genetic technology but essential collaboration between clinical and academic communities.

Cross-platform variant curation, universal specialist knowledge & metadata limited

# Summarising everything





*"It's one small step for man, one giant leap for mankind."*

- Neil Armstrong



# Skills (Proteomics & transcriptomics)

## ILRI/BecA (CGIAR)

- East Coast fever is a lymphoproliferative disease caused by the tick-borne protozoan parasite *Theileria parva*.
- Sporozoite proteome defined by LC–MS/MS analysis.
- Raw MS/MS data files were analysed with Peaks software (Bioinformatics solutions) using a database containing all Uniprot database entries for *R. appendiculatus* and the re-annotated proteome of *T. parva*
- Proteins were classified according to their putative localization in the sporozoite using TargetP. Trans-membrane domains were predicted by TMHMM Server, signal peptides by SignalP and glycosylphosphatidylinositol (GPI)-anchor signal by GPI-SOM.
- We clustered all sporozoite protein coding sequences of *T. parva* with the complete predicted proteome of *Plasmodium falciparum* into putative orthologous.
- We identified proteins predicted to be orthologs of *Plasmodium falciparum* sporozoite surface molecules and invasion organelle proteins, and proteins that may contribute to the phenomenon of bovine lymphocyte transformation.

James Nyagwange <sup>a, b</sup>, Edwin Tijhaar <sup>b</sup>, Nicola Ternette <sup>c</sup>, Fredrick Mobergi <sup>f</sup>, Kyle Tretina <sup>d</sup>, Joana C. Silva <sup>d, e</sup>, Roger Pelle <sup>a</sup>, Vishvanath Nene <sup>a</sup>  

Show more 

+ Add to Mendeley  Share  Cite 

---

<https://doi.org/10.1016/j.ijpara.2017.09.007>

Under a Creative Commons license

[Get rights and content](#)

 [Open access](#)

## Highlights

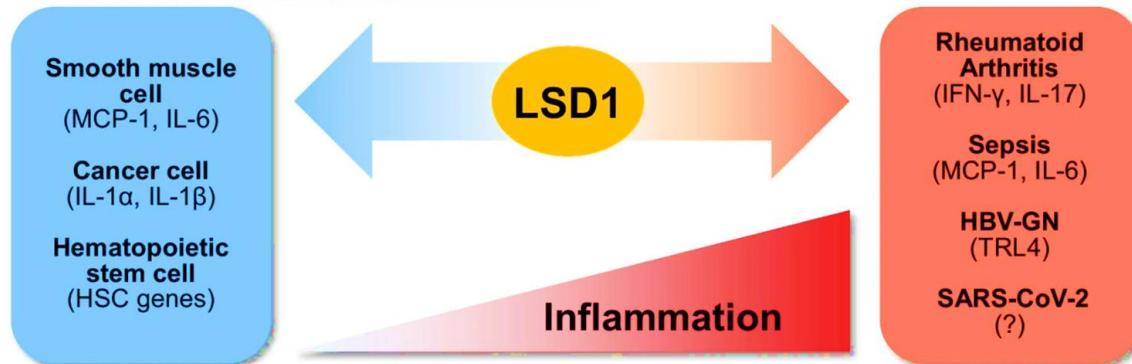
- 2007 *Theileria parva* proteins expressed in the sporozoite were identified.
- Proteins include known *T. parva* antigens targeted by antibodies and cytotoxic T cells.
- Proteins predicted to be orthologs of *Plasmodium falciparum* sporozoite surface molecules were identified.
- Proteins predicted to be orthologs of *P. falciparum* invasion organelle proteins were identified.
- Proteins that may contribute to the phenomenon of bovine lymphocyte transformation were identified.

# Skills (Epigenetics & Transcriptomics)

## Netherlands Cancer Institute & GSK: The dichotomy of LSD1

- Lysine specific demethylase 1 (LSD1) suppresses gene expression through demethylation of H3K4.
- LSD1 acts as a transcriptional coactivator when bound to the androgen receptor (AR) via demethylation of repressive histone marks at H3K9, resulting in the de-repression of AR target genes .

From: [Roles of lysine-specific demethylase 1 \(LSD1\) in homeostasis and diseases](#)



LSD1 functions as a key epigenetic regulator in inflammatory disease. LSD1 controls gene expression in two ways during inflammation. LSD1 increases the expression levels of inflammatory response genes by acting as a positive regulator in inflammatory diseases such as rheumatoid arthritis (RA), sepsis, hepatitis B virus-associated glomerulonephritis (HBV-GN), and SARS-CoV-2 infection. By contrast, LSD1 functions as a negative regulator of the inflammatory response that decreases the expression of cytokine genes in smooth muscle cells (SMCs), cancer cells, and hematopoietic stem cells (HSCs).

[Haematologica](#). 2019 Jun; 104(6): 1156–1167.

Prepublished online 2018 Dec 4. doi: [10.3324/haematol.2018.199190](https://doi.org/10.3324/haematol.2018.199190)

PMCID: PMC6545850

PMID: [30514804](#)

Lysine specific demethylase 1 inactivation enhances differentiation and promotes cytotoxic response when combined with all-trans retinoic acid in acute myeloid leukemia across subtypes

Kimberly N. Smitheman,<sup>1</sup> Tesa M. Severson,<sup>2</sup> Satyajit R. Rajapurkar,<sup>1</sup> Michael T. McCabe,<sup>1</sup> Natalie Karpinich,<sup>1</sup> James Foley,<sup>1</sup> Melissa B. Pappalardi,<sup>1</sup> Ashley Hughes,<sup>3</sup> Wendy Halsey,<sup>3</sup> Elizabeth Thomas,<sup>3</sup> Christopher Traini,<sup>3</sup> Kelly E. Federowicz,<sup>1</sup> Jenny Laraio,<sup>1</sup> Fredrick Mobegi,<sup>2</sup> Geraldine Ferron-Brady,<sup>4</sup> Rabinder K. Prinjha,<sup>1</sup> Christopher L. Carpenter,<sup>1</sup> Ryan G. Kruger,<sup>1</sup> Lodewyk Wessels,<sup>2,5</sup> and Helai P. Mohammad<sup>1</sup>

► Author information ► Article notes ► Copyright and License information ► [Disclaimer](#)

- LSD1 is a known regulator of normal haematopoiesis.
- LSD1 inhibition promotes differentiation of leukemic cells.
- Therapeutic potential of a selective, potent inhibitor of LSD1, GSK2879552, in combination with all-trans retinoic acid in relapsed refractory AML patients was explored.
  - RNA-seq (Active treatments vs DMSO) :: DESeq2
  - ChIP-seq (H3K4me1, H3K4me2, H3K4me3, H3K9ac, H3K27ac vs DMSO) :: Peaks call (MACS2); differential analysis (DESeq2); nearest-gene + peak intersection (BEDTools).

# Skills (WGS; Metagenomics; Comparative)



**Journal of Cystic Fibrosis**  
Volume 20, Issue 3, May 2021, Pages 413-420

The cystic fibrosis gut as a potential source of multidrug resistant pathogens

Steven L. Taylor <sup>a, b, 1</sup>, Lex E.X. Leong <sup>c, 1</sup>, Sarah K. Sims <sup>a, b</sup>, Rebecca L. Keating <sup>d</sup>, Lito E. Papanicolas <sup>a, b</sup>, Alyson Richard <sup>a, b</sup>, Fredrick M. Mobegi <sup>a, b</sup>, Steve Wesselingh <sup>b</sup>, Lucy D. Burr <sup>d, e, 2</sup>, Geraint B. Rogers <sup>a, b, 2</sup>

**GUT MICROBES**  
2019, VOL. 10, NO. 3, 367–381  
<https://doi.org/10.1080/19490976.2018.1534512>

**RESEARCH PAPER/REPORT** **OPEN ACCESS** 

**Opportunistic bacteria confer the ability to ferment prebiotic starch in the adult cystic fibrosis gut**

Yanan Wang <sup>a,b</sup>, Lex E.X. Leong <sup>a,b</sup>, Rebecca L. Keating <sup>c</sup>, Tokuwa Kanno <sup>d</sup>, Guy C.J. Abell <sup>a</sup>, Fredrick M. Mobegi <sup>a,b</sup>, Jocelyn M. Choo <sup>a,b</sup>, Steve L. Wesselingh <sup>a</sup>, A. James Mason <sup>a,d</sup>, Lucy D. Burr <sup>a,c,e,f</sup>, and Geraint B. Rogers <sup>a,b</sup>

<sup>a</sup>Infection and Immunity Theme, South Australia Health and Medical Research Institute, Adelaide, Australia; <sup>b</sup>SAHMRI Microbiome Research Laboratory, Flinders University School of Medicine, Adelaide, Australia; <sup>c</sup>Department of Respiratory Medicine, Mater Health Services, South Brisbane, Australia; <sup>d</sup>King's College London, Institute of Pharmaceutical Science, London, UK; <sup>e</sup>Mater Research, University of Queensland, South Brisbane, Australia



**RESEARCH ARTICLE**  

**Intestinal microbiology shapes population health impacts of diet and lifestyle risk exposures in Torres Strait Islander communities**

Fredrick M Mobegi<sup>1,2</sup>, Lex EX Leong<sup>1</sup>, Fintan Thompson<sup>1,3</sup>, Sean M Taylor<sup>3</sup>, Linton R Harriss<sup>3</sup>, Jocelyn M Choo<sup>1,2</sup>, Steven L Taylor<sup>1,2</sup>, Steve L Wesselingh<sup>4</sup>, Robyn McDermott<sup>3,5</sup>, Kerry L Ivey<sup>1,6,7†\*</sup>, Geraint B Rogers<sup>1,2†\*</sup>

**Mechanisms linking low-calorie sweeteners to impaired glycaemic control**

D Kreuch, K Ivey, F M Mobegi, L Leong, N J Isaacs, N Pezos, M Horowitz, C K Rayner, G B Rogers, R L Young

Microbiome And Host Health, Intestinal Nutrient Sensing

Research output: Contribution to journal > Article > peer-review

**Long-Term Azithromycin Reduces *Haemophilus influenzae* and Increases Antibiotic Resistance in Severe Asthma**

Steven L. Taylor <sup>1,2</sup>, Lex E. X. Leong <sup>1,2</sup>, Fredrick M. Mobegi <sup>1,2</sup>, Jocelyn M. Choo <sup>1,2</sup>, Steve Wesselingh <sup>1,2</sup>, Ian A. Yang <sup>3,4</sup>, John W. Upham <sup>3,5</sup>, Paul N. Reynolds <sup>6,7</sup>, Sandra Hodge <sup>6,7</sup>, Alan L. James <sup>8,9</sup>, Christine Jenkins <sup>10,11</sup>, Show All...  
+ Author Affiliations

 58  5,216

<https://doi.org/10.1164/rccm.201809-1739OC> PubMed: 30875247

Received: September 22, 2018 Accepted: March 13, 2019

# Skills (WGS & Transcriptomics, NGS, TGS)

## Curtin University: Characterization fungal isolates

- Comprehensive characterization of the biology and survival of fungal pathogens is a prerequisite to developing more effective disease management strategies.
- At the CCDM, I lead the bioinformatics analyses (in two projects) focusing on Ascochyta blight of pulses.
- We employ whole-genome sequencing using nanopore, PacBio, and Illumina technologies, RNA-seq, and targeted genotyping to study Ascochyta interactions with the hosts and identify novel mechanisms of resistance and pathogenicity.
- These datasets provide unmatched resolution at a molecular level into various mechanisms of pathogenicity and resistance, allowing us to identify genetic strengths and vulnerabilities that can be exploited in breeding programs.

Research | Open Access | Published: 08 May 2021

Analysis of differentially expressed *Sclerotinia sclerotiorum* genes during the interaction with moderately resistant and highly susceptible chickpea lines

[Virginia W. Mwape](#)✉, [Fredrick M. Mobegi](#), [Roshan Regmi](#), [Toby E. Newman](#), [Lars G. Kamphuis](#)✉ & [Mark C. Derbyshire](#)

*BMC Genomics* **22**, Article number: 333 (2021) | [Cite this article](#)

**1063** Accesses | **1** Citations | **5** Altmetric | [Metrics](#)

# Skills (Molecular epidemiology)

**RESEARCH** **Open Access** 

## Investigating potential transmission of antimicrobial resistance in an open-plan hospital ward: a cross-sectional metagenomic study of resistome dispersion in a lower middle-income setting

Anushia Ashokan<sup>1,2,3</sup>, Josh Hanson<sup>4,5</sup>, Ne Myo Aung<sup>5,6</sup>, Mar Mar Kyi<sup>5,6</sup>, Steven L. Taylor<sup>1,2</sup>, Jocelyn M. Choo<sup>1,2</sup>, Erin Flynn<sup>1,2</sup>, Fredrick Mobegi<sup>1,2</sup>, Morgyn S. Warner<sup>3,7</sup>, Steve L. Wesselingh<sup>8</sup>, Mark A. Boyd<sup>3</sup> and Geraint B. Rogers<sup>1,2\*</sup>

**OPEN**

## The post-vaccine microevolution of invasive *Streptococcus pneumoniae*

Amelieke J. H. Cremers<sup>1,\*</sup>, Fredrick M. Mobegi<sup>1,2,\*</sup>, Marien I. de Jonge<sup>1</sup>, Sacha A. F. T. van Huijum<sup>2</sup>, Jacques F. Meis<sup>3,4</sup>, Peter W. M. Hermans<sup>1,7</sup>, Gerben Ferwerda<sup>1</sup>, Stephen D. Bentley<sup>5</sup> & Aldert L. Zomer<sup>1,2</sup>

Received: 29 June 2015  
Accepted: 10 September 2015  
Published: 23 October 2015

The 7-valent pneumococcal conjugated vaccine (PCV7) has affected the genetic population of *Streptococcus pneumoniae* in pediatric carriage. Little is known however about pneumococcal population genomics in adult invasive pneumococcal disease (IPD) under vaccine pressure. We sequenced and serotyped 349 strains of *S. pneumoniae* isolated from IPD patients in Nijmegen between 2001 and 2011. Introduction of PCV7 in the Dutch National Immunization Program in 2006 precluded substantial alterations in the IPD population structure caused by serotype replacement. No

**OPEN**

## Deciphering the distance to antibiotic resistance for the pneumococcus using genome sequencing data

Received: 12 October 2016  
Accepted: 16 January 2017  
Published: 16 February 2017

Fredrick M. Mobegi<sup>1,2,3</sup>, Amelieke J. H. Cremers<sup>1</sup>, Marien I. de Jonge<sup>1</sup>, Stephen D. Bentley<sup>4</sup>, Sacha A. F. T. van Huijum<sup>2</sup> & Aldert Zomer<sup>1,2,5</sup>

# Skills (Genotype-phenotype association)

## Phage-Derived Protein Induces Increased Platelet Activation and Is Associated with Mortality in Patients with Invasive Pneumococcal Disease

Rahajeng N. Tunjungputri,<sup>a,b</sup> Fredrick M. Mobegi,<sup>c</sup> Amelieke J. Cremers,<sup>c,d</sup> Christa E. van der Gaast-de Jongh,<sup>c</sup> Gerben Ferwerda,<sup>c</sup> Jacques F. Meis,<sup>d,e</sup> Nel Roeleveld,<sup>f,g</sup> Stephen D. Bentley,<sup>h</sup> Alexander S. Pastura,<sup>c</sup> Sacha A. F. T. van Hijum,<sup>i</sup> Andre J. van der Ven,<sup>a</sup> Quirijn de Mast,<sup>a</sup> Aldert Zomer,<sup>i,j</sup> Marien I. de Jonge<sup>c</sup>

Department of Internal Medicine, Radboud Institute for Molecular Life Sciences, Radboud University Medical Center, Nijmegen, The Netherlands<sup>a</sup>; Center for Tropical and Infectious Diseases (CENTRID), Faculty of

## The Contribution of Genetic Variation of *Streptococcus pneumoniae* to the Clinical Manifestation of Invasive Pneumococcal Disease

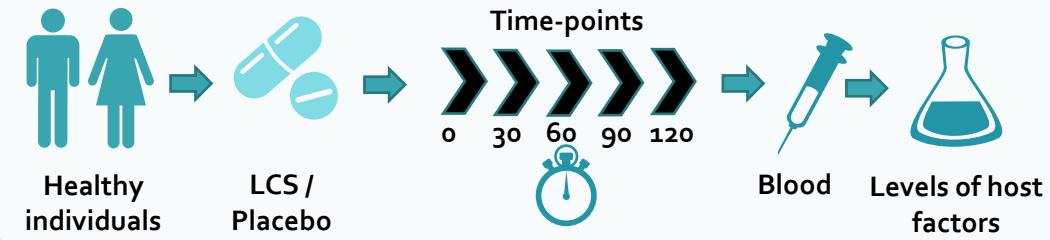
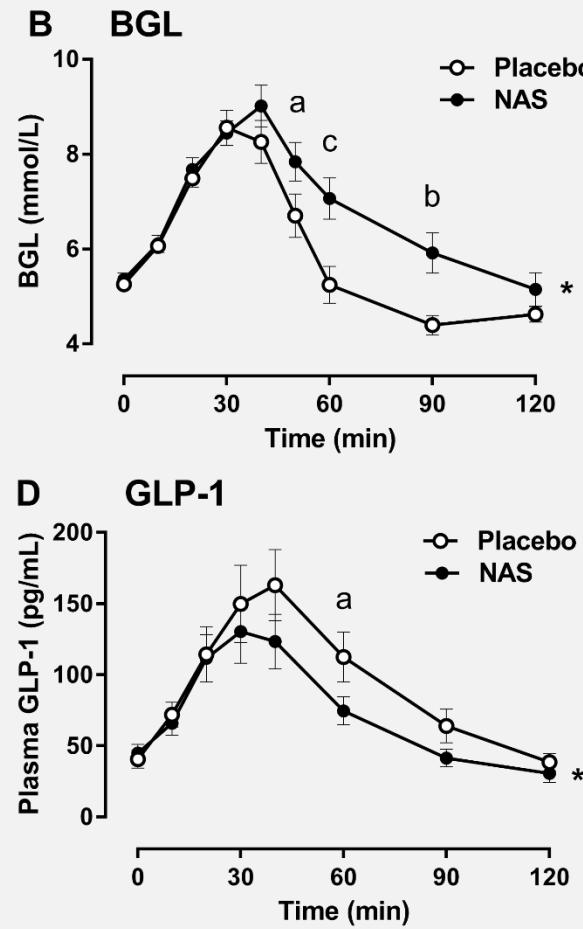
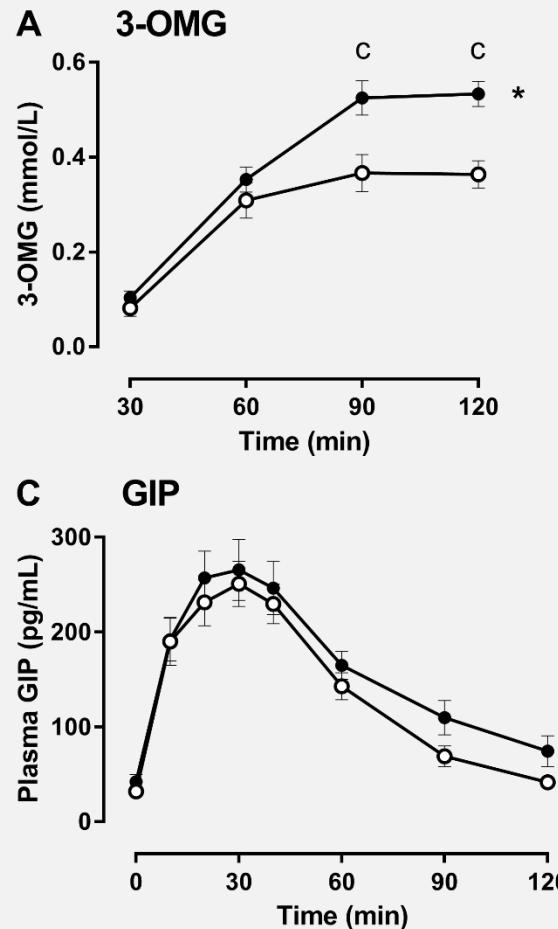
Amelieke J. H. Cremers,<sup>1,2,3</sup> Fredrick M. Mobegi,<sup>1,2,4</sup> Christa van der Gaast-de Jongh,<sup>1,2</sup> Michelle van Weert,<sup>1,2</sup> Fred J. van Opzeeland,<sup>1,2</sup> Minna Vehkala,<sup>5</sup> Mirjam J. Knol,<sup>6</sup> Hester J. Bootsma,<sup>6</sup> Niko Välimäki,<sup>5</sup> Nicholas J. Croucher,<sup>7</sup> Jacques F. Meis,<sup>8</sup> Stephen Bentley,<sup>9</sup> Sacha A. F. T. van Hijum,<sup>2,4,10</sup> Jukka Corander,<sup>5,9,11</sup> Aldert L. Zomer,<sup>12</sup> Gerben Ferwerda,<sup>12</sup> and Marien I. de Jonge<sup>1,2</sup>

<sup>1</sup>Section of Pediatric Infectious Diseases, Laboratory of Medical Immunology, Radboud Institute for Molecular Life Sciences, <sup>2</sup>Radboud Center for Infectious Diseases, <sup>3</sup>Department of Medical Microbiology, and <sup>4</sup>Bacterial Genomics Group, Center for Molecular and Biomolecular Informatics, Radboudumc, Nijmegen, The Netherlands; <sup>5</sup>Department of Mathematics and Statistics, University of Helsinki, Finland; <sup>6</sup>Centre for Infectious Disease Control, National Institute for Public Health and the Environment (RIVM), Bilthoven, The Netherlands; <sup>7</sup>Medical Research Council Centre for Outbreak Analysis and Modelling, Department of Infectious Disease Epidemiology, Imperial College London, United Kingdom; <sup>8</sup>Department of Medical Microbiology and Infectious Diseases, Canisius-Wilhelmina Hospital, Nijmegen, The Netherlands; <sup>9</sup>Wellcome Trust Sanger Institute, Pathogen Genomics Group, Hinxton, Cambridge, United Kingdom; <sup>10</sup>NIZO, Ede, The Netherlands; <sup>11</sup>Department of Biostatistics, University of Oslo, Norway; and <sup>12</sup>Department of Infectious Diseases and Immunology, Faculty of Veterinary Medicine, Utrecht University, The Netherlands

## From microbial gene essentiality to novel antimicrobial drug targets

Fredrick M Mobegi<sup>1,2</sup>, Sacha AFT van Hijum<sup>2,3\*</sup>, Peter Burghout<sup>1</sup>, Hester J Bootsma<sup>1</sup>, Stefan PW de Vries<sup>1,4</sup>, Christa E van der Gaast-de Jongh<sup>1</sup>, Elles Simonetti<sup>1</sup>, Jeroen D Langereis<sup>1</sup>, Peter WM Hermans<sup>1,5</sup>, Marien I de Jonge<sup>1</sup> and Aldert Zomer<sup>1,2\*</sup>

# Artificial sweeteners: the microbiome perspective



## Contribution of LCS to dysglycaemia.

- Regular high intake of beverages sweetened with **low-calorie sweeteners (LCS)** increase the risk of developing type 2 diabetes mellitus (T2DM)
- Underlying mechanisms remain unknown.
- 2-week intervention (randomised double-blind experiment) supplementation with a LCS combination (92 mg sucralose + 52 mg acesulfame-K, equivalent to ~1.5L of diet beverage consumption/day) or placebo

## RBG and Insulin

Intake of LCS caused increased blood sugar and insulin levels in healthy non-diabetic subjects.

## 3-OMG

Intake of LCS caused augmented glucose absorption (3-OMG) in healthy non-diabetic subjects.

## GLP-1 and GLP-2

Intake of LCS caused attenuated release of glucagon-like peptide-1 (GLP-1) levels in healthy non-diabetic subjects.

Does the gut microbiota play a role for LCS effects in humans just like in rodents?





# RNA sequencing data analysis

---

Fredrick M. Mobegi, PhD  
BHKi hybrid seminar series



[linkedin.com/fmobegi](https://linkedin.com/fmobegi)



[twitter.com/mobeginomics](https://twitter.com/mobeginomics)



[github.com/fmobegi](https://github.com/fmobegi)

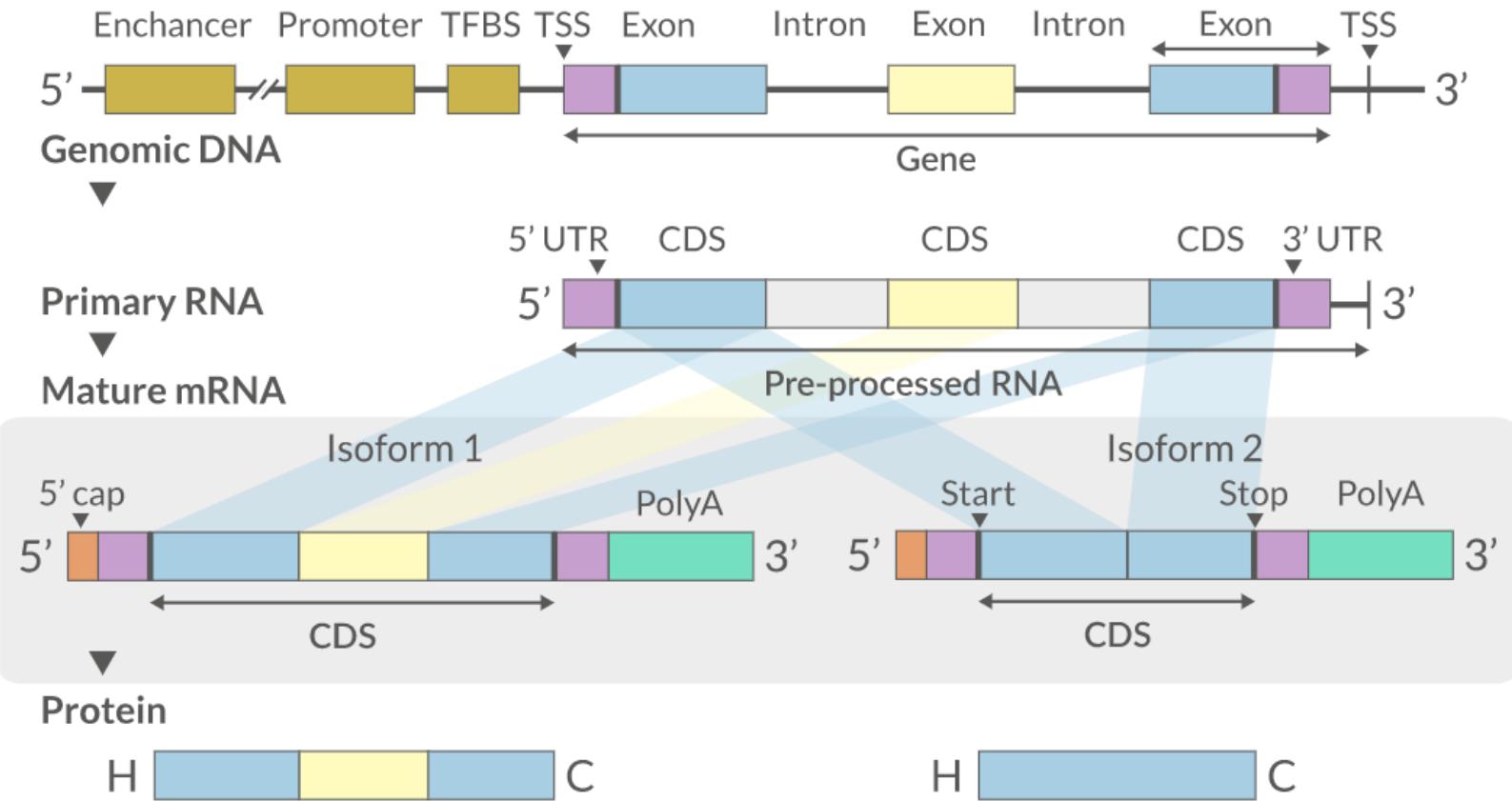


[rpubs.com/fmobegi](https://rpubs.com/fmobegi)

- RNA sequencing
- Workflow
- DGE analysis pipeline/workflow
- Read Quality assessment
- Read mapping
- Alignment Quality Check
- Quantification of transcripts
- Exploratory analyses
- Differential Gene Expression analysis
- Functional analyses
- Summary

## Contents

# RNA sequencing



- The transcriptome is spatially and temporally dynamic
- Data comes from functional units (coding regions)
- Only a tiny fraction of the genome

Image credits: scilifelab

# RNA-seq experiment design

## Scotty: a web tool for designing RNA-Seq experiments to measure differential gene expression FREE

Michele A. Busby, Chip Stewart, Chase A. Miller, Krzysztof R. Grzeda, Gabor T. Marth ✉

[Author Notes](#)

*Bioinformatics*, Volume 29, Issue 5, 1 March 2013, Pages 656–657,

<https://doi.org/10.1093/bioinformatics/btt015>

Published: 12 January 2013 Article history ▾

[PDF](#) [Split View](#) [Cite](#) [Permissions](#) [Share ▾](#)

### Abstract

**Motivation:** A common question arises at the beginning of every experiment where RNA-Seq is used to detect differential gene expression between two conditions: How many reads should we sequence?

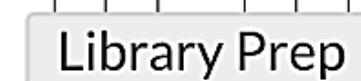
- Balanced experiment
- Biological replicates
- Statistical power
- Technical replicates

Scotty performs Power analysis with cost estimates

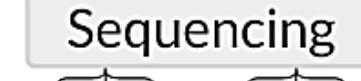
### Confounding



### RNA extraction



### Library Prep



### Sequencing



### Balanced

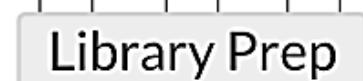
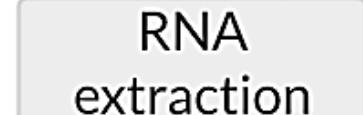
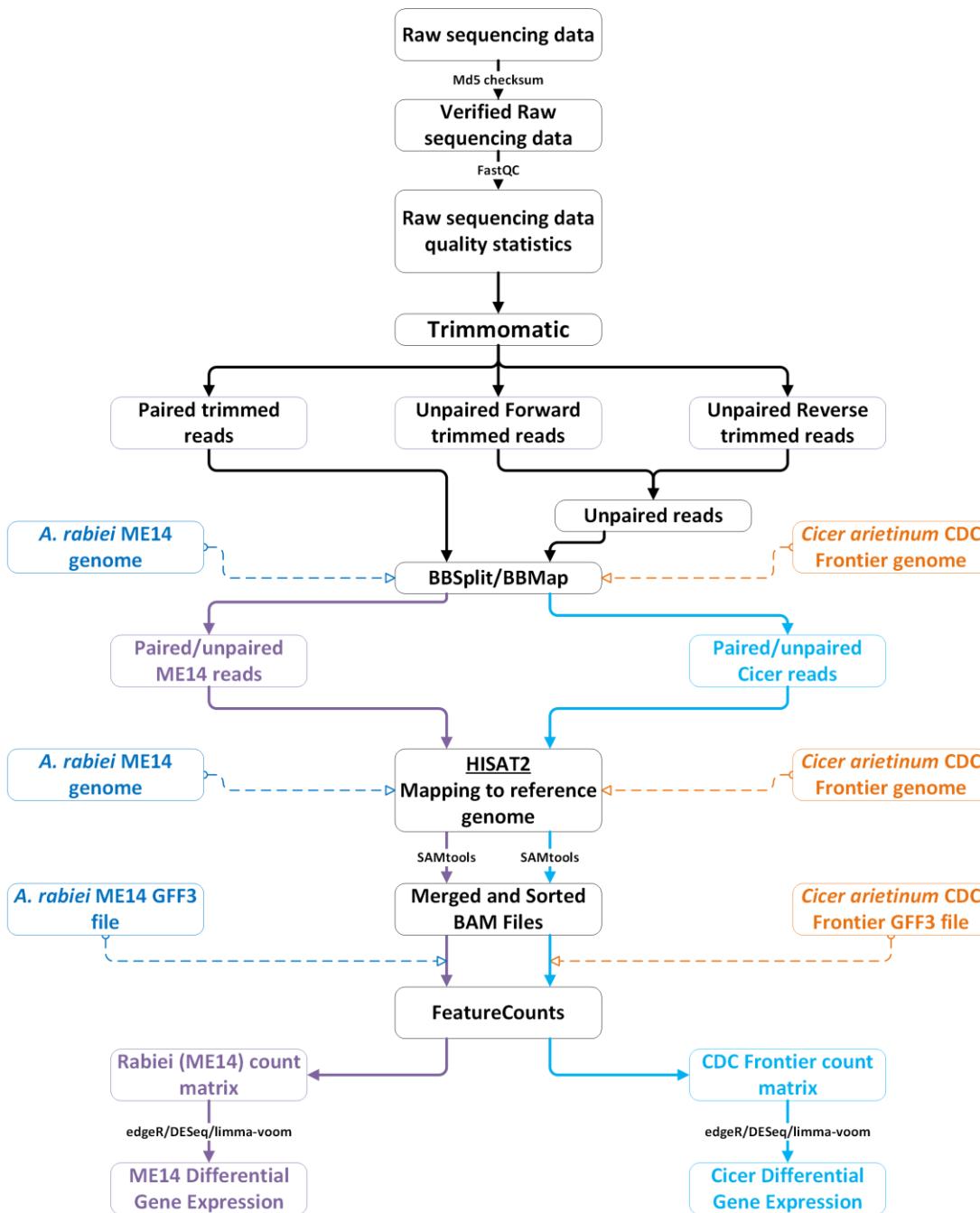


Image credits: scilifelab

# RNA-seq data analysis pipeline



# RNA-seq data analysis pipeline

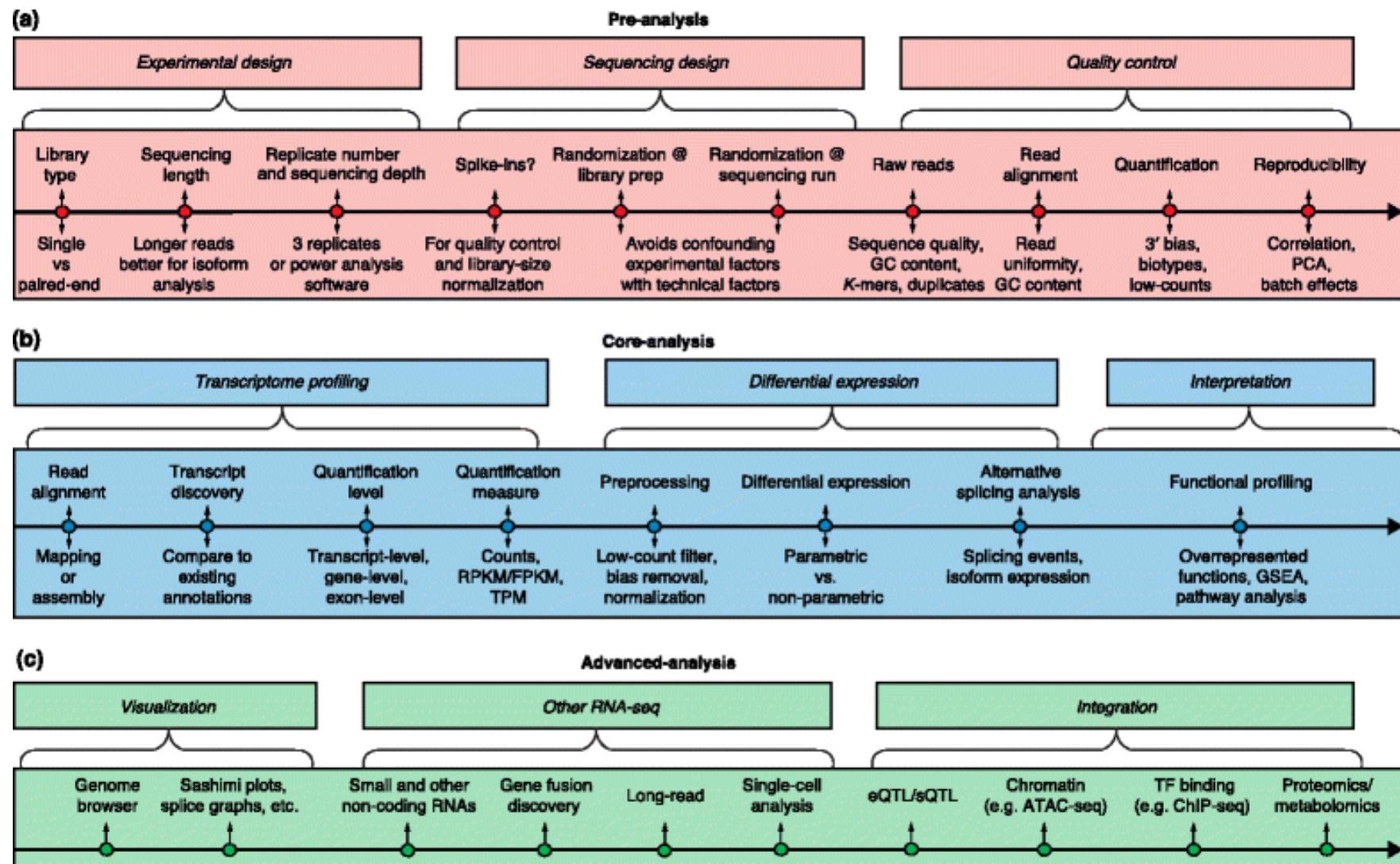
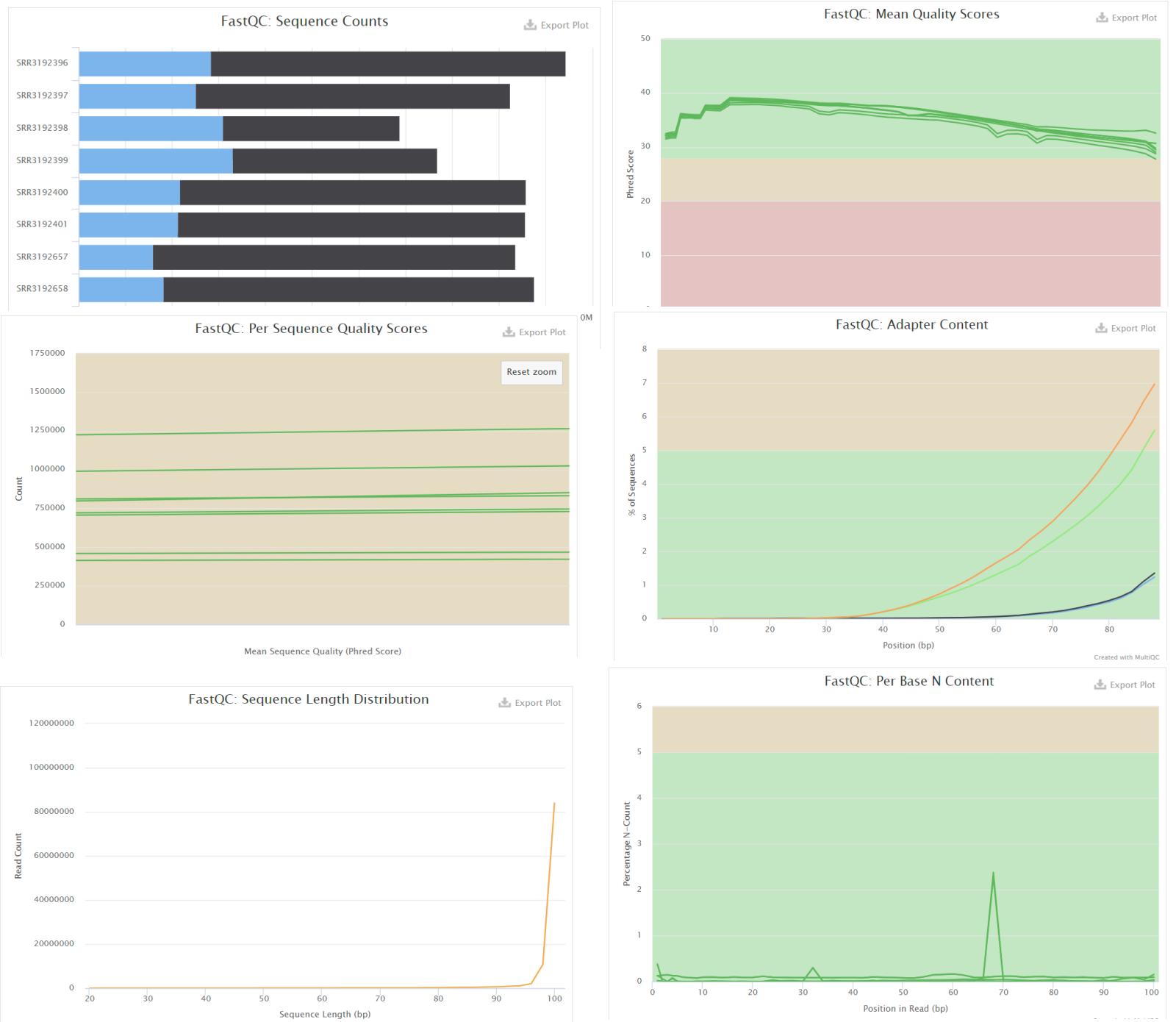


Figure 1: Conesa *et al.* [A survey of best practices for RNA-seq data analysis](#). Genome Biology volume 17, Article number: 13 (2016)

# Quality assessment

- MultiQC
- FastQC
- Trimmomatic
- CutAdapt



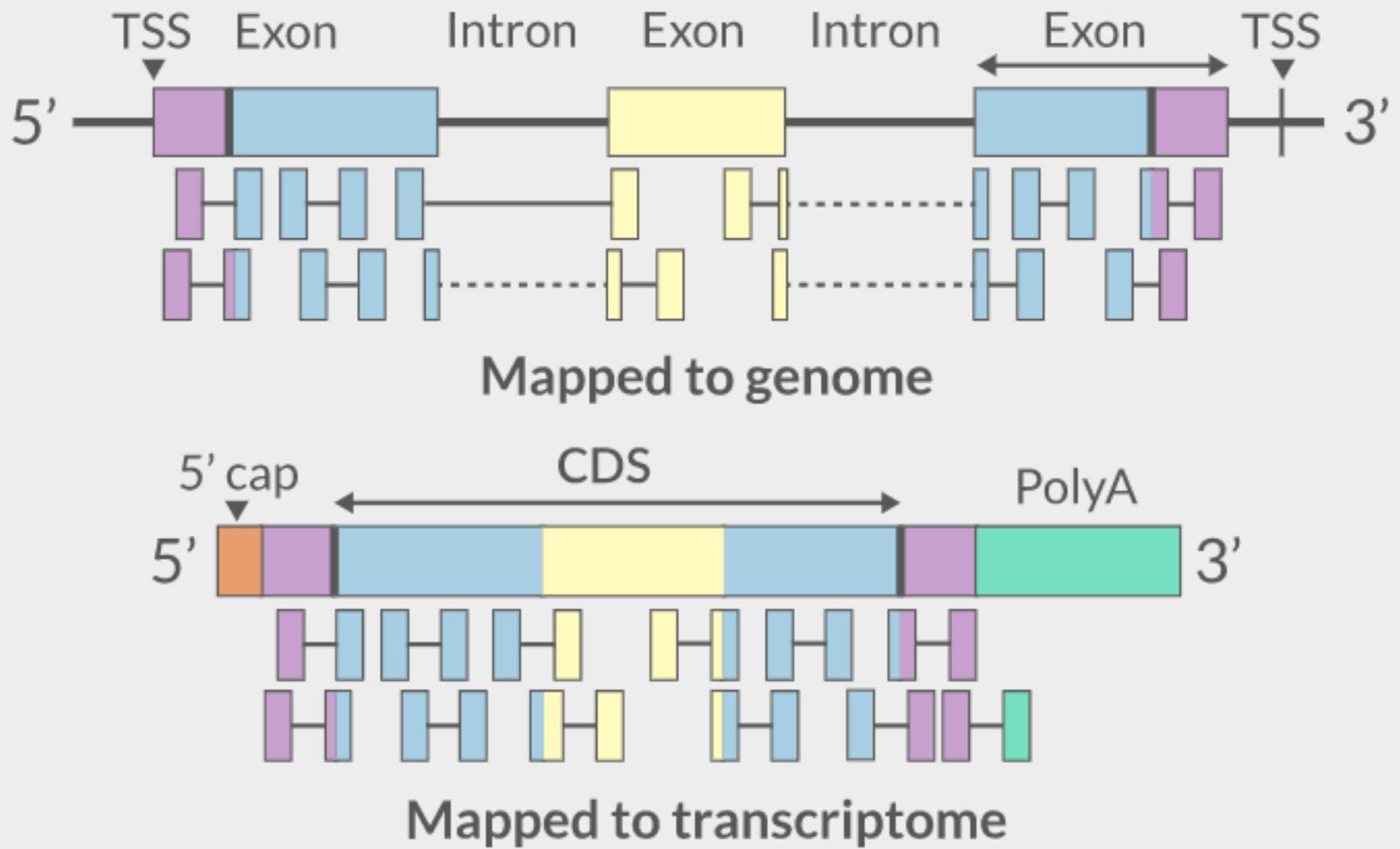
# Quality assessment (FastQC and MultiQC)

The screenshot shows the MultiQC software interface. At the top, there is a navigation bar with tabs: MultiQC Example Reports, RNA-Seq, Whole-Genome Seq, Bisulfite Seq, Hi-C, and MultiQC\_NGI. The RNA-Seq tab is selected. On the left, a sidebar lists various QC modules: General Stats, Quality Histograms, Sequence Quality Scores, Per Base Sequence Content, Per Sequence GC Content, Per Base N Content, Sequence Length Distribution, Sequence Duplication Levels, and Overrepresented sequences. The main content area features the MultiQC logo and a brief description: "A modular tool to aggregate results from bioinformatics analyses across many samples into a single report." Below this, a message says "Report generated on 2022-02-08, 23:02 based on data in: /Users/phil/GitHub/MultiQC/website/public\_html/examples/rna-seq". A welcome message at the bottom left says "Welcome! Not sure where to start? Watch a tutorial video (6:06)" with a "don't show again" button. The central part of the screen displays a table titled "General Statistics" with data for five samples (SRR3192396, SRR3192397, SRR3192398, SRR3192399, SRR3192400). The table includes columns for Sample Name, % Assigned, M Assigned, % Aligned, M Aligned, % BP Trimmed, % Dups, % GC, and M Seq. The data is presented in a grid format with color-coded cells.

Sample Name	% Assigned	M Assigned	% Aligned	M Aligned	% BP Trimmed	% Dups	% GC	M Seq
SRR3192396	67.5%	71.9	93.7%	97.8	4.0%	72.8%	50%	104.4
SRR3192397	66.6%	63.0	94.7%	87.1	3.5%	72.8%	48%	92.5
SRR3192398	50.9%	36.5	88.2%	58.7	5.0%	55.0%	47%	68.8
SRR3192399	52.3%	42.3	88.2%	65.6	5.0%	57.1%	47%	76.8
SRR3192400	70.3%	63.4	77.3%	73.4	7.2%	77.3%	45%	95.8

## Mapping and splice- aware alignment

- STAR
- HiSat2
- GSNAP
- SOAPsplice
- TopHat2
- Novoalign™
- CLC™



- Some tools like Kallisto can quantify reads directly without a reference genome
- Viable aligners must (i) align reads across splice junctions, (ii) handle paired-end reads, (iii) handle strand-specific data, and (iv) run efficiently.
- Ability to align reads across unannotated splice junctions is also a plus

# Sequence Alignment Map (SAM) QC

- samtools stats
- bamtools stats
- QoRTs
- RSeQC
- Qualimap

## Size matters

Format	Size_GB
SAM	7.4
BAM	1.9
CRAM lossless Q	1.4
CRAM 8 bins Q	0.8
CRAM no Q	0.26

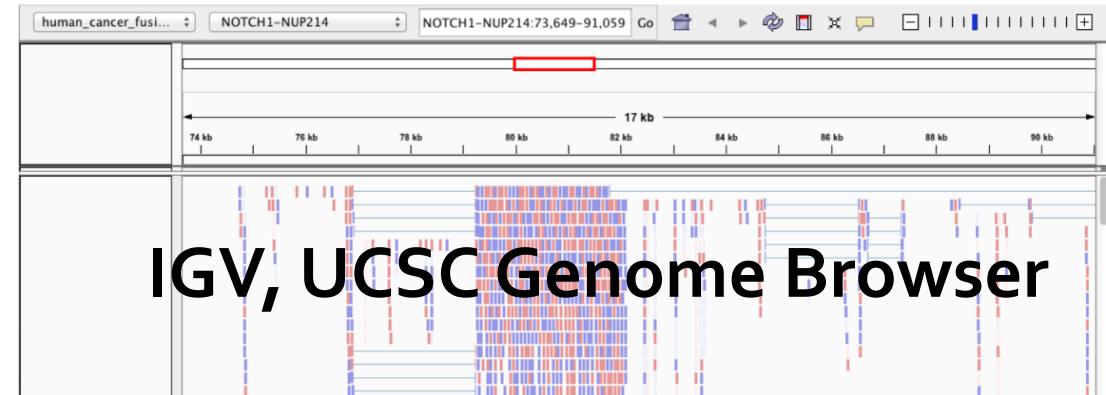
- # of reads mapped/unmapped
- Uniquely mapped
- Insert size distribution
- Coverage + Gene body coverage
- Biotype counts / Chromosome counts
- Counts by region: gene/intron
- Sequencing saturation
- Strand specificity

## Looks matter too..

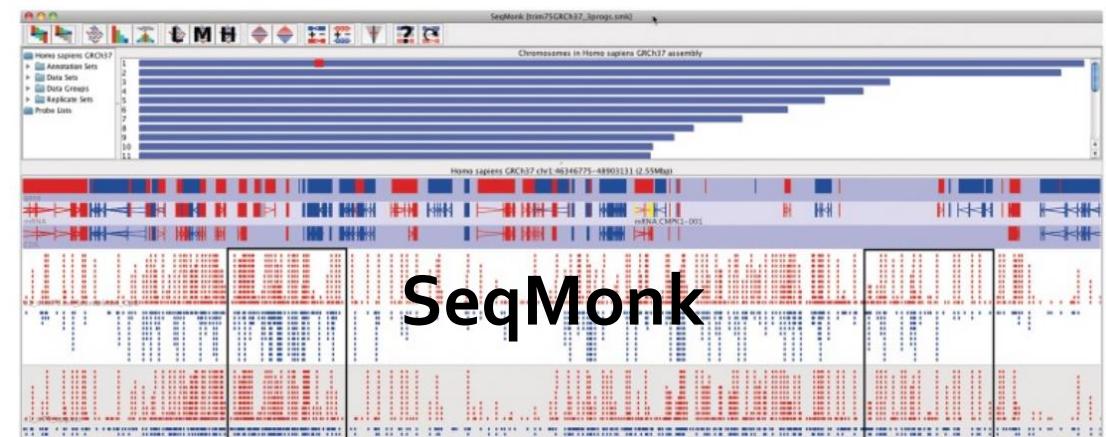
Samtools tview alignment.bam genome.fasta

The screenshot shows a command-line interface for Samtools tview. The command is "samtools tview alignment.bam genome.fasta". The output displays a grid of sequence alignments. The columns represent genomic positions from 911 to 1851, and the rows represent individual reads. Each cell contains the sequence of a read aligned against a specific genomic region. The alignments show various patterns of matches and mismatches.

## SAMTOOLS



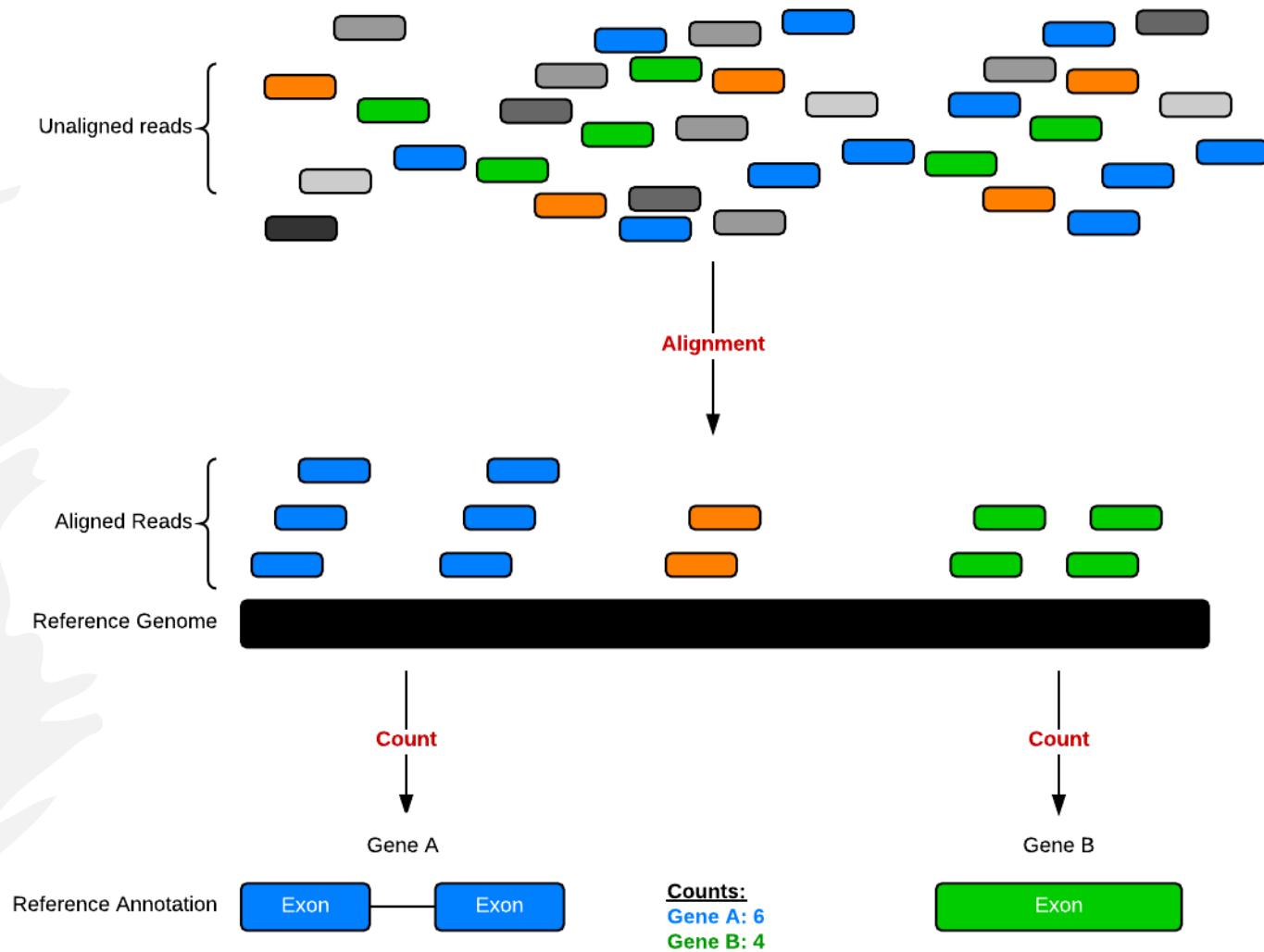
## IGV, UCSC Genome Browser



## SeqMonk

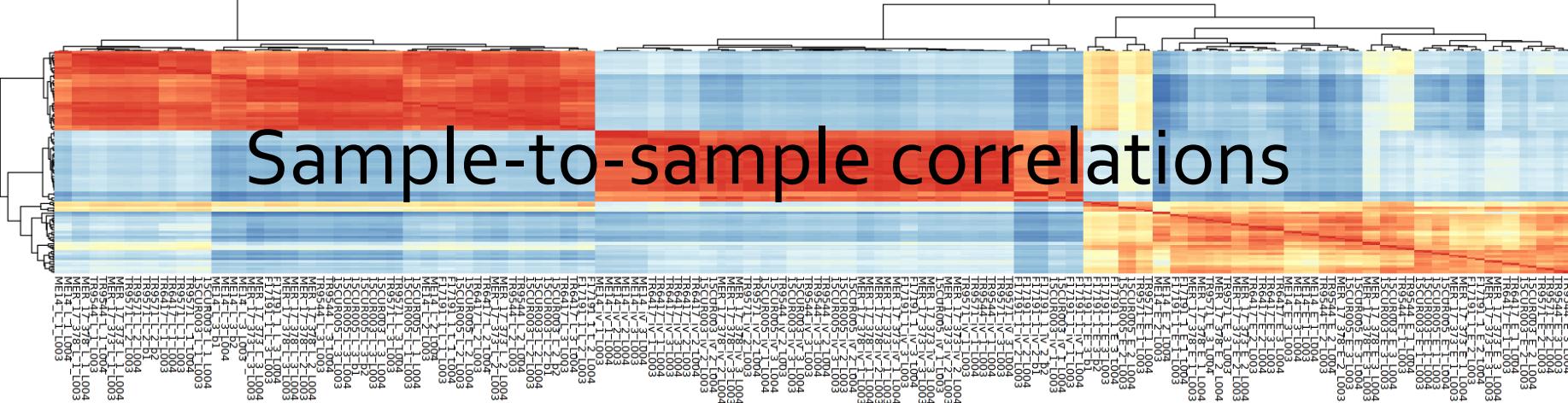
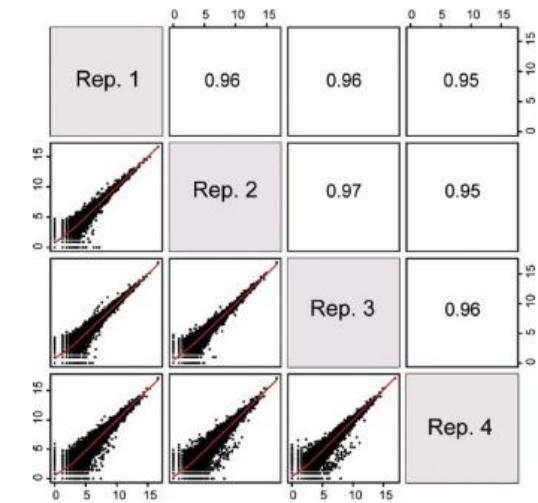
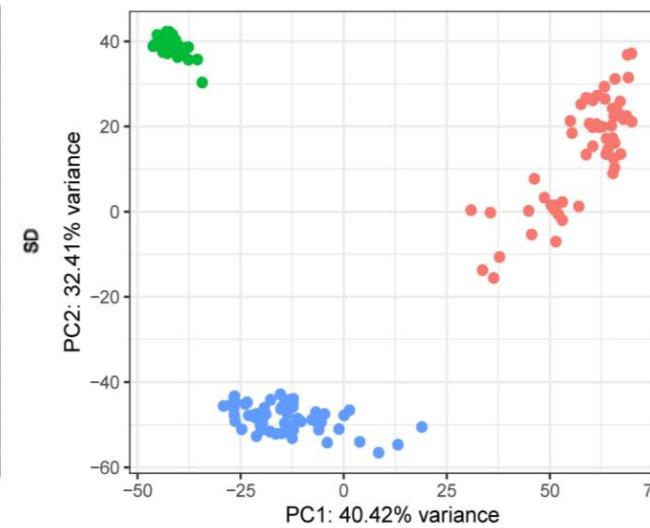
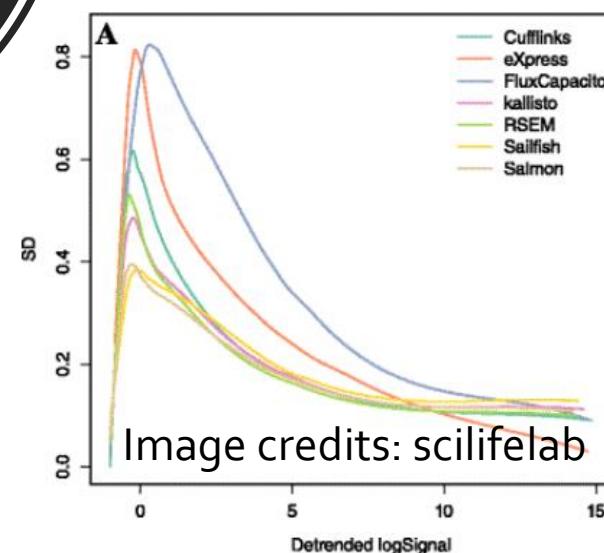
# Quantifying counts

- Htseq, StringTie, RSEM, Cufflinks, and Trinity are some of the popular tools.
- Reference annotation tells counting software the assembly position/location of features to quantify.
- Reads can be quantified on any feature in the reference annotation (gene, transcript, exon).
- Gene and transcript level are most common.



- RNASeqComp
- MultiQC (reports)
- R data wrangling

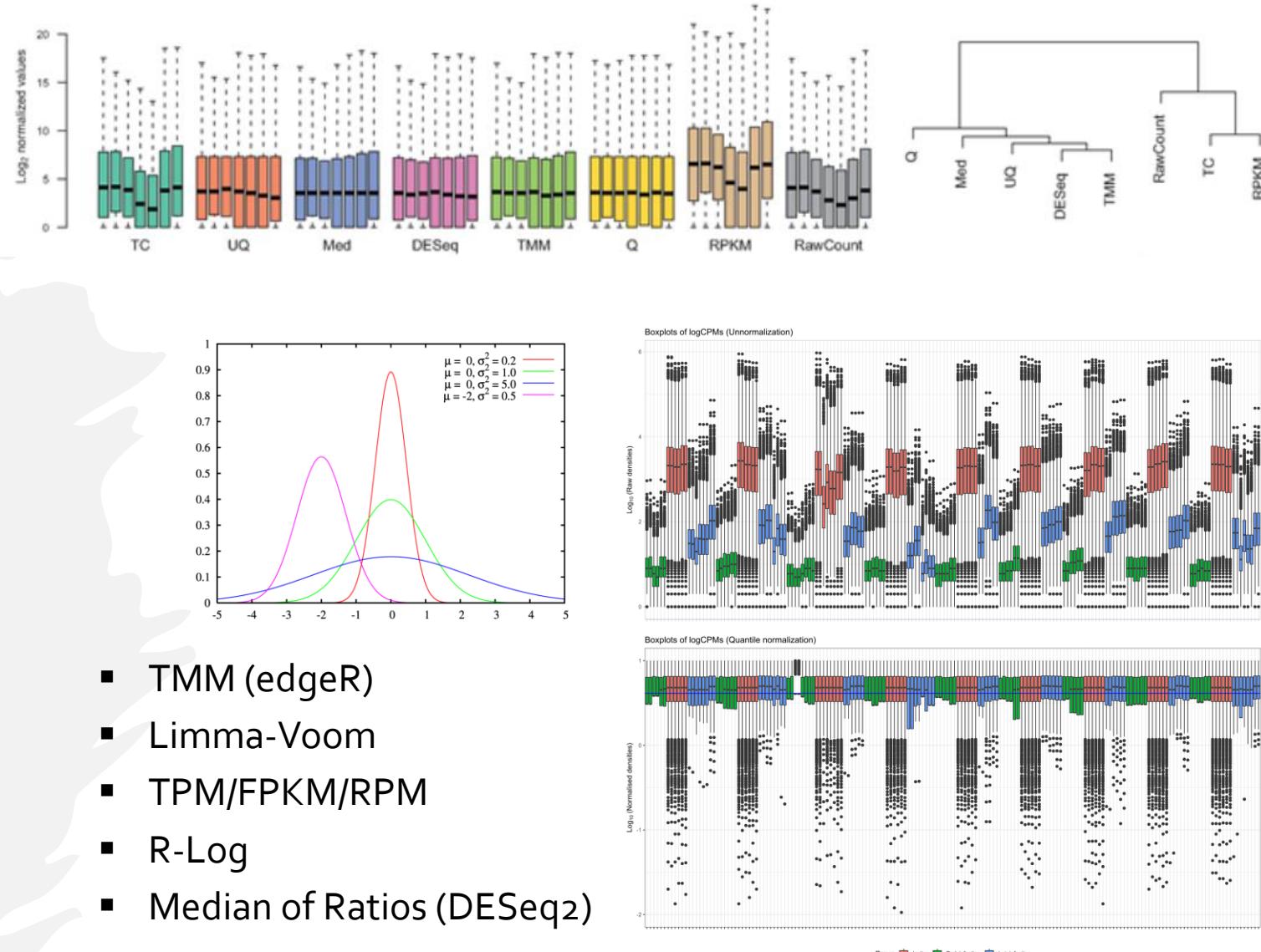
## QC on Counts



ENSG000000000003	140	242	188	143	287	344	438	280	253
ENSG000000000005	0	0	0	0	0	0	0	0	0
ENSG000000000419	69	98	77	55	52	94	116	79	69
ENSG000000000457	56	75	104	79	157	205	183	178	153

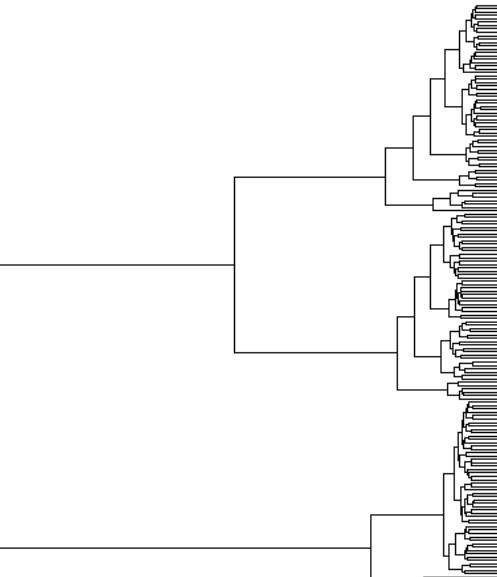
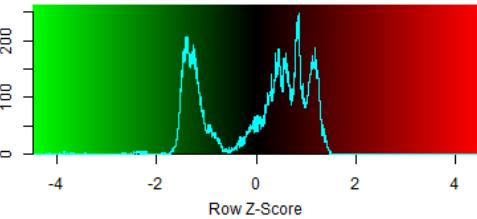
# Normalizing counts

- Normalization controls for compositional bias and sequencing/library depth.
- Own analysis and plots: use normalized counts like TPM/FPKM.
- Heatmaps and clustering: use VST, R-log, Voom.
- Differential Gene Expression analysis: use raw counts.

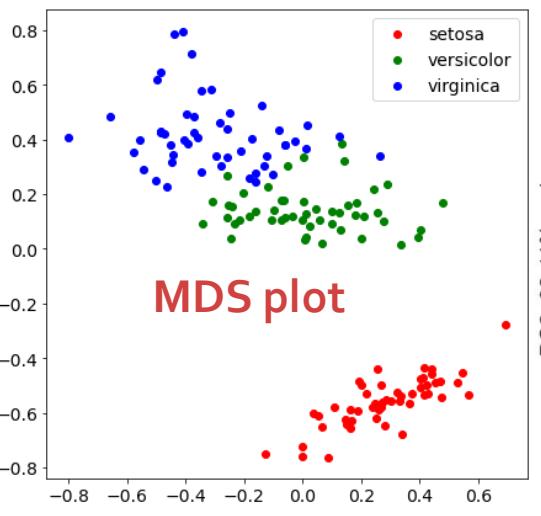
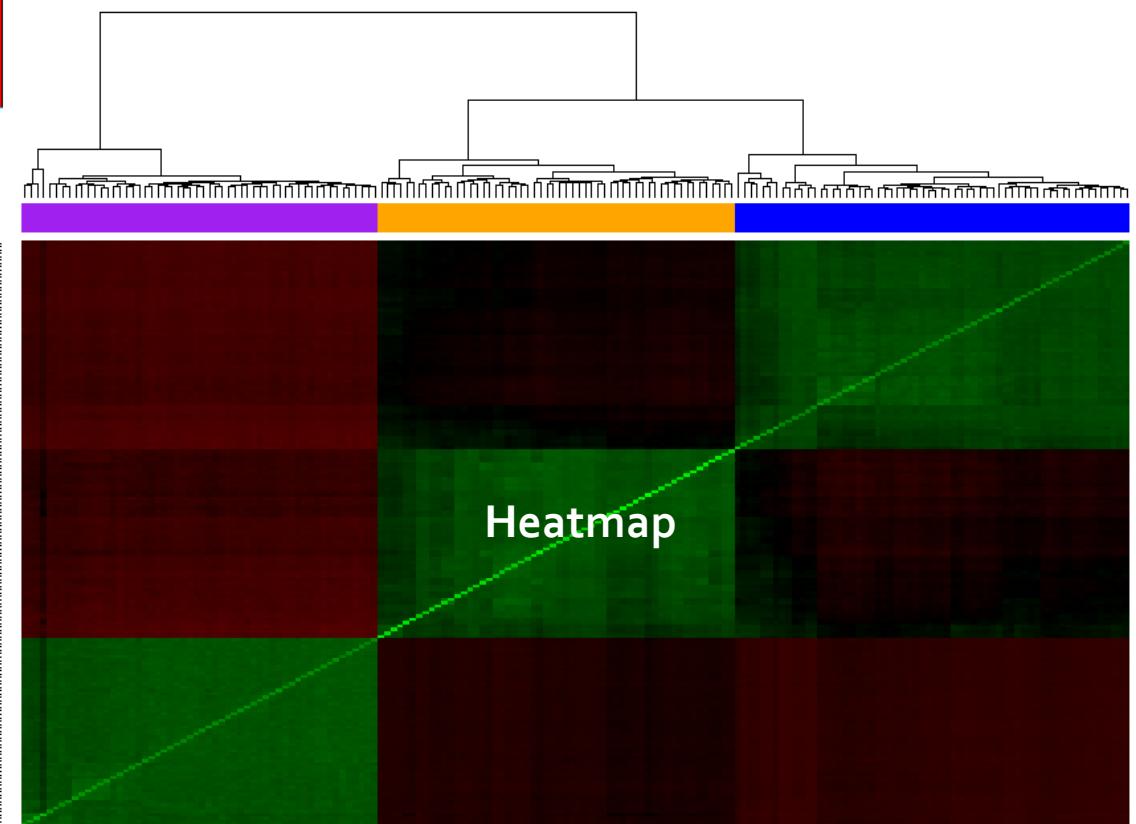


# Exploration

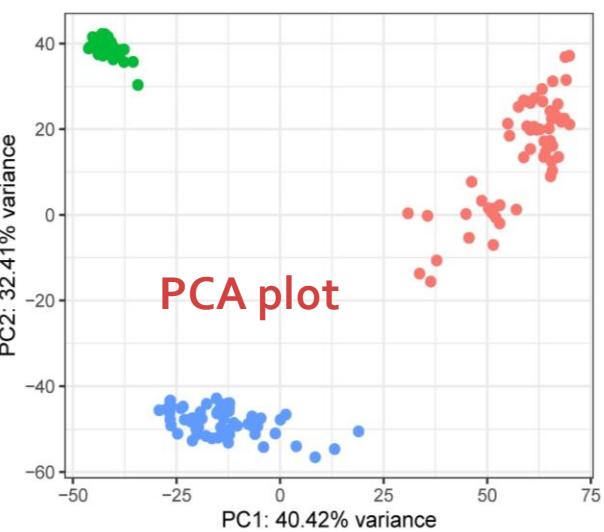
- SVA
- PVCA
- BatchQC
- ComBat
- R



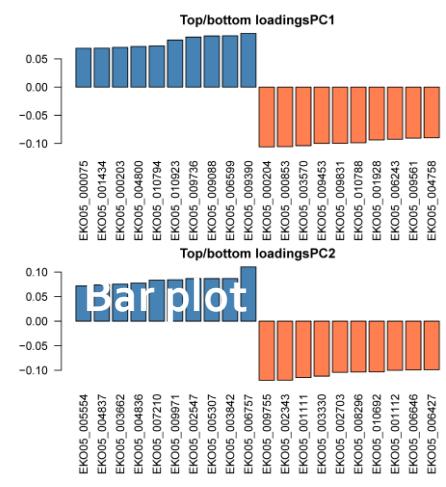
Top 500 most variable genes across samples



MDS plot



PCA plot



Bar plot

# DGE analysis

## countData

	ctrl_1	ctrl_2	exp_1	exp_1
geneA	10	11	56	45
geneB	0	0	128	54
geneC	42	41	59	41
geneD	103	122	1	23
geneE	10	23	14	56
geneF	0	1	2	0
...	...	...	...	...
...	...	...	...	...
...	...	...	...	...

## colData

	treatment	sex
ctrl_1	control	male
ctrl_2	control	female
exp_1	treatment	male
exp_2	treatment	female

Sample names:

**ctrl\_1, ctrl\_2, exp\_1, exp\_2**

**countData** is the count matrix  
(number of reads mapping to each gene for each sample)

**colData** describes metadata about the *columns* of countData

**colnames(countData) == rownames(colData)**

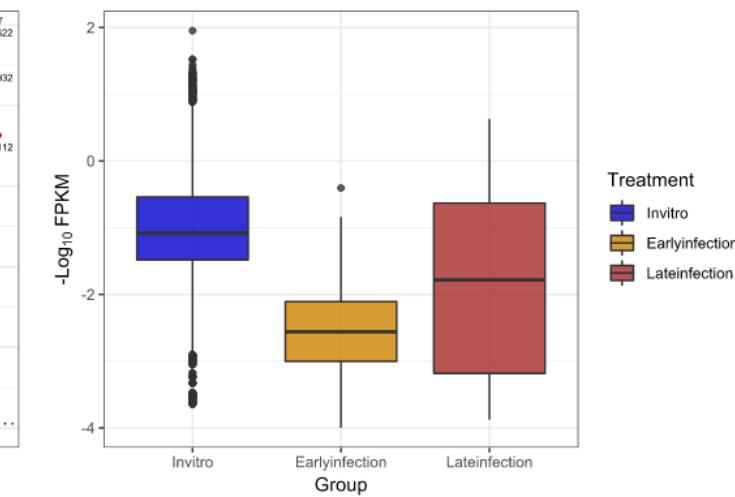
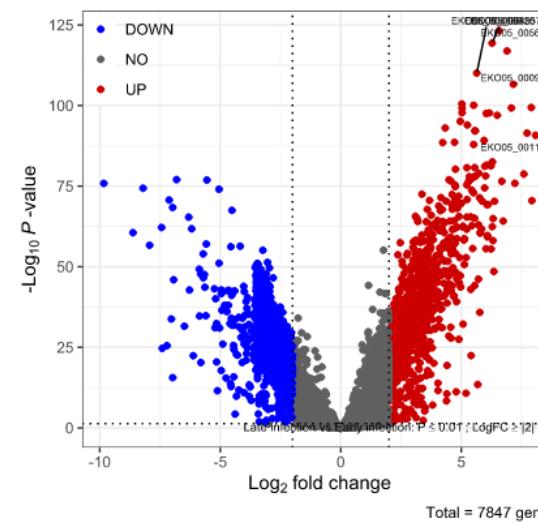
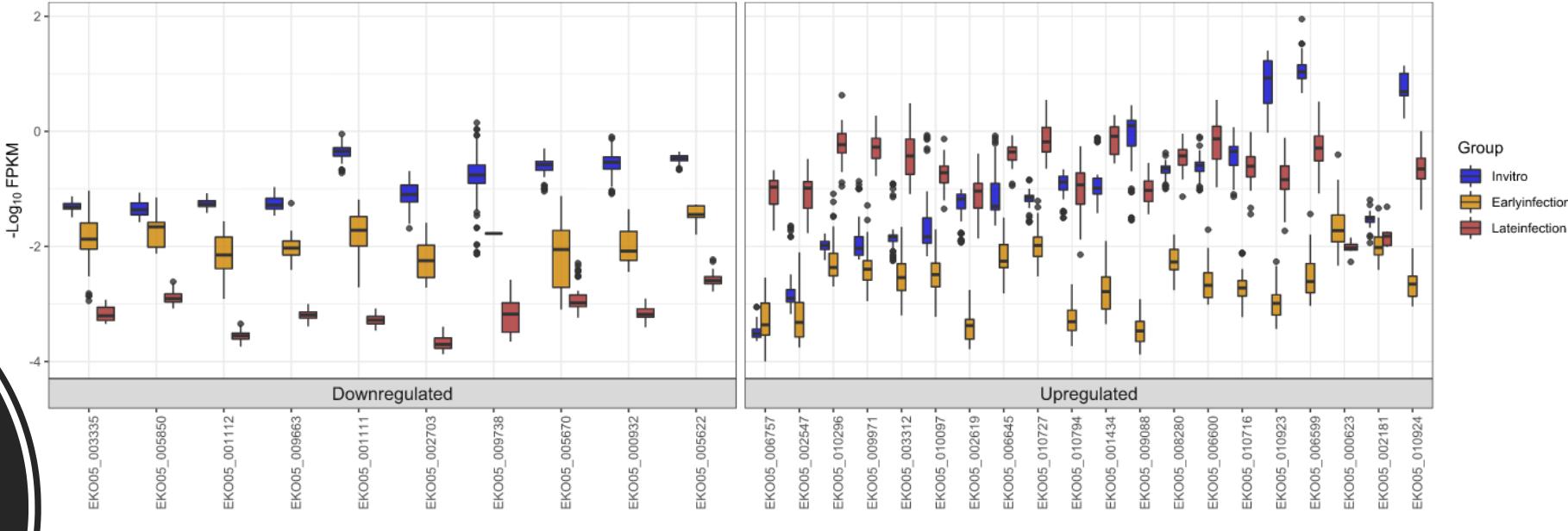
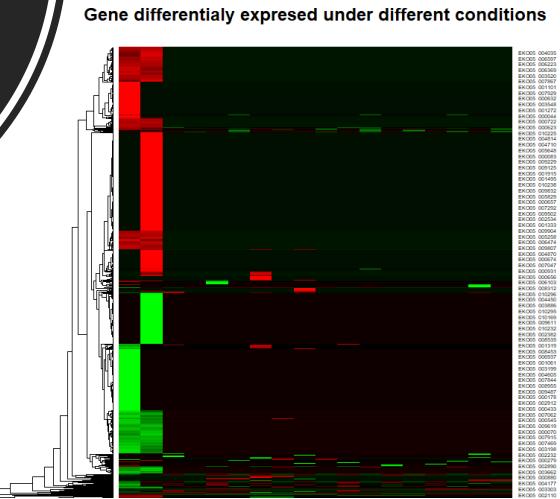
## Experimental design

- Factors always ordered alphabetically by default
- The first variable is used as the base (intercept) upon which all other variables are compared
- Set variable levels manually to ensure the model is consistent with the hypothesis.

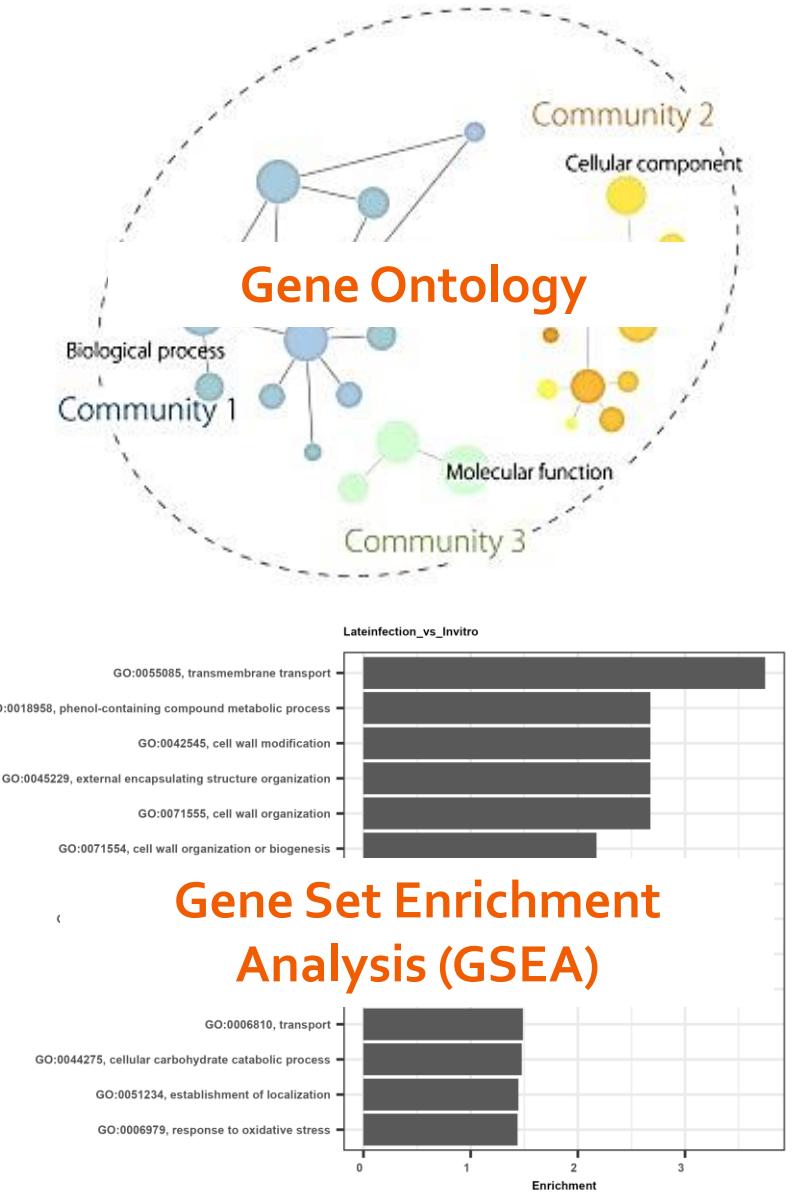
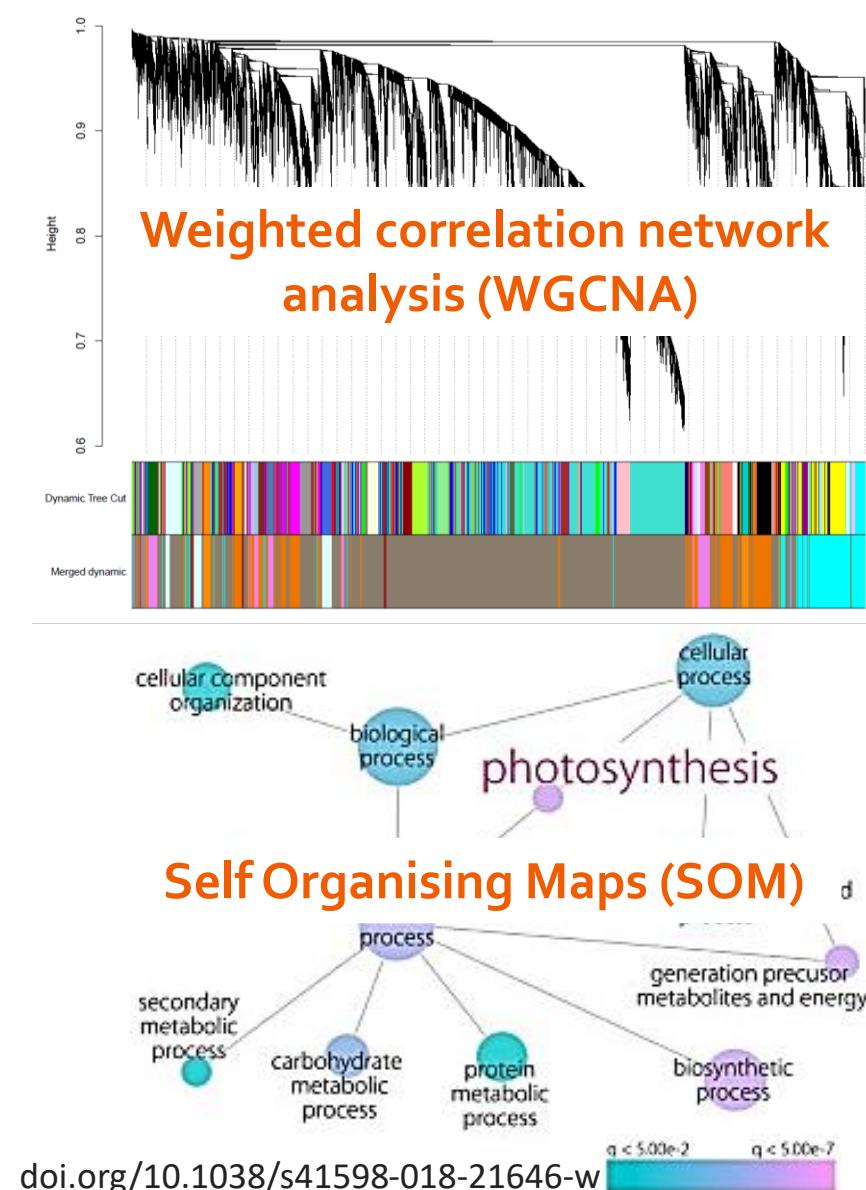
`relevel(colData$treatment, 'control')`  
`x$col <- factor(x$col, levels = c(a, b..))`

design = ~treatment :::: compare treated-vs-control using control as reference.  
design = ~treatment+sex :::: a + M vs F with female as reference.  
design = ~o+sex+treatment :::: a and b but with sex not contrasted.  
design = ~treatment+sex+treatment:sex :::: a+b+interaction between a & b.  
design = ~treatment+sex+batch+treatment:sex :::: d + model for batch effects.

# Visualize



# Functional analysis

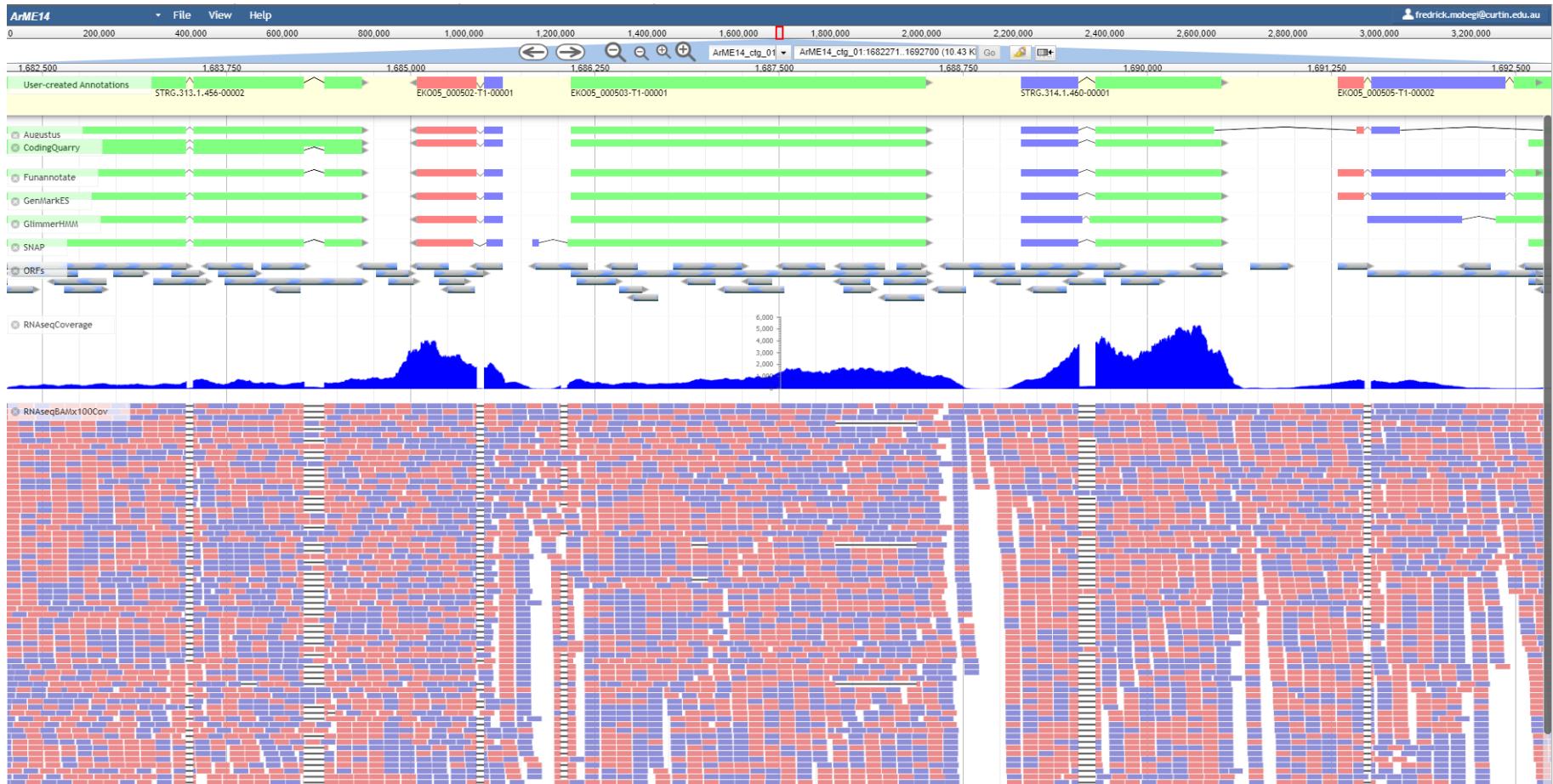


The aim is to identify important pathways or functional modules represented in your expression data..

# Web Apollo: a web-based genomic annotation editing platform

[Eduardo Lee](#), [Gregg A Helt](#), [Justin T Reese](#), [Monica C Munoz-Torres](#), [Chris P Childers](#), [Robert M Buels](#),  
[Lincoln Stein](#), [Ian H Holmes](#), [Christine G Elsik](#) & [Suzanna E Lewis](#) 

[Genome Biology](#) 14, Article number: R93 (2013) | [Cite this article](#)



Genome  
annotation



## Further reading

- Fu, Yu *et al.*, "Elimination of PCR duplicates in RNA-seq and small RNA-seq using unique molecular identifiers." *BMC genomics* 19.1 (2018): 531.
- Parekh *et al.*, "The impact of amplification on differential expression analyses by RNA-seq." *Scientific reports* 6 (2016): 25533.
- Klepikova *et al.*, "Effect of method of deduplication on estimation of differential gene expression using RNA-seq." *PeerJ* 5 (2017): e3091.
- Soneson *et al.*, "Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences." *F1000Research* 4 (2015).
- Zhang *et al.*, "Evaluation and comparison of computational tools for RNA-seq isoform quantification." *BMC genomics* 18.1 (2017): 583.
- Teng *et al.*, "A benchmark for RNA-seq quantification pipelines." *Genome biology* 17.1 (2016): 74.
- Conesa *et al.*, "A survey of best practices for RNA-seq data analysis." *Genome biology* 17.1 (2016): 13.
- Hsieh *et al.*, "Effect of de novo transcriptome assembly on transcript quantification" 2018 *bioRxiv* 380998.
- Wang, Sufang, and Michael Gribkov. "Comprehensive evaluation of de novo transcriptome assembly programs and their effects on differential gene expression analysis." *Bioinformatics* 33.3 (2017): 327-333.



## Summary

- Quality Control is KING in everything transcriptomics.
- A good experimental design helps avoid/minimize confounding effects.
- Biological replicates rule above all other replicates and/or sequencing methods.
- Carefully planning sample prep and sequencing ensures you capture the intended hypothesis in your experiment.
- Always discard low quality bases, reads, genes and samples.
- Experiment with different pipelines and software to identify biases.
- Corroborate assumptions about the data with methods and tools used for analysis.
- No two tools will give you identical results. Choose wisely
- You can never have too many illustrations in RNA-seq.

Questions

