# A model high level analysis, and correlation analysis

Harshita Gaur
Hossameldin Ali

# A **model** high level analysis, and **correlation analysis**

- What is model analysis?

- What is correlation analysis?
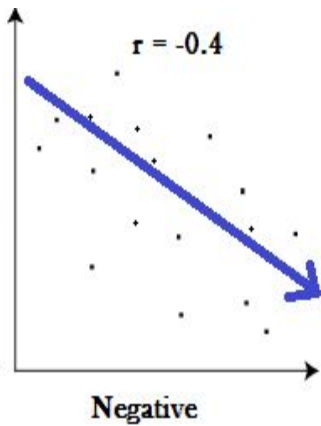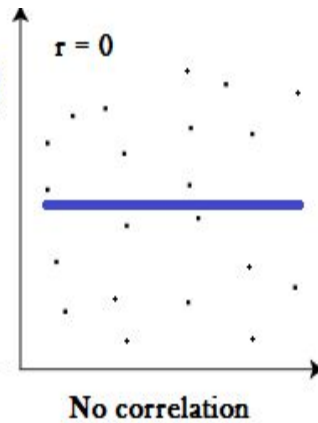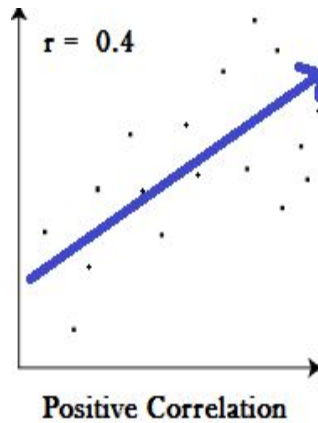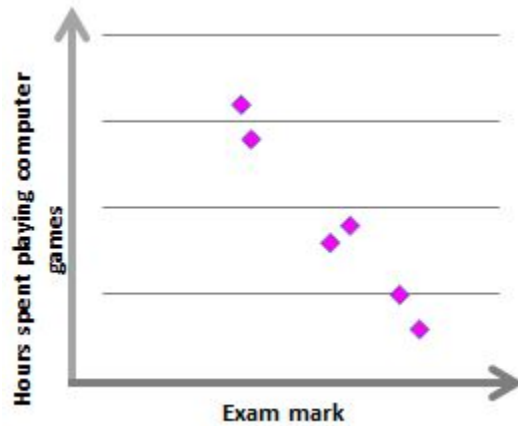
- What is a correlation network?

# Model high level analysis

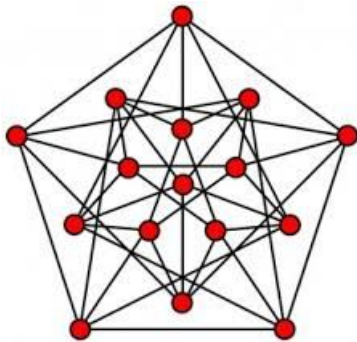- Ex: Jaguar doors testing
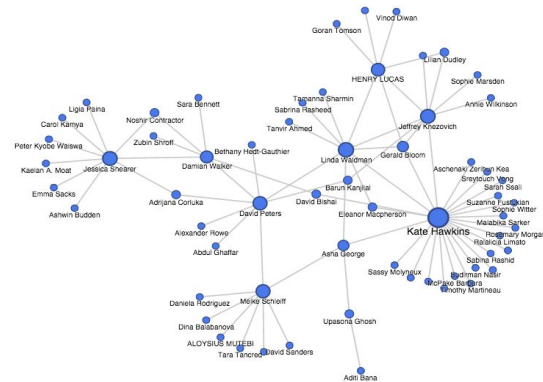- We are more into computer models

# Correlation analysis



- Range -1 to 1
- Stronger correlation >0.5 or > -0.5

# What is a network



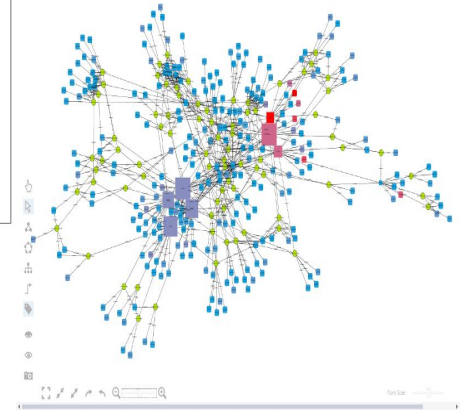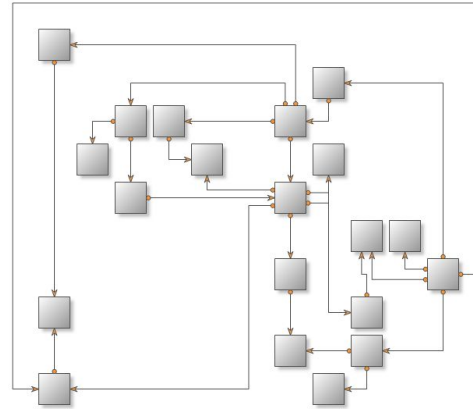



Graphs: mathematical structures used to model pairwise relations between objects

# Biological network



- Nodes are usually genes, transcripts or proteins.
- Nodes: shape, size and color
- Edges: directed or undirected
- Edges: can be weighted
- Aesthetically pleasing: number of edges crossing is minimized

# BioLayout Express 3D

# Biolayout Express 3D



- 30,000 nodes and 2 – 3 million edges
- Appropriate data from any source, biological or otherwise
- Analysis gene expression data
- Optimized layout algorithm and graph clustering

# The graph paradigm for microarray expression data

- Makes no earlier assumptions regarding the experimental design, normalization method, microarray platform or questions being addressed by the study

- However, all of these factors do influence the properties of a network, and therefore it is worth spending some time to discuss them.

# Experimental design

- The size of a data set is a function of the number of probes on the array and the number of samples analyzed
  - number of samples is small ( < 10): higher correlation between data derived from different probes
  - ~ 20 samples or more: no need to filter data

# Correlation network (WGCNA)

- Weighted Gene Co-expression Network Analysis
- A method used to discover patterns among variables based on correlation
- Modules (Clusters)
- WGCNA is linear:
  - Faster computations
  - No overfitting problems
  - Suitable for small datasets

# Identification of a human neonatal immune-metabolic network associated with bacterial infection

# Power Calculations



Sample size to obtain 90% power for a given proportion of genes on array, at $p \leq 0.01$

- 30 control pre-vaccinated 9 month old Gambian infant samples was used to perform the power calculations
- A high-performance computing model for host RNA demonstrated that as few as 24 whole blood RNA markers would be sufficient for predicting infection.

# what is the data

**a** PCA Mapping (37.2%)

**b** PCA Mapping (79.7%)

Category
• control
• infected

Sequence of study analyses prior to validating 52-gene set as a classifier

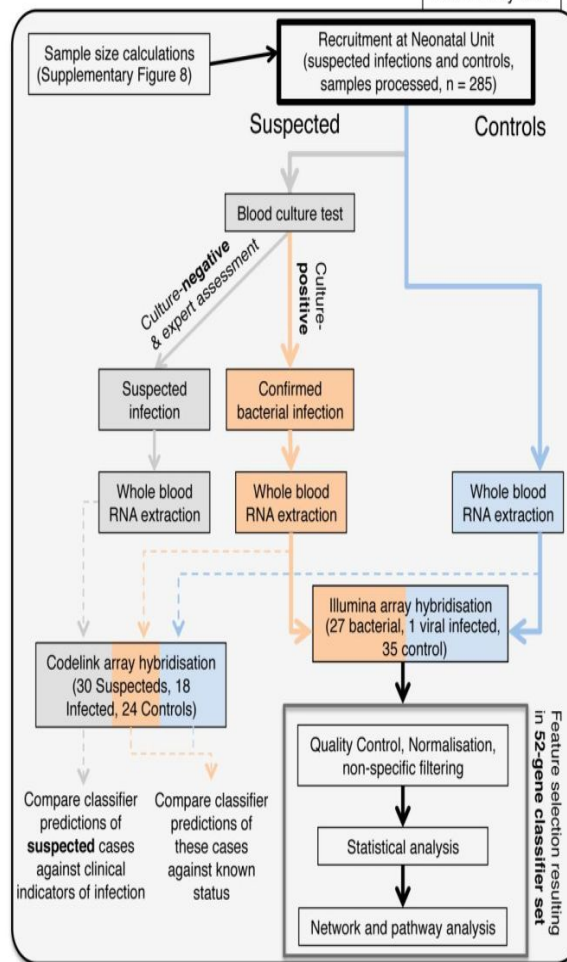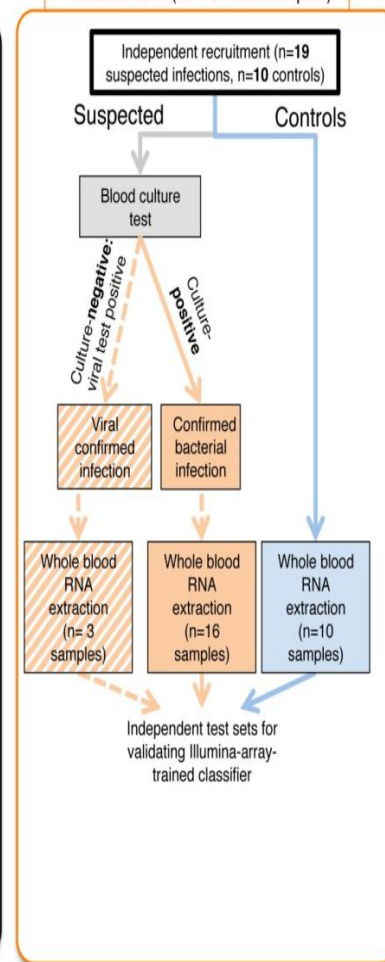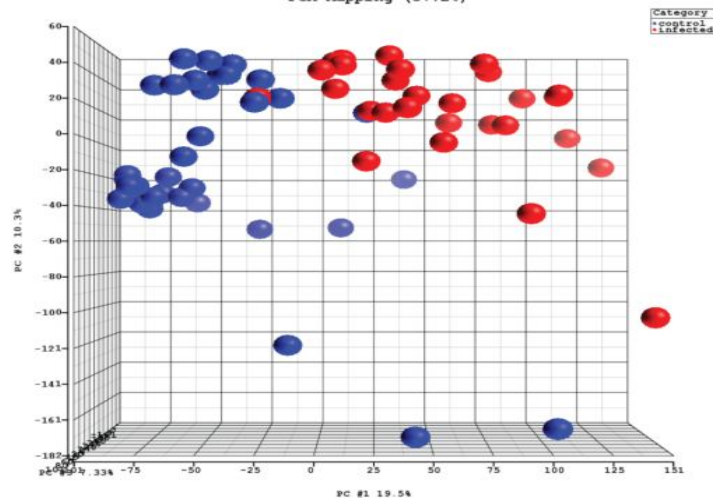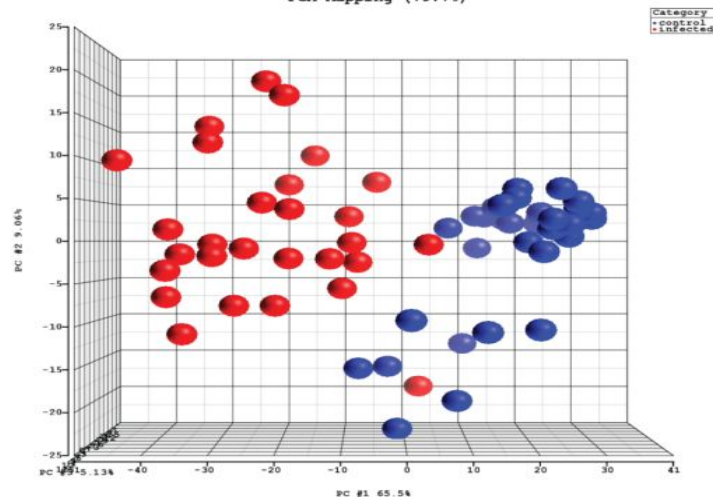Illumina array data on 84 neonatal samples (n=48,802 gene probes) → Sample and data quality control → 63 samples and n=48,802 gene probes → Normalisation and non-specific filtering → n=23,342 gene probes → Statistical inference testing between study groups → Multiple-testing adjusted statistical p-values and differential expression

Select statistically significant (adj. p≤0.01) and biologically relevant (2-fold up/down regulation) genes

(CV ≥ 10%)

Hierarchical clustering of samples (n=10,206 gene probes)
Fig. 1a

Analysis of co-stimulatory/ inhibitory genes
Fig. 5d

Immune cell compartment analysis
Fig. 3a,b,c,e,f
Fig 6c

Clinical marker regression analysis
Fig. 5a

n=824 gene probes

Hierarchical clustering of genes = 3 expression clusters
Fig. 1b

Identification of patient-specific responses with Biolayout
Fig. 6a,b

Network analysis with integration of InnateDB gene relationships
Fig. 4 and Suppl. Fig. 2

Investigation motivates removal of one cluster of genes and increased stringency of statistical/biological thresholds (adj. p ≤ 10^-5, 4-fold change)

n=52 gene probes

Use as classifier biomarker set, test and validate
Fig. 7 and Suppl. Fig. 7

Network and pathway analysis confirmation of biomarker set
Fig. 2 and Suppl. Fig. 3

# Heat map showing hierarchical clustering

- using eBayes with Benjamini-Hochberg correction adj.P<= 0.01, absolute fold change >=2

- 824 propes
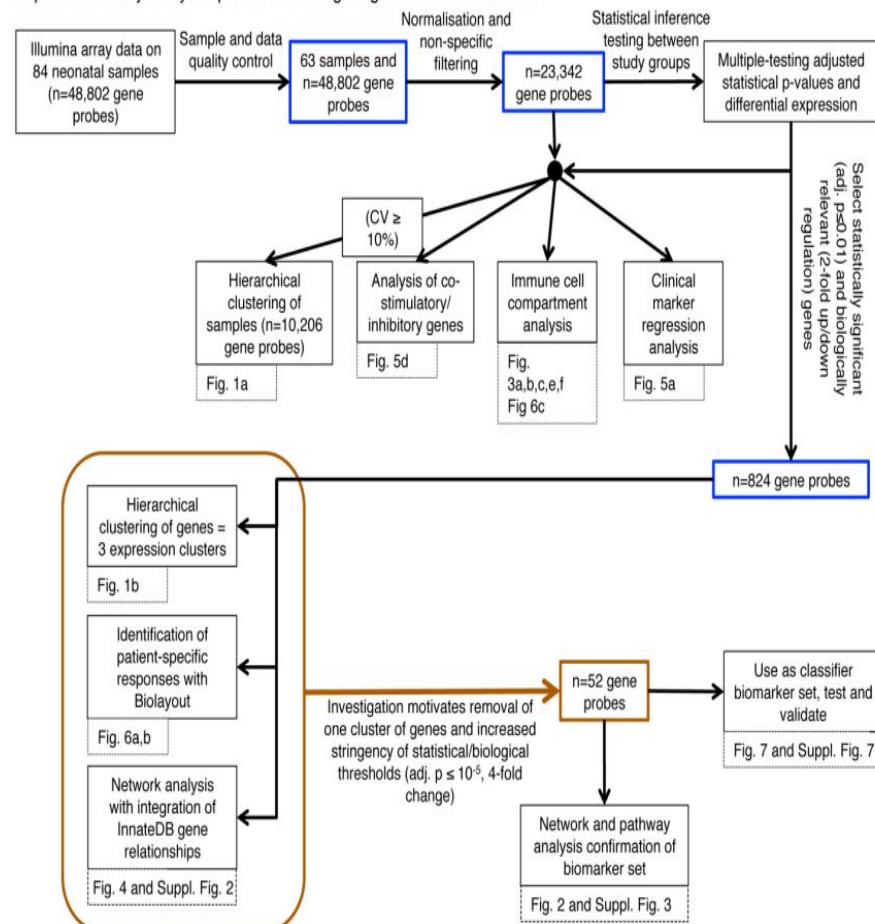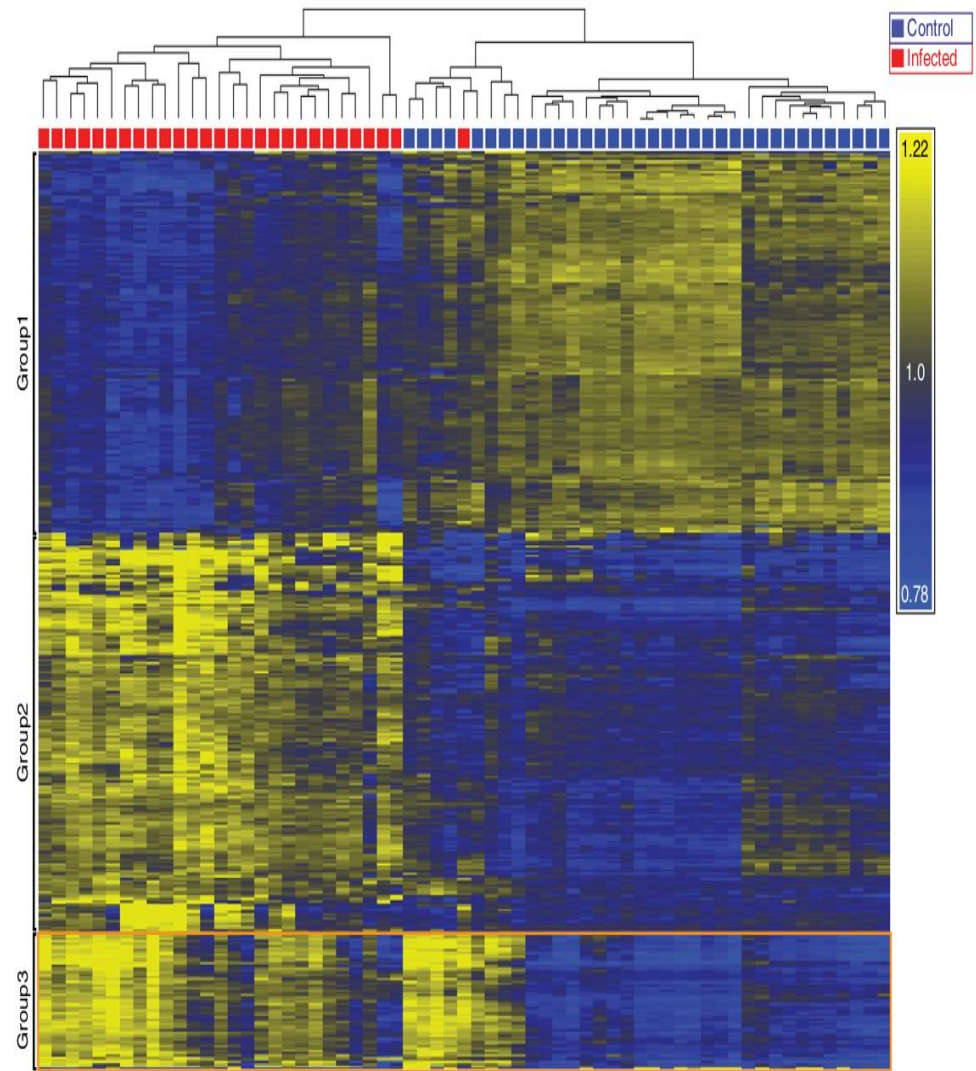
**a** PCA Mapping (37.2%)

**b** PCA Mapping (79.7%)

Sequence of study analyses prior to validating 52-gene set as a classifier

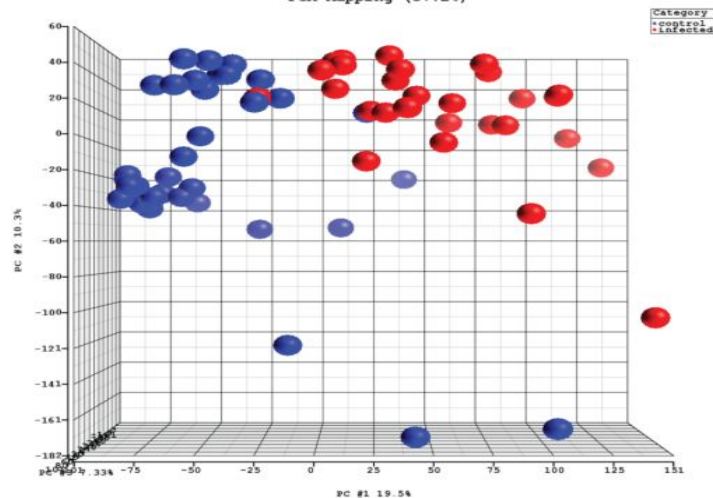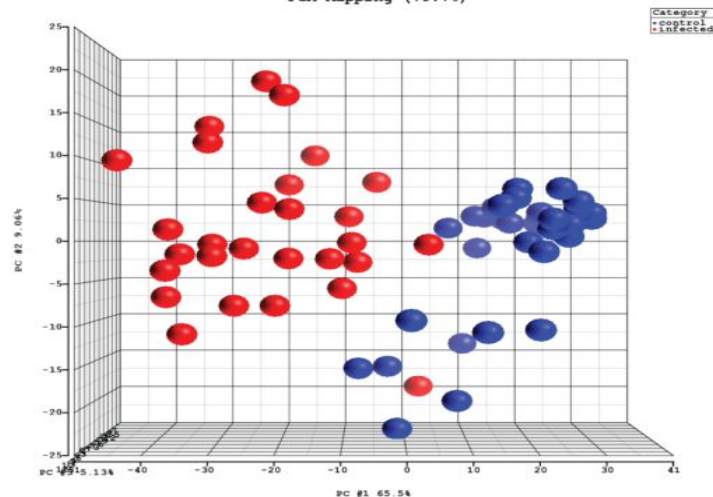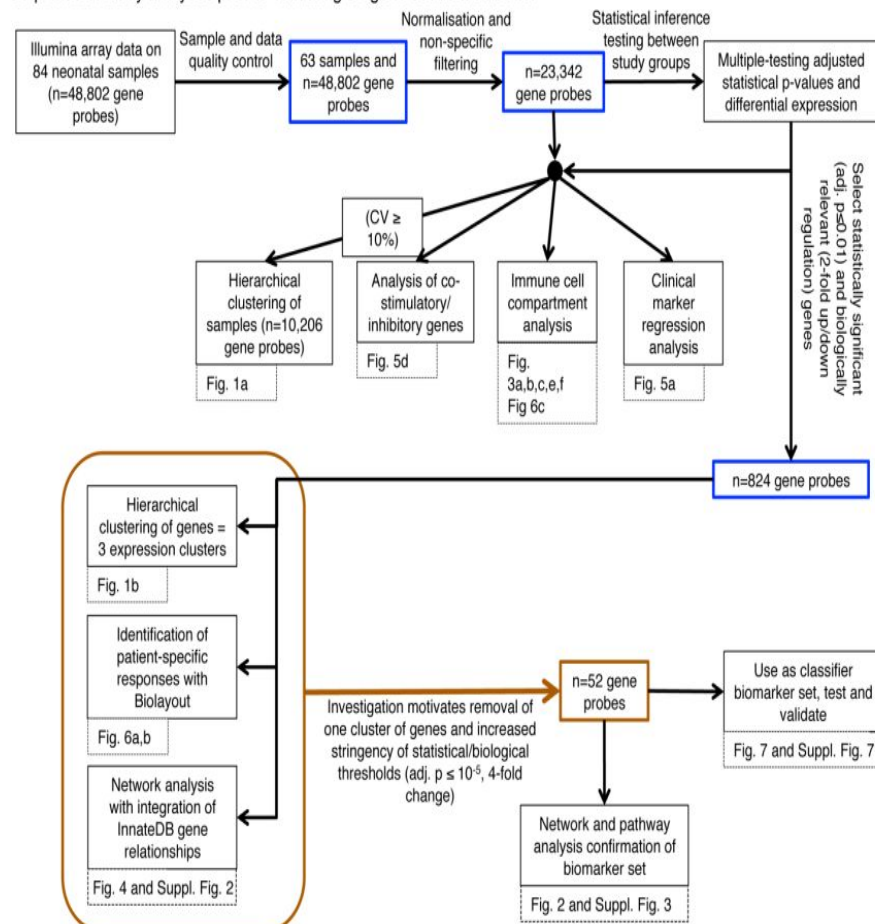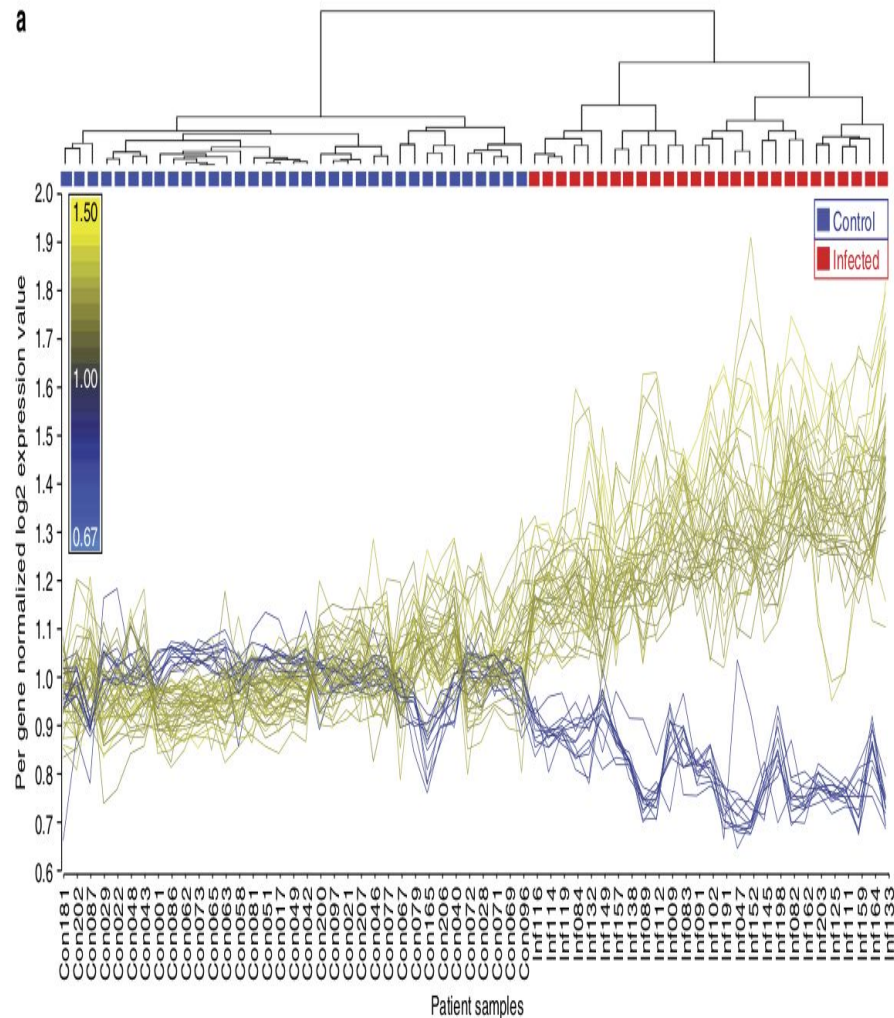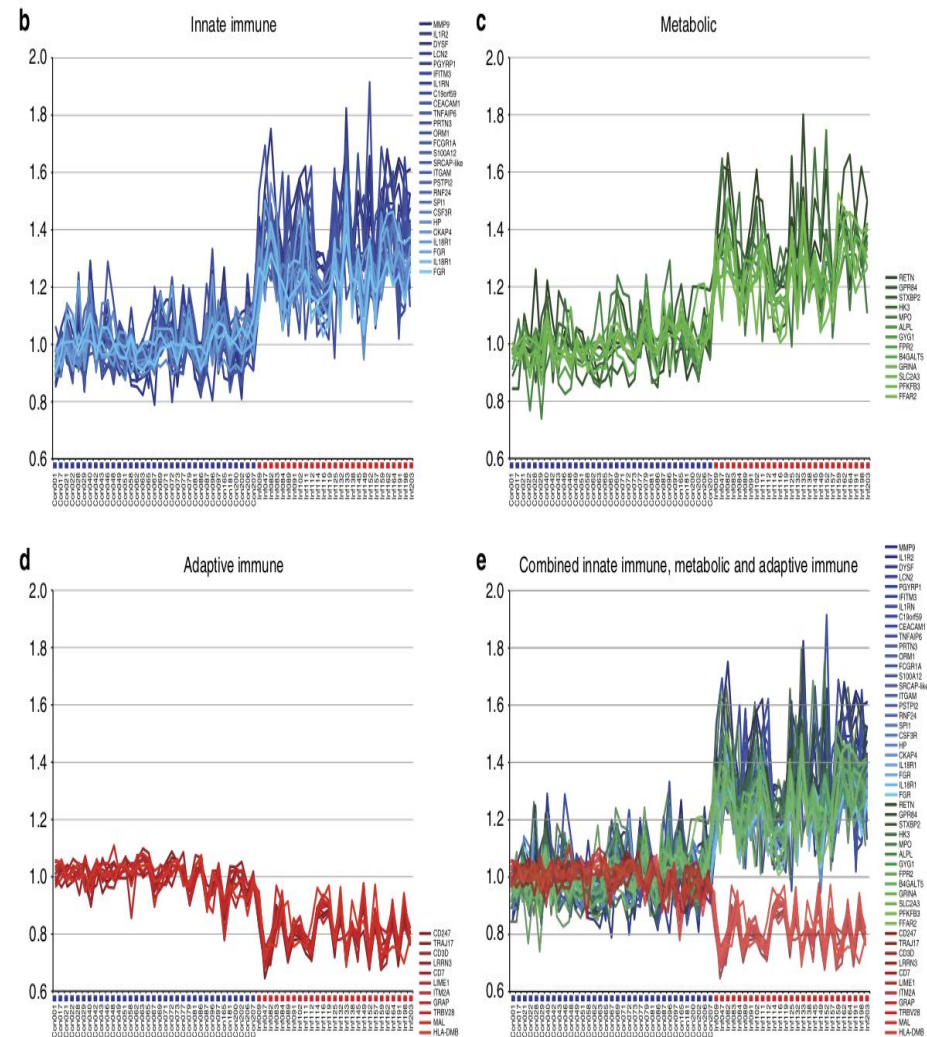# Pathway biology responses to infection
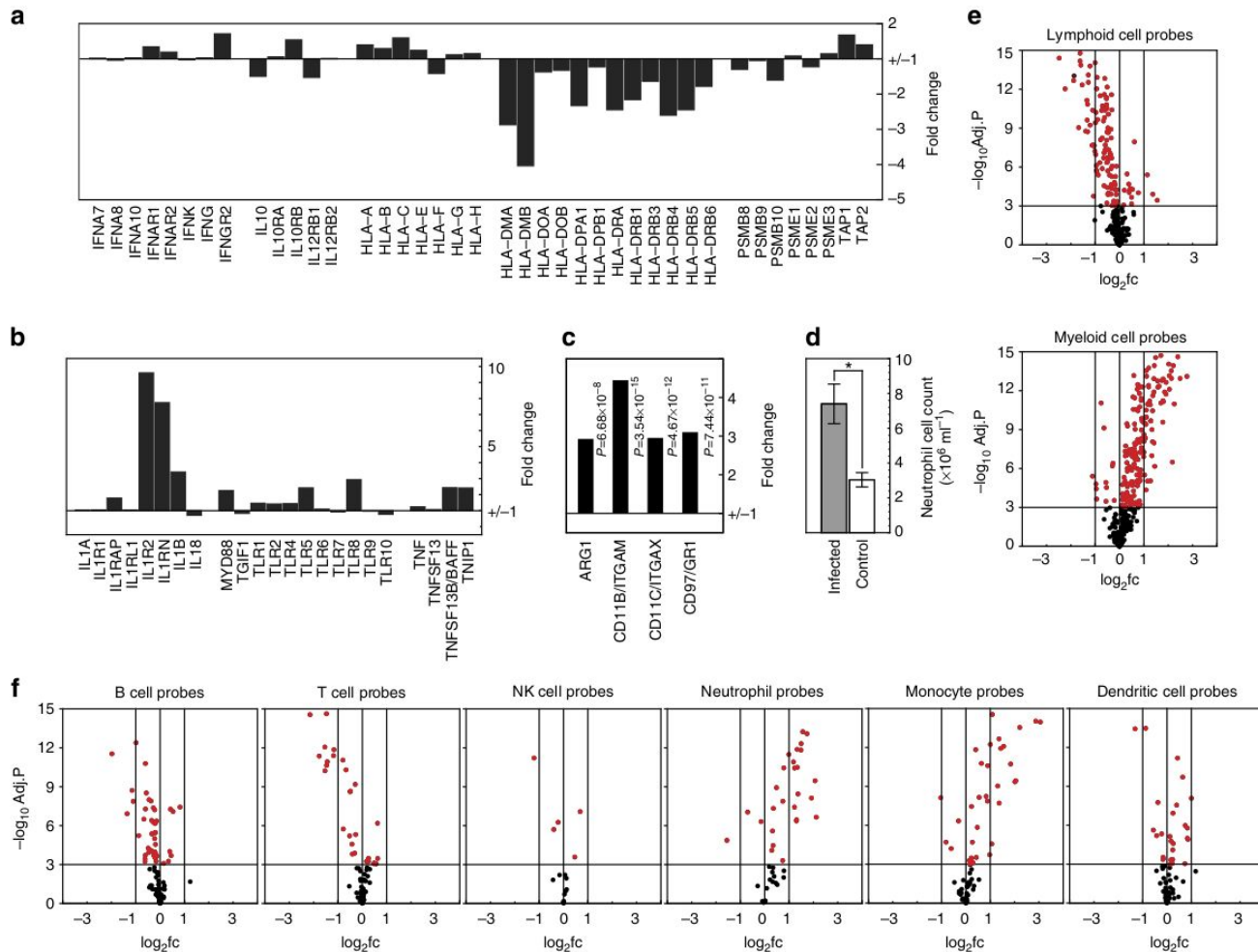
additional filtering to groups 1 and 2 using more stringent cutoffs

Analysis of the resulting 52 genes revealed sub-networks

# Three functional pathway classes

- Innate immunity
- Adaptive immunity
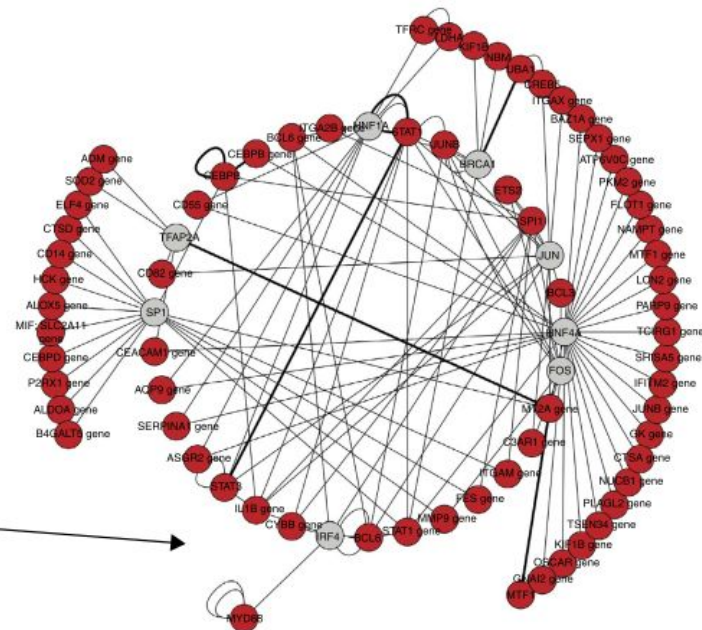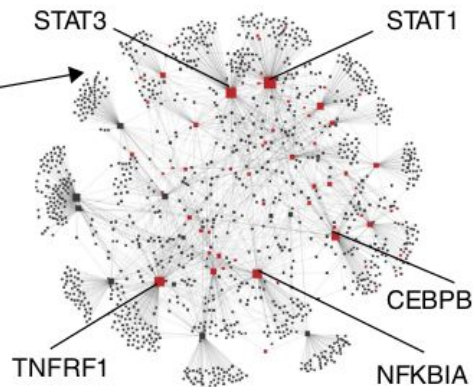- Unexpectedly metabolic pathways for sugar and lipid metabolism

# Changes in immune cell compartment

"We conclude that in sepsis the heightened innate immune cellular response is driven by monocytes/macrophages and neutrophils and is counter-balanced by inhibitory pathways resulting in a net suppression of the adaptive immune arm, especially associated with the T-cell compartment."

a

InnateDB resource filtered for a high confidence manually annotated data set of approximately 2,500 human molecular interactions is used

# Regulatory nodes

we conclude that the overriding pathophysiological signal associated with neonatal infection is one of the increased innate immune metabolic responses with an unbalanced homeostatic regulation of adaptive immunity. The specific and intense activation of innate immune signalling, moderated via inhibitory pathways is consistent with the notion of an elevated set point in neonates in comparison with adults for guiding a suppressed adaptive immune response

# Metabolic pathway biology and immune homeostasis

Observed changes in specific metabolic pathways including glucose, energy and cholesterol metabolism.
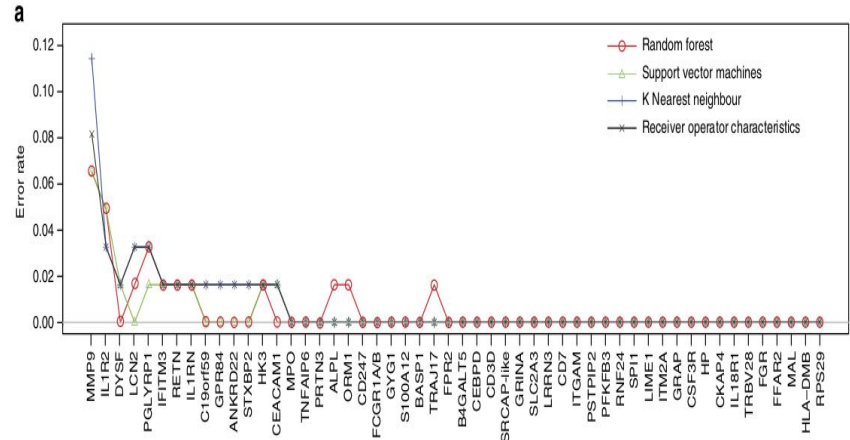
Including change in B4GALT5 and lipocalin 2, involved in innate response to bacterial infection is also involved in regulating lipid and glucose metabolism.

Indicating theses metabolic processes are likely to be linked with innate immune response

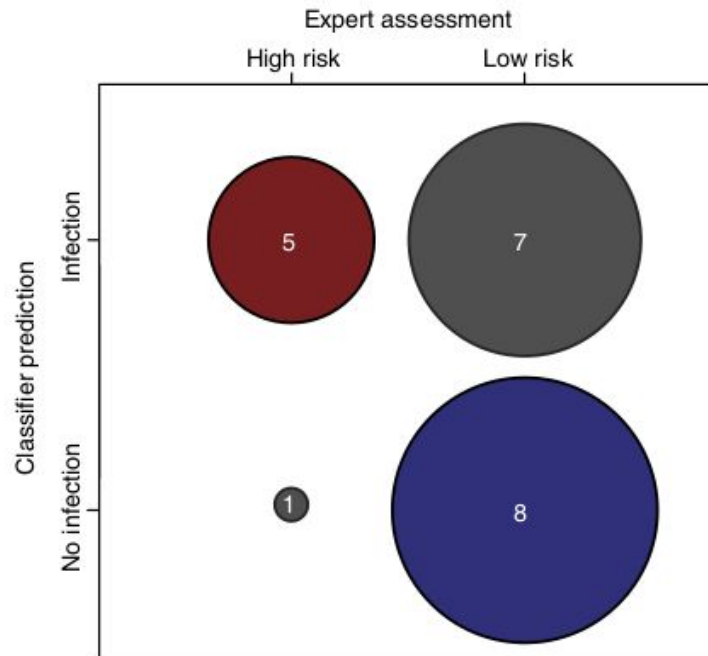# 52-Gene dual-network accurately predicts neonatal infection

- Four distinct machine-learning algorithms
  - Random forest
  - support vector machine
  - k-nearest neighbour
  - Receiver operator characteristics (ROC)
- When the number of genes included was 19 or more the error rate was consistently 0%*



*consistent with the in-silico model

# 52-gene classifier vs Expert assessment

- These findings highlight the difficulty faced by clinicians in determining cases of blood culture-negative sepsis and strongly support the future clinical utility of the 52-gene classifier

# Using other platforms

A ROC classifier when applied to genes with a completely different microarray platform correctly assigned 100% of samples to control or bacterially infected groups (sensitivity 100%, specificity 100%).