# Gene enrichment analysis

**Most of today's material courtesy of Terry Speed**

---

# Are sets of genes differentially expressed?

**The sets** we refer to here are all the *outcomes* of analyses. Later we discuss sets specified *a priori*.

**Examples of sets**. They could be the list of all genes whose differential expression (e.g. average *M*-value) exceeds a given threshold, typically a liberal one, which would not correspond to any real "significance", e.g. *1.5*-fold. They might be clusters.

**What do we mean** by a set being differentially expressed. Here it is a convenient shorthand for being unusual in relation to all the genes represented on the array, for example, by being functionally enriched, in the sense of having more genes of a given category than one would expect, *by chance.*
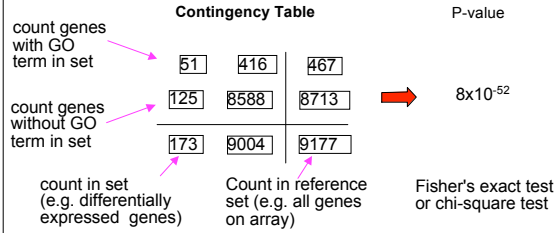
---

# GO and microarray gene sets

**Hypothesis:** Functionally related, differentially expressed genes should accumulate in the corresponding GO-group.

**Problem:** to find a method which scores accumulation of differential gene expression in a node of the GO.

We describe the calculation from the program Gostat. For all the genes analysed, it determines the annotated GO terms and all splits. It then counts the # of appearances of each GO term for the genes in the set, as well as the # in the reference set, which is typically all genes on the array. Then a 2×2 table is formed, see over page, and a *p*-value calculated.

# Is a GO term is specific for a set?

**Contingency Table**       P-value

count genes with GO term in set →

| 51 | 416 | 467 |

count genes without GO term in set →

| 125 | 8588 | 8713 |

→ $8 \times 10^{-52}$

| 173 | 9004 | 9177 |

count in set (e.g. differentially expressed genes)

Count in reference set (e.g. all genes on array)

Fisher's exact test or chi-square test

---

# The multiple testing problem

Naturally one doesn't test a single GO term or split, but many, perhaps 1000s. As with testing of single genes, we need to deal with the multiple testing problem. Many of the solutions from there carry over: Bonferroni, Holm, step-down minP, FDR, and so on. But there are also special problems here, deriving from the nesting relationships between splits. In my view, these are not easily dealt with, and require more research.

**Related questions**. How can we compare the results of different lists being compared? And, rather than select a set of genes using a cut-off, can we make use the gene abundances or p-values for differential expression?

---

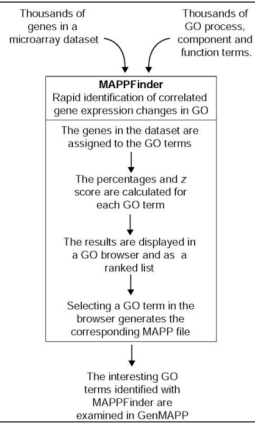**GOstat:** Tool for finding significant GO terms in a list of genes
http://gostat.wehi.edu.au



List of significant GO terms

List of genes and GO associations

2

# There are many similar tools

Here are a few.

GenMAPP, and MAPPFinder
EASE (DAVID)
FunSpec
FatiGO
.....

---

**Outline of MAPPfinder:**
**MAPP = MicroArray**
**Pathway**
**Profiler**

Thousands of genes in a microarray dataset

Thousands of GO process, component and function terms.

**MAPPFinder**
Rapid identification of correlated gene expression changes in GO

The genes in the dataset are assigned to the GO terms

The percentages and $z$ score are calculated for each GO term

The results are displayed in a GO browser and as a ranked list

Selecting a GO term in the browser generates the corresponding MAPP file

The interesting GO terms identified with MAPPFinder are examined in GenMAPP

---

# Analyzing microarray data by functional gene sets defined *a priori*

Analysis at the level of single gene:
- Identifying differentially expressed genes becomes a challenge when the magnitude of differential expression is small.
- For some differences, many genes are involved.

Analysis at the level of functional group: why?
By incorporating biological knowledge, we can hope to detect modest but coordinate expression changes of sets of functionally related genes.

**PGC-1α-responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes**

Mootha *et al*, Nature Genetics July 2003

Data: Affymetrix microarray data on 22,000 genes in skeletal muscle biopsy samples from 43 males, 17 with normal glucose tolerance (NGT), 8 with impaired glucose tolerance and 18 with Type 2 diabetes (DM2).

In their single gene analysis, a *t*-statistic was calculated for each gene. No significant difference found between NTG and DM2 after adjusting for multiple testing.

Their idea: test 149 *a priori* defined gene sets for association with disease phenotypes.

---

# 149 gene sets

Sets of metabolic pathways:
- manually curated pathways (standard textbook literature reviews, and LocusLink)
- Netaffx annotations using GenMAPP metabolic pathways
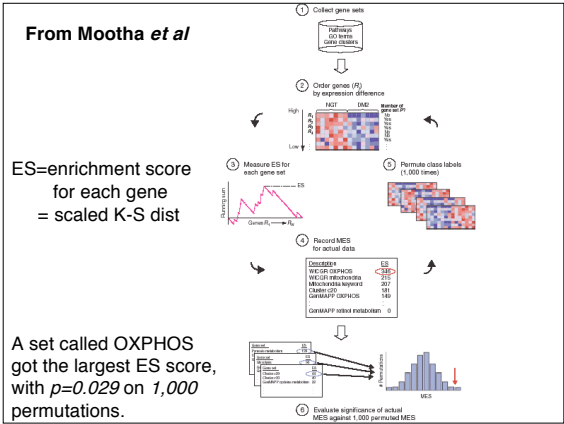
Sets of coregulated genes:
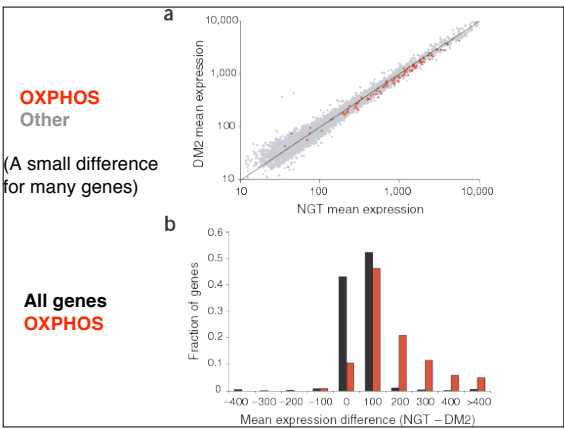- SOM clustering of the mouse expression atlas

**How do you compare group of interest to rest?**

---

**Two sample Kolmogorov-Smirnov test**



To compare two empirical cdfs, $S_M(x)$ and $S_N(x)$ based on samples of size M and N, resp, the Kolmogorov-Smirnov (K-S) test uses the K-S distance $D_{MN} = \max_x |S_M(x) - S_N(x)|$. This is normalized by multiplying by $\sqrt{(M^{-1} + N^{-1})}$. It has a complicated null distribution, which can be approximated by permuting.

**From Mootha _et al_**

ES=enrichment score
   for each gene
   = scaled K-S dist

A set called OXPHOS
got the largest ES score,
with _p=0.029_ on _1,000_
permutations.



---



OXPHOS
Other

(A small difference
for many genes)

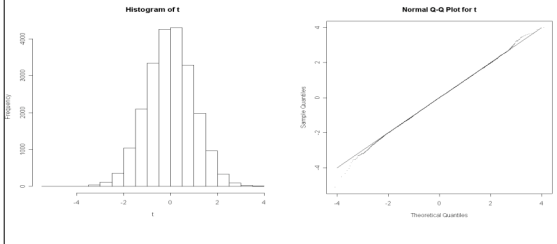All genes
OXPHOS

---

# Simplification

Mootha _et al_ did a two sample K-S test to compare genes
in a specific gene set with genes not in that set.

Instead of doing this, why don't we simply do a one
sample test, comparing each gene set to the whole
(population) directly?
Each gene set is small w.r.t. the entire set of genes, so all
other genes ≈ all genes.

If we have approximate normality, a _z_-test should work
for shift alternatives. A chi-squared test for scale changes
also works.

# Mootha's ts are approx normal

**Histogram of t**

**Normal Q-Q Plot for t**



---

# One sample z-test

Assumption: the (population of) *t*-statistics of all genes follow normal distribution. Denote the mean by $\mu$ and the SD by $\sigma$.
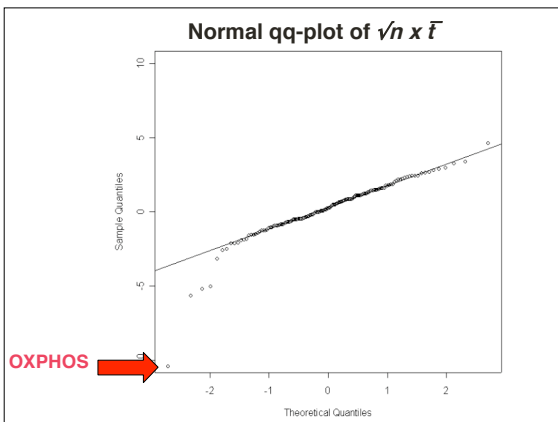
If this is the case, the best test of the null hypothesis that a sample $t_1$, $t_2$, …..,$t_n$ is from this distribution, with alternative a shift of the original distribution is based on $\bar{t}$. Specifically, it uses

$$z = (\bar{t} - \mu)/\sigma/\sqrt{n}.$$

In general, we'd expect $\mu=0$ and $\sigma=1$, and this is the case for Mootha's ts. Thus we test the null hypothesis that our sample comes from the same population using

$$z = \sqrt{n}\,\bar{t}.$$

Let's do a normal *qq*-plot of the 149 *z*-statistics of this form.

---

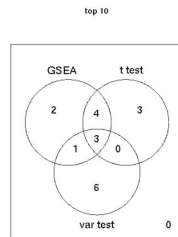## Normal qq-plot of $\sqrt{n}\ x\ \bar{t}$



**OXPHOS** ➡

## Result from one sample z-test

- OXPHOS is easily identified as ≈ -10.
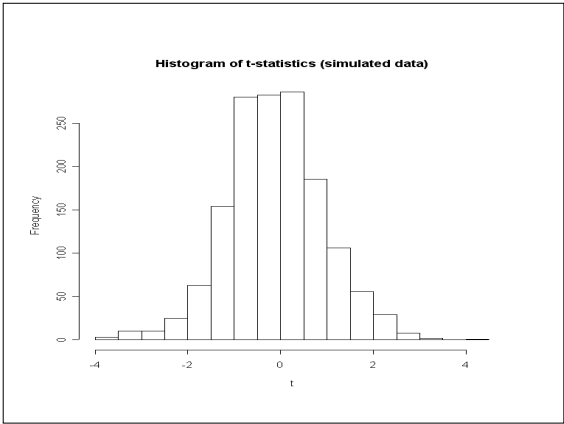- The next three sets on the top ranking list are all related to oxidative phosphorylation.

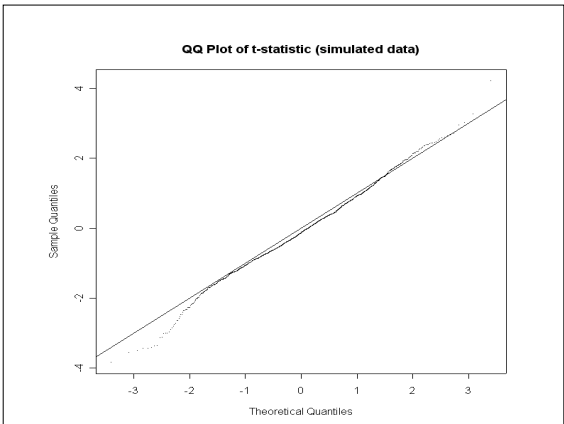|  | z | n | # overlapping w/ OXPHOS |
|---|---|---|---|
| OXPHOS_HG-U133A_Probes | -10.4 | 114 | 114 |
| Human_mitoDB_6_2002_HGU133A_ probes | -5.6 | 594 | 106 |
| Mitochondr_HG-U133A_probes | -5.2 | 615 | 103 |
| MAP00190_Oxidative_phosphorylation | -5.0 | 75 | 29 |

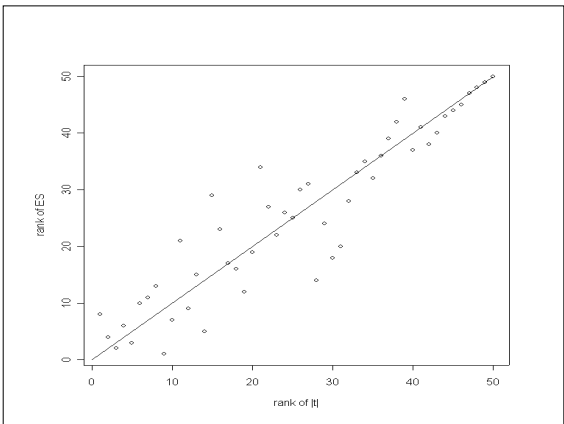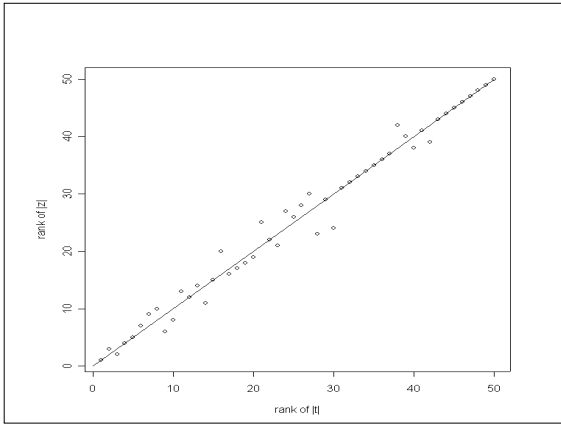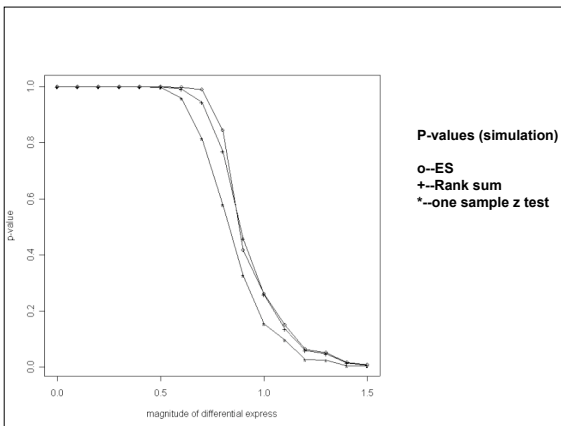## Similar Results

top 10



## Simulation 1

- 1500 × 29 gene expression values are generated from N(0,1), representing 1500 genes for 9 cases and 20 controls.

- The 1500 genes are divided into 50 gene sets, each with 30 genes. The genes are correlated within each gene set.

- We manipulate the gene expression level of the cases of the first gene set so that the magnitude of difference is known.

**Histogram of t-statistics (simulated data)**

**QQ Plot of t-statistic (simulated data)**

**P-values (simulation)**

**o--ES**
**+--Rank sum**
**\*--one sample z test**

# Conclusion

- When the population follows a normal distribution, the one-sample *z*-test is most powerful for shift alternatives (no surprise: theory says it has to be).
- From the simulation study, the one sample z-test is seen to be more powerful than the two sample *K-S* test for shift alternatives (even less of a surprise).
- The new method is not as compute intensive as the K-S test.
- Similar results can be given for the following test statistic, for scale change alternatives: for a set of n genes

$$z' = \sum_{i=1,...n} [(t_i - \bar{t})^2 - (n-1)] / \sqrt{(2(n-1))}.$$

(A test of no scale change might locate a set of genes that was split, with some having larger and others having smaller ts than average.)