

MODULE 8 - PART 2: PATHWAY ANALYSIS LAB

Lab Part 1: Enrichment Test

1.1. Let's get a gene list for analysis

ARTICLE

OPEN

doi:10.1038/nature12634

Mutational landscape and significance across 12 major cancer types

Cyriac Kandoth^{1*}, Michael D. McLellan^{1*}, Fabio Vandin², Kai Ye^{1,3}, Beifang Niu¹, Charles Lu¹, Mingchao Xie¹, Qunyuan Zhang^{1,3}, Joshua F. McMichael¹, Matthew A. Wyczalkowski¹, Mark D. M. Leiserson², Christopher A. Miller¹, John S. Welch^{4,5}, Matthew J. Walter^{4,5}, Michael C. Wendt^{1,3,6}, Timothy J. Ley^{1,3,4,5}, Richard K. Wilson^{1,3,5}, Benjamin J. Raphael² & Li Ding^{1,3,4,5}

The Cancer Genome Atlas (TCGA) has used the latest sequencing and analysis methods to identify somatic variants across thousands of tumours. Here we present data and analytical results for point mutations and small insertions/deletions from 3,281 tumours across 12 tumour types as part of the TCGA Pan-Cancer effort. We illustrate the distributions of mutation frequencies, types and contexts across tumour types, and establish their links to tissues of origin, environmental/carcinogen influences, and DNA repair defects. Using the integrated data sets, we identified 127 significantly mutated genes from well-known (for example, mitogen-activated protein kinase, phosphatidylinositol-3-OH kinase, Wnt/ β -catenin and receptor tyrosine kinase signalling pathways, and cell cycle control) and emerging (for example, histone, histone modification, splicing, metabolism and proteolysis) cellular processes in cancer. The average number of mutations in these significantly mutated genes varies across tumour types; most tumours have two to six, indicating that the number of driver mutations required during oncogenesis is relatively small. Mutations in transcriptional factors/regulators show tissue specificity, whereas histone modifiers are often mutated across several cancer types. Clinical association analysis identifies genes having a significant effect on survival, and investigations of mutations with respect to clonal/subclonal architecture delineate their temporal orders during tumorigenesis. Taken together, these results lay the groundwork for developing new diagnostics and individualizing cancer treatment.

Highlights:

- Using the integrated data sets, we identified 127 significantly mutated genes
- Genes under positive selection, either in individual or multiple tumour types,
- tend to display higher mutation frequencies above background.
- Our statistical analysis identified 127 such genes
- The mutational significance in cancer (MuSiC) package was used to identify
- significant genes for both individual tumour types and the Pan-Cancer collective. [Dees et al. MuSiC: Identifying mutational significance in cancer genomes. *Genome Res.* 2012]
- These significantly mutated genes are involved in a wide range of cellular processes, including transcription factors/regulators, histone modifiers, genome integrity, receptor tyrosine kinase signalling, cell cycle, mitogen-activated protein kinases (MAPK) signalling, phosphatidylinositol-3-OH kinase (PI(3)K) signalling, Wnt/ β -catenin signalling, histones, ubiquitin-mediated proteolysis, and splicing (Fig. 2).

Supplementary Data, Table 4

- globally significant, frequency $\geq 1\%$ for glioblastoma multiforme (GBM): 46

- globally significant, frequency $\geq 1\%$ for kidney renal clear cell carcinoma (KIRC): 53

1.2. Let's use g:Profiler to obtain enrichment results

<http://biit.cs.ut.ee/gprofiler/>

First set the parameters and filter gene sets to be analysed

J. Reimand, M. Kull, H. Peterson, J. Hansen, J. Vilo: g:Profiler -- a web-based toolset for functional profiling of gene lists from large-scale experiments (2007) NAR 35 W193-W200 [PDF]
 J. Reimand, T. Arak, J. Vilo: g:Profiler -- a web server for functional interpretation of gene lists (2011 update) Nucleic Acids Research 2011; doi: 10.1093/nar/gkr378 [PDF]

Organism
 Homo sapiens

Query (genes, proteins, probes)
 PTEN
 TP53
 EGFR
 PIK3R1
 PIK3CA
 NF1
 RB1
 ATRX

Options

- Significant only
- Ordered query
- No electronic GO annotations
- Chromosomal regions
- Hierarchical sorting
- Hierarchical filtering
- Show all terms (no filtering) [v]
- Output type**
 Generic Enrichment Map (TAB)
- Hide advanced options
- Evidence codes in txt output
- Measure underrepresentation
- Gene list as a stat. background
- 1.00 User p-value
- Size of functional category
 10 [v] 1000 [v]
- Size of query / term intersection
 3 [v]
- Numeric IDs treated as
 AFFY_HUGENE_2_0_ST_V1 [v]
- Significance threshold**
 g:SCS threshold
- Statistical domain size**
 Only annotated genes
- Download g:Profiler data as GMT:
 ENSG, name

Gene Ontology

- Biological process**
- Cellular component
- Molecular function
- Inferred from experiment [IDA, IPI, IMP, IGI, IEP]
- Direct assay [IDA] / Mutant phenotype [IMP]
- Genetic interaction [IGI] / Physical interaction [IPI]
- Traceable author [TAS] / Non-traceable author [NAS] / Inferred by curator [IC]
- Expression pattern [IEP] / Sequence or structural similarity [ISS] / Genomic context [IGC]
- Biological aspect of ancestor [IBA] / Rapid divergence [IRD]
- Reviewed computational analysis [RCA] / Electronic annotation [IEA]
- No biological data [ND] / Not annotated [NA]
- Biological pathways**
- KEGG
- Reactome
- Regulatory motifs in DNA
- TRANSFAC TFBS
- miRBase microRNAs
- Protein databases
- Human Protein Atlas
- CORUM protein complexes
- Human Phenotype Ontology (sequence homologs in other species)
- Online Mendelian Inheritance in Man
- BioGRID protein-protein interaction

g:Profile!

g:Profiler version 1.536_e83_eg30

Then paste in the gene list and press g:Profile to perform the analysis; also download the GMT file with gene symbols, which will be necessary to use Enrichment Map for visualization

J. Reimand, M. Kull, H. Peterson, J. Hansen, J. Vilo: g:Profiler -- a web-based toolset for functional profiling of gene lists from large-scale experiments (2007) NAR 35 W193-W200 [PDF]
 J. Reimand, T. Arak, J. Vilo: g:Profiler -- a web server for functional interpretation of gene lists (2011 update) Nucleic Acids Research 2011; doi: 10.1093/nar/gkr378 [PDF]

Organism
 Homo sapiens

Query (genes, proteins, probes)
 PTEN
 TP53
 EGFR
 PIK3R1
 PIK3CA
 NF1
 RB1
 ATRX

Options

- Significant only
- Ordered query
- No electronic GO annotations
- Chromosomal regions
- Hierarchical sorting
- Hierarchical filtering
- Show all terms (no filtering) [v]
- Output type**
 Generic Enrichment Map (TAB)
- Hide advanced options
- Evidence codes in txt output
- Measure underrepresentation
- Gene list as a stat. background
- 1.00 User p-value
- Size of functional category
 10 [v] 1000 [v]
- Size of query / term intersection
 3 [v]
- Numeric IDs treated as
 AFFY_HUGENE_2_0_ST_V1 [v]
- Significance threshold**
 g:SCS threshold
- Statistical domain size**
 Only annotated genes
- Download g:Profiler data as GMT:
 ENSG, name

Gene Ontology

- Biological process**
- Cellular component
- Molecular function
- Inferred from experiment [IDA, IPI, IMP, IGI, IEP]
- Direct assay [IDA] / Mutant phenotype [IMP]
- Genetic interaction [IGI] / Physical interaction [IPI]
- Traceable author [TAS] / Non-traceable author [NAS] / Inferred by curator [IC]
- Expression pattern [IEP] / Sequence or structural similarity [ISS] / Genomic context [IGC]
- Biological aspect of ancestor [IBA] / Rapid divergence [IRD]
- Reviewed computational analysis [RCA] / Electronic annotation [IEA]
- No biological data [ND] / Not annotated [NA]
- Biological pathways**
- KEGG
- Reactome
- Regulatory motifs in DNA
- TRANSFAC TFBS
- miRBase microRNAs
- Protein databases
- Human Protein Atlas
- CORUM protein complexes
- Human Phenotype Ontology (sequence homologs in other species)
- Online Mendelian Inheritance in Man
- BioGRID protein-protein interaction

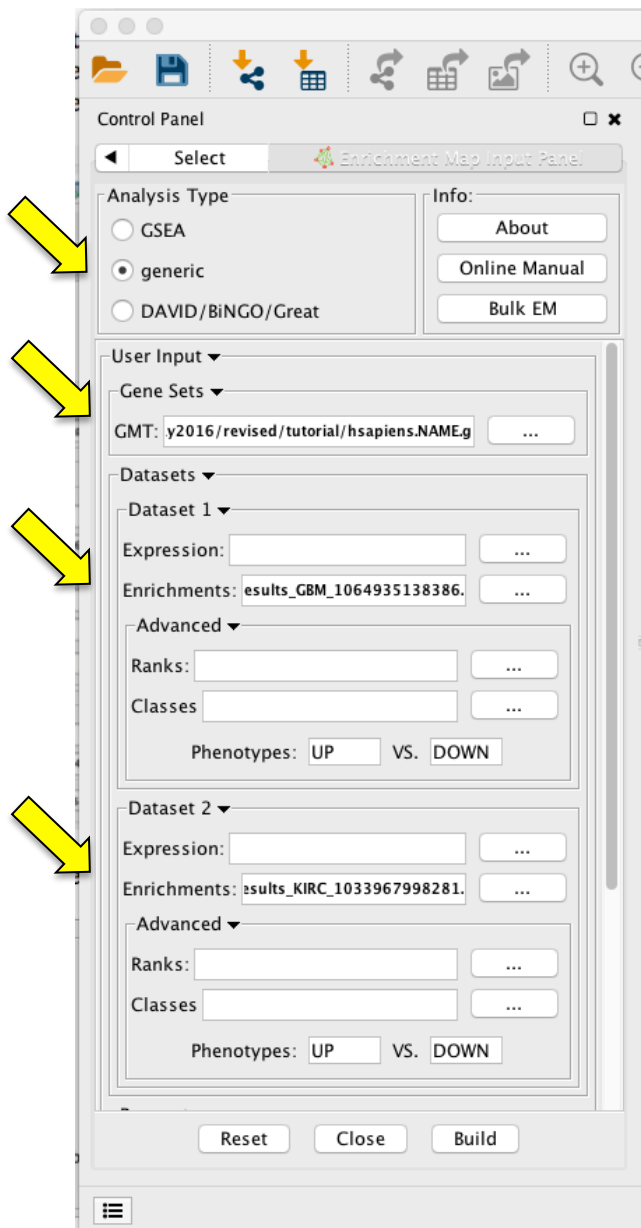
g:Profile!

g:Profiler version 1.536_e83_eg30

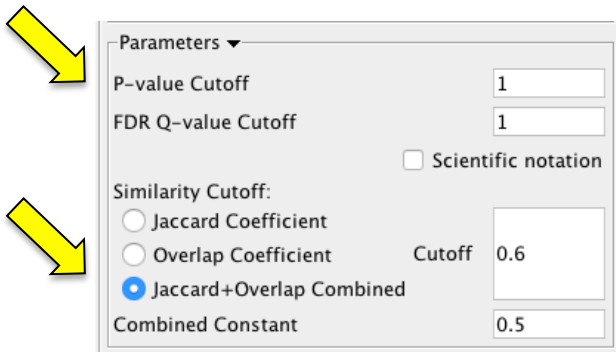
Open Cytoscape 3.2.1, Apps > Enrichment Map > Create enrichment map

First load all the files:

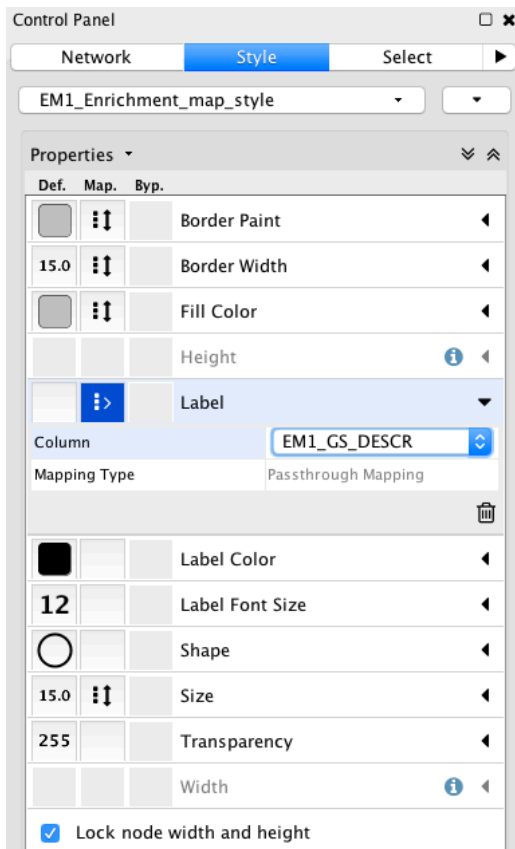
- GMT file: has gene-set definitions
- Enrichment of data-set 1: has enrichment statistics for GBM
- Enrichment of data-set 2: has enrichment statistics for KIRC



Then set the analysis parameters



From the Cytoscape control panel, select the style tab, and map the EM1_GS_DESCR to graphic attribute Label



Then, for both data-sets, reset the color based on FDR.

First, change the attribute mapped to the graphical attribute:

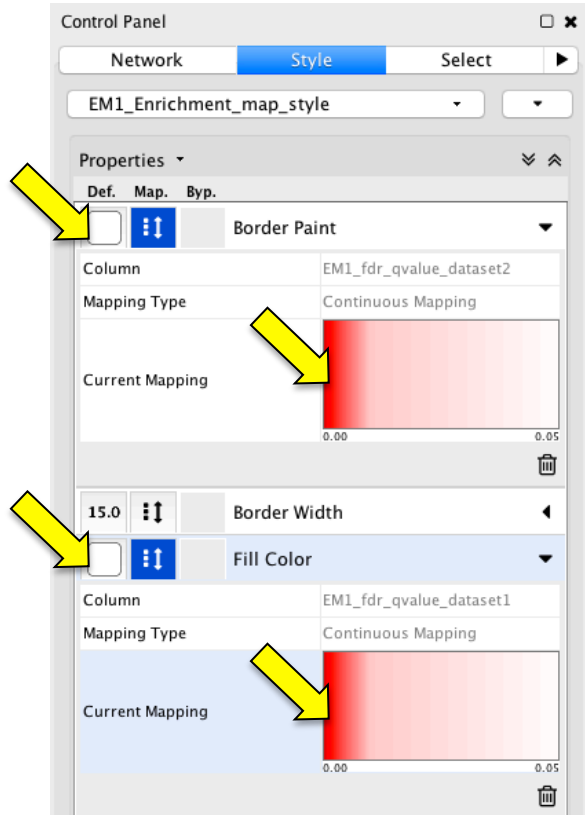
- fill color: EM1_fdr_qvalue_dataset1
- border paint: EM1_fdr_qvalue_dataset2

Second, set the color notches to these values:

- 0.0001: dark red
- 0.001: red

- 0.005: lighter red
- 0.01: lightest red
- 0.05: white

Third, set the default colors of nodes and borders to white.



Finally, we are going to improve the layout.

1. From: Layout > Settings
2. Select: Prefuse Force Directed Layout
3. And reset default spring coefficient to: 1E-6

Tips for interpreting results for these data-sets:

- Entirely red circles represent pathways enriched both in GBM and KIRC
- Red cores with white borders represent pathways only seen in GBM
- White cores with red borders represent pathways only seen in KIRC
- Intensity of red maps to strength of pathway enrichment p-value
- Singletons at the bottom are often redundant with larger clusters; however they sometimes include additional, unique pathways.
- Browse groups of pathways and identify major functional themes characteristic of these groups.