

Comparative analysis of microbial genomes

Focus on *Mycobacterium
tuberculosis*

Nicky Mulder

University of Cape Town

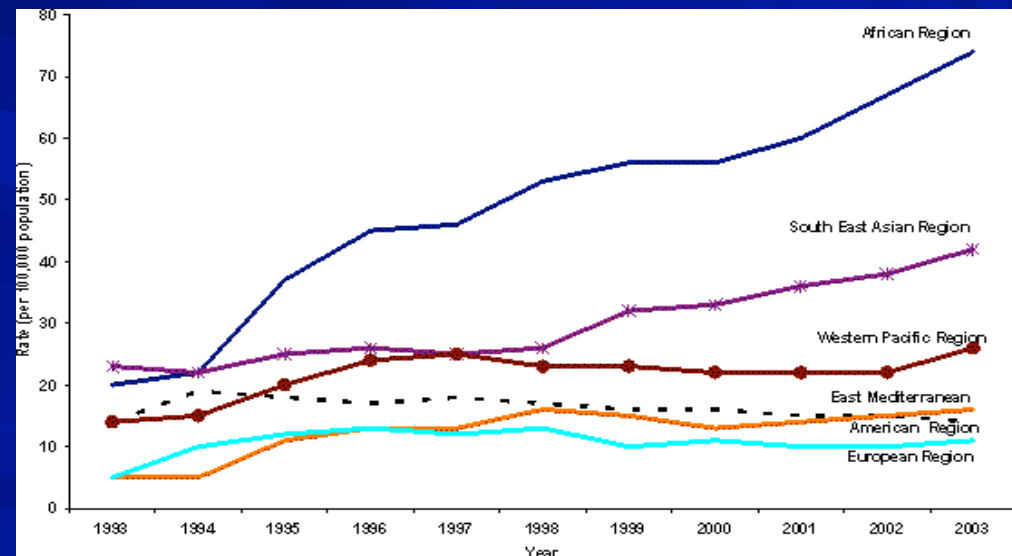
European Bioinformatics Institute

Outline

- Tuberculosis background
- Mycobacterial genomics
- Comparative analysis of microbial genomes
- Identification of expanded families in *M. tuberculosis*
- Future work

Background

- 1/3 world's population is infected with tuberculosis
- ~1000 people die of TB in South Africa daily
- Increased incidence due to AIDS
- Caused by *Mycobacterium tuberculosis*



World TB notification rate per 100,000 population size
(<http://www.afro.who.int/tb/notificationtrend.html>)

Mycobacterial genome sequencing

- Hundreds of microbial genomes completely sequenced
- Mycobacteria sequenced:
 - Mycobacterium leprae
 - Mycobacterium tuberculosis (H37Rv, Oshkosh, F11)
 - Mycobacterium bovis (& BCG)
 - Mycobacterium paratuberculosis
 - Mycobacterium (sp. MCS, KMS, JLS)
 - Mycobacterium smegmatis
 - Mycobacterium avium
 - Mycobacterium ulcerans
 - Mycobacterium vanbaaleni

Mycobacterium tuberculosis H37Rv

- 4.4Mb genome
- ~4000 proteins
- 51% have predicted functions
- Approximately half the genome is made up of duplicated families
- Has many lipid and fatty acid metabolism pathways
- Still so much not known about the remaining genes and pathogenesis

Comparative genomics

- Aim: generate phylogenetic profiles to identify potential virulence genes
- 84 genomes, ~250,000 genes
- 56 pathogens 28 non-pathogens
- BLASTClust at different score densities
- Score density 0.5 and 1.0

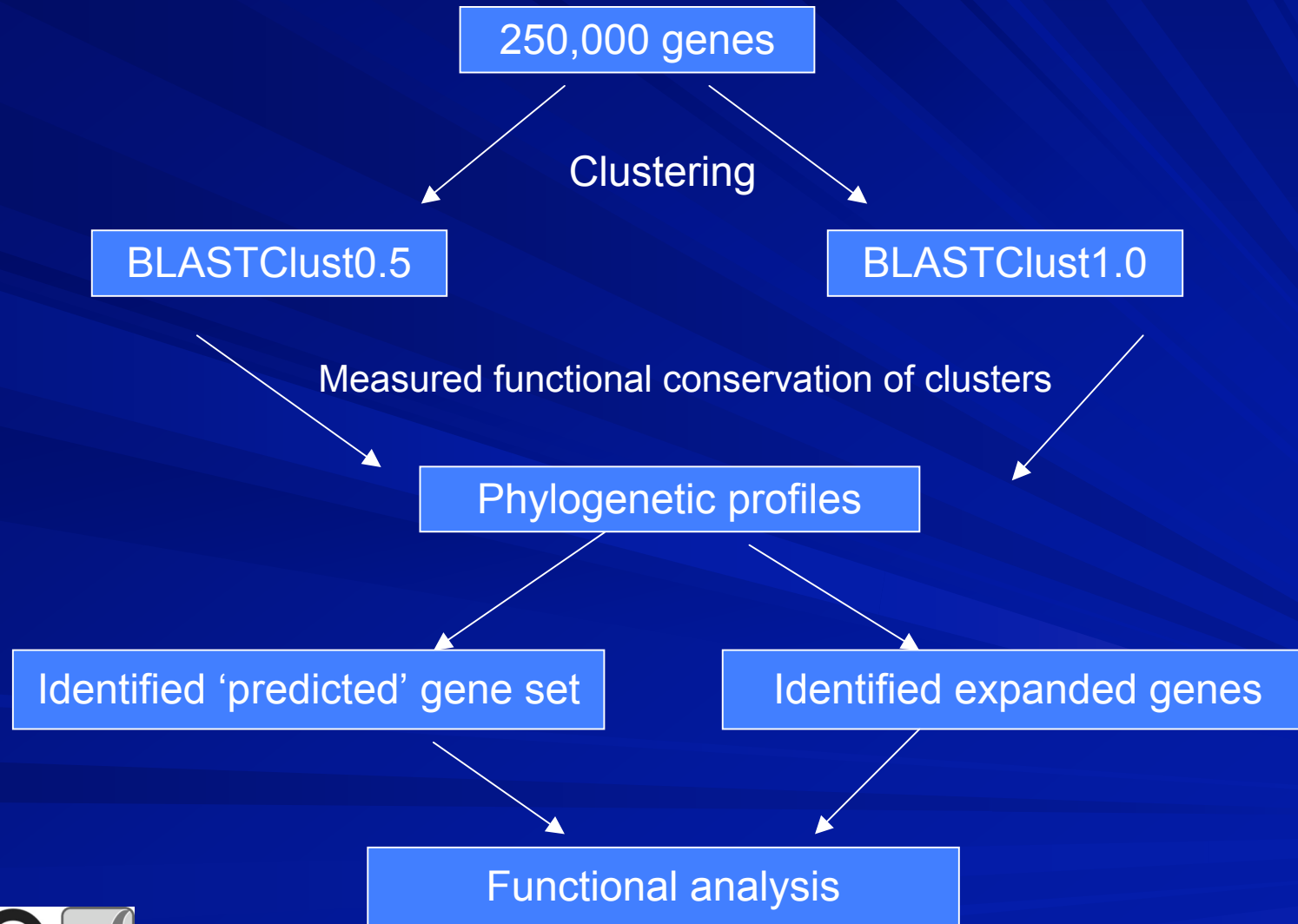
$$S = -N/AL'$$

L' is length of sequence in the alignment, L

N is the number of identical residues

A is the total alignment length

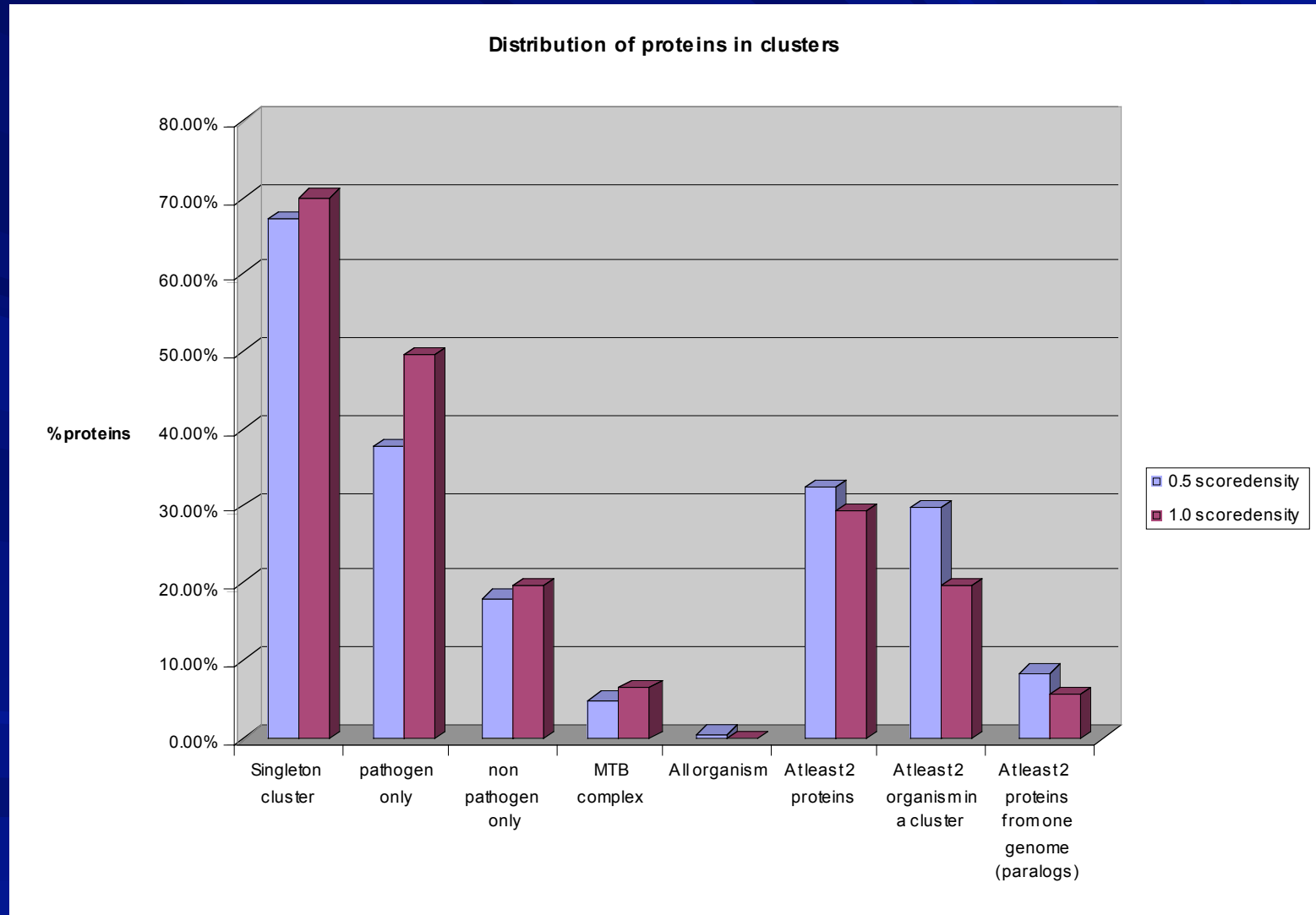
Methods



Summary of clustering results

Experiments	0.5 score density	1.0 score density
Total clusters created	63819	121905
Total Singletons (one protein families)	42770	85735
Total more than one protein clusters	20849	36170
Total Pathogen only clusters	7904	18022
Total non Pathogen only clusters	3824	7165
Total MTB Complex only clusters	1060	2362
Number with at least two genomes per cluster	19056	24110
Number of clusters with at least 2 proteins (paralogs) in one genome	1794	2061
All organisms (84 genomes)	10	2
At least 80 genomes	135	39
All organism (at least 60 genomes)	448	130

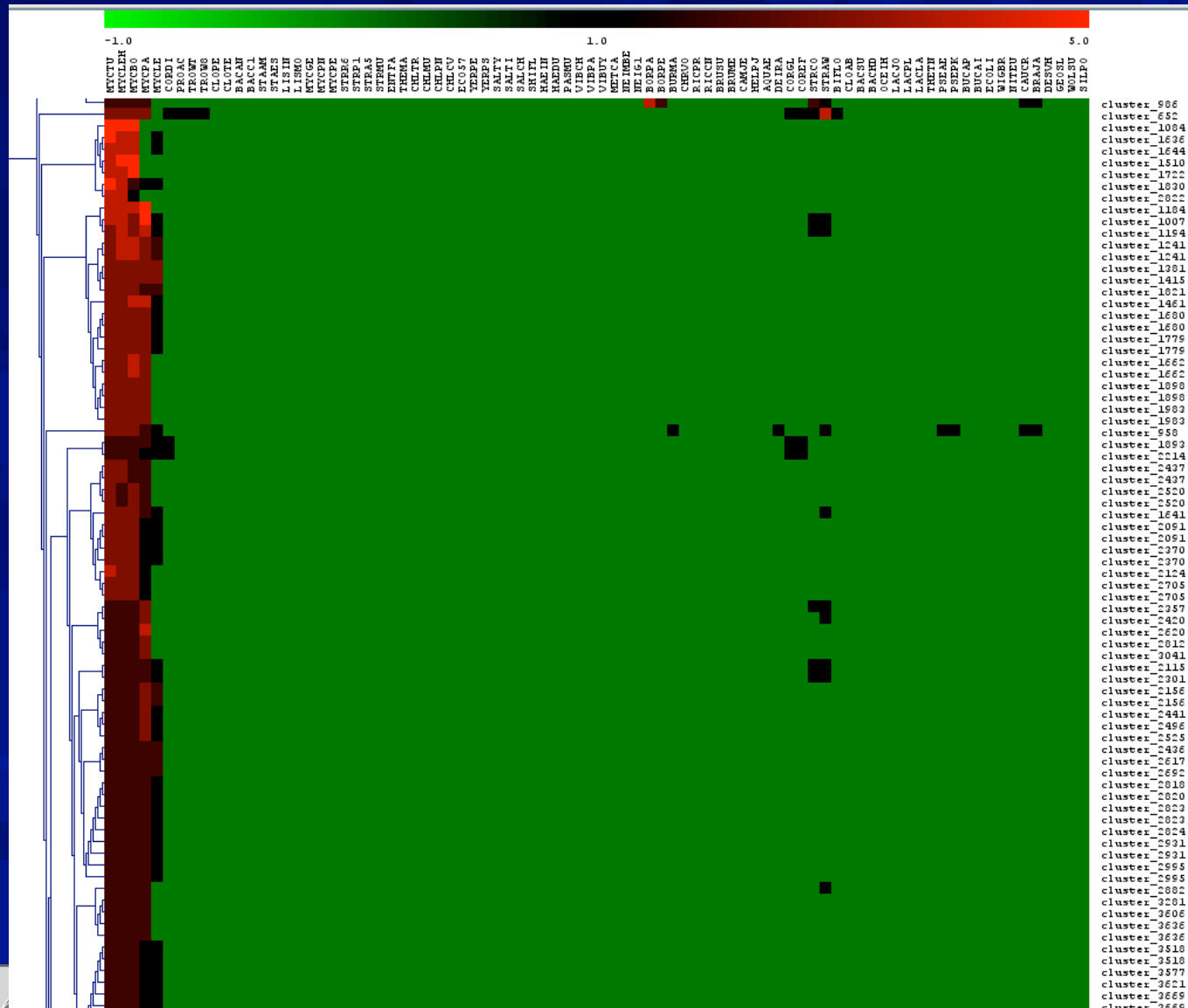
Distribution of proteins in cluster types



Measuring functional diversity within clusters

- Used DE lines, InterPro matches and GO Slim
- Can only be done for those covered:
~70% InterPro, 50% GO slim
- BlastClust1.0 greatest conservation –most
>90% conservation of InterPro domains
and annotation
- In 1.0 PE/PPE families agreed with
previous reports

Heatmap of results 2

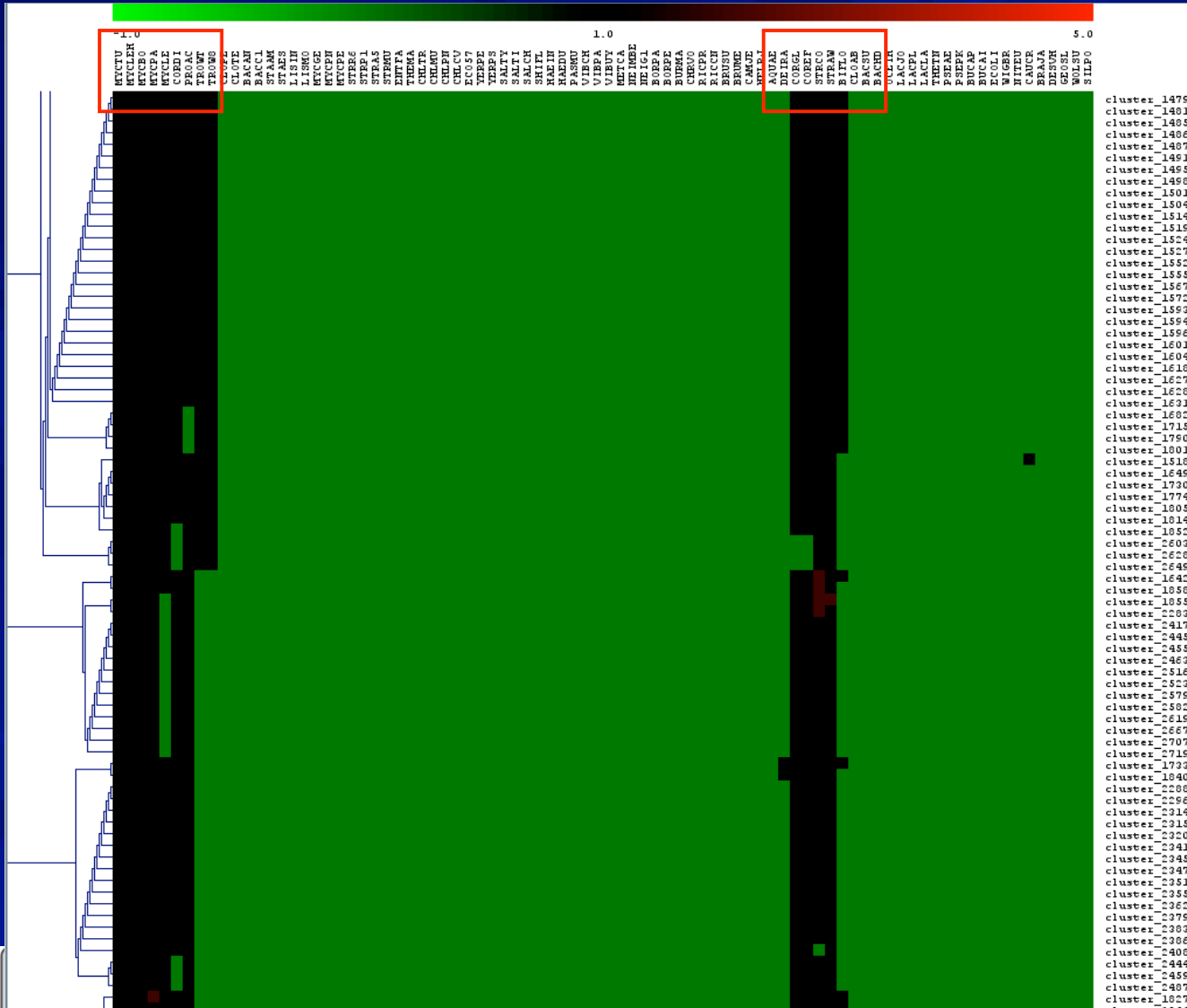


Heatmap of results 3

Mycobacteria

Corynebacteria

Tropheryma



Streptomyces
Corynebacteria

Clustering by organism

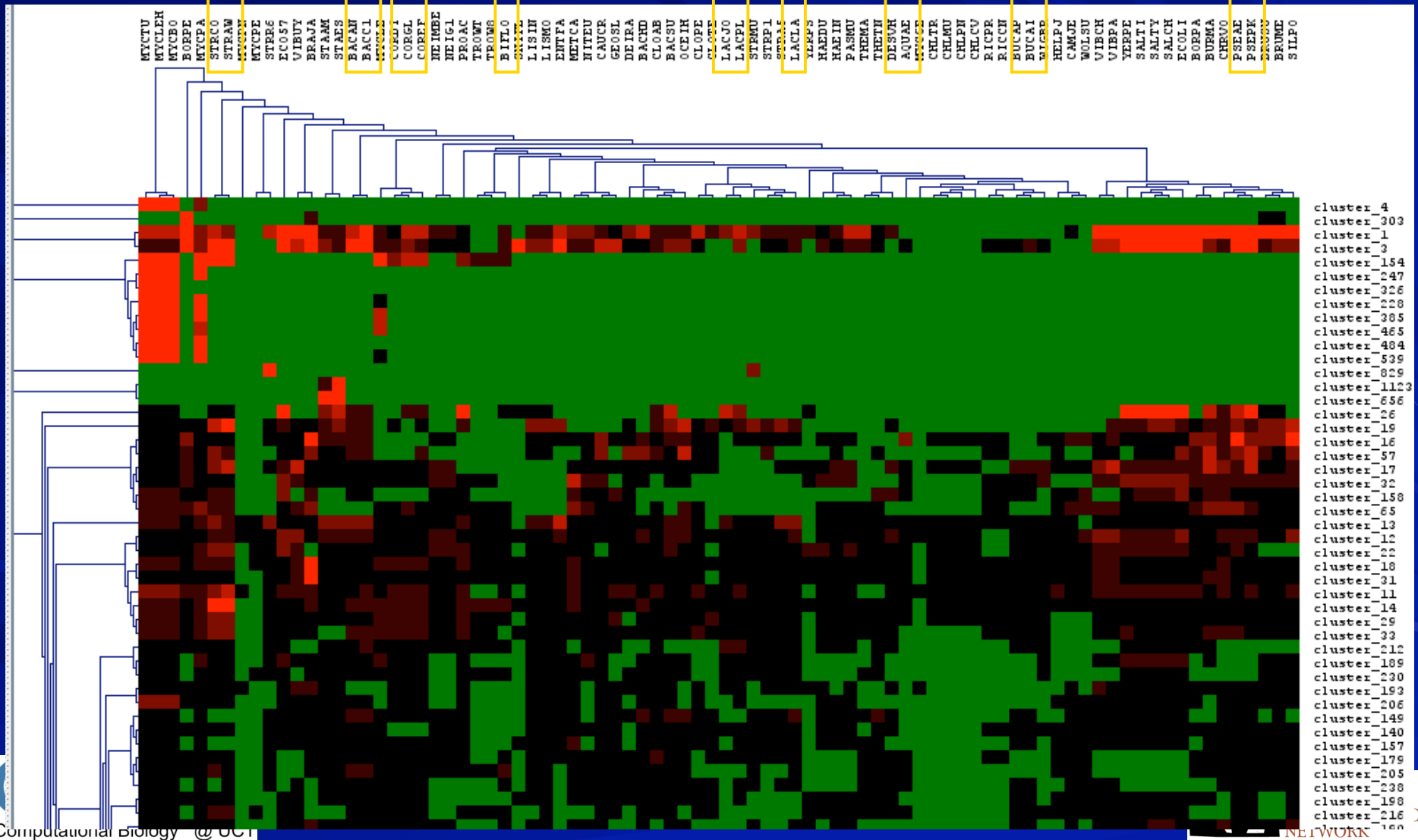
High GC G+

Gamma proteo

High GC G+

Low GC G+

Gamma proteo

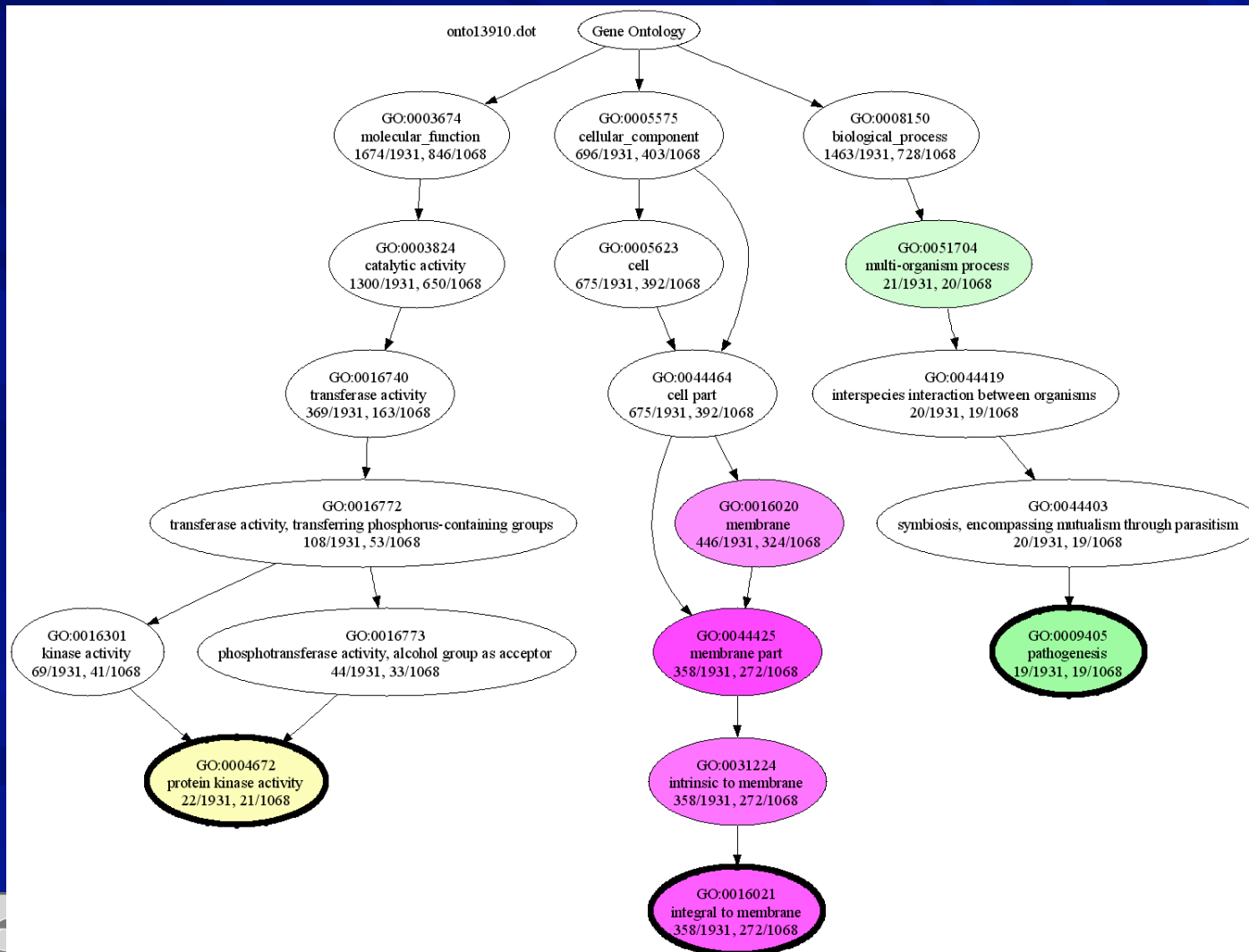


Known virulence genes

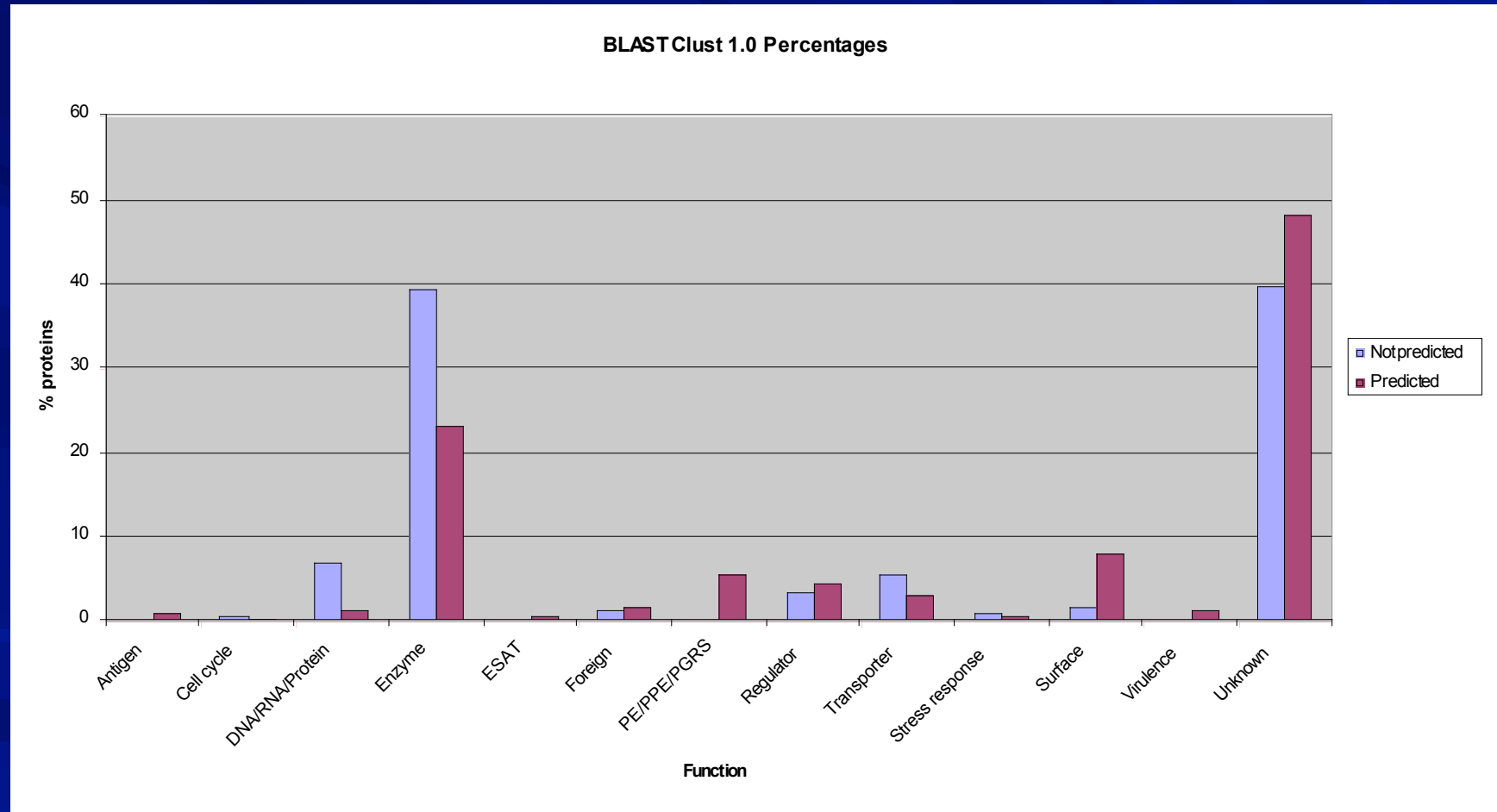
Cluster No	Proteins in pathogens (no. organisms)	Proteins in non pathogens (no. organisms)	Common DE line
652	16 (8)	8 (5)	DNA-binding_response_regulator
1583	7 (7)	7 (4)	Mycobacterial persistence regulator MRPA
2441	10 (5)	0	Virulence_factor/mce-family_protein
2613	6 (6)	4 (4)	Cytotoxin_/haemolysin_homologue
2818	9 (5)	0	Virulence_factor_mce_family_protein
2820	9 (5)	0	Mce-family_protein_mce2d
2822	9 (3)	0	Phospholipase_C_1_precursor_(EC_3.1.4.3)(MTP40_antigen)
2823	9 (5)	0	Virulence_factor_mce_family_protein
2931	9 (5)	0	Mce-family_protein_mce1b_
3332	8 (5)	0	Virulence_factor_mce_family_protein
7120	5 (5)	0	Heparin-binding_hemagglutinin_(Adhesin)
8232	4 (4)	0	Sulfatase family protein
8234	4 (4)	0	Virulence_factor_mce_family_protein
8935	4 (4)	0	Exported_repetitive_protein_precursor_(Cell_surface_protein_pirG)
11724	3 (3)	0	Virulence_factor_mce_family_protein
10914	3 (3)	0	Virulence-regulating_(arac/xyls_f

Functional analysis

Ontologizer –statistically over-represented GO terms in predicted set

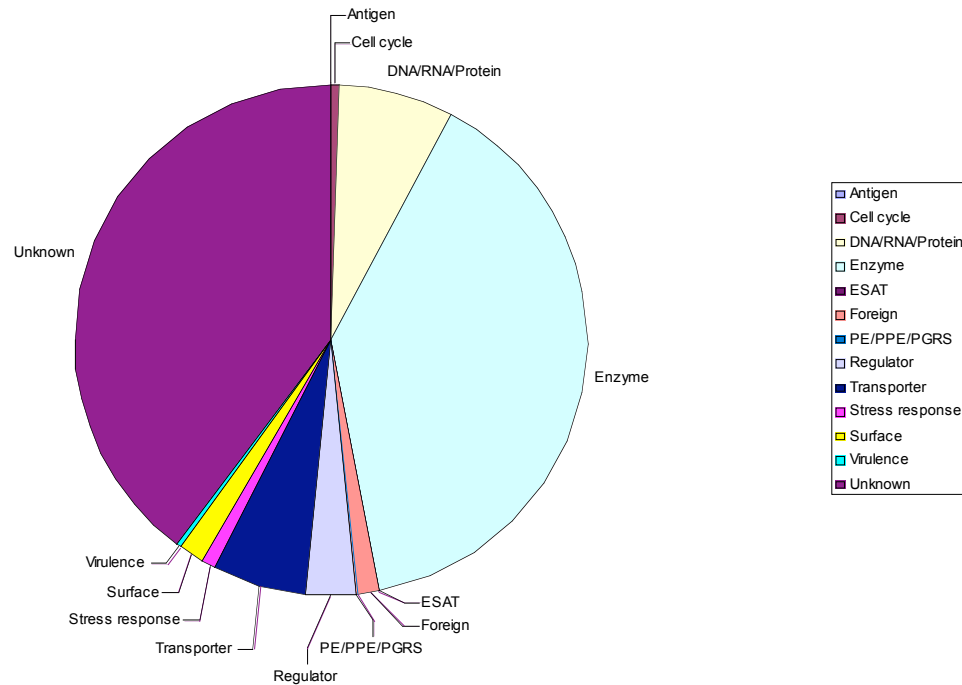


Functional analysis continued



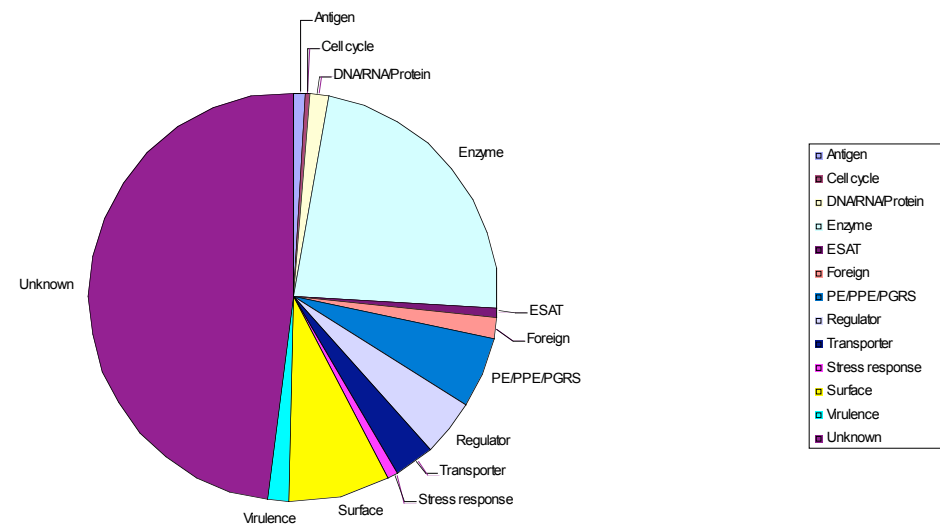
Functional composition of protein sets

BLASTClust 1.0 Not predicted % proteins per function



Both have high proportion of hypothetical proteins

BLASTClust 1.0 Predicted % proteins per function



Hypothetical proteins

InterPro Signatures

IPR001054 PF00211 Guanylate_cyc
 IPR001054 PS5012 GUANYLATE_CYCLASE_2
 IPR001054 SM00044 CYCc
 IPR003660 PF00672 HAMP
 IPR003660 PS50885 HAMP
 IPR003660 SM00304 HAMP

Structural Predictions

MB_P63527

Hypothetical protein Rv1320c/MT1362

InterPro Signatures

IPR001054 PF00211 Guanylate_cyc
 IPR001054 PS5012 GUANYLATE_CYCLASE2
 IPR001054 SM00044 CYCc
 IPR003660 PF00672 HAMP
 IPR003660 PS50885 HAMP
 IPR003660 SM00304 HAMP

Structural Predictions

MB_Q10633

Adenylate cyclase, putative

InterPro Signatures

IPR001054 PF00211 Guanylate_cyc
 IPR001054 PS5012 GUANYLATE_CYCLASE2
 IPR001054 SM00044 CYCc
 IPR003660 PF00672 HAMP
 IPR003660 PS50885 HAMP
 IPR003660 SM00304 HAMP

Structural Predictions

MB_Q10633

InterPro Signatures

IPR003455 PF02409 Omt_N
 IPR011610 TIGR00027 mthyl_TIGR00027

Structural Predictions

MB_053686

InterPro Signatures

IPR003455 PF02409 Omt_N
 IPR011610 TIGR00027 mthyl_TIGR00027

Structural Predictions

MB_053841

InterPro Signatures

IPR003455 PF02409 Omt_N
 IPR011610 TIGR00027 mthyl_TIGR00027

Structural Predictions

MB_086359

InterPro Signatures

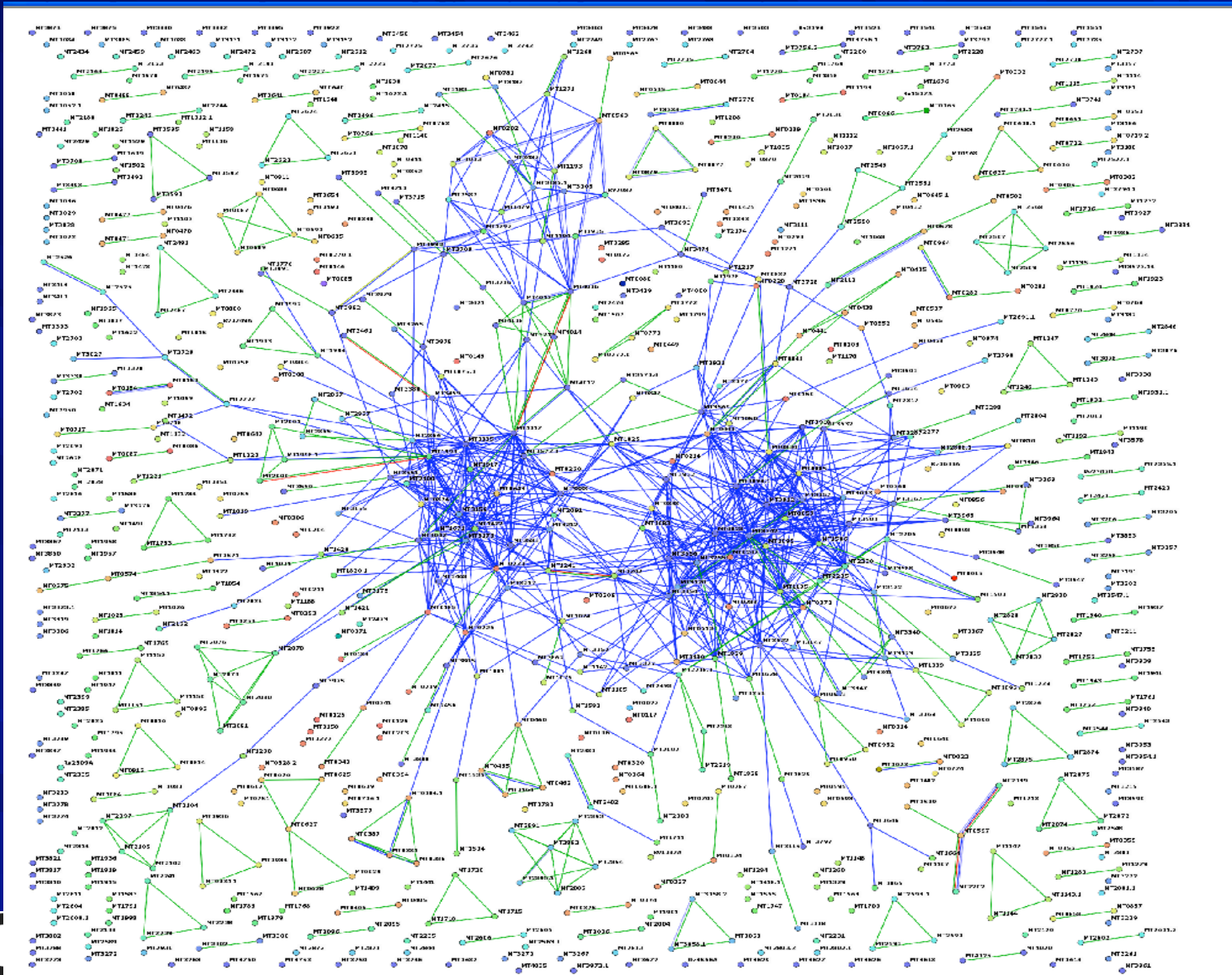
IPR003455 PF02409 Omt_N
 IPR011610 TIGR00027 mthyl_TIGR00027

Structural Predictions

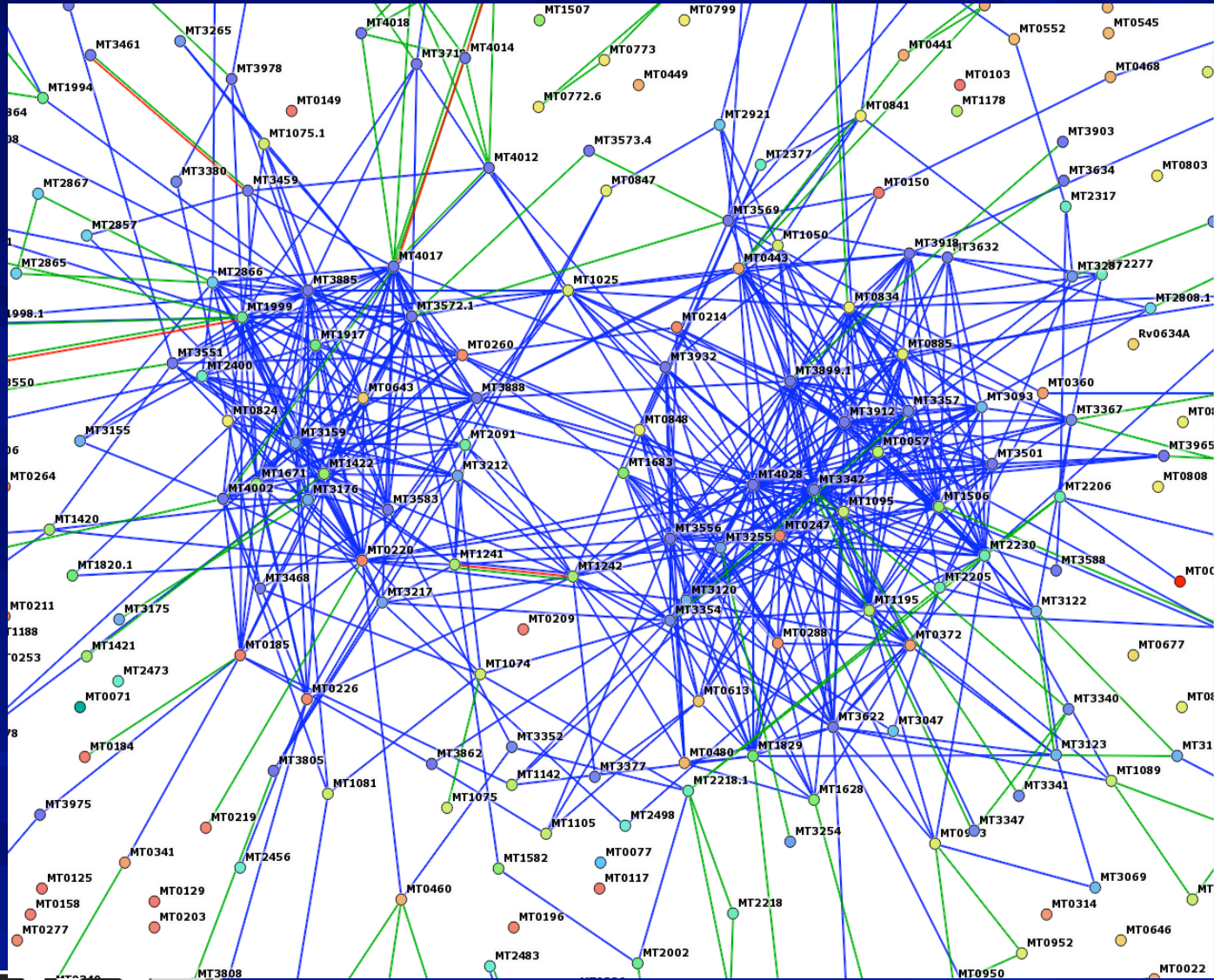
MB_P64747

InterPro Signatures

Connections between hypotheticals using STRING



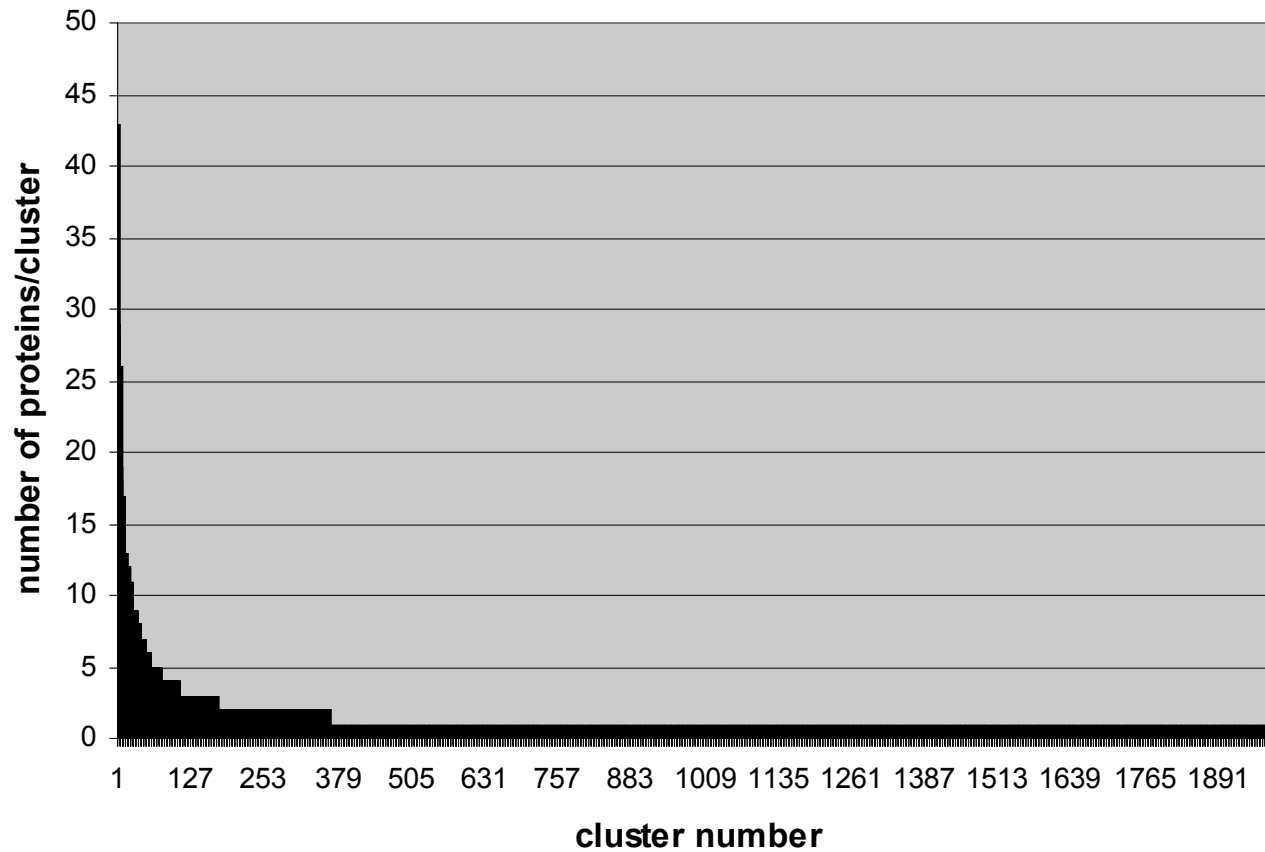
Zooming in



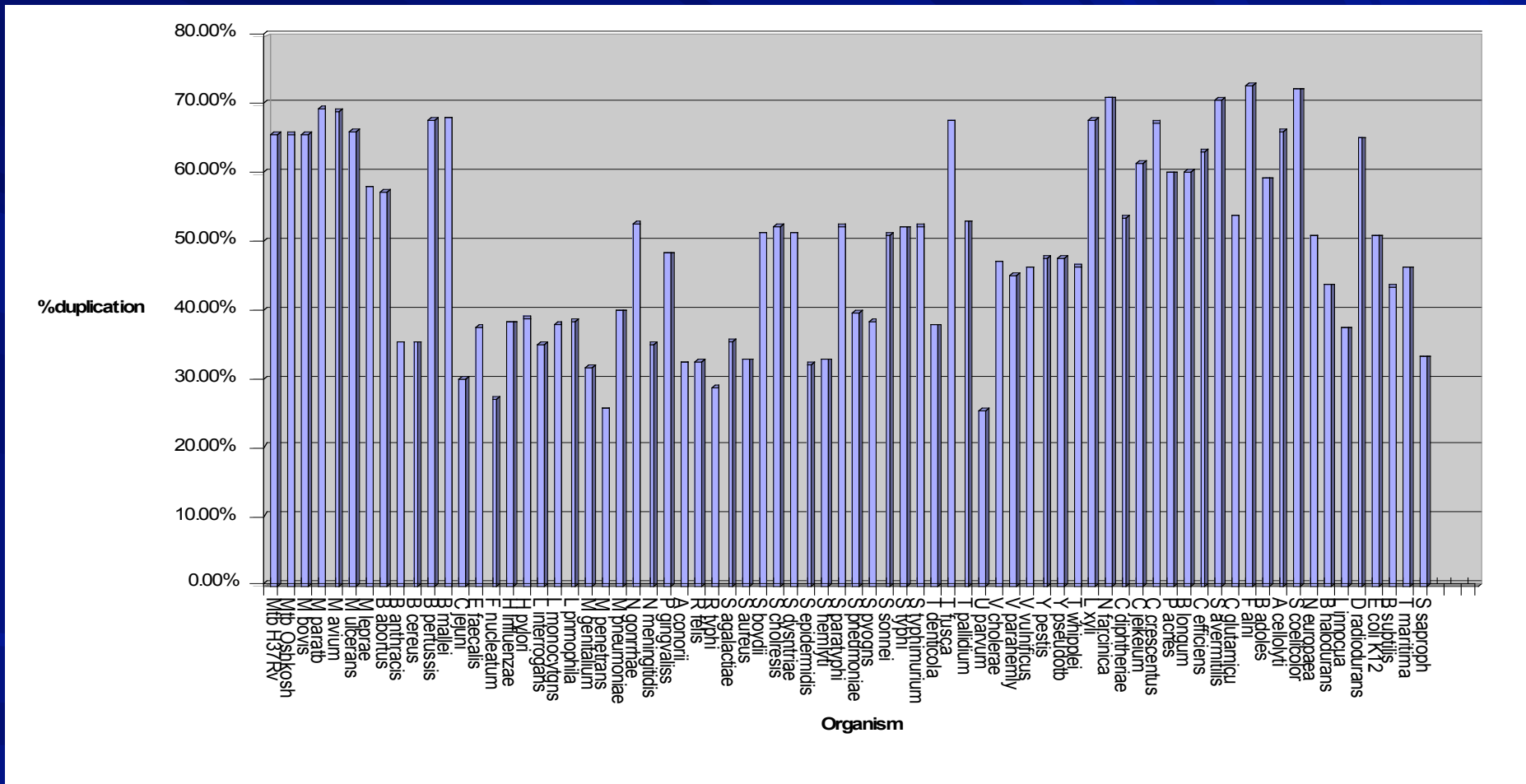
- Neighborhood
- Gene Fusion
- Cooccurrence
- Coexpression
- Experiments
- Databases
- Textmining
- [Homology]

Expanded families in H37Rv

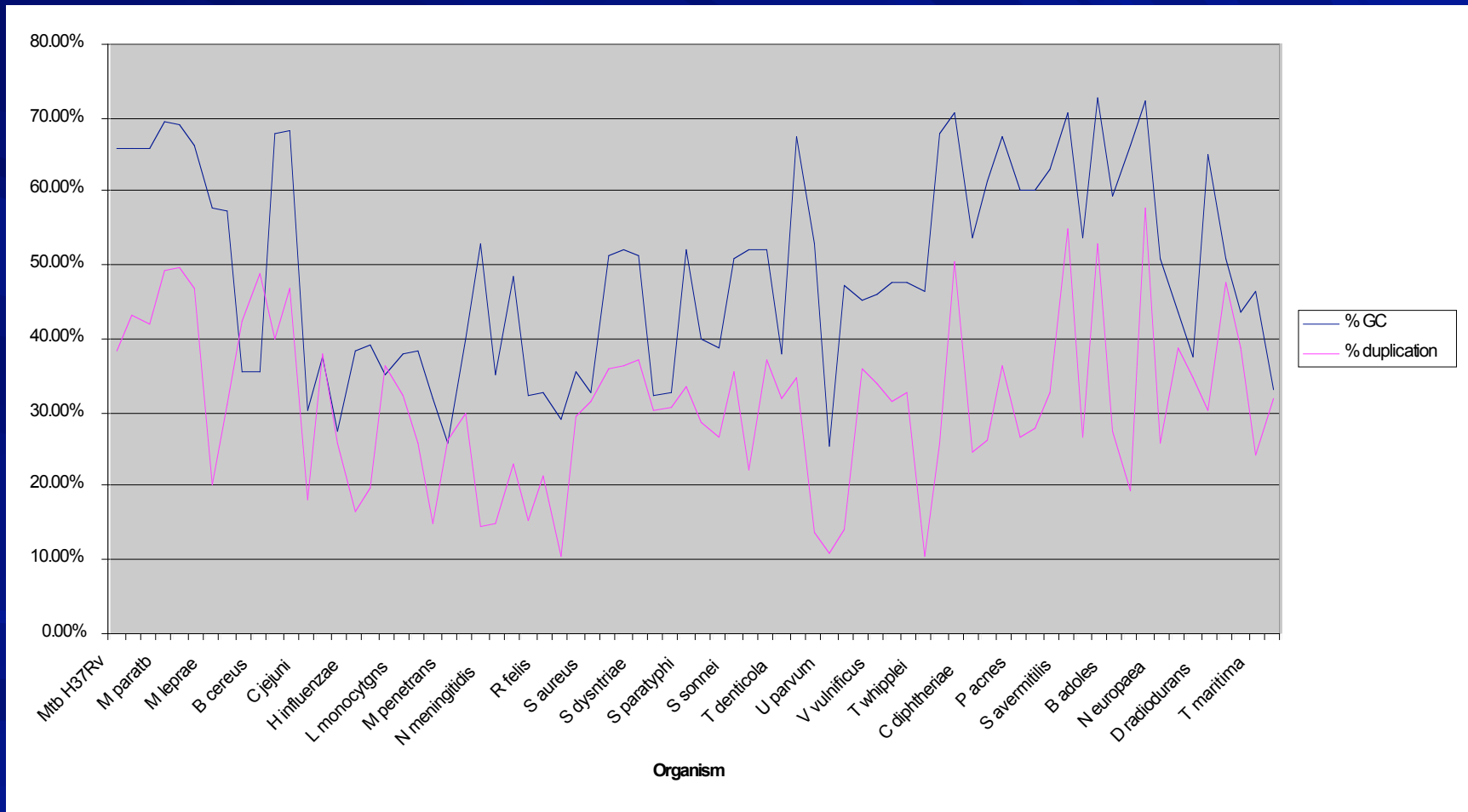
Paralogs from the BLAST results grouped into clusters for *Mycobacterium tuberculosis* H37RV



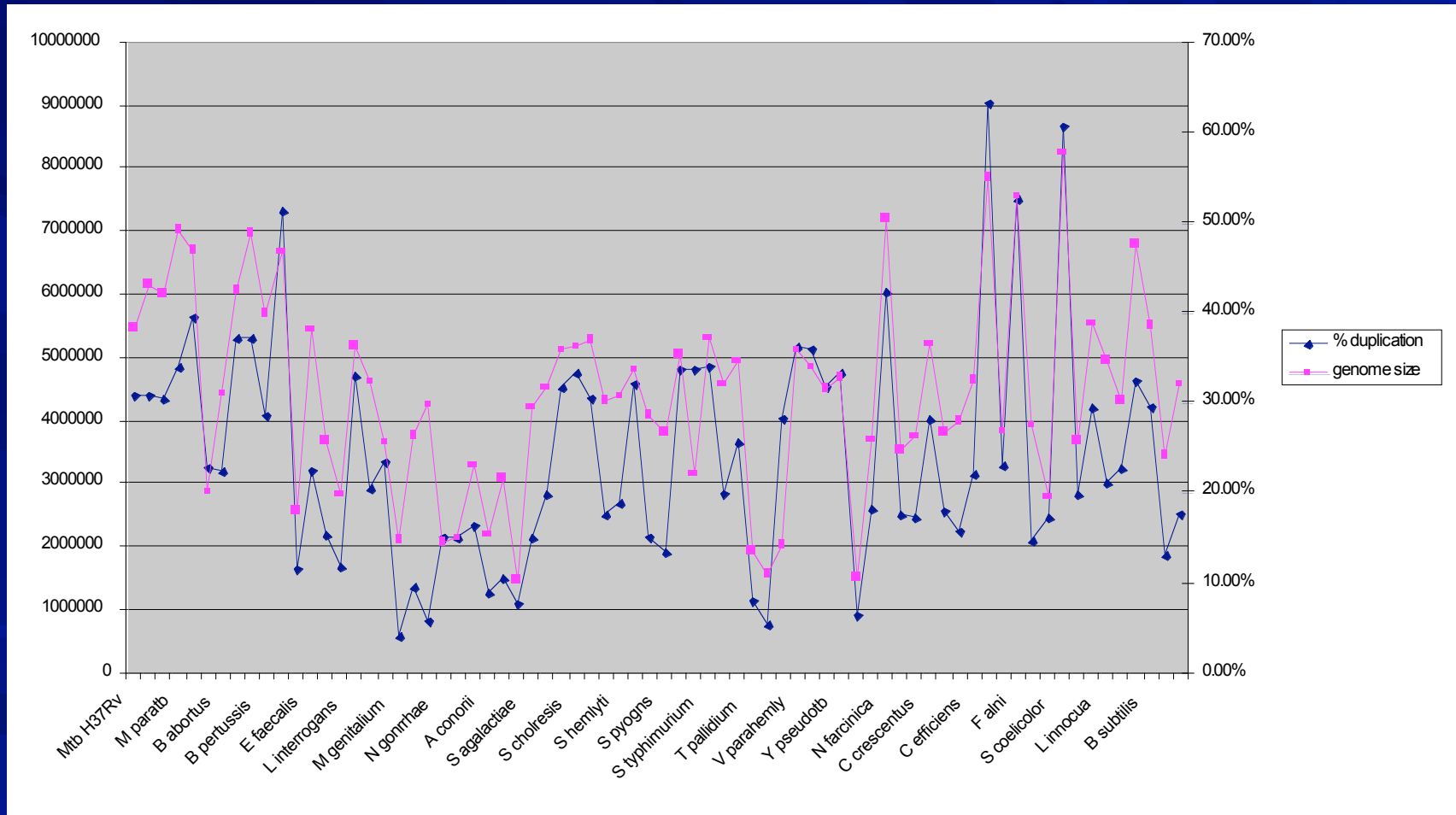
Duplication in multiple genomes



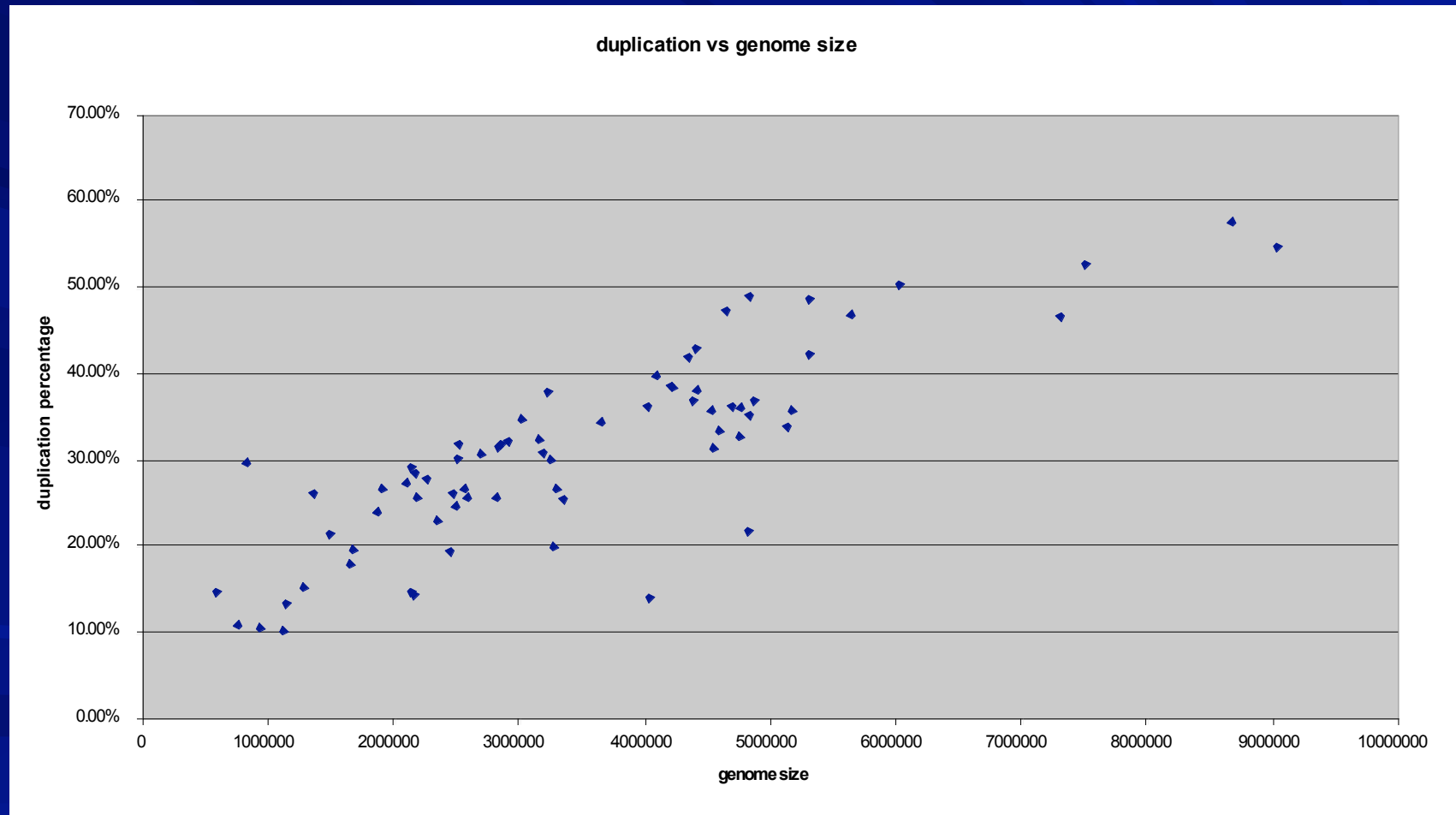
Relationship between GC content and % duplication



Relationship between genome size and % duplication



Relationship between genome size and % duplication



Duplicated families in *M.tb* and *M. leprae*

Combined InterPro and BLAST results gave 354 clusters

354 clusters

Duplicated families in *M.tb* and *M. leprae*

Combined InterPro and BLAST results gave 354 clusters

More in *M. leprae* than TB (hypothetical, WhiB, conserved membrane, inv)



Duplicated families in *M.tb* and *M. leprae*

Combined InterPro and BLAST results gave 354 clusters

More in *M. leprae* than TB (hypothetical, WhiB, conserved membrane, inv)



Future work

- Investigate hypothetical protein families further
- Determine functional fate of expanded families that are unique to TB or TB complex
- Measure diversity vs cluster size
- Compare diversity in upstream and coding regions
- Look at expression correlation of expanded families

Acknowledgements

- Halimah Rabiou
- Venu Vuppu
- Gordon Jamieson
- NBN funding

