

DATA INTEGRATION AND LINKING OF GENOMIC, METABOLOMIC AND TRANSCRIPTOMIC DATA USING MAPMAN

14. FEB. 2024 | BJÖRN USADEL

IBG-4, FORSCHUNGSZENTRUM JUELICH, GERMANY



Forschungszentrum Jülich

1.7 km² research campus 11 institutes and over 80 institute departments,
7,100 people



HHU Düsseldorf Plant Excellence cluster

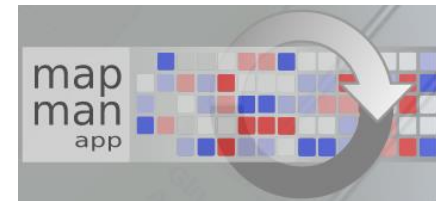


Slide 2



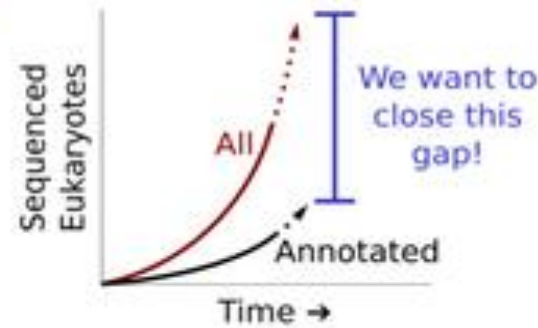
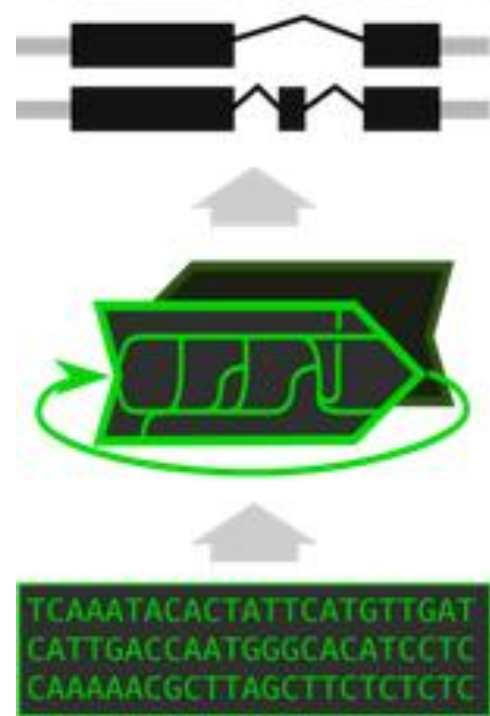
GWAS, GENOMES AND THEN?

- Gene Finding with **Helixer**
- Functional Gene Prediction with **Mercator4**
- Visualisation of Omics Data with **MapMan**

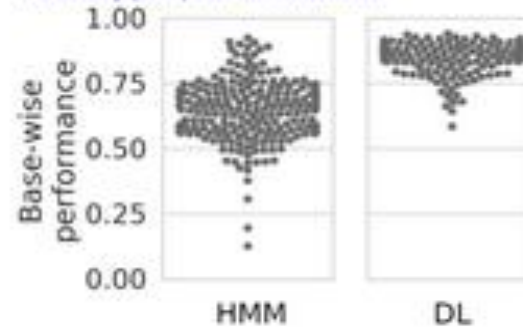


FINDING GENES IS STILL HARD

Gene prediction with Deep Learning



Prototype performance



Who has not had to deal with V1, V2 etc of gene annotation? Gene finding is hard.

Helixer uses AI to help – and often outperforms simple “full length” RNA sequencing endeavors

We use both but nothing is perfect. We can still find and improve genes in Arabidopsis now

HELIXER

about

gene annotation

Job submission (beta version)

- upload a file with nucleotide sequences (one or many records in a valid FASTA format)
 - minimum sequence length of a single record: 25 kbp
 - maximum file size (including all records): 1 GByte
 - optionally, compress the file ('.gz', '.zip' and '.bz2' are supported)
- select the organism lineage 'land plants', 'vertebrates', 'invertebrates' or 'fungi' (required for lineage-specific annotation)
- optionally, provide an email address (where a link to the GFF3-formatted gene annotation results will be sent)
- optionally, specify a label (used as prefix in the GFF3-formatted gene annotation results to specify the source of the nucleotide sequence)

Helixer beta version

Please note:

- We are frequently applying improvements which may impact the functionality of Helixer. If you are experiencing difficulties, please try again later.
- We currently support compressed (.gz .zip .bz2) and uncompressed input files.
- As this is a beta version, we may update aspects of the tool which could change the output generated.

Thank you for using Helixer

The Helixer team

Upload nucleotide sequence file (FASTA format) Keine Datei ausgewählt.

Use demo file - Arabidopsis_lyrata.v.1.0.dna.chromosome.8.fa.gz

Select Lineage-specific mode

Enter label for GFF feature prefix ⓘ



Meant as a bioinformatics downloadable tool

but “small” genomes ca 1 GBase i.e. most berry crops can be done online



WHAT DO ALL THE GENES/PROTEINS DO?

Job submission

[Result tree viewer](#)

[Result Heatmap viewer](#)

Sequence type Protein DNA 

Include Prot-scriber annotations (Beta version) 

Include Swissprot annotations 

Upload FASTA file Keine Datei ausgewählt.

Use demo FASTA file

Job name 

Email address 



COMPARATIVE VIEW OF GENE NUMBERS

▼ 9 | Secondary metabolism

▼ 9.1 | terpenoids

↳ mevalonate (MVA) pathway

- ▣ 3 2 2 2 acetyl-CoA C-acyltransferase *(ACAT1/2)
- ▣ 3 1 1 1 3-hydroxy-3-methylglutaryl-CoA synthase *(HMGS)
- ▣ 3 2 3 1 3-hydroxy-3-methylglutaryl-CoA reductase *(HMGR)
- ▣ 0 1 1 1 mevalonate kinase *(MVK)
- ▣ 2 1 1 1 phosphomevalonate kinase *(PMK)
- ▣ 1 2 1 1 mevalonate diphosphate decarboxylase *(MVD1/2)
- ▣ 1 2 2 1 isopentenyl diphosphate isomerase *(IDI1/2)

▼ methylerythritol phosphate (MEP) pathway

↳ 1-deoxy-d-xylulose 5-phosphate import

- ▣ 1 1 1 1 D-xylulose kinase
- ▣ 1 1 1 0 D-xylulose 5-phosphate transporter
- ▣ 7 3 4 2 1-deoxy-D-xylulose 5-phosphate synthase *(DXS)
- ▣ 2 1 1 2 1-deoxy-D-xylulose 5-phosphate reductase *(DXR)
- ▣ 0 1 1 1 4-diphosphocytidyl-2-C-methyl-D-erythritol synthase
- ▣ 1 1 1 1 4-diphosphocytidyl-2-C-methyl-D-erythritol kinase
- ▣ 1 1 1 1 2-C-methyl-D-erythritol 2,4-cyclodiphosphate synthase



CLASSES FOR ALL LAND PLANTS

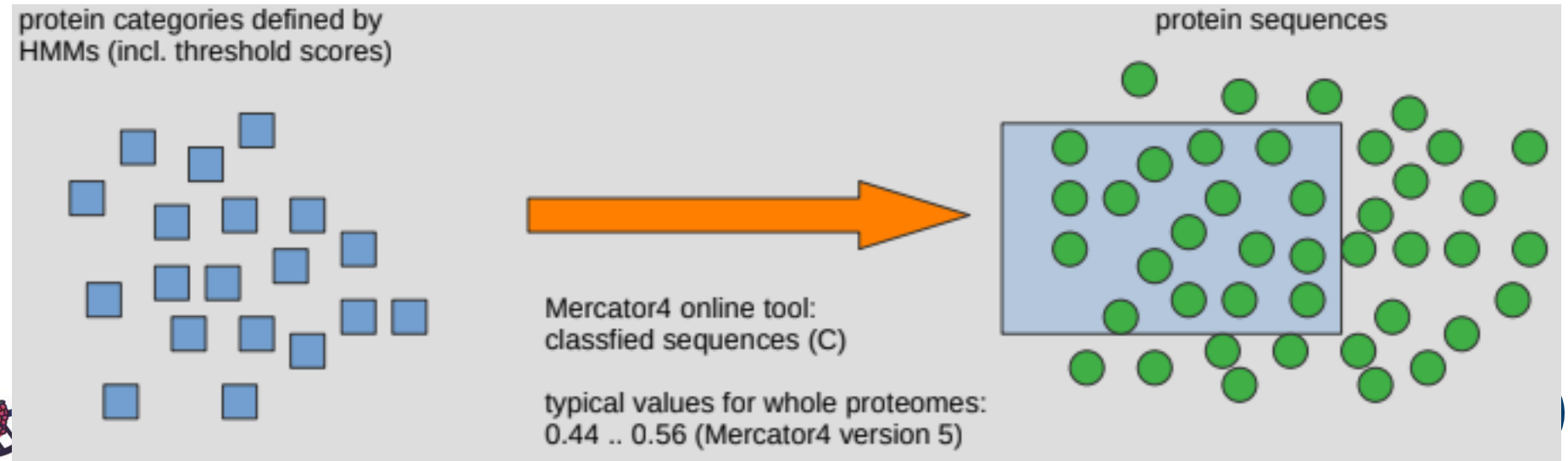
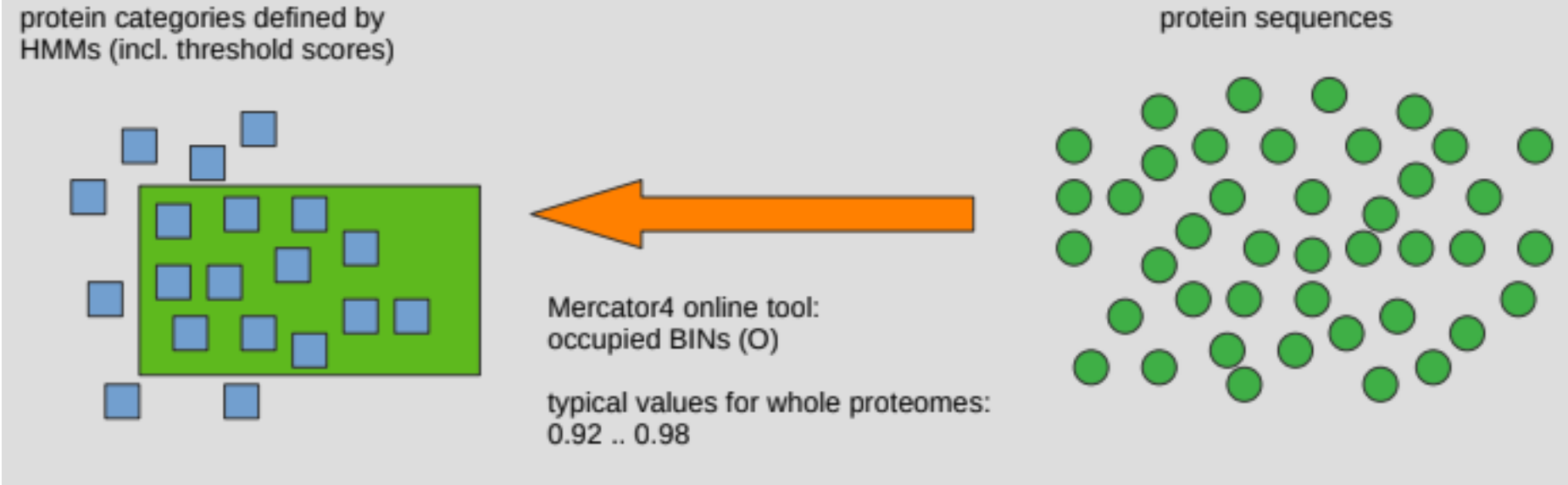
Hierarchical framework consisting of 31-top-level categories (also called *BINs*).

- 1 Photosynthesis
- 2 Cellular respiration
- 3 Carbohydrate metabolism
- 4 Amino acid metabolism
- 5 Lipid metabolism
- 6 Nucleotide metabolism
- 7 Coenzyme metabolism
- 8 Polyamine metabolism
- 9 Secondary metabolism
- 10 Redox homeostasis
- 11 Phytohormone action
- 12 Chromatin organisation
- 13 Cell division
- 14 DNA damage response
- 15 RNA biosynthesis
- 16 RNA processing
- 17 Protein biosynthesis
- 18 Protein modification
- 19 Protein homeostasis
- 20 Cytoskeleton organisation
- 21 Cell wall organisation
- 22 Vesicle trafficking
- 23 Protein translocation
- 24 Solute transport
- 25 Nutrient uptake
- 26 External stimuli response
- 27 Multi-process regulation
- 28 Plant reproduction
- 30 Clade-specific metabolism
- 50 Enzyme classification

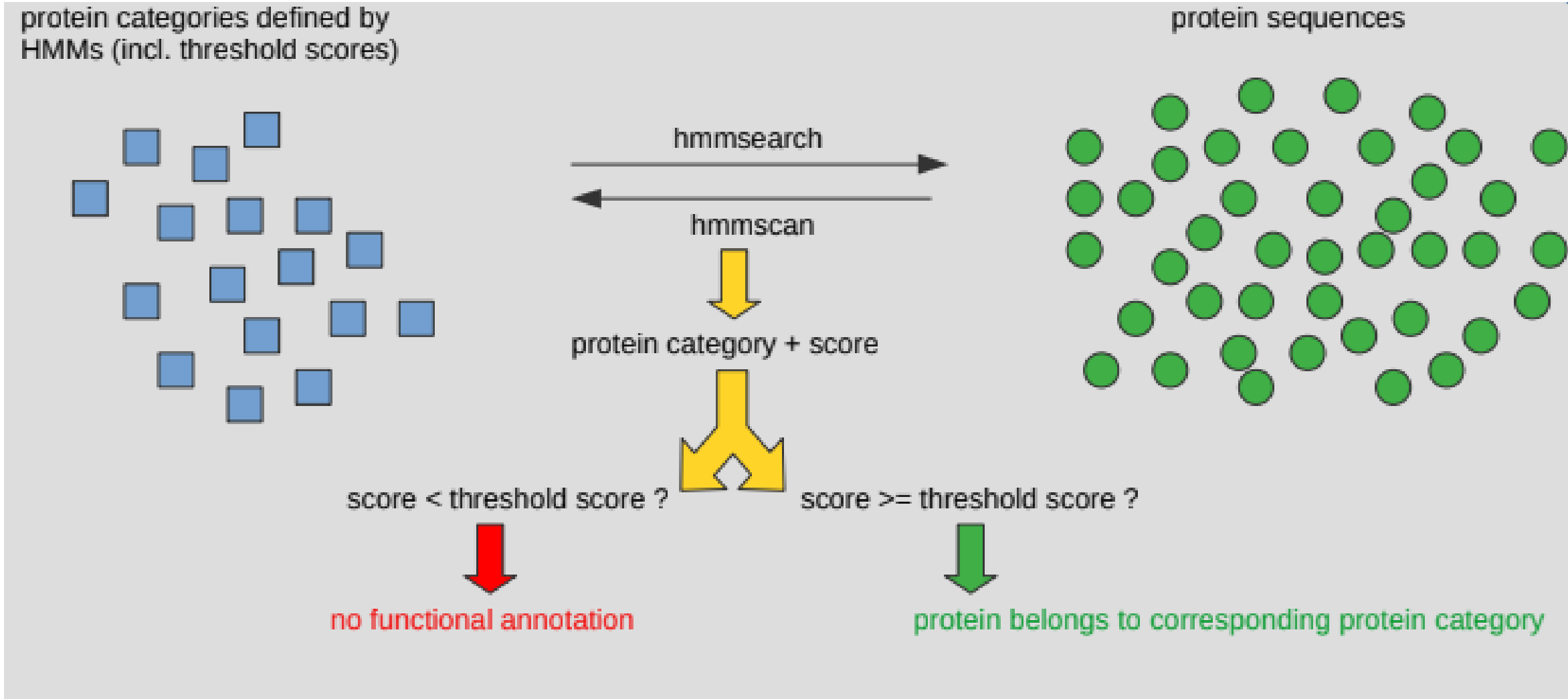
- ▶ 11 Phytohormone action
 - ▶ 11.1 abscisic acid
 - ▶ 11.1.1 biosynthesis
 - ▶ 11.1.1.1 zeaxanthin epoxidase *(ABA1)
 - ▶ 11.1.1.2 neoxanthin synthase *(ABA4)
 - ▶ 11.1.1.3 neoxanthin biosynthesis cofactor *(NXD1)
 - ▶ 11.1.1.4 9-cis-epoxycarotenoid dioxygenase *(NCED)
 - ▶ 11.1.1.5 xanthoxin oxidase *(ABA2)
 - ▶ 11.1.1.6 xanthoxin oxidase molybdopterin sulfurase *(ABA3)
 - ▶ 11.1.1.7 abscisic aldehyde oxidase *(AAO)
 - ▶ 11.1.2 perception and signalling
 - ▶ 11.1.3 conjugation and degradation
 - ▶ 11.1.4 transport
 - ▶ 11.2 auxin
 - ▶ 11.3 brassinosteroid
 - ▶ 11.4 cytokinin
 - ▶ 11.5 ethylene
 - ▶ 11.6 gibberellin
 - ▶ 11.7 jasmonic acid
 - ▶ 11.8 salicylic acid
 - ▶ 11.9 strigolactone
 - ▶ 11.10 karrikins
 - ▶ 11.11 signalling peptides



QUECK QUALITY OF PROTEOME



FUNCTIONAL ANNOTATION OF PLANT PROTEOMES



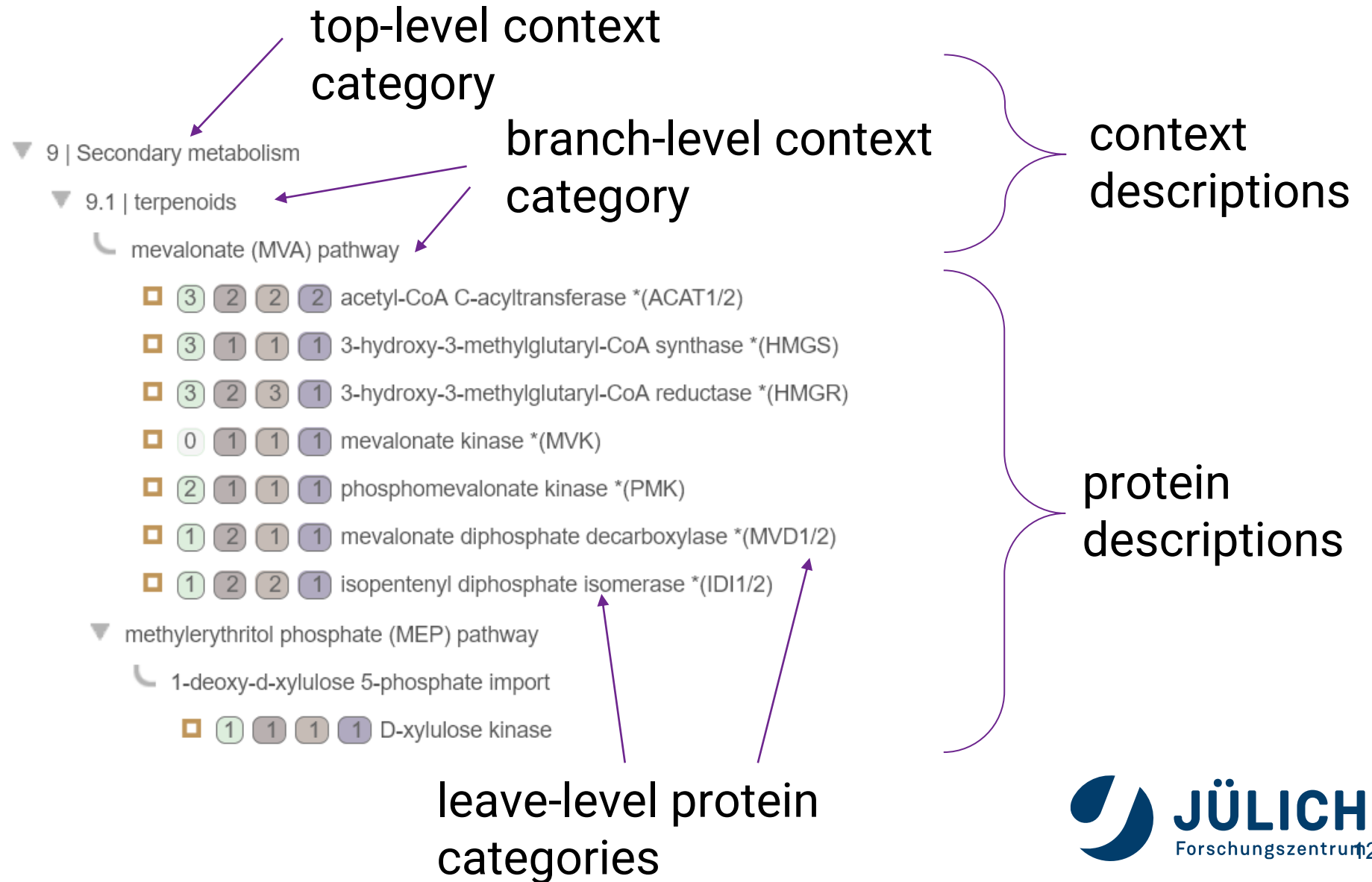
hmmsearch (protein HMM query vs protein database) & *hmmscan* (protein query vs protein HMM database)



BIN STRUCTURE

Each child node is more specialised than its parent *BIN*. Protein sequences are assigned to leaf level.

Soon... Enzyme classes linked to metabolic reactions



VERSION HISTORY

Annual release cycle (late summer -Autumn)

version	0.6 (2017)	1.0 (2018)	2.0 (2019)	3.0 (2020)	4.0 (2021)	5.0 (2022)	6.0 (2023)
protein categories + context nodes	3395 + 1068	4145 + 1339	4500 + 1491	4869 + 1608	5251 + 1702	5783 + 1817	~ 6180 + 1940



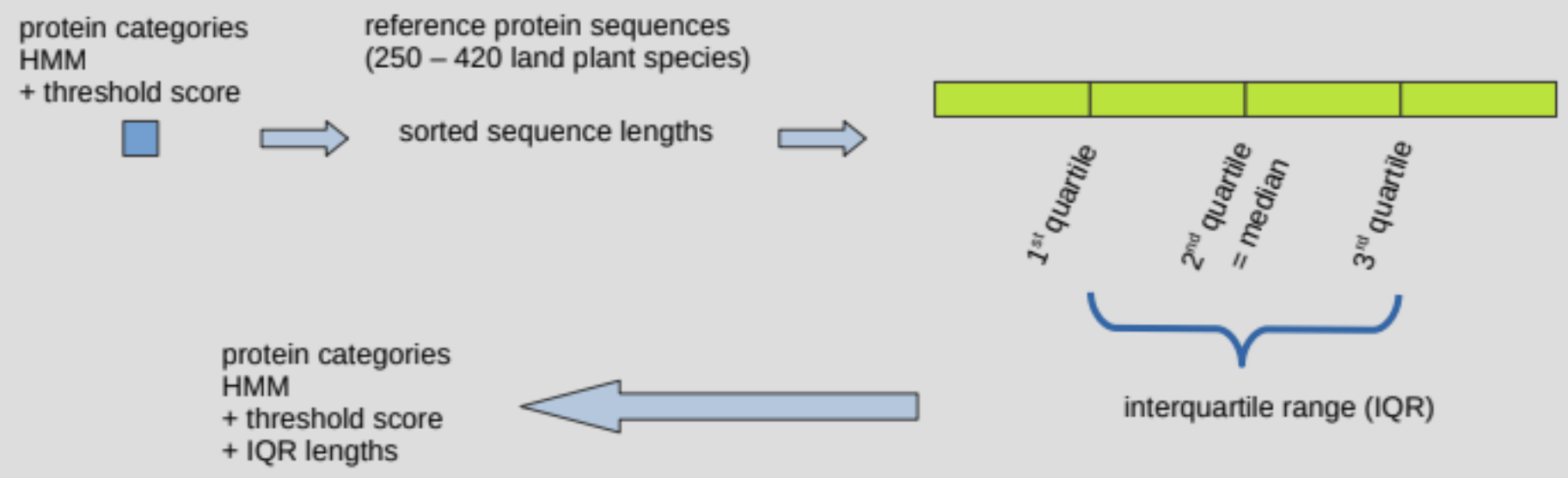
VERSION HISTORY

Annual release cycle (late summer -Autumn)

version	0.6 (2017)	1.0 (2018)	2.0 (2019)	3.0 (2020)	4.0 (2021)	5.0 (2022)	6.0 (2023)
protein categories + context nodes	3395 + 1068	4145 + 1339	4500 + 1491	4869 + 1608	5251 + 1702	5783 + 1817	~ 6180 + 1940
Classified		47.88	49.73	52.49	54.42	56.08	57.83
Annotated		72.0	72.76	73.53	74.60	74.93	94.54



FRAGMENTED SEQUENCES SINCE VERSION 6



example values:

species	below IQR	within IQR	above IQR
Arabidopsis thaliana	0.150	0.728	0.122
Glycine max	0.138	0.707	0.155
Trifolium pratense	0.126	0.721	0.152
Solanum pennellii	0.146	0.719	0.135
Cuscuta campestris	0.202	0.594	0.204
Cuscuta pentagona (from transcriptome)	0.447	0.436	0.117
Peperomia fraseri (from transcriptome)	0.656	0.327	0.017

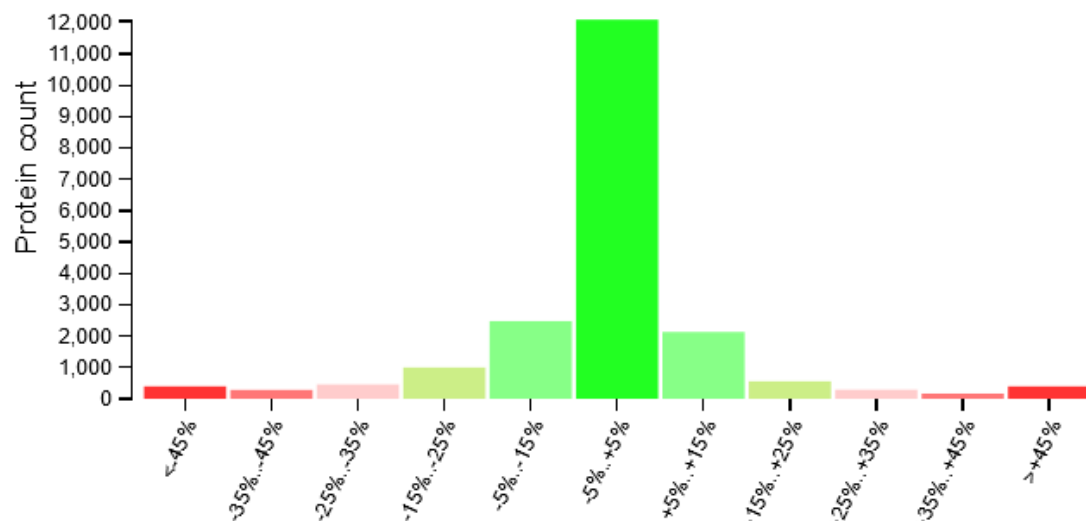
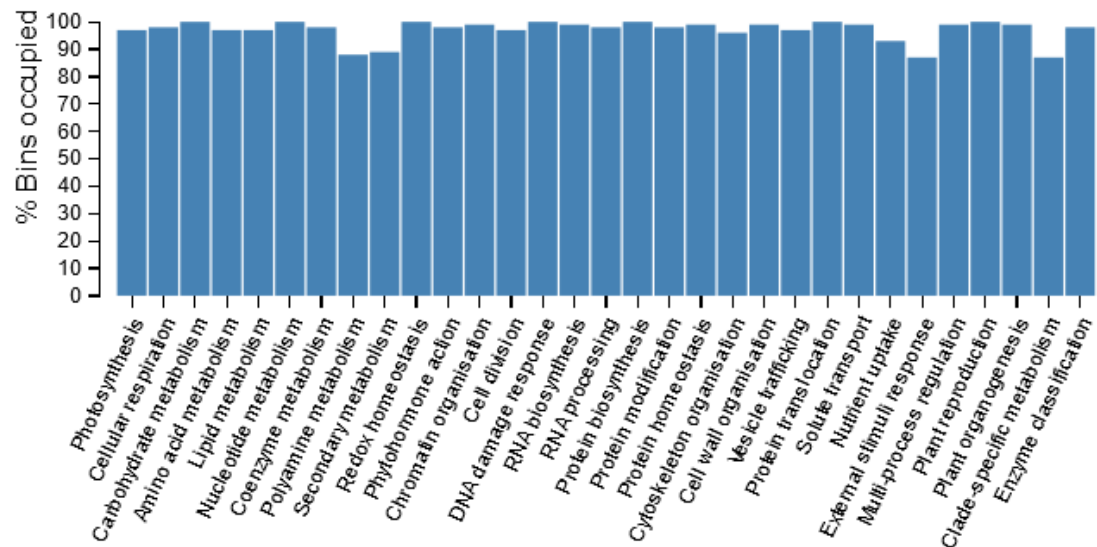
fragmented proteome



Job name: pep6
Job ID: GFA-76bd35bbcdb46999f2fc89f9b6b4b212
Number of sequences: 27416
Sequence type: Protein
Prot-scriber: False
Swissprot: False
Submitted to cluster: 13.2.2024, 21:37:07
Job status: FINISHED
Last update: 13.2.2024, 21:47:08

Submitted sequences (S): 27416
Annotated sequences (A): 15854
Classified sequences (C): 15854
Occupied bins (O): 6115
Bins available (B): 6223
Summary ⓘ: S:27416,A:57.83%,C:57.83%,O:98.26%,B:6223

Mapping file for MapMan: [MapMan mapping file](#)
Annotated FASTA file: [Mercator4 annotated FASTA file](#)



PRESENCE ABSENCE AND COPY NUMBERS

▼ 9 | Secondary metabolism

▼ 9.1 | terpenoids

↳ mevalonate (MVA) pathway

- ▣ 3 2 2 2 acetyl-CoA C-acyltransferase *(ACAT1/2)
- ▣ 3 1 1 1 3-hydroxy-3-methylglutaryl-CoA synthase *(HMGS)
- ▣ 3 2 3 1 3-hydroxy-3-methylglutaryl-CoA reductase *(HMGR)
- ▣ 0 1 1 1 mevalonate kinase *(MVK)
- ▣ 2 1 1 1 phosphomevalonate kinase *(PMK)
- ▣ 1 2 1 1 mevalonate diphosphate decarboxylase *(MVD1/2)
- ▣ 1 2 2 1 isopentenyl diphosphate isomerase *(IDI1/2)



▼ methylerythritol phosphate (MEP) pathway

↳ 1-deoxy-d-xylulose 5-phosphate import

- ▣ 1 1 1 1 D-xylulose kinase
- ▣ 1 1 1 0 D-xylulose 5-phosphate transporter
- ▣ 7 3 4 2 1-deoxy-D-xylulose 5-phosphate synthase *(DXS)
- ▣ 2 1 1 2 1-deoxy-D-xylulose 5-phosphate reductase *(DXR)
- ▣ 0 1 1 1 4-diphosphocytidyl-2-C-methyl-D-erythritol synthase
- ▣ 1 1 1 1 4-diphosphocytidyl-2-C-methyl-D-erythritol kinase
- ▣ 1 1 1 1 2-C-methyl-D-erythritol 2,4-cyclodiphosphate synthase



Search by position on chromosome

chromosome (available position range)

Chr1 (1 .. 23207424)

1

100000

search reset

19 protein result(s)

	Plant species Sequence details	Description Functional context (based on MAPMAN4)	Orthologous proteins
1	Camelina sativa Csa01g001010	putative component of lumen subcomplex of chloroplast NDH <i>Photosynthesis: photophosphorylation</i>	protein family
2	Camelina sativa Csa01g001020	(yeast MED11)-like component of Mediator transcriptional regulatory complex <i>RNA biosynthesis: RNA polymerase II-dependent transcription</i>	protein family
3	Camelina sativa Csa01g001030	alpha-dioxygenase <i>Lipid metabolism: lipid degradation</i>	protein family
4	Camelina sativa Csa01g001040	putative ribonuclease H1 involved in R-loop homeostasis <i>RNA biosynthesis: RNA polymerase II-dependent transcription</i>	protein family
5	Camelina sativa Csa01g001050	protein of unknown function	protein family
6	Camelina sativa Csa01g001060	G-type subunit of vacuolar H(+)-ATPase peripheral V1 subcomplex <i>Solute transport: primary active transport</i>	protein family
7	Camelina sativa Csa01g001070	putative PIG-N-type protein of GPI biosynthetic pathway <i>Protein modification: lipidation</i>	protein family
8	Camelina sativa Csa01g001080	putative plastidial CRS/CFM-type RNA intron splicing factor <i>RNA processing: organelle machinery</i>	protein family
9	Camelina sativa Csa01g001090	putative plastidial CRS/CFM-type RNA intron splicing factor	protein family



Mercator4 BIN enrichment analysis

Upload Mercator4 mapping file [Use Example Dataset](#)

Perform One-sided Fisher's exact test
 Over-representation analysis
 Under-representation analysis
 Two-sided Fisher's exact test

FDR-adjusted p-value cutoff

Genes of Interest

Background Genes

[Submit](#) [Reset](#) [Download Result Table as TSV](#)

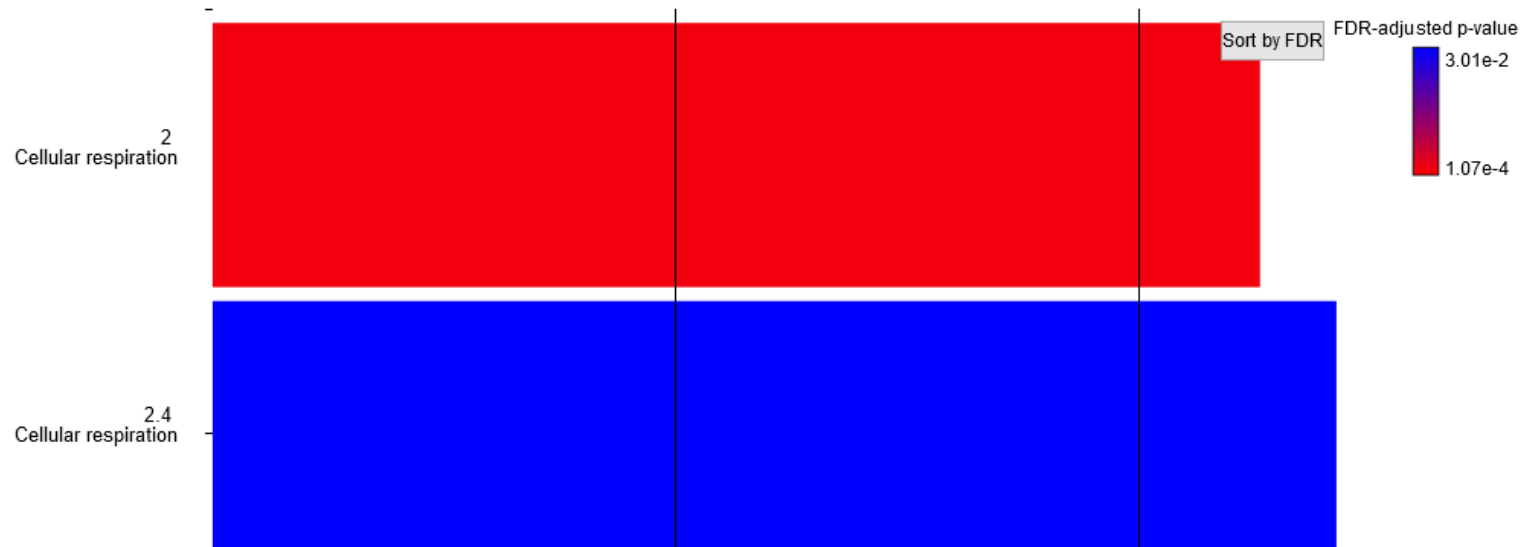


Toggle column: [Genes of Interest \(List\)](#)

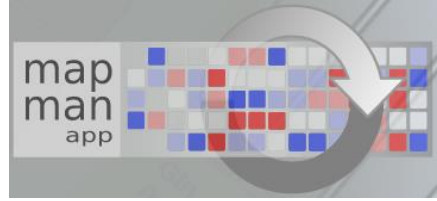
Search:

MapMan4 category number	Context of Protein Function	#Genes of Interest IN MapMan4 category	#Genes of Interest NOT IN MapMan4 category	#Background Genes IN MapMan4 category	#Background Genes NOT IN MapMan4 category	Enrichment Factor	p-value	FDR-adjusted p-value
2	Cellular respiration	18	184	29	1441	4.52	0.0000027040082208663597	0.001616996916078083
2.4	Cellular respiration.oxidative phosphorylation	12	190	18	1452	4.85	0.00008388516891239334	0.03009799860576673
35	No Mercator4 annotation	40	162	569	901	0.51	5.988643766125384e-8	0.00010743626916428939
35.1	No Mercator4 annotation.other annotation available	40	162	505	965	0.58	0.000020617733553228634	0.009247053498623042
50	Enzyme classification	62	140	219	1251	2.06	1.8643729410822603e-7	0.00016723425281507876

Showing 1 to 5 of 5 entries



MAPMAN VISUALISATIONS



Protein function annotation



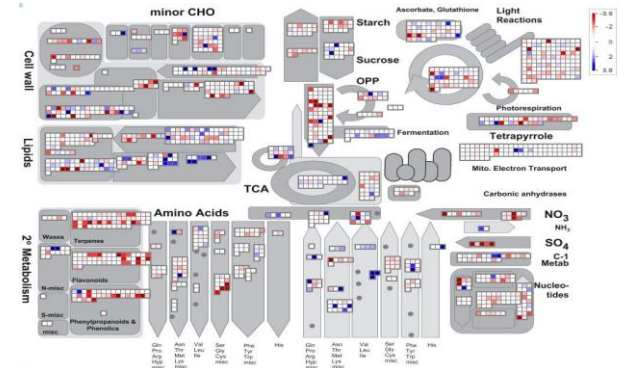
Mapman4 category (BIN)		protein name	protein description
BINCODE	NAME	IDENTIFIER	DESCRIPTION
"1.3.1"	"Photosynthesis, photorespiration, phosphoglycolate phosphatase"	"ab0g12345"	"phosphoglycolate phosphatase"
"2.3.1"	"Cellular respiration, tricarboxylic acid cycle, citrate synthase"	"ab0g23456"	"citrate synthase"
"18.9.1"	"Protein modification, hydroxylation, prolyl hydroxylase"	"ab0g34567"	"prolyl hydroxylase"

Gene expression data



	condition 1		condition 2	
	expression change		expression change	
	log2 fold	significance	log2 fold	significance
<i>example</i>	"cond1 vs. normal"	"cond1 vs. normal"	"cond2 vs. normal"	"cond2 vs. normal"
"ab0g12345"	-0.25	0	0.85	1
"ab0g23456"	0.63	1	1.25	0
"ab0g34567"	0.46	0	1.62	1

Pathway diagram



Run Mercator4 on genome assembly



Differential gene expression (snRNA, RNA-Seq, Microarray data), Metabolites, proteins

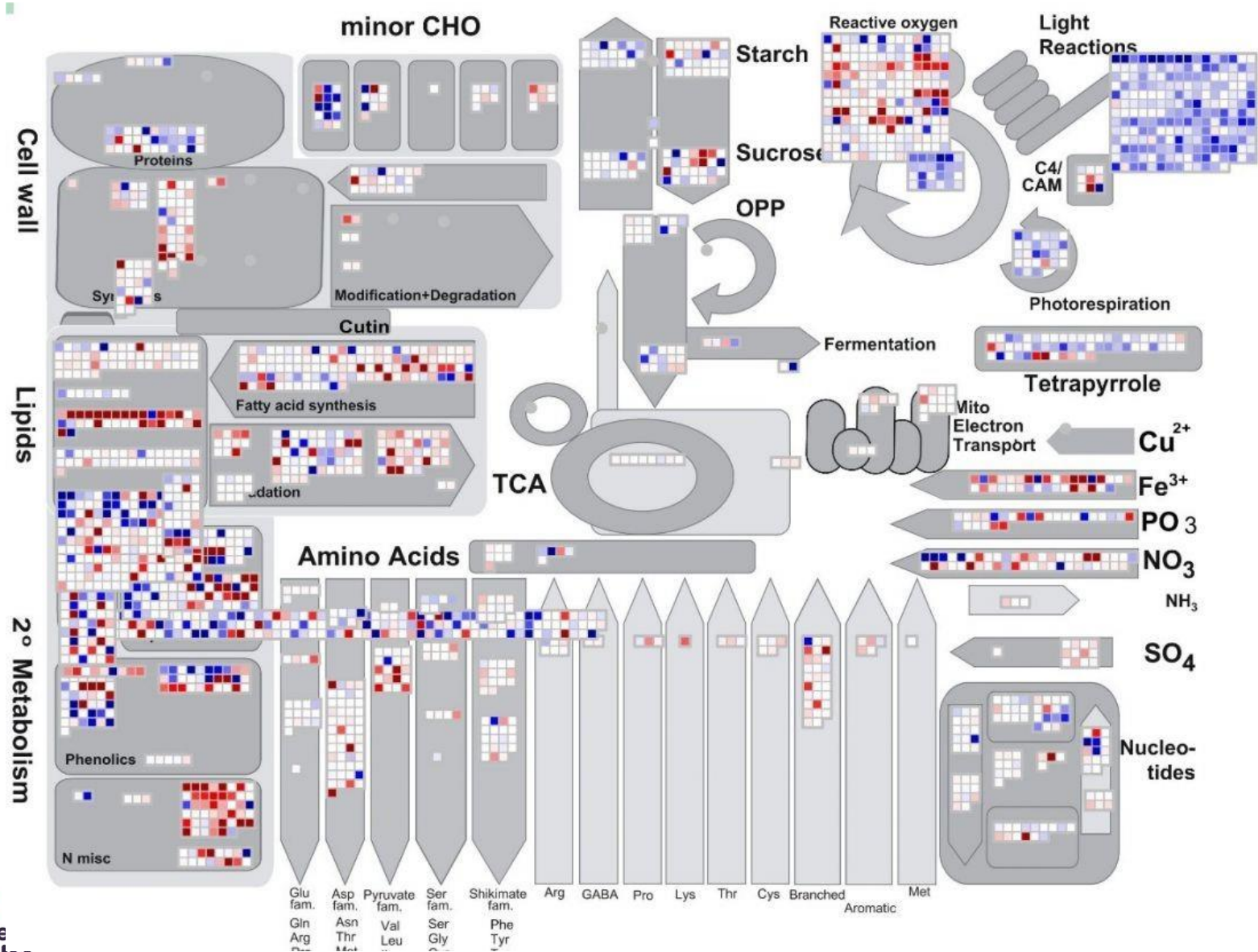
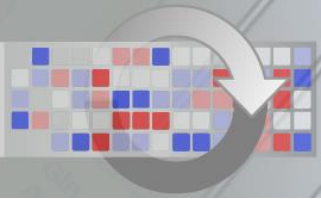


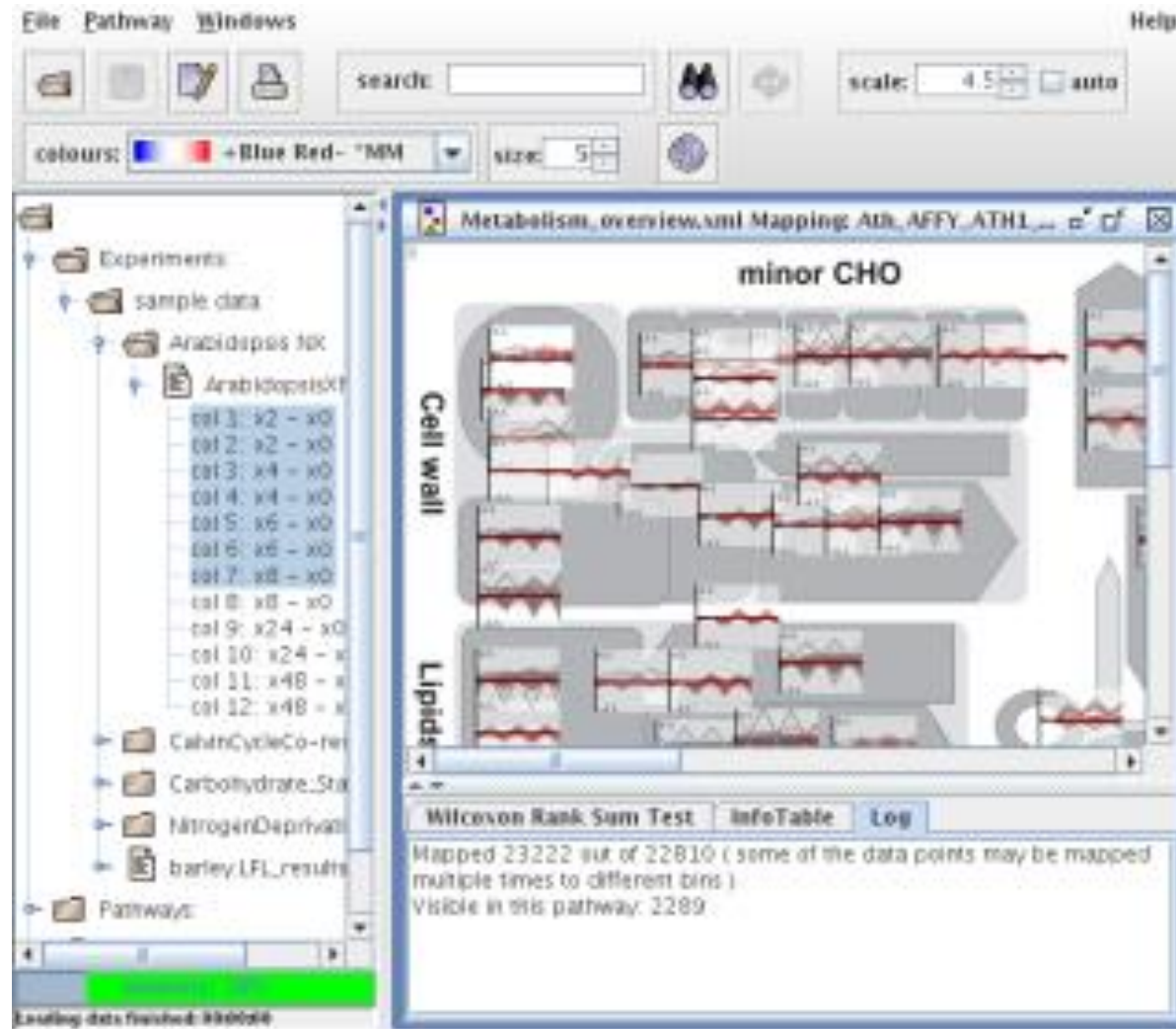
Pathway you are interested in, available on the plabipd website



Important Note: Identifiers must be identical







- Dr. Alisandra Denton CEPLAS Düsseldorf (HELIXER)
- Dr. Rainer Swacke, Sebastian Beier, Marie Bolger

DFG

