# Analysis and visualization of biological experiment data in the context of relevant networks

## Christian Klukas

Plant Bioinformatics Group
Leibniz Institute of Plant Genetics and Crop Plant Research

Progress Seminar Department Molecular Genetics

**IPK-Gatersleben, January 23rd 2009**

**IPK**
**GATERSLEBEN**

# Outline

1. Motivation

2. Methods

   ◦ Definition and combination of data models for experiment data and networks

   ◦ Network-integrated visualization of experiment data

   ◦ Interaction-, layout- and navigation-techniques

   ◦ Statistical analysis

3. Implementation

4. Example use-cases (online-demo)

# Motivation

# Methods' Requirements



1. Experiment Data

Methods
· Data mapping
· Visualization
· Interaction
· Layout, Navigation
· Statistical Analysis

3. Identifier Knowledge

2. Networks or classification hierarchies

# Methods: Data model for experiment data

- Different standards for biological domains:
  - Genomics:  MIAME / MAGE
  - Proteomics:  PEDRo
  - Metabolomics: ArMet
- Existing models differ greatly, contain all kinds of annotations, new model:
  - Flexible model for all of the domains
  - Only data for visualization, analysis and identification

| experiment | | condition | | sample | | measurement | | substance |
|---|---|---|---|---|---|---|---|---|
| name : String<br>coordinator : String<br>importDate : Date<br>importBy : String<br>comment : String | 1    1..* | species : String<br>genotype : String<br>treatment : String | 1    0..* | time : Float<br>timeUnit : String | 1    0..* | value : Float<br>unit : String<br>replicateID : int | 1..*    1 | name : String<br>synonyms : String |

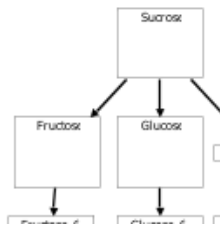# Methods: Data model for biological networks and classification hierarchies

- (Generic) Mapping-Graph MG=(V, E, l, $t_V$, $t_E$, z)
    - V – set of nodes
    - E – set of edges (directed or undirected)
    - l – function node/edge → label (L)
    - $t_V$, $t_E$ – function node/edge → node/edge type ($T_V$, $T_E$)
    - z – function node/edge → mapped experiment data

# Methods: Data model for biological networks and classification hierarchies

- ## Protein-Protein Network - $MG_{PPI}$
  - Nodes represent proteins, $T_V$={protein}
  - Undirected edges represent interaction between two proteins, $T_E$={interaction}
- ## KEGG Pathway – $MG_{KEGG}$
  - $T_V$={Orthologe, Enzyme, Gene, Gene-Group, Metabolite, Map-Link}
  - $T_E$={ECrel, PPrel, GErel, PCrel, rProd, rSub, link} (enzyme-enzyme relation, protein-protein relation, gene-expression, protein-metabolite relation, reaction product, reaction substrate, link to pathway)
- ## (extended) Pathway-Overview Graph – $MG_O$, $MG_{OE}$
- ## Gene Ontology Hierarchy – $MG_{GO}$
- ## KEGG BRITE Hierarchy - $MG_{BRITE}$

# Methods: Data-mapping – combination of network and experiment data

**Data**

| Graph (Nodes & Edges) | Dataset | Synonyms | Alternative Identifiers |
|---|---|---|---|



Graph (Nodes & Edges)

Dataset

Synonyms
- KEGG Compounds
- Expasy Enzymes
- KEGG KO

Alternative Identifiers

**Pseudo code data m.**

```
for each substance in dataset do
    Set of Substance IDs = expandWithDBids(substance.name,
                              substance.synonyms)
    for each graphelment in targetGraph do
        Set of Graphelement IDs = expandWithDBids(label)
        if intersection(substance IDs, Graphelement IDs).length > 0 then
            graphelment.assignData(dataset.getSubsetFor(substance))
    end for
end for
```
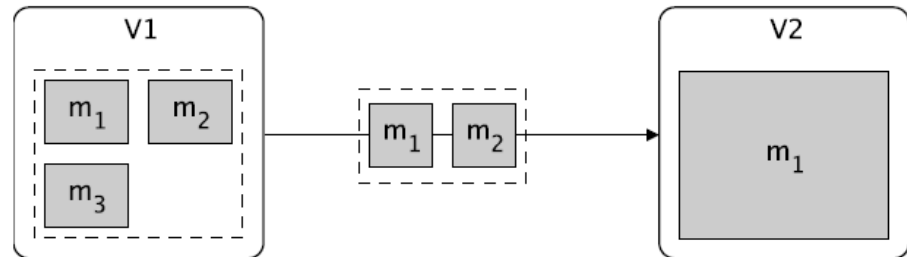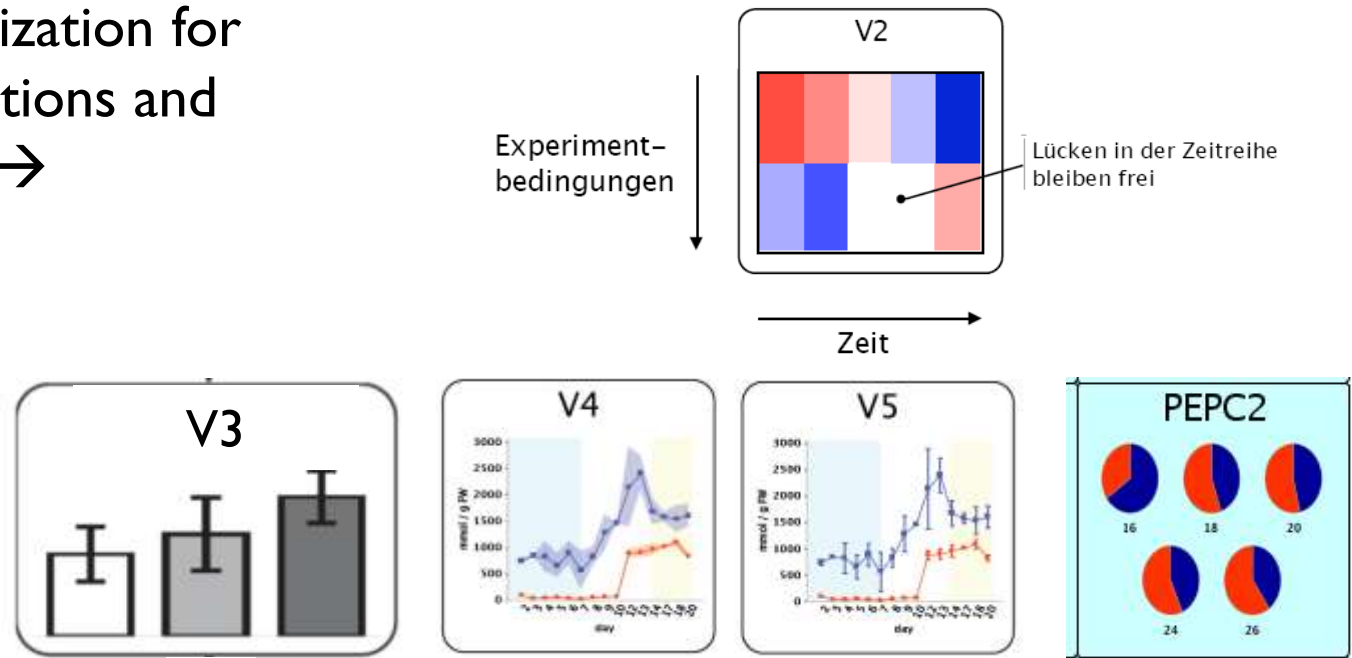
**Result**



8

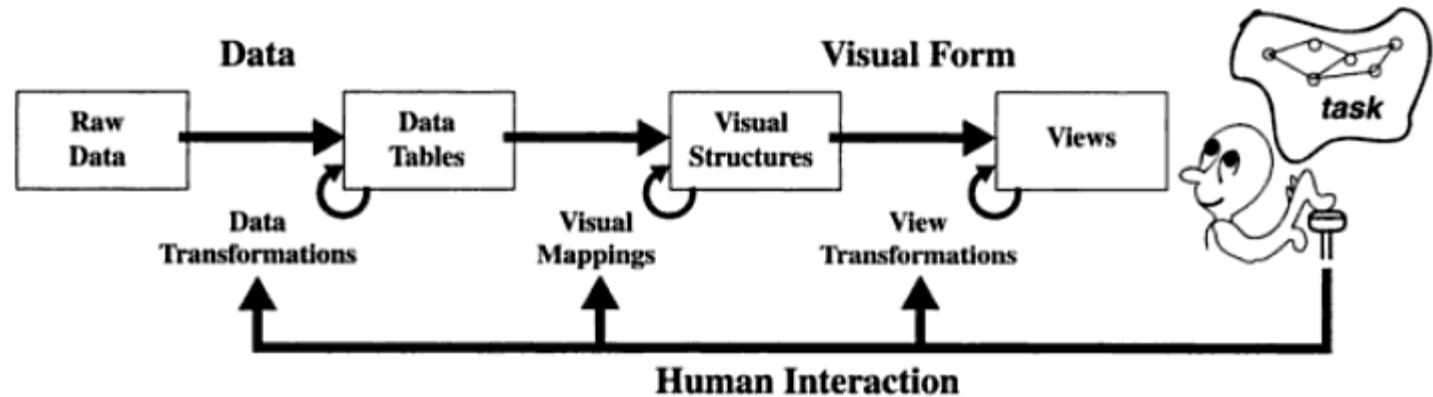# Methods: Network integrated data visualization

Processing of multiple mappings →

Data visualization for conditions and time →
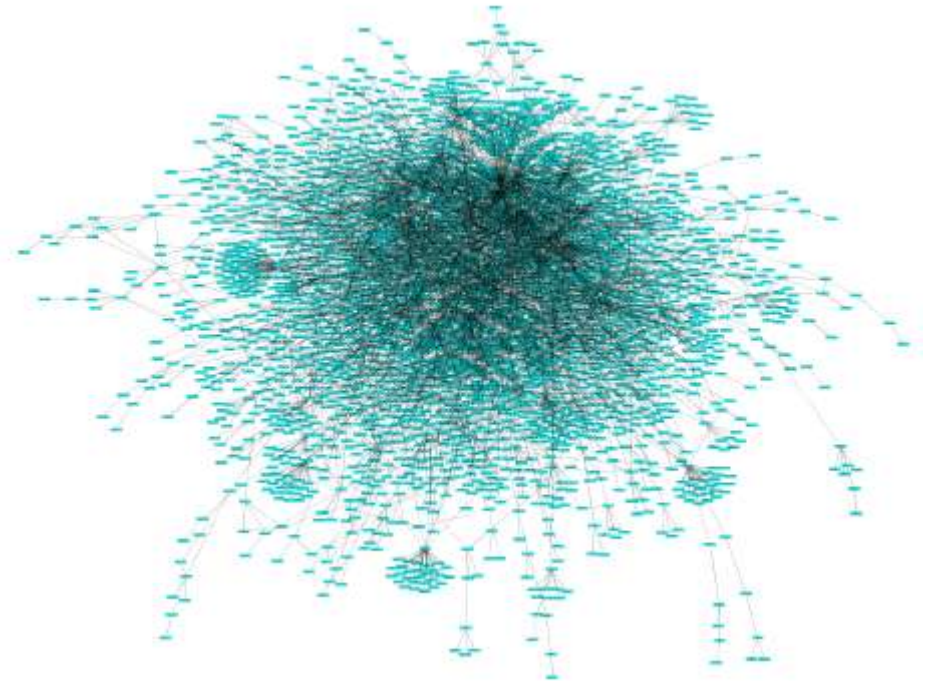
# Methods: Interaction techniques



Visualization pipeline  (Card, Mackinlay, Shneiderman, 1999)

- Data transformations
  - Dynamic queries (select conditions, search graph elements)
  - Direct walk (pathway navigation)
  - Attribute walk (select similar elements)
  - Details-on-demand (show/hide details)
  - Direct manipulation (node position, label)
- Visual mappings
  - Experiment data display: color-coded, size-/width-coded, diagrams
- View transformations
  - Direct selection (click & select)
  - Overview-and-Detail (multiple views with varying detail)

# Methods: Layout of specific mapping-graphs
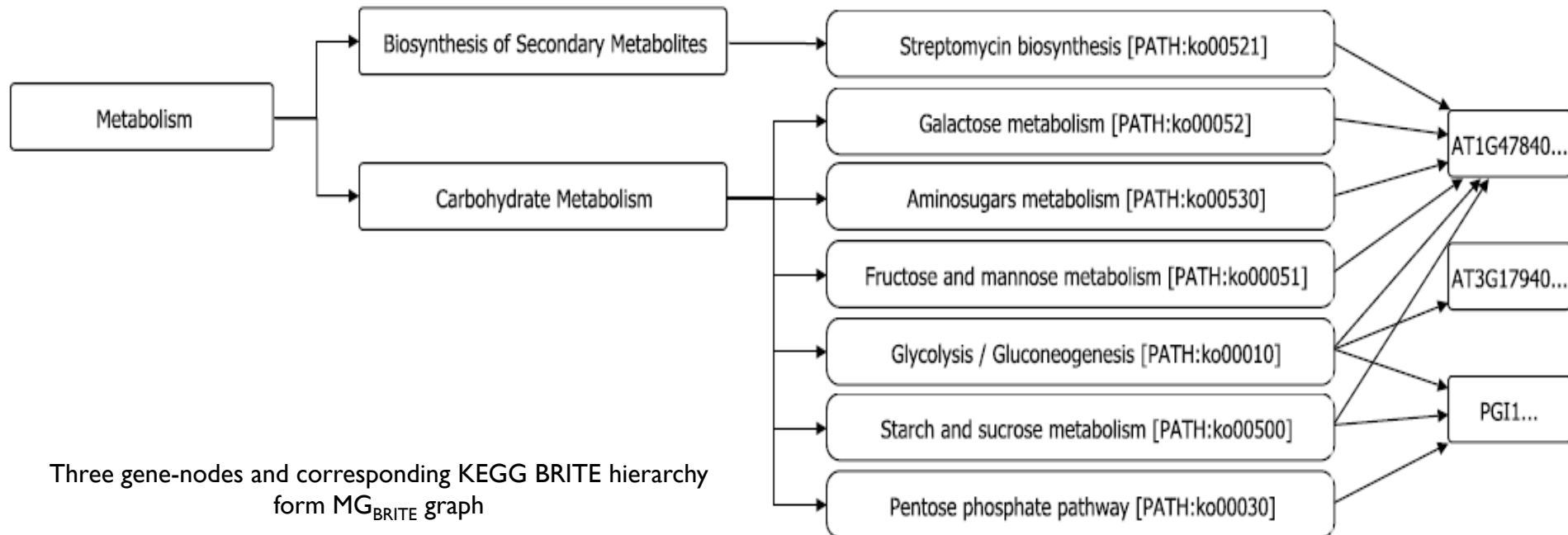
- $MG_{PPI}$: Force-Directed (Spring-Embedder)

- $MG_{GO}$, $MG_{BRITE}$: Hierarchical (Sugiyama)

- $MG_{KEGG}$: Manual layout is given



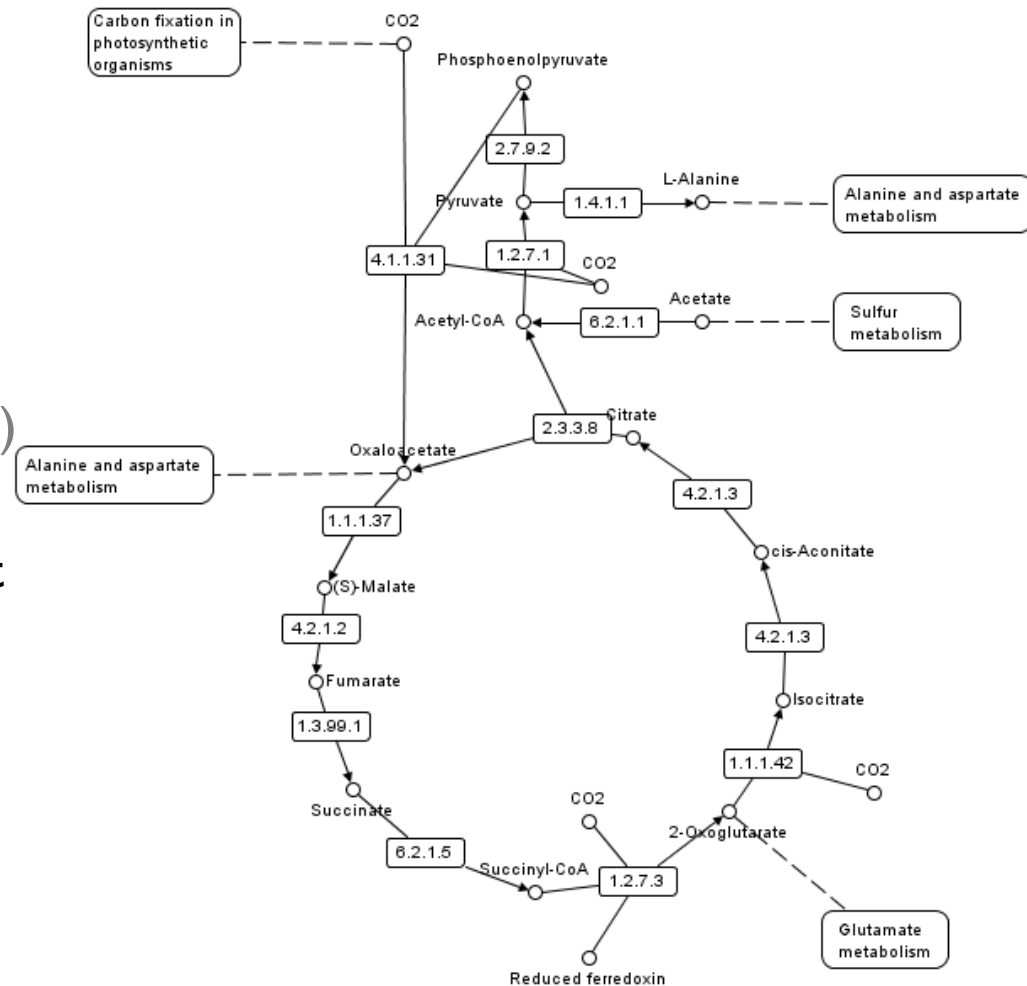PPI network of worm C.elegans (5418 interactions, 2992 proteins). Datasource: GraphWeb

# Methods: Layout of specific mapping-graphs

- $MG_{PPI}$: Force-Directed (Spring-Embedder)

- $MG_{GO}$, $MG_{BRITE}$: Hierarchical (Sugiyama)



Three gene-nodes and corresponding KEGG BRITE hierarchy form $MG_{BRITE}$ graph
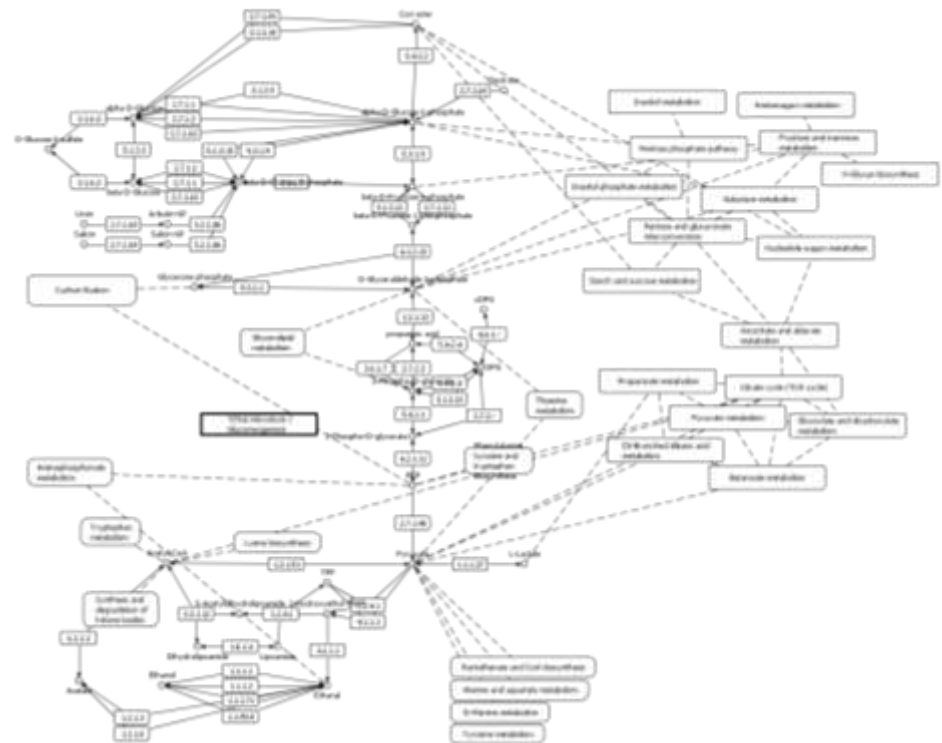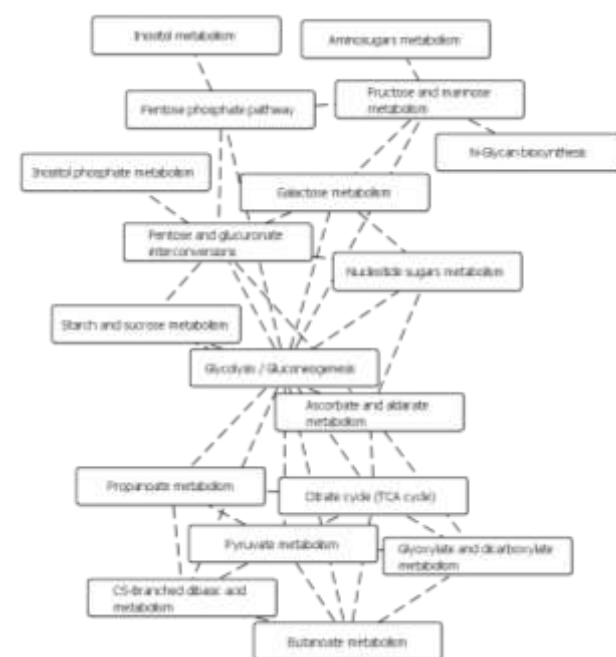
# Methods: Layout of specific mapping-graphs

- $MG_{PPI}$: Force-Directed (Spring-Embedder)

- $MG_{GO}$, $MG_{BRITE}$: Hierarchical (Sugiyama)

- $MG_{KEGG}$: Manual layout is given



KEGG Pathway „Reductive carboxylate cycle (CO2 fixation)" , visualization: VANTED
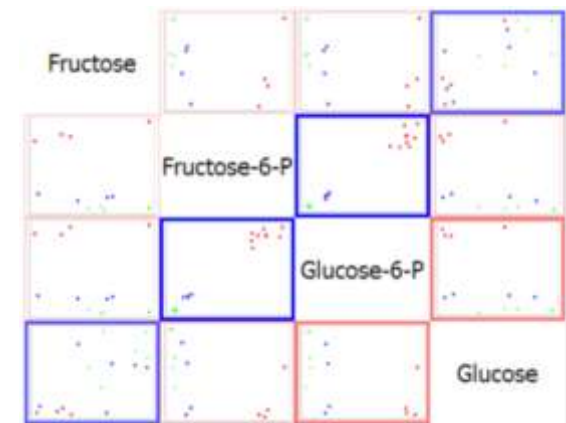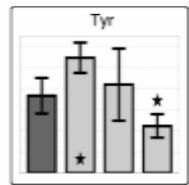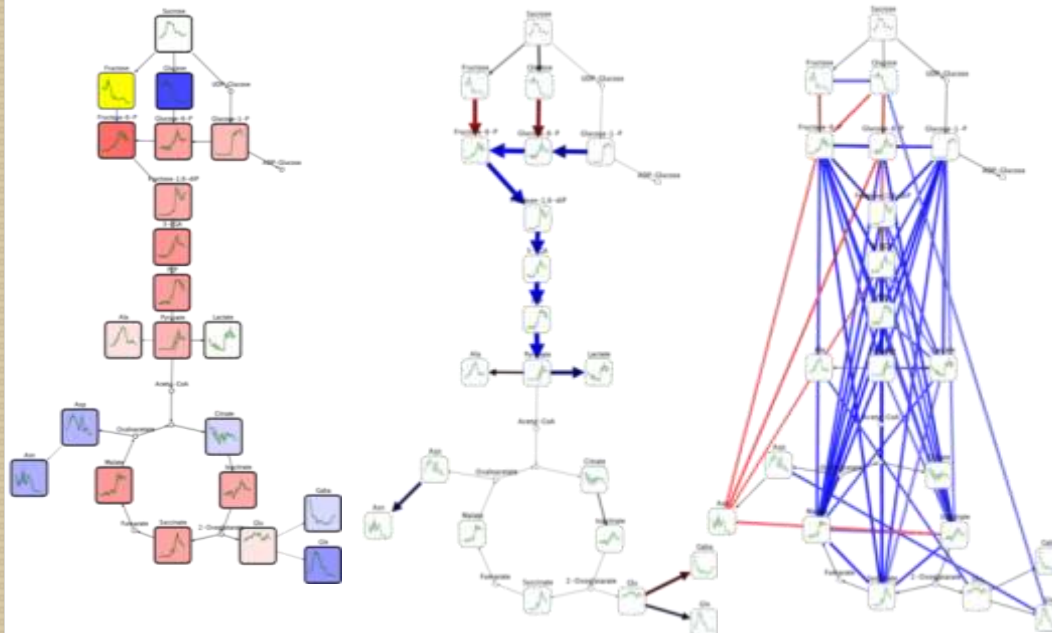
# Methods: Pathway-Navigation

- Extending the overview
- Collapsing pathways
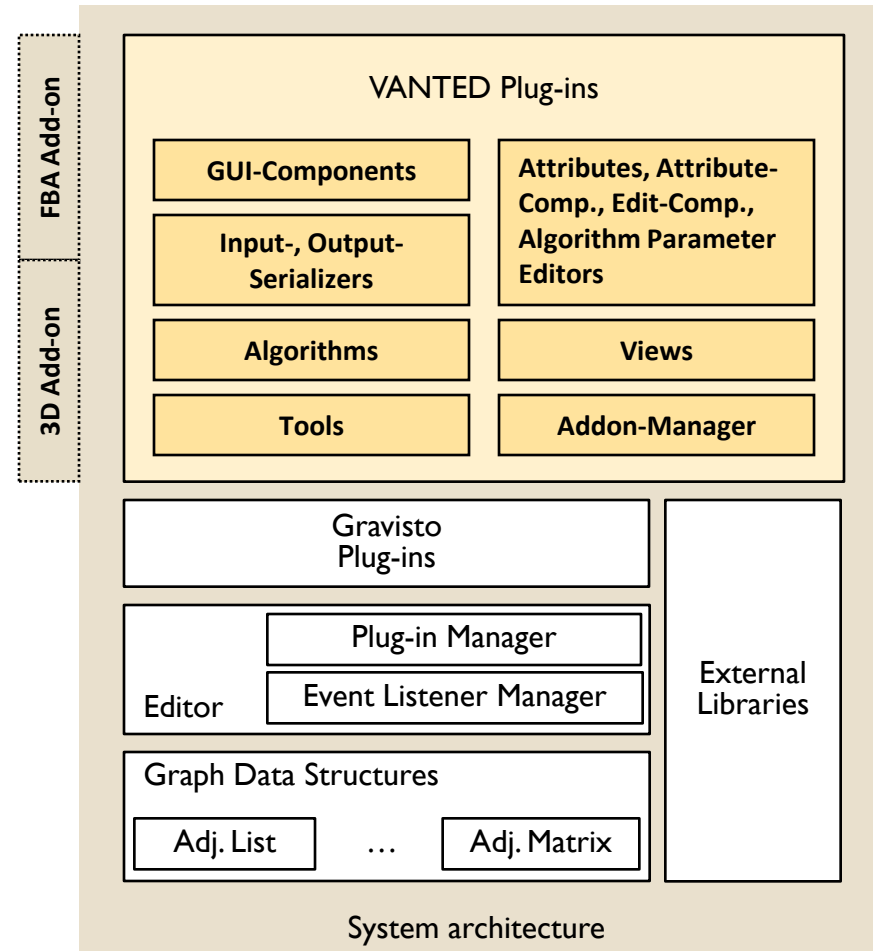- Arranging pathways
- Stepwise pathway-navigation

# Methods: Statistical analysis

- Statistic tests: t-Test (2 variants), Grubbs-Test, David-Quicktest
- Flexible calculation of correlations
- Significance analysis for $MG_{KEGG}$ and $MG_{BRITE}$ (Fisher's exact test)
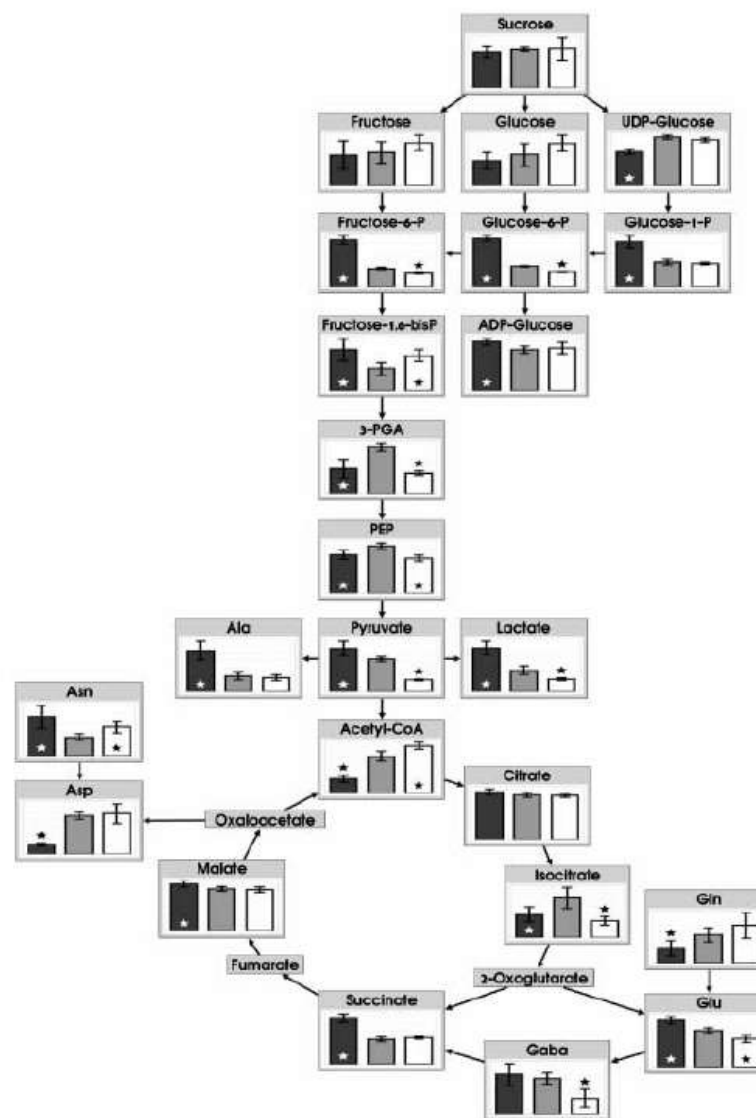- Scatter-Plots

# Implementation

- Based on the extensible, plugin-based graph visualization toolkit **Gravisto** (developed at the University of Passau and at the IPK)

- MVC pattern

- Event management (observer design pattern)

- Java application (Windows exe, Mac OS X image, platform-neutral ZIP, Java WebStart)

- External plug-ins ("Add-ons")
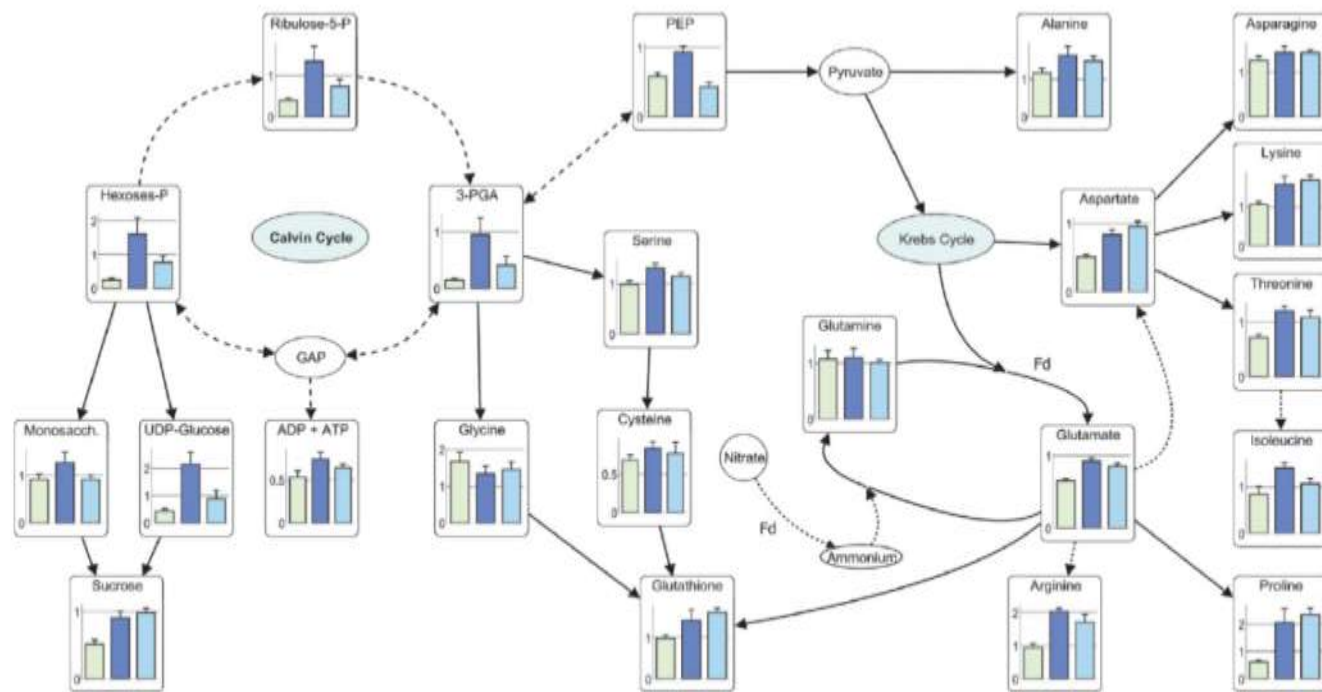


System architecture

# Visualization of relative metabolite changes under different stress situations

Rolletschek et al. (2005): Regulation of lipid biosynthesis in soybean seeds: evidence for a key role of photosynthetic oxygen release. **New Phytologist**

**Fig. 5** Simplified schematic representation of primary metabolites in soybean (*Glycine max*) seeds and their response to changing $O_2$ supply. Intermediates of glycolysis and the citrate cycle are shown, as well as branch paths to related sugars and free amino acids. Arrows indicate probable direction of predominant carbon flow. Vertical bars show the level of each metabolite in seeds after a 6 h treatment with 2.1, 21 and 42 kPa oxygen (black, grey and white bars, respectively). Data are given in relative units (mean ± SD). *, Significant differences vs control treatment (21 kPa), $t$-test $P < 0.05$.
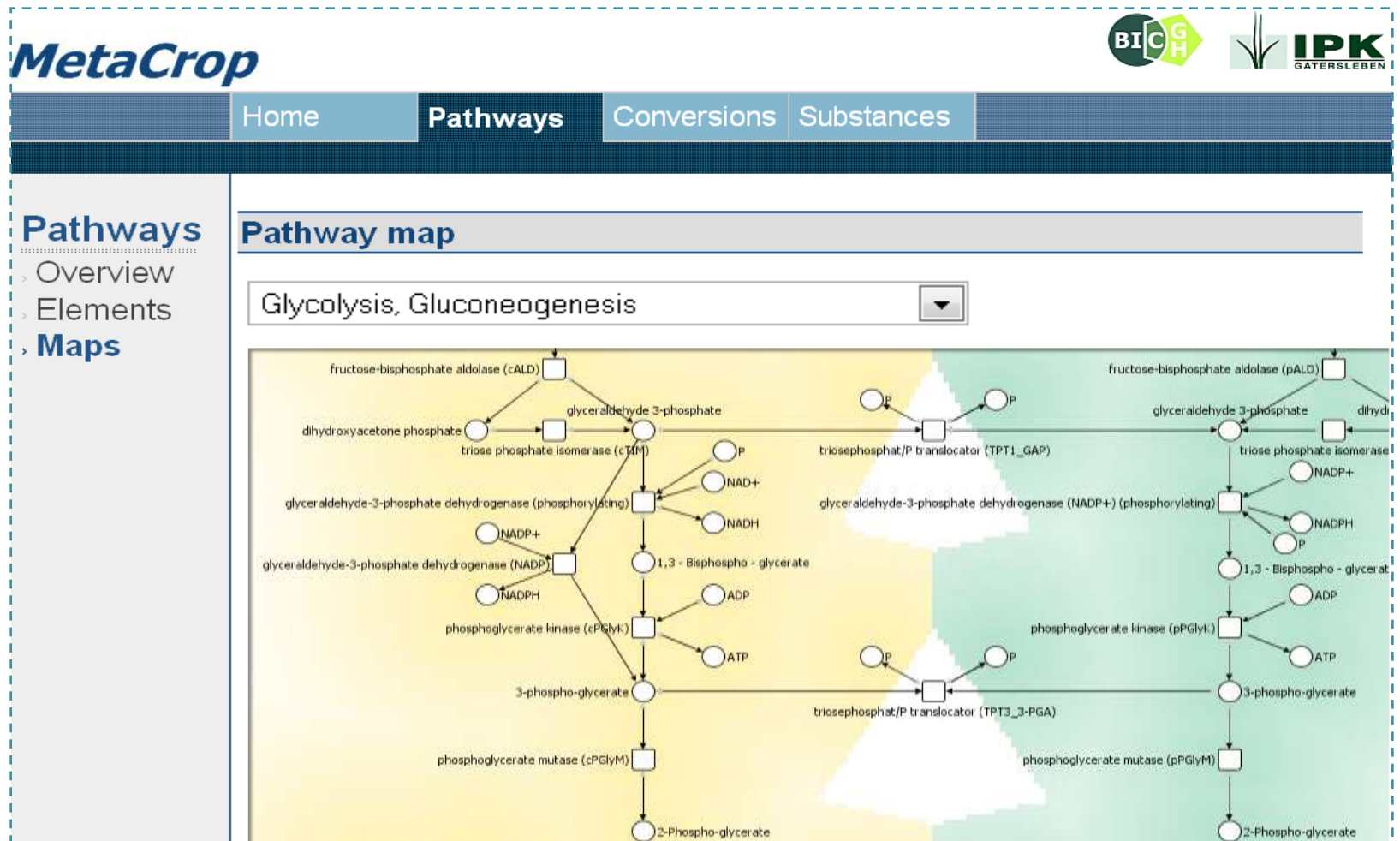
# Visualization of relative metabolite changes under different stress situations



Fig. 4. Relative metabolite changes in iron-starved and control plants. Four-week-old plants were transferred to hydroponic Hoagland solution supplemented with either FeSO4-EDTA or CaCO3 (pH 8.0). Leaf material was harvested after 29 days, and the corresponding metabolites were measured as described in *Materials and Methods*. Depicted are the ratios ± SE of metabolite contents between Fe-starved and -replete plants of WT (green bars), *pfld*5-8 (blue bars), and *pfld*4-2 (light blue bars) lines (*n* = 8–10 independent plants). The graph was created by using the visualization system Vanted (38).

Tognetti et al. (2007): *Enhanced plant tolerance to iron starvation by functional substitution of chloroplast ferredoxin with a bacterial flavodoxin.* **Proc. Natl. Acad. Sci. USA**
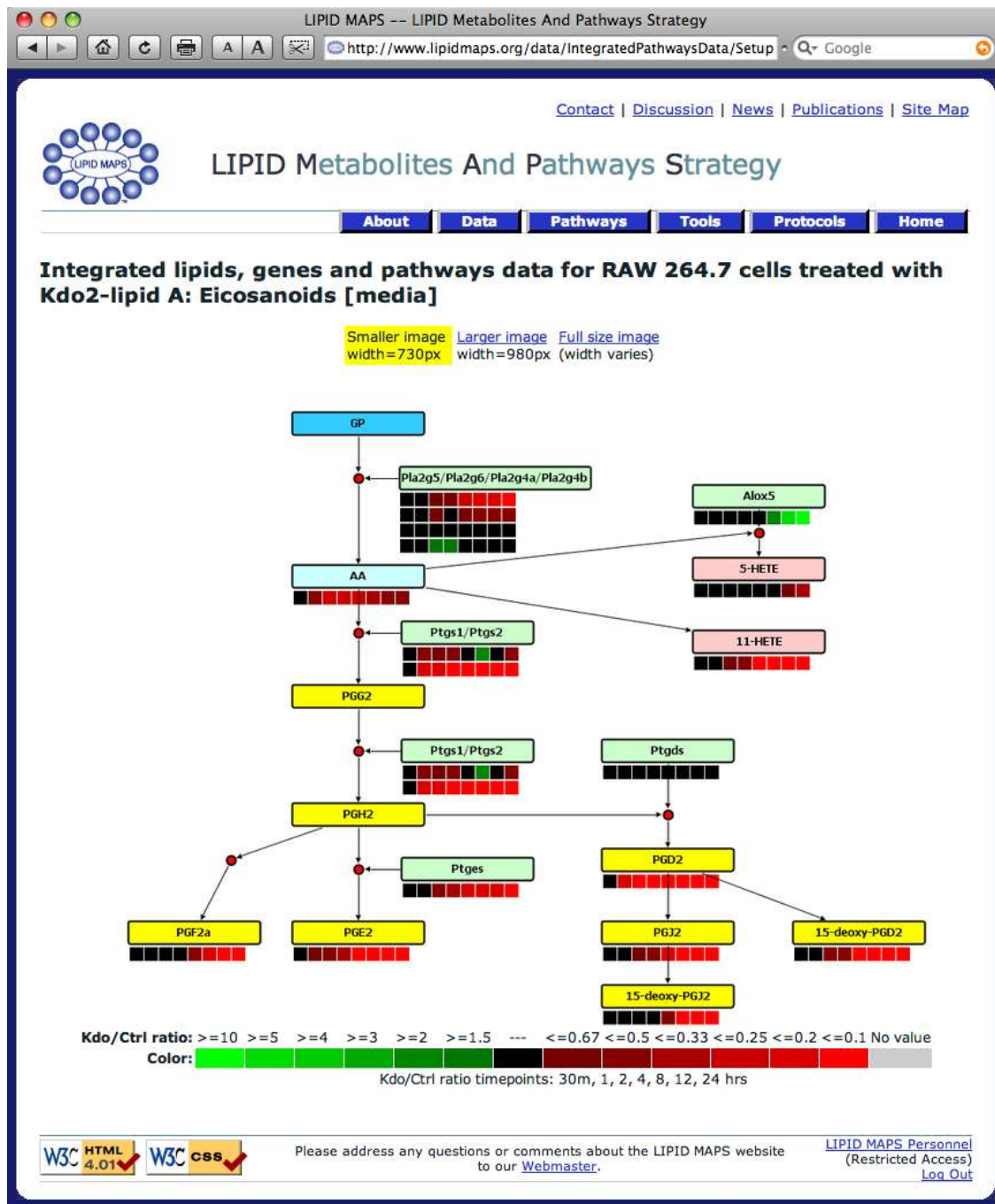
# Interactive network visualization for internet database systems



Grafahrend-Belau et al. (2008): *MetaCrop – A detailed database for crop plant metabolism.* **Nucleic Acids Research**
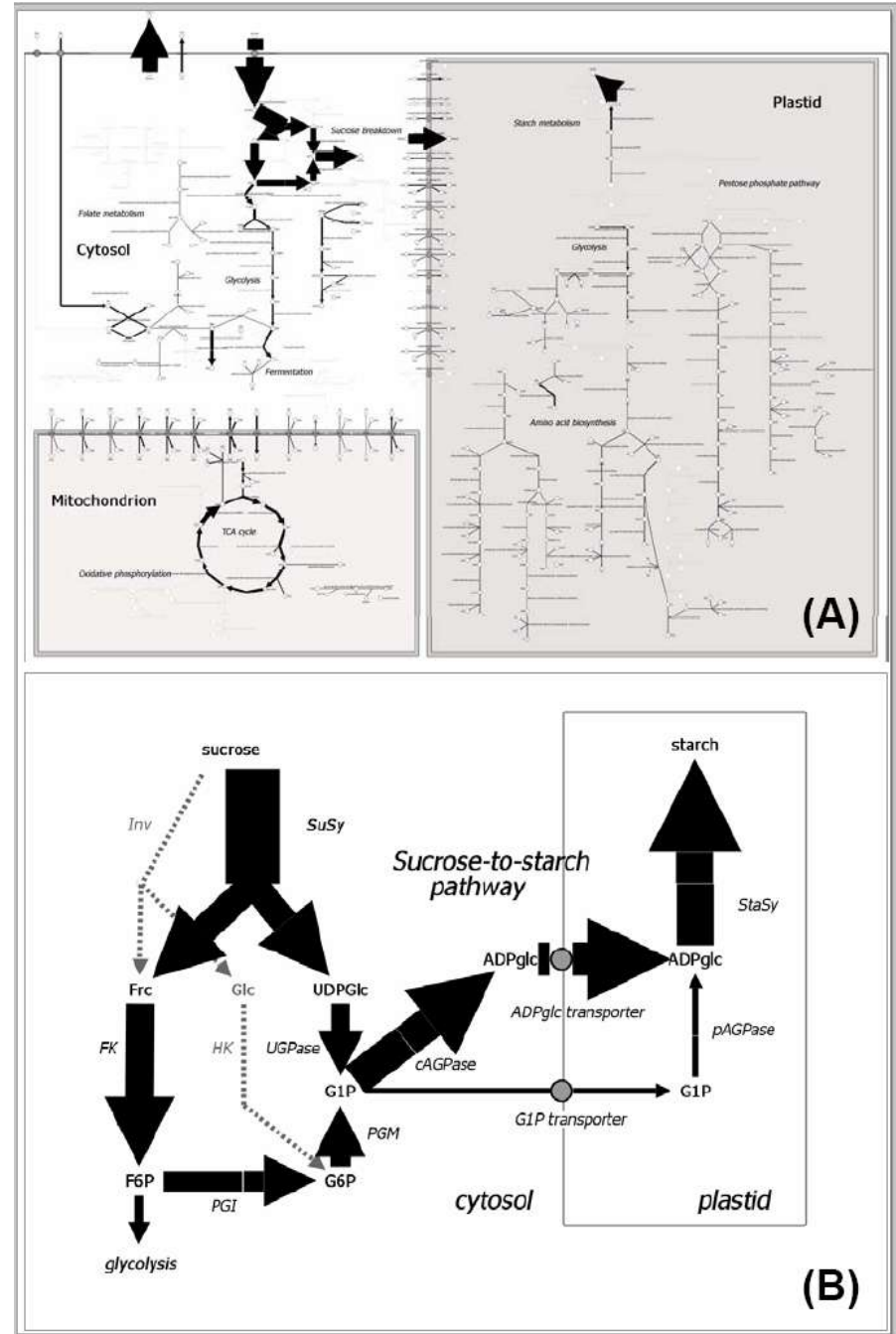
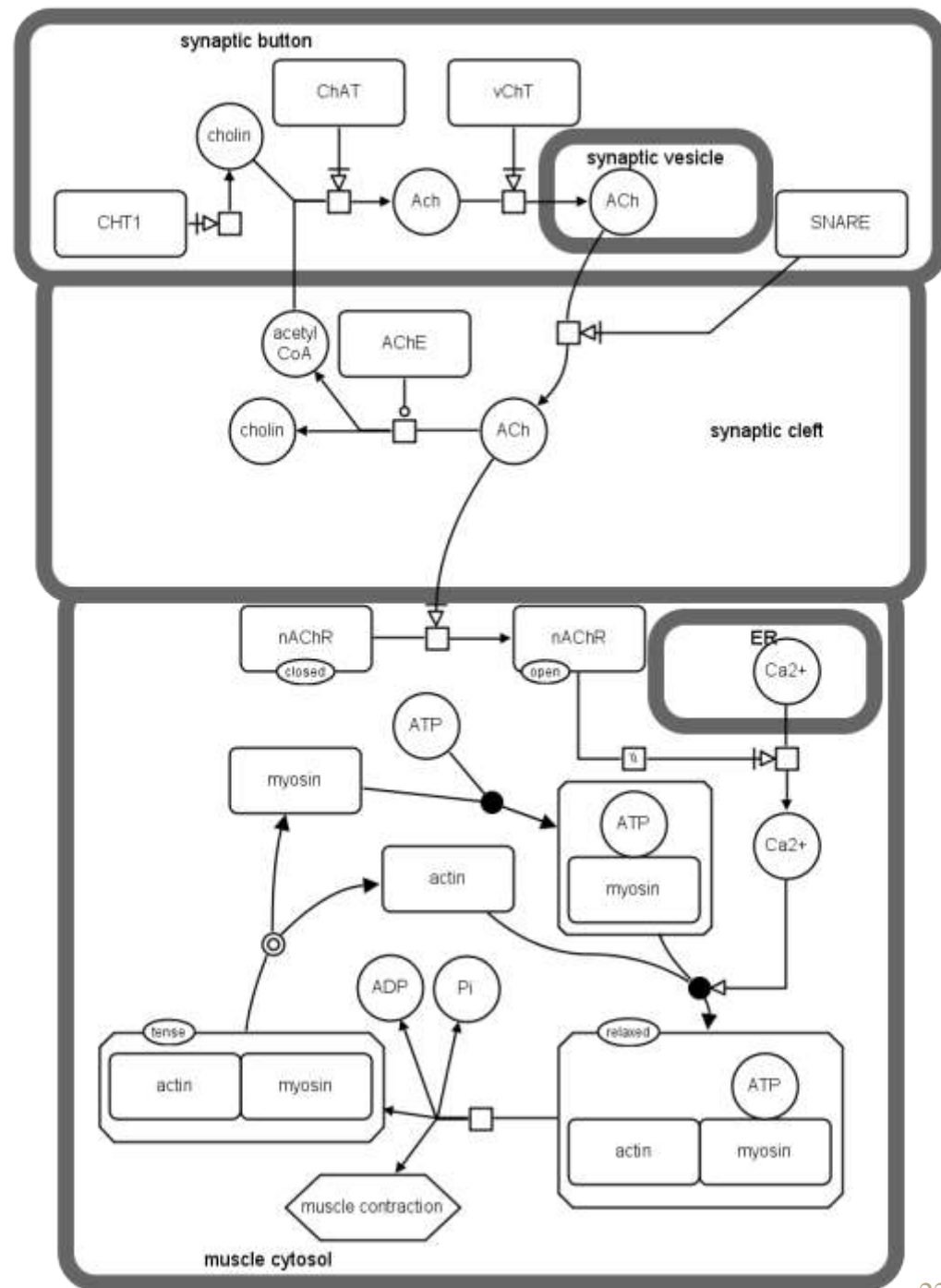# Pathway editing, data visualization

www.lipidmaps.org

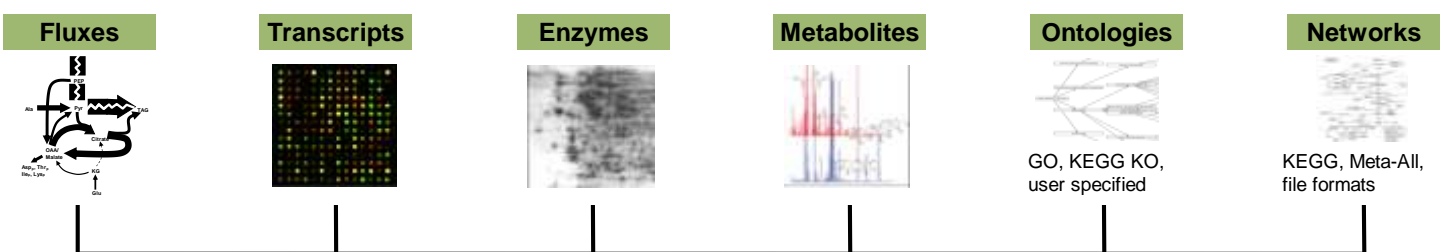# Visualization of metabolic flux steady state simulation results
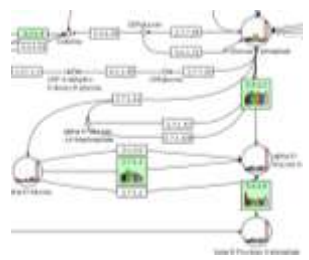


Grafahrend-Belau et al. (2008)

# Support for editing of SBGN diagrams

**Fluxes** **Transcripts** **Enzymes** **Metabolites** **Ontologies** **Networks**

GO, KEGG KO, user specified
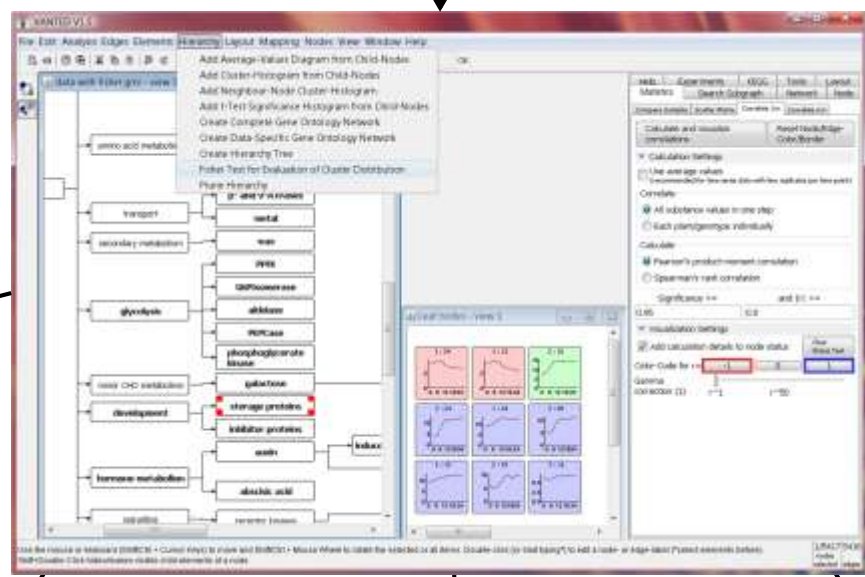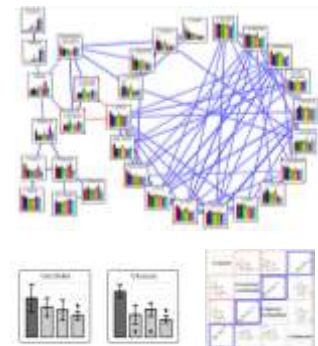
KEGG, Meta-All, file formats
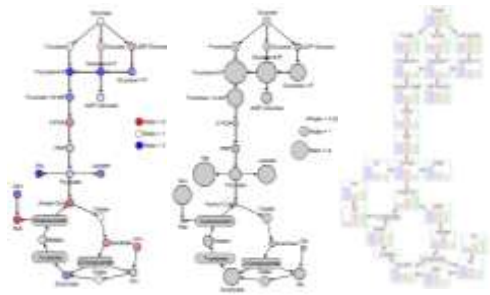
**Combination of data from different omics-areas**

Integration of enzyme and metabolite data into KEGG pathways

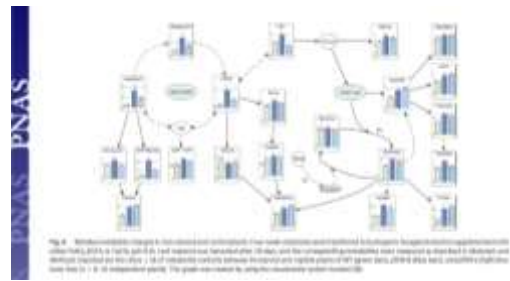**Data analysis, statistics, correlation networks**

**Flexible data visualization**

Color-coding    Shape-coding    Diagrams

**Figures for publications**

Tognetti *et al* (2007), PNAS 104: 11495-11500

**Processing hierarchical data**

http://vanted.ipk-gatersleben.de