# Knowtator, a framework to annotate phenotype-genotype relationships relevant to Arabidopsis leaf growth and development

## Pierre Hilson

pierre.hilson@versailles.inra.fr

pierre.hilson@versailles.inra.fr

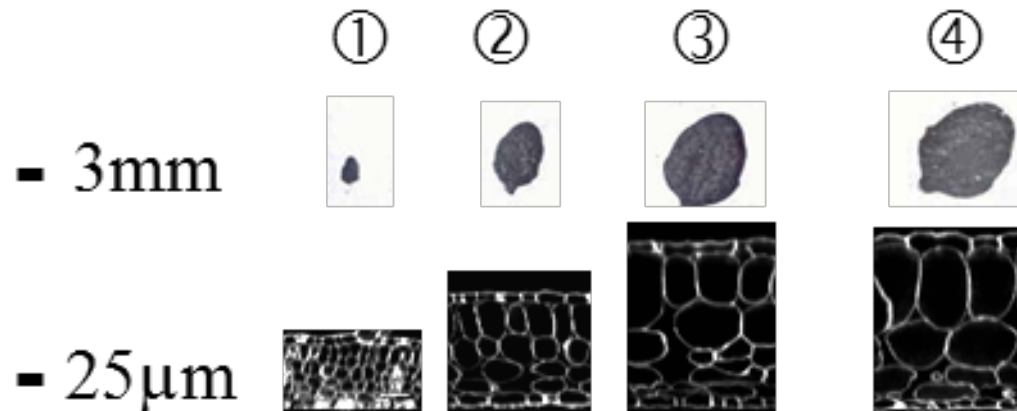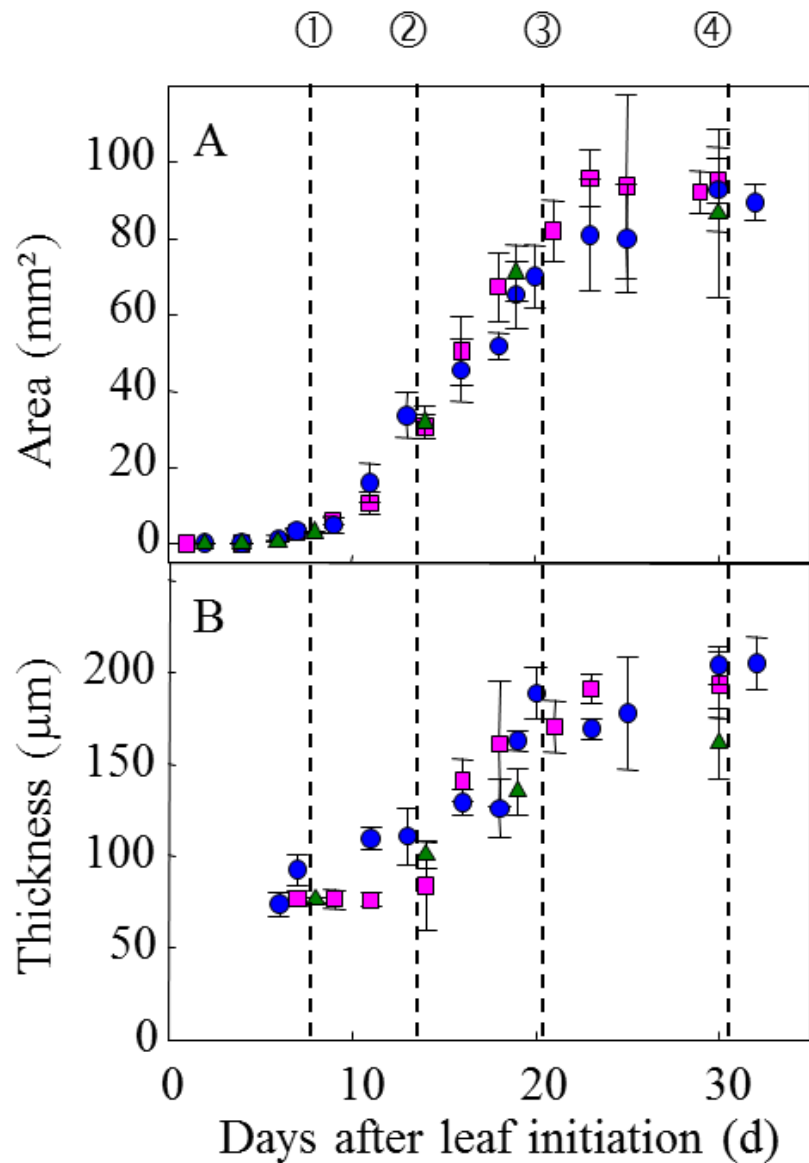# AGRON-OMICS

> = **A**rabidopsis **GRO**wth **N**etwork integrating **OMICS** technologies

Focus: **leaf growth**

## 2. Development of tools for the research community

- **AGRONOMICS1 tiling array** covering both genome strands
- **pep2pro** for comprehensive proteome data analysis, **MASCP Gator** for integration and visualization of Arabidopsis proteomics data, ORFeome resources and cloning resources
- Protein-protein interaction networks- **Arabidopsis Interactome Mapping**
- Improved methods for protein localization studies
- Enzymatic and metabolic networks involved in biomass production
- Plant structure visualization by high-resolution X-ray computed tomography
- Data mining and integration, **CORNET , Arabidopsis Reactome**

**http://www.agron-omics.eu/**

① ② ③ ④

A

Area (mm²)
100
80
60
40
20
0

■ 3mm

■ 25μm

B

Thickness (μm)
200
150
100
50
0

0    10    20    30

Days after leaf initiation (d)

① ② ③ ④
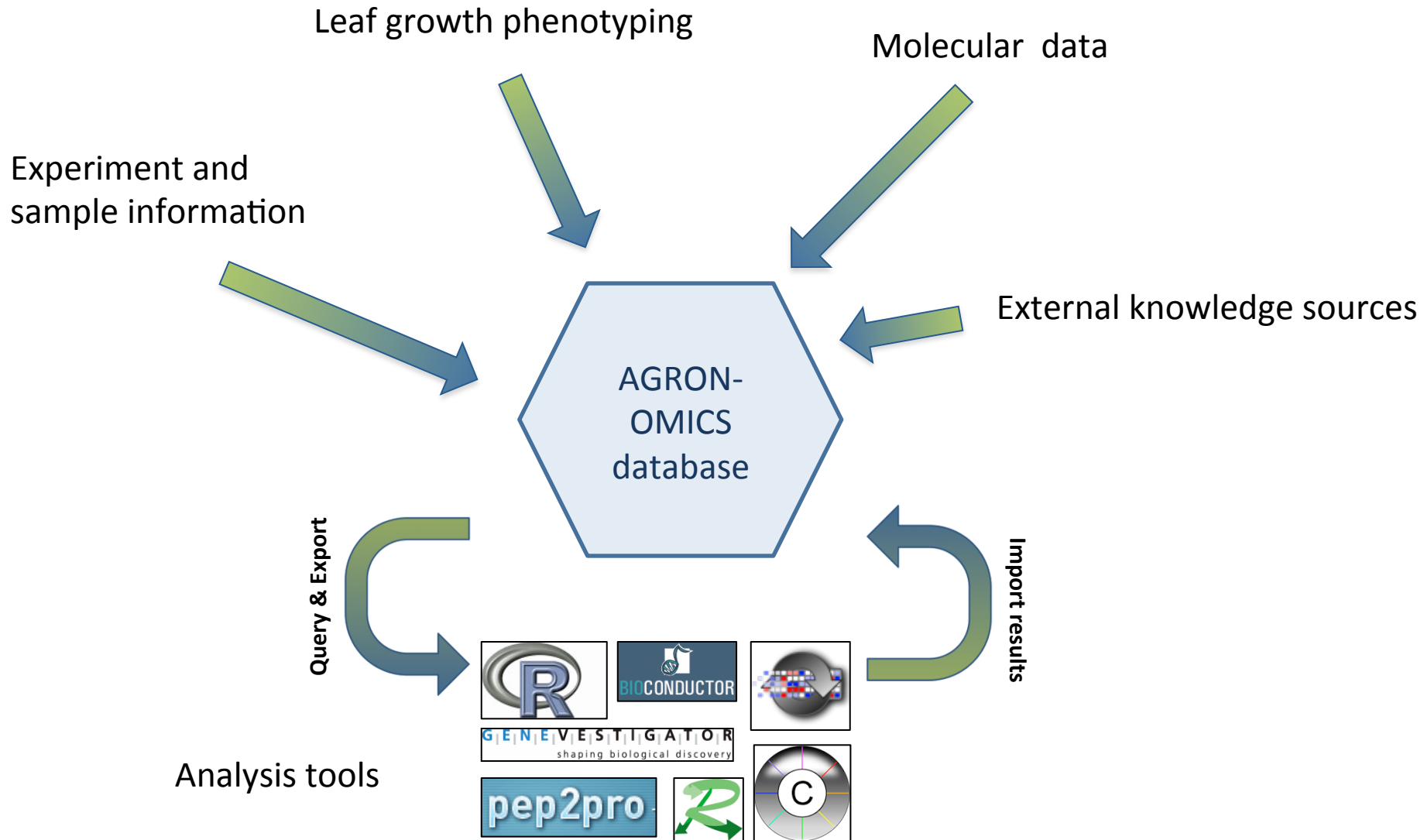
Stage characterised by
(Boyes identifiers in SOW)

Stage 1
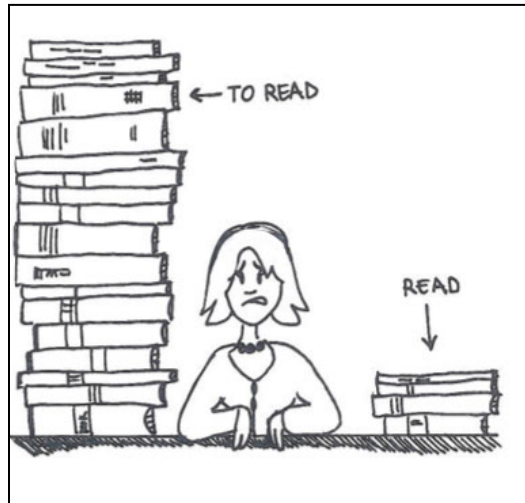active cell division (1.07)

Stage 2
rapid cell expansion (1.11)

Stage 3
decreasing cell expansion (1.18)

Stage 4
mature leaf, no expansion and well before
senescence (> 1.28)

# The AGRON-OMICS database



Sean Walsh, ETH Zurich

# Need to access data in primary literature

There are no good comprehensive literature databases for plant research.

<u>AIM:</u> To capture and mine knowledge from literature

<u>PREMISES:</u>

- Papers about leaf growth and development
- Information collected with biomedical ontologies
- Structured statements that can be simplified and imported into relational databases

<u>TASKS:</u>

- Develop a method (test library, 174 papers)
- Annotate papers with volunteers (adding 109 papers)
- Merge with public molecular resources
- Deploy query and visualization tools

# Further rules and restrictions

- Only Arabidopsis

- Only leaf data (also partly cotyledon, meristem, embryo)

- Restricted to the *Results* section

- Exclude reviews

- Information structure, one-to-one relations

# What sort of information to record?

| Relation | Example |
| --- | --- |
| Phenotype | The rot3-2 allele causes enlarged leaf blades |
| Gene expression | ANT mRNA accumulated in leaf |
| Feature | AtCPL2 contains one dsRNA-binding domain |
| DNA-protein interaction | ARF2 ... bound to the promoter region of GH3.1 |
| Genetic interaction | hyl1 ... appeared to suppress the as2 phenotypes |
| Protein-protein interaction | AN3 interacted strongly with ... AtGRF9 |
| Process | RHL2 ... involved during endocycles |
| Regulation of gene expression | AtCPL1 ... negative regulators of RD29A expression |
| Regulation of process | AN3 ... promoting ... cell proliferation |
| Regulation of phenotype | PHABULOSA ... influence leaf shape |

# Structured statement

| Phenotype | The reduced leaf area in the *hub1-1* mutant was confirmed by morphological measurements of the fully expanded leaves 1 and 2 |
|---|---|

| Slot | Original text | Ontology |
|---|---|---|
| Developmental stage | fully expanded leaves | 3 leaf fully expanded_PO:0001053 |
| Factuality | | |
| Genotype | hub1-1 | mutated gene_MI:0804<br>RDO4 HUB1_AT2G44950<br>loss of function_APO:0000011<br>homozygous diploid _APO:0000229 |
| Growth condition | | |
| Localisation | | |
| Methodology | | |
| Plant part | leaf | leaf_PO:0025034 |
| Process | | |
| Property | area | area_PATO:0001323 |
| Value | reduced | decreased area_PATO:0002058 |

# Biomedical ontologies

Defined terms
Avoiding redundancy and confusion
Established parent-child relationships
Community endorsed

| Ontology | Acronym | URL | Reference |
|---|---|---|---|
| BRENDA tissue / enzyme source | BTO | http://www.brenda-enzymes.info | Gremse *et al.* 2011 |
| Gene Ontology | GO | http://www.geneontology.org/ | Ashburner *et al.* 2000 |
| Molecular Interaction | MI | http://psidev.sf.net<br>http://psidev.sourceforge.net/molecular_interactions/xml/doc/user/index.html | Hermjakob *et al.* 2004 |
| Phenotype, Attribute and Trait Ontology | PATO | http://obofoundry.org/wiki/index.php/PATO:Main_Page | |
| Plant environmental conditions | EO | http://www.gramene.org/plant_ontology/ontology_browse.html#eo | Liang *et al.* 2008 |
| Plant Ontology | PO | http://www.plantontology.org/ | Jaiswal *et al.* 2005 |
| The Arabidopsis Information Resource | TAIR | http://arabidopsis.org/ | Lamesch *et al.* 2012 |

# Annotation flow chart

**Paper selection**
Web browser-html
Acrobat reader-pdf

**Annotation**
Protégé/Knowtator-pprj,
pins, pont

**Quality control**
Computer algorithm-txt

**Database update**
Parsing algorithm
MySQL

# Knowtator, a custom annotation interface

- "general-purpose text annotation tool"

- is a plug-in to Protégé, free, open-source platform to construct domain models and knowledge-based applications with ontologies

- flexible, can be adapted to the project

- easy to import and handle existing ontologies

- export to xml format

- small project files

- easy to share

- mistakes can be corrected

- relatively user friendly

# Knowtator interface



**ORIGINAL TEXT**

**SLOTS**

**Ontologies in use**

**Uploaded external ontologies**

**Relationship categories**

# Text tagging and annotation

# Monitoring relation consistency

- Rigorous guidelines and training of community curators (hands-on sessions, documentation

- Records quality checked with scripts designed detecting different types of errors
  - completeness of relation annotations (i.e. were required slots filled)
  - consistency of ontology terms
  - report of orphan annotations or seemingly undefined ontology terms

- Logs examined by curators, relations adjusted when necessary

Relations produced by reference annotator highly consistent and complete 174 curated articles, the quality control script reported on average only **3.1** missing slots reported per article (not expected to be zero, missing textual info)
19,267 required slots in total, on average of 111 per article; 2.8% missing slots

Relations encoded originally by twelve community annotators contained on average **4.9** missing slots per article, dropped to **2.8**
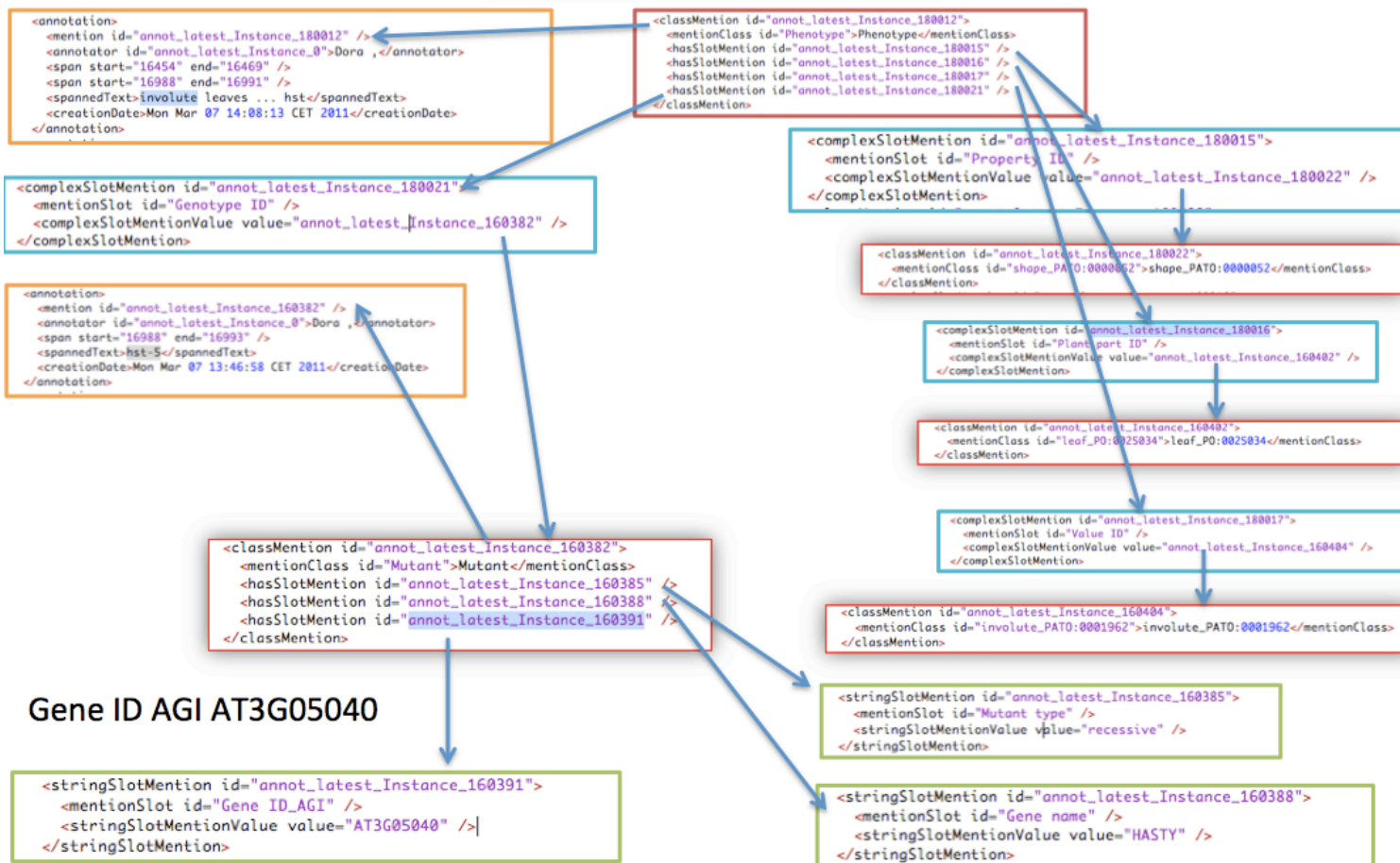
# Knowtator output format

Protégé/Knowtator exports XML representing objects and pointers to objects

# Parsing the data

```
Record number 78
                    file : 10430960.txt.knowtator.xml
                   class : Phenotype
             spannedText : The rot3-2 allele causes ... short petioles
          Annotated_Text : Plant part ID=petioles|Genotype ID=rot3-2|Property Slot=NULL|Value ID=short
   Growth Condition Slot :
      Developmental stage :
           Plant part ID : petiole_PO:0020038
         Localisation ID :
           Property Slot : length_PATO:0000122
              Process ID :
                Value ID : decreased length_PATO:0000574
           Regulation ID :
       Gene expression ID :
            Gene studied :
         Interaction type :
          Protein studied :
       Interactor protein :
              Protein ID :
             Gene target :
       Genetic interactor :
           DNA target ID :
             Genotype ID : mutated gene_MI:0804
                Genotype : Gene ID=ROT3_AT4G36380 | Genotype_Zygosity=homozygous diploid _APO:0000229 | Mutant LOF_GOF ID=gain of function_APO:0000010
            Factuality ID :
```

# Parsing the data

| Information type | Record |
|---|---|
| file | 10430960.txt.knowtator.xml |
| class | Phenotype |
| SpannedText | The rot3-2 allele causes ... short petioles |
| span | 10725\|10749,10776\|10790 |
| Annotated_Text | Plant part ID=petioles\|Genotype ID=rot3-2\|Property Slot=NULL\|Value ID=short |
| Developmental stage | |
| Factuality | |
| Genotype ID | mutated gene_MI:0804 |
| Mutant info | Genotype_Zygosity=homozygous diploid _APO:0000229\|Mutant LOF_GOF ID=gain of function_APO:0000010 |
| AGI | Gene ID=ROT3_AT4G36380 |
| Growth condition | |
| Localisation | |
| Methodology | |
| Plant part | petiole_PO:0020038 |
| Process | |
| Property | length_PATO:0000122 |
| Value | decreased length_PATO:0000574 |

# KnownLeaf MySQL database

Now in a position to begin asking questions of the data as a whole
e.g. give list of AGIs involved in leaf epidermal phenotypes

# Overview of KnownLeaf database content

| Relation category | # AGI | # unique AGI | Ratio |
| --- | --- | --- | --- |
| Phenotype | 5608 | 381 | 14.72 |
| Gene expression | 4767 | 704 | 6.77 |
| Genetic interaction | 658 | 186 | 3.54 |
| Feature | 462 | 175 | 2.64 |
| Protein-protein interaction | 310 | 121 | 2.56 |
| Process | 235 | 140 | 1.68 |
| Regulation of gene expression | 204 | 70 | 2.91 |
| Regulation of process | 178 | 85 | 2.09 |
| DNA-protein interaction | 92 | 47 | 1.96 |
| Regulation of phenotype | 20 | | |
| Total | 12534 | | |

# Representation of collected annotations merged with public knowledge sources

- Interactive graph in Cytoscape

- Selected Knowtator relations part of a larger network of objects to interpret and further expand these relations

- Represented edges
  - Co-expression (ATTED-II mutual rank score ≥ 25)
  - Protein-protein interactions (Y2H, AI-1 interactome)
  - Phenotype relations (gene-mutation-phenotype)

- Seeded nodes
  - Genes/proteins annotated via Knowtator (AGI code)
  - Proteins that vary significantly across leaf development (Bärenfaller et al., 2012, Mol Syst Biol)
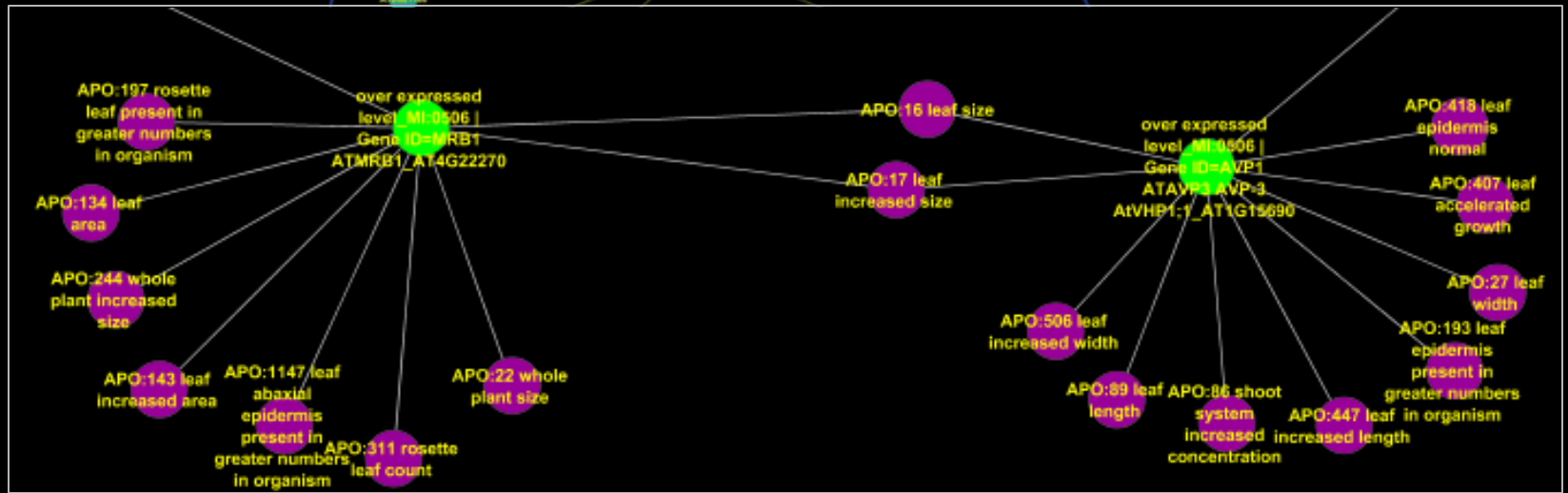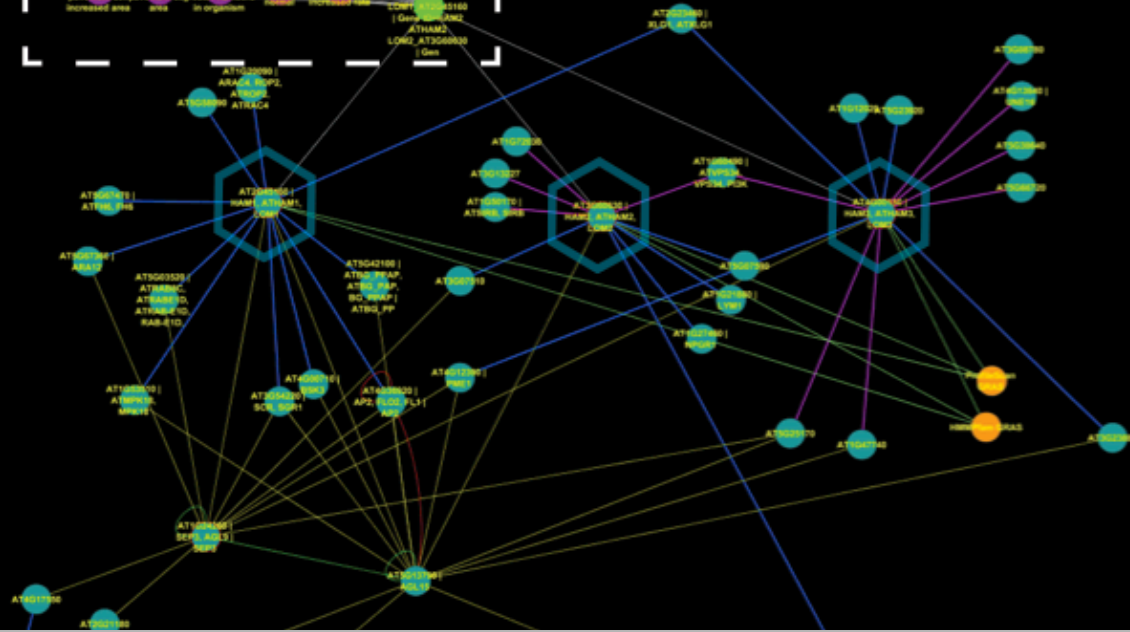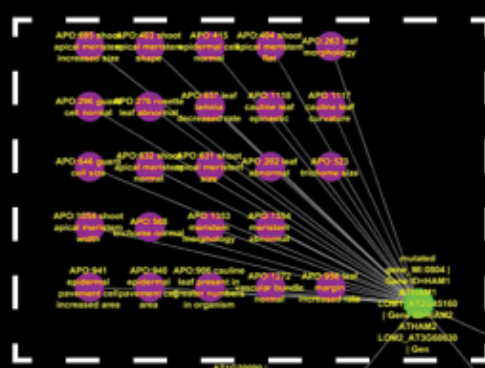  - Connected protein/gene (co-expression, PPI)
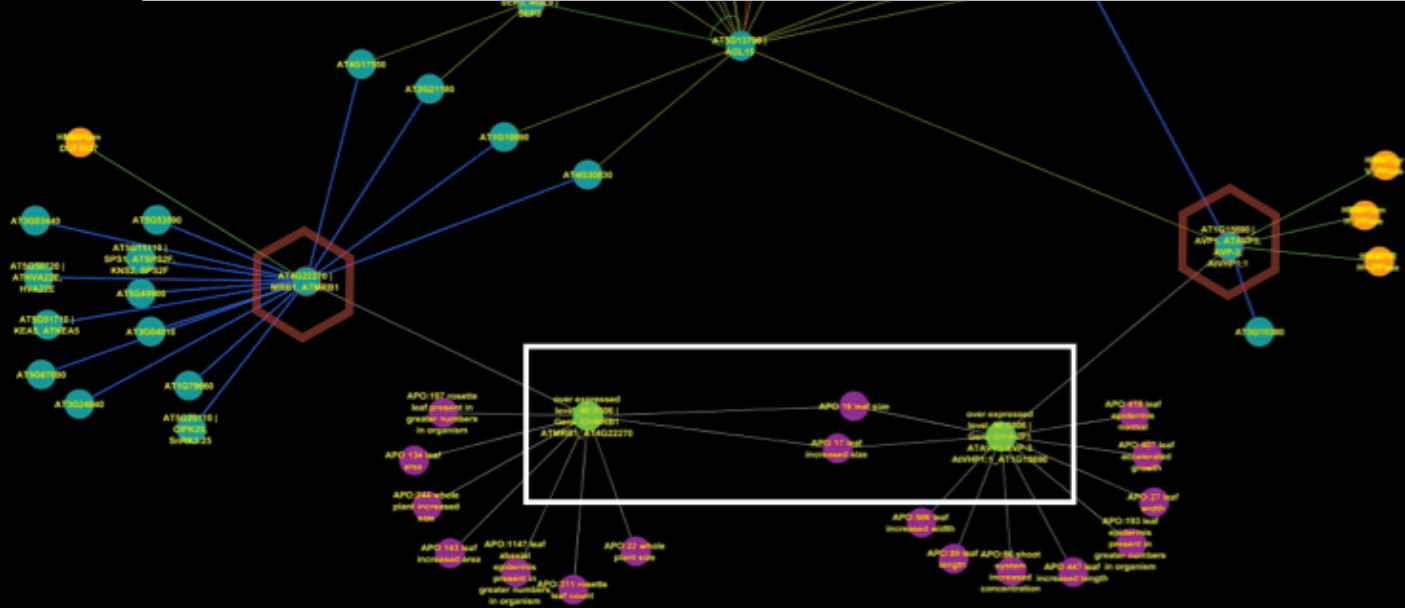
**-> LeafNet: 19,055 nodes connected by 39,649 edges**

LeafNet

LeafNet

# LeafNet

LeafNet

LeafNet

# Knowtator is useful for literature annotation

Keeping the original text → For humans → Literature database

Working with ontologies → Machine readable → Modeling / Automated text mining efforts

3 papers/day

Training to use the program: 3 hours intro, few days of practice

Easy to share projects (6 relatively small files + textsources)

Program is free

# Engage your community

- Biological systems are complex, many genes/proteins are involved in any given process

- Many scientific articles are published... or will be soon

- Few willing annotators

- Long term benefits obvious, short term investments not so

- Practical tools are required

# Perspectives

- Mining of the KnownLeaf database

- Exploration of LeafNet graphs to generate hypotheses

- Exhaustive annotation of all relevant papers in the "leaf growth and development" domain (now about 1/3)

- Implementation of the Leaf Knowtator annotation system in other domains

- Training of machine learning algorithms embedded in automated text mining tools with the Knowtator-generated data set

- All tools and data will soon be made public
  … except spanned text because of copyright restrictions

# Contributors

Fabio Fiorani, initiator (Phenote)

**Dóra Szakonyi, reference annotator and developer of Leaf Knowtator**
**Sofie Van Landeghem, text mining expert**
**Sean Walsh, KnownLeaf database and LeafNet**
Pierre Hilson, molecular biologist

**Community annotators**: Katja Bärenfaller, Lieven Baeyens, Jonas Blomme, Rubén Casanova-Sáez, Stefanie De Bodt, David Esteve-Bruna, Nathalie Gonzalez, Jesper Grønlund, Richard G.H. Immink, Sara Jover-Gil, Asuka Kuwabara, Tamara Muñoz-Nortes, Aalt-Jan van Dijk, David Wilson-Sánchez

PIs: Vicky Buchanan-Wollaston, Gerco C. Angenent, Yves Van de Peer, Dirk Inzé, José Luis Micol, Wilhelm Gruissem,

Researchers groups:
Department of Plant Systems Biology, VIB, Ghent University
Department of Biology, ETH Zurich
División de Genética, Universidad Miguel Hernández, Elche
Warwick Systems Biology Centre, University of Warwick
Plant Research International (PRI), Bioscience, Wageningen