



# Information Extraction Challenge

## *Gene Regulation Network*

### *in Arabidopsis thaliana (GRNA)*



L. Lepiniec  
B. Dubreucq  
A. Fatihi  
D. Valsamou

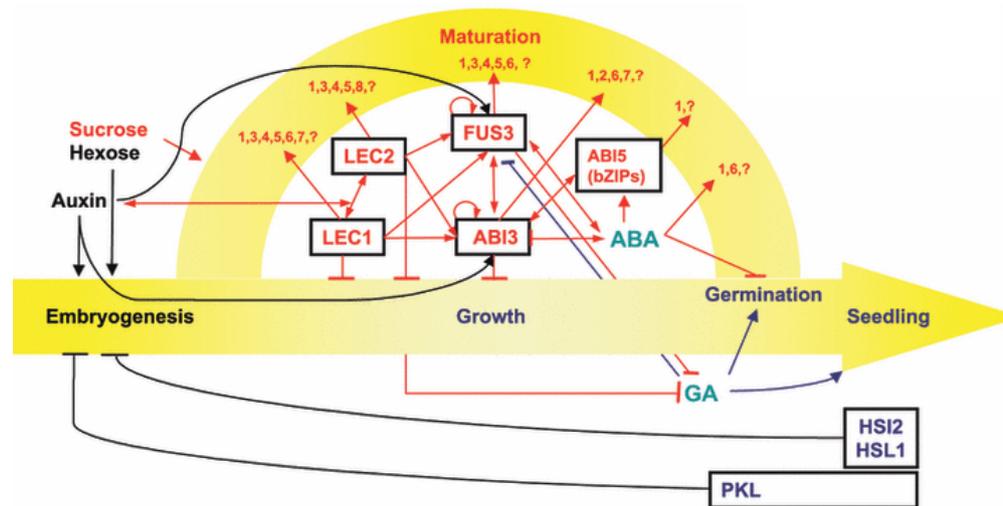
C. Nédellec  
P. Bessières  
R. Bossy  
L. Deléger  
E. Chaix  
D. Valsamou

P. Zweigenbaum  
D. Valsamou



# Biological motivations

- *A.thaliana*, a model plant
  - From *A.th* to other species
- Seed development: industrial, agricultural and fundamental research interest
- Complex regulatory networks that involve genetic, molecular and physiological mechanisms



A model of genetic (framed) and molecular (in blue cyan) interactions involved in the control of seed development and maturation in *Arabidopsis thaliana* (from Santos-Mendoza *et al.*, 2008) (Numbers are biological processes).

# Biological motivations

The screenshot shows the NCBI PubMed search interface. The search query is 'arabidopsis genetic regulation'. The results page displays 'Search results' with 'Items: 1 to 20 of 22877'. The first result is a link to a paper titled 'SKIP Interacts with the Paf1 Complex to Regulate Flowering via the Activation of FLC Transcription in Arabidopsis' by Cao Y, Wen L, Wang Z, Ma L. The paper is from Mol Plant, 2015 Sep 14, pii: S1674-2052(15)00365-2, doi: 10.1016/j.molp.2015.09.004. The page also shows navigation options like '<< First', '< Prev', 'Page 1 of 1144', 'Next >', and 'Last >>'. On the left side, there are filters for 'Article types', 'Text availability', 'PubMed Commons', and 'Publication dates'.

- Crop improvement: improved selection of novel species with a better knowledge of genetic interactions
- Knowledge harvest of seed development from text data

# Information extraction motivations

- Text mining: extraction of scarce and critical information

Genetic Regulation Network: some studies but not at this level of complexity

- Requirements
  - ✓ A rich knowledge model
  - ✓ Annotated corpus for relation extraction
- Towards Systems Biology
- International visibility
- Shared task: text-mining community building

# BioNLP Shared Task, information extraction in the biomedical domain

BioNLP-ST'09: 1 task (GENIA)

BioNLP-ST'11: 5 tasks (2 by INRA)

BioNLP-ST'13: 6 tasks (1 by INRA)

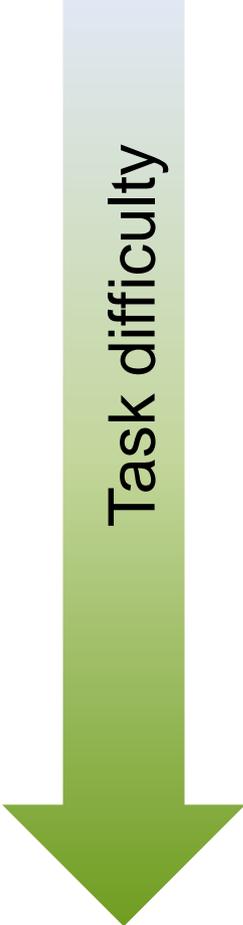
➤ **BioNLP-ST'16: 4 tasks (2 by INRA)**

*Bacteria Biotope*

*Plant Seed Development (SeeDev)*

- Knowledge-based construction from text
- More direct application to Life Science

Task difficulty



# Information Extraction Challenge

## Step 1

Annotation model

Annotated text  
by human (Training set)

Training of learning methods

Classifier

## Step 2

Classifier

Test text

Annotation prediction

Automatically  
annotated text

Evaluation

Test text annotated  
by human (Test set)

## Step 3

Community workshop: 11-12 August 2016  
(Association for Computational Linguistics)

# Information Extraction Challenge

## Step 1

Annotation model

Annotated text  
by human (Training set)

Training of learning methods

Classifier

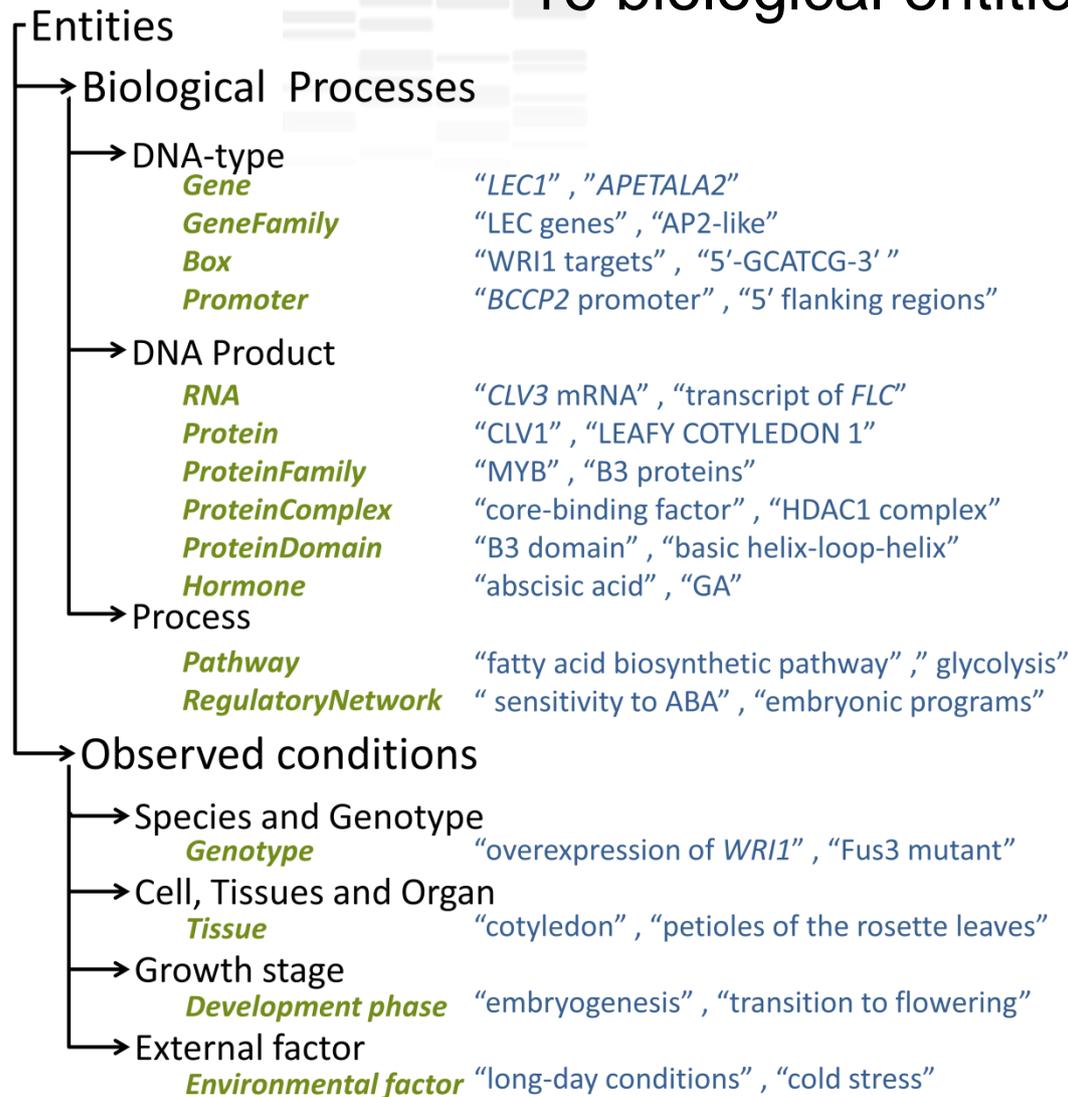
- Critical stage, long and complex
- Involved biology experts
- Assisted by knowledge engineer (CDS funding)
- Guidelines

# BioNLP Shared Task Schedule

November 2015	1st announcement (with sample data set and detailed plans)
January 2016	Training data release and on-line evaluation service open
(mid) March 2016	Evaluation on test data open
(end) March 2016	Evaluation done, notification to participants
8 May 2016	Workshop paper due
5 June 2016	Notification of acceptance
22 June 2016	Camera-ready papers due
11-12 August 2016	Workshop

# Knowledge model

Fined-grained and complex model  
16 biological entities and 10 events



10 relation types were defined :

→ Regulation :

- *RegulatesActivityOf*
- *RegulatesAccumulationOf*
- *RegulatesExpressionOf*

→ Interaction :

- *InteractWith*
- *BindTo*

→ Localisation :

- *IsFoundIn*
- *IsFoundDuring*

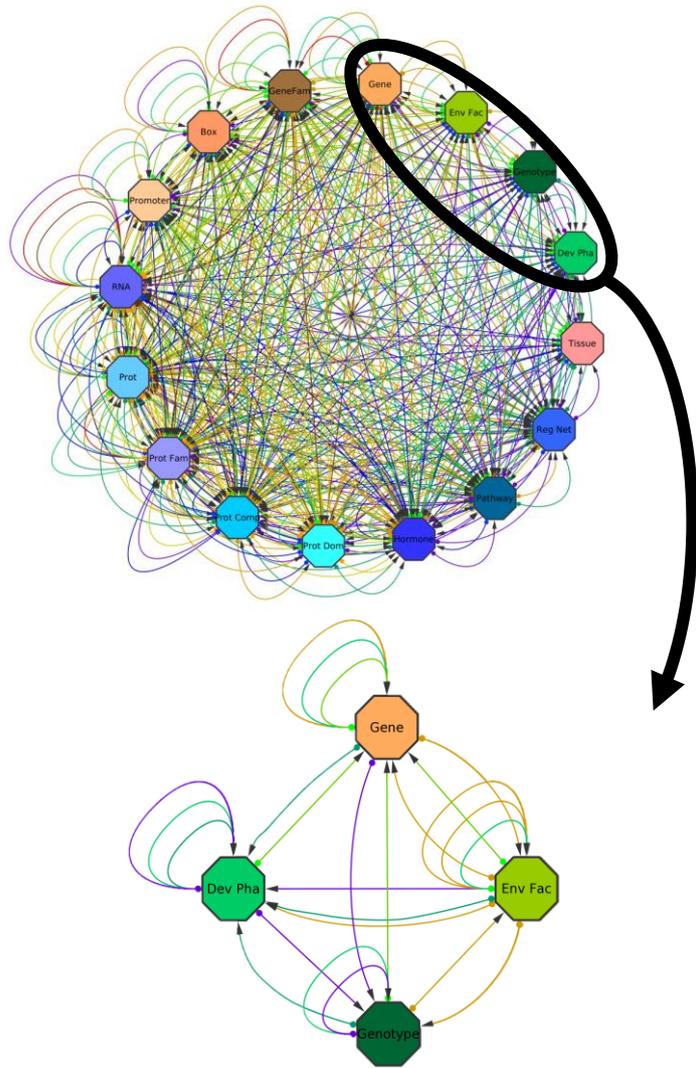
→ Similarity :

- *Comparison*
- *Belongs to*
- *Encodes*

1 relation to define n-ary events

- *Condition*

# Gene Regulation Network in *A.th*



- Integration of the knowledge model in the corpus annotation editor, AlvisAE
- Annotation of a corpus of scientific articles by biology experts (on-going)  
Currently 4,444 entities and 1,421 n-ary events.

# AlvisAE : corpus annotation editor

- A module of Alvis suite, an integrated software library that provides all the necessary tools for semantic analysis and terminology and ontology acquisition for Ontology-based Data Integration.

manual-annotation : [1] Regulates the Stem

Arguments

- Agent: ecd02...9705a AP2
- Target: 8cd9e...a3d0d AG

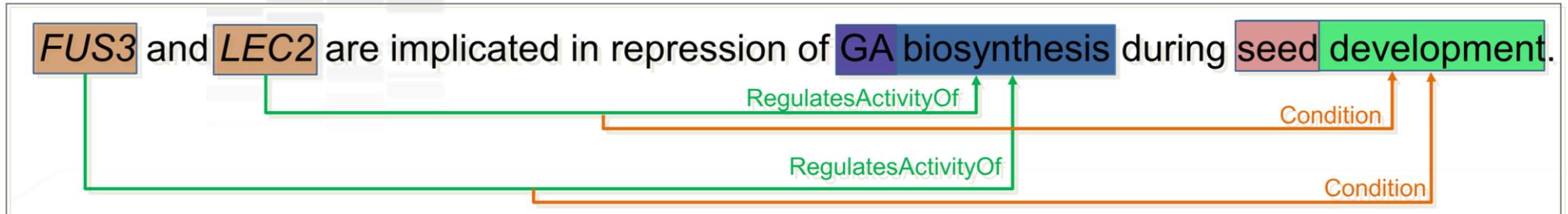
Properties

- modality: inhibition
- speculation
- negation

Id	Annotation	Ki	Type	Details	Visi
estelle 10dbf	@manual-annotation		Protein	AP2	
estelle 149b5	@manual-annotation		Protein	AP2	

(Papazian *et al.*, 2012)

# Example of annotation



## Entities

Gene *FUS3*, *LEC2*

Hormone GA

Pathway GA biosynthesis

Tissue seed

Developmental Phase seed development

## Relation

Agent (Gene) regulates\_activity\_of Target (Pathway)

Agent (Event) has\_as\_condition Target (Developmental Phase)

- ***FUS3* regulates activity of *GA biosynthesis* with condition *seed development***  
(-)

# Further use of *Arabidopsis* corpus



Model and reference corpus of regulation network for *A.th* seed development

BioNLP Shared Task 2016 : *Plant seed development*

Training data and evaluation data for

- relation extraction (D. Valsamou, PhD)
- ontology based annotation (A. Ferré, PhD)

Iindex Paris-Saclay grants

Use case plant in text-mining European project **openMIN7ED**

Extension to species of agronomic interest (wheat)

Integration of knowledge extracted in plant models of system biology (on-going project with IMSV)