# Text-mining methods used for information extraction in plant scientific papers

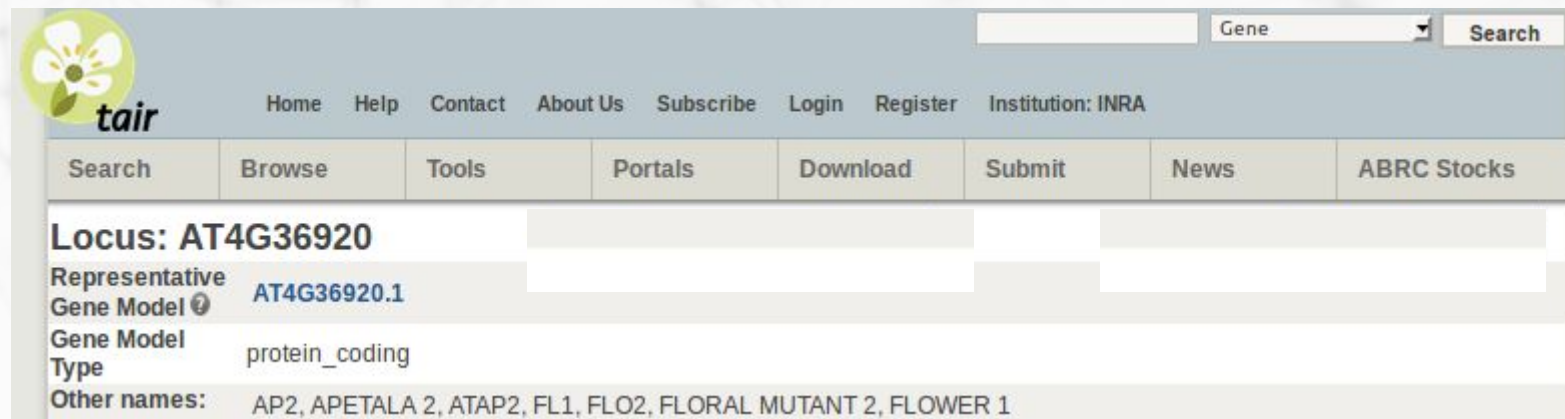## 4. Entity categorization

## Categorization

Different entities predicted could be similar to others, and be linked to a single category.

Categorization allows to extract the information and to go back to a superior level of conceptualization compared to the form of the raw text.

# Categorization

In the case of synonyms: explains that AP2 and APETALA 2 are the same concept, and could be linked to a unique category such as AT4G36920 (locus of the gene)

➢ Improves the query of TDM results
➢ Allows link with external resources

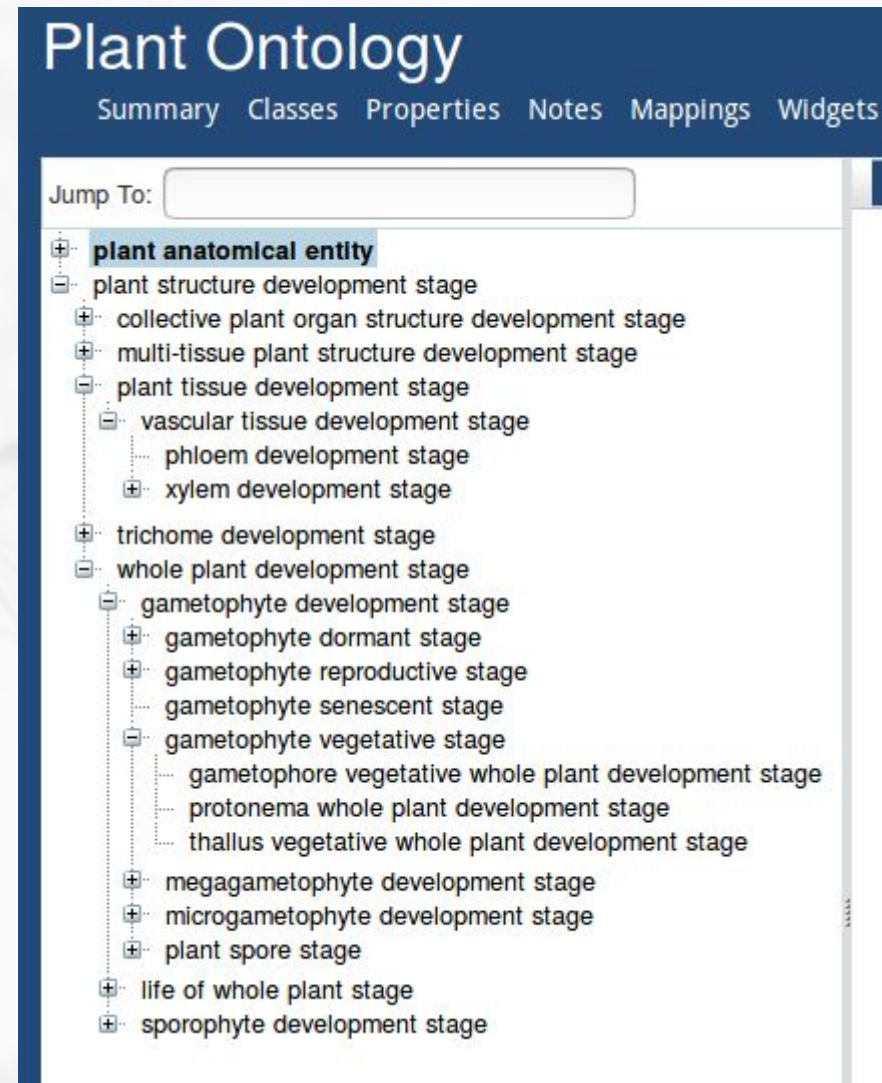# Ontology for Entity recognition & Categorization

To detect and categorize some complex entities such as Development Phase or Tissue, we use ontologies.

An ontology is a semantic representation of knowledge. Concepts are organized and linked together by different relations, for example "is_a", "part_of", "has_part" ...

# Example of an ontology

## Plant Ontology

We can reason
and infer knowledge
from ontologies.

# A plant ontology for TDM use

It will be built to allow detections of complex entities, such as Tissue, Developmental Phase, Pathway ...

It will be used for :
- Entity detection
- Categorization of predicted entities
- The interrogation of the final application