

# The Cell Cycle Ontology: an application ontology for data integration



Erick Antezana<sup>1</sup>, Vladimir Mironov<sup>1,4</sup>, Mikel Egaña<sup>2</sup>, Robert Stevens<sup>2</sup>, Ward Blondé<sup>3</sup>, Bernard De Baets<sup>3</sup>, and Martin Kuiper<sup>1,4</sup>

<sup>1</sup> VIB Dept. of Plant Systems Biology, Ghent University, Technologiepark 927, B-9052 Ghent, Belgium. E-mail: erick.antezana@gmail.com

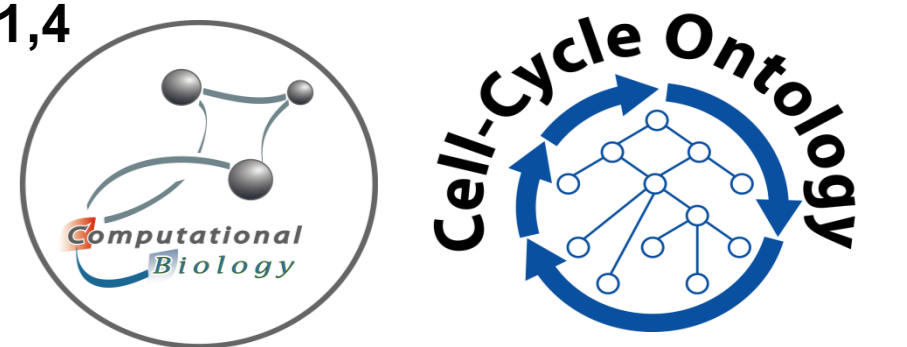
<sup>2</sup> The University of Manchester, School of Computer Science, Oxford Road, M13 9PL Manchester, UK. E-mail: {egaanarm,stevenr}@cs.man.ac.uk

<sup>3</sup> Ghent University, Department of Applied Mathematics, Biometrics and Process Control, Coupure links 653, 9000 Ghent, Belgium. E-mail: {ward.blonde,bernard.debaets}@ugent.be

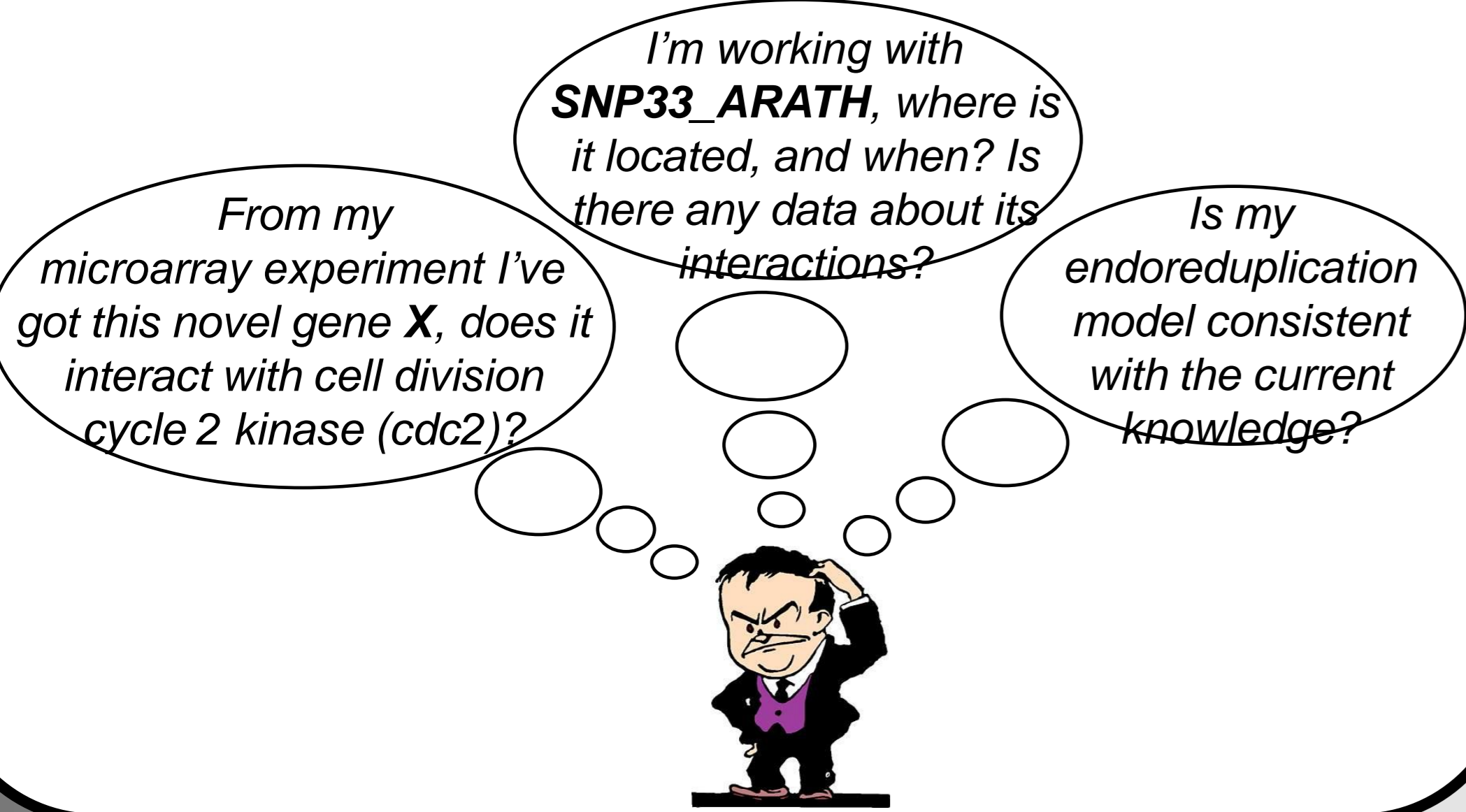
<sup>4</sup> Norwegian University of Science and Technology, Department of Biology, Høgskoleringen 5, 7491 Trondheim, Norway. E-mail: {mironov,kuiper}@nt.ntnu.no

<http://www.CellCycleOntology.org>

Plant Systems Biology



## Motivating scenarios



## Objective

To capture the knowledge about the **cell cycle** process (particularly its dynamic facets) and to promote sharing, reuse and enable better computational integration with existing resources (semantic web). The ultimate aim is to support evaluation and generation of hypotheses via reasoning services about cell-cycle regulation. **Target organisms:** *S. cerevisiae*, *S. pombe*, *A. thaliana* and *H. sapiens*.

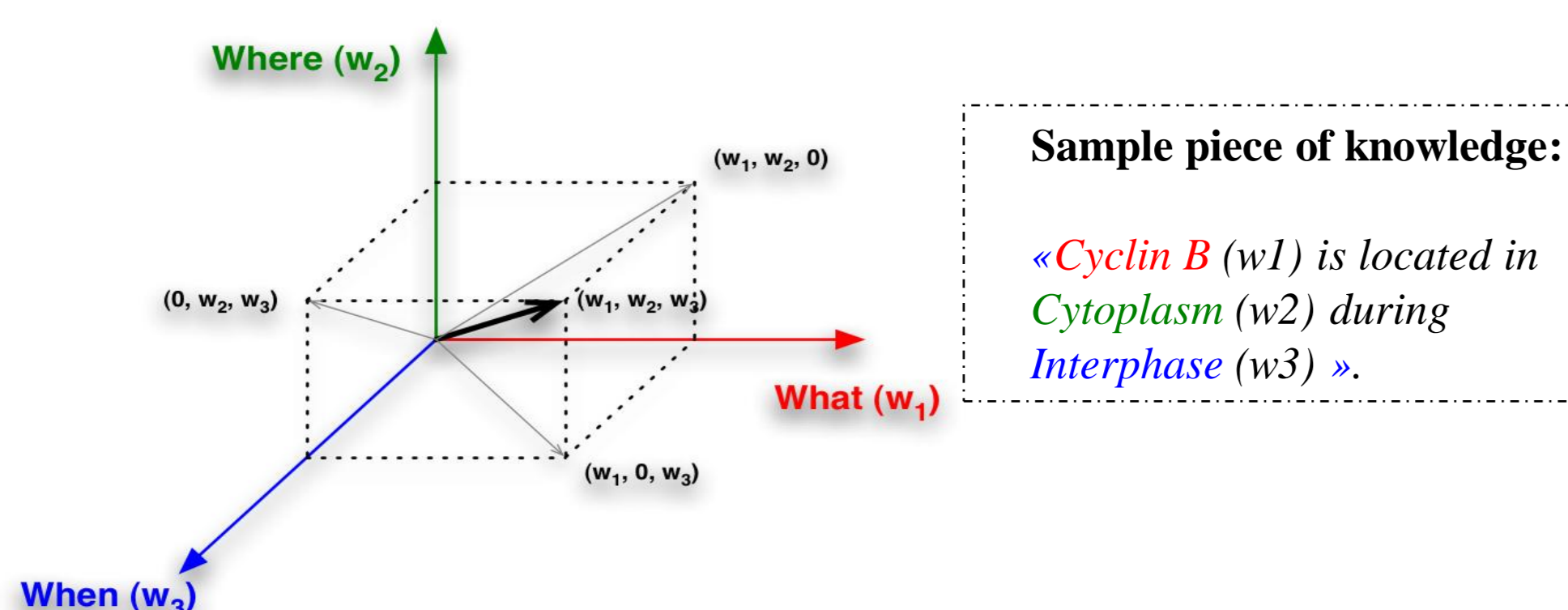


Fig 1. W3 paradigm.

## Data

CCO should capture the semantics and spatio-temporal relationships (Fig 1) of cell-cycle components (proteins, genes, cellular locations, phases, ...). The data sources are:

- GO (Cellular Location branch, and branches for 'cell cycle' (GO:0007049), 'cell division' (GO:0051301), 'cell proliferation' (GO:0008283), 'DNA replication' (GO:0006260) and 'chromosome segregation')
- RO
- MI (IntAct ontology)
- GOA files
- PPI: IntAct
- NCBI taxonomy
- UniProt
- Cell cycle functional data
- Data obtained with bio-tools (e.g. OrthoMCL)



OBO and OWL-DL formats have been chosen for representing the knowledge. RACER is mainly used for checking the data consistency and for doing classifications.

## Data integration pipeline

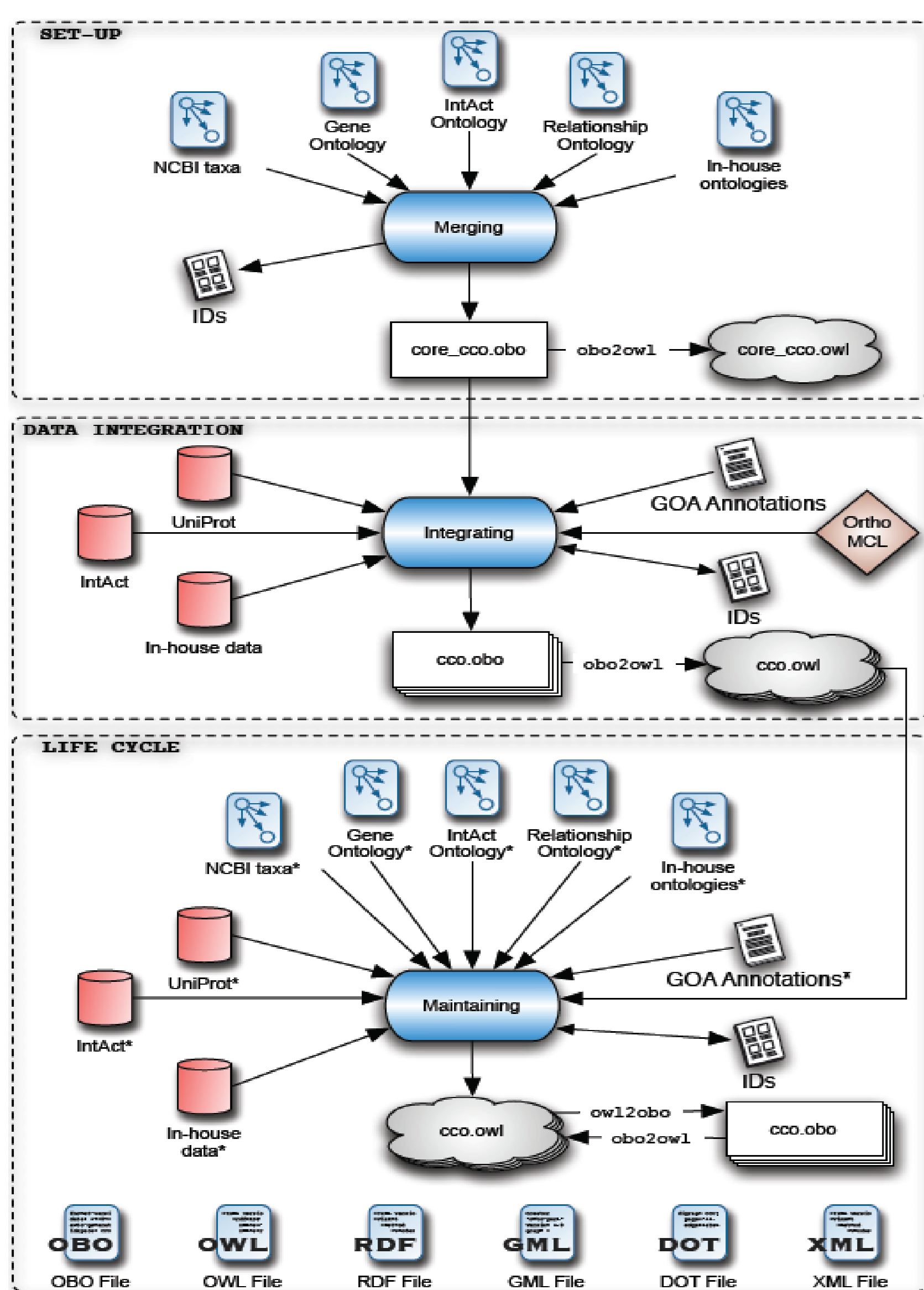


Fig 2. Data integration pipeline.

## Exploring CCO

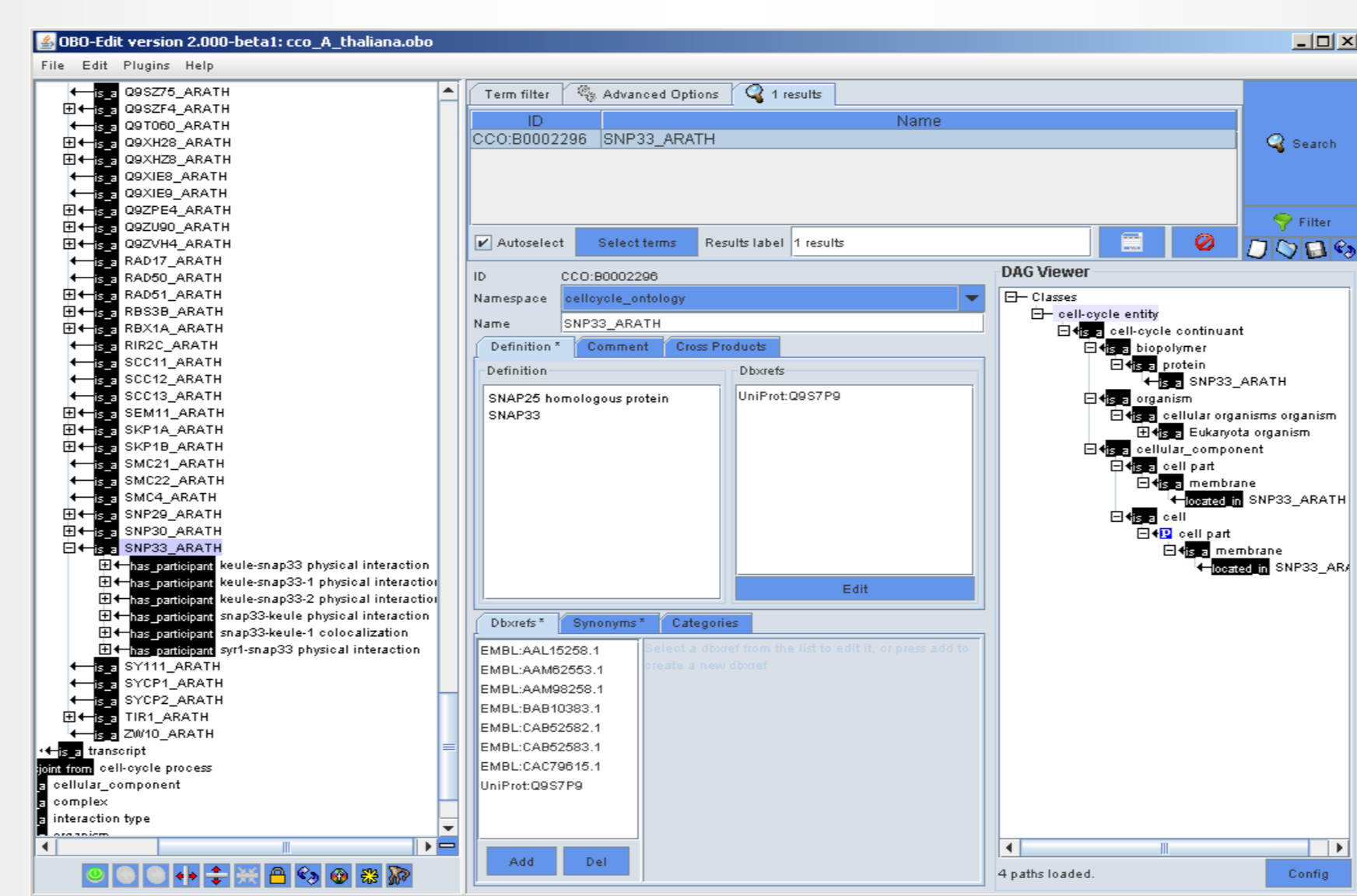


Fig 3. CCO in OBO-Edit.

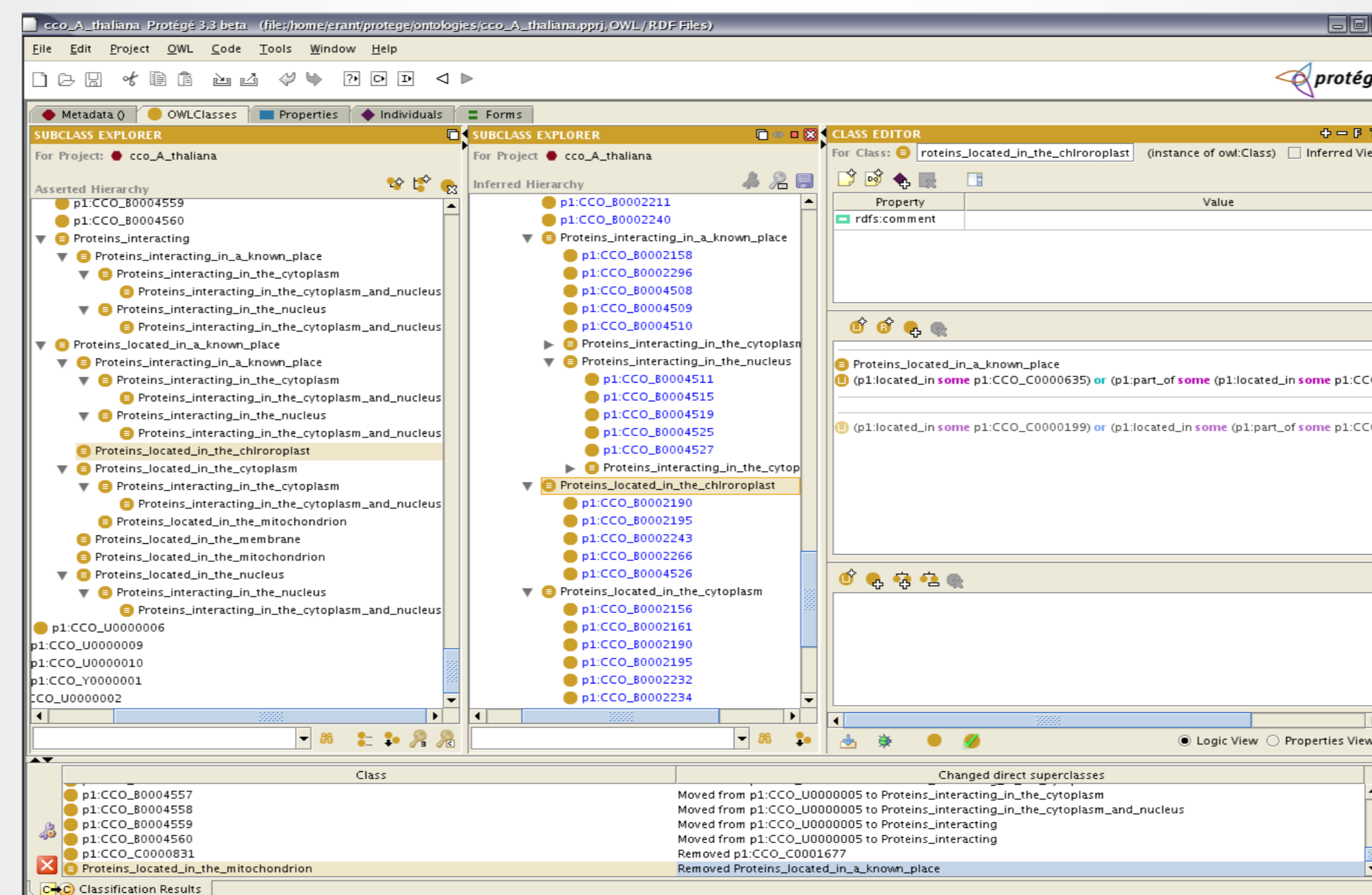


Fig 4. CCO in Protégé.

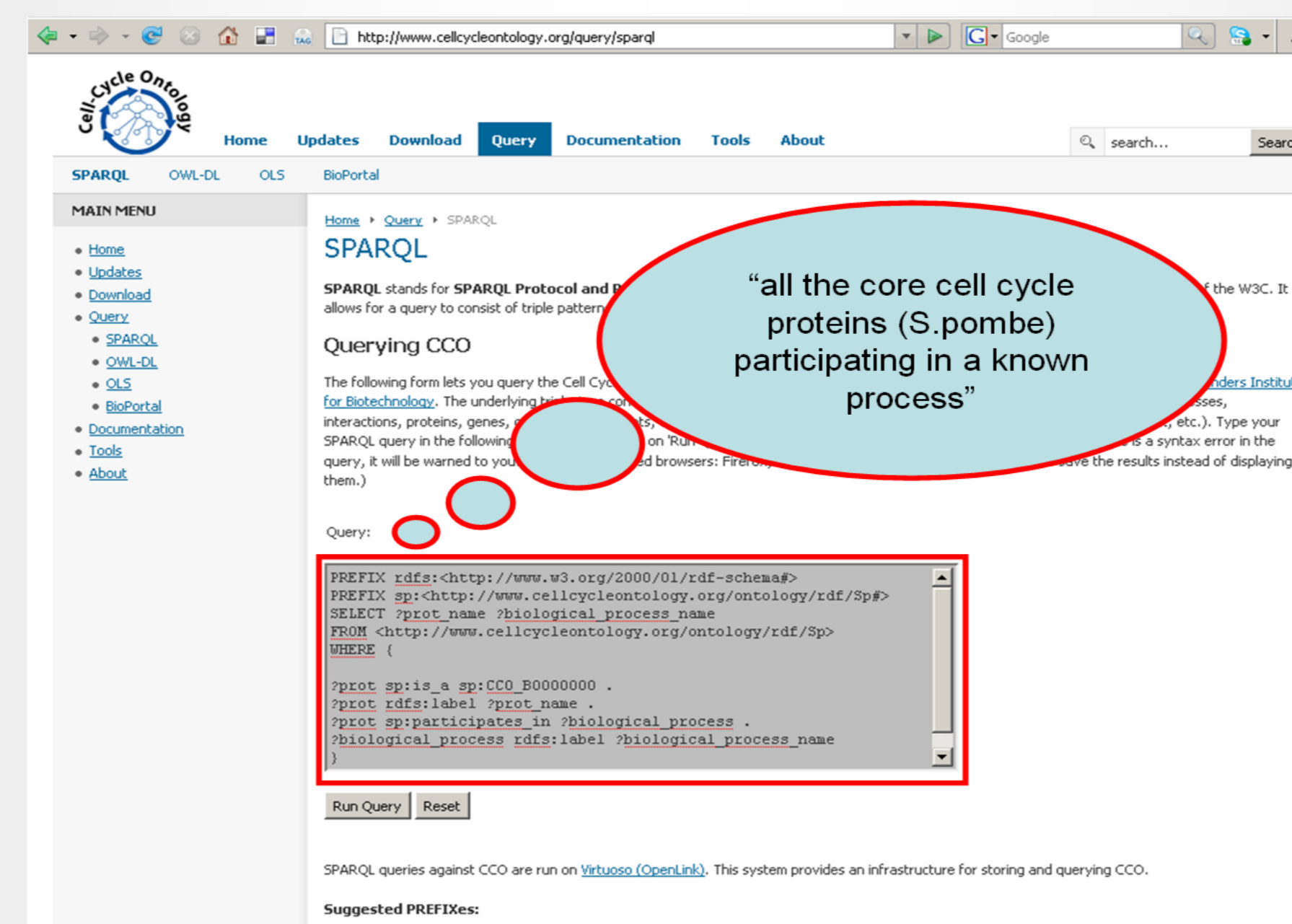


Fig 5. CCO web site (online).

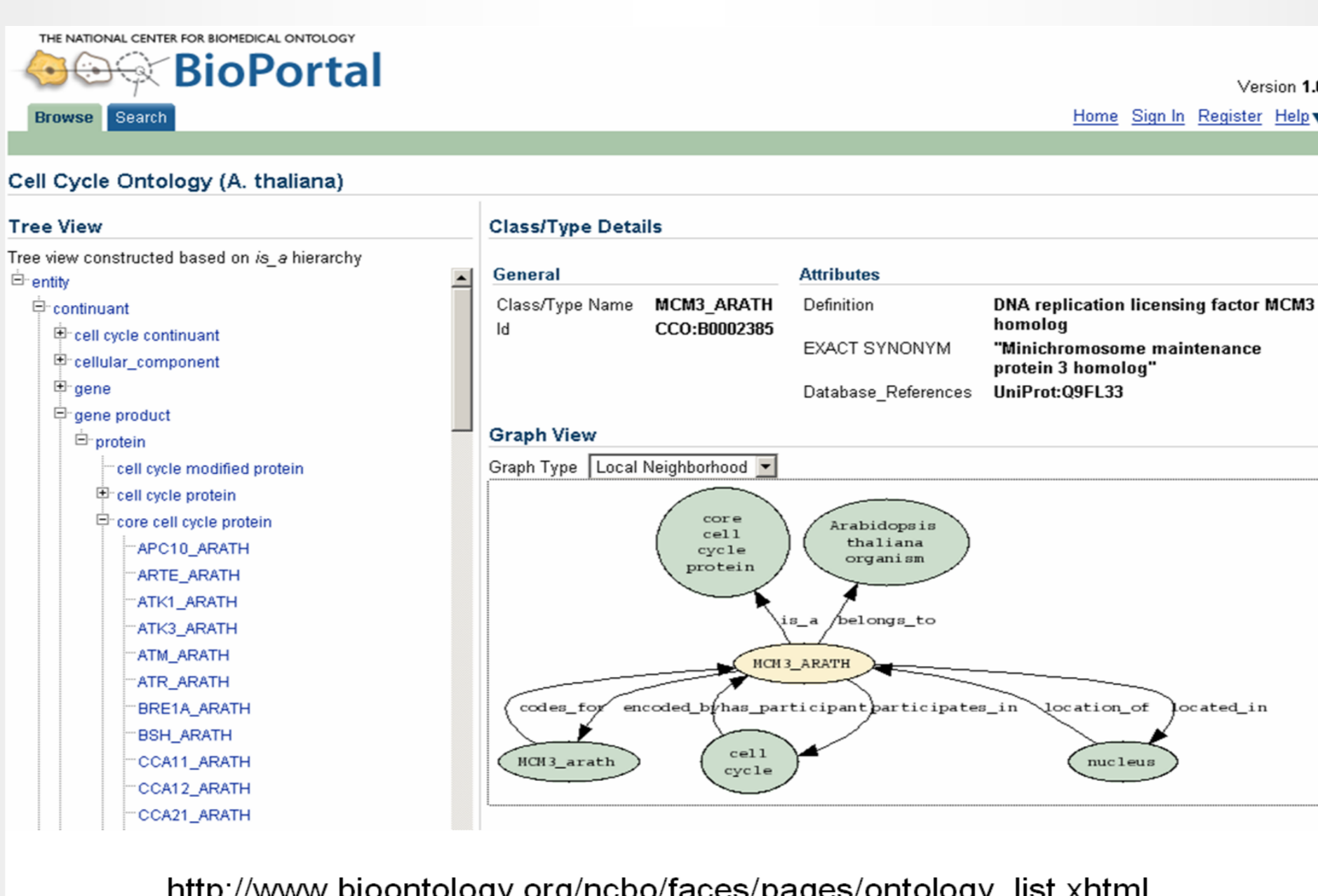


Fig 6. CCO in the BioPortal.

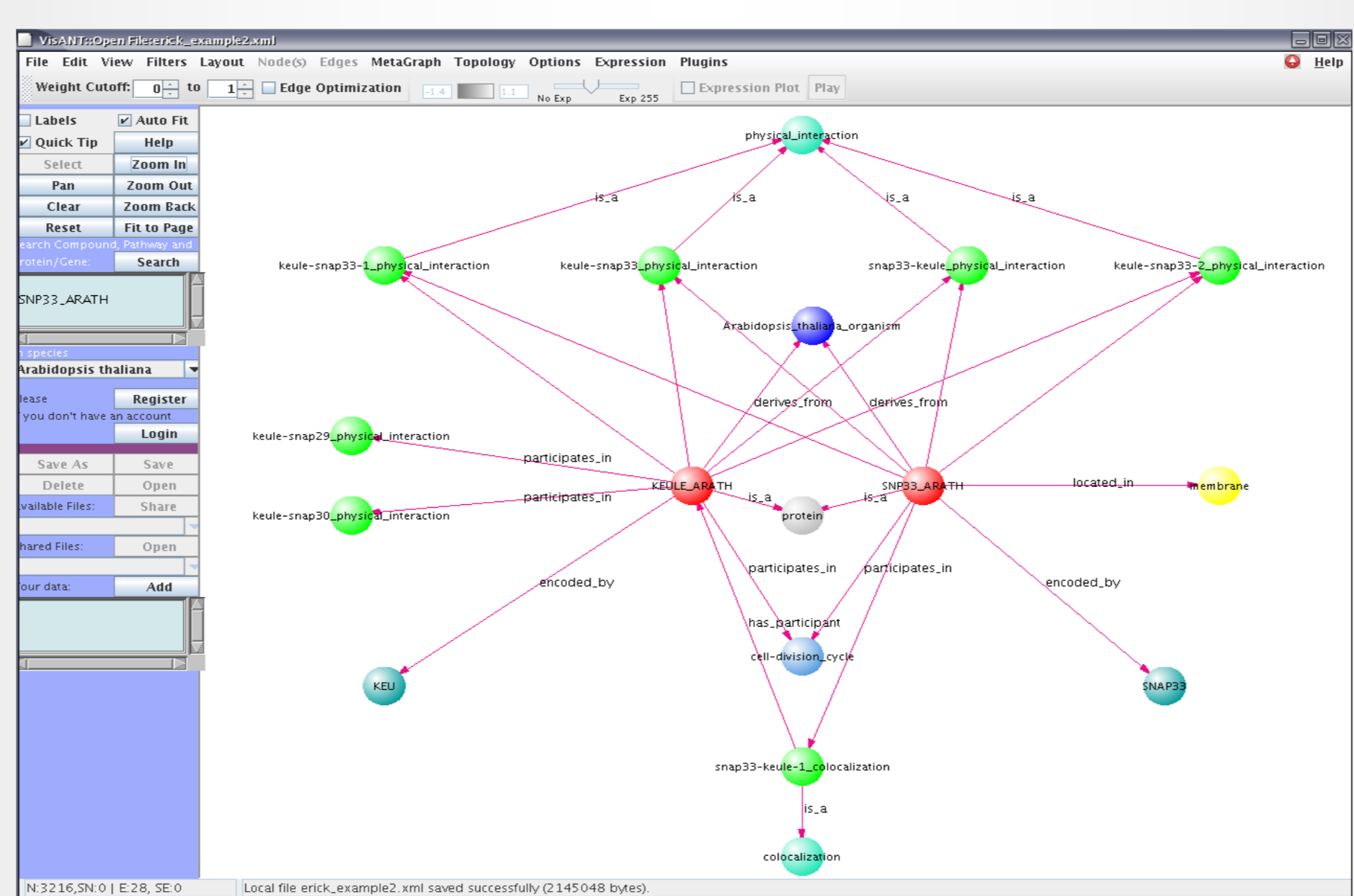
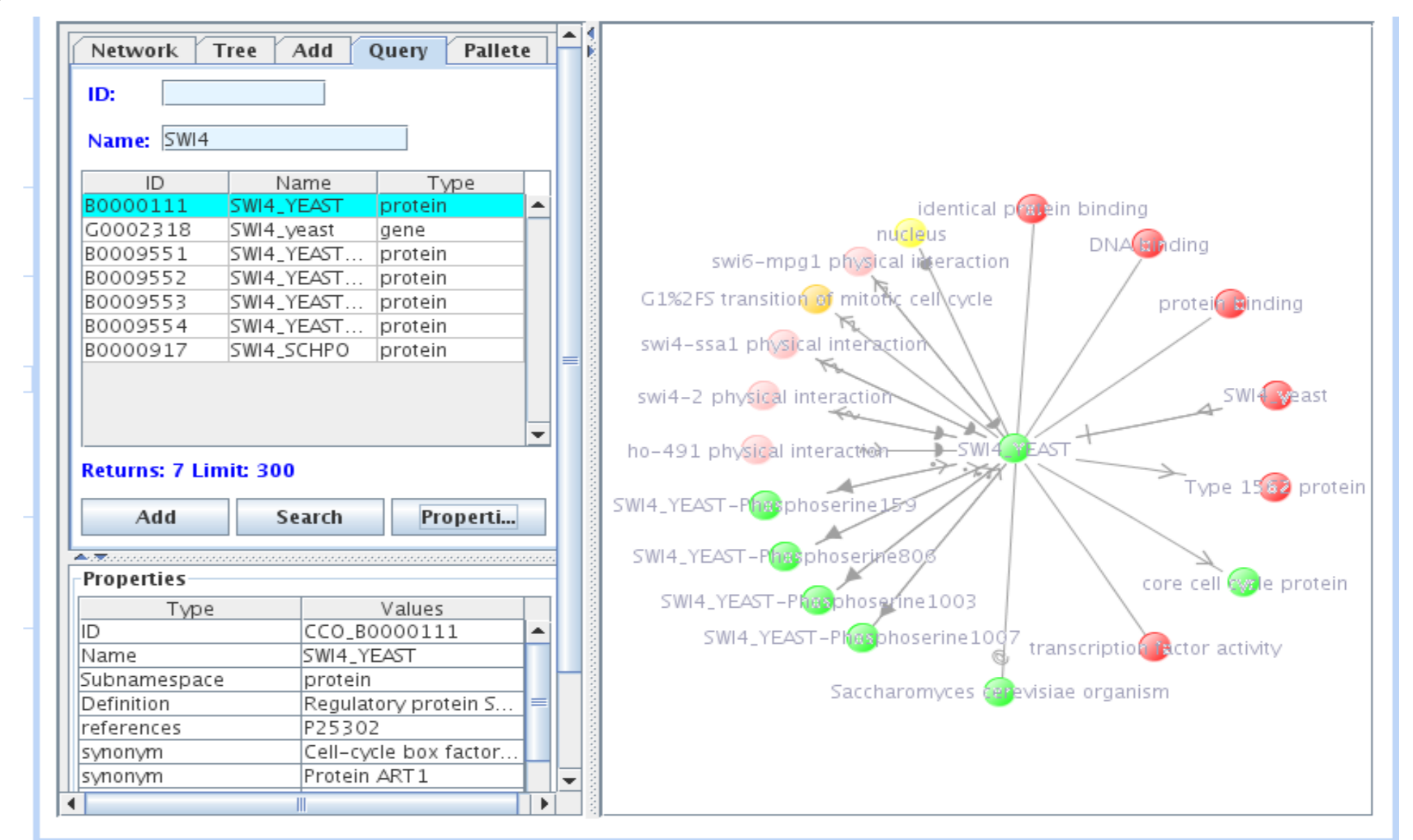


Fig 7. CCO in visANT.

## DIAMONDS Platform



- Within the EU FP6 project DIAMONDS (LSHG-CT-2004-512143) one of the objectives was to build a data integration platform dedicated to cell cycle biology.
- The Cell Cycle Ontology was chosen as data integration paradigm.
- NorayBio, a software company in Bilbao, Spain (people involved: Iñaki Bilbao, Aitzol Illarramendi, Marta Acilu), developed a software platform that allows querying and visualization of the CCO knowledge, and it provides links to Expression Profiler and ArrayExpress data. (<http://www.semantic-systems-biology>)

## Reasoning results

- There are a number of relationships in GO (core source of CCO) that might have been better annotated as *part\_of* instead of *is\_a*.
- The results inspired the GO team to make some amendments to the process part of the GO (e.g. regulation of cell cycle).
- Inconsistencies found in the data about the cellular localizations and protein-protein interactions.

## Conclusions and Results

- A fully automated **data integration pipeline** (nightly launched) was developed (Fig 2).
- Concrete problems and results related to the implementation of automatic format mappings (OWL, XML, DOT, GML) between ontologies and inconsistency checking issues have been identified.
- Several exports in commonly used formats have been developed (Fig 3-7).
- Existing integration obstacles due to the diversity of data formats and lack of formalization approaches as well as the trade-offs that are common in biological sciences.

## Future work

- **Knowledge** will be **weighted** (e.g. evidence codes) expressing the support media similar to those implemented in GO (experimental, electronically inferred, and so forth).
- **Ontolome analysis** (e.g. hypothesis generation by ontology alignments).
- **An advanced query system** will be developed (DL-based).
- **More data** to be integrated (upon feedback).

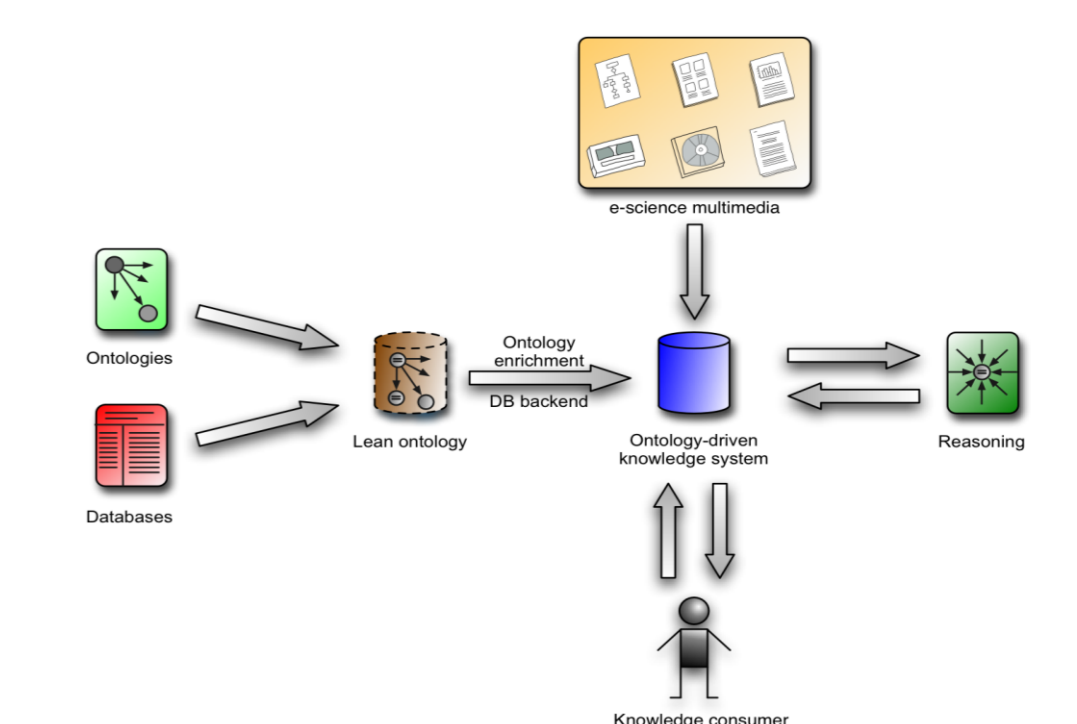


Fig 8. Outlook into the future

## Acknowledgements

This work has been funded by the EU (the DIAMONDS project, contract number LSHG-CT-2004-512143), The Manchester University, EPSRC, and Marie Curie EST.