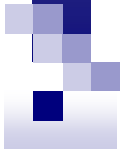


Slicing through the Scientific Literature

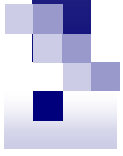
Christopher Baker (1), Patrick Lambrix (2), Jonas Laurila Bergman (2),
Rajaraman Kanagasabai (3), Wee Tiong Ang (3)

(1) University of New Brunswick, (2) Linköpings Universitet,
(3) ASTAR Singapore



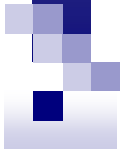
Literature search

- n Huge amount of scientific literature.
- n Need to integrate a spectrum of information to perform a task.



Literature search

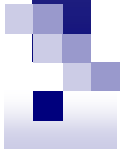
- n How to know what is in the repository
 - α Lack of knowledge of the domain
- n How to compose an expressive query
 - α Lack of knowledge of search technology



Example scenario

“Lipid”

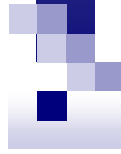
- n Keyword search returns all documents containing lipid.
 - α No knowledge; terminology problem
- n Relationships: use of multiple keywords with/without boolean operators, e.g. *lipid and disease*



Example scenario

“Lipid”

- n Keyword search returns a list of relevant questions concerning lipid. User selects question and retrieves knowledge and provenance documents.
- n Multiple search terms: requirement that there are relevant connections between the keywords.



The screenshot shows a Mozilla Firefox browser window with the title "KnowleFinder - Mozilla Firefox 3.0 Beta 5". The address bar contains the URL "http://localhost:8080/ESTKnowleFinder". The browser's menu bar includes "File", "Edit", "View", "History", "Bookmarks", "Tools", and "Help". The toolbar contains icons for "Smart Bookmarks", "Home", "Back", "Forward", "Stop", "Reload", "Print", "Email", "EST Live", and "demo (localhost)".

The main content area of the browser displays the "KnowleFinder" application. At the top, the word "KnowleFinder" is written in a large, bold, black serif font. Below the title is a search interface consisting of a text input field containing the word "lipid" and a "Search" button to its right. The search results area is a large, empty rectangular box. At the bottom of the browser window, the status bar shows the word "Done".

KnowleFinder - Mozilla Firefox 3.0 Beta 5

File Edit View History Bookmarks Tools Help

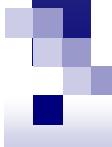
Smart Bookmarks email est EST Live demo (localhost)

Search

KnowleFinder

1. Which lipid has a broad synonym
2. Which lipid has a lipid KEGG_ID and has a broad synonym
3. Which lipid is implicated in a disease
4. Which lipid interacts with proteins
5. Which lipid is implicated in a disease and interacts with proteins
6. Which lipid is implicated in a disease and interacts with proteins involved in signal pathways
7. Which lipid is found in a sentence is implicated in a disease and interacts with proteins involved
8. Which document contains a sentence in which lipid is implicated in a disease and interacts with proteins involved

Done

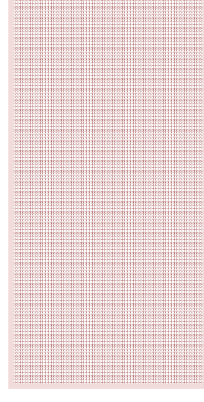


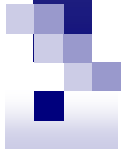
Question

NLG: Which lipid is implicated in a disease and interacts with proteins involved in signal pathways ?

Result

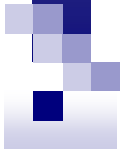
Protein	Lipid	Disease	Signal Pathway
P53	Unsat. Fatty Acid	Ovarian Cancer	Apoptosis





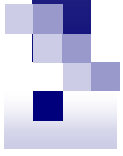
Example scenario

- n Requirements
 - ✕ Natural language interface
 - ✕ Ontology-driven query
 - ✕ Context of query terms
 - ✕ Cross-domain queries



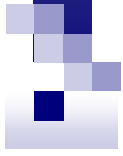
Outline

- n Relevant queries
- n Framework for slicing through the scientific literature
- n Algorithms and example
- n Conclusion & Future Work



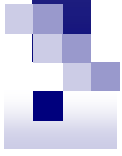
Outline

- n **Relevant queries**
- n Framework for slicing through the scientific literature
- n Algorithms and example
- n Conclusion & Future Work



Ontologies

“Ontologies define the basic terms and relations comprising the vocabulary of a topic area, as well as the rules for combining terms and relations to define extensions to the vocabulary.”

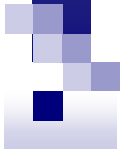


Ontologies

GENE ONTOLOGY (GO)

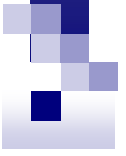
immune response
i- acute-phase response
i- anaphylaxis
i- antigen presentation
i- antigen processing
i- cellular defense response
i- cytokine metabolism
i- cytokine biosynthesis
synonym cytokine production
...
p- regulation of cytokine biosynthesis
...
...
i- B-cell activation
i- B-cell differentiation
i- B-cell proliferation
i- cellular defense response
...
i- T-cell activation
i- activation of natural killer cell activity
...

We will assume a graph representation.

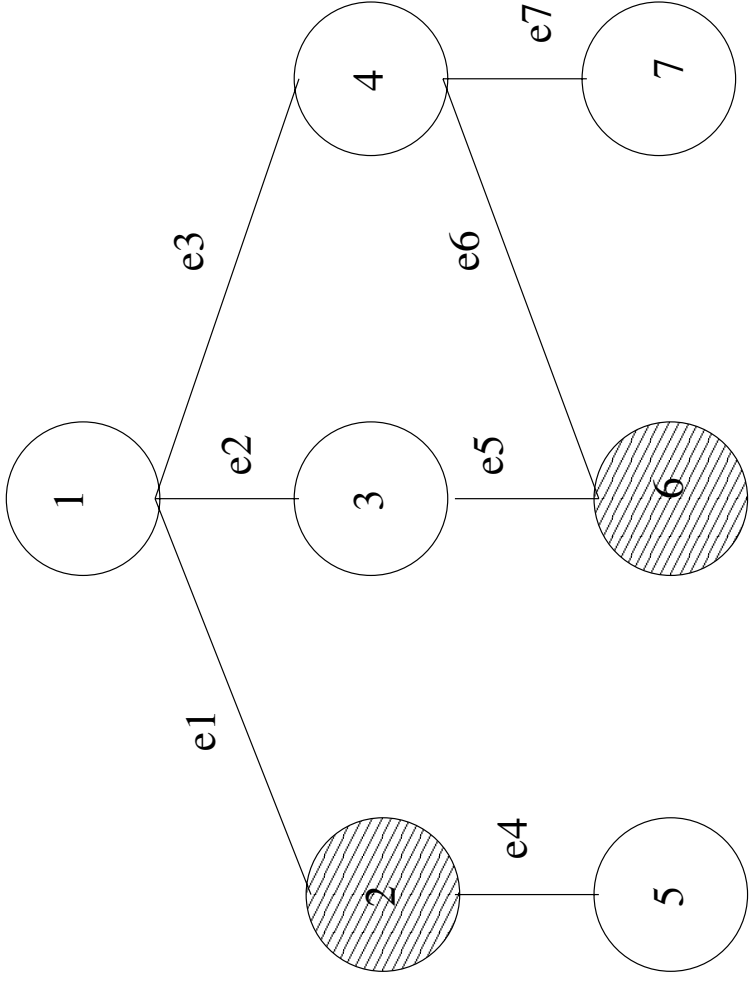


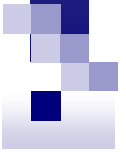
Relevant queries

- Relevant query including a number of concepts and relations from an ontology
- connected sub-graph of the ontology that includes the concepts and relations.
(query graph based on the concepts and relations; slice is set of all query graphs based on the concepts and relations)

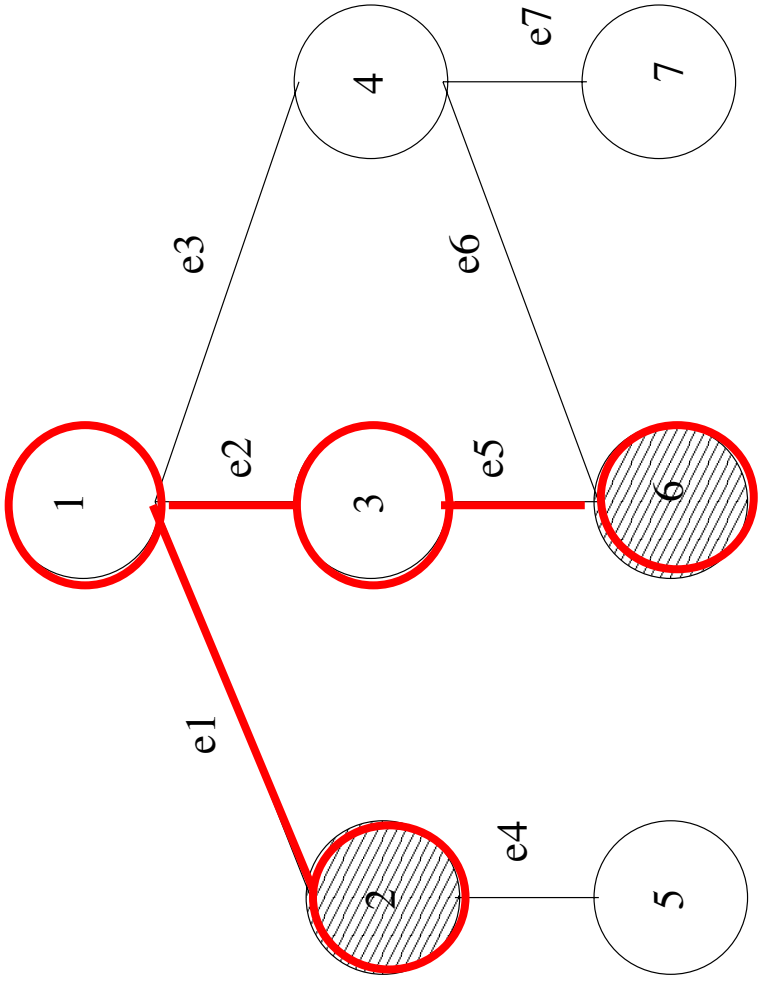


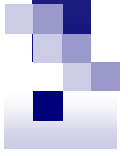
Query graph



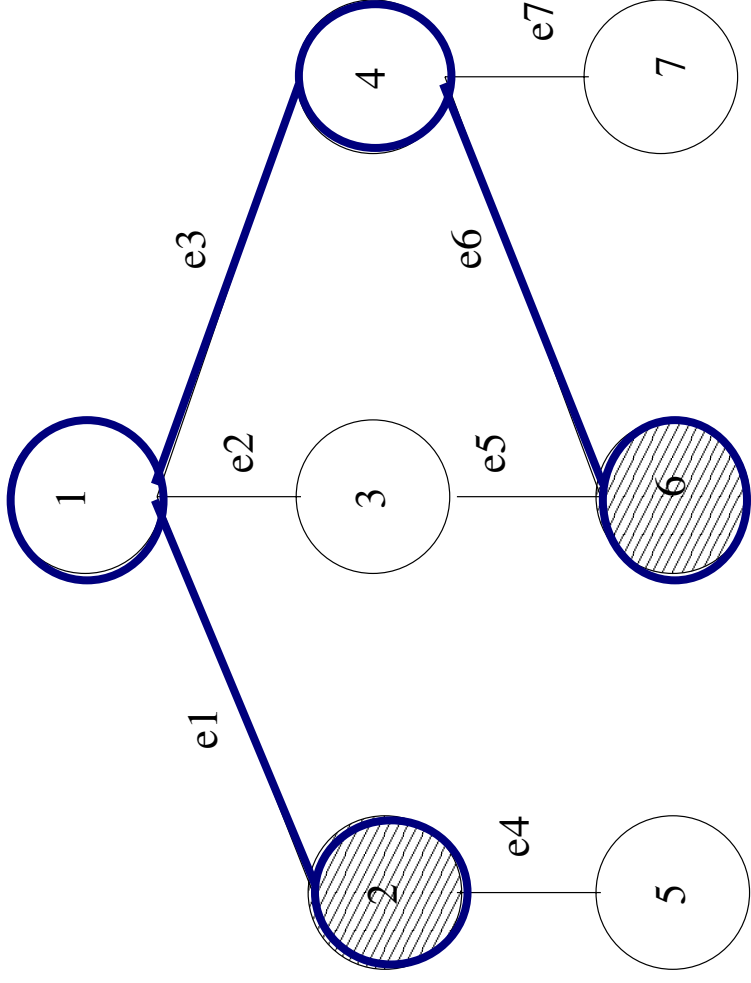


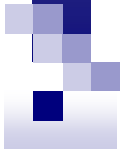
Query graph





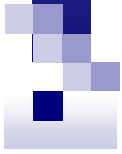
Query graph





Special cases

- n No relations, several concepts
 - ✗ Relevant queries regarding concepts; relations are suggested by the system.
 - ✗ Difference with traditional techniques: extra requirement that search terms need to be connected in the ontology.
- n No relations, one concept
 - ✗ Relevant queries including a specific query term.
 - ✗ Computes the ontological environment of the query term.



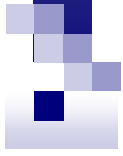
Multiple ontologies

GENE ONTOLOGY (GO)

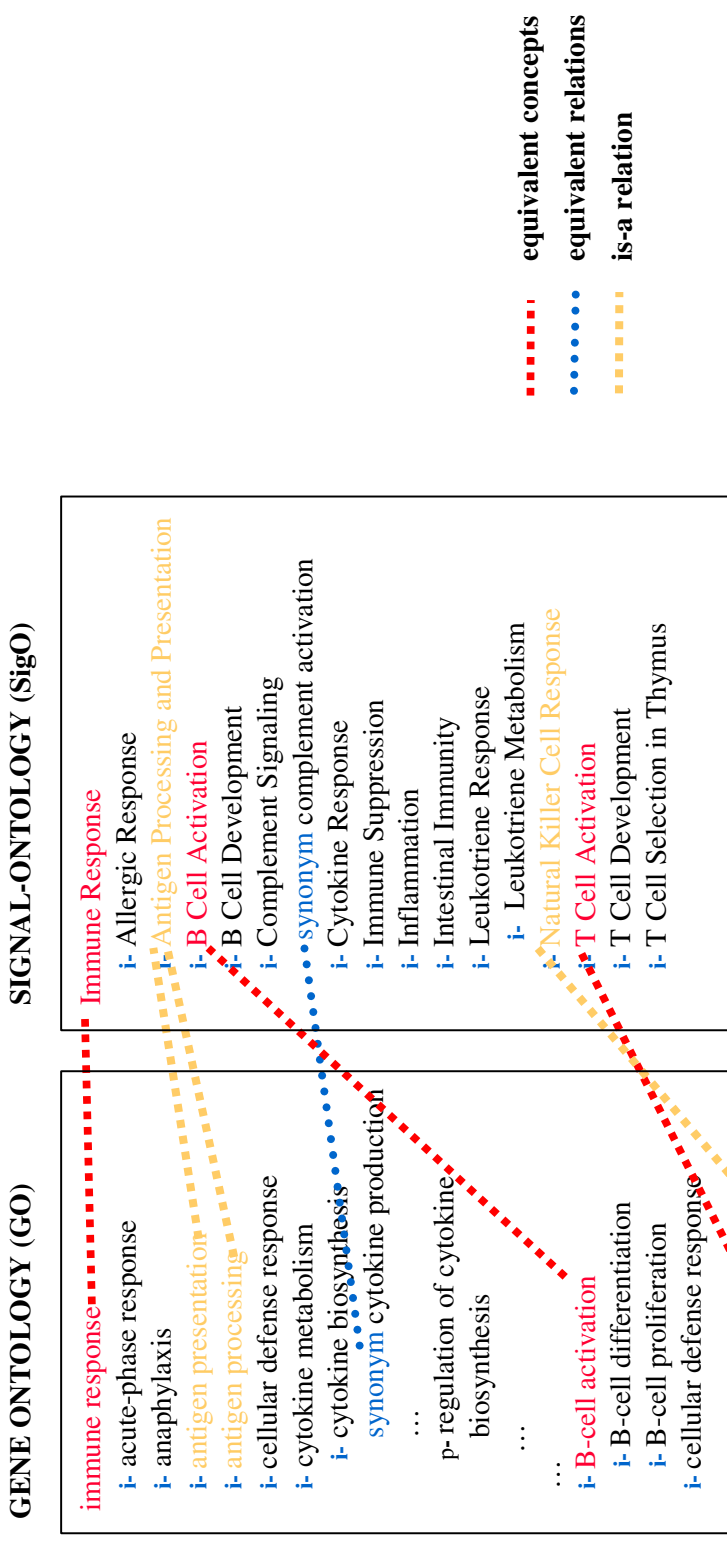
immune response
i- acute-phase response
i- anaphylaxis
i- antigen presentation
i- antigen processing
i- cellular defense response
i- cytokine metabolism
i- cytokine biosynthesis
synonym cytokine production
...
p- regulation of cytokine biosynthesis
...
...
i- B-cell activation
i- B-cell differentiation
i- B-cell proliferation
i- cellular defense response
...
i- T-cell activation
i- activation of natural killer cell activity
...

SIGNAL-ONTOLOGY (SigO)

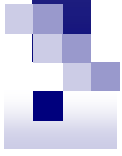
Immune Response
i- Allergic Response
i- Antigen Processing and Presentation
i- B Cell Activation
i- B Cell Development
i- Complement Signaling
synonym complement activation
i- Cytokine Response
i- Immune Suppression
i- Inflammation
i- Intestinal Immunity
i- Leukotriene Response
i- Leukotriene Metabolism
i- Natural Killer Cell Response
i- T Cell Activation
i- T Cell Development
i- T Cell Selection in Thymus



Ontology Alignment



Alignment is a set of mappings between terms in the ontologies.

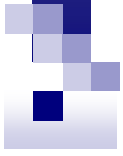


Relevant queries – multiple ontologies

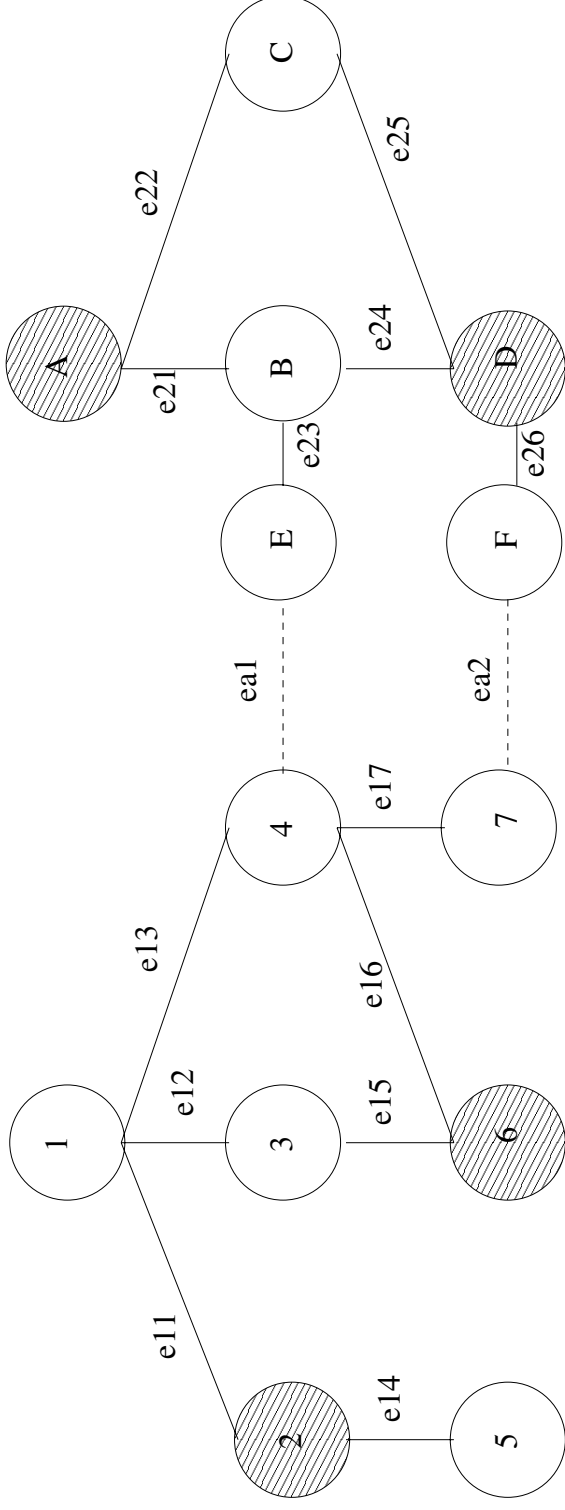
ⁿ Relevant query including a number of concepts and relations from multiple ontologies

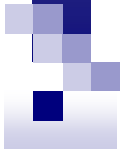
Query graphs connected by a path going through a mapping in the alignment.

*(aligned query graph based on query graphs;
aligned slice is set of all aligned query graphs
based on the query graphs)*

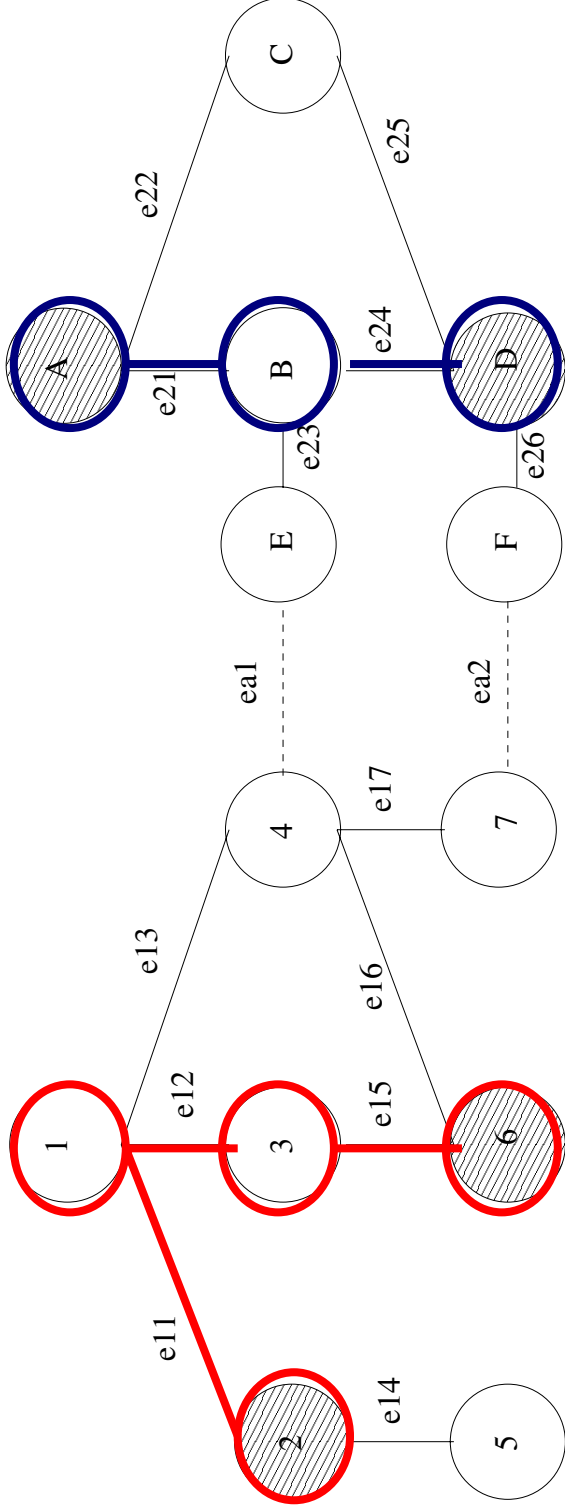


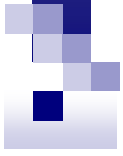
Aligned query graph



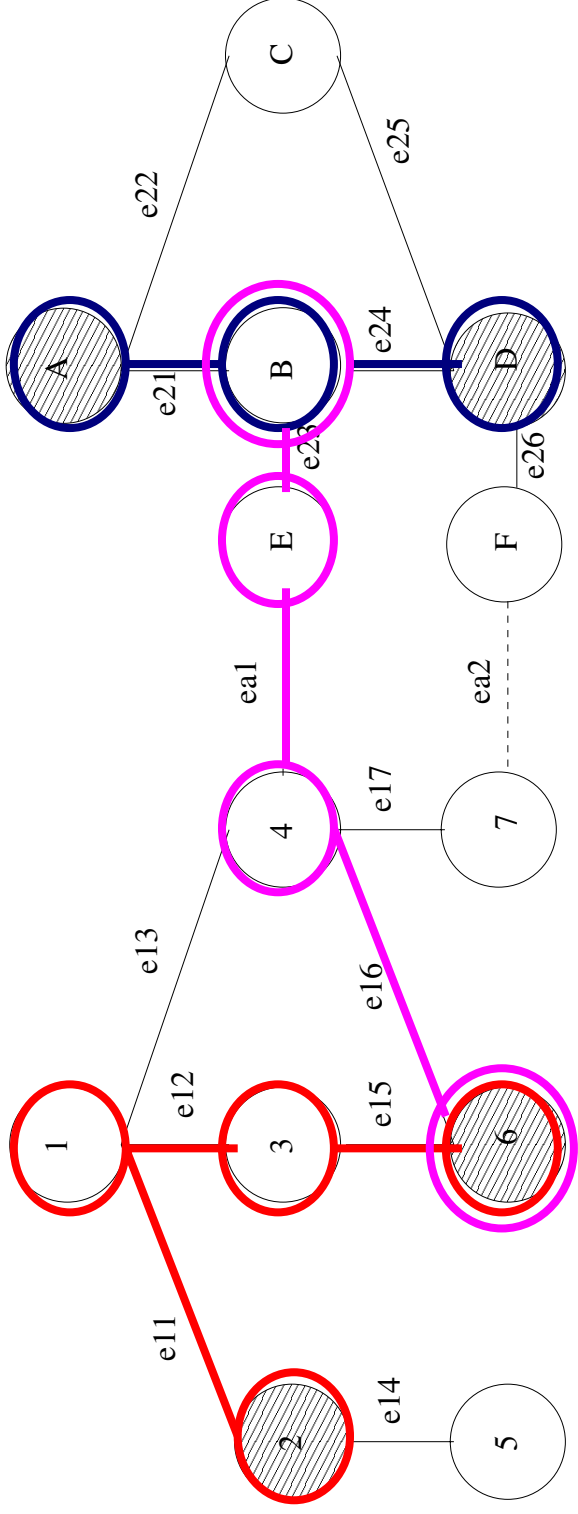


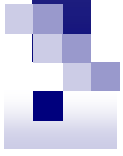
Aligned query graph





Aligned query graph

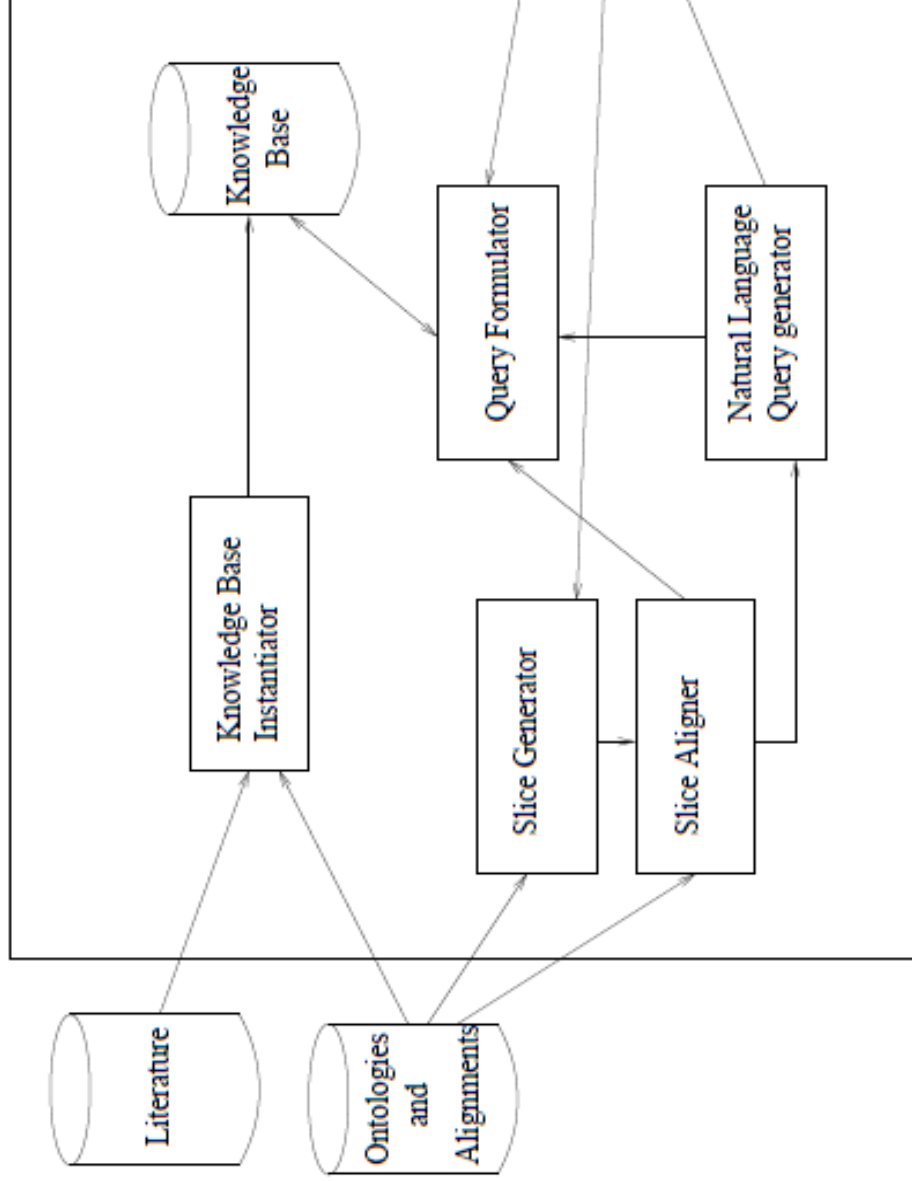


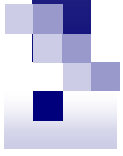


Outline

- n Relevant queries
- n **Framework for slicing through the scientific literature**
- n Algorithms and example
- n Conclusion & Future Work

Framework

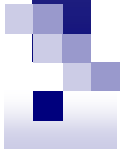




Framework

- n External resources
 - ✕ Literature document base
 - ✕ Ontology and ontology alignment repository

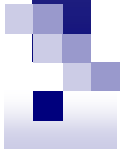
- n Computed resources
 - ✕ Knowledge base



Framework

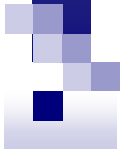
- n Process and translation
 - ✘ Knowledge base instantiation
 - ✘ Slice generation and alignment
 - ✘ Natural language query generation

n Query



Outline

- n Relevant queries
- n Framework for slicing through the scientific literature
- n **Algorithms and example**
- n Conclusion & Future Work



External resources

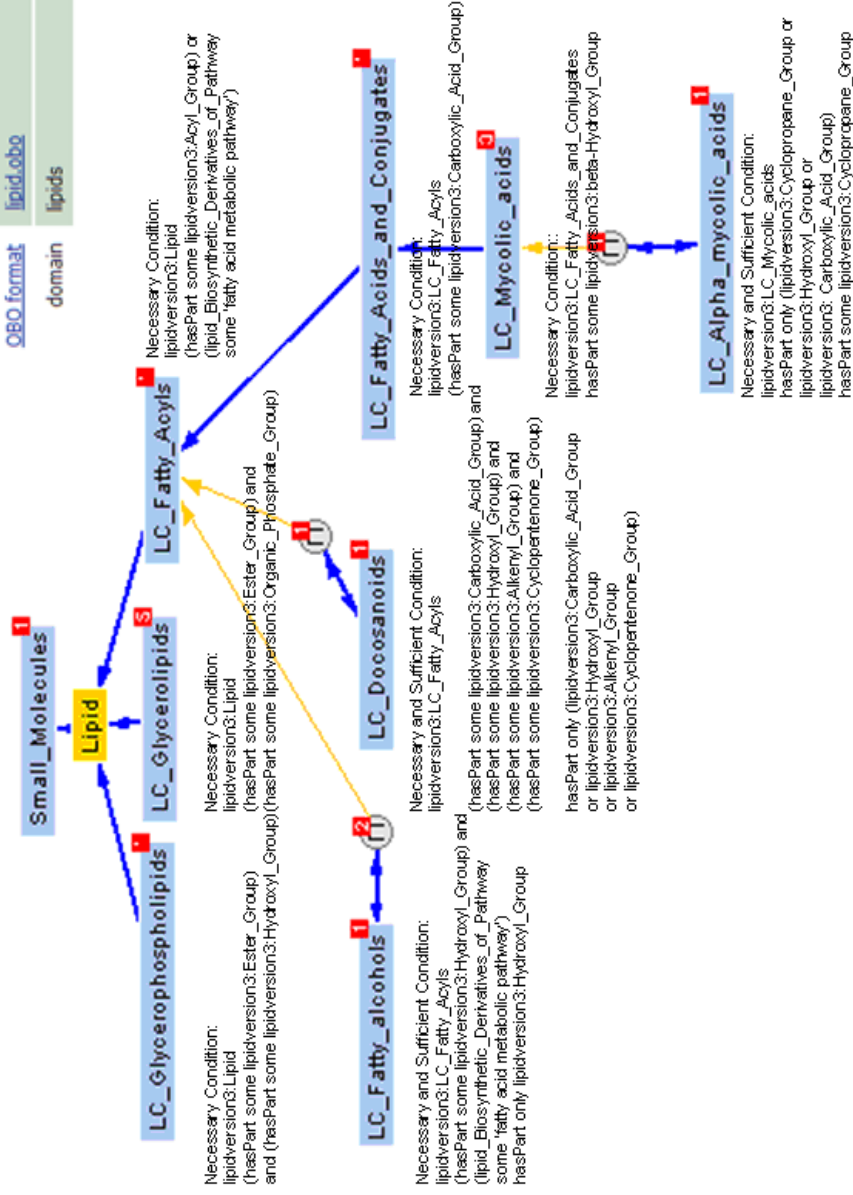
- n Literature document base
 - ✕ Generated from a collection of 7498 PubMed abstracts relevant for Ovarian Cancer. 683 papers included lipid names from which 241 full papers were downloadable.
- n Ontology and ontology alignment repository
 - ✕ Lipid ontology
 - ✕ Signal ontology
 - ✕ Alignment using SAMBO

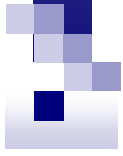


Lipid Ontology

Lipid research is increasingly integrated within systems level biology such as lipidomics where lipid classification is required before appropriate annotation of chemical functions can be applied. The ontology describes the LIPIDMAPS nomenclature classification explicitly using description logics (OWL-DL). Lipid classes are organized hierarchically with the super-classes restricted by generic necessary conditions. More specific necessary conditions are used to define membership requirements for sub classes of lipid according to appropriate functional groups.

namespace	LIPRO
current activity	Active
contact	Christopher Baker
OWL format	LIPRO.owl
OBO format	lipid.obo
domain	lipids





SAMBO (1)

n SAMBO (System for Aligning and Merging Biomedical Ontologies)

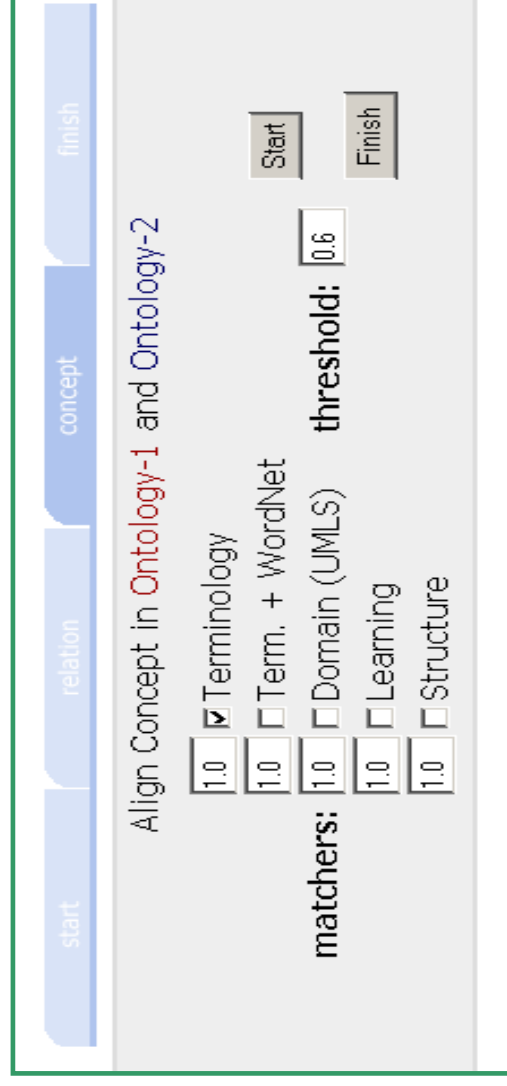
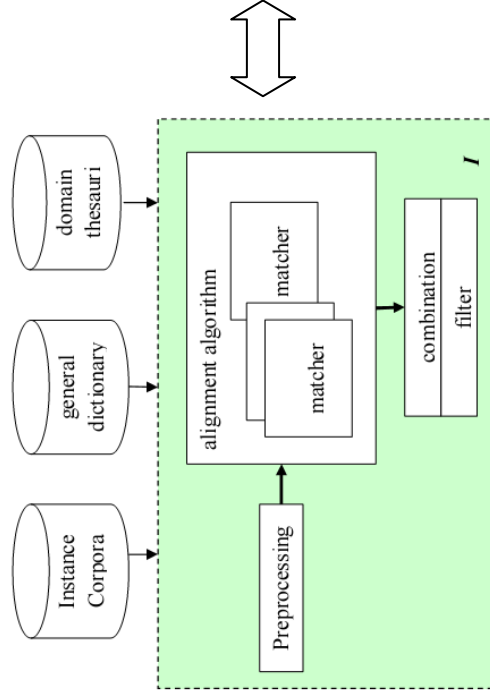
α Phase I

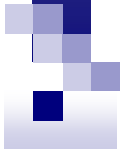
Winner Anatomy Track of OAEI 2008

n Matchers

n Weighted sum combination of matcher results

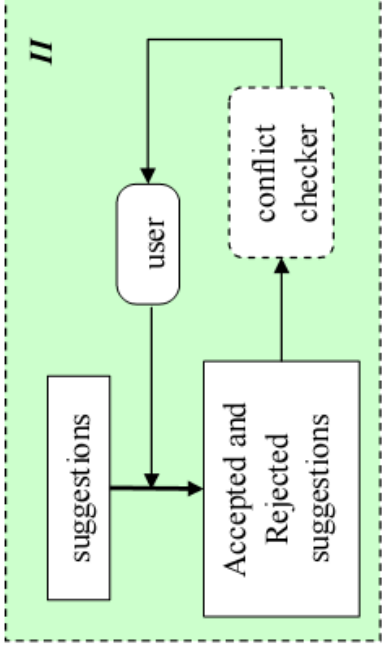
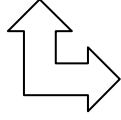
n Single threshold filtering





SAMBO (2)

Phase II:

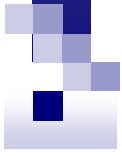


nose_MA	nose_MESH
nasal cavity epithelium	nasal_mucosa
Id: MA_0001324	Id: MESH_A.04.531.520
definition: nasal mucosa	definition: nasal_epithelium
Synonym: nasal cavity	Synonym: nasal_mucosa
Part of: nasal cavity	Part of:

comment on the alignment

new name for the alignment

warning



Knowledge base instantiation

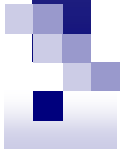
BioText toolkit

- n Named entity recognition
 - ✕ Gazetteer approach
 - ✕ Termlists
 - n Lipids: LIPIDMAPS, LipidBank, KEGG, IUPAC
 - n Proteins: Swiss-Prot
 - n Diseases: disease ontology of Center of Genetic Medicine
 - ✕ Tags found entities with ontology concepts

n Normalization and grounding

n Relation detection

- ✕ Based on co-occurrence in sentence
 - + rule set developed with domain expert



Knowledge base instantiation

- 1) Document Content
- 2) Sentence Extraction
- 3) Sentence Detection: **lipid interaction protein**
- 4) Entity Recognition: term identification / assign **lipid** class
- 5) Normalization: collapse **lipid** synonyms
- 6) Relation Extraction: **Lipid-Protein** or **Lipid Disease**
"TLR4 binds to POPC", tagged as
"<term category=**protein**> TLR4</term>
binds to
<term category=**lipid**>POPC</term>"
- 7) Classification: Identify ontology classes and specify relations for all sentences, proteins, **lipid** subclasses.
- 8) Populate OWL ontology (JENA -API)

Term List DB's:

Lipid names,
LIPIDMAPS, Lipid Bank,
KEGG classifications,
Disease names,
Protein names
Stemmed Interactions

Document and
sentence meta data

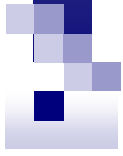
Complete
Instantiated
OWL-DL
Ontology



Mined Interactions

Table 1. Interactions mined from the ovarian cancer bibliome. OC and AP represent a cancer and apoptosis pathway proteins respectively.

Interaction Type	Abstract (7498)	Full Paper (241)
OC-OC	223	13
OC-AP	505	195
OC-Lipid	11	14
OC-Hormone	8	1
AP-AP	113	59
AP-Lipid	10	8
Lipid-Lipid	3	23
Lipid Hormone	2	18
Protein Hormone	9	2
Hormone-Hormone	2	6



Knowledge base instantiation

n Population

⌘ Concepts

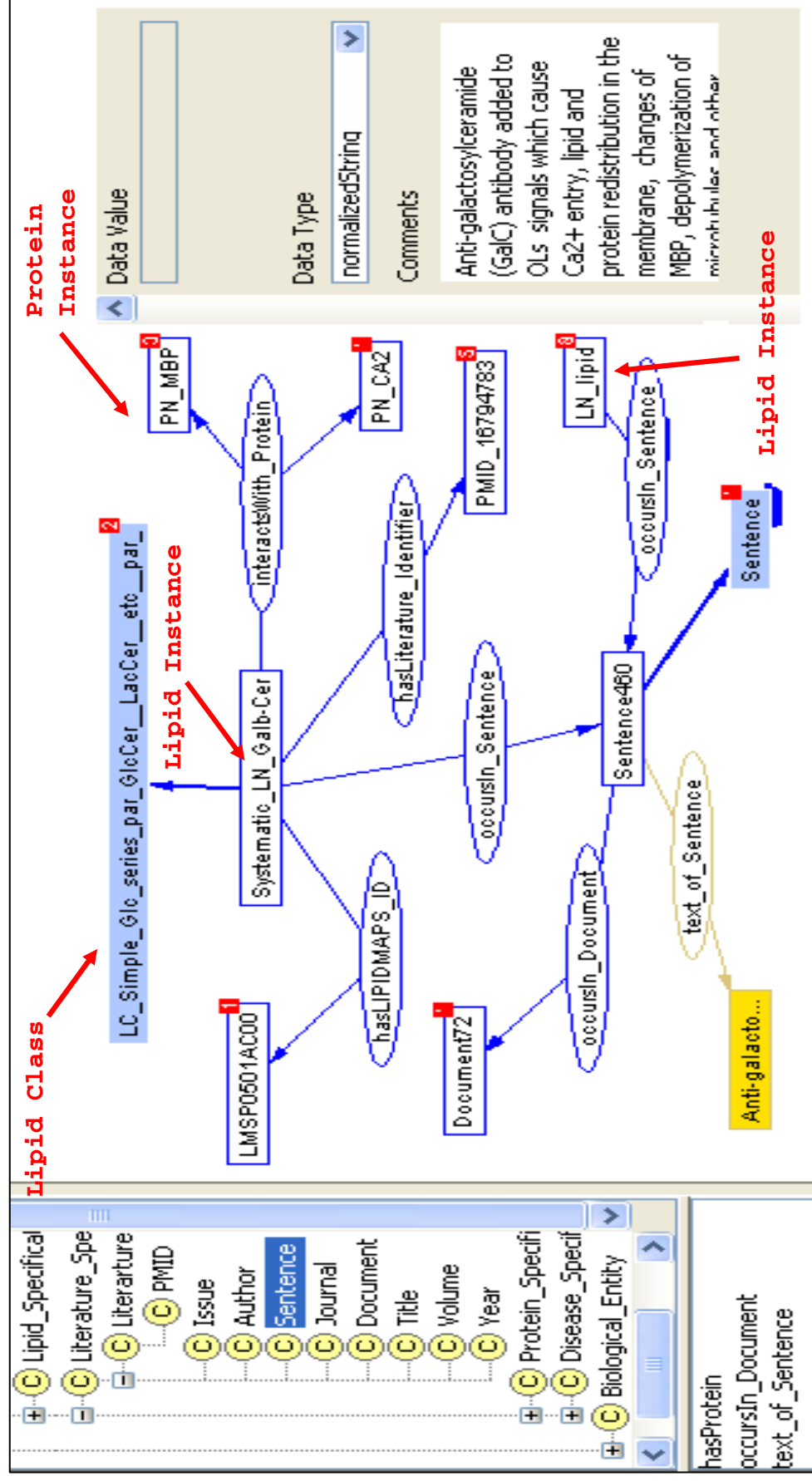
⌘ Properties

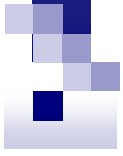
Knowledge

Provenance documents



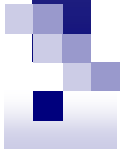
Knowledge base instantiation





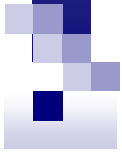
Slice generation

- n Current implementation focuses on slices based on concepts.
- n Depth-first traversal of ontology to find paths between given concepts; paths can be put together to find slices/query graphs.



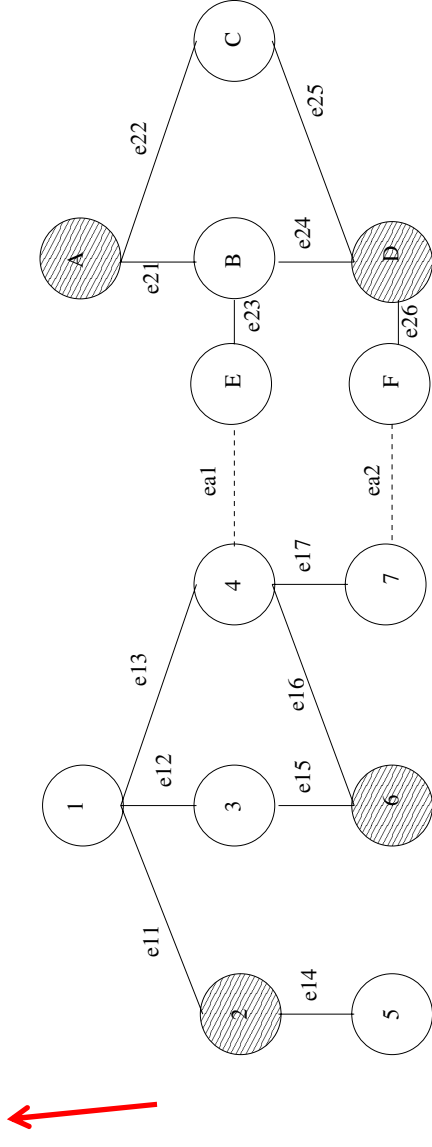
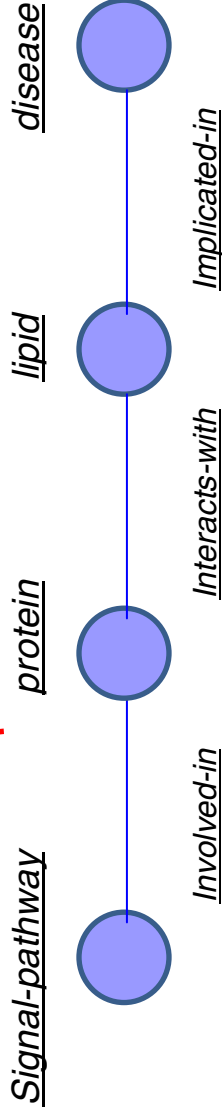
Slice alignment

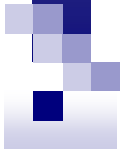
- n Algorithm computes subset of aligned slice.
- n Assumption: shorter paths represent closer relationships.
- n Algorithm connects slices using shortest paths from given concepts in one ontology to given concepts in other ontology.



Slicing through the literature

```
nROL: (RETRIEVE (?X ?Y ?Z ?W)  
(AND (?X Protein) (?Y Lipid) (?Z Disease) (?W SignalPathway)  
(?X ?Y Interacts_with) (?Y ?Z Implicated_in) (?X ?W Involved_in)))
```





Natural language query generation

- n Triple representation:
 - <*lipid, interacts-with, protein*>
- n Rule base to generate NL statements.
 - What *lipid interacts with proteins*?
 - ⌘ Learned from examples.
- n Aggregation of statements from different triples, grammar checking.



File Edit View History Bookmarks Tools Help

Smart Bookmarks email est EST Live demo (localhost)

KnowleFinder

1. Which lipid has a broad synonym
2. Which lipid has a lipid KEGG ID and has a broad synonym
3. Which lipid is implicated in a disease
4. Which lipid interacts with proteins
5. Which lipid is implicated in a disease and interacts with proteins
6. Which lipid is implicated in a disease and interacts with proteins involved in signal pathways
7. Which lipid is found in a sentence is implicated in a disease and interacts with proteins involved
8. Which document contains a sentence in which lipid is implicated in a disease and interacts with proteins involved

Done



Query

n Send nRQL query to RACER.

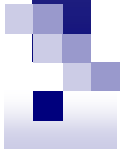
Question

NLG: Which lipid is implicated in a disease and interacts with proteins involved in signal pathways ?

```
nRQL: (RETRIEVE (?X ?Y ?Z ?W)
(AND (?X Protein) (?Y Lipid) (?Z Disease) (?W signalPathway)
(?X ?Y Interacts_with) (?Y ?Z Implicated_in) (?X ?W Involved_in)))
```

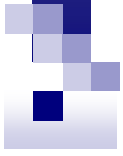
Result

Protein	Lipid	Disease	Signal Pathway
P53	Unsat. Fatty Acid	Ovarian Cancer	Apoptosis



Outline

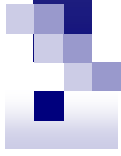
- n Relevant queries
- n Framework for slicing through the scientific literature
- n Algorithms and example
- n **Conclusion & Future Work**



Conclusions

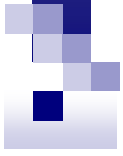
- n Framework for literature search dealing with
 - ⌘ Lack of knowledge of the domain.
 - ⌘ Lack of knowledge of the search technology.

- n Proof of concept implementation.



Future Work

- n Heuristics and their influences for slice generation.
- n Tradeoff in query generation between completeness and information overload.
- n Relevance measure and query ranking.
- n Optimization of slice generation.



Future Work

- n Integration in larger system
- n Integrated implementation.
- n Scalability testing.