

YeastMed: An XML-Based System for Biological Data Integration of Yeast

Abdelaali Briache, Kamar Marrakchi, Amine Kerzazi, Ismael Navas-Delgado, Jose F Aldana Montes, Badr D. Rossi Hassani and Khalid Lairini

LABIPHABE, Department of Biology, F. S. T. Of Tangier , Morocco.

KHAOS, Departamento de Lenguajes y Ciencias de la Computación, Universidad de Málaga, Spain.

Index

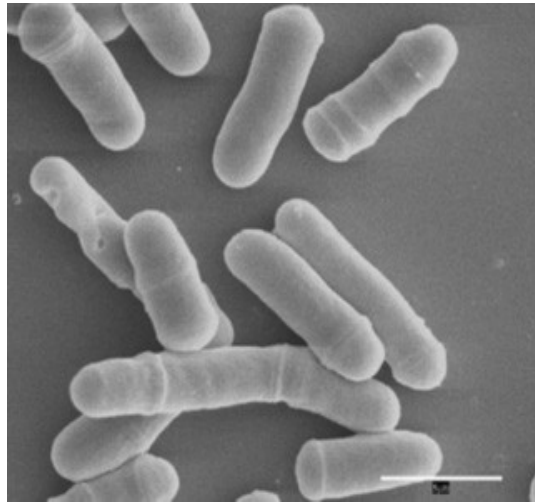
- Introduction.
 - Yeasts.
 - Data integration.
- YeastMed System.
 - Architecture.
 - Data sources.
 - Schemas.
 - Ontology.
 - Mappings.
 - SB-KOM.
- Conclusions.

Introduction: Yeasts

- ▶ Tiny forms of fungi;
- ▶ Microorganisms visible only with a microscope;
- ▶ Cell cycle similar to that of Humans,
- ▶ rapid growth,
- ▶ dispersed cells,
- ▶ a well-defined genetic system,
- ▶ genome can be easily manipulated,
- ▶



filament-shaped *Candida albicans*



Elongated
schizosaccharomyces pombe



Egg-shaped
Saccharomyces cerevisiae

Introduction: Yeasts

Yeast-Specialized sources

AGD	<i>Ashbya gossypii</i> genome database
CandidaDB	<i>Candida albicans</i> genome database
Candida Genome	<i>Candida albicans</i> genome database
CYGD	MIPS Comprehensive yeast genome database
Génolevures	A comparison of <i>S.cerevisiae</i> and 14 other yeast species
PROPHECY	Profiling of phenotypic characteristics in yeast
SCMD	<i>Saccharomyces cerevisiae</i> morphological database: micrographs of budding yeast mutants
SCPD	<i>Saccharomyces cerevisiae</i> promoter database
SGD	<i>Saccharomyces</i> genome database
TRIPLES	Transposon-insertion phenotypes, localization and expression in <i>Saccharomyces</i>
YDPM	Yeast deletion project and mitochondria database
Yeast Intron Database	Ares lab database of splicing introns in <i>S.cerevisiae</i>
Yeast snoRNA Database	Yeast small nucleolar RNAs
yMGV	Yeast microarray global viewer
YRC PDR	Yeast resource center public data repository

General Sources

NCBI Protein database	All protein sequences: translated from GenBank and imported from other protein databases
Swiss-Prot	Now UniProt/Swiss-Prot: expertly curated protein sequence database, section of the UniProt knowledgebase
UniProt	Universal protein knowledgebase: merged data from Swiss-Prot, TrEMBL and PIR protein sequence databases
PROSITE	Biologically significant protein patterns and profiles
PRINTS	Hierarchical gene family fingerprints
Pfam	Protein families: multiple sequence alignments and profile hidden Markov
DDBJ—DNA Data Bank of Japan	All known nucleotide and protein sequences
EMBL Nucleotide Sequence Database	All known nucleotide and protein sequences

Introduction: Data Integration

Different data models and schemas,

Different model constructs that can be used to describe the same object, even if the same model is used,

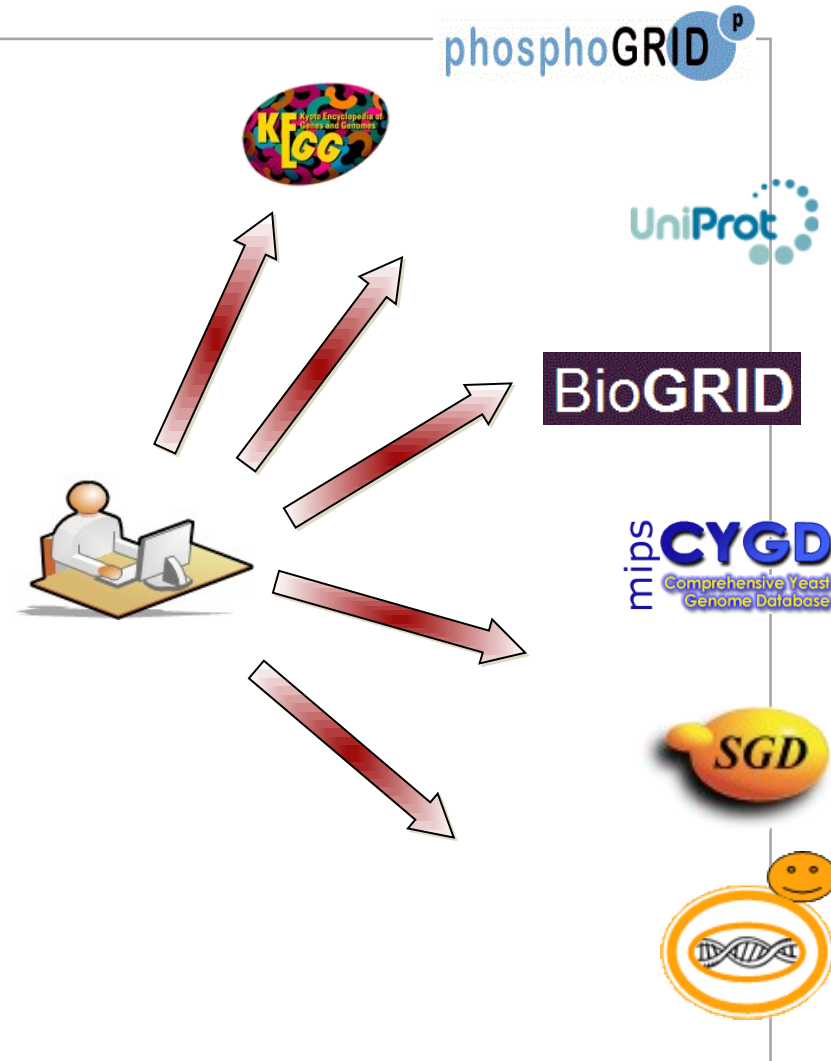
wide variety of formats that have been used for data representation : ASN.1 (*Abstract Syntax Notation One*), XML, HTML...etc.,

Data sources make their data available in different ways,

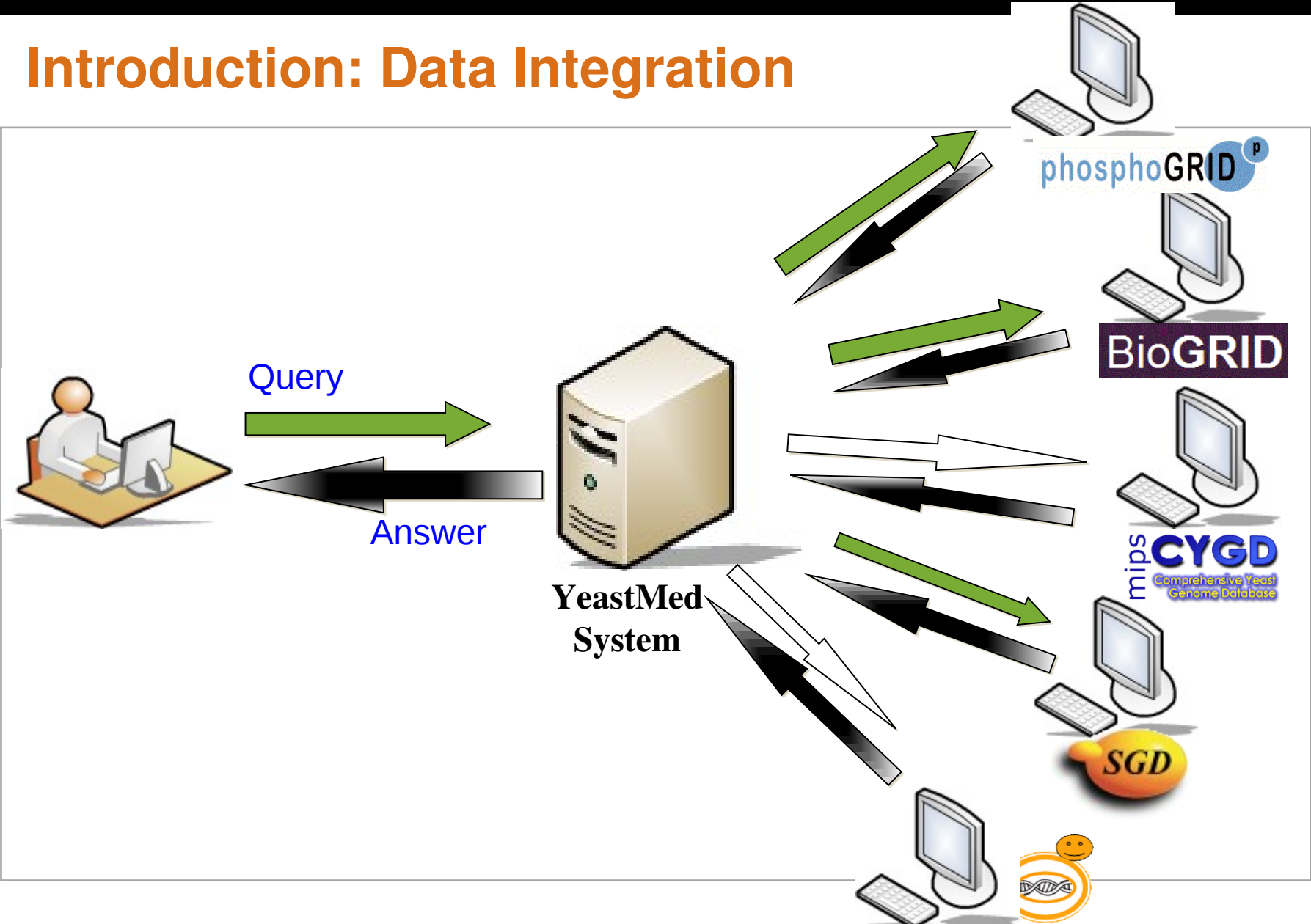
The query language used to interrogate data sources,

inconsistent use of nomenclature,

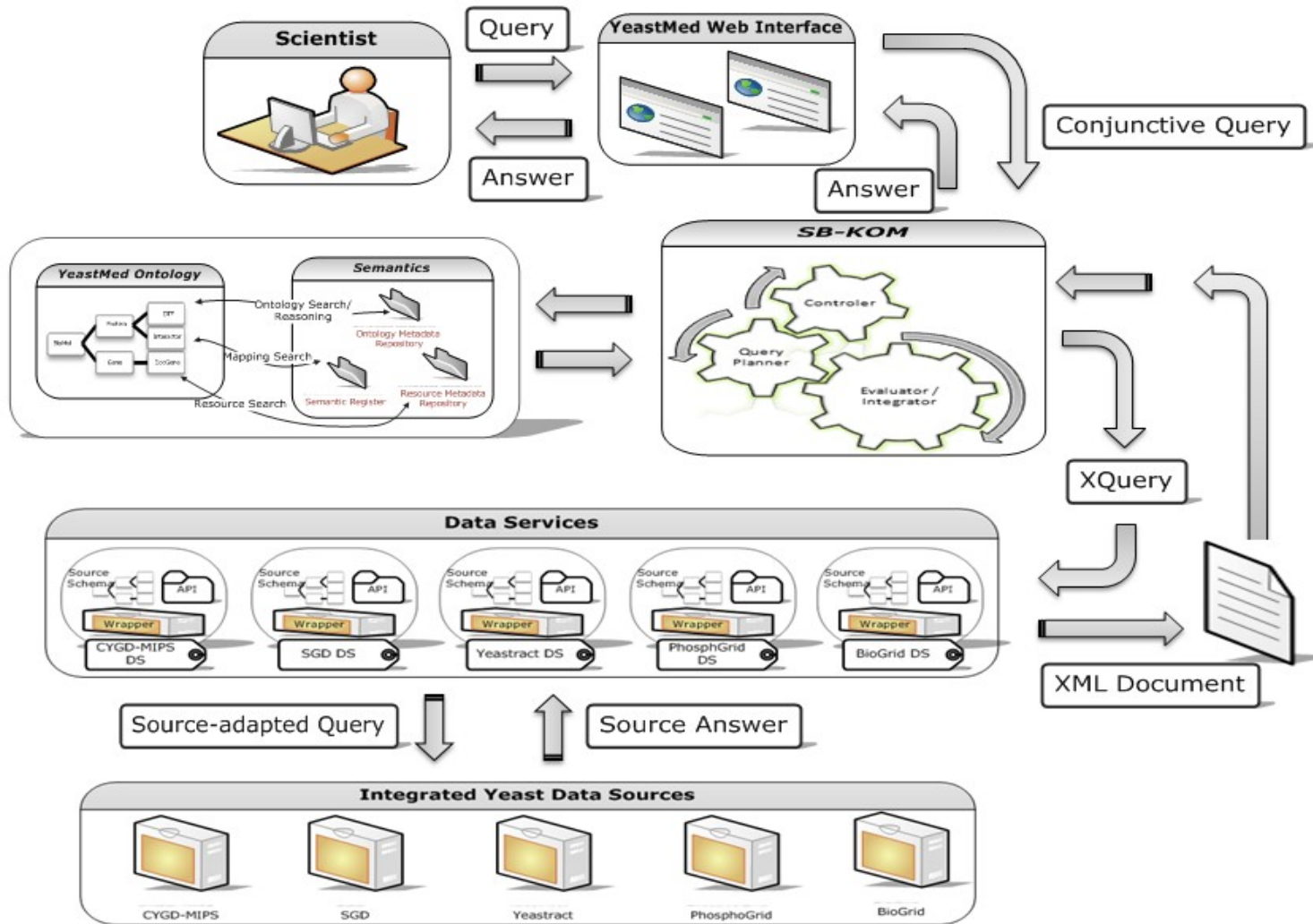
...



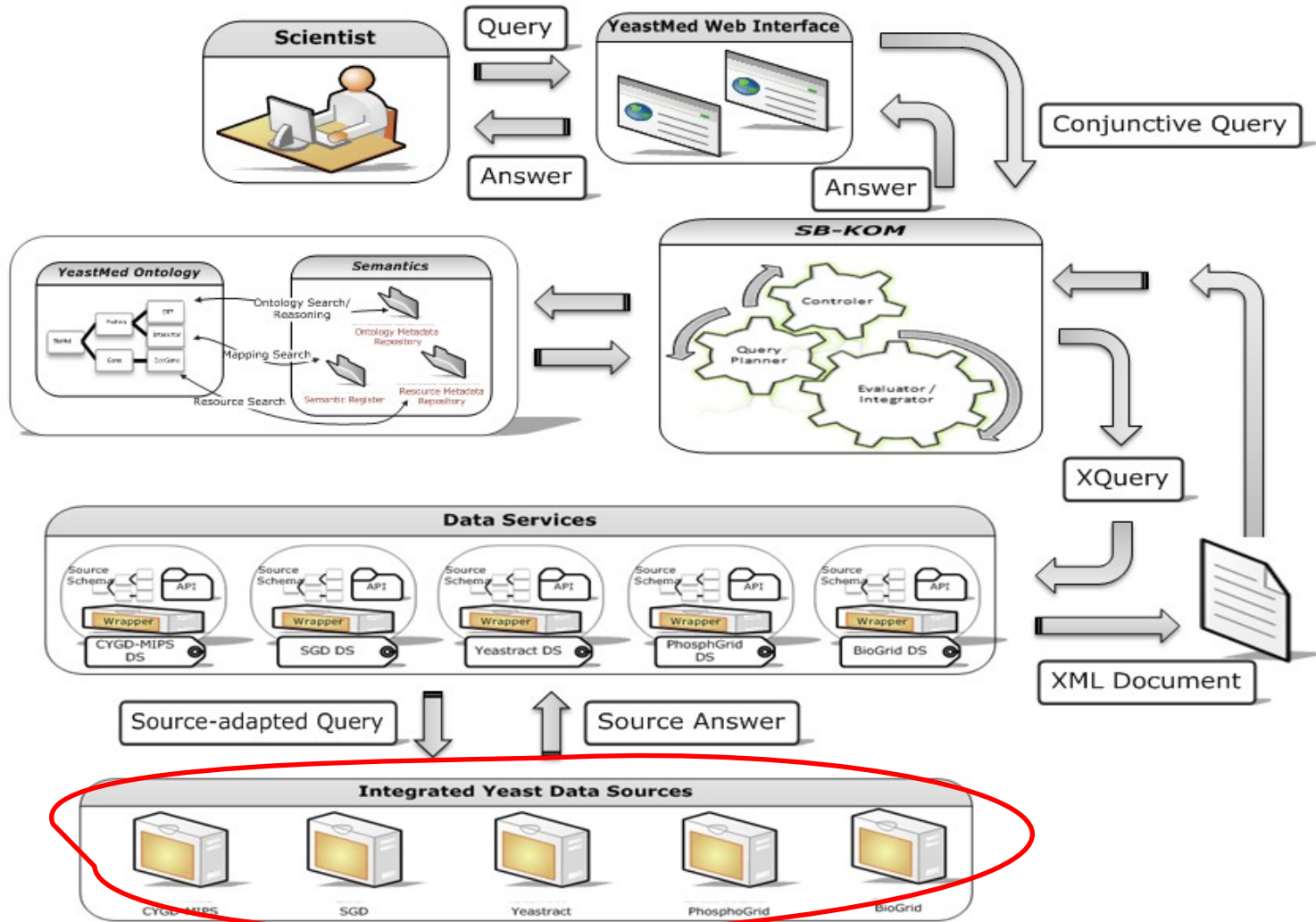
Introduction: Data Integration



YeastMed : Architecture



YeastMed : Architecture

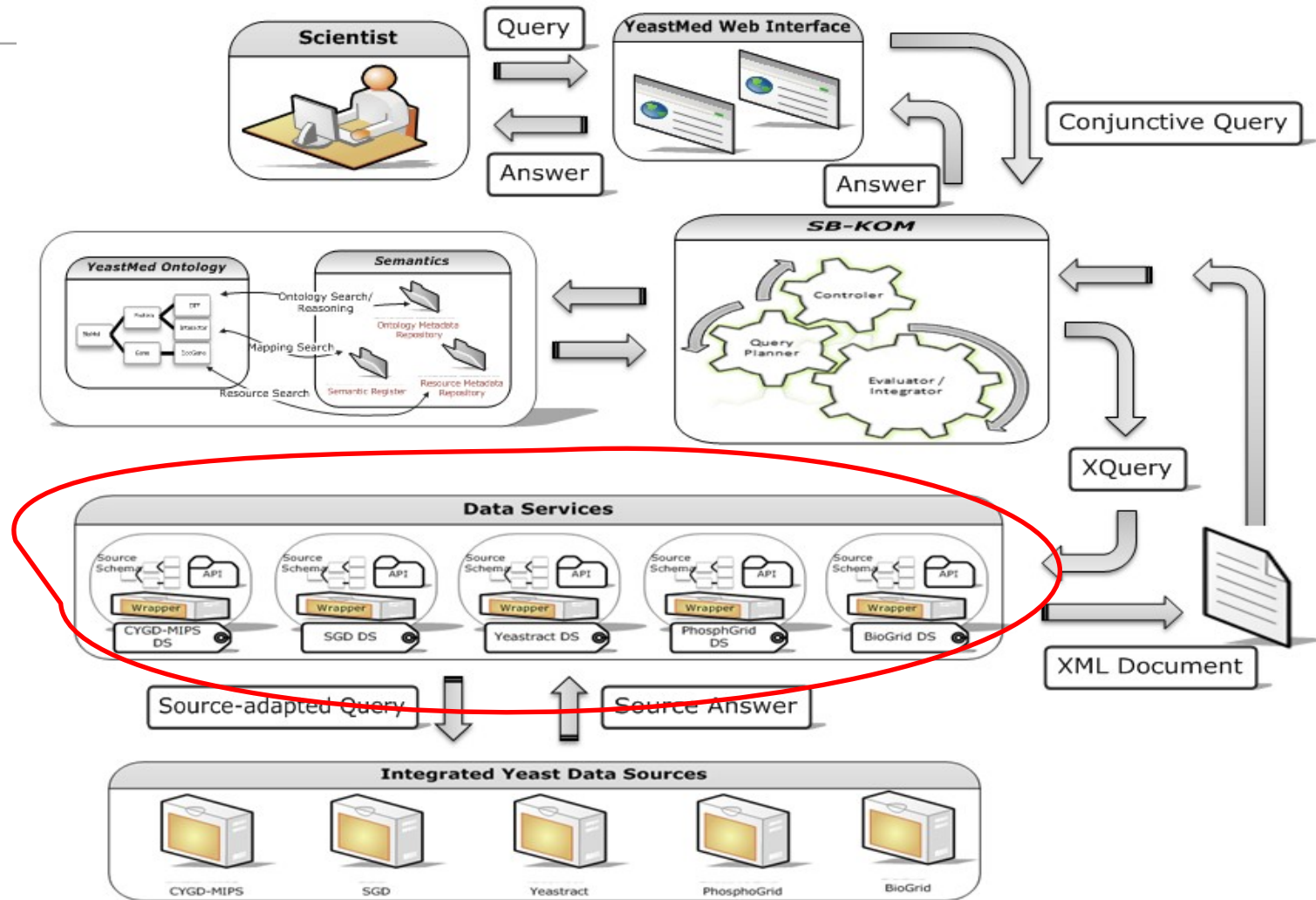


YeastMed: Data Sources

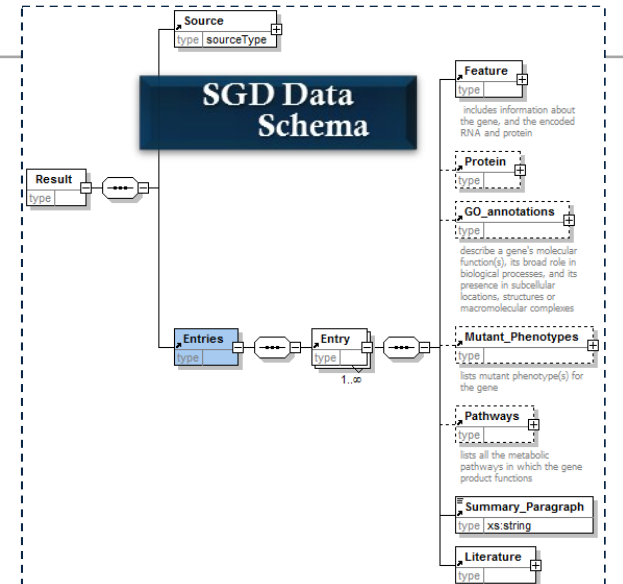
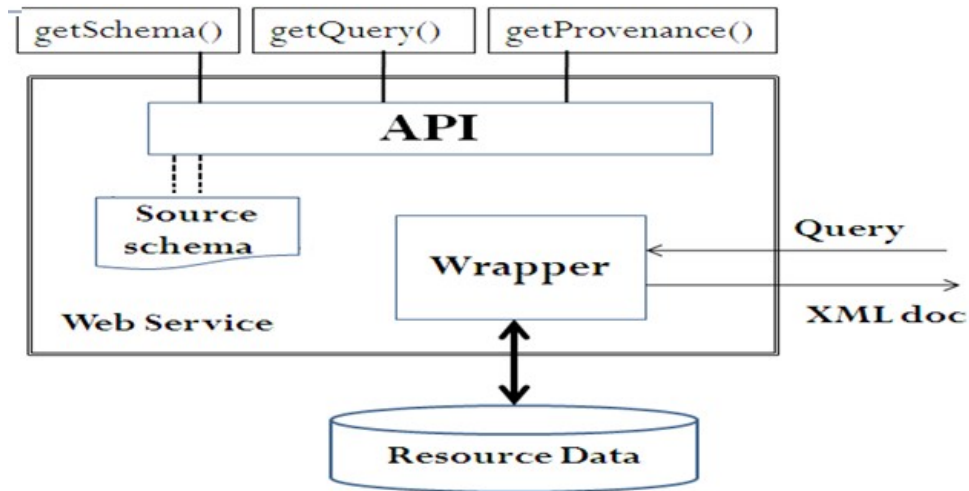
- **SGD:** Collection of genetic and molecular biological information about Sc.
- **MIPS-CYGD:** Information on the molecular structure and functional network of Sc.
- **Yeasttract:** A repository of regulatory associations between transcription factors and target genes, based on experimental evidence.
- **PhosphoGrid** records the positions of specific phosphorylated residues on gene products.
- **BioGrid:** An online interaction repository with data compiled through comprehensive curation efforts.



YeastMed: Architecture



YeastMed: Data Services



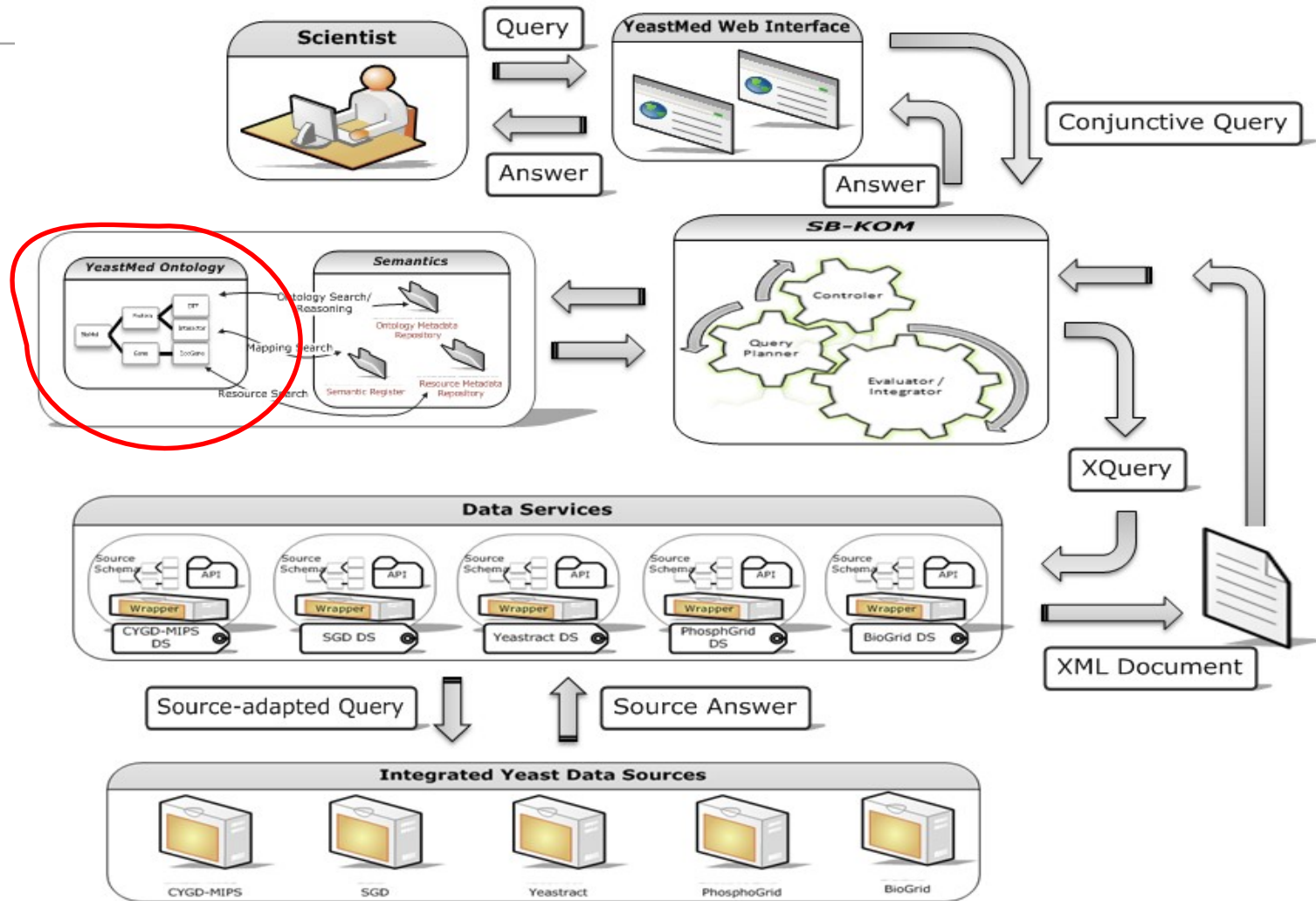
Operations

Name	
getSchema	--
getDataProvenance	--
Query	--

Endpoints [Add](#) [Remove](#)

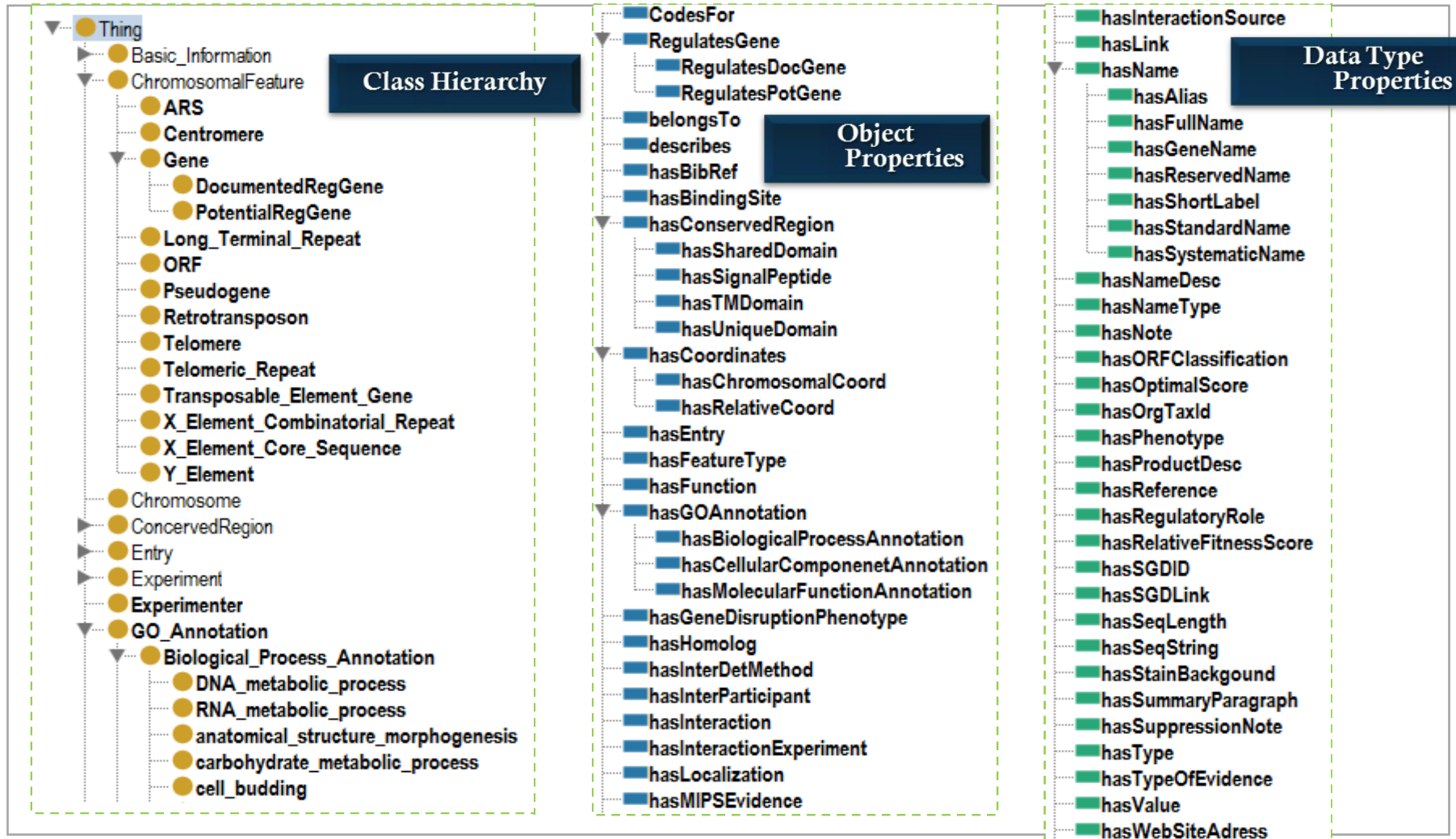
<input type="checkbox"/>	
<input type="checkbox"/>	http://172.16.51.3:8080/SGD/services/SGD

YeastMed: Architecture

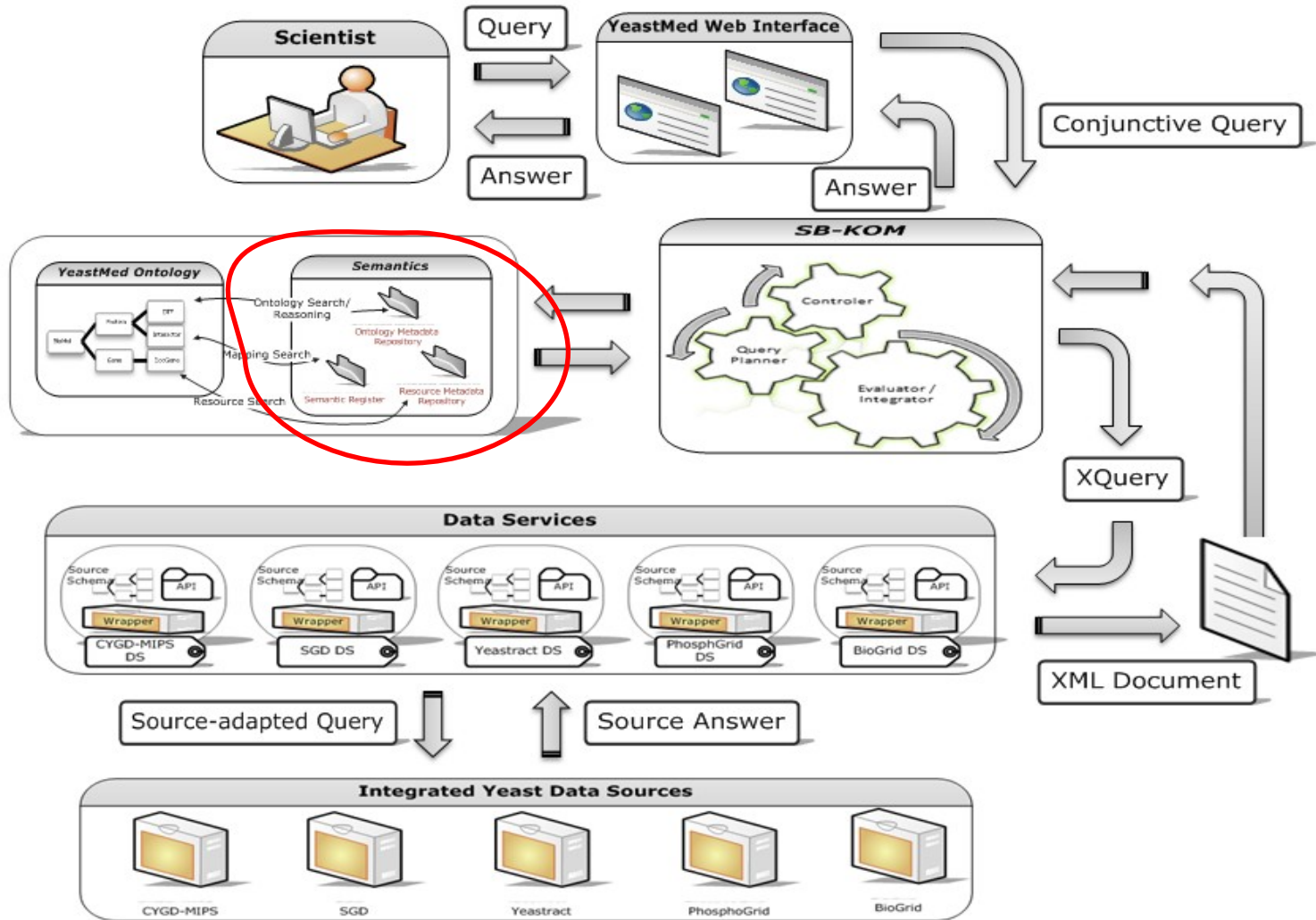


YeastMed: Ontology

Nature Precedings : doi:10.1038/npre.2010.5396.1 : Posted 15 Dec 2010



YeastMed: Architecture

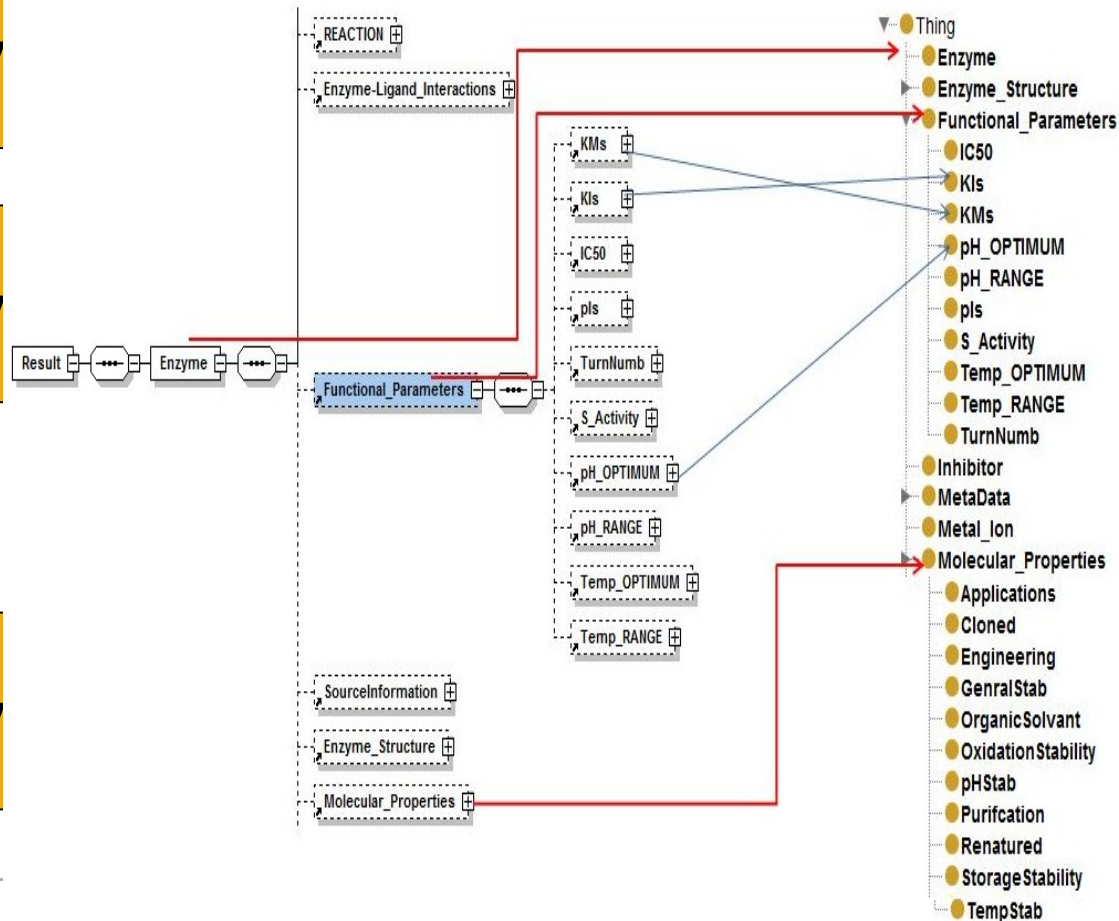


YeastMed: Mappings

Data Serv

Data Serv

Data Serv



ontology

YeastMed: Mappings

- **Class Mapping:** *it maps an Ontology class to the source schema.*

XPath-Element-Location, Ontology-Class-Name, Correspondence-index

- **Datatype Property Mapping:** *it maps an Ontology datatype property to the source schema.*

Example:

Result/Entries/Entry/Protein, Protein, 100

XPath-Domain-Location; XPath-value-Location, Ontology-Domain-Name; Property-Name, correspondence-index

- **Object Property Mapping:** *it maps an Ontology object property to the source schema.*

XPath-Domain-Location; XPath-Range-Location, Ontology-Domain-Name; Ontology-Range-Name;

Property-Name, correspondence-index

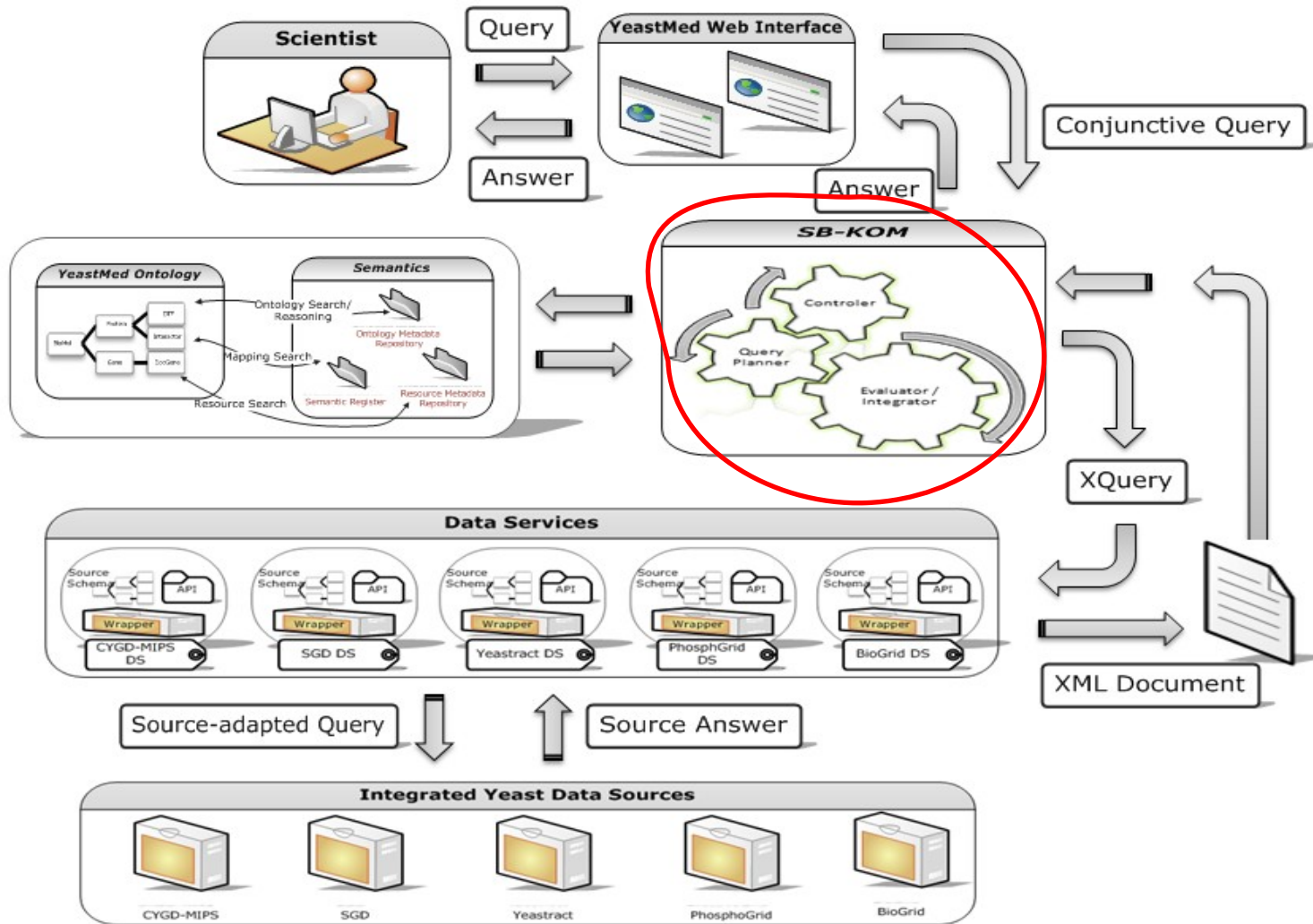
Example:

Result/Entries/Entry/Protein;

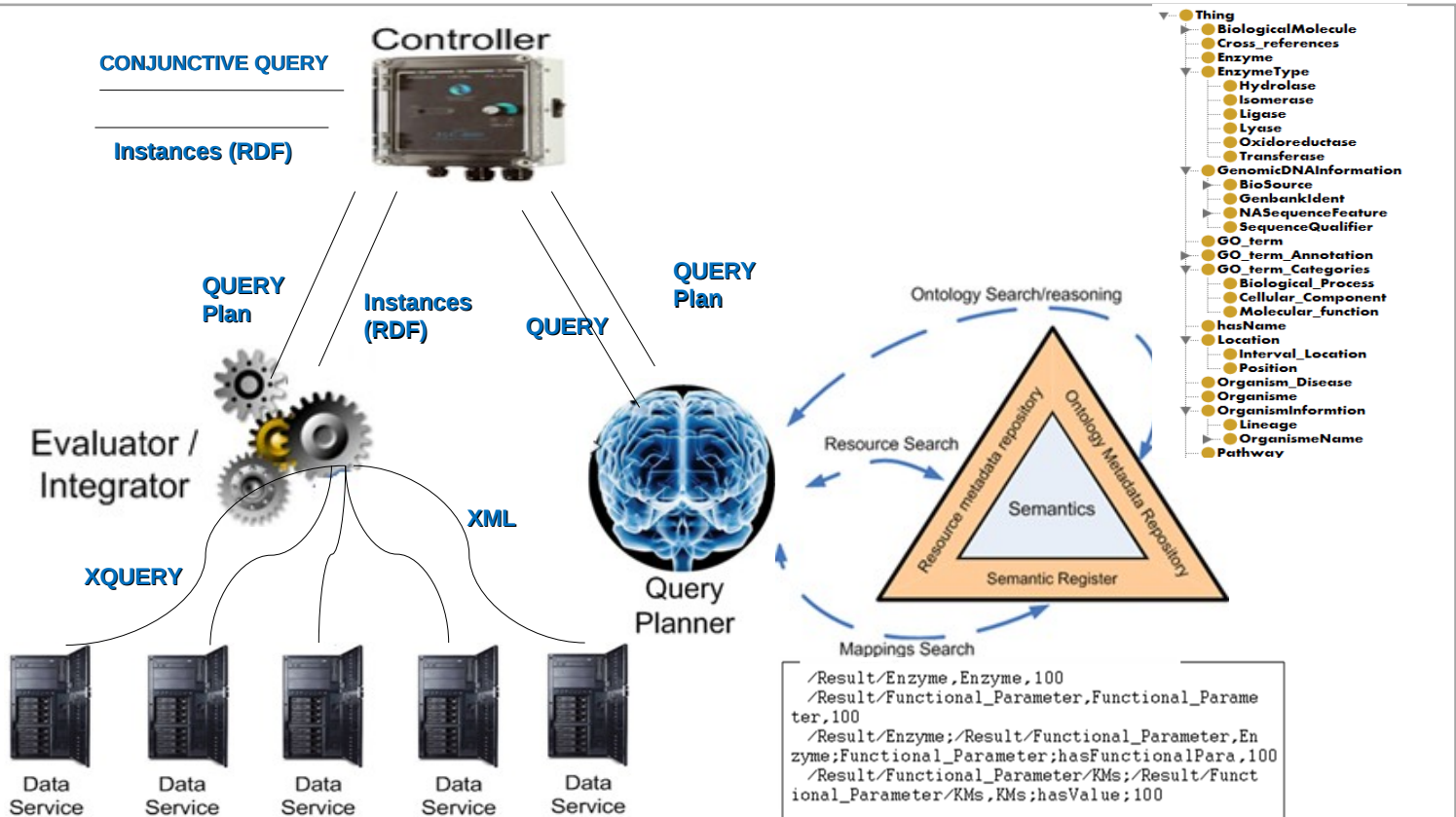
Example:

Result/Entries/Entry/Protein/SysName, TranscriptionFactor; hasName, 100
Result/Entries/Entry/Protein; Result/Entries/Entry/Literature, Protein; Bib
Ref; hasBibRef, 100

YeastMed: Architecture



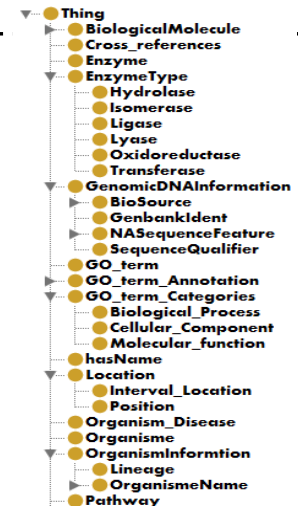
YeastMed: SB-KOM



YeastMed: SB-KOM

“find all the **bibliographic References** of **topoisomerase III** and also all the information about the **phosphorylation sites** that contain the **transcription factors** of DNA Topoisomerase III, and specially the one (or ones if exist) whose gene is located on the **Chromosome XVI**.”

- **Classes** : *Protein, BibRef, TranscriptionFactor, Chromosome* and *PhosphoSite*.
- **Datatype properties**: *hasDescription, hasSystematicName* and *hasName*
- **Object properties**: *hasBibRef, regulatedBy, belongsTo* and *hasPhosphoSite*

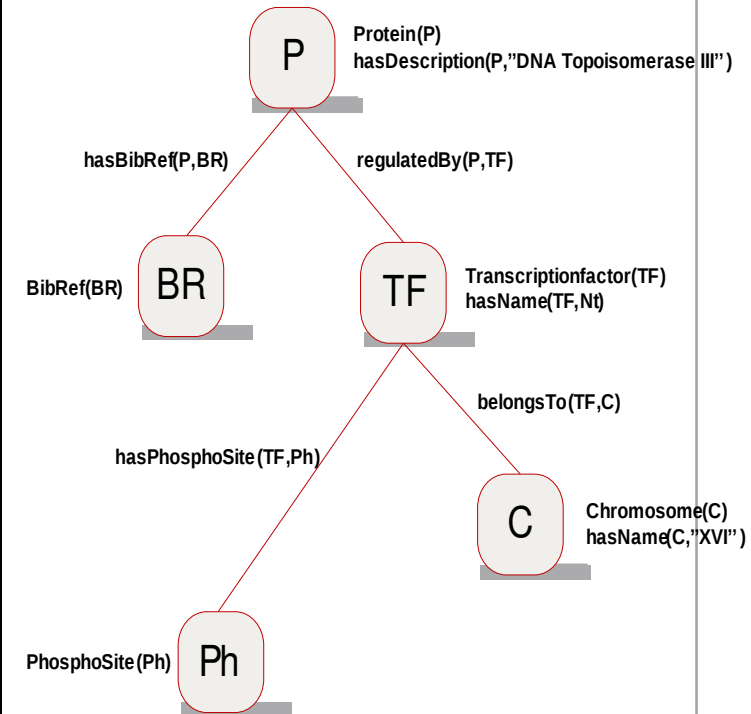


$$\text{Ans}(BR, Ph) := \text{Protein}(P), \text{hasDescription}(P, \text{"DNA Topoisomerase III"}), \text{BibRef}(BR), \text{hasBibRef}(P, BR), \\
 \text{hasSystematicName}(P, SN), \text{regulatedBy}(P, TF), \text{hasName}(TF, Nt), \text{TranscriptionFactor}(TF), \\
 \text{Chromosome}(C), \text{hasName}(C, \text{"XVI"}), \text{BelongsTo}(TF, C), \text{PhosphoSite}(Ph), \text{hasPhosphoSite}(TF, Ph);$$

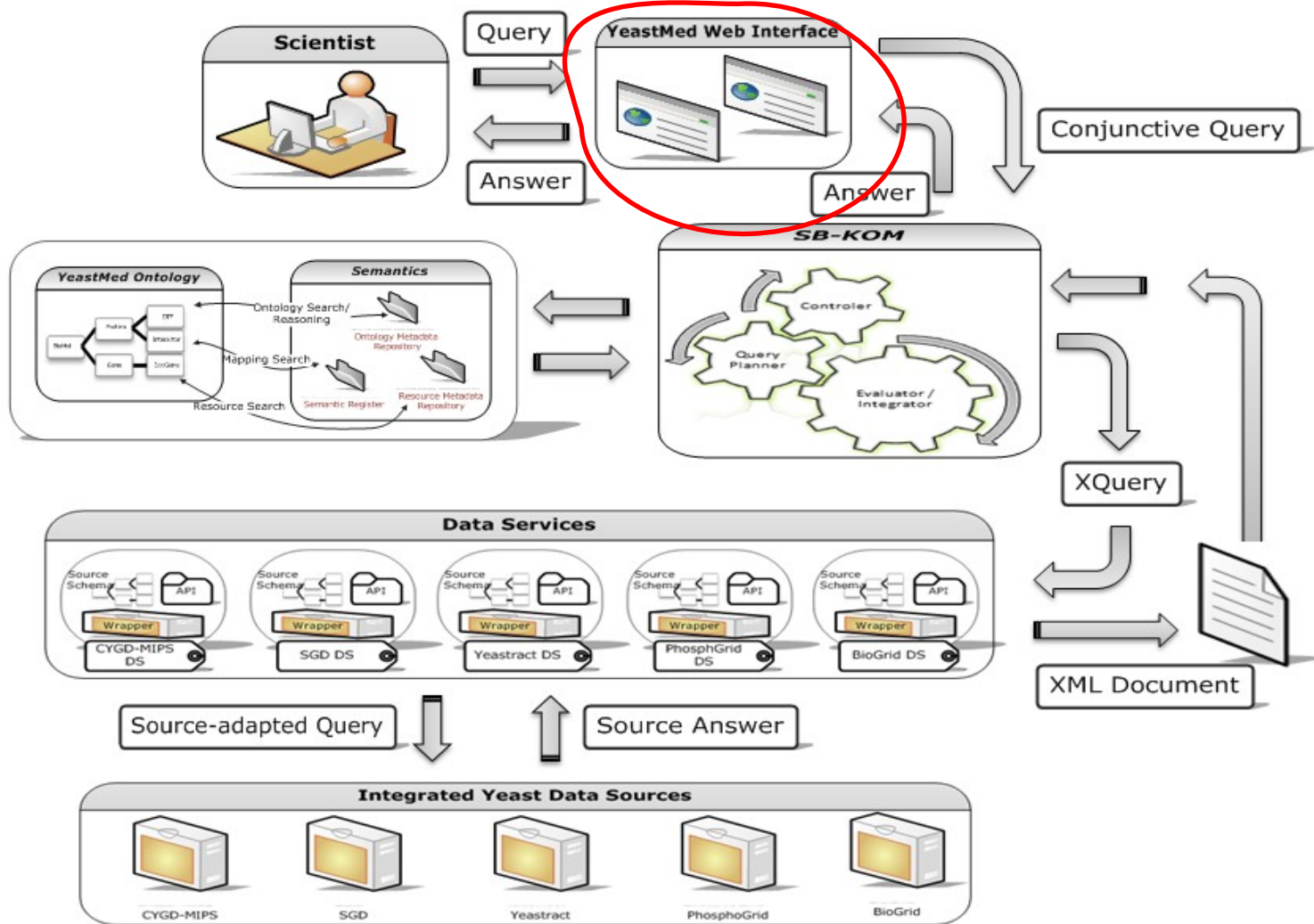
YeastMed: SB-KOM

Involved Data Sources : **SGD, Yeasttract and PhosphoGrid**

Group	Query	Mapping source
G1	Protein(P), hasBibRef(P, BR)	SGD
G2	Protein(P), hasDescription(P, "DNA Topoisomerase III")	SGD
G3	Protein(P), hasSystematicName(P, SN)	Yeasttract
G4	Protein(P), RegulatedBy(P, TF)	Yeasttract
G5	TranscriptionFactor(TF), hasName(TF, Nt)	Yeasttract
G6	TranscriptionFactor(TF), belongsTo(TF, C)	Yeasttract
G7	TranscriptionFactor(TF), hasPhosphorylationSite(TF, Ph)	PhosphoGrid
G8	Chromosom(C), hasName(C, "XVI")	Yeasttract
G9	regulatedBy(P, TF)	Yeasttract
G10	hasBibRef(P, BR)	SGD
G11	belongsTo(TF, C)	Yeasttract
G12	hasPhosphoSite(TF, Ph)	PhosphoGrid
G13	Protein(P)	SGD Yeasttract PhosphoGrid
G14	TranscriptionFactor(TF)	Yeasttract PhosphoGrid
G15	BibRef(BR)	SGD
G16	Chromosome(C)	Yeasttract
G17	PhosphoSite(Ph)	PhosphoGrid



YeastMed Architecture



YeastMed: Web Interface

[Organism](#)
[Gene](#)
[Protein](#)
[Enzyme](#)
[Pathway](#)

Protein Information

Protein Names :

Recommended Name :

Full Name : Acetyl-coenzyme A carboxylase carboxyl transferase subunit beta

Protein Existence :

Protein Existence :

inferred from homology

Protein Annotation :

function

catalytic activity

similarity

subunit

pathway

collected

SubCellular Location :

Employment

To use *YeastMed* follow the two steps below.

Step 1 : Selection of databases

This step allows you to select Databases wished to be interrogated

All possible databases

CYGD-MIPS

YEASTRACT

SGD

BioGrid

PhosphoGrid

Step 2 : Query formulation

This step allows you to formulate queries in order to retrieve chromosomal features that match the selected criteria.

Here you can specify instances of selected concepts (e.g. by providing Gene and Protein names, IDs,...etc.)

This allows you to select more criteria for your queries.

If you need to

Gene
Protein
 Pseudogene
 ORF
 rRNA
 snRNA

located on
 Coding for
 having Sequence
 having function
 Coded by
having description

Gene
 Protein
 Pseudogene
 ORF
 rRNA
 snRNA

DNA Topoisomerase III

and

Coding for
 having Sequence
 having function
 Coded by
 having Description
regulated by

Gene
 Protein
 Pseudogene
 ORF
Transcription Factor
 rRNA

Conclusions

- ▶ Allowing a transparent and simultaneous access to several autonomous and heterogeneous Yeast sources;
- ▶ Helping biologists to find convenient data to interpret results of their experiments;
- ▶ Avoiding biologists to confront technical and structural problems in data retrieving process.
- ▶ We're working on extending the set of the integrated sources.

Thanks

- Contact a_briache@lcc.uma.es
- Available soon on
<http://www.yeastmed.uma.es>