

SWAT4LS, Berlin, 10 Dec. 2010

A renaissance for the point mutation:
from legacy data to semantic web services

Christopher J. O. Baker
University of New Brunswick



Outline

- Problem Statement
- Mutation Databases
- Text Mining
- Grounding
 - Mutations
 - Molecular Functions
- Mutation Impact Extraction
 - NAR Subset
- Performance Metrics
- Reuse Scenarios
- Impact Prediction Data
- 3 DEMOS reusing mutation annotations and impacts
- Knowlegator / mSTRAP
- PubMed Semantic Assistant
- SADI Web Services
 - 3 sample queries



Data Integration Challenges:

- (i) publishing of salient mutation impact descriptions in unstructured text.
- (ii) prevalence of boutique databases of mutation information with a many years of latency.
- (iii) errors within manually populated mutation databases.
- (iv) mining of mutations from scientific documents for denovo database creation.



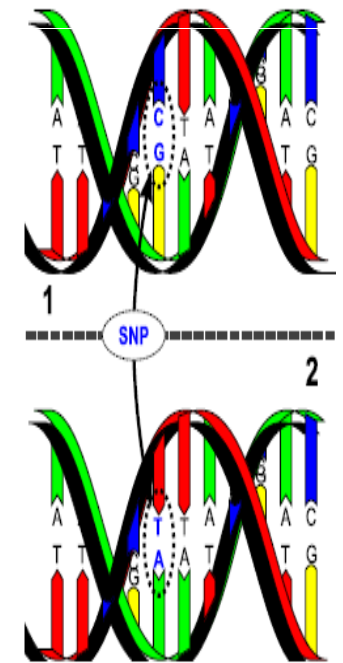
Types of Mutations

- Genotype
- Missense/nonsense
- Splicing
- Regulatory
- Small deletions
- Small insertions
- Gross deletions
- Gross insertions
- Repeat variations
- Single Amino Acid polymorphisms (exp.)
- SNP / nsSNP

Single Nucleotide Polymorphism (SNP)

Wikipedia¹:

- "DNA sequence variation occurring when a single nucleotide—A, T, C, or G—in the genome differs between members of a species"
- "Successful" point mutations in a population: > 1%



¹http://de.wikipedia.org/wiki/Single_Nucleotide_Polymorphism



Data Types Stored

- Residues
- Coordinates
- DNA / Protein
- Publication Links
- Experimental conditions
- Organisms
- Impact Annotations
 - Phenotype
 - Modified Properties
 - Stability,
 - Protein Interaction
 - Molecule conformation
 - Pathway Dynamics
 - Units of measurement



Relevance of non-synonymous SNPs

- Changes to the protein lead to changes in
 - molecular function (e.g., impaired signaling)
 - metabolism (e.g., cystic fibrosis)
 - cellular phenotype (e.g., neurofilament aggregation)
 - physiological phenotype (e.g., QT prolongation)
 - the human being
- SNPs make human beings different
 - The collection of SNPs form the haplotype
 - SNPs and haplotypes are the basis to understand human variability



Future Relevance of SNPs

- Increasingly single individuals scanned for their genotypic variability
- Comparison against standard genomes:
 - Human genome project
 - 1,000 genomes project (www.1000genomes.org)
- Commercialization of sequencing
 - High speed and cost efficiency
 - Sequencing of individuals for a reasonable amount of money

• =>

..... How do we make sense of the data and reuse it



400 + Mutation Databases

Web [Images](#) [Maps](#) [News](#) [Video](#) [Gmail](#) [more](#) ▼



mutation database

Search

[Advanced Search](#)
[Preferences](#)

Search: the web pages from Canada



Search help

Statistics

New genes

What's new

Background

Publications

Contact us

Register

Gene symbol

Go!

The Human Gene Mutation Database at the Institute of Medical Genetics in Cardiff

Copyright © Cardiff University 2008. All rights reserved.

This database is maintained by D.N. Cooper, E.V. Ball, P.D. Stenson, A.D. Phillips, K. Howells and M.E. Mort with the assistance of N.S.T. Thomas.



*Please note that the less up-to-date public version of our database is free only for [registered](#) users from academic institutions/non-profit organisations. Commercial users are required to purchase a license from BIOBASE, commercial and academic non-profit users wishing to access the most up-to-date version of the database (see [example](#) HGMD Professional entry). Read more about how HGMD is [funded](#).



CENTER FOR COMPUTATIONAL BIOLOGY AND BIOINFORMATICS
CUTTING EDGE RESEARCH THROUGH COLLABORATION
MOONEY LAB



MutDB

The goal of MutDB is to annotate human variation data with protein structural information and other functionally relevant target or disease, search below (the search may take up to a minute).

Universal Mutation Database

Software and databases for mutations in human genes

UNIVERSAL MUTATION DATABASE SOFTWARE

UMD softwares have moved to the [UMD website](#)

DISEASE OR GENE SPECIFIC DATABASES

Links to gene/disease specific web sites:



International Agency for Research on Cancer
Centre International de Recherche sur le Cancer

IARC TP53 Mutation Database

DNA Mutation Database

Familial Hypertrophic Cardiomyopathy



This database contains mutations in various genes known to cause familial hypertrophic cardiomyopathy, a genetic disorder associated with defects in the sarcomere [1]. Only gene symbols approved by [HUGO](#) are used and mutations are reported in accordance with guidelines recommended by the [Mutation Database Initiative](#) of HUGO and [EBI](#). Clinical phenotypes are classified as 'typical' or 'atypical'. A typical phenotype means clinical manifestations described in the WHO criteria [2] and/or McKenna's criteria [3].



Mutation Databases

- **Specialized databases:**
 - focus on an individual disease / phenotype
- **General databases:**
 - Omim (NCBI): monogenetic diseases, keeping track of genetic variability and disease implication, gathered from the literature
 - GAD: SNPs from associations studies, gathered from the literature – identifies medically relevant polymorphism from the large volume of polymorphism and mutational data
 - dbSNP (NCBI): broad range of SNPs - 51,312,474 variations for 43 different organisms



Different in scope and scale

	Diseases	Distinct genes	Associations
MIM	1864	2708	4851
GAD	2303	1740	7945

- Omim makes reference to less diseases and to more genes
- GAD makes reference to more diseases, gene-disease associations, but monitors less genes



Mining Mutation Annotations

- From databases
 - High number of specialized databases
 - No uniform format => data mining is difficult
- From the literature:
 - Advantage:
 - Contextual information supports interpretation
 - Universal resource => different types of SNPs available, i.e. disease related and experimental SNPs
 - Disadvantage:
 - Normalisation of SNPs is not straight forward
 - Language variability



Maintenance of Mutation Databases

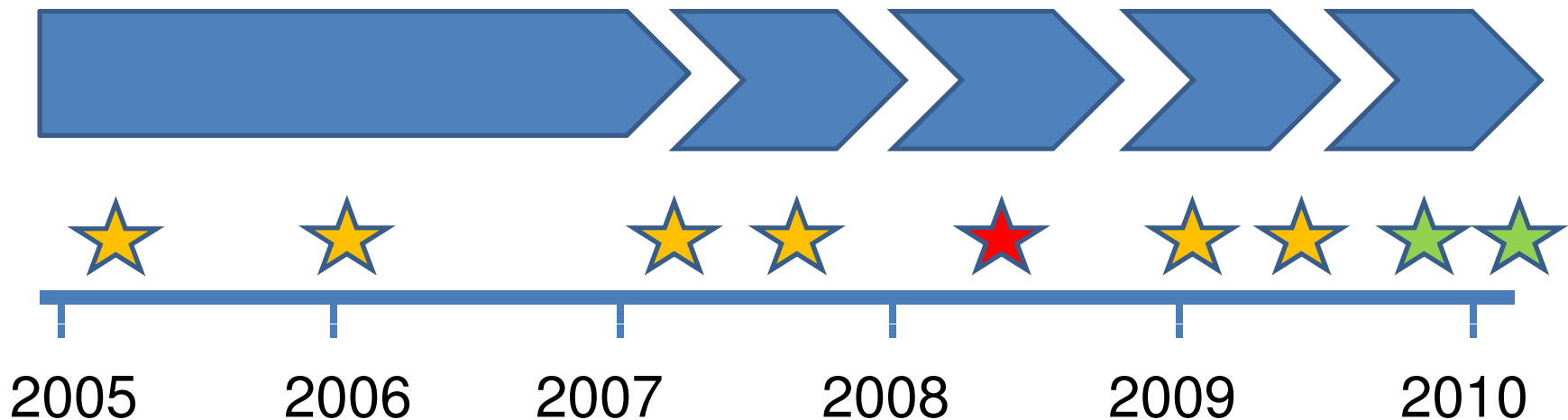
- Update
 - Scientist reading papers and entering data (PMD 1999)
 - Authors depositing information
- Curation
 - Often out of date (backlog of unprocessed papers)
 - Funding dries up ?
 - PDB ~40 % inaccurate wrt mutations*
 - What about other techniques

[Intrinsic Evaluation of Text Mining Tools May Not Predict Performance on Realistic Tasks](#)

- *J. Gregory Caporaso, Nita Deshpande, J. Lynn Fink, Philip E. Bourne, K. Bretonnel Cohen, Lawrence Hunter; Pacific Symposium on Biocomputing 13:640-651(2008)*



Mutation Mining: a history



Front ISR 2005	- Baker and Witte - Mutation Mining
JBCB 2007	- Kanagasabai et al – mSTRAP Workflow
JBCB 2007	- Witte and Baker - Systematic Evaluation Mutation Extraction
IJBRA. 2007	- Witte, Kappler, Baker – <i>En. semantic access</i> to protein eng. Lit.
ISMB 2008	- Baker et al - Towards Ontology-driven Navigation of Mutation Lit.
BMC Bioinf. 2009	- Baker and Rebholz-Schuhmann Special Issue
AMIA 2009	- ISCB Workshop - Interpretation of <i>Mutations</i> with <i>Semantic Supp</i>
DILS 2010	- Laurila et al - Algorithm for Mutation Grounding
BMC Geno 2010	- Laurila et al - Semantic Infrastructure for Mutant Impact Extraction

Single Nucleotide Polymorphism—examples

BACKGROUND AND PURPOSE: The collagen alpha2(I) gene (COL1A2) on chromosome 7q22.1, a positional and functional candidate for intracranial aneurysm (IA), was extensively screened for susceptibility in Japanese IA patients. **METHODS:** Twenty-one single nucleotide polymorphisms (SNPs) of COL1A2 were genotyped in genomic DNA from 260 IA patients (including 115 familial cases) (mean age, 59.9 years) and 293 controls (mean age, 61.6 years). Differences in allelic and genotypic frequencies between the patients and controls were evaluated with the chi(2) test. Circular dichroism spectrometry was monitored with collagen-related peptides that mimic triple-helical models of type I collagen with Ala-459 and Pro-459 to estimate the conformation and stability of alterations. **RESULTS:** Significant genotypic association in the dominant model was observed between an exonic SNP of COL1A2 and familial IA patients (chi(2)=11.08; df=1; P=0.00087; odds ratio=3.19; 95% CI, 2.22 to 6.50). This SNP induces Ala to Pro substitution at amino acid 459, located on a triple-helical domain. Circular dichroism spectra showed that the Pro-459 peptide had a higher thermal stability than the Ala-459 peptide. **CONCLUSIONS:** The variant of COL1A2 could be a genetic risk factor for IA patients with family history.

maps to rs number

states, locations, genes, types

PMID: 14739420, T. Yoneyama et. al: "Collagen type I alpha2 (COL1A2) is the susceptible gene for intracranial aneurysms".

Roman Klinger et al . Identifying gene specific variations in biomedical text. Journal of Bioinformatics and Computational Biology, Special Issue: Making Sense of Mutations requires Knowledge Management, December 2007.

Text Mining Systems for Mutations: Performance

2004 – 2009

- **MuteXt** (Protein Point Mutation) (P=87.9% R=85.8%)
(P=49.3% R= 64.5%)
- **MEMA** (Regex DNA / Protein, HUGO) (P=75% R=98%)
- **MutationFinder** (Regex) + rules (P=98%, R=81%)
- **ProMiner** SNPExtraction and normalization / grounding (P=78%, R=67%)
- **mSTRAP** RegEx plus protein or organism name, (P=94.5% R=79.6%)
- **mSTRAP** Grounding / Normalization to db (P=91.8% R=80.9%)
- **VTag**: (CRF approach) in special context of cancer, no mapping to database
- **OSIRIS**: Query expansion: for all SNPs of a found gene: PubMed query)
slow, limited to results of PubMed search engine (P=99% R=82 %)

Ongoing Challenges

- ❑ Normalization (Manual nom to dbSNP R=61%)
- ❑ Defining appropriate metrics and definitions for individual tasks for system benchmarking
- ❑ Mining Impacts / Causality from sentences

- Protein Engineering – Impact
- Disease Studies - Causality

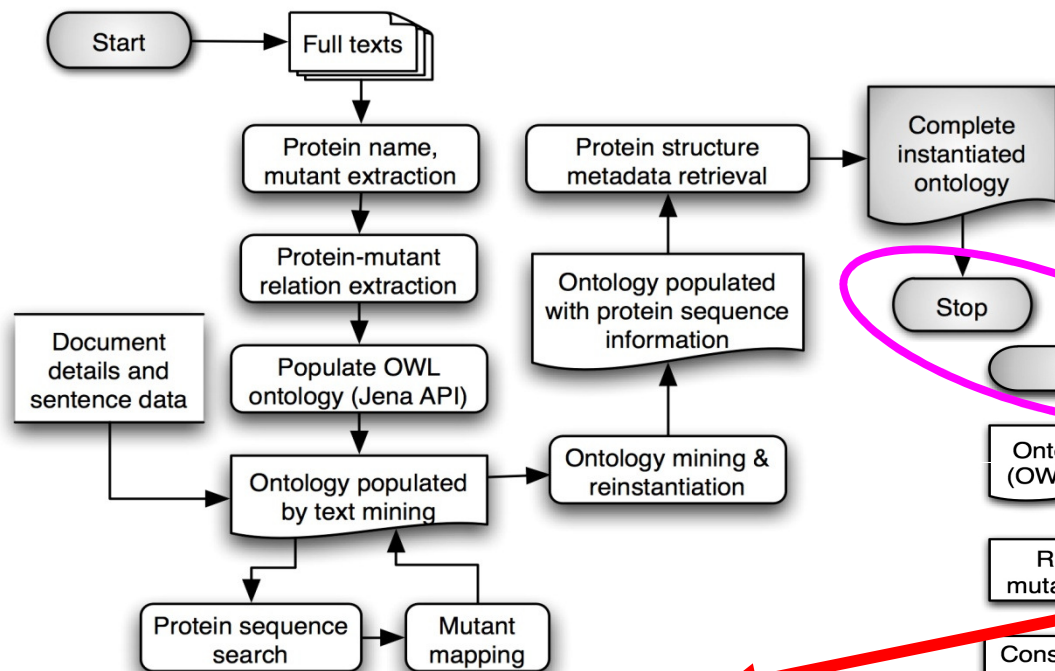
<hasContent>It led to the amino acid sequence change of **H1047R**, which was found to be a **gain-of-function** mutation at the **kinase domain**. (d) FISH analysis showing DNA copy number gains and amplifications of **PIK3CA** locus on the primary tumor cases T10 and T24. </hasContent></Sentence>

- ❑ Cottage Industry ? – Scale up to publishing of mutation annotations according to Semantic Metadata for incorporation into systems level approaches and prediction tools e.g.
 - pathway analyses based on SNP annotations
 - impact prediction for un annotated residues
 - reuse scenarios



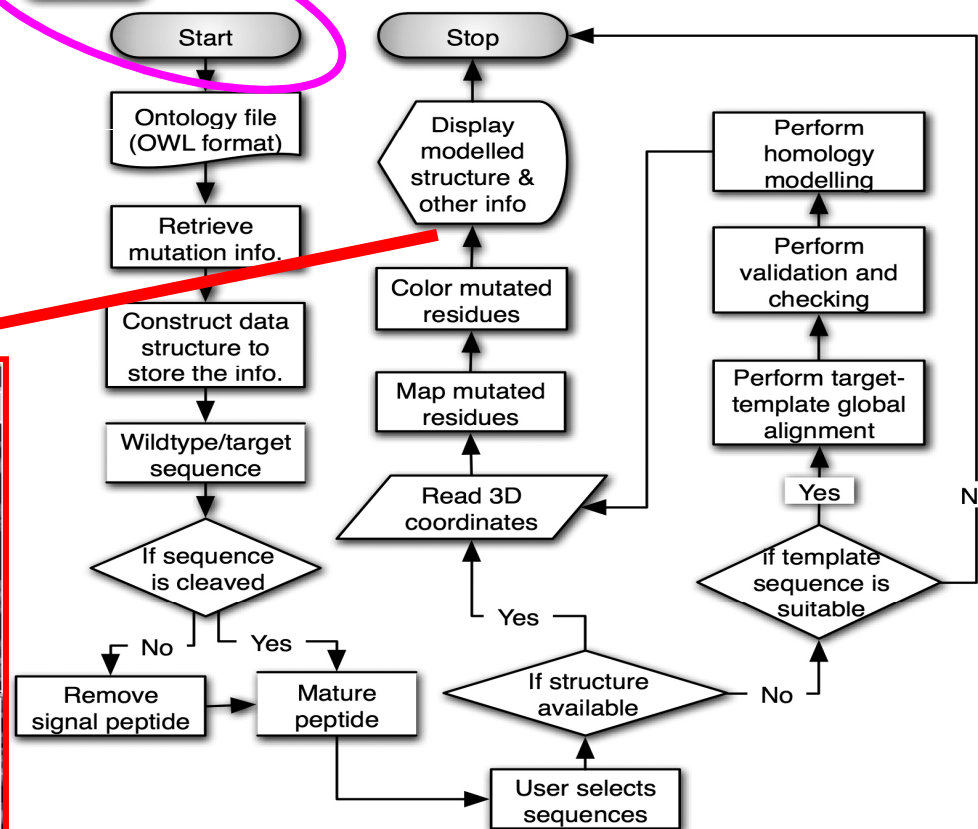
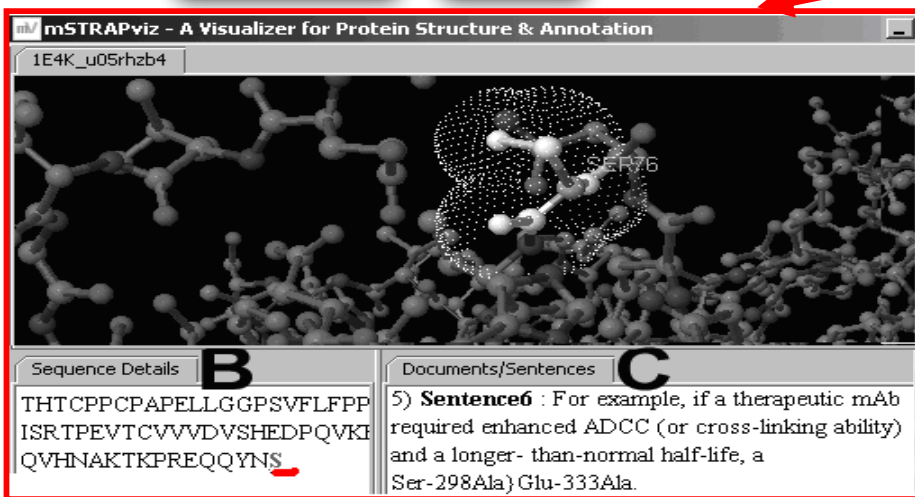
Mutation Reuse Scenarios

- Annotation
 - Cancer Genome (Forbes et al 2010)
 - Pathway (Bauer Mehren et al 2009)
 - **Pubmed Abstracts (Laurila et al 2010)**
 - **Protein Structure (Kanagasabai et al 2007)**
- Impact Prediction
 - SNAP (Bromberg and Rost 2008)
 - Membrane protein stability (Winnenburg et al 2009)
- Federated Query over Mutation Triplestores
 - **SADI Semantic Web Services (Riazanov et al 2010)**



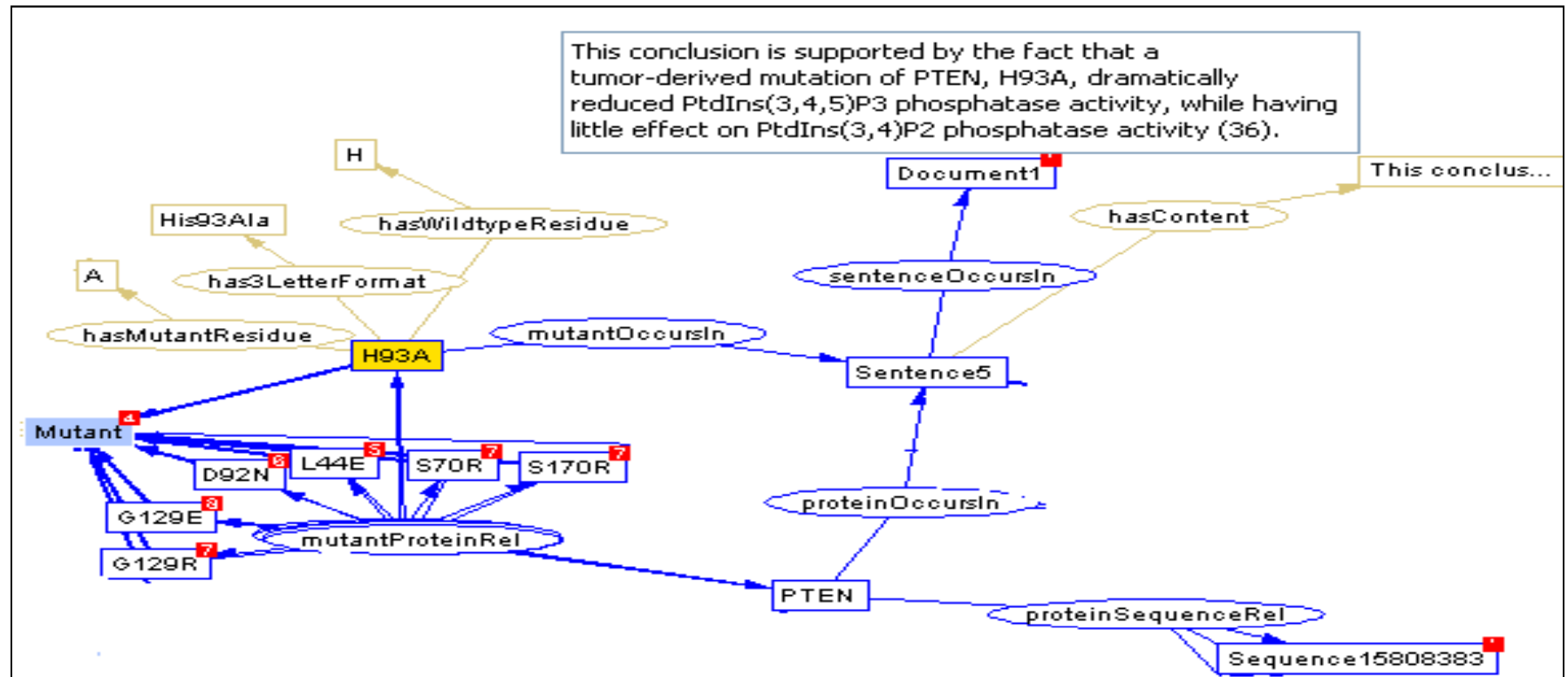
E214F, The mutation increased temp stability by 7 °C

R211G, 100 fold lower activity resulted from disruption of disulphide bonds ...

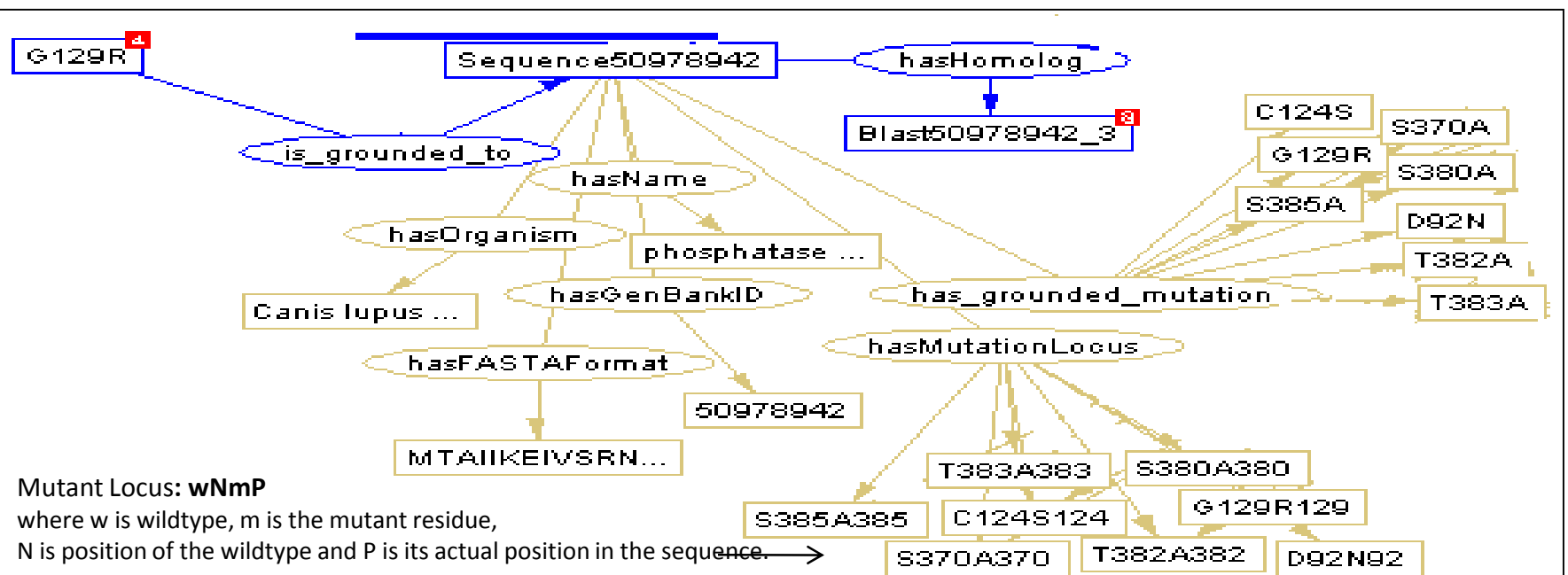


- COSMIC database content used for building a gold standard corpus for evaluation of mutation extraction and grounding
- Target families:
 1. PIK3CA - 42 papers out of 78 online
 2. FGFR3 - 37 papers out of 59 online
 3. MEN1 - 19 papers out of 72 online
- Protein–Mutation Tuple (extraction, normalisation)
- Mutation Grounding to sequences
- Sequences then checked by pair wise sequence alignment with gold standard protein sequence with >99% homology

A

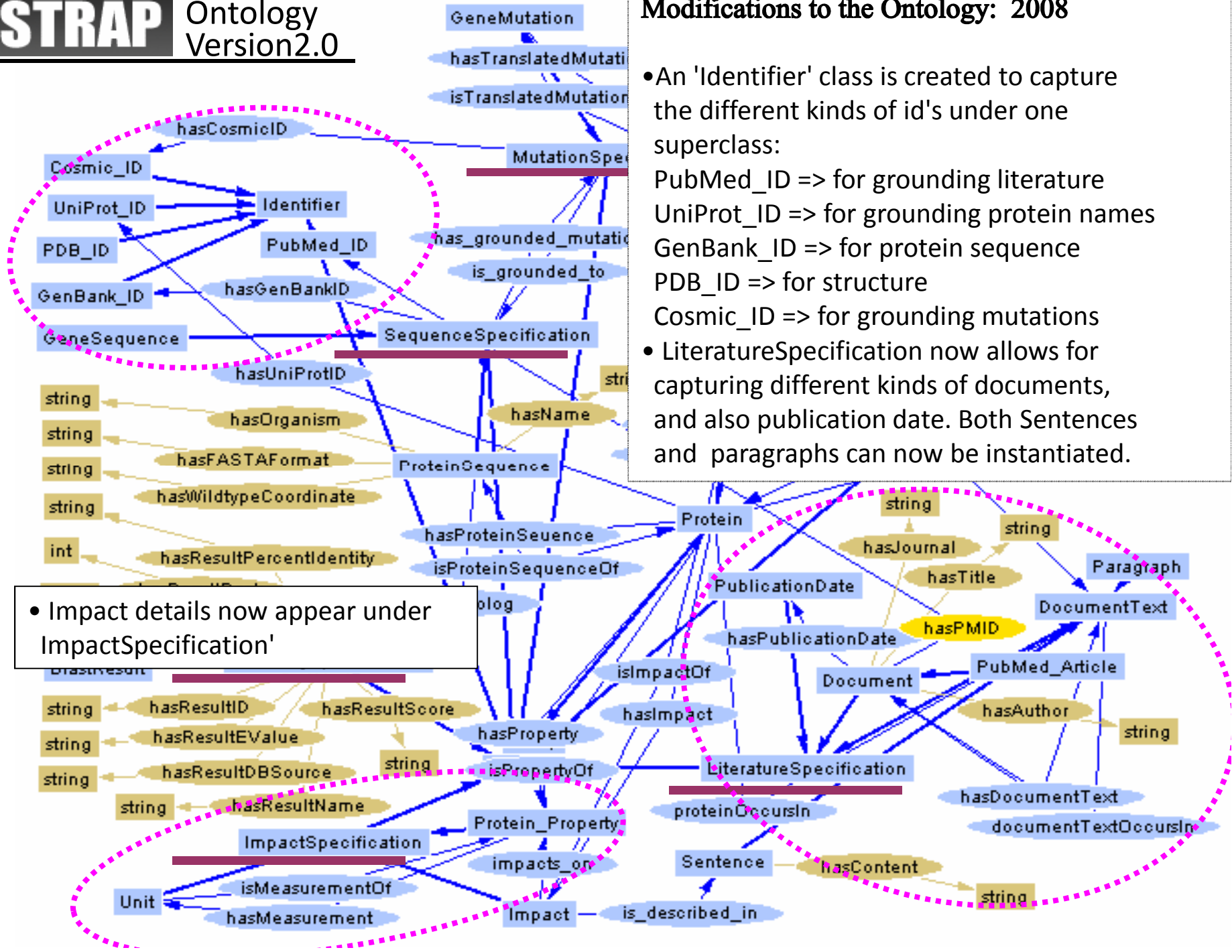


B



Modifications to the Ontology: 2008

- An 'Identifier' class is created to capture the different kinds of id's under one superclass:
 PubMed_ID => for grounding literature
 UniProt_ID => for grounding protein names
 GenBank_ID => for protein sequence
 PDB_ID => for structure
 Cosmic_ID => for grounding mutations
- LiteratureSpecification now allows for capturing different kinds of documents, and also publication date. Both Sentences and paragraphs can now be instantiated.



• Impact details now appear under ImpactSpecification'

Toward a Richer Representation of Sequence Variation In the Sequence Ontology (2010)

Michael Bada and Karen Eilbeck

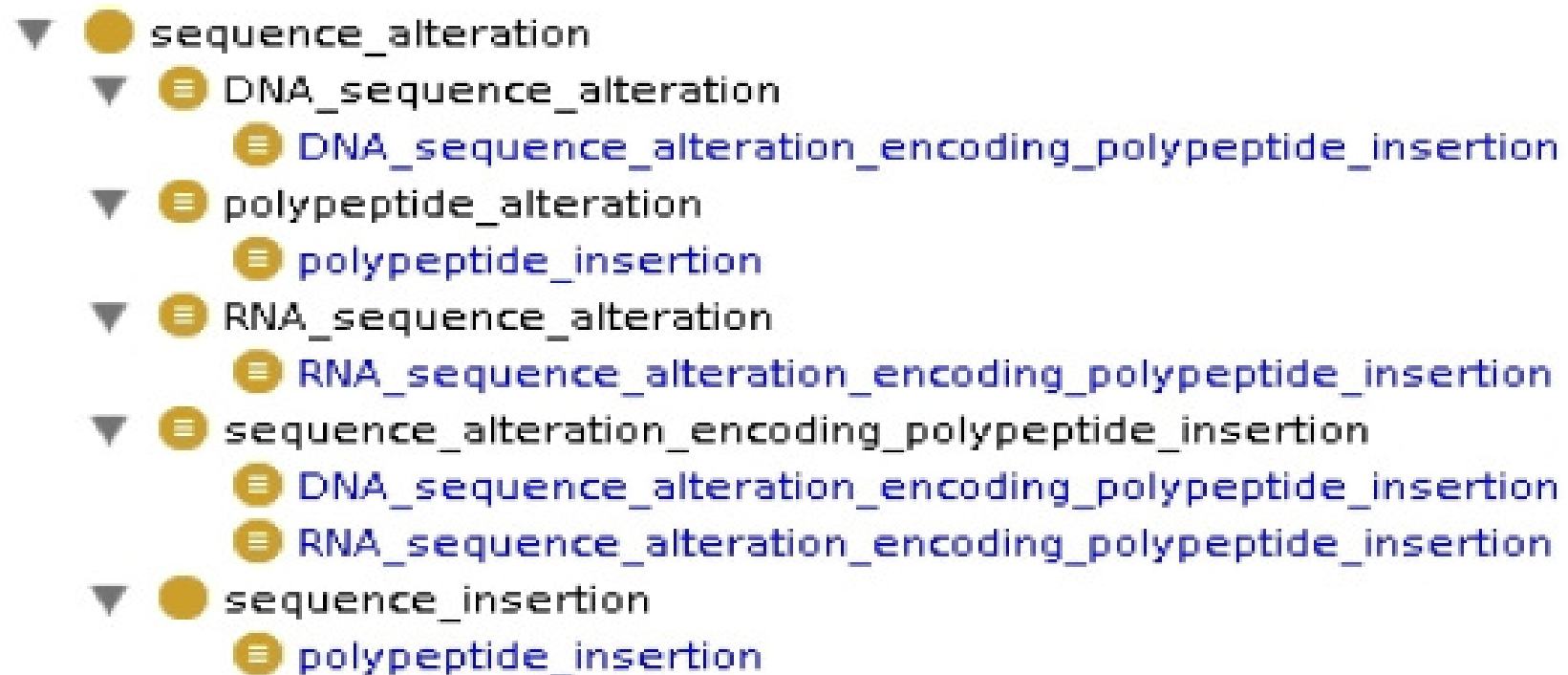


Figure 5. Automatically classified hierarchy of sequence alterations.

Proceedings of the Workshop on Annotation, Interpretation and Management of Mutations ([AIMM-2010](#)), Annotation, Interpretation and Management of Mutations A workshop at [ECCB10](#). Ghent, Belgium, September 26th, 2010.

Enterprise KnowleGator with RacerPro (nRQL) Version r0.9.0 Beta

File mSTRAPviz

Navigation

- KnowleGator
 - Ask a question...
 - with a Concept
 - BlastResult : (36)
 - Cosmic_ID : (4)
 - Document : (6)
 - DocumentText : (165)
 - Entity : (5)
 - GenBank_ID : (12)
 - Gene : (0)
 - GeneMutation : (0)
 - GeneSequence : (0)
 - GeneSynonym : (0)
 - HomologSpecification :
 - Identifier : (34)
 - Impact : (2)
 - ImpactSpecification : (1)
 - LiteratureSpecification :
 - MutationSpecification :
 - MutationSynonym : (21)
 - Paragraph : (0)
 - PDB_ID : (12)
 - Protein : (5)
 - ProteinMutation : (30)
 - ProteinSequence : (12)
 - ProteinSynonym : (0)
 - Protein_Property : (4)

Editor

Welcome Ask a question... mSTRAPviz

Type your free-text question here...

Question 1 : (6)

```

    graph TD
      A[Protein : PN_PTEN] --> B(hasProteinMutation)
      B --> C[ProteinMutation : ?Y1]
      C --> D(is_grounded_to)
      D --> E[SequenceSpecification : ?Y4]
      C --> F(hasImpact)
      F --> G[Impact : ?Y2]
      G --> H(impacts_on)
      H --> I[Protein_Property : lipid_phosphatase_activity]
  
```

Search

Advanced Simple

Output

Get the answer... Status details

Question 1 (6 results found)

ProteinMutation : ?Y1	Impact : ?Y2	SequenceSpecification : ?Y4
G129R	negative	Sequence73765544
G129R	negative	Sequence13928830
G129R	negative	Sequence50978942
C1245	negative	Sequence73765544
C1245	negative	Sequence13928830
C1245	negative	Sequence50978942

01:21

Enterprise Knowlegator with RacerPro (nRQL) Version r0.9.0 Beta

File mSTRAPviz

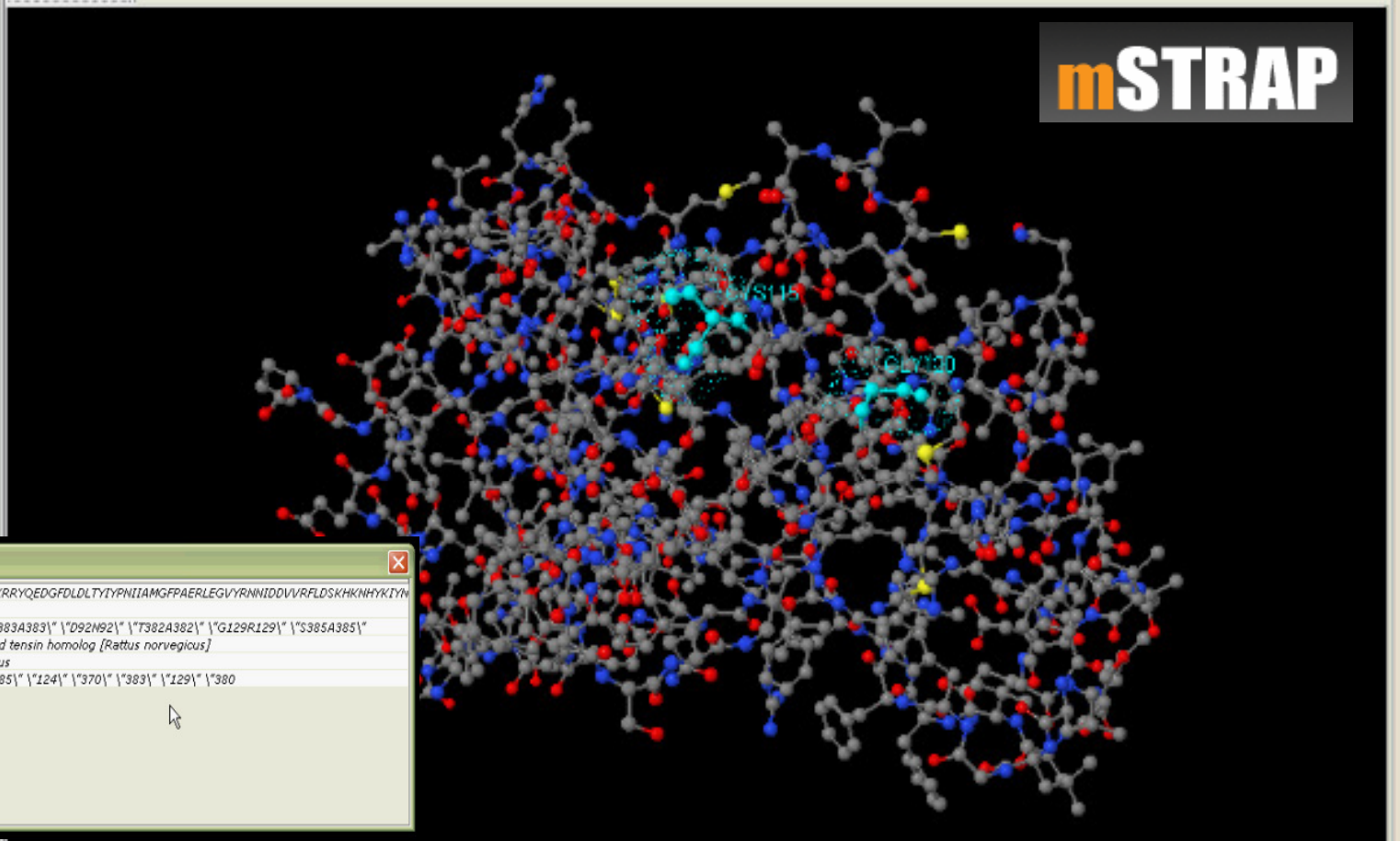
Navigation

- Knowlegator
 - Ask a question...
 - with a Concept
 - described_in : (4)
 - hasBlastResult : (72)
 - hasImpact : (6)
 - hasMeasurement : (5)
 - hasMutantSynonym : (30)
 - hasPMID : (6)
 - hasProperty : (4)
 - hasProteinSynonym : (21)
 - hasSentence : (108)
 - Impacts_on : (4)
 - isImpactOf : (6)
 - isMeasurementOf : (5)
 - isMutantSynonymOf : (30)
 - isPropertyOf : (4)
 - isProteinSynonymOf : (21)
 - mutantOccursIn : (246)
 - mutantProteinRel : (99)
 - mutantSequenceRel : (95)
 - proteinMutantRel : (99)
 - proteinOccursIn : (145)

Editor

Welcome Ask a question... mSTRAPviz

:2RTR_uir2np9v:



'Sequence13928830' Details...

hasFASTAFormat	MTAIKEIVSRNRRRYQEDGFDLDTYYPNIIAMGFPAERLEGVYRNNIDVVRFDSKHKHNYKIYN
hasGenBankID	I3928830
hasMutationLocus	S370A370 T383A383 D92H92 T382A382 G129R129 S385A385
hasName	phosphatase and tensin homolog [Rattus norvegicus]
hasOrganism	Rattus norvegicus
hasWildtypeCoordinate	92 382 385 244 370 383 129 380

Search

Advanced Simple

Sequence Details Alignment

```

S...
HAML TGLKIAVIGGDARQLEIIRKLTQQADIYLVGFDQLDHC
DSIILPVSATTEGGVSTVPSNKEEVVLRKQDHLDRTPAHCVIPS
LVKLFERDDIALYNSIPTVEGTIMEAIQHTDYTIHGSQVAVLC
EVKVGARSSAHLARITENGLVPTHTDELKEHVKDIDICINTIE
ILDLASRPGGTDPKYAEKQGIKALLAPGLPGIVAPKTAGQILF
+
  
```

Documents/Sentences

- Sentence1** : However, because translocation of GLUT4 in cells overexpressing a dominant inhibitory PTEN mutant (C124S) was similar to that of control cells, we conclude that endogenous PTEN may not modulate metabolic functions of insulin under normal physiological conditions. 2001 Academic Press Key Words: metabolism; signal transduction; insulin resistance; phosphatase; glucose.
- Sentence11** : C124S is a dominant inhibitory PTEN mutant (36).
- Sentence13** : Using an Akt phosphorylation assay to assess effects of PTEN to

News

*Kanagasabai R., Choo K. H., Ranganathan S. and *Baker CJO, Workflow for Mutation Extraction and Structure Annotation, Journal of Bioinformatics and Computational Biology, (2007), vol 5, 6:1319-1337*

- mSTRAPviz™ is now integrated into [Knowlegator™](#) ([pic1](#), [pic2](#)), an ontology-driven knowledge integration and navigation tool. This system will be presented at Toronto during the [ISMB2008 Highlights Track](#)

ISMB 2008 Highlights Track: HL17

Towards Ontology-driven Navigation of the Lipid Bibliosphere and Mutation Literature

Monday, July 21 - 10:45 a.m. - 11:10 a.m.

Room: 718A

Presented by: Christopher J. O. Baker, Institute for Infocomm Research, SG

Session Chair: Andrey Rzhetsky

DEMO clips of the Knowlegator™+mSTRAPviz™ system:

[Clip #1 : Query & Navigation in Knowlegator](#) (25MB)

[Clip #2 : Results from query fed into mSTRAPviz for automated homology modeling](#) (17MB)

- Release of v1.1 (9Jul08) - skip repeated modeling of same sequences when reselected

Mutation Impact Extraction

Laurila *et al.* *BMC Genomics* 2010, **11**(Suppl 4):S24
<http://www.biomedcentral.com/1471-2164/11/S4/S24>



PROCEEDINGS

Open Access

Algorithms and semantic infrastructure for mutation impact extraction and grounding

Jonas B Laurila¹, Nona Naderi², René Witte², Alexandre Riazanov¹, Alexandre Kouznetsov¹, Christopher JO Baker^{1*}

From Asia Pacific Bioinformatics Network (APBioNet) Ninth International Conference on Bioinformatics (InCoB2010)
Tokyo, Japan. 26-28 September 2010

A semantic assistant for mutation mentions in PubMed abstracts.



Jonas B Laurila¹, Alexandre Kouznetsov¹ and Christopher J O Baker^{*1}

AIMM-2010 @ ECCB 2010

Annotation, Interpretation and Management of Mutations 2010 Ghent, Belgium.

Nature Precedings : doi:10.1038/npre.2010.5443.1 : Posted 27 Dec 2010

Replacement of tryptophan residues... +

NCBI Resources ▾ How To ▾ My NCBI 9

PubMed.gov
U.S. National Library of Medicine
National Institutes of Health

Search: PubMed ▾ Limits Advanced search Help

Search Clear

Display Settings: ▾ Abstract

Send to: ▾



Eur J Biochem. 1995 Mar 1;228(2):403-7.

Replacement of tryptophan residues in haloalkane dehalogenase reduces halide binding and catalytic activity.

Kennes C, Pries F, Krooshof GH, Bokma E, Kingma J, Janssen DB.

Department of Biochemistry, Groningen Biomolecular Sciences and Biotechnology Institute, University of Groningen, The Netherlands.

Haloalkane dehalogenase catalyzes the hydrolytic cleavage of carbon-halogen bonds in short-chain haloalkanes. Two tryptophan residues of the enzyme (Trp125 and Trp175) form a halide-binding site in the active-site cavity, and were proposed to play a role in catalysis. The function of these residues was studied by replacing [Trp125 with phenylalanine](#), glutamine or arginine and Trp175 by glutamine using site-directed mutagenesis. All mutants showed a more than 10-fold reduced kcat and much higher Km values with 1,2-dichloroethane and 1,1,1-trichloroethane. Fluorescence quenching experiments showed a decrease in the affinity of the mutants for the substrate. The isotope effect observed with the wild-type enzyme in deuterium oxide was lost in the mutants. The results indicate that both tryptophans are involved in stabilizing the transition state of the substitution reaction that causes carbon-halogen bond cleavage.

Point Mutations

PointMutation

- hasmentionedposition - 125
- hascorrectposition - 125
- isgroundedto - Q6Q3H0
- haswildtyperesidue - W
- hasmutantresidue - F

PMID: 7705355 [PubMed - indexed for MEDLINE] **Free Article**

+ Publication Types, MeSH Terms, Substances, Secondary Source ID

+ LinkOut - more resources

Related citations

[Kinetic analysis and X-ray structure of haloalkane dehalogenase](#) [Biochemistry]

[Repositioning the catalytic triad aspartic acid and glutamic acid in the active site of haloalkane dehalogenase](#) [Biochemistry]

[Kinetic characterization and X-ray structure of a mutant of haloalkane dehalogenase](#) [Biochemistry]

[Improved catalytic properties of haloalkane dehalogenase by modification of the active site](#) [Biochemistry]

Review [Evolving haloalkane dehalogenase](#) [Curr Opin Chem Biol]

See re

S

Cited by 5 PubMed Central articles

[Functionally relevant motions of haloalkane dehalogenases occur in the active site](#) [Protein Science]

[Persistently conserved positions in structurally similar, sequence dissimilar haloalkane dehalogenases](#) [Protein Science]

[The importance of reactant positioning in the active site of haloalkane dehalogenase](#) [Protein Science]

Nature Precedings : doi:10.1038/npre.2010.5443.1 : Posted 27 Dec 2010

PubMed.gov
U.S. National Library of Medicine
National Institutes of Health

Search: PubMed

[Limits](#) [Advanced search](#) [Help](#)

[Display Settings:](#) Abstract

[Send to:](#)

Point Mutations

Biochemistry. 1999 Sep 14;38(37):12052-61.

Crystallographic and kinetic evidence of a collision complex formed during halide import in haloalkane dehalogenase.

Pikkemaat MG, Ridder IS, Rozeboom HJ, Kalk KH, Dijkstra BW, Janssen DB.

Laboratory of Biochemistry, BIOSON Research Institute, Groningen Biomolecular Sciences and Biotechnology Institute, University of Groningen, The Netherlands.

Haloalkane dehalogenase (DhIA) converts haloalkanes to their corresponding alcohols and halide ions. The rate-limiting step in the reaction of DhIA is the release of the halide ion. The kinetics of halide release have been analyzed by measuring halide binding with stopped-flow fluorescence experiments. At high halide concentrations, halide import occurs predominantly via the rapid formation of a weak initial collision complex, followed by transport of the ion to the active site. To obtain more insight in this collision complex, we determined the X-ray structure of DhIA in the presence of bromide and investigated the kinetics of mutants that were constructed on the basis of this structure. The X-ray structure revealed one bromide ion firmly bound in the active site and two bromide ions weakly bound on the surface of the enzyme. One of the weakly bound ions is close to Thr197 and Phe294, near the entrance of the earlier proposed tunnel for substrate import. Kinetic analysis of bromide import by the [Thr197Ala](#) and [Phe294Ala](#) mutants of DhIA at high halide concentration showed that the rate constants for halide binding no longer displayed a

PointMutation

hasmentionedposition - 294

hascorrectposition - 294

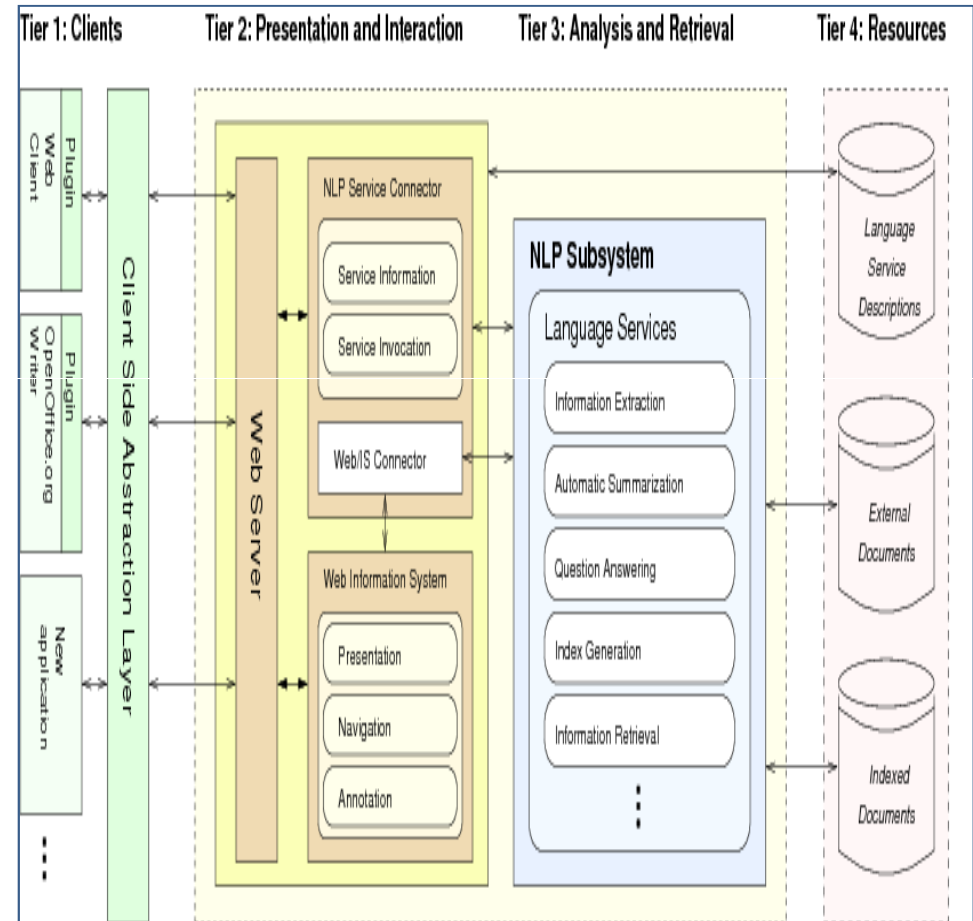
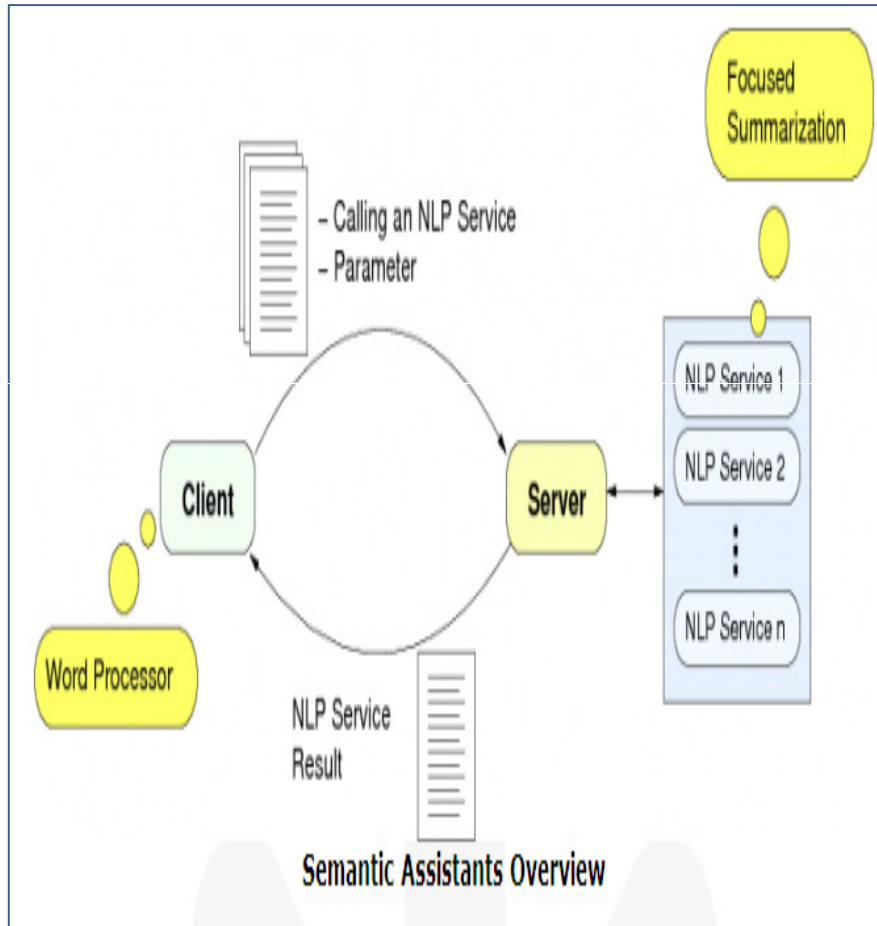
isgroundedto - Q6Q3H0

haswildtyperesidue - F

hasmutantresidue - A

increase with increasing bromide concentrations. This is in agreement with an elimination or a surface-located halide-binding site. Likewise, chloride binding kinetics of the mutants indicated with wild-type enzyme. The results indicate that Thr197 and Phe294 are involved in the formation of an for halide import in DhIA and provide experimental evidence for the role of the tunnel in substrate and

Semantic Assistant Framework



René Witte and Thomas Gitzinger.

[**Semantic Assistants – User-Centric Natural Language Processing Services for Desktop Clients.**](#)

3rd Asian Semantic Web Conference (ASWC 2008), February 2–5, 2009, Bangkok, Thailand.

Springer LNCS 5367, pp. 360–374. (Acceptance rate: 31%)



Mutation Grounding

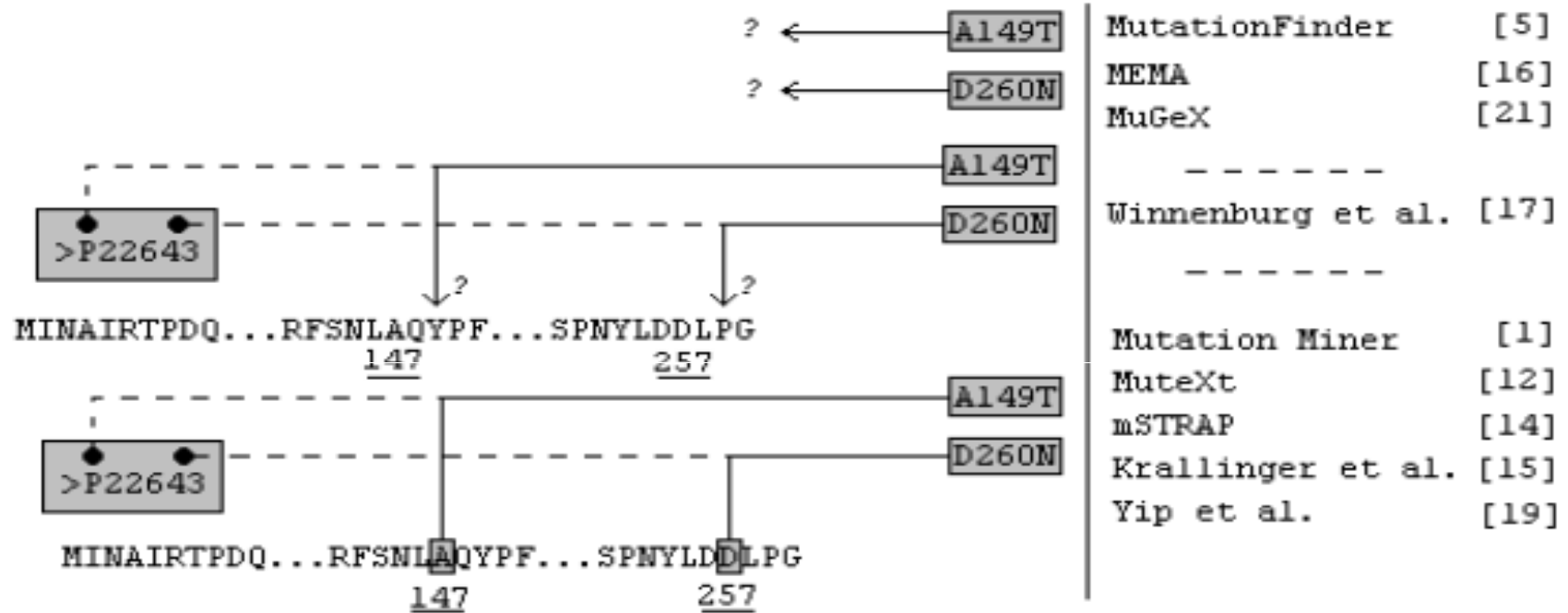
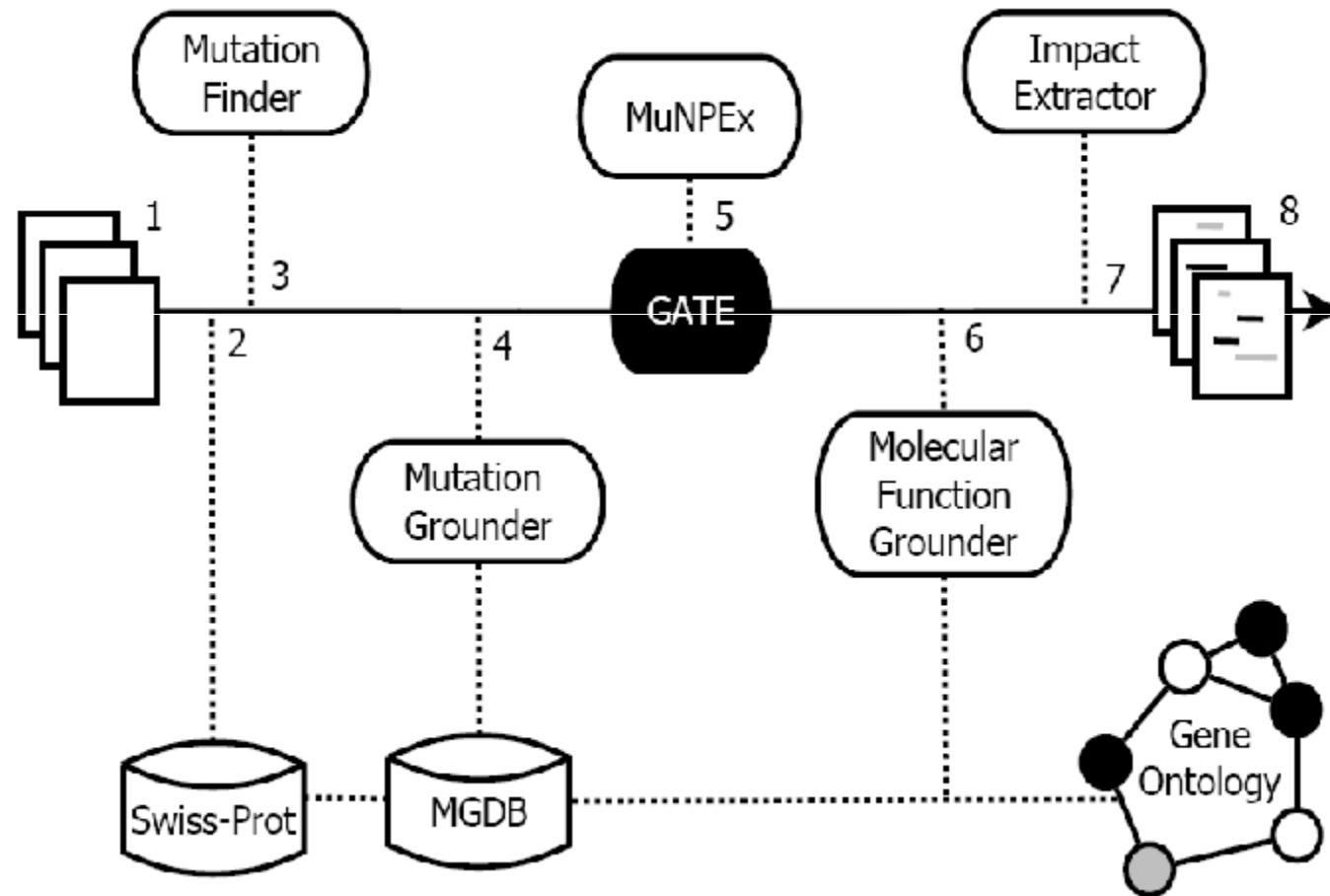


Fig. 1. Degrees of Mutation Grounding. *Uppermost*, mutation mentions are extracted but their relation to the appropriate protein sequence is not (no grounding). *Middle*, the related protein is found and mappings are established to the sequence if both text and database follow the same numbering scheme. *Lowermost*, the mutations are properly grounded, i.e. mapped to the correct position on the amino acid sequence of the related protein. Systems performing and not performing mutation grounding are displayed to the right.

Implementation A GATE pipeline





Sample sentences

“As expected, complete loss in activity of **W109L** and sustained activity of **F151W** were observed.”

“In order to further understand the catalytic mechanism we constructed an **Asp-124->Asn** mutant enzyme.”

“DhIA shows only a small decrease in activity when **Trp-125 is replaced with phenylalanine.**”

“The **W125F** mutant showed only a slight reduction of activity (Vmax) and a larger increase of Km with 1,2-dibromoethane.”

"Haloalkane dehalogenase (DhIA) from **Xanthobacter autotrophicus GJ10** hydrolyses terminally chlorinated and brominated n-alkanes to the corresponding alcohols."

Mutation	Description	Protein name	Gene name	Organism name
-----------------	--------------------	---------------------	------------------	----------------------



Sample sentences

“As expected, complete **loss** in **activity** of W109L and **sustained activity** of F151W were observed.”

“In order to further understand the catalytic mechanism we constructed an Asp-124->Asn mutant enzyme.”

“DhIA shows only a small **decrease** in **activity** when Trp-125 is replaced with phenylalanine.”

“The W125F mutant showed only a slight **reduction** of **activity** (V_{max}) and a larger **increase** of **K_m** with 1,2-dibromoethane.”

“Haloalkane dehalogenase (DhIA) from Xanthobacter autotrophicus GJ10 hydrolyses terminally chlorinated and brominated n-alkanes to the corresponding alcohols.”

Direction

Protein Property



Named entities of interest

- Mutations
- Proteins
- Genes
- Organisms
- Protein properties
 - Protein functions (*activity, binding etc*)
 - Kinetic variables (*K_m, k_{cat} etc*)
 - (stability)

Mutation Grounding

1. Retrieve and normalize mutations
2. For each candidate sequence
 1. For each pair of mutations
 1. Make regexp $w_1.(N_2-N_1)w_2$
 2. Match regexp to sequence
 3. Check remaining residues at corrected positions.
3. Ground proteins and mutations to the same AC / sequence

Mutation Grounding: Example

Candidate sequences

```
1 MGAKACYGAKCVAVAIIVAGASSESLGKEQY
2 MAPEALFDRKYTYGKKVWSFGVLLWEIFTL
3 MQVSLESYGSMSSNTPLVRIARLSSGEGPT
```

Candidate mutations

```
K5Q
Y8C
G9C
K11E
V35Q
```

Mutation Grounding: Example

Candidate sequences

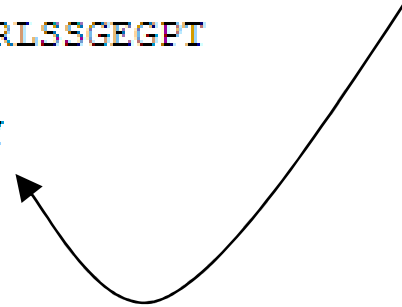
```
1 MGA KACYGAKCVAVAI VAGASSESLGKEQY
2 MAPEALFDRKYTYGKKVNSFGVLLWEIFTL
3 MQVSLESYGSMSSNTPLVRIARLSSGEGPT
```

Candidate mutations

```
{ K5Q
  Y8C
  G9C
  K11E
  V35Q
```

K..Y

Compute regular expression



Mutation Grounding: Example

Candidate sequences

1 MGA**K****A**Y**E**AKCVAVAIVAGASSESLG**K****E****C****Y**
2 MAPEALFDRKYTYGKKVWSFGVLLWEIFTL
3 MQVSLESYGSMSSNTPLVRIARLSSGEGPT

K..Y

Candidate mutations

{ K5Q
Y8C
G9C
K11E
V35Q

Match with sequence 1

Mutation Grounding: Example

Candidate sequences

1 MGA**K**AC**Y**E**K**CVAVAIVAGASSESLGKEQY
2 MAPEALFDRKYTYGKKVNSFGVLLWEIFTL
3 MQVSLESYGSMSNTPLVRIARLSSGEGPT

Candidate mutations

~~K5Q~~
Y8C
G9C
K11E
V35Q

Extend match A

Mutation Grounding: Example

Candidate sequences

```
1 MGAKACYGAKCVAVAIVAGASSESLGKKEQYXX  
2 MAPEALFDRKYTYGKKVNSFGVLLWEIFTL  
3 MQVSLESYGSMSSNTPLVRIARLSSGEGPT
```

Candidate mutations

```
K5Q  
Y8C  
G9C  
K11E  
V35Q
```

Extend match B

Mutation Grounding: Example

Candidate sequences

```
1 MGA[K]A[Q]Y[Q]K[V]VAVAI VAGASSESLGKEQY
2 MAPEALFDRK[Y]I[Y]K[W]WSFGVLLWEIFTL
3 MQVSLES[Y]E[S]MSSNTPLVRIARLSSGEGPT
```

Candidate mutations

```
K5Q
Y8C
G9C
K11E
V35Q
```

Match with sequence 2 & 3

Mutation Grounding: Example

Candidate sequences

Mutations / Offset

1	MGAKAC Y GA K CVAVVAIVAGASSESLGKEQY	4 / -1
2	MAPEALFDR K Y I Y G K KWWSFGVLLWEIFTL	4 / 5
3	MQVSLES Y E SMSSENTPLVRIARLSSGEGPT	2 / 0

Choose best candidate sequence:

1. Most grounded mutations
2. Least absolute offset

Evaluation

- COSMIC
 - Catalogue Of Somatic Mutations In Cancer
 - PIK3CA, FGFR3, MEN1
 - 63 documents
- Haloalkane dehalogenases
 - Protein engineering literature
 - 13 documents



Evaluation

Corpus	Precision	Recall	Corpus size
PIK3CA	0.86	0.70	30
FGFR3	0.89	0.66	26
MEN1	0.54	0.32	7
Haloalkane Dehalogenases	0.83	0.73	13
Average	0.84	0.65	76



Grounding of named entities

- Protein grounding
 - Assign the correct UniProt id to each detected protein entity.
- Mutation grounding
 - Verify and, if necessary, correct each mutation location to match its corresponding protein's sequence as obtained from UniProt.
- Protein function grounding
 - Assign the correct gene ontology id to detected protein functions



Protein & mutation grounding

- Combined into one method
 1. A pool of accession numbers is created based on occurrence of protein and gene names
 2. Mutations are matched to candidate sequences, going from min to max amount of mutations.
 3. Sequence with most grounded mutations is considered correct for the entire paper



Protein function grounding

1. Retrieve go:mf concepts related to previously grounded proteins
2. ground noun phrases, with *activity*, *binding*, *affinity* or *specificity* as head nouns, to retrieved go:mf concepts.
3. score them according to lexical similarity with the retrieved go:mf concepts.
4. use scores to solve contradictions in output protein function grounding and impact information

Protein function grounding

1. Retrieve go:mf concepts related to previously grounded proteins
2. ground noun phrases, with *activity*, *binding*, *affinity* or *specificity* as head nouns, to retrieved go:mf concepts.
3. score them according to **lexical similarity** with the retrieved go:mf concepts.
4. use scores to solve contradictions in output protein function grounding and impact information

Protein function grounding Example

Grounded protein

(P22643) Haloalkane dehalogenase

Related gene ontology molecular function concepts

(GO:0018786) haloalkane dehalogenase activity

Found noun phrases

anhydrase activity

dehalogenase activity

activity

First stem!

Protein function grounding Example

$$\text{compare} - \text{similarity} = \frac{|N \cap G|^2}{|N| |G|}$$

haloalk dehalogen activ

anhydr activ $s = 1 / (2 * 3) = 1/6$

dehalogen activ $s = 4 / (2 * 3) = 4/6$

activ $s = 1 / (1 * 3) = 2/6$

Protein function grounding

Example:

ground to a certain degree

(GO:0018786) haloalk dehalogen activ

anhydr activ	→((GO:0018786) , 0.17)	(3)
dehalogen activ	→((GO:0018786) , 0.67)	(1)
activ	→((GO:0018786) , 0.33)	(2)



Impact direction term lists

Positive	Negative	(cont.)	Neutral	Negation	Non-Neutral
increase	abolish	loose	identical	without	affect
-increases	decrease	defect	similar	no	effect
-increased	reduce	disrupt	full	not	alter
-increasing	lower	diminish			differ
enhance	inhibit				
higher	impair				
improve					

Relation detection

1. Impacts

1. Directionality+Protein property
2. «**sustained activity** of F151W were observed»

2. Mutants

1. Set of grounded mutations

3. Mutant+Impact

- Relations found in text by the use of rules



Impact rules - an example

#negative-impact-rule

If (**Sentence contains Protein Function and Sentence contains Negative Direction and Sentence not contains Positive Direction**)



Markup Sentence as Negative Impact on Protein Function

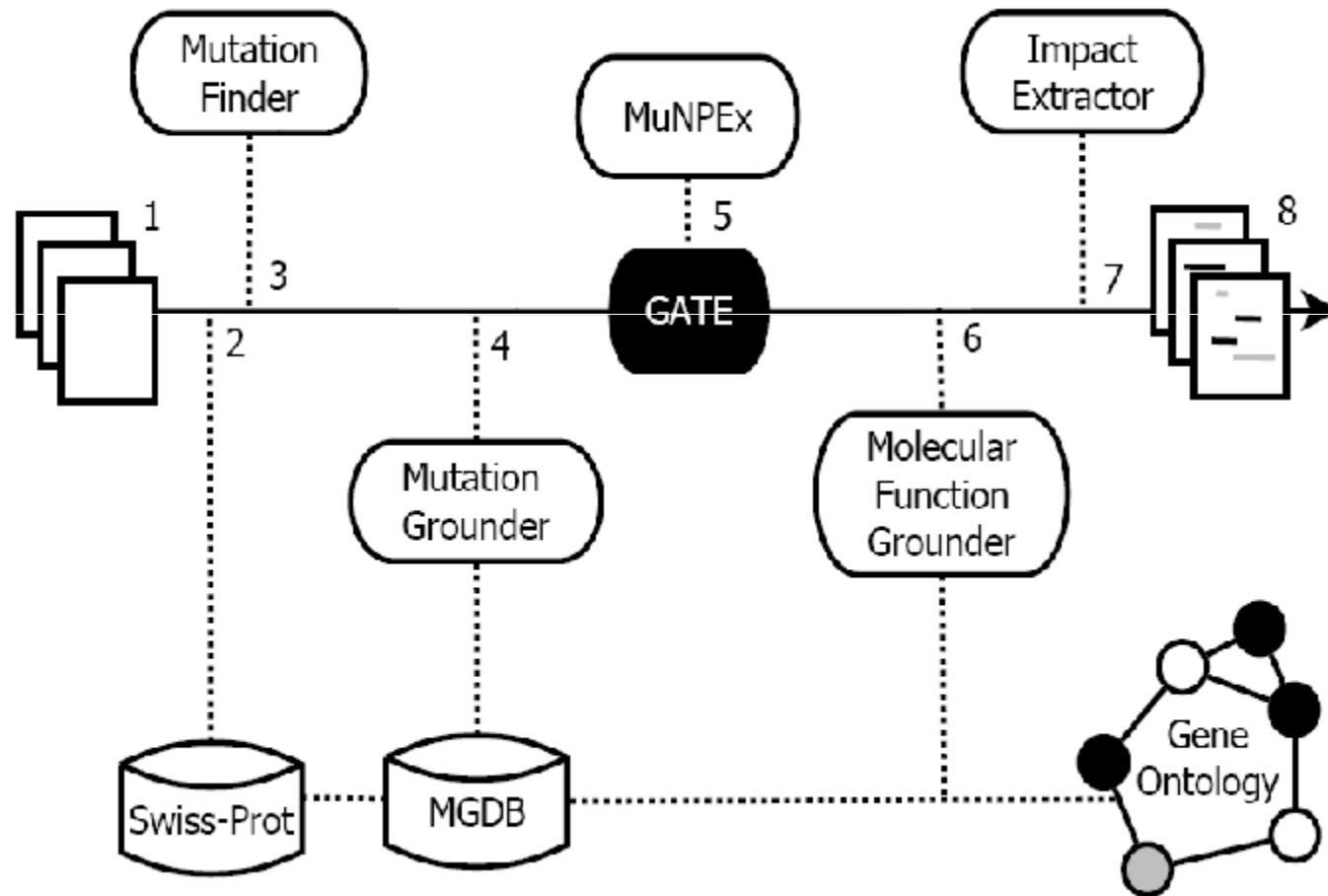


DhIA shows only a small **decrease** in **activity** when Trp-125 is replaced with phenylalanine.



A larger **increase** of **Km** with 1,2-dibromoethane.

Implementation: a GATE pipeline



Access to mutation information

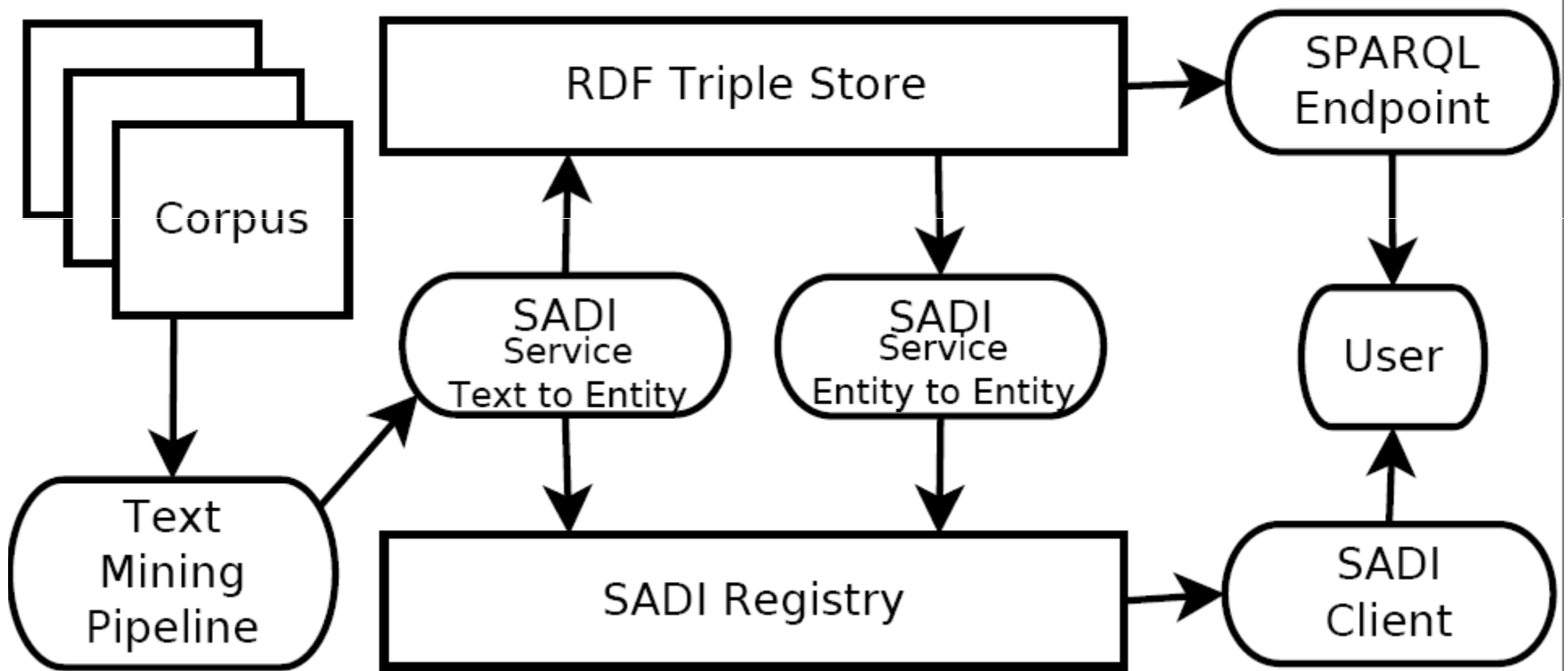
1. Export mutation impact information from text to an RDF triple store
2. Provide a SPARQL endpoint as a query interface to the triple store
3. Make both the pipeline and triple store available through semantic web services (SADI)



Example queries

- Retrieve all reported mutations and their impacts on haloalkane dehalogenase.
- Find all mutants with an increased Ca^{2+} affinity.
- Has the F137T mutation of carbonic anhydrase II previously been studied?

Access to mutation information - summary





The NAR Subset

PubMed

PubMed Central

PubMed Central Open Access subset

Nucleic Acids Research

**Subset where *MutationFinder*
recognizes more than one
point mutation →
Σ 1146 full text articles**



Statistics of NAR subset

- Mutations could be grounded in 733 of the 1146 relevant documents by our system.
 - A total of 4008 *GPM-occurs-in-Document* triples
 - 2977 *unique/distinct GPMs*
 - 759 *Protein-occurs-in-Document* triples
 - 995 *Impacts* (191 positive, 120 neutral, 684 negative)
 - 110 *unique/distinct* Gene Ontology Terms...were extracted from these documents.

Impact Extraction for Reuse in SNP Prediction

#	Mut	SNAP _{MAT} Pred	Reference	Func Change/ Impact	GO ID	Agreement	GO Descrip
	R165W		34	positive	GO_0004016		adenylate cyclase activity
17	F51L	Non-neutral	16	negative	GO_0004016	Y	adenylate cyclase activity
36	I102T	Neutral	16	negative	GO_0004016	N	adenylate cyclase activity
70	A154D	Non-neutral	16	negative	GO_0004016	Y	adenylate cyclase activity
	S127L		43	positive	GO_0004144		diacylglycerol O-acyltransferase activity
	T11S		13	negative	GO_0004144		diacylglycerol O-acyltransferase activity
	R18C		13	negative	GO_0004144		diacylglycerol O-acyltransferase activity
	F121R		19	negative	GO_0004977		melanocortin receptor activity
44	N123A	Non-neutral	38	negative	GO_0004978	Y	adrenocorticotropin receptor activity
45	N123D	Non-neutral	38	negative	GO_0004978	Y	adrenocorticotropin receptor activity
	D126A		38	neutral	GO_0004978		adrenocorticotropin receptor activity
156	I316T	Non-neutral	33	negative	GO_0005000	Y	Vasopressin receptor activity
15	E49A	Non-neutral	29	negative	GO_0005184	Y	neuropeptide hormone activity
26	D90A	Non-neutral	29	negative	GO_0005184	Y	neuropeptide hormone activity
	E100A		29	negative	GO_0005184		neuropeptide hormone activity
	E100C		1	negative	GO_0005515		protein binding

In silico mutagenesis: a case study of the melanocortin 4 receptor.

[Bromberg Y](#), [Overton J](#), [Vaisse C](#), [Leibel RL](#), [Rost B](#) – FASEB 2009



Impact Extraction: Deployed as Semantic Web Service

SADI, SHARE, and the *in silico Scientific Method* - Mark D Wilkinson^{1§*}, Luke McCarthy^{1*}, Benjamin Vandervalk^{1*}, David Withers^{1*}, Edward Kawas^{1*}, Soroush Samadian^{1*} BMC Bioinformatics (2010)

Deploying the Mutation Impact mining pipeline with SADI: an exploratory case study, Alexandre Riazanov, Jonas Bergman Laurila and Christopher J O Baker

Proceedings of the Workshop on Annotation, Interpretation and Management of Mutations ([AIMM-2010](#))
Annotation, Interpretation and Management of Mutations. A workshop at [ECCB10](#).
Ghent, Belgium, September 26th, 2010.

AIMM2010

Annotation, Interpretation and Management of Mutations.
A workshop at ECCB10.

HOME

SCHEDULE

WORKSHOP THEMES

KEYNOTE SPEAKERS

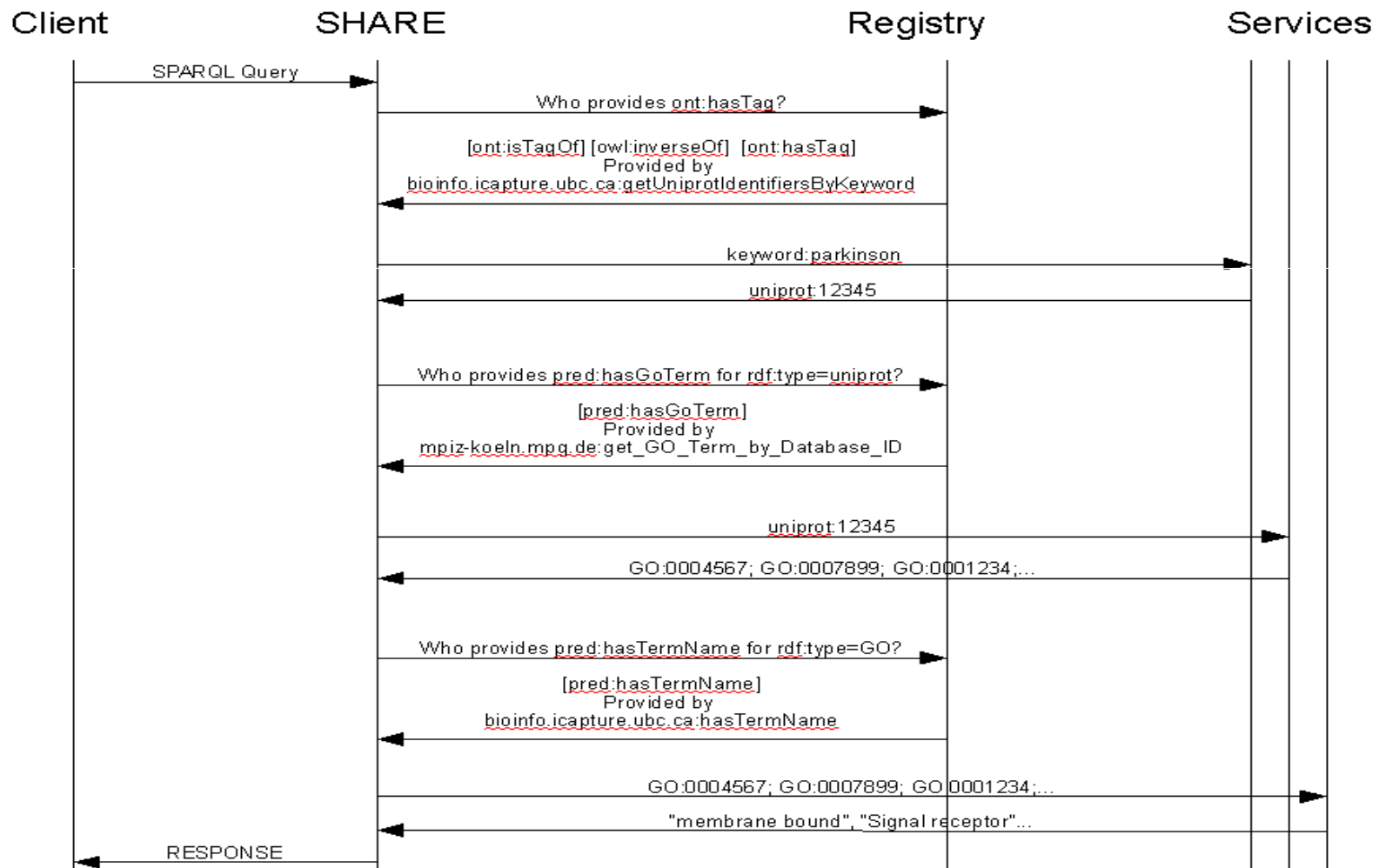
AIMM publications

2008 CEUR: <http://sunsite.informatik.rwth-aachen.de/Publications/CEUR-WS/Vol-429/>

2008 BMC special issue: <http://www.biomedcentral.com/bmcbioinformatics/10?issue=98>

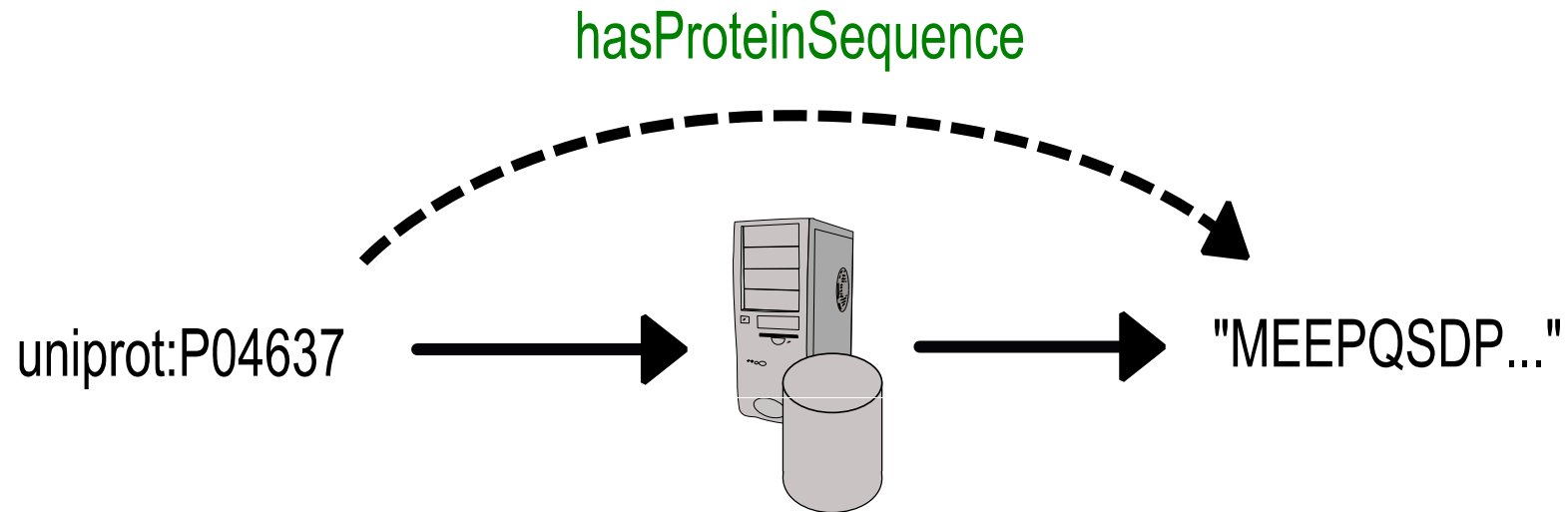
2010 CEUR: <http://sunsite.informatik.rwth-aachen.de/Publications/CEUR-WS/Vol-645/>

Semantic Health And Research Environment (SHARE) prototype.



Slide by
Michel
Dumontier

hasProteinSequence



Input Data:	BRCA1	rdf:type	Gene ID
Output Data:	BRCA1	hasDNASequence	AGCTTAGCCA...
Registry Index:	Service provides "hasDNASequence" property to Gene IDs		

Predicate-based web service invocation. Using the `hasProteinSequence` predicate in a query automatically invokes a web service capable of obtaining the amino acid sequence for UniProt entry P04637.

Nature Precedings : doi:10.1038/npre.2010.5443.1 : Posted 27 Dec 2010

SADI registry — registered services - Windows Internet Explorer

http://sadiframework.org/registry/services/

SADI framework predicate map

Register Services

Registered services <http://sadiframework.org/>

Service URL	Input Class	Output
- http://sadiframework.org/services/getGOTerm Name <code>getGOTerm</code> Description gets the text-label for a GO Term Properties attached http://sadiframework.org/ontologies/predicates.owl#hasTermName (with values from http://www.w3.org/2000/01/rdf-schema#Literal)	GO Record	getGO
+ http://sadiframework.org/examples/uniprot2go	UniProt Record	Annotate
+ http://sadiframework.org/examples/uniprot2pdb	UniProt Record	Annotate
+ http://dev.biordf.net/~kawas/cqi-bin/getGOTermDefinitions	GO Record	getGO
+ http://dev.biordf.net/~kawas/cqi-bin/getGeneInformation	GeneID Record	getGene
+ http://sadiframework.org/services/getDrugBankByUniProt	UniProt Record	getDrug
+ http://sadiframework.org/services/getKEGGIDFromUniProt	UniProt Record	getKEGG
+ http://dev.biordf.net/~kawas/cqi-bin/getKeqgPathwaysByEnzyme	ENZYME Record	getKeqg
+ http://sadiframework.org/services/getKEGGPathwayDiagram	KEGG PATHWAY Record	getKEGG
+ http://sadiframework.org/services/getDrugNames	DrugBank Record	getDrug

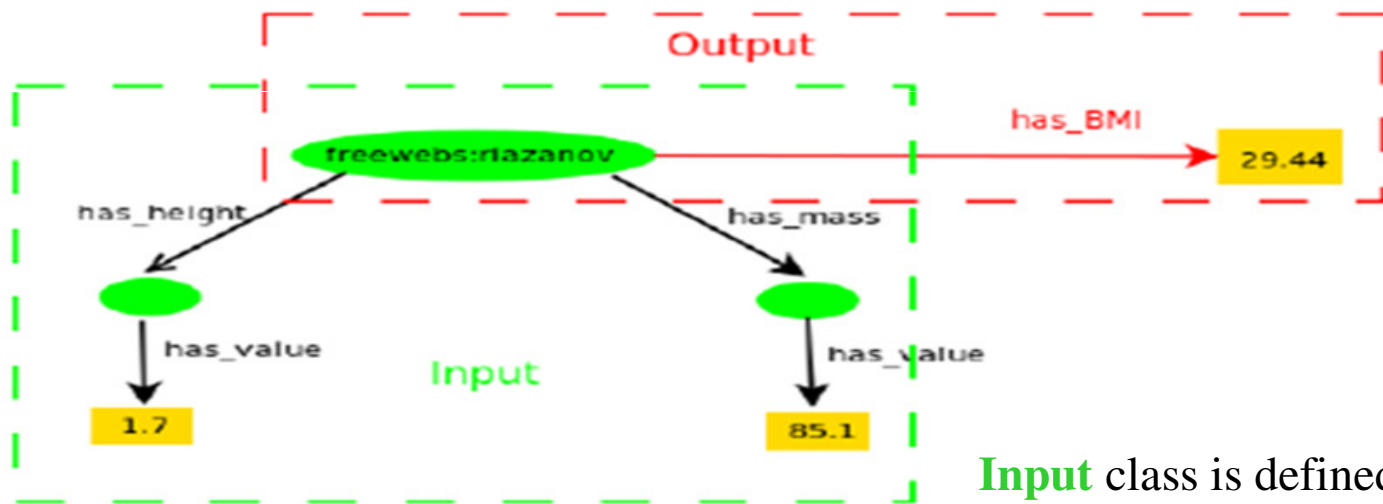
Internet 100%

start 3 Internet Explorer C-BRASS-Atlantic Se... Alex-presentation wo... C-BRASS NEP2 SoW-... 10:36 PM

Input / Output Semantics: Example

OWL class of the main input node specifies the service output.

- Uses existential property restrictions to define the predicates of what the service does.
- Service computing BMI has this output class: (has_BMI exactly 1 xsd:float) advertising that it computes the Body Mass Index.
- Service annotates its input with the predicate has BMI with typical **output:** `<http://www.freewebs.com/riazanov> has_BMI float"29.44"`



Input class is defined as by

(has_mass some (has_value some xsd:float))

(has_height some (has_value some xsd:float))

the service expects something like this in the input:

`<http://www.freewebs.com/riazanov> has_height [has_value float"1.7"] .`

`<http://www.freewebs.com/riazanov> has_mass [has_value float"85.1"] .`



Mutation Impact Pipeline

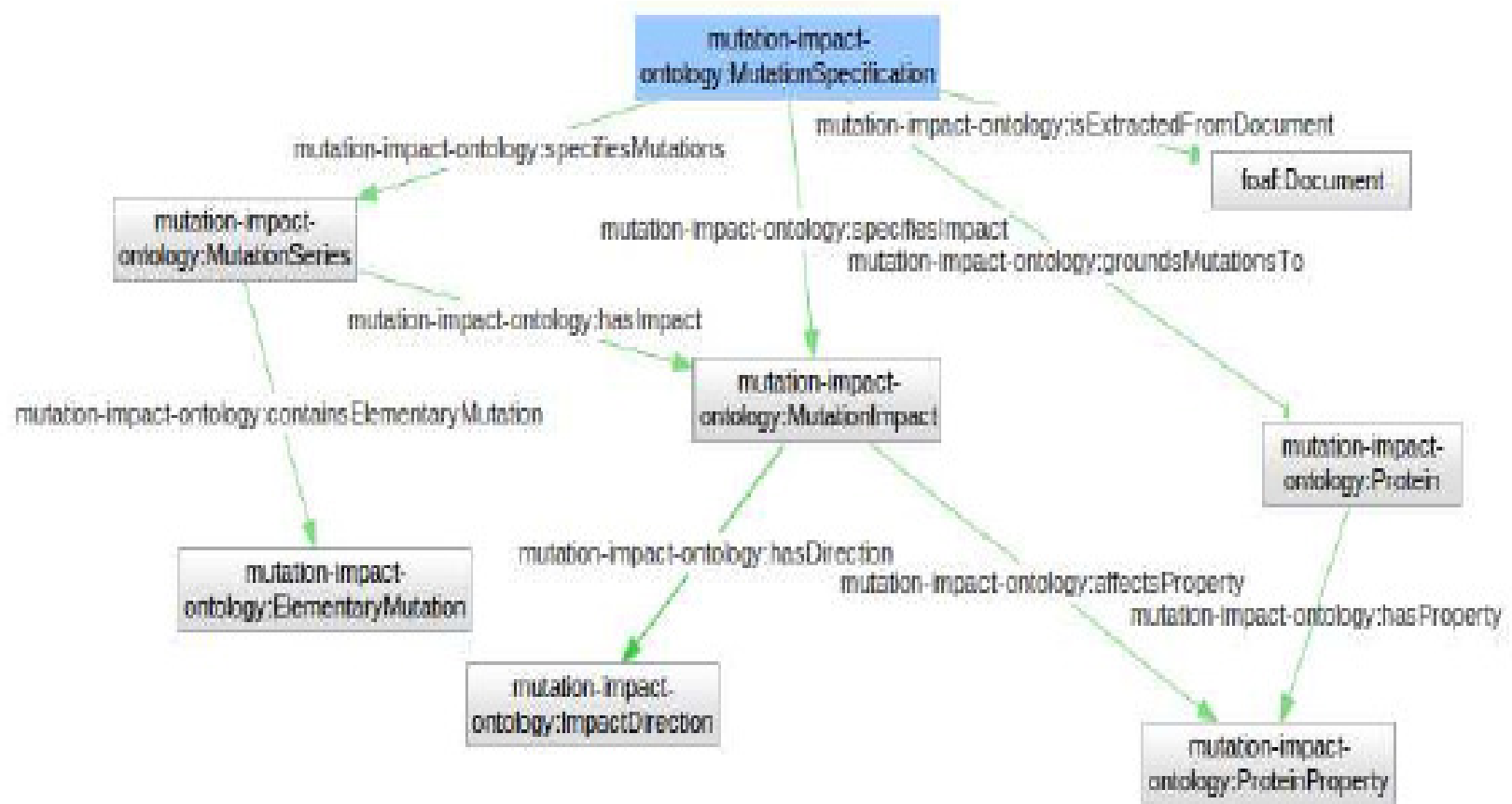
The haloalkane dehalogenase from the nitrogen-fixing hydrogen bacterium *Xanthobacter autotrophicus* GJ10 (DhlA) prefers 1,2-dichloroethane (DCE) as substrate and converts it to 2-chloroethanol and chloride ... DhlA shows only a small decrease in activity when Trp-125 is replaced with phenylalanine

- “haloalkane dehalogenase” is a protein
- its UniProt Id is P22643
- “Trp-125 is replaced with phenylalanine” is the point mutation W125F
- “activity” is the protein property being affected, GO_00188786 in Gene Ontology
- “decrease” means the impact is negative

UNB-Concordia joint work. GATE pipeline wrapped as a Java library. Outputs Java objects representing mutation specifications: statements about mutations and their impacts on protein properties.



Mutation Impact Ontology





Text -> Mutation Specifications

- Just wraps the Mutation Impact mining pipeline as a SADI service
- Identifies input text in the service input, submits it to the pipeline and converts the results to an RDF graph

Input:

```
<http://example.com/text1>    rss:link    anyURI"http://example.com/text1"
```

Output:

```
<http://example.com/text1>    foaf:topic                midb:mutationSpec243
miodb:mutationSpec243         mio:groundsMutationsTo    miodb:protein528
miodb:protein528              mio:hasSwissProtId        "P22643"
miodb:mutationSpec243         mio:specifiesImpact       miodb:mutationImpact624
miodb:mutationImpact624      mio:hasDirection          mio:Positive
miodb:mutationImpact624      mio:affectsProperty       miodb:proteinProperty326
miodb:proteinProperty326     mio:hasType                GO:GO_0018786
```

.



Services: Mutation Impact DB

- Wildtype protein → mutation specifications (complete description of)
- Mutant protein → mutation specifications
- Specific property of a specific protein → known mutation impacts
- Mutation impact (on a specific property of a specific protein) → mutation specifications
- Bio entity type (e.g., `mio:Protein` or `mio:PointMutation`) → known instances
- Set of elementary mutations → subsets described in the literature (with links to the corresponding `mio:MutationSpecification`)



Use Case Queries

Samples use cases queries are online:

http://unbsj.biordf.net/mutation-impact/queries_Dec9.html

and can be run at

<http://138.119.1.172:8080/cardioSHARE-mutations/>

Use Case 1

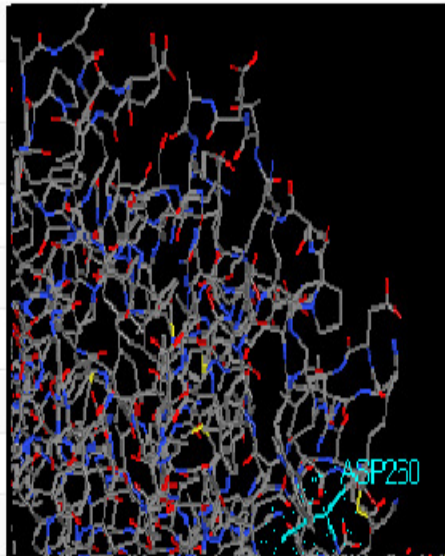
- A protein engineer is looking for mutations that can improve catalytic activity of an enzyme
- Query: find all mutations and the structure images of wild type proteins that were mutated, where the impact of the mutation is an enhanced haloalkane dehalogenase activity (GO_0018786)
- Predicates from our ontology take us from GO_0018786 to mutations and proteins: *proteinPropertyHasType + affectsProperty + hasDirection + specifiesImpact + containsElementaryMutation + groundMutationsTo*
- Two external SADI services provide *has3DStructure* to link proteins to their structure descriptions, and *hasJmol3DStructureVisualization* to link the structure to a 3D image file

Use Case 1 Query

```

PREFIX mio:<http://unbsj.biordf.net/ontologies/mutation-impact-ontology.owl#>
PREFIX mioe:<http://unbsj.biordf.net/ontologies/mutation-impact-ontology-extras.owl#>
PREFIX go:<http://purl.org/obo/owl/GO#>
PREFIX sio:<http://semanticscience.org/resource/>
PREFIX props:<http://sadiframework.org/ontologies/properties.owl#>
PREFIX objects:<http://sadiframework.org/ontologies/service_objects.owl#>
SELECT DISTINCT ?NormalizedMutation ?Protein ?Visualisation
FROM <http://unbsj.biordf.net/mutation-impact/service-data/protein_property_types.rdf>
WHERE { # impact <-- property instance ?Impact mio:affectProperty
?Property .
# protein property instance <-- GO_0018786 ?Property
mioe:proteinPropertyHasType go:GO_0018786 .
# check that the impact is positive
?Impact mio:hasDirection mio:Positive .
# grounded mutation
<-- impact ?MutationSpecification mio:specifiesImpact ?Impact .
# grounded mutation --> wildtype protein ?MutationSpecification
mio:groundMutationsTo ?Protein .
# grounded mutation --> point mutation series
?MutationSeries mio:mutationSeriesIsSpecifiedBy
?MutationSpecification .
# point mutation series --> separate point mutations
?MutationSeries mio:containsElementaryMutation ?Mutation .
?Mutation mio:hasNormalizedForm ?NormalizedMutation .
# grounded mutation
--> Web page with Jmol applet call ?MutationSpecification
objects:hasJmol3DStructureVisualization ?Visualisation . }

```

Normal	Protein	Visualisation
D260N	http://biordf.net/moby/UniProt/P22643	http://unbsj.biordf.net/mut-vis/visualiseMut
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	
D260N	http://biordf.net/moby/UniProt/P22643	

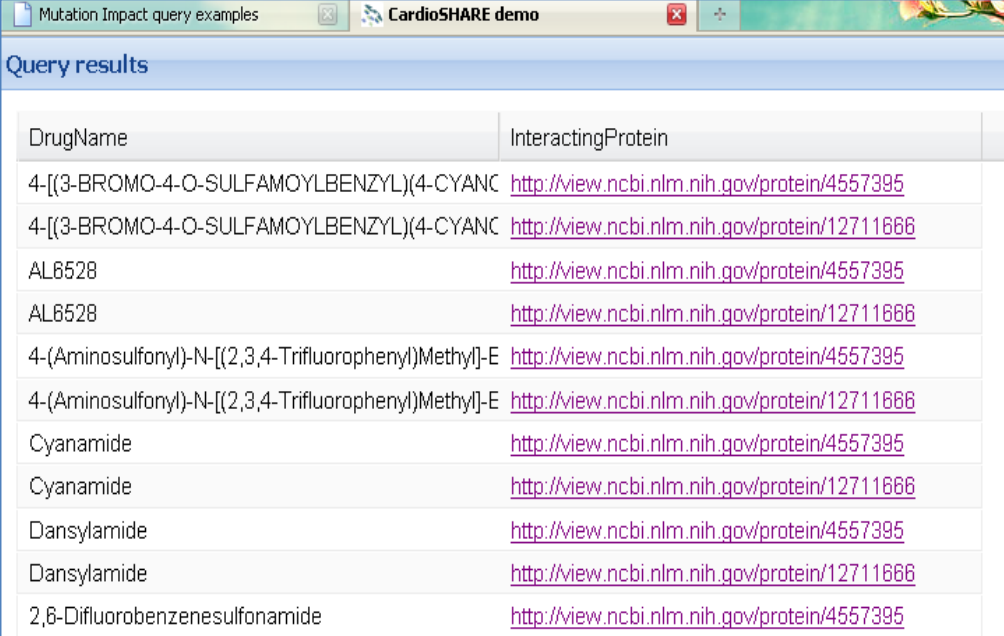


Use Case 3

- A researcher in drug discovery is looking for existing drugs targeting a new disease condition
- Query: find all drugs related to mutated proteins, together with their interaction partners, where the mutation impact is a increased carbonic anhydrase activity (GO_0008270)
- Our predicates link GO_0008270 to proteins via the instances of this protein properties, positive impacts and mutation specifications
- External ontologies facilitate the linking of proteins to the IDs of drugs affecting them
- Another external predicate, *hasMolecularInteractionWith*, links the proteins to proteins they interact with

Use Case 3 Query

```
PREFIX mio:<http://unbsj.biordf.net/ontologies/mutation-impact-ontology.owl#>
PREFIX mioe:<http://unbsj.biordf.net/ontologies/mutation-impact-ontology-extras.owl#>
PREFIX go:<http://purl.org/obo/owl/GO#>
PREFIX objects:<http://sadiframework.org/ontologies/service_objects.owl#>
PREFIX pred: <http://sadiframework.org/ontologies/predicates.owl#>
SELECT ?DrugName ?InteractingProtein
FROM <http://unbsj.biordf.net/mutation-impact/service-data/protein_property_types.rdf>
WHERE {# enumerate known instances of go:GO_0008270
?Property mioe:proteinPropertyHasType go:GO_0008270 .
# impact <-- protein property instance ?Impact mio:affectProperty
?Property .
# check that the impact is positive ?Impact mio:hasDirection
mio:Positive .
# grounded mutation <-- impact ?MutationSpecification
mio:specifiesImpact ?Impact .
# grounded mutation --> wildtype protein
?MutationSpecification mio:groundMutationsTo ?Protein .
# wildtype protein --> drug ?Protein objects:isTargetOfDrug
?Drug . ?Drug objects:hasDrugGenericName ?DrugName .
# wildtype protein --> interacting proteins ?Protein
pred:hasMolecularInteractionWith ?InteractingProtein }
```



DrugName	InteractingProtein
4-[(3-BROMO-4-O-SULFAMOYL)BENZYL](4-CYANC	http://view.ncbi.nlm.nih.gov/protein/4557395
4-[(3-BROMO-4-O-SULFAMOYL)BENZYL](4-CYANC	http://view.ncbi.nlm.nih.gov/protein/12711866
AL6528	http://view.ncbi.nlm.nih.gov/protein/4557395
AL6528	http://view.ncbi.nlm.nih.gov/protein/12711866
4-(Aminosulfonyl)-N-[(2,3,4-Trifluorophenyl)Methyl]-E	http://view.ncbi.nlm.nih.gov/protein/4557395
4-(Aminosulfonyl)-N-[(2,3,4-Trifluorophenyl)Methyl]-E	http://view.ncbi.nlm.nih.gov/protein/12711866
Cyanamide	http://view.ncbi.nlm.nih.gov/protein/4557395
Cyanamide	http://view.ncbi.nlm.nih.gov/protein/12711866
Dansylamide	http://view.ncbi.nlm.nih.gov/protein/4557395
Dansylamide	http://view.ncbi.nlm.nih.gov/protein/12711866
2,6-Difluorobenzenesulfonamide	http://view.ncbi.nlm.nih.gov/protein/4557395



Use Case 4

In this use case, a genomics researcher asks for all known mutations reported in the literature for a protein containing a non-synonymous SNP. Here the researcher is primarily looking for any literature describing impacts of a nsSNP on a protein. By retrieving all known mutations for the protein in which the nsSNP is reported, the researcher can find out if any of these reported mutations corresponds to the location of the SNP in question. Minimally the researcher can retrieve the full set of mutations to the protein based on reported experimental analysis and their impacts, together with references to the supporting literature.

Use Case 4 Query

Query: Find all documented mutations of the protein with SNP rs2305178.

PREFIX mio:<http://unbsj.biordf.net/ontologies/mutation-impact-ontology.owl#>

PREFIX mioe:<http://unbsj.biordf.net/ontologies/mutation-impact-ontology-extras.owl#>

PREFIX dbsnp:<http://lsrn.org/dbSNP:>

PREFIX objects:<http://sadiframework.org/ontologies/service_objects.owl#>

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>

PREFIX foaf: <http://xmlns.com/foaf/0.1/>

PREFIX rss: <http://purl.org/rss/1.0/>

PREFIX pred: <http://sadiframework.org/ontologies/predicates.owl#>

PREFIX sio: <http://semanticscience.org/resource/>

SELECT DISTINCT ?NormalizedMutation ?DocumentURL

WHERE { # SNP --> gene (Entrez)

'is variant of' dbsnp:rs2305178 sio:SIO_000272 ?EzGene .

enumerate known proteins ?Protein mioe:biologicalEntityHasType
mio:Protein .

proteins --> genes (KEGG) ?Protein pred:isEncodedBy?KeggGene

gene (KEGG) --> reference sequence ?KeggGene

objects:hasRefSeqTranscript ?RefSeq .

reference sequence --> gene (Entrez) ?RefSeq

objects:correspondsToEntrezGene ?EzGene .

protein --> mutation info ?MutationSpecification

mio:groundMutationsTo ?Protein . ?MutationSeries

mio:mutationSeriesIsSpecifiedBy ?MutationSpecification .

?MutationSeries mio:containsElementaryMutation ?Mutation .

?Mutation mio:hasNormalizedForm ?NormalizedMutation .

mutation --> literature reference ?Document foaf:topic

?MutationSpecification . ?Document rss:link ?DocumentURL }

Query results

DocumentURL	NormalizedMutation
http://www.freewebs.com/riazanov/15880580.pdf.txt	G697C

Int. J. Cancer: 117, 166168 (2005) ' 2005 Wiley-Liss, Inc.

SHORT REPORT Constitutive activating mutation of the FGFR3b in oral squamous cell carcinomas
Yan Zhang¹ , Yoshiko Hiraishi¹ , Hua Wang¹ , Ken-saku Sumi¹ , Yasutaka Hayashido¹ , Shigeaki Toratani¹ ,



Summary

- Rule based mutation and impact extraction
- Methods for grounding of mutations to protein sequences and protein functions to Gene Ontology
- Algorithms deployed with Semantic Assistant
- Algorithms and mutation information exposed with semantic metadata and as semantic web services (SWS)
- Reuse of SWS Mutation services with multiple use cases

Acknowledgements

Jonas B. Laurila

Alexandre Riazanov

Alexandre Kouznetsov

*** Christopher J.O. Baker**

Computer Science & Applied Statistics

University of New Brunswick

Saint John, Canada

Nona Naderi

René Witte

*Computer Science &
Software Engineering*

Concordia University

Montreal, Canada

This research was funded in part by:

- New Brunswick Innovation Foundation, New Brunswick, Canada
- NSERC, Discovery Grant, Canada
- Quebec – New Brunswick Co-operation in Advanced Education – Research Program, Government of New Brunswick, Canada

Acknowledgements

Wee Tiong Ang

Justin Choo

Rajaraman Kanagasabai

Institute for Infocomm Research

A-STAR, Singapore

Luke McCarthy

Mark Wilkinson

University of British Columbia

Vancouver, Canada

This research was also funded in part by:

- Agency for Science, Technology and Research, Singapore
- CANADA'S Advanced Research And Innovation Network, CANARIE, Canada