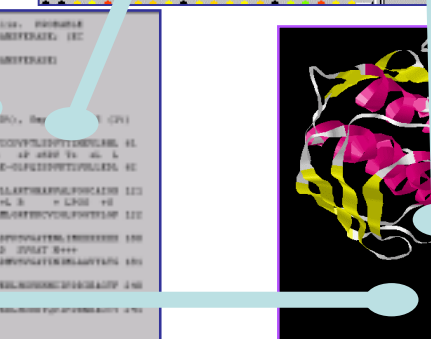
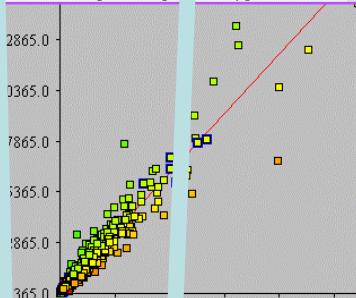
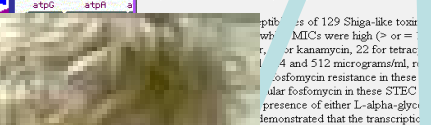
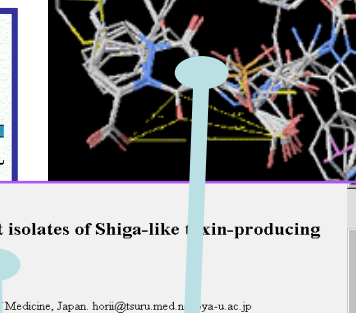
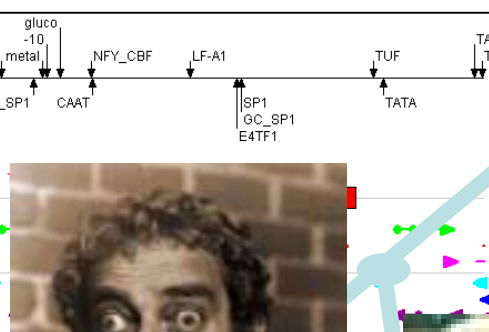


Treating “shimantic web” syndrome with ontologies

AKT workshop on Semantic Web Services

**Duncan Hull, Robert Stevens, Phillip Lord,
Chris Wroe and Carole Goble**

University of Manchester



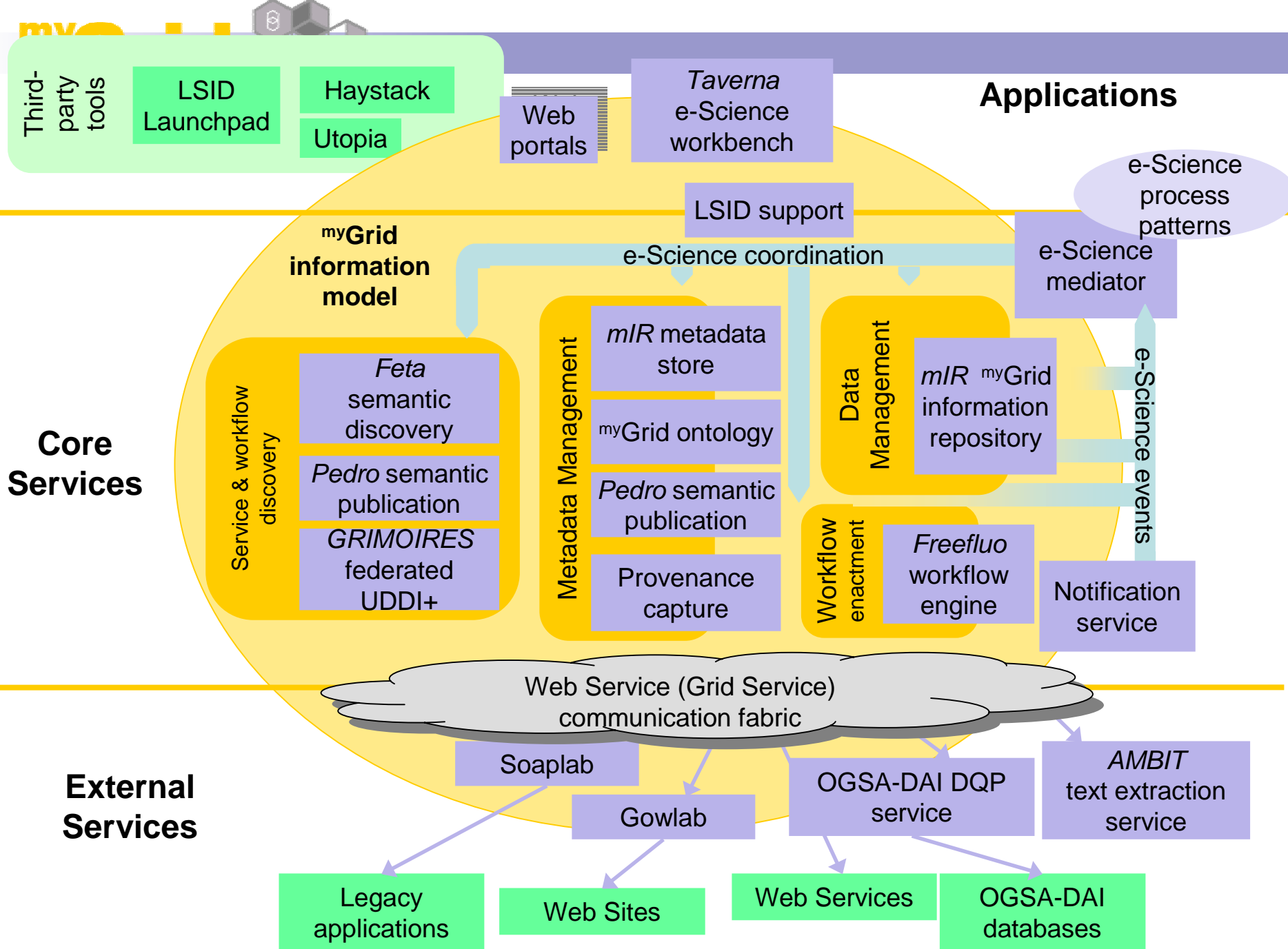
[54]	PURIFIED HOMOGENEOUS UDP-GLCNAC (GALNAc) PYROPHOSPHORYLASE	[56]	References Cited
[75]	Inventor: Alan D. Elbein , Little Rock, Ark.		U.S. PATENT DOCUMENTS
[73]	Assignee: University of Arkansas , Little Rock, Ark.	4,340,340 8/1986 Smo et al. 5/1979 1/1979	552/8
[21]	Appl. No.: 437,140		Primary Examiner —Blaine Lundford
[22]	Filed: May 5, 1995		Attorney Agent, or Firm —Benjamin A. Adick
			ABSTRACT
[21]			The present invention provides the enzyme UDP-N-acetylglucosamine pyrophosphorylase in purified and purified form. Also provided is a various method of using and preparing this purified, homogeneous enzyme
[54]	Int. Cl.⁶	C12N 9/16; C12N 9/14	
[55]	U.S. CL.	435/196; 435/194	
[58]	Field of Search	435/196, 233, 435/194	
		3 Claims, 9 Drawing Sheets	

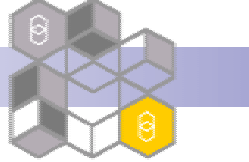
```

DE PROBABLE UDP-N-ACETYLGLUCOSAMINE 1-CARBOXYVINYLTRANSFERASE
DE (EC 2.5.1.7) (ENOYLPRUVATE TRANSFERASE) (UDP-N-ACETYLGLUCOSAMINE
DE ENOYLPRUVYL TRANSFERASE) (EPT).
GN MURA OR MURZ.
OS BACILLUS SUBTILIS.
OC BACTERIA; FIRMICUTES; BACILLUS/CLOSTRIDIUM GROUP; BACILLACEA
OC BACILLUS.
KW PEPTIDOGLYCAN SYNTHESIS; CELL WALL; TRANSFERASE.
FT ACT_SITE 116 116 BINDING (BY SIMILARITY).
FT CONFLICT 374 374 S -> A (IN REF. 3).
SQ SEQUENCE 429 AA; 46016 MW; 02018C5C CRC32;
MEKLNIAAGD SLNGTVHIS APLNSAVALPI ATLANSYVME IEGLEPEISDI ETLRLDLK
GNVHPFENG MVVDPTSMIS MKPDGKVRK LRIASYLVMG MLGRFKQAVI GLPGGCHL
RPIDQHIKGF ELGAELVINE QGAYILRAER VVARIYLDV VSVGATINIM LAAVLAEGH
IITENAKKEP IUDVATILIS MGAKIKGAG

```

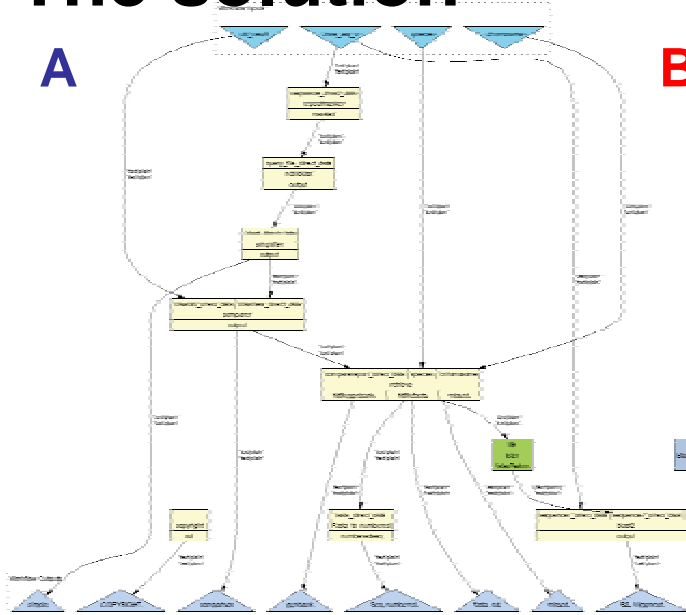
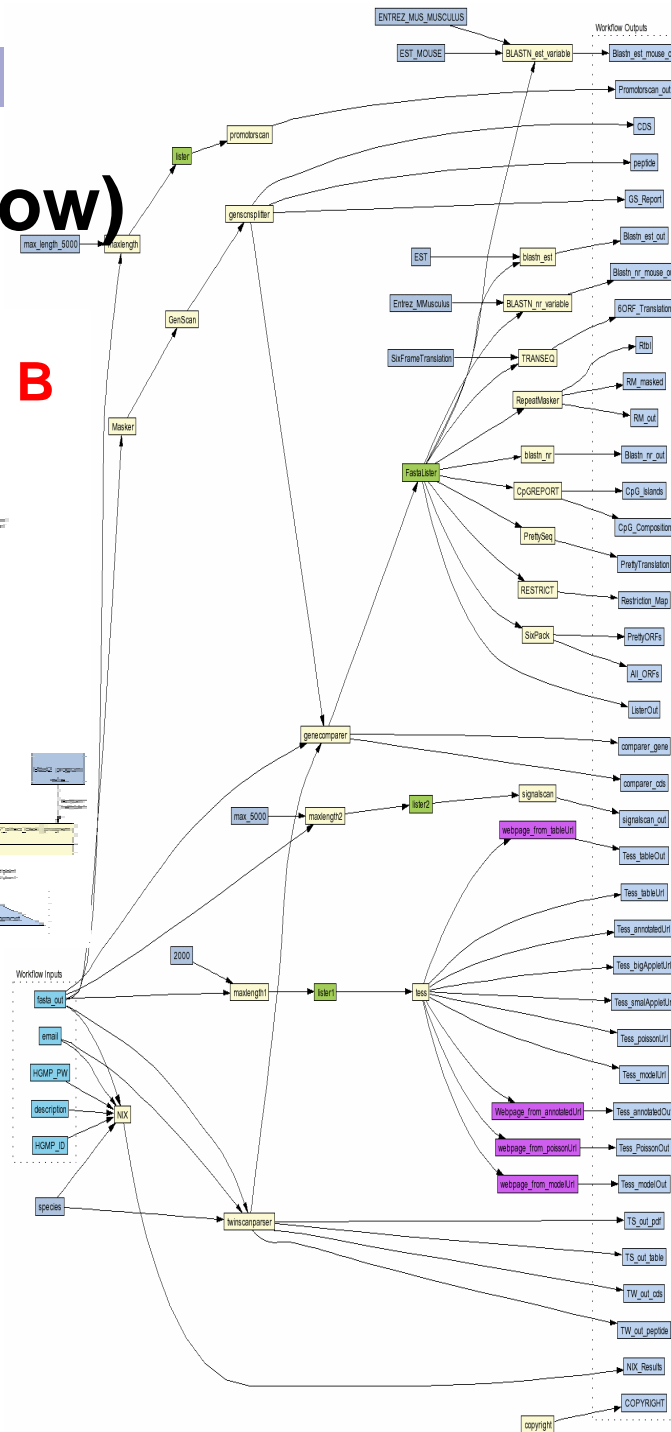
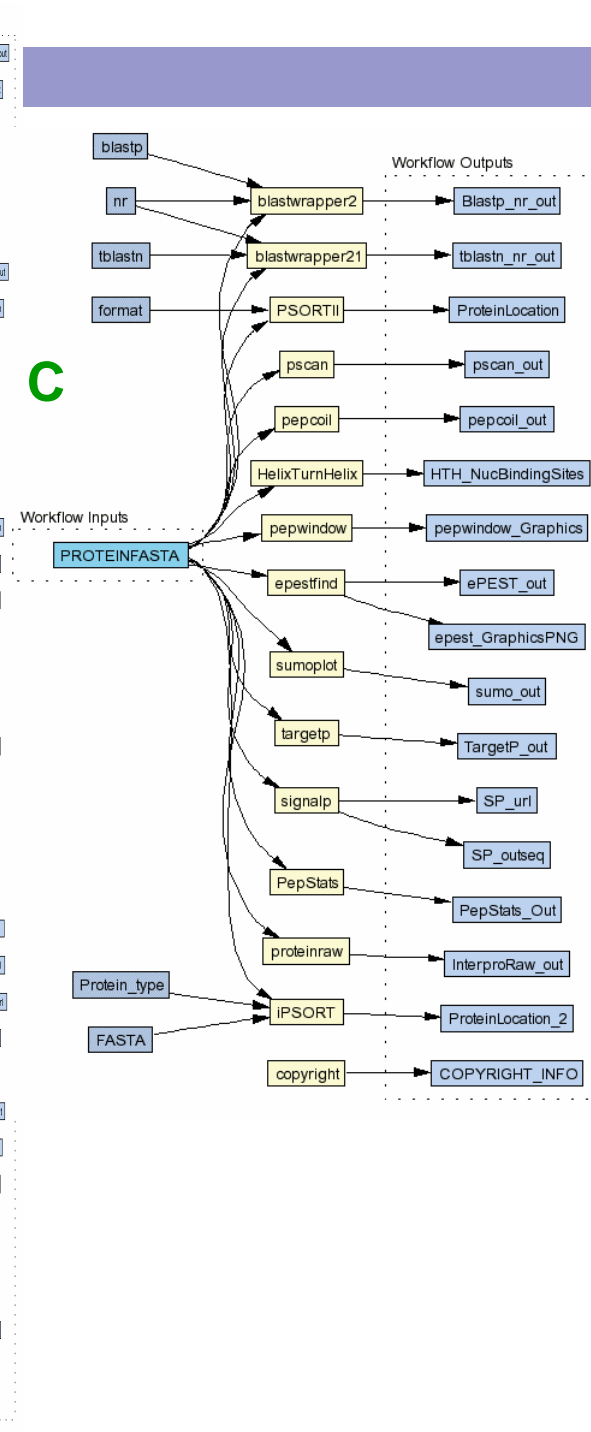
[illegible][illegible]

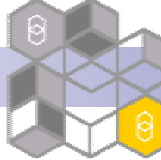




Workflows (dataflow)

The solution

A**B****C**



ID	Q99NI3	PRELIMINARY	unique identifiers
AC	Q99NI3;		
DT	01-JUN-2001	(TrEMBLrel. 17, Created)	
DT	01-JUN-2001	(TrEMBLrel. 17, Last sequence update)	

We don't (always) have XML

This is

xsd:string!

```

OC   Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
OX   NCBI_TaxID=10090;
RN   [1]
RP   SEQUENCE FROM N.A.
RC   STRAIN=C57BL/6J;
RX   PubMed=15100712; DOI=10.1038/sj.ejhg.5201174;
RA   Tipney H.J., Hinsley T.A., Brass A., Metcalfe K., Donnai D.,
RA   Tassabehji M.;
RT   "Isolation and characterisation of GTF2IRD2, a novel fusion gene and
RT   member of the TFII-I family of transcription factors, deleted in
RT   Williams-Beuren syndrome.";
RL   Eur. J. Hum. Genet. 12:551-560(2004).
RN   [2]
RP   SEQUENCE FROM N.A.
RA   Bayarsaihan D.;
RL   Submitted (MAY-2002) to the EMBL/GenBank/DDBJ databases.
DR   EMBL; AY014963; AAG41674.1; -.
DR   EMBL; AY116023; AAM48282.1; -.
DR   MGD; MGI:2149780; 1700012P16Rik.
DR   InterPro; IPR008938; ARM.
DR   InterPro; IPR004212; GTF2I.
DR   Pfam; PF02946; GTF2I; 2.
SQ   SEQUENCE 936 AA; 104577 MW; 28CB2994C0ABB1A9 CRC64;
      MAQVAVTTQP TDEPSDGRMV VTFLMSALES MCKELAKSKA KACIAVYET DVYVVGTERG
      CAFVNARQDL QKDFAHQCQG EGLPEEKPLC LGNGEACPGE AQLLRRRAVD HFCLCYRKAL
      GTTAMVPVPY EOMLODEAAV VVRGLPEGLA FOHPDNYSLA TLKWILENKA GISFAVKRPF
  
```

Gene name

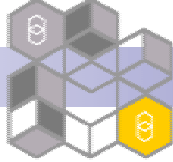
Protein sequence

...no type system either

- To add insult to injury, we don't have a type system either
 - Probably never will do
 - Getting people to describe their data types using a global model (schema or ontology) is *hard*
- Especially when there is no obvious immediate benefit of doing so (see www.BioMOBY.org)
- myGrid accommodates 3rd party Web Services because bioinformatics is an *open* world

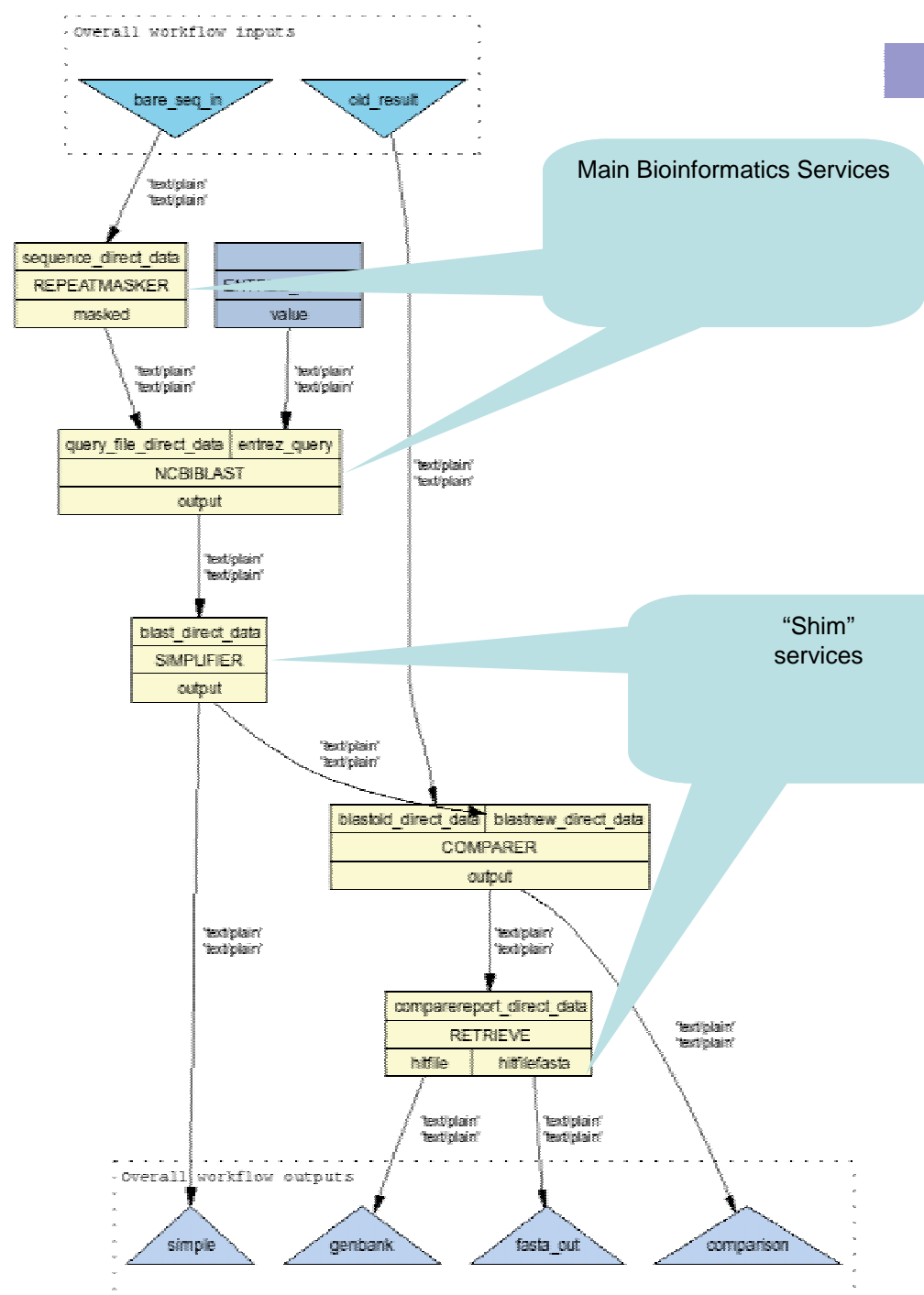
Web Services don't quite fit together

- Many bioinformatics Web Services *nearly* fit together
- This is where shims come into play
 - Thin strips of metal used to align pipes or rails
 - Software components that transform between closely related data (either syntactically or semantically)



Shims

- ~10 / 30 Web Services are shims
- ready made WS
- homemade WS wrap legacy applications



Problems with shim solution

- Finding the right shim at the right time
 - Resource intensive: knowledge, time etc
 - Shims re-invent the wheel
 - “Every biologist has written a BLAST parser”
 - For n services there are n^2 shims
 - In mygrid 600 services, therefore
worst case scenario: 360,000 shims
- Like it or not, we're stuck with shims
- What do we do about it?

Research goals

- Characterise the shims that are out there
- Classify and describe shims (and services they mediate between) using an ontology
- Create a library (not a factory) of shims
 - Reuse don't rewrite
- Model identifies mismatches
- Facilitate automation? It's not always safe...
- ...or at least support and guide user selection of shims during workflow construction

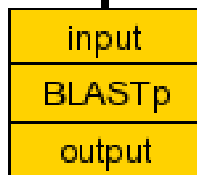
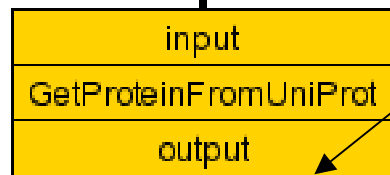
Shim description and classification

Shim type	Input	Operation	Output	Description
Dereferencer	Identifier or Pointer	Dereference	The dereferenced resource	GenBank ID replaced with GenBank record
Syntax translator	Data represented in concrete representation, x	Translate	Data represented in an alternate concrete representation, y	SeqRet translates between representations of sequence data.
Semantic translator	DNA sequence	Translate	Protein sequence	Translate DNA into protein
Mapper	Unique identifier	Map	Unique Identifier	Maps between IDs. E.g. GenBank to EMBL
Parser	Record	Parse	Abstract syntax tree	Parse BLAST report.
Iterator	A set	Iterate	A member of a set	Iterate over members of a given set
Comparer	Two or more sets	Diff	Set of differences	Comparing BLAST reports notifies of new sequences
Accessor	Record	Access	Subset of record	Access a subset

Table 1. Examples of shims taken from the *myGrid* project. This partial classification is based on inputs, outputs and the task or operation performed.

Workflow Inputs

MyFavouriteProtein



Workflow Outputs

BLAST Reports

Example 1

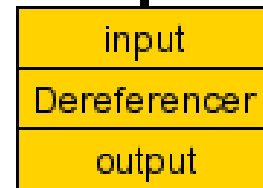
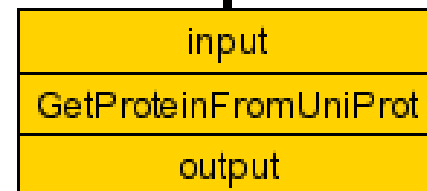
UniProt_ID

identifies

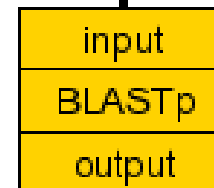
UniProt_Record

Workflow Inputs

MyFavouriteProtein



Shim



Workflow Outputs

BLAST reports

✗ Mismatch

✓ Match

Workflow Inputs

myFavouriteProtein

Example 2

UniProt record

hasPart

protein_sequence

input

GetProtein

output

input

BLASTp

output

Workflow Outputs

BLASTreports

Workflow Inputs

myFavouriteProtein

input

GetProtein

output

input

Accessor

output

Shim

input

BLASTp

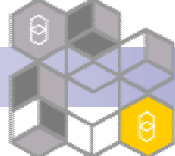
output

Workflow Outputs

BLASTreports

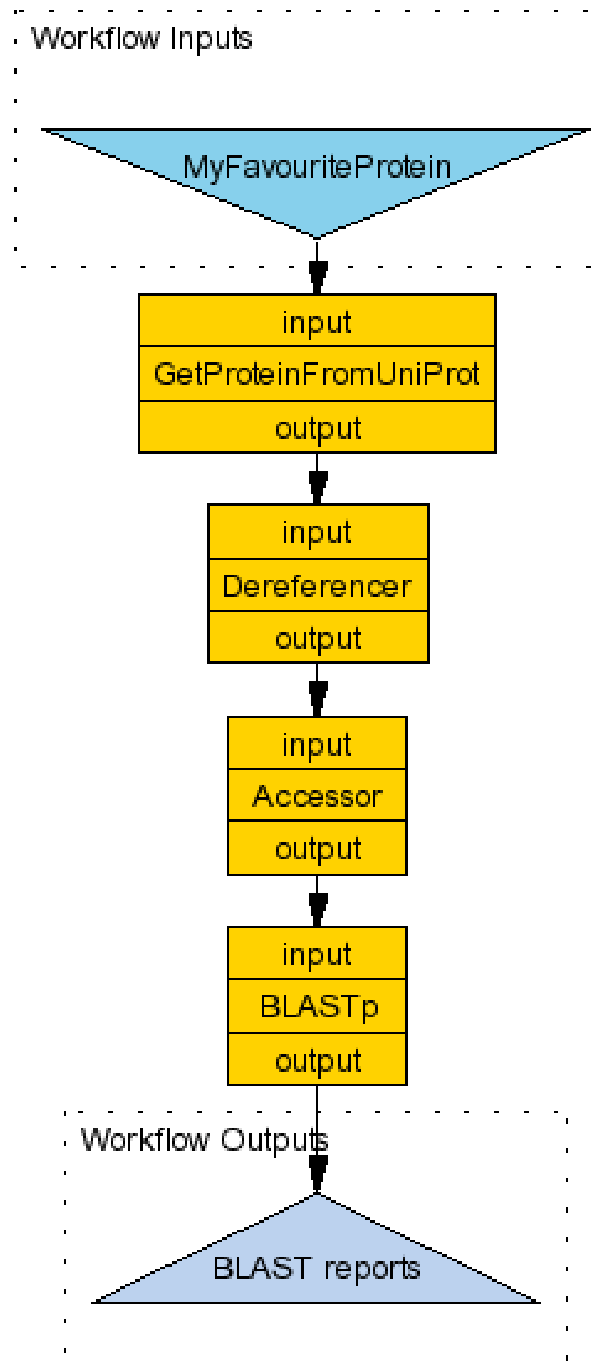
✗ Mismatch

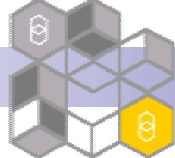
✓ Match



Example 3

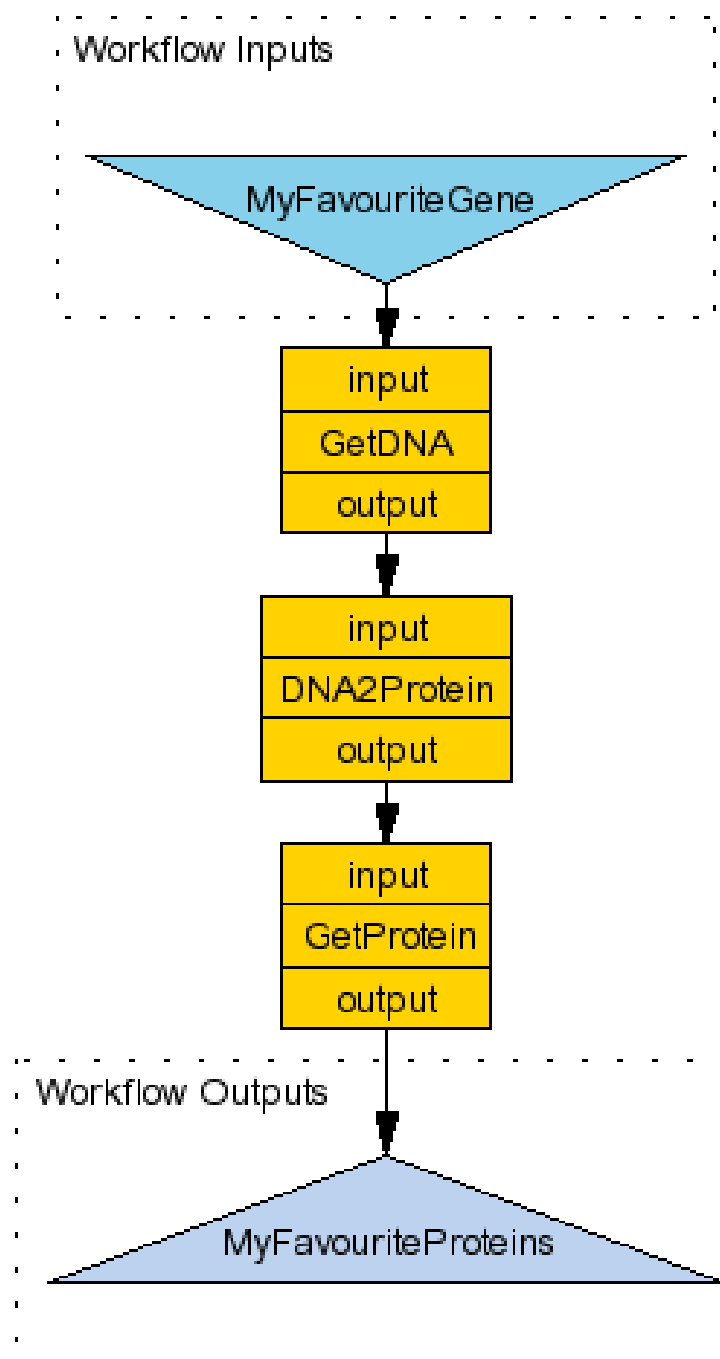
- Some shims are multistep
 - Depending on how you describe them
- Investigating AI planning techniques for this





It's not always safe!

- “Just because you can, doesn't mean you should”
- Automation is context dependent
- change workflow, change the experiment (seriously annoy the biologist)



Conclusions

- Web Services in bioinformatics leave integration problems: Shims mediate
- Tractable number of shim *types*
- Semantically annotating *all* services reduces mediation to service discovery
- Some shims maybe be automated but it can be dangerous

Acknowledgements

PhD supervisors/advisors

Robert Stevens, Carole Goble, Phillip Lord, Chris Wroe, Alvaro Fernandes

myGrid users

Hannah Tipney, May Tassabehji, Medical Genetics team at St Marys Hospital, Manchester, UK. Andy Brass

Simon Pearce and Claire Jennings, Institute of Human Genetics School of Clinical Medical Sciences, University of Newcastle, UK

myGrid core

Matthew Addis, Nedim Alpdemir, Tim Carver, Rich Cawley, Neil Davis, Alvaro Fernandes, Justin Ferris, Robert Gaizaukaus, Kevin Glover, Carole Goble, Chris Greenhalgh, Mark Greenwood, Yike Guo, Ananth Krishna, Peter Li, Phillip Lord, Darren Marvin, Simon Miles, Luc Moreau, Arijit Mukherjee, Tom Oinn, Juri Papay, Savas Parastatidis, Norman Paton, Terry Payne, Matthew Pocock Milena Radenkovic, Stefan Rennick-Egglestone, Peter Rice, Martin Senger, Nick Sharman, Robert Stevens, Victor Tan, Anil Wipat, Paul Watson and Chris Wroe.

Postgraduates

Martin Szomszor, Jun Zhao, Pinar Alper, John Dickman, Keith Flanagan, Antoon Goderis, Tracy Craddock, Alastair Hampshire

Industrial

Dennis Quan, Sean Martin, Michael Niemi, Syd Chapman (IBM), Robin McEntire (GSK)

