

Linking Proteins, Mutations, and Drugs

Applications of Data Integration using RDF

Emily Doughty

Collaborator: Brian Kirk

Mentor: Olivier Bodenreider

August 11, 2010

Outline

- Background on pharmacogenomics and our contribution to the field
- Brief overview of the framework and data integration of our system
- Validation applications and system traversals
- System limitations and future directions

Pharmacogenomics

Definition:

“The study of how variations in the human genome affect the response to medications”¹

- Mutations are one of the variations existing between individual human genomes
- How does the existence of a mutation impact drugs or vice versa?
- How does the association with a mutation impact drugs on proteins?



1. <http://www.medterms.com/script/main/art.asp?articlekey=15313>

Relating Drugs, Mutations, and Proteins

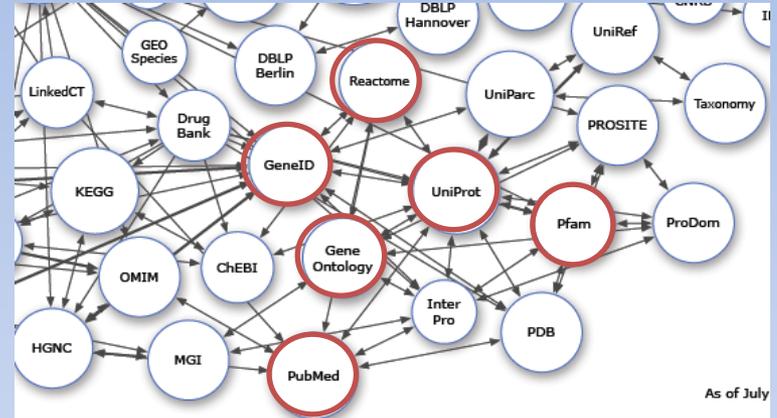
- Want to relate clinical drugs and mutations to protein structure and function
- Specifically, want to look at protein domains and pathways and relate to drug classes and properties
- Formulate new hypotheses for validation and future research

Framework

From Brian Kirk:

- Database Integration:

- UniProt
- NDF-RT
- PharmGKB
- Reactome and BioCyc
- Gene Ontology
- Protein to Domain table provided by Dr. Maricel Kann's lab at UMBC

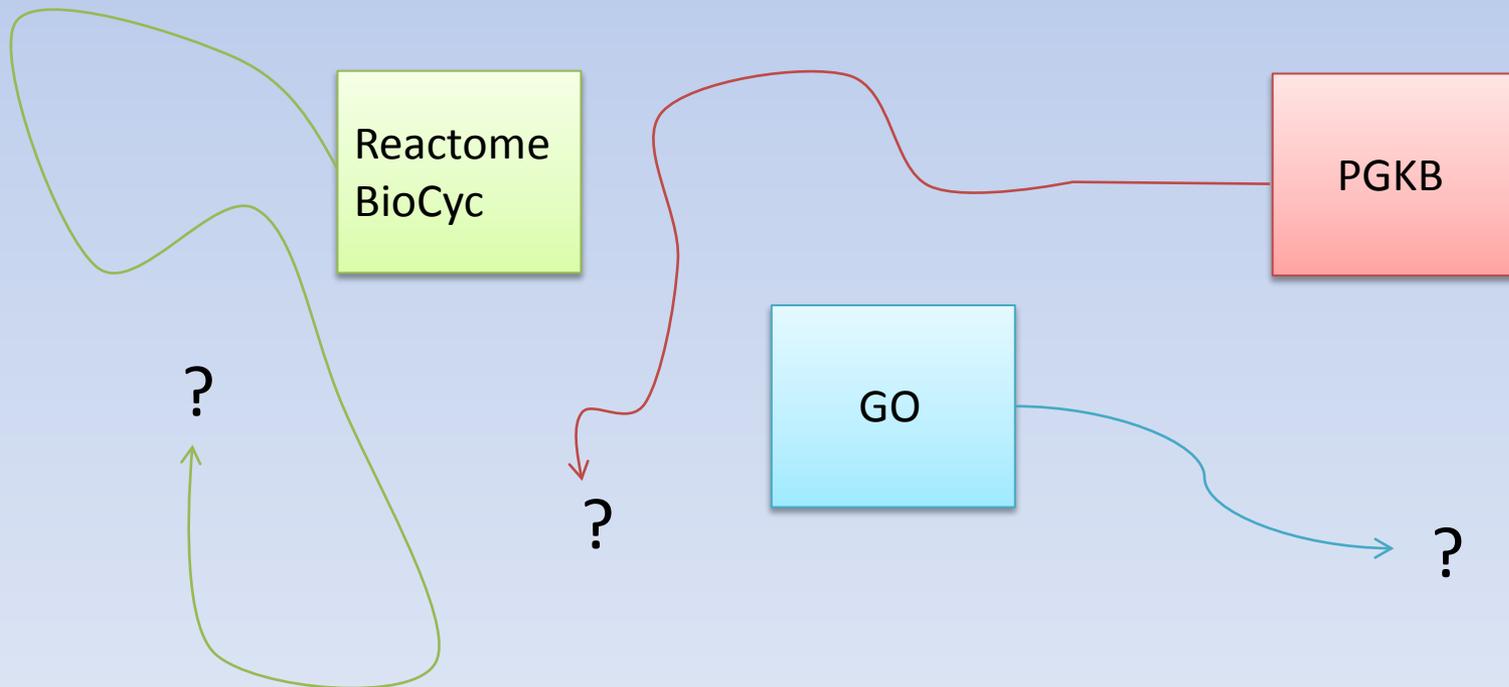


- High-throughput method for text mining for drug-related mutations

- Extractor of MUtations (EMU) developed at Dr. Kann's lab

Making Connections: Data Integration

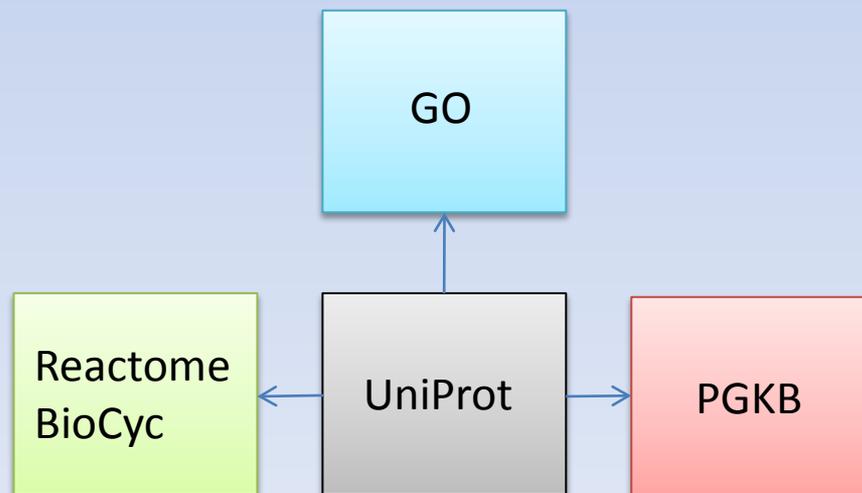
- Need to be able to connect database “pieces” together
- Solution: database identifiers and UniProt



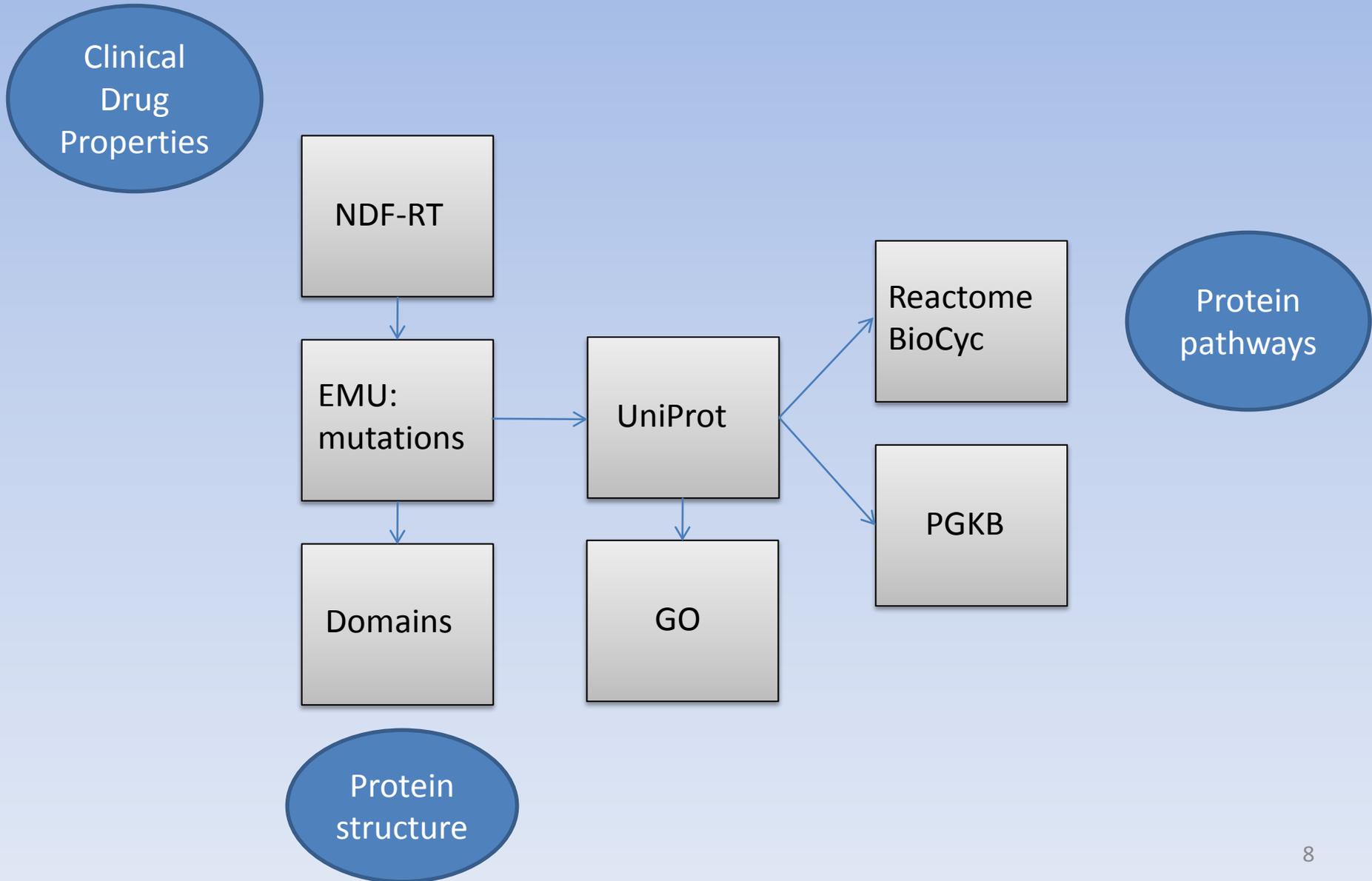
Making Connections: Data Integration

UniProt

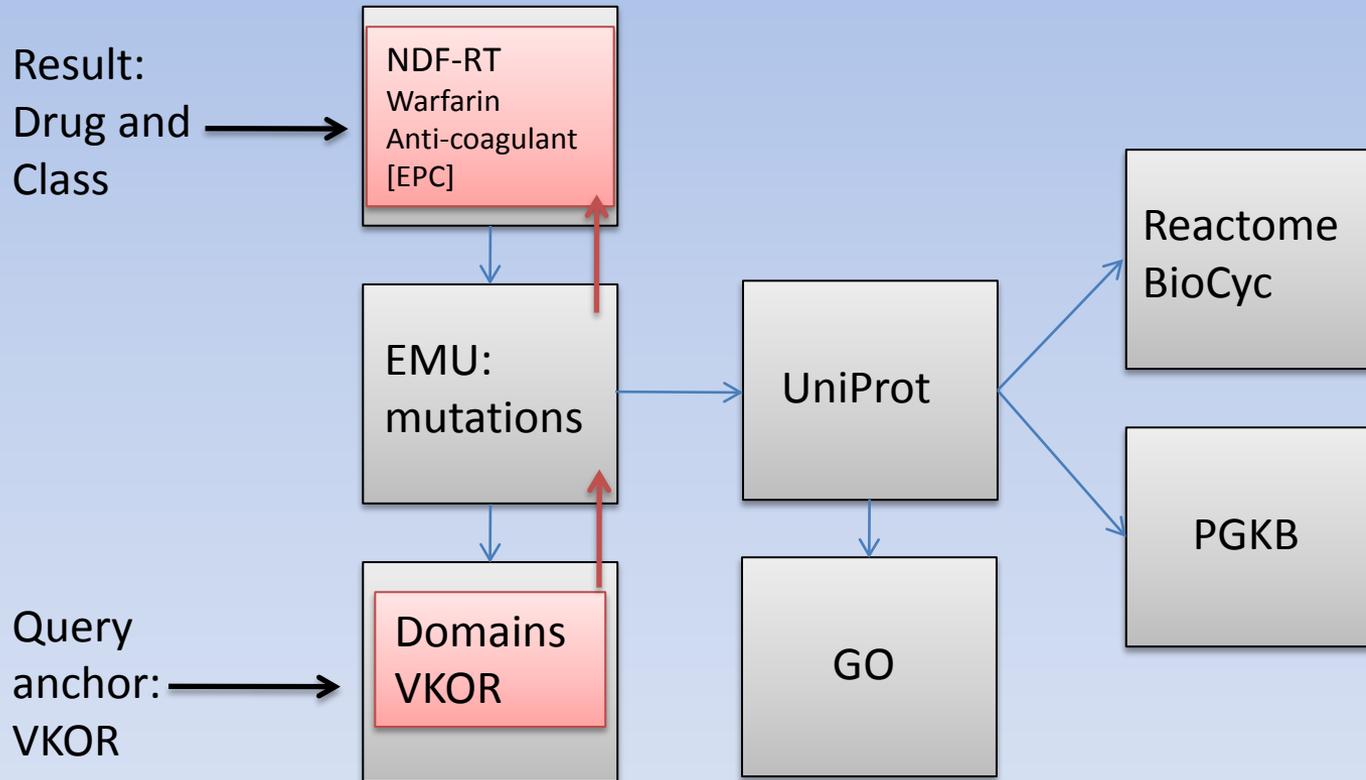
- For proteins and genes, contains database identifiers to other databases, including PharmGKB, GO, Reactome, BioCyc
- Connect to the mutational information through genes (or UniProt identifiers if available)



Schema



Query Strategy: an Example



What Can The System Accomplish?

- Able to link drugs to mutations to protein domains, pathways, and drug properties
- Can help with drug validation (ex. PharmGKB)
- Dynamic structure: can integrate new datasets easily for other drug-mutation questions (Ex. Drug adverse events)

Analysis: Warfarin

PharmGKB
Pharmacogenomics Knowledge Base

Search PharmGKB

Home ▾ Search ▾ Submit ▾ Download Help ▾ PGRN ▾ Contributors ▾ Clinical PGx

DRUG/SMALL MOLECULE:
warfarin

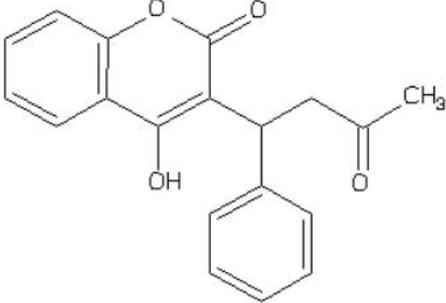
Overview Properties Genetics Related Genes Related Drugs Related Diseases Datasets Downloads/LinkOuts

Overview

Generic Names: Warfarin sodium

Trade Names: Athrombin; Athrombin-K; Athrombine-K; Brumolin; Co-Rax; Coumadin; Coumafen; Coumafene; Coumaphen; Coumaphene; Coumarins; Coumefene; D-Con; Dethmor; Dethnel; Dicusat E; Frass-Ratron; Jantoven; Kumader; Kumadu; Kumatox; Kypfarin; Latka 42; Mar-Frin; Marevan; Maveran; Panwarfin; Place-Pax; Prothromadin; RAX; Rosex; Sofarin; Solfarin; Sorexa Plus; Temus W; Tintorane; Tox-Hid; Vampirinip II; Vampirinip III; Waran; Warf 42; Warfarat; Warfarin Plus; Warfarin Q; Warfarine; Warficide; Warfilone; Zoocoumarin

PharmGKB Accession Id: PA451906



Description
(source: Drug Bank)

Indication
(source: Drug Bank)

ATC Therapeutic Category

- B01AA: Vitamin K antagonists

Analysis: Warfarin

DRUG/SMALL MOLECULE:

warfarin

Overview Properties Genetics Related Genes Related Drugs Related Diseases Datasets Downloads/LinkOuts

Pharmacology, Interactions, and Contraindications

Mechanism Of Action

(source: Drug Bank)

Pharmacology

(source: Drug Bank)

Food Interactions

(source: Drug Bank)

Absorption, Distribution, Metabolism, Elimination & Toxicity

Biotransformation

(source: Drug Bank)

Protein Binding

(source: Drug Bank)

Toxicity

(source: Drug Bank)

Isomeric SMILES Code:

CC(=O)CC(c1ccccc1)c2c(c3ccccc3oc2=O)O (source: Drug Bank)

1 [sparql PREFIX rdf: <ht...> Messages

drug_name	propName	diseaseName
Warfarin	has_Ingredient	Warfarin
Warfarin	has_MoA	Vitamin K Epoxide Reductase Inhibitors
Warfarin	has_PE	Decreased Coagulation Factor Activity
Warfarin	has_PE	Decreased Coagulation Factor Concentration
Warfarin	may_prevent	Cerebrovascular Accident
Warfarin	may_prevent	Coronary Thrombosis
Warfarin	may_prevent	Ischemic Attack, Transient
Warfarin	may_prevent	Myocardial Infarction
Warfarin	may_prevent	Postoperative Complications
Warfarin	may_prevent	Pulmonary Embolism
Warfarin	may_prevent	Thromboembolism
Warfarin	may_prevent	Venous Thrombosis
Warfarin	may_treat	Atrial Fibrillation
Warfarin	may_treat	Cerebrovascular Accident
Warfarin	may_treat	Postoperative Complications
Warfarin	may_treat	Pulmonary Embolism
Warfarin	may_treat	Thromboembolism
Warfarin	may_treat	Thrombophlebitis
Warfarin	site_of_metabolism	Hepatic Metabolism

Query executed in 1351 ms. Number of rows returned: 19

Analysis: Warfarin

- Found two genes associated with mutations: CYP2C9 and VKORC1
- Both genes annotated with mutations in PharmGKB
- Reconstruction of known information

The screenshot displays the PharmGKB (Pharmacogenomics Knowledge Base) website interface. At the top, there is a search bar and navigation links for Home, Search, Submit, Download, Help, PGRN, Contributors, and Clinical PGx. The main content area is titled 'DRUG/SMALL MOLECULE: warfarin' and features a tabbed interface with 'Genetics' selected. Under the 'Genetics' tab, there are three 'In-Depth Annotations' (indicated by three stars):

- rs1799853 at chr10:96692037 in CYP2C9**
This variant has been shown to influence warfarin dose as well as affecting the clearance of several other drugs.
Variant Name: CYP2C9*2; CYP2C9:144Arg>Cys
Related Drugs: fluvastatin, glipizide, phenytoin, tolbutamide, warfarin
Evidence: <http://www.pharmgkb.org/.../variant.jsp#ImportantVariantInformationforCYP2C9-111>
- rs1057910 at chr10:96731043 in CYP2C9**
This variant has been shown to correlate significantly with warfarin dose as well as affecting the clearance of several other drugs.
Variant Name: CYP2C9*3; CYP2C9:359Ile>Leu
Related Drugs: fluvastatin, glipizide, phenytoin, tolbutamide, warfarin
Evidence: <http://www.pharmgkb.org/.../variant.jsp#ImportantVariantInformationforCYP2C9-222>
- rs7294 at chr16:31009822 in VKORC1**
May be associated with a higher warfarin dose phenotype.
Variant Name: VKORC1:G9041A; VKORC1:3730G>A
Related Drugs: warfarin
Evidence: <http://www.pharmgkb.org/.../variant.jsp#ImportantVariantInformationforVKORC1-9041>

Analysis: Heparin

PharmGKB
Pharmacogenomics Knowledge Base

Search PharmGKB

Home ▾ Search ▾ Submit ▾ Download Help ▾ PGRN ▾ Contributors ▾ Clinical PGx

DRUG/SMALL MOLECULE:
heparin

Overview Properties Related Genes Related Drugs Related Diseases Downloads/LinkOuts

Pharmacology, Interactions, and Contraindications

Mechanism Of Action
(source: Drug Bank)

Pharmacology
(source: Drug Bank)

Food Interactions
(source: Drug Bank)

Absorption, Distribution, Metabolism, Elimination & Toxicity

Biotransformation
(source: Drug Bank)

Protein Binding
(source: Drug Bank)

Absorption
(source: Drug Bank)

Toxicity
(source: Drug Bank)

1 [sparql PREFIX ndf: <ht...>] Messages

drug_name	propName	diseaseName
Heparin	has_Ingredient	Heparin
Heparin	has_MoA	Antithrombin Activators
Heparin	has_MoA	Thrombin Inhibitors
Heparin	has_PE	Decreased Coagulation Factor Activity
Heparin	may_prevent	Postoperative Complications
Heparin	may_prevent	Pulmonary Embolism
Heparin	may_prevent	Thromboembolism
Heparin	may_prevent	Venous Thrombosis
Heparin	may_treat	Angina, Unstable
Heparin	may_treat	Brain Ischemia
Heparin	may_treat	Cerebral Infarction
Heparin	may_treat	Coronary Thrombosis
Heparin	may_treat	Myocardial Infarction
Heparin	may_treat	Postoperative Complications
Heparin	may_treat	Pulmonary Embolism
Heparin	may_treat	Thromboembolism
Heparin	may_treat	Thrombophlebitis

Query executed in 1381 ms. Number of rows returned: 17

Analysis: Heparin

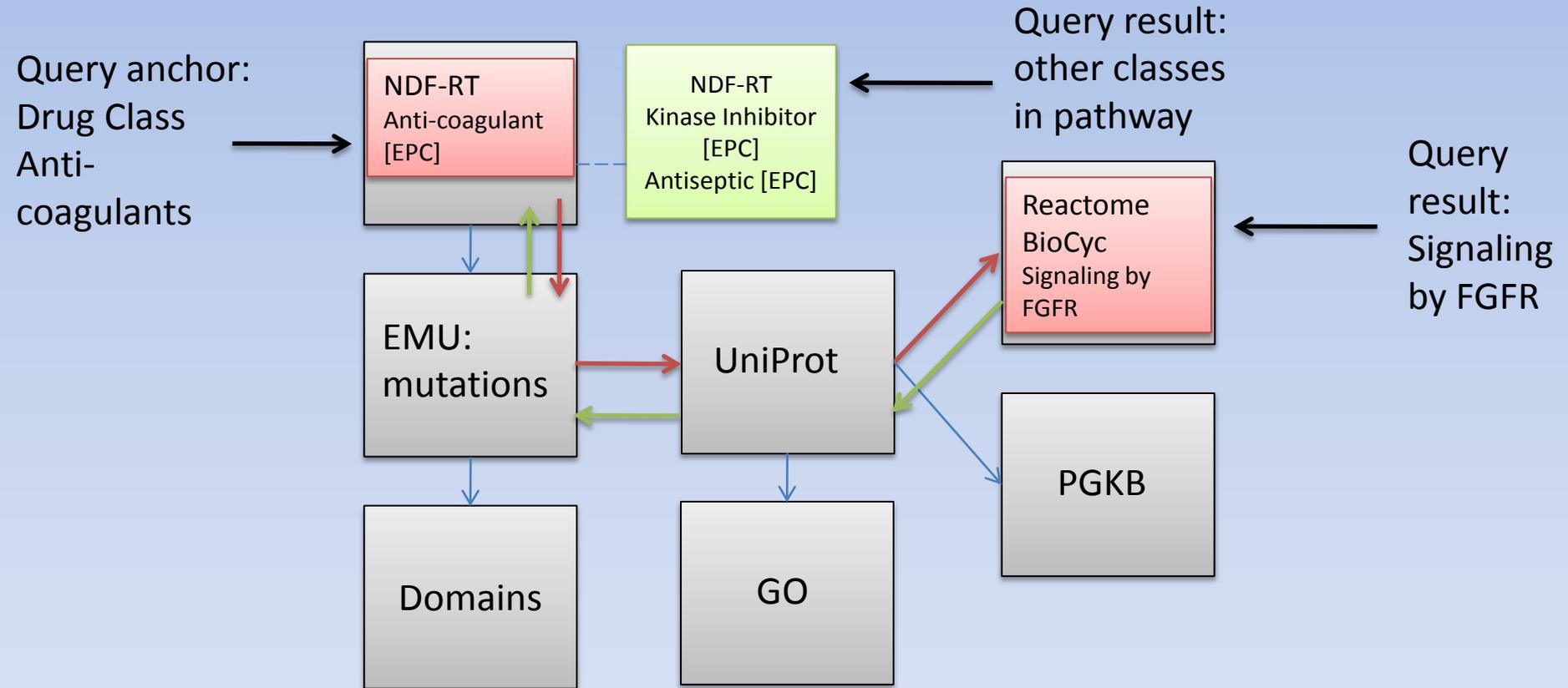
- Uncovered 239 domains from two sources (name variations possible)
- 29 citations found with associations to mutations
- PharmGKB has no information on genetics and thus no mutational links easily available
- Recommendation: validation of Heparin towards genetics using our list of citations as a starting point

Further System Applications

For a given class (ex. anti-coagulants):

- What are the pathways associated with that class through mutations
- What are the other classes associated with that pathway

Query Classes to Pathways

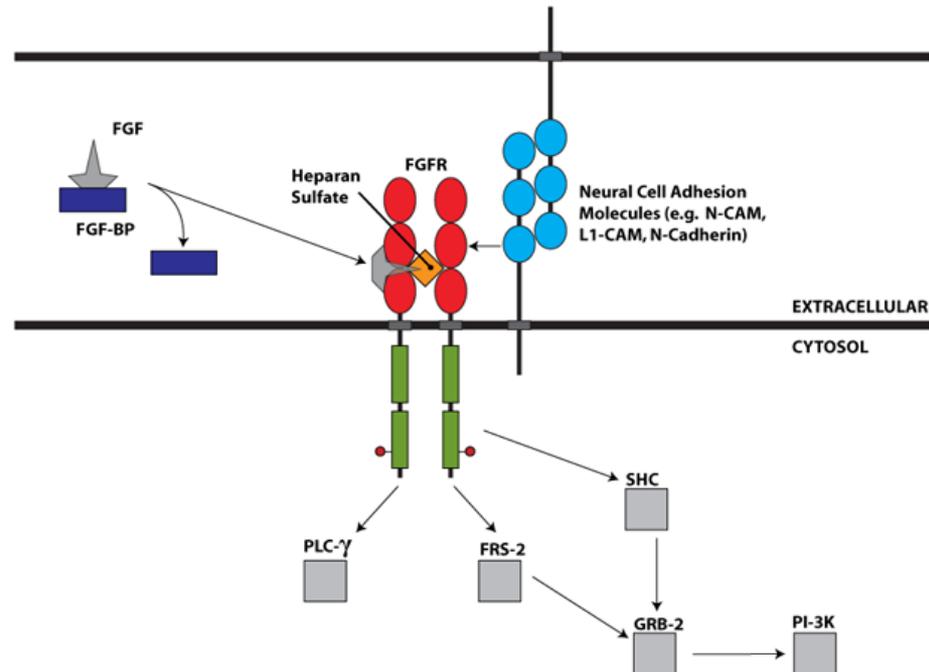


Demonstrating the System

Signaling by FGFR

Stable identifier	REACT_9470.2
Authored	de Bono, B, 2007-01-10
Reviewed	Mohammadi, M, 2007-02-06
Your feedback	Let us know what you think of this article (click here)

The 22 members of the fibroblast growth factor (FGF) family of growth factors mediate their cellular responses by binding to and activating the different isoforms encoded by the four receptor tyrosine kinases (RTKs) designated FGFR1, FGFR2, FGFR3 and FGFR4. These receptors are key regulators of several developmental processes in which cell fate and differentiation to various tissue lineages are determined. Unlike other growth factors, FGFs act in concert with **heparin** or heparan sulfate proteoglycan (HSPG) to activate FGFRs and to induce the pleiotropic responses that lead to the variety of cellular responses induced by this large family of growth factors. An alternative, FGF-independent, source of FGFR activation originates from the interaction with cell adhesion molecules, typically in the context of interactions on neural cell membranes and is crucial for neuronal survival and development. Upon ligand binding, receptor dimers are formed and their intrinsic tyrosine kinase is activated causing phosphorylation of multiple tyrosine residues on the receptors. These then serve as docking sites for the recruitment of SH2 (src homology-2) or PTB (phosphotyrosine binding) domains of adaptors, docking proteins or signaling enzymes. Signaling complexes are assembled and recruited to the active receptors resulting in a cascade of phosphorylation events. This leads to stimulation of intracellular signaling pathways that control cell proliferation, cell differentiation, cell migration, cell survival and cell shape, depending on the cell type or stage of maturation. [Eswarakumar *et al* 2005, Dailey *et al* 2005, Schlessinger 2000]



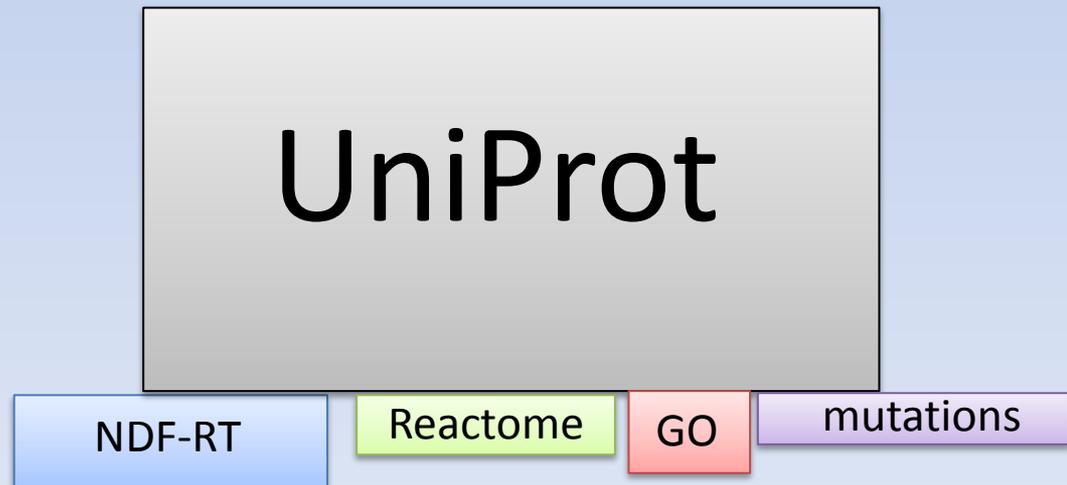
Preliminary Quantitative Results

		Drug Characteristics	
		Class (EPC)	Properties
Protein Characteristics	Pathway	5,724 (pathway, gene, drug, class)	In Progress...
	Domain	55,930 (domain, drug, class)	In Progress...

Technical Challenges

■ UniProt

- massive database with ~1.5 billion triples (majority of data)
- needs a robust server and efficient queries in order to be used effectively
- however, contains useful information, not overwhelming and still query-able



Technical Challenges

- Interoperability
 - No way to link Reactome and BioCyc together
 - Solution: go through UniProt
- Smart and efficient query strategies
 - Due to size of the entire system and complexity, need efficient query strategies to get results
 - Need an entry point for quick querying (starting from specific drug, domain, etc)
 - Reducing UniProt will help with this problem, though not completely

System Limitations

- With current level of knowledge, difficult to make generalized hypotheses
- Difficult to make direct connections between properties of drugs and protein levels due to interaction complexity
- System is weakened overall by its weakest links, but difficult to fully exploit in 2.5 months
- Post processing difficult with the high volume of data

Additional Limitations

- Mutations
 - EMU mutation extractions were not curated for drug-mutation relationships
 - All output automatically generated
 - EMU was evaluated on mutation extraction on a corpus of disease-related abstracts, and not evaluated on a mutation-drug corpus
- Mutation-Drug Association
 - Due to lack of curation, only say that drugs are associated with mutations, and cannot mention drugs specifically causing or affecting a mutation

Future Work

- Add additional functionality such as drug adverse events
 - easy to do if identifiers to other databases are in UniProt or if additional databases use Entrez geneID
- Curate a random portion of EMU output to create a gold standard
- Use only a subset of UniProt to improve query runtime
- Submit our findings for the Pacific Symposium on Biocomputing workshop on Mining the Pharmacogenomics Literature

Acknowledgements

- Brian Kirk, collaborator
- Jonathan Mortensen
- Olivier Bodenreider, mentor
- Dr. Maricel Kann's lab, UMBC
- May Cheh
- Dr. McDonald