

Investigating plant systems using data integration and network analysis

Plant Bioinformatics, Systems and Synthetic Biology Summer School

Nottingham, July 30th 2009 Chris Rawlings Catherine Canevet, Michael Defoin-Platel







Examples of Plant BioInformatics simple questions requiring data integration

 Which metabolic pathways is a gene or set of genes involved in?

- Information distributed across several pathway databases
- What transcription factors are involved in that pathway?
 - Information in pathway databases as well as transcription factor databases





ROTHAMSTED



- Data integration
 - Needs for data integration
 - Benefits of data integration
 - Challenges
 - Solutions
 - Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation





- Data integration
 - Needs for data integration
 - Benefits of data integration
 - Challenges
 - Solutions
 - Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation



Plant databases many locations, many formats

o TAIR

(The Arabidopsis Information Resource)

o AraCyc

(Arabidopsis Thaliana metabolic pathways)

AtRegNet

(Arabidopsis Thaliana Regulatory Networks)

o Gramene

(Comparative Grass Genomics)

o Grassius

(Grass Regulatory Information Services)

o TIGR

(Rice Genome Annotation Resource)



Plant & crop databases (2/2)

- The Brachypodium distachyon Information Resource
- o GrainGenes
- Maize GDB
- Maize Mapping Project
- SGN (Sol Genomics Network)
- SoyBase
- UrMeLDB (Medicago trunculata)
- o BeanGenes
- o Legume Base
- o JGI Poplar database



Solo -

Other bioinformatics databases

- Bcsdb (Bacterial Carbohydrate Structures)
- BioCyc (Pathway/Genome Databases)
- BioGRID (General Repository for Interaction Datasets)
- BRENDA (Enzymes)
- GOA (Gene Ontology Annotation to UniProtKB proteins)
- KEGG (Kyoto Encyclopedia of Genes and Genomes)
- PDB (proteins, nucleic acids and complexes)
- Pfam (protein families)
- SGD (Saccharomyces Genome)
- TRANSFAC (transcription factors)
- TRANSPATH (pathways)
- UniProt (proteins)







Spa Spa	rseness o	of plan	t d	lata	a (2)	
🔰 🔰 Jackson e	t al.		Order	Family	Species	Represented Families
S Plant Cell 3	2006	Eurosid I (Fabids)	8	64	Glycine max, Lotus japonicus, Medicago, Populus trichocarpa	2
		Eurosid II (Malvids)	З	34	Arabidopsis lyrata, A. thaliana, Brassica oleracea, B. rapa, Capsella rubella, Carica papya	2
	}	• Other Rosids	5	35	Eucalyptus globulus, Vitis vinifera	2
		Euasterid I (Lamiids)	4	40	Lycopersicon esculentum, Mimulus guttatus, Solanum tuberosum, Triphysaria versicolo	3 r
Model		Euasterid II (Campanulids)	4	35		0
nlanta		Other Asterids	2	30		0
Few and far between		Other core eudicots	4	24		o
		Early diverging dicots	2 ^a	12		0
		Magnoliids	4	18		0
	-	Commelinids	3 ^b	29	Musa acuminata, Oryza sativa, Sorghum bicolor, Triticum aestivum, Zea mays	2
		• Other monocots	; 8	53		ο
		Basal Angiosperms	1 ^c	5		0
		Gymnosperms	4	11		0
		Ferns	16	55		0
		Lycophytes	3	3	Selaginella moellendorfii	1
		Bryophytes	13	158	Physcomitrella patens	1
Convright @2006 American Society of Plant Biologiste		Green Algae			Critamydomonas reinhardtii, Micromonas pusilla, Ostreococci tauri	us
ROTHAMSTED RESEARCH		Total	84	606		13

Sparseness of plant knowledge (4)



Model plant genomes small and "simple"

TABLE I

Nuclear genome size in different species

Common name	Scientific name	Nuclear genome size ⁽¹⁾
Wheat	Triticum aestivum	15,966
Onion	Allium cepa	15,290
Garden pea	Pisum sativum	3,947
Com	Zea mays	2,292
Aspairagus	Asparagus officinalis	1,308
Tomato	Lycopersicum esculentum	907
Sugarbeet	Beta vulgaris	758
Apple	Malus X domestica	743
Common bean	Phaseolus vulgaris	637
Cantaloupe	Cucumis melo	454
Grape	Vitis vinifera	483
Man	Homo sapiens	2,910



1: Expressed in Megabases (1Mb:1,000,000 bases)



Binary Interactions Proteins



Conclusion

 Plant and crop scientists, more than others need to pake most of all their data

 Exploit a from get eater variety of sources

mention other challenges • Not to

- Plany in the environment
- Crop systems interactions with other organisms





Data integration

- Needs for data integration
- Benefits of data integration
- Challenges
- Solutions
- Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation











Data integration

- Needs for data integration
- Benefits of data integration
- Challenges
- Solutions
- Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation



Syntactic integration challenge



Over **1000 databases** freely available to public

Over **60 million sequences** in GenBank

Over **870 complete** genomes and many ongoing projects

Over **17 million citations** in PubMed

PubMed growth by 600,000 publications each year

Integration of Life Science data sources is essential for Systems Biology research



Syntactic integration example

- Databases widely different formats for same information
 - KEGG (Kyoto Encyclopedia of Genes and Genomes, <u>www.genome.jp/kegg/</u>)
 - TAIR (The Arabidopsis Information Resource, <u>www.arabidopsis.org</u>)
- Screenshots for same entry
 - Protein FUMARASE 1
 - Arabidopsis AT2G47510







Arabidopsis thaliana (thale cress): AT2G47510

Help



TAIR - The Arabidopsis Information Resource

*	un Unio Contra	A About II.	Lanin			0		
Search Bro		le	Stocks	Dortals	Download	Submit	Nows	arch
Search Bio	100	15	SLUCKS	Fortais	Download	Submit	news	
Locus: AT2G	47510					ID		
Date last modified	2003-05-02							
TAIR Accession	Locus:2061966							
Gene Model @	AT2G47510.1							
Other names:	FUM1, FUMAR/	ASE 1, T30E	322.19			G	ENE NAM	ME
Description 0	fumarase (FUM	1)						
Other Gene Models	FUM1 AT2G4 (splice	7510.2 variant)						
Annotations 0	Category		Re	ationship Type 0	Key	word 0		
	GO Biological F	rocess	inv	olved in	resp	onse to oxidativ	/e stress	
	GO Cellular Co	mponent	loc	ated in	mito	chondrion		
	GO Molecular F	unction	ha	в —	fuma	arate hydratase	activity	
				Annotation	Detail			
RNA Data								
					•			
Two-channel Arrays	array element name Ø	avg (st). log d. err		Concep	ot		
	103P1	-0.0	018 (I	4700				
			<u>ID:</u> Lo	bcus: AI2G	47510			
One-channel Arrays	array element		a <u>Anno</u>	<u>tation:</u>				
	name 🛛		ⁱⁿ Conc	eptClass: F	PROTEIN			
	248461_S_AT		¹¹ Data	Source: TA	IR			
	16604_S_AT		5 Conc	ontNamo:	FLIM1 (pre	forrod)		
13 AL MARKS						ereneu),		
Associated	type			ARASE I, I	30BZZ.19			
Transcripts V	EST		(35 Conc	eptAccessio	<u>on:</u> AT2G4	/510 (1/	AIR)	
	cDNA		(5) Conc	eptGDS:				
Chromosome	2		TAXII	D (String)	: 3702			,
Nucleotide Sequence Ø	full length CDS	full length	genomic full I	ength cDNA				
Protein Data 0	name	Length (aa)	molecular weight	isoelectric point	domains(# of	domains)		

Semantic Integration challenge

Same concept different names

- synonyms
- ontologies
- Same name different concepts
 - homographs



 Will expand on these issues later

antAmiGO DATABASE: PO_0205

inflorescence

Accession: PO:0009049 Aspect: plant structure Synonyms: cob (sensu sorghum) corymb cyme dichasium drepanium helicoid cyme monochasium panicle raceme rhipidium scorpioid cyme spike (sensu Triticeae) umbel verticillaster

Definition:

That part of the axial system of plants above the uppermost foliage leaf/pair of foliage leaves that bears flowers. Comment:

Some plants have only solitary flowers, e.g. Magnolia









Data integration

- ➤ Needs for data integration
- Benefits of data integration
- Challenges
- Solutions
- Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation





- List of standardized terms or descriptors whose meanings are specifically defined or authorized by a standards organization
- EC nomenclature, Gene nomenclature
 - <u>http://www.arabidopsis.org/nomencl.html</u>
- Compendium of nomenclature resources
 - <u>http://www.expasy.ch/cgi-bin/lists?nomlist.txt</u>
- CVs help tackle consistency within species/group/taxa





Underpinning Data Integration

Ontologies

- an ontology is a formal representation of a set of concepts within a <u>domain</u> and the relationships between those concepts.
- A semantic structure
- Generally hierarchical
- Provides mechanisms for standardizing <u>across</u> species/groups/taxa





ROTHAMSTED

Underpinning Data Integration Example: The Gene Ontology





Plant Ontologies

- o <u>http://www.gramene.org/plant_ontology/</u>
- Plant structure (PO), Growth stage (GRO)
- Trait (TO), Environment (EO)

GRAMENE Search Image: Search Search Genomes Species Download Resources About Help Image: Search Browse Ontology Submission I Tutorial FAQ Help Image: Find: GRO:0007057 Image: Gene Image: G			GRAMENE Search Search Genomes Species Download Resources About Help Ontologies Search Browse Ontology Submission Tutorial FAQ Help Find: TO:0000179 Image: Contrology: Gene (GO) Plant Structure (PO) Growth Stage (GRO) Trait (TO) Options: Exact Match Include Obsolete Terms				
■ ■ Growth Stage Term "1.06-le	af just at coleoptile tip" (GRO:0007057)	11	🗷 🗖 Trait Term "biotic stress tra	ait" (T(0:0000179)		
Term Name	1.06-leaf just at coleoptile tip	11	Term Name		biotic stress trait		
Term Accession	SRO:0007057		Term Accession	(TO:0000179		
Aspect	Growth Stage		Aspect		Trait		
Synonyms (0)	None.		Synonyms (0)		None.		
Definition	The leaf has just emerged.		Definition		Response by the plant in terms of resistivity or sensitivity to bi		
Comment	Zadok scale-9, Haun scale-0		Comment		None		
Derivation			Derivation				
 all (all) #858384 E [i] cereal plant growth stage ontology (GRO:0007199) #14835 E [i] wheat, barley and oat growth stage (GRO:0007156) #9 E [i] 01-germination (GRO:0007051) #0 E [i] 1.06-leaf just at coleoptile tip (GRO:0007057) #0 			all (all) #858384 ⊞				



 Extensible Markup Language (XML) facilitates the sharing of data across heterogeneous computer systems and improves the consistency

Markup Language	Purpose
AGAVE	Genomic annotation and visualization
BioML	Experimental information for biopolymers
<u>BSML</u>	Genomic sequences and biological function
CML	Management of molecular information
MAGE-ML	Microarray gene expression data exchange
<u>SBML</u>	Systems biology and biochemical networks





Minimum Information for Biological and Biomedical Investigations

- Minimum Information Standards for Metadata
 - Data about data
- MIBBI Initiative
- o <u>http://www.mibbi.org/</u>
 - MIAME, MIAPE, MIQAS etc
 - Microarray, Proteomics, QTL and Association Experiment





Standards and Underpinnings

- Nomenclatures, Controlled Vocabularies
- Ontologies
- o Metadata
- All help tame the complexity of data integration – but they are not enough
 - Legacy information
 - Not widely followed
 - Competing approaches
 - Data integration solutions needed
 - Has been major Bioinformatics challenge



- Annal

Solutions for Data Integration (1/4)

Workflows

- techniques for enacting series of linked processes, useful to systematically automate protocols in bioinformatics
- e.g. InfoSense, Taverna
- generally hard to write almost programming

Mashups

- Data from different Web services or RSSfeeds are selected by the user and then "mashed" to form a new Web application
- e.g. UniProt DASTY
- more aggregation than integration



House

Solutions for Data Integration (2/4)

Service oriented architectures

- interconnection of data sources
- e.g. CORBA, Web services
- only "plumbing" techniques
- often poorly documented and constructed and only for programmers

Link integration

- cross-reference data entry from different data sources
- e.g. SRS (40% EMBL-EBI traffic), Entrez, Integr8
- problems with name clashes, ambiguities and updates





Solutions for Data Integration (3/4)

Data Warehousing

- data sources are processed and combined into a new data model to form a new data source
- e.g. eFungi, ATLAS, GIMS
- toolkit: GMOD, BioMART, Intermine, BioWarehouse
- high building and maintenance costs

View integration ("virtual warehouse")

- data are kept in the original sources but appear to be in a single database
- e.g. BioZon
- complex and slow for users and developers





Solutions for Data Integration (4/4)

Model-driven service oriented architecture

- data resources and tools are obliged to adhere to a designed model
- e.g. caBIG
- only possible in tightly coupled systems often difficult to achieve consensus

Integration applications

- built specifically to integrate data
- e.g., ToolBus, Youtopia, Ondex
- more often for a single application domain




Ondex as an example data integration system

• Three closely coupled aspects:

- Data input technical integration and transformation into data domain graph
- Mapping across graphs semantic integration
- Data visualisation and quantitative analysis
- In use in all BBSRC-funded systems biology Centres



o Open source



Everything is a network











Data integration in Ondex







Properties: compound name, protein sequence, protein structure, cellular component, KM-value, PH optimum ...



Ontology of Concept Classes, Relation Types and additional Properties

Concepts and relations (2/2)





Semantic Integration by Graph Alignment

- Going beyond technical data integration
 - Normalising data formats
- Create relations between equivalent entries from different data sources
- Identified by *Mapping methods*
 - Concept name (gene name), synonyms
 - Concept accessions (UniProt ID)
 - Graph neighbourhood
 - Sequence methods







Data integration – accession matching

- Matching of accession
- Within the context of data source

Concept	Concept
ID: AT2G47510 Annotation: catalytic/ fumarate hydratase [EC:4.2.1.2] ConceptClass: PROTEIN DataSource: KEGG ConceptName: FUM1 (preferred), FUMARASE 1 ConceptAccession: AT2G47510 (TIGR), AT2G47510 (TAIR), AT2G47510 (MIPS), 15226618 (NCBI-GI), 819364 (NCBI-GeneID), P93033 (UniProt) ConceptGDS: TAXID (String) : 3702 AAseq (String) : MSIYVASRRKSGGTTVTALRY NAseq (String) : atgrcgatttacgtcgcgtcgcgacggct	ID: Locus:AT2G47510 <u>Annotation:</u> <u>ConceptClass:</u> PROTEIN <u>DataSource:</u> TAIR <u>ConceptName:</u> FUM1 (preferred), FUMARASE 1, T30B22 19 <u>ConceptAccession: AT2G47510 (TAIR)</u> <u>ConceptGDS:</u> TAXID (String) : 3702





Data integration – name matching

• At least two concept names have to match

Concept ID: AT2G47510 Appetation: aptalutia/ fumarata hudrataga [EC:4.2.1.2]	Concept
<u>ConceptClass:</u> PROTEIN DataSource: KEGG <u>ConceptName: FUM1 (preferred), FUMARASE 1</u> <u>ConceptAccession:</u> AT2G47510 (TIGR), AT2G47510 (TAIR), AT2G47510 (MIPS), 15226618 (NCBI-GI), 819364 (NCBI-GeneID), P93033 (UniProt) <u>ConceptGDS:</u> TAXID (String) : 3702 AAseq (String) : MSIYVASRRKSGGTTVTALRY NAseq (String) : atgtcgatttacgtcgcgacggct	ID: Locus:AT2G47510 <u>Annotation:</u> <u>ConceptClass:</u> PROTEIN DataSource: TAIR <u>ConceptName:</u> FUM1 (preferred), FUMARASE 1, T30B22.19 <u>ConceptAccession:</u> AT2G47510 (TAIR) <u>ConceptGDS:</u> TAXID (String) : 3702





Principle of alignment of concepts as nodes in a data graph

- How can we map 2 concepts (proteins)
 - from 2 different data sources (KEGG & TAIR)
- Examine attributes (name, accession, graph structure, sequence)







Data integration – sequence matching

- Sequence similarity methods
 - e.g. BLAST
- o Link to a new resource
 - Without any concept attributes in common
- Other similarity matches include
 - EC2GO
 - Pfam2GO





Ondex data integration scheme

Data input Visualisation Data integration & transformation ONDEX Clients/Tools Heterogeneous Integration Methods data sources ONDEX Visualization Generalized Object Data Model Accession Parser UniProt **Tool Kit** Name based Database Web Client Parser AraCyc Transitive Taverna Layer Parser GO Blast **ProteinFamily** Parser Data Exchange Pfam Pfam2GO OXL/RDF 4 Lucene Parser PDB Text mining WebService





Importing data into Ondex

- What databases to import
- What format these are in
- Ondex parsers already written
 - Generic
 - o OBO, PSI-MI, SBML, Tab-delimited, Fasta
 - Database-specific
 - Aracyc, AtRegNet, BioCyc, BioGRID, Brenda, Drastic, EcoCyc, GO, GOA, Gramene, Grassius, KEGG, Medline, MetaCyc, Oglycbase, OMIM, PDB, Pfam, SGD, TAIR, TIGR, Transfac, Transpath, UniProt, WGS, WordNet







Data integration

- Needs for data integration
- Benefits of data integration
- Challenges
- Solutions
- Case study

Tutorial

Network analysis

- Definition
- Visualising biological networks
- Annotation





- Reference database of virulence and pathogenicity genes validated by gene disruption experiments
 - Literature mining
 - http://www.phi-base.org/
- Sequence comparison orthology and gene cluster analysis



http://www.phi-base.org/

PHI- base	Pathogen Host Interactions	This database contains expertly curated molecular and biological information on genes proven to affect the outcome of pathogen-host interactions. Information is also given on the target sites of some anti-infective chemistries.				
Search Abo	out Release notes	Download	d Disclaimer	Errors & contributions	Help Consortium	
Quick Searc	h					
Search all	Y for			order by Gene	name 💌 Go Clear	
e.g. 'ACE*', 'Ca	andida a*' or 'PHI:441'					
Advanced Se	earch					
Search	Gene	for	all	*	Go Clear	
⊙ and ○ or	Disease	for	all	*	• List of "hot" targ	
⊙ and ○ or	Host	for	all	*	genes curated fr	
💿 and 🔘 or	Pathogen	for	all		literature	
⊙ and ○ or	Anti-Infective	for 💿	all		interatore	
		0	Allylamines	~	-Loss of	
			Carboxylic acids	~	pathogenicity	
⊙ and ○ or	Phenotype	for 💿	all		– Reduced virule	
		0	Loss of pathogeni	city 🔼	Only genes	
			Reduced virulence Unaffected pathor	e senicity	validated by gen	
⊙ and ○ or	Experimental evidence	for 💿	all	,,	valluated by gen	
		Õ	gene disruption		- alsruption	
			+ gene mutation	have stavia ad	experiments	



Integrated phenotype and comparative genome information



Fusarium graminearum (*microscopic fungus*) genome

101000	18 2 4	ad Water ort Charges Spitzeen			
A1 *	& CLA	for by cond			
A	81	C. C.	D	E	
Janfer fer nend	18383 [BC	10_14516(
	20448 68	s CoAlt gelated PTH22	[042628 (UPHOT) PH 584 (PH) A4889887 (EMBL)]	Loss of pathogenecity (Pheno)	
	19666 [P1	H2, 90	[CK2828 (UPHOF) PHE120 (PH) A4688887 (EMBL)]	Loss of pathogeneoty ("tienc)	
Cluster for sould	171001100	10 14030			
	2010223 04	11 102 1020	C PARETO JENE AAUTEON JENES CONNER A RECOTO	Deduced strakenes (Thereit)	
	20126 (21	R. PHERIC	[PHI 33 (PHI PADENI A PROTI ANAGENIC (EMEL)]	Polycel visionce (Phone)	
Cluster for seet	20061.00	19,00046]			
	10291 ((0)	19_032201			
	20563 [11	19, uni006263	[Parts22 [Part]]	Unaffected pathogenicity (Phena)	
	1000 000	19_13386			
Claster for need	12378 / 00	10 14527			
	20120 / 308	L ACA1	LAMONTO (EMIL) OFOR YO LUPIO'S PHEAD (MED)	Loss of authorshills (Phenel)	
	20002 [15]	9, CAP1	(AAD42019 (EMBL) FHE213 (FHE06ABF9 (UFRO?))	Loss of pathoganicity (Phane)	
	-	an nature			
ACCESSION FOR Dealth	2000 00	414 8000	C BALLETO BEAU OFFICE CARDON A REPORTS BEAUTY	Loss of a thready in filtered	
		0.000	A CAREFOR A REPORT THAT FOR THE A REAL PROPERTY AND A REAL PROPERTY.	Loss of a shortened by randy	
	20000 1 00	C2, 3941	LAADDOSD (EMDL) PHEADL (PHE OLIVE) (UPPOT)	Underted pathematical Press)	
Cluster for poed.	6497 [80	30,09245			
	202031.56	r.PDEN	L 09C214 (LPROT) AM07743 (EMBL) PHL435 (PH0)	Peducet visitince (Ptena)	
Cluster for seed	SP (B)	19 (F19)			
	20084 (27		(PRE304 (PH) AAP12366 (EMBL) OBE2P4 (UFROT))	Reduced visilence (Pheno)	
	20082 [28	P, CPSN	[ORC322 (LERIOT) AASICSTREE (EMEL) FEE 203 (FEE)	Reduced visikince (Ftens)	
	2045 I M	00 00443 30871	LED. WHEN MARKS AMERICA LIPPORT PROPERTY AND	Deduced studence (Therei)	
		a part of parts	Construction of the second sec		
Sunfer for need	15683 (60	10_11096			
	76752 (60	10_12901		la ne constante a series	
	2009 [M	30_11671,1070]	I E1403349 (EMBL) A4ROOB (URROT) PHEB76 (PHD)	Feduced visikince (Pheno)	
Cluster for oread	14150 (00	10 091531			
	20243 (40	A MOSA	(PHLMF (PHLCAG1947 (EMBL) 05020 ((PR07))	Peduced visikance (Phene)	
Chester for soled	15.61 EC	10,14495			
	24,61,83	01.6/5	Exercise bank even point (Eweld) oppose (Period)	reduced visitince (rhend)	
Cluster for seed	961100	15 (1980)			
	30378 (568	5,504, HLAP1]	[PRL4D1 (PHI) AAU14030 (EMEL) 0501P9 (LEWOT))		
Classies for need	8795 [80	10 00000			
	79000 (51	P7951	[PRI 121 [PHI CA2621 (LERIOT) AA888888 (EMBL)]	Loss of pathagenicity (Phene)	

tab separated text file of clusters loaded in Excel



Ondex front-end





Data integration - Summary

- Choose and import data of interest
- Address semantic and syntactic challenges
- Integration methods generate new merged data
- Need to analyse the results of integration
- $\circ \rightarrow$ Network Analysis





- Data integration
 - Needs for data integration
 - Benefits of data integration
 - Challenges
 - Solutions
 - Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation





Everything is a network











- Data integration
 - Data types and sources
 - Benefits of data integration
 - Challenges
 - Solutions
 - Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation





Biological networks can be described using Graph Theory

- \circ A graph G = (V, E) is
 - a set of vertices V (nodes)
 - a set of edges *E* (links between the nodes)



- $V = \{1, 2, 3, 4, 5, 6, 7\}$
- $E = \{ (1,2), (1,3), (2,3), (3,5), (5,4), (4,6) \}$





Biological networks can be described using Graph Theory

O Graphs can be :
Undirected

Directed



 Graphs can also be more complex: multigraphs, bipartite graphs, hypergraphs





o Node degree

Is the number of edges of a node. In a directed graph, we can also define the in-degree and out-degree.

Adjacent nodes

Are nodes joined by an edge

o Path

Is a sequence of adjacent nodes between two nodes







Connected graph

If there is a path between all nodes in a graph

o Distance

Between two nodes *u* and *v* is the length of the shortest path between *u* and *v* (can be infinite)

o Graph diameter

is the maximum distance in a graph



Biological networks can be described using Graph Theory

Clustering coefficient

characterizes the overall tendency of nodes to form clusters or groups

o Centrality measures

Identify "most important" nodes in a graph (hubs). The importance can be assessed for example using :

- the number of connections (degree centrality)
- the proximity (closeness centrality)
- the number of paths (betweenness centrality)





Social Networks

Zachary's karate club; J. Anthropological Res. 33, 452-473 (1977)

 34 members of a karate club

 Visual inspection shows some substructure



 Some very 'popular' individuals (hubs)





Social Networks

Zachary's karate club; J. Anthropological Res. 33, 452-473 (1977)



Degree: 34, 1, 33, 2



Closeness: 1, 3, 34, 32





Data from http://www-personal.umich.edu/~mejn/netdata/

Biological networks can be described using Graph Theory

o Clique

is a subset of nodes such that every pair of nodes in the subset is joined by an edge







- Data integration
 - Needs for data integration
 - Benefits of data integration
 - Challenges
 - Solutions
 - Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation





Existing data visualisation tools

Review

Open Access

A survey of visualization tools for biological network analysis Georgios A Pavlopoulos*, Anna-Lynn Wegener and Reinhard Schneider

o <u>http://www.biodatamining.org/content/1/1/12</u>

 Medusa, Cytoscape, BioLayout Express 3D, Osprey, Proviz, Ondex, PATIKA, PIVOT, Pajek





Data integration

- Network analysis
 - Definition
 - Visualising biological networks
 - Layouts
 - Case study
 - ➤ Filters
 - Annotation





Data integration

- Network analysis
 - Definition
 - Visualising biological networks
 - Layouts
 - Case study
 - ➤ Filters
 - Annotation





Visualising biological networks e.g. PPI network




- topology \rightarrow shape \rightarrow geometry
- o Draw 2-dimensional representations
 - Trees
 - Directed acyclic graph (DAG)
 - General graph (e.g. force-directed)
- Case study in Ondex
 - bioenergy crop improvement





Data integration

- Network analysis
 - Definition
 - Visualising biological networks
 - Layouts
 - Case study
 - ➤ Filters
 - Annotation





Bioenergy Crop Improvement Candidate Gene Identification







ROTHAMSTED

List of candidate genes linked to biological processes



Integrating Relevant Data Sources

- Poplar (JGI) and Arabidopsis (TAIR) genomes
- Linking genes to KEGG, AraCyc, GO, MEDLINE
- Arabidopsis Hormone Database for trait/hormones/genes relations







Command console v1.0 Type 'help;' for more informatio ROT ONDEX>

RESEARCH



ROTHAMSTED

Genomic View in Ondex





Network View in Ondex



ROTHAMSTED



Data integration

- Network analysis
 - Definition
 - Visualising biological networks
 - Layouts
 - Case study
 - Filters
 - Annotation





Integrating different datasets

 → large resulting graph
 Need to narrow down
 Select meaningful areas of the graph

o Example in Ondex

protein-protein interaction network







Protein protein interactions measured using quantitative techniques → Relations on graph have confidence values (confidence) → Threshold filter







- Data integration
 - ➤ Needs for data integration
 - Benefits of data integration
 - Challenges
 - Solutions
 - Case study
- Network analysis
 - Definition
 - Visualising biological networks
 - Annotation





Annotators (1/3)

 Visualise concepts and relations using their attributes/properties

- Colour
- Shape
- Size









ROTHAMSTED

Wild-type FGSG_9908 PH-1 Pkar PH-1



Annotators (2/3)

Virtual Knock-out

- Annotator to see how important a single concept is to all possible paths contained in a network
- Ondex resizes the concepts based on this score
- Scale Concept by Value
 - Pie charts
 - Up/down regulation is indicated in red/green



Mapping microarray expression data to integrated pathways







Annotators (3/3)

Run network statistics such as:
Connectivity
Centrality
Clustering
Network diameter

 \rightarrow Add annotation to the graph







Arabidopsis PPI network

- Calculate
 betweenness and
 degree centralities
- Identify '*influential*' nodes with low degree and high betweenness
 - yellow, green nodes
 - communication between functional modules?









Dynamic process







Investigating plant systems Conclusions (1/2)

Data integration

- Integrate various heterogeneous data sources
- Unify and consolidate knowledge
- Overcome sparseness of data
- Is difficult
 - \rightarrow semantic and syntactic challenges





Investigating plant systems Conclusions (2/2)

Network analysis

- Once data is integrated \rightarrow large volume
 - ➔ need to filter down to regions of interest
- Use various layouts to study merged data
- Annotate data with network statistics to analyse further
- \circ Ondex = 1 solution









o Catherine Caneveto Michael Defoin Platel

Rothamsted members:

- Keywan Hassani-Pak
- Matthew Hindle
- Shao Chih Kuo
- Artem Lysenko
- Chris Rawlings
- Mansoor Saqi
- Andrea Splendiani
- Jan Taubert

Former member:

Jacob Köhler

Biological collaborators:

- Kim Hammond-Kosack
- Martin Urban
- Dimah Habash
- David Wild
- Katherine Denby
- Roxane Legaie

