Visualization and Analysis of FLAREX Gene Expression Data with VANTED

Christian Klukas

Network Analysis Group

Leibniz Institut für Pflanzengenetik und Kulturpflanzenforschung (IPK)

> PGRC Progress Seminar (2005-11-24)





VANTED - Visualization and Analysis of Networks containing Experimental Data Overview

- Motivation
- System overview
- VANTEDs Features
- Demo
- Discussion

Motivation

Massively-parallel techniques generate more and more data

- > A top-down view on the biochemistry of a organism is made possible
- The amount of work needed to evaluate the data increases
 - New tools need to be evaluated or developed
- Goals
 - ☑ Show large amounts of data in a readable and understandable form
 - ☑ Consider related networks
 - Fast data evaluation with the help of statistic functions like t-test or correlation analysis, and clustering algorithms

System overview Roots of VANTED



System overview Related Tools and Devices



SOAP Access to FLAREX

- API designed together with Andreas Stephanik and Matthias Lange, implemented by Karl Spies
- Contains 19 methods
- Sequence Diagram (Example for access to not-normalized spot-intensities):



Key Features

Data sources

- Measurement data
 - DBE-Database (→ VANTED-DB)
 - FLAREX (Array experiment database at the IPK)
 - Excel Files (VANTED-template)
 - Text Files (J-Express format)
- Pathway data (GML, Pajek-.NET, SBML)
- Data transformation and evaluation
 - t-test, U-test, Pearson- and Spearman correlation, SOM-data clustering, various layout commands, search and filter operations, extensible with script commands
- Data export
 - Image files (JPG, PNG, PDF, SVG)
 - Print out
 - Graph files (GML, Wilmascope-.XWG, DOT)

Loading of experimental data



Data-Visualization



Statistic Tests

- Analyze data samples...
 - Check for normal distribution
 - ☑ David et al. quick test
 - Chi-square test
 - Detect/Remove outliers
 - Grubbs test
 - Detection of significant mean differences with
 - ✓ t-test (2 variants)
 - ✓ U-test (rank-sum test)





Correlation Analysis (1/2)

- Calculation of the Pearson (linear) or Spearman (rank-order) correlation
- Detection of correlations, shifted in time:
 - Repeated correlation calculation (r_i) for multiple timeoffsets (i=-3...3, t₋₃...t₊₃)
 - Using max|r_i| for data visualization
- Test of significance with approximation to the tdistribution





Correlation Analysis (2/2)



Summary & Outlook

Website

- http://vanted.ipk-gatersleben.de
- Publications
 - Borisjuk, Hajirezaei, Klukas, Rolletschek, Schreiber: Integrating data from biological experiments into metabolic networks with the DBE information system. In Silico Biology (2004)
 - Rolletschek, Radchuk, Klukas, Schreiber, Borisjuk: Oil storage in soybean seeds: evidence for a key role of photosynthetic oxygen release. New Phytologist (2005)
 - BMC Bioinformatics ?
- Outlook
 - Improve analysis and visualization of array data
 - Based on discussions with colleagues and feedback from users of the system
 - Display of simulated experimental data, generated with SyBME

Thanks to Colleagues

 Ljudmilla Borisjuk, Mohammad-Reza Hajirezaei, Björn Junker, Hardy Rolletschek, Nese Sreenivasulu, Winfriede Weschke, Ruslana Radchuk

> Discussion of system features and data provision

- Matthias Lange, Uwe Scholz, Andreas Stephanik, Karl Spies
 Database services and SOAP access to FLAREX
- Dirk Koschützki, Falk Schreiber
 Support and Inspiration

Software Demo

1. Loading of Array Data

- 23 Hybridizations:
 - Wild type / trangenic line
 - 6 time points
 - ~ 146000 measurement values
 - ~ 3800 clones
- 2. Filter Dataset
 - **Remove 50% of the data with highest average standard deviations**
 - Remove 90% of remaining nodes with low differences between wild type and transgenic line
- 3. Detect Common Developmental Patterns (SOM-Clustering)





NB V	INTE	D BE IA															کالاتار
Eile	Edit	Analysis	Edges (Graph	Layout	Nodes	<u>W</u> indow	Help									
₿.	6	06	× •		ə e	Q	€ @	Q	<u>001</u> -0}-	바 후 뤽 팬	Search:				ок		
												1	Experiments	KEGG	Statistics		
												8	DBE-Databas	se .			4
													Flarex (in de	velopme	ent, not working)		4
															List Providers		Login
													User 19: Ruslana Ra CITY: Gatersl COUNTRY: G EMAIL: ruslai HYBRIDISAT STREET: Cor TEL: ++49 39 Hybridise NAME: 24_ Hybridise NAME: 34_ Hybridise NAME: 24_ Hybridise NAME: 24_ Hybridise NAME: 24_ Hybridise NAME: 24_	adchuk eben iermany har@ipk IONS: 3: rensstr. 1 482 5264 tion 19_1 tion 19_2 tion 15_1 tion 15_2 tion 15_2 tion 17_1 tion 17_2 Gene Ir t Data f	-gatersleben.de -gatersleben.de 8 85: 86: 87: 88: 89: 90: htensities for Hybridisation IDs 85, 86, 87, 88, 89, 91), 91, 92,	
													Please Walt	(up to	a rew minutes)		iton
Get D	ata fr	or Hybridi	sation ID	85 (1))	38)												otop
				(_					

10

Ve V	ANTE	D BETA																			×
File	<u>E</u> dit	Analysis	Edges G	raph	Layou	ut No	odes <u>V</u>	<u>/indow</u>	Help												
[≱	6	86	20	0	9	e (e e		Q	<u>001</u>	0) <u>0</u> 1	冒	후 릐	Search:				ок			
															1	Experiments	KEGG	Statistics			
																DBE-Databa	se			*	7
																Flarex (in de	evelopme	ent, not workin	g)		7
																Load Input F	⁼ile		(0	7
																					_
				00 - H.M.			a an as		141		-	105									_

Welcome to VANTED - Visualization and Analysis of Networks containing Experimental Data! In the Help menu you find a tutorial section which quickly gives an overview of the various features of this application. Furthermore you will find [?] buttons throughout the system which point directly to the topics of interest. If you experience problems or would like to suggest enhancements, feel free to use the Send feedback command in the Help menu!

Compose Experiment Dataset 🛛 🔀										
Experiment Name*	[Enter Experiment-Name]									
Remark		1								
Coordinator*	Ruslana Radchuk									
Experiment Started	•									
Time Points	Time Points									
Time Unit	days after pollination									
Plant Names*	Plant Names 🔷 💎									
Genotypes*	Genotypes 🔷									
Select Identifier*	IPK Gene Identifiers 🛛 👻									
	* values must be set									
OK Cance	9									

Compose Experi	ment Dataset		×	
Experiment Name*	Enter Experiment	-Nam	e]	
Remark				
Coordinator*	Ruslana Radchuk			
Experiment Started				
	Time Points	•	4	
	Hybridisation 85	13		
	Hybridisation 86	13		
	Hybridisation 87	15		
	Hybridisation 88	15		
	Hybridisation 89	17		
	Hybridisation 90	17		
	Hybridisation 91	19		
	Hybridisation 92	19		
	Hybridisation 93	11		
Time Points	Hybridisation 94	11		
	Hybridisation 95	13		
	Hybridisation 96	13		
	Hybridisation 97	15		
	Hybridisation 98	15		
	Hybridisation 99	17		
	Hybridisation 100	17		
	Hybridisation 101	19		
	Hybridisation 102	19		
	Hybridisation 103	21		
	Hybridisation 104	21		
Time Unit	days after pollina	tion		
Plant Names*	Plant Names		*	
Genotypes*	Genotypes		٠	
Select Identifier*	IPK Gene Identifie	rs	~	
	* values must be s	et		
	el 📄			

VB V	ANTE	D BETA												$[\times]$
<u>F</u> ile	<u>E</u> dit	Analysis	Edges	Graph	Layout	Nodes	<u>W</u> indo	∾ <u>H</u> el	D					
Ľ,	6	88	% (9 e	Q	•		<u>00+</u> -00		🗐 Search:			Ок
								Exp	eriments	KEGG Statis	itics			
								DE	E-Databa	se				*
								Fk	irex (in de	velopment, no	t working)			*
								Lo	ad Input F	File			۲	*
								Cr	ate Data:	set				
								Pro	cessing Da	ata				
								Ple	ase wait	,				
_												l	Stop	
Proce	essin	g data:72	948 valu	les fron	n 19000	sample	es (50%), 3823	l substar	nces				

[Enter Experiment-Name]

Close this Tab

Perform Data Mapping

 Δ

~

~

(if no graph window is open, a node grid will be created in a new editor window)

Experiment Info

Experiment-Name: [Enter Experiment-Name] Remark: Coordinator: Ruslana Radchuk Measurement values: 146092 Import time: Tue Nov 22 11:19:57 CET 2005 Experiment started: Tue Nov 22 11:19:43 CET 2005

Specify Mapping-Data

Time Points

days after pollination 11 days after pollination 13 days after pollination 15 days after pollination 17

days after pollination 19

Plants/Genotypes

34 ([Enter Genotype]) id=1

WT ([Enter Genotype]) id=2

Specify Mapping-Options

Map empty datasets

Useful if more than one (incomplete) dataset will be mapped onto the nodes)

Create new nodes for datasets that can otherwise not be mapped onto the active graph/nodes

Ask for user-given mapping if no automatic mapping is possible.

Data Mapping Task

Map XML data for substance PSC12K03 to node Node ID=1129

Add new graph node finished

Stop

4 min

2

.

[* 🗃 🖫 💑 🖸 👔 🤰 C 🔍 C 🔍 C 🔍 C 🔍 U 🕂 町 臣 冬 릐 Search:

ALANAL

OK

>

BSH Script

Calculate Average Standard Deviation of Samples

- Goal:
 - Remove 50% of the data with highest average standard deviations
- 1. Calculate average standard deviation for all clones (graph nodes)

```
//@Nodes:Calculate Average Sample StdDev§
series = node.getMappedSeriesData();
stddevs = new ArrayList();
for (SeriesData sd : series)
    stddevs.addAll(sd.getStdDevValues());
double sum = 0;
int i=0;
for (Double stddev : stddevs)
    sum += stddev;
node.setAttributeValue("script", "avg_stddev",
    new Double(sum/stddevs.size()));
```

 $s = \sqrt{\frac{\sum_{i=1}^{n} (X_i - \overline{X})^2}{1}}$

The standard deviation value represents the average distance of a set of scores from the mean.

BSH Script

Calculate Average Ratio of Sample-Mean Differences

- Goal:
 - Identify clones with high differences between wild type / transgenic line
- 1. Calculate ratio-difference for all clones (graph nodes)

BSH Script Remove 50% of clones with highest average standard deviation

BSH Script Remove 90% of clones with low average ratio between WT / transgenic line

VANTED BETA





23

₹ª

- P 🛛

<u>File E</u>dit Analysis Edges Graph Layout Nodes <u>W</u>indow <u>H</u>elp

[4 🗃 🖪 🖁 👗 № 💼 ② C C C C C C C Q C L O F F P キ 引 Search:

👙 demo_pgrc_all_data.gml<u>* - view 1</u>



OK



2

E.

🕜 🔶

@ 4

4 🗘

0 🗘

10 🗘

Eile Edit Analysis Edges Graph Layout Nodes <u>Wi</u>ndow <u>H</u>elp

PSCNLI

PSCIIKI

PSC2963

PSCHO

AVE

[🔄 🖪 🖶 👗 🖪 🔵 🧶 🔍 🔍 🔍 🔍 🖳 🖉 🖷 🗦 🚔 Search:

ME

PSC21L

PGC23

👙 demo_pgrc_all_data.gml* - view 1

- 7 🛛

OK

PSC27P

Experiments KEGG Statistics Network Node ٩ Graph Directed Edges 🔽 Charting (all nodes) IHMM Bar-Outline/Line Thickness Horizontal/Vertical ~ Label Rotation (degree) Series Colors N 0 2 4 6 8 10 Show Category Labels ~ da 2,5 Show Range Labels V T-Test-Marker Size Use SE instead of SD for Error-Bar

Network Attributes

Cluster-Colors

Charting (all line-charts)

Error-Bar Line-Thickness

Shape-Size

Show Error as Fill-Range

Show Error as vert. Line

Image: Show Error as vert. Line<

Apply & Redraw

>



V B	NTED BETA																
Eile	<u>E</u> dit Analysis B	Edges Graph	Layout N	odes <u>W</u> ind	w <u>H</u> elp)		12									
₿.	a 8 B	% 🗅 🌒	9 ¢	q q 6	l Q	. 에 제 바 비) 립 Sea	rch:		ок							
23	🛔 demo_pgro	c_final.gml	view 1								- 7 🛛	Experimer	nts KEGG Stati	istics Ne	twork Node		
₹"					K K						^	Graph					0 🍝
				10	1							Directe	d Edges 🔽				
												Chartin	g (all nodes)				0 4
												Bar-Ou	tline/Line Thickne	ess	H	1h	4 🧘
				1								Horizor	ital/Vertical		~		
				40								Label R	otation (degree)	i.			0 🗘
				10								Series	Colors		~	~	
												Show C	ategory Labels		0246	5 8 10 day	
												Show F	ange Labels		2,5		
												T-Test-	Marker Size		* *		10 🗘
									4			Use SE	instead of SD for	r Error-Ba	ır 🔲 🗌		
		6					5	AL	TA.			Networ	k Attributes				0 🍝
												Cluster	-Colors				
												Cluster	graph Graph I	ID=7			
			(1 ⁴⁴) (1 ⁴⁴)									Chartin	a fall line-charts	ો			0
					~ ~							c.ndi cili	g (an me churc	-	— т		
					A M							Error-B	ar Line-Thickness				4 🗸
			ie ie		a is							Shape-	Size		, "The law		6 🗘
			LA LA	200								Show E	rror as Fill-Range	• 🗹 🕨	\sim		
				I II		Ā Ā Ā						Show E	rror as vert. Line	• 🗖 I	T.A		
			1			I I I						Show L	ines		J!		
			N	4	the second	Contraction of Contraction						Show S	hapes	•	1 /		
											~		De de un				
s - 18	<	_	_	_	_		_			_	>	Apply &	Redraw				160

160 nodes no edges

<

	PSC31N14	PSC30K22	PSC27814	PSC30F17	PSC34F09	PSC23L24	PSC28K19	PSC27409	PSC24011	F.
	\wedge	A		1	X	F	A		A	
	PSC27812	PSC28P08	PSC23322	PSC34815	PSS15P04	P5C33H05	PSC25114	PSC30003	PSC23110	
-	X	X	N	A	1	17	F	17	X	
	PSC35307	PSC31C11	PSC25L1C	PSC28002	PSC24L17	PSC34015	PSC25H03	PSC2582C	PSC29N11	1
	A	1			X	A	1	1	A	
	PSC28001	PSC34N1:	P5508H0:	P\$\$05322	PSC22H10	PSC31K05	PSC33M23	PSC34G24	PSC28N1:	
	Ne	-	15-		A		>	A	A	
	PSC34J15	PSC26M18	PSC24L16	PSC29124	P5C21L15	PSC22C22	PSC22P01	PSC27P13	PSC26C24	5
	H	F	JE			A		H	A	
	PSC28E01	PSC33F05	PSC27N04	PSC22M12	PSC30H24	PSC26E07	PSC27K18	PSC33M17	PSC31K12	
	1×	A	1	A		1	A	1	1	
	PSC34004	PSC30P22	PSC28004	PSC24016	PSC25P03	PSC28L23	PSC21813	PSC24012	PSC33C10	
	1	M	1	1	1	×	×		1	
	PSC26L04	PSC32P06	PSC23022	PSC21824	PSC23320	PSC33007	PSC33M14	PSC25L18	PSC32112	1
	1	1	1		X	A	A	F	1	
	PSC27H08	PSC33E02	PSC23A24	PSC31105						
	1	T	1	A						







