



LabKey Software™
PARTNERS IN SCIENCE

LabKey Software: Past Present & Future

Mark Igra

Partner, LabKey Software
marki@labkey.com

- Why Scientific Data Integration
- Evolution of LabKey
 - Software
 - Company
 - Common requirements for scientific data integration
- Where we are now
- Future Directions

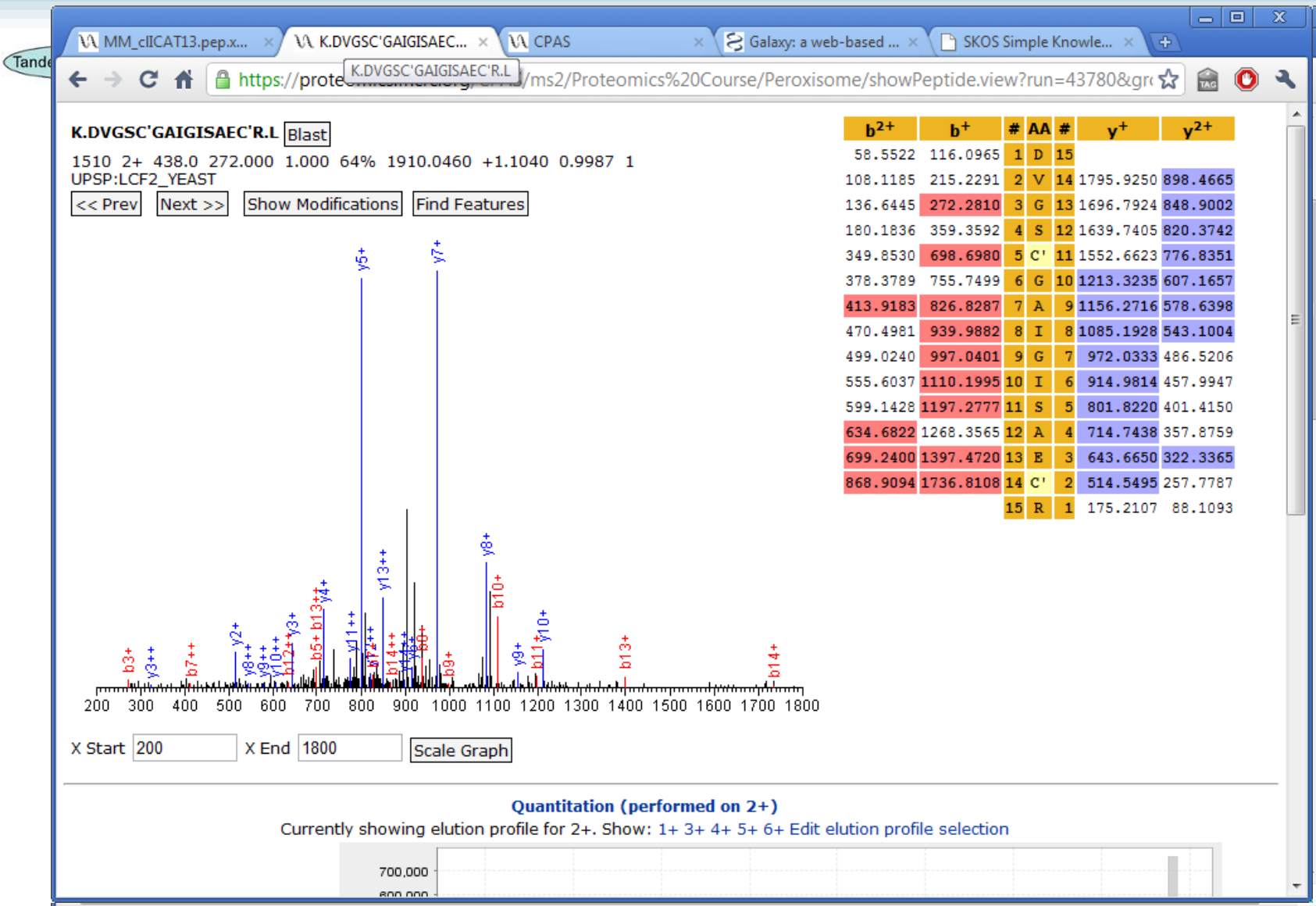
- Data Volume
 - Hundreds of millions of results from thousands of high throughput assay runs
- Data Variety
 - Clinical, Demographic, Assay & Specimen Data
 - Scientific data annotation
- Collaboration
 - Investigators aren't all in the same place
 - Need secure data sharing
- Broad Range of Analyses
 - Queries, Reports, Domain Specific Tools

- Find Cancer Biomarkers using Proteomics
- Problems Faced
 - Number of reads from MS2-based proteomics assay were exploding & difficult to handle on existing tools
 - Analysis Pipeline on cluster was difficult to optimize



- Solution
 - The original LabKey Server (CPAS)
 - Pipeline to run MS2 analysis jobs via web browser and load results into Database
 - Web based analysis tools to view, combine & share results
 - Lots of Data
 - More than 90,000 MS2 Runs
 - More than 700,000,000 peptide identifications
- Rauch et al, Journal of Proteome Research, 5/2006

Proteomics Examples



Software Requirement Summary

	Proteomics				
Data Size	700M Peptides 90K Runs				
Data Pipeline	✓				
Collaboration	✓				
Specimens	✓				
Security	✓				

- LabKey Software becomes independent company



- Measure intracellular cytokines in many samples
- Per-run quantity of data increasing
- Need consistent analysis within an experiment
- Cross-run analyses
- Each experiment type may have different statistics
















- Solution: LabKey Flow
 - High-throughput flow analysis engine
 - Loading of flow statistics
 - Adaptable data model for varying analysis types
 - Query tools for analyzing data
- Shulman et al, Cytometry A, Sep 8 2008

Flow Dashboard > labkey-analysis >

labkey-demo.xml analysis

Run Comment:

SHOW GRAPHS ▶

QUERY ▼	VIEWS ▼	EXPORT ▼	PRINT	PAGE SIZE ▼					
	Name	Flag	S:Count	S:%P	4+:Count	8+:Count	4+/(!IFNg+&!IL2+):Count	4+/(!IFNg+&!IL2+&IL4+&!TNFa+):Count	4+
DETAILS ▶	119142.fcs		9,764	97.64	3,713	1,917	3,711		1
DETAILS ▶	118813.fcs		9,770	97.7	2,357	1,888	2,357		1
DETAILS ▶	118981.fcs		9,807	98.07	2,338	1,848	2,337		47
DETAILS ▶	118947.fcs		9,555	95.55	2,065	1,703	2,065		2
DETAILS ▶	119199.fcs		9,728	97.28	2,336	1,907	1,977		0
DETAILS ▶	119154.fcs		9,545	95.45	2,541	1,629	2,541		0
DETAILS ▶	118836.fcs		9,496	94.96	504	466	503		0
DETAILS ▶	119010.fcs		9,817	98.17	3,804	1,824	3,801		0
DETAILS ▶	118762.fcs		9,743	97.43	5	131	5		0
DETAILS ▶	118756.fcs		9,753	97.53	0	0	0		0
DETAILS ▶	119171.fcs		9,662	96.62	2,200	2,196	1,661		5
DETAILS ▶	118801.fcs		9,706	97.06	2,251	2,087	2,251		2
DETAILS ▶	118754.fcs		9,787	97.87	1	0	1		0

View Graphs

Flow Dashboard > labkey-analysis >

labkey-demo.xml analysis

Run Comment:

HIDE GRAPHS ▾ [Large Graphs] [Medium Graphs] [Small Graphs]

QUERY ▾	VIEWS ▾	EXPORT ▾	PRINT	PAGE SIZE ▾	1 - 68 of 68					
	Name	Flag	S:Count	S:%P	4+:Count	8+:Count	4+/(!!IFNg+&!IL2+):Count	4+/(!!IFNg+&!IL2+&IL4+&!TNFa+):Count	4+/(!!IFNg+&!IL2+&IL4+&TNFa+):Count	4+/(!!IFNg+&IL2+
DETAILS ▸	119142.fcs		9,764	97.64	3,713	1,917	3,711	1		0
	<div> <div>S/Lv/L</div> <div>S/Lv/L/3+</div> <div>S/Lv/L/3+/4+</div> <div>S/Lv/L/3+/4+</div> </div>									
	<div> </div>									
DETAILS ▸	118813.fcs		9,770	97.7	2,357	1,888	2,357	1		0
	<div> <div>S/Lv/L</div> <div>S/Lv/L/3+</div> <div>S/Lv/L/3+/4+</div> <div>S/Lv/L/3+/4+</div> </div>									
	<div> </div>									

- Sort Ascending
- Sort Descending
- Clear Sort
- Filter...
- Clear Filter

Customize View

Flow Dashboard > labkey-analysis >

labkey-demo.xml analysis

Run Comment:

HIDE GRAPHS ▶ [Large Graphs] [Medium Graphs] [Small Graphs]

QUERY ▼ VIEWS ▼ EXPORT ▼ PRINT PAGE SIZE ▼

1 - 68 of 68

Columns Available Fields

- Filter
- Sort
- ☒ Flag
 - ☐ Created
 - ☒ Run
 - ☐ Analysis Script
 - ☐ Compensation Matrix
 - ☐ Statistic
 - ☐ Graph
 - ☐ (<APC-A>)
 - ☐ (<Alexa 680-A>)
 - ☐ (<FITC-A>)
 - ☐ (<PE Cy55-A>)
 - ☐ (<PE Cy7-A>)

☐ Show Hidden Fields


Editing an unsaved view.

DELETE

REVERT

VIEW GRID

SAVE

	Name	Flag	S:Count	S:%P	4+:Count	8+:Count	4+/(!IFNg+&!IL2+):Count	4+/(!IFNg+&!IL2+&IL4+&!TNFa+):Count	4+/(!IFNg+&!IL2+&IL4+&TNFa+):Count	4+/(!IFNg+&IL2+&IL4+&TNFa+):Count	4+/(!IFNg+&IL2+
DETAILS ▶	119142.fcs		9,764	97.64	3,713	1,917	3,711		1	0	
		S/Lv/L			S/Lv/L/3+		S/Lv/L/3+/4+		S/Lv/L/3+/4+		
		<input type="text" value=""/>			<input type="text" value=""/>		<input type="text" value=""/>		<input type="text" value=""/>		

Software Requirement Summary

	Proteomics	Flow			
Data Size	700M Peptides 90K Runs	50M Statistics			
Data Pipeline	✓	✓			
Collaboration	✓				
Specimens	✓	✓			
Security	✓	✓			
Data Variety		✓			
Query		✓			

- Combine many data types for HIV Vaccine studies
- Clinical Response Forms (CRF), Specimens, Many Assays
- Enable secure collaboration for scientists worldwide
- Allocate & distribute valuable specimens



- Solution
 - Secure web portal for HIV Vaccine Enterprise Data
 - Used by several networks to share data
 - CHAVI, CAVD, HVTN, HPTN (3000 Users Worldwide)
 - Core software was written by LabKey
 - SCHARP runs Atlas
 - Defines available data and relationships
 - Manages security and permissions
 - Manages data loading
 - Builds custom modules
- Nelson et al, BMC Bioinformatics, March 2011
- Piehler et al, BMC Immunology, May 2011

Participant View

CHAVI 001 > Study Overview > Dataset: Binding Antibody, All Visits >

Participant - [REDACTED]

[NEXT PARTICIPANT](#) ▶ [SEARCH FOR \[REDACTED\]](#) ▶

[All Datasets](#)

[Specimen Timeline](#)

[Clinical Data](#)

[Requested Specimens](#)

[All Specimens](#)

[Binding Antibodies](#)

[Cytokines](#)

Clinical Data

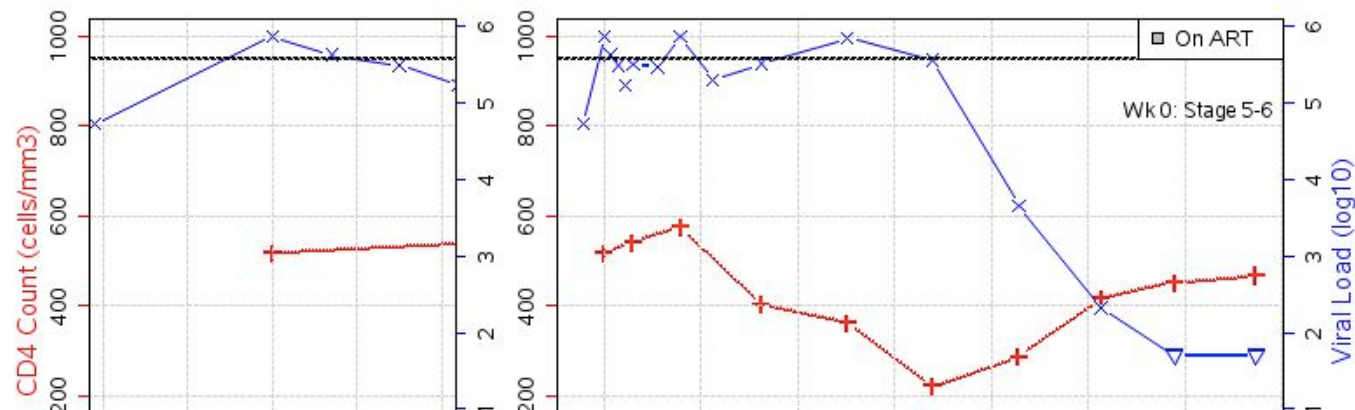
CRF Data

[EXPORT TO EXCEL](#)

Cohort:	Confirmed Acute
Age:	18
Gender:	Female
Race:	White
Tribe:	
Clade:	

▽ Result 'Less Than' Value --- Est Viral Load Set Point
△ Result 'Greater Than' Value

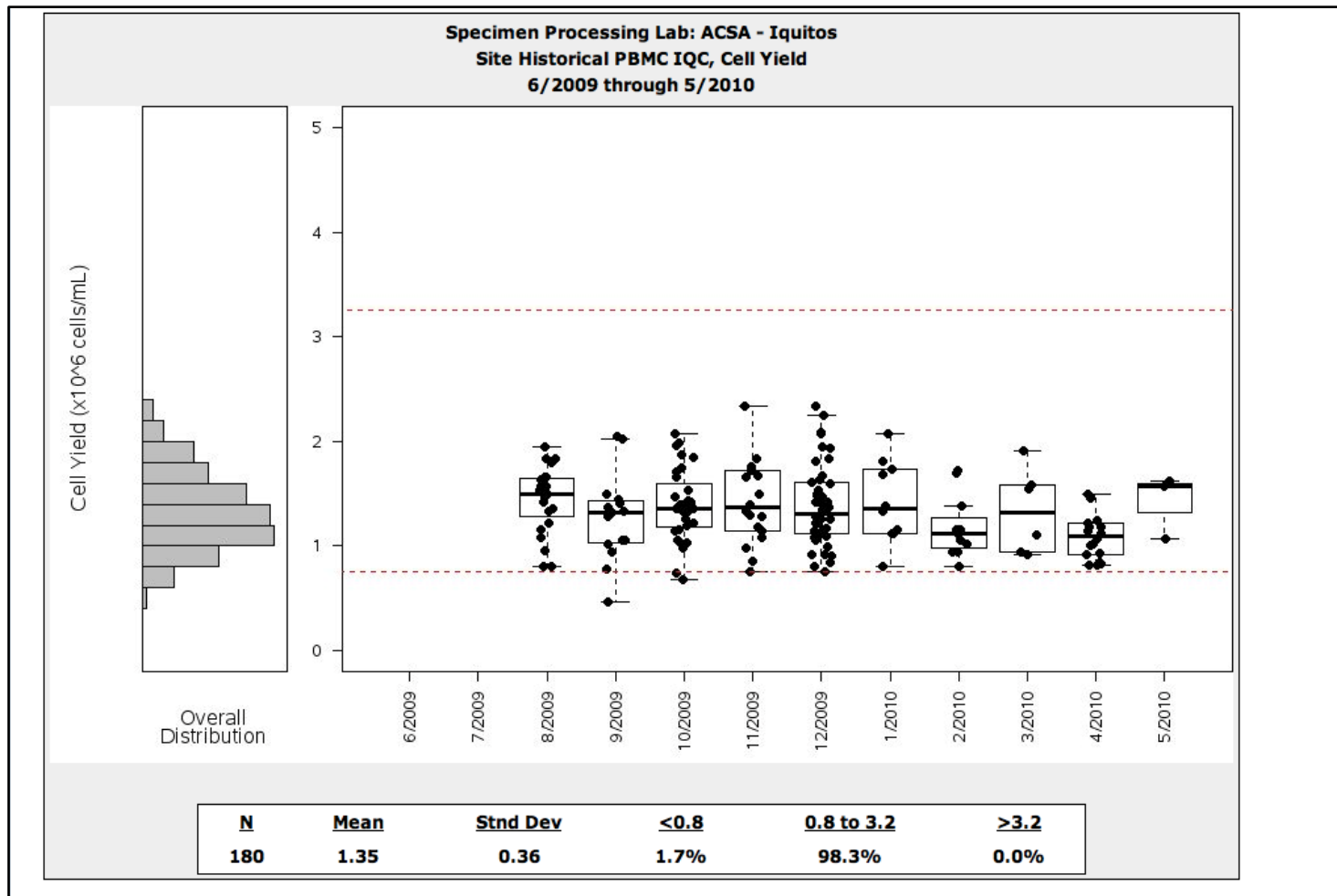
[REDACTED] - Confirmed Acute



- Dozens of small applications needed for different groups
 - Most used data already in LabKey System
 - But custom workflows, reports & analysis required
- SCHARP has detailed knowledge of the requirements
- LabKey provided an API and simple application building tools



- Internal Quality Control for Sample Processing
 - Cell Yield & Processing Time



RV144 Followup Study Tracking System

- Tools for tracking information about studies

Study Tracking System >

Admin

Study Data Entry: Borrow cytokine chemokine

[Study Overview](#) [Assay\(s\)](#) [Team Members](#) [SCHARP](#) [Atlas Admin](#) [Other Tools](#)

Stat Task List

Study Tasks

	Task Description	Date Completed
[edit]	Sampling Plan from SCHARP	28-Dec-2009

Sample Set Tasks

Task Description	Sample Set(s) Used	
Sample set(s) for Borrow Chemokine		Add Sample Set
Sample set(s) for Borrow Cytokine		Add Sample Set

LDO Task List

Assay Tasks: Borrow Chemokine

	Task Description	Date Completed
[edit]	Data Uploaded for Borrow Chemokine	
[edit]	Data Imported for Borrow Chemokine	
[edit]	Data QCd for Borrow Chemokine	
[edit]	Data Sent to Stat for Borrow Chemokine	
[edit]	Analysis Complete for Borrow Chemokine	

Assay Tasks: Borrow Cytokine

	Task Description	Date Completed
[edit]	Data Uploaded for Borrow Cytokine	
[edit]	Data Imported for Borrow Cytokine	
[edit]	Data QCd for Borrow Cytokine	
[edit]	Data Sent to Stat for Borrow Cytokine	
[edit]	Analysis Complete for Borrow Cytokine	

Software Requirement Summary

	Proteomics	Flow	Atlas (Studies)	Atlas (Labs)	
Data Size	700M Peptides 90K Runs	50M Statistics	20K Subjects 800K Specimens 1200 Datasets	30K Assay Runs	
Data Pipeline	✓	✓	✓	✓	
Collaboration	✓		✓	✓	
Specimens	✓	✓	✓	✓	
Security	✓	✓	✓	✓	
Data Variety		✓	✓	✓	
Query		✓	✓	✓	
Study Model			✓		
Custom Reports			✓	✓	
API			✓	✓	
Auditing			✓	✓	

- 30 years of daily data on thousands of animals
- Clinical staff & vets need health record
- Researchers need scientific data
- Colony Management is an ongoing problem

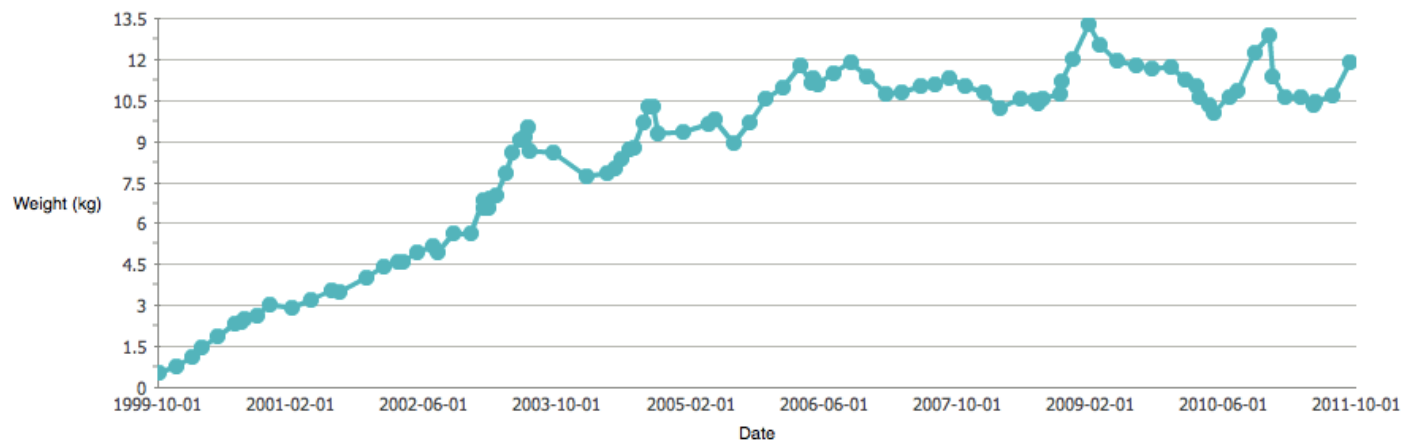


- EHR is a a specialized, ongoing “study”
- LabKey enhanced scalability of study solution
- LabKey enhanced API
- LabKey wrote tools to transfer data into new EHR within minutes of entry into old EHR
 - Now LabKey EHR is only one in use
- Wisconsin Primate Center built custom, forms, views & reports to enter & analyze the data
- Solution extending to Oregon National Primate Center

Problem List:

VIEWS ▾EXPORT ▾PRINTPAGE SIZE ▾MORE ACTIONS ▾											1 - 12 of 12
<input type="checkbox"/>		Id	Current Room	Current Cage	Problem Number	Date Observed ▾	Date Resolved	Category	Remark	Code	
<input type="checkbox"/>	DETAILS ▸				12	2011-10-04	2011-10-07				
<input type="checkbox"/>	DETAILS ▸				11	2011-05-16	2011-05-20				
<input type="checkbox"/>	DETAILS ▸				10	2010-11-22	2010-11-30				
<input type="checkbox"/>	DETAILS ▸				9	2008-05-22	2008-05-23				
<input type="checkbox"/>	DETAILS ▸				8	2007-09-05	2007-09-07				
<input type="checkbox"/>	DETAILS ▸				7	2007-06-05	2007-06-08				
<input type="checkbox"/>	DETAILS ▸				6	2005-01-12	2005-01-26				
<input type="checkbox"/>	DETAILS ▸				5	2003-12-03	2003-12-03				
<input type="checkbox"/>	DETAILS ▸				4	2003-03-17	2003-04-10				
<input type="checkbox"/>	DETAILS ▸				3	2002-11-19	2002-11-19				
<input type="checkbox"/>	DETAILS ▸				2	2002-07-22	2002-07-22				
<input type="checkbox"/>	DETAILS ▸				1	2002-03-07	2002-03-19				
VIEWS ▾EXPORT ▾PRINTPAGE SIZE ▾MORE ACTIONS ▾											1 - 12 of 12

Weight:



Min Date:

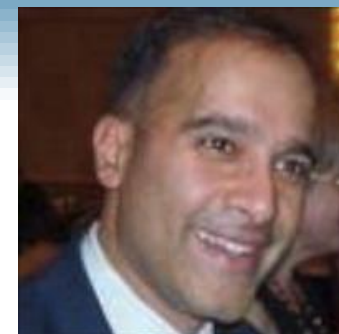
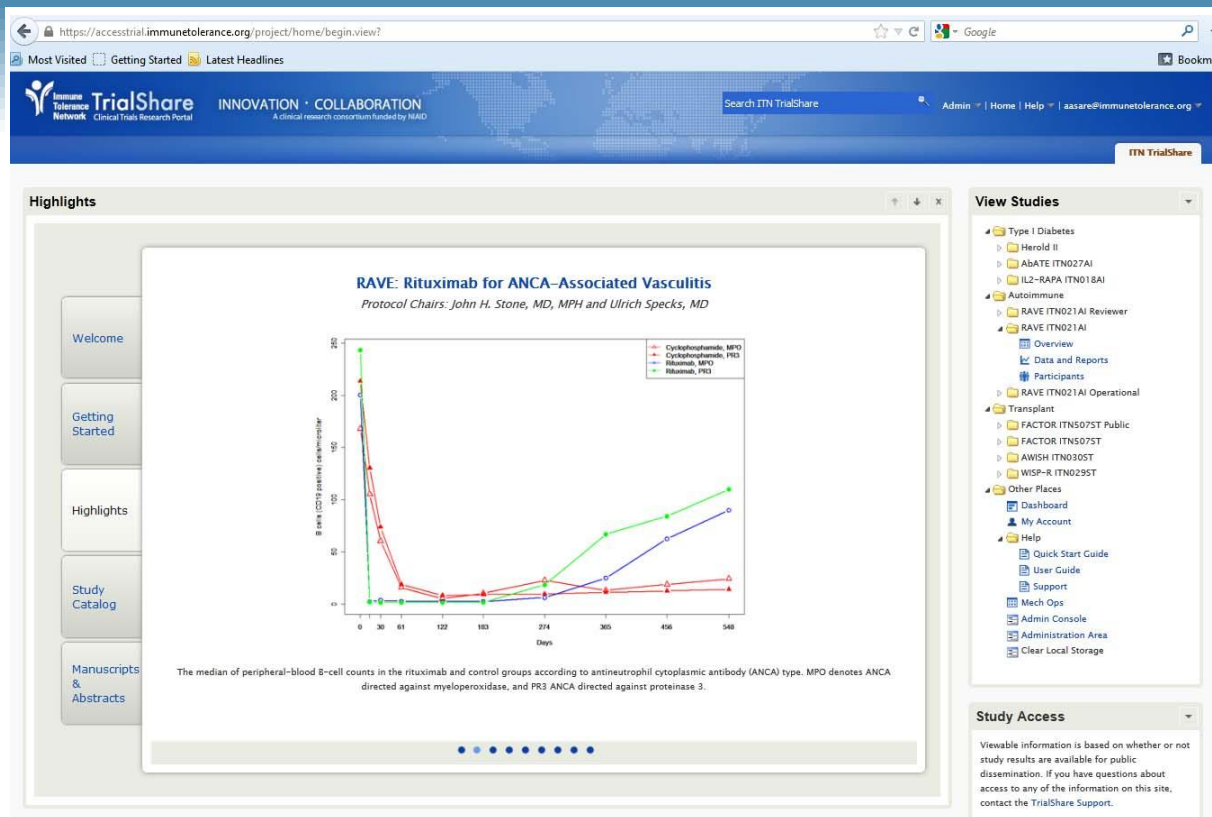
Max Date:

REFRESH

Software Requirement Summary

	Proteomics	Flow	Atlas (Studies)	Atlas (Labs)	Primate EHR
Data Size	700M Peptides 90K Runs	50M Statistics	20K Subjects 800K Specimens 1200 Datasets	30K Assay Runs	13K Animals 1.2M Drug Doses
Data Pipeline	✓	✓	✓	✓	
Collaboration	✓		✓	✓	✓
Specimens	✓	✓	✓	✓	
Security	✓	✓	✓	✓	✓
Data Variety		✓	✓	✓	✓
Query		✓	✓	✓	✓
Study Model			✓		✓
Custom Reports			✓	✓	✓
API			✓	✓	✓
Auditing			✓	✓	✓

2011 – ITN Polish, Analysis



- Two teams
- Every month
 - Demo progress
 - Gather requirements & priorities
 - Write specs
 - Detailed cost estimates
 - Development, automated & manual testing, bug fixing
- Every day
 - Stand up meeting
 - Schedule progress
 - Quality metrics
 - Blocking Issues

What the process looks like

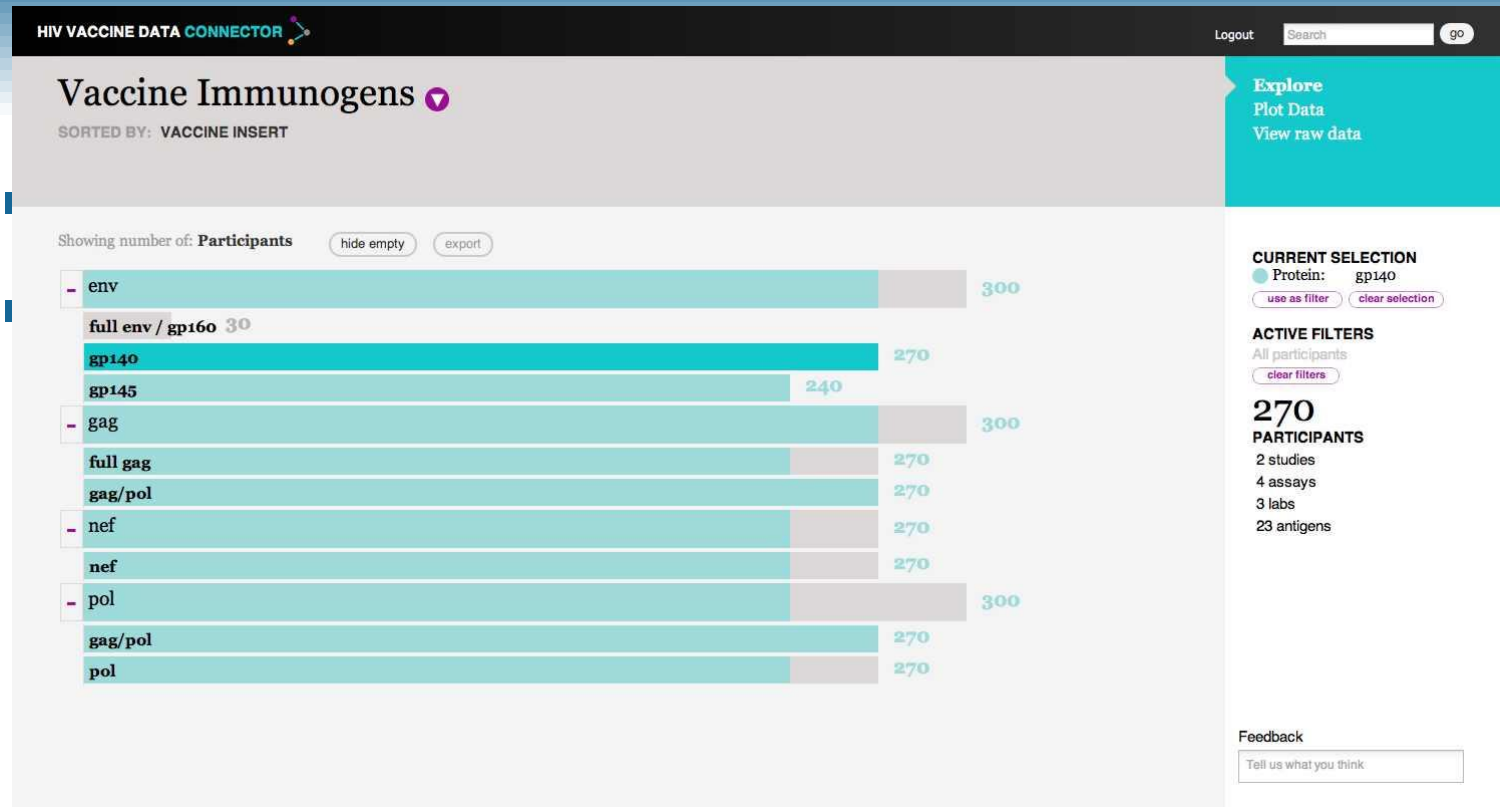
Current Sprint 23: **12.3-2** (10 Sep - 28 Sep)

[Add a Story](#)

Show work for: Everyone

Story	Ignore	Pending	Blocked	Started	Done	+ unstack
2 pts started CDS - Design/investigation for ELISA assay		13.0h KL Placeholder				
5 pts started Dataspace - Reflecting grid filters in the part...		1.0h KF Buddy testing			W f A st G NA Bug fixing	
2 pts started Dataspace - improve measure picker usability		5.0h NA Animation for axis selector panel display/collapse	5.0h Automated tests (add marker class for animations)	5.0h NA Refactor axis selectors to render separately	R NA Hook up axis buttons	
			5.0h NA Bug fixing			
			3.0h Buddy testing			
2 pts started HVTN - Spec for folders as tabs		5.0h BP Figure out plan				
		3.0h MI Write spec				
3 pts started HVTN - Menu Web Part - Customizable			5.0h Automated tests	3.0h DGB Create customize view page	DGB Create new webpart	
			2.0h DGB Bug fixing	1.0h DGB Persist webpart settings		
			2.0h Buddy testing			

feedback



Software Requirement Summary

	Proteomics	Flow	Atlas (Studies)	Atlas (Labs)	Primate EHR
Data Size	700M Peptides 90K Runs	50M Statistics	20K Subjects 800K Specimens 1200 Datasets	30K Assay Runs	13K Animals 1.2M Drug Doses
Data Pipeline	✓	✓	✓	✓	
Collaboration	✓		✓	✓	✓
Specimens	✓	✓	✓	✓	
Security	✓	✓	✓	✓	✓
Data Variety		✓	✓	✓	✓
Query	✓	✓	✓	✓	✓
Study Model			✓	✓	✓
Custom Reports	✓	✓	✓	✓	✓
API	✓	✓	✓	✓	✓
Auditing	✓	✓	✓	✓	✓

- Data integration is the unifying theme
 - Samples to Assays
 - Assays to Subjects
 - Subjects to Studies
 - Translational Medicine involves integrating data from “Molecules to Populations” (Kuhn et al 2008)
- Data types to integrate are constantly evolving
 - Both file and “structured data types”
 - Can’t rebuild new systems for new data types
- Tools grow as fast as data types
- Scientists need to share data
 - Central lab doing work for distributed clients
 - Distributed group of labs contributing
 - Need to manage many users

- Our collaborators guide our direction
 - But we have a few ideas of what's coming
- Networks/Consortia
 - Cross protocol data integration
 - Ancillary Studies
 - Deidentification
- Labs
 - Sequencing
 - Improved Throughput
 - Pipelines

- Developers
 - Dependency & library management
 - Rserve – Build web tools with R back ends
 - Flexibility in display & navigation
- For Everyone
 - Integration with more systems (e.g. IMMPort)
 - Visualization & Analysis
 - Terminology and Identity Management

- Number & complexity of engagements has grown
 - We have evolved from 3 – 22 people in 9 years
 - Processes have evolved to scale
 - Can apply what we've learned to new engagements
- Open source model works
 - Inspires Trust
 - Fits the research market
- We're happy with our market, our growth & our business
- Thank You

Acknowledgements

- **Partners & Funders**
- Martin McIntosh (FHCRC)
 - NCI
 - Canary Foundation
- Steve Self (SCHARP)
 - CHAVI (NIH)
 - HVTN (NIH)
 - CAVD (Gates Foundation)
- David O'Connor (Wisconsin)
 - Primate EHR (ARRA)
 - Genotyping Tools (NIAID)
- Parag Mallick (USC)
- Michael Katze (UW)
- Immune Tolerance Network
- **LabKey Development**
 - Matthew Bellew
 - Adam Rauch
 - Britt Piehler
 - Josh Eckels
 - Kevin Krause
 - Brendan MacLean (now UW)
 - Nick Shulman (now UW)
 - Karl Lum
 - Nick Arnold
 - Cory Nathe
 - Ben Bimber (now ONPRC)
 - Alan Veniza
 - Dax Hawkins
 - Dave Bradlee
- **Documentation**
 - Elizabeth Nelson
 - Steve Hanson
- **Test**
 - Trey Chaddick
 - Elizabeth Van Nostrand
- **Management & Operations**
 - Britt Piehler
 - Kristin Fitzimmons
 - Peter Hussey
 - Ren Lis
 - Ben Hackett

Mark Igra
marki@labkey.com



LabKey Software™
PARTNERS IN SCIENCE

If you use LabKey Server for your research, please reference one of these publications about the platform:

General Use: Nelson EK, Piehler B, Eckels J, Rauch A, Bellew M, Hussey P, Ramsay S, Nathe C, Lum K, Krouse K, Stearns D, Connolly B, Skillman T, Igra M. [LabKey Server: An open source platform for scientific data integration, analysis and collaboration](#). BMC Bioinformatics 2011 Mar 9; 12(1): 71.

Proteomics: Rauch A, Bellew M, Eng J, Fitzgibbon M, Holzman T, Hussey P, Igra M, Maclean B, Lin CW, Detter A, Fang R, Faca V, Gafken P, Zhang H, Whitaker J, States D, Hanash S, Paulovich A, McIntosh MW: [Computational Proteomics Analysis System \(CPAS\): An Extensible, Open-Source Analytic System for Evaluating and Publishing Proteomic Data and High Throughput Biological Experiments](#). Journal of Proteome Research 2006, 5:112-121.

Flow Cytometry: Shulman N, Bellew M, Snelling G, Carter D, Huang Y, Li H, Self SG, McElrath MJ, De Rosa SC: [Development of an automated analysis system for data from flow cytometric intracellular cytokine staining assays from clinical vaccine trials](#). Cytometry 2008, 73A:847-856.



LabKey User Conference 2011
NOVEMBER 14-15, SEATTLE