BUSINESS TECHNOLOGY



Pfizer Informatics for NGS

. . .

Giles Day & Adam Kraut Pfizer BioTherapeutics R&D – BioTeam, Inc. XGEN Mar 15th 2010





- Small biotech purchased by Pfizer 2006
- Located in South San Francisco
- Antibody Therapeutics
- 454 Sequencing Team co-located
- Approx 100 people





- Capture of knowledge & decisions
 - Why did we do that, who did that, when did they do it
- Integration of structured & un-structured information
- Discovery of relevant information
- Foster collaboration



Overview of NGS Strategy



- New target discovery
 - GWAS
 - In-house & external
- Support of existing therapeutic programs
 - Phage display library analysis
 - Patient stratification



Solutions Come and Go



- In-house relational platforms
- SharePoint
- Semantic Solutions
- Google/Search Devices
- Pathway Platforms



What They Are



- Inflexible
 - Unable to adapt to changing research processes & organisations
 - Can never capture everything
 - Difficult to administer
- Culture Defiant
 - I'm not ready to share
 - What benefit do I get
 - It's not my solution
- Traditional
 - Software Development cycles are too long even agile ones
 - Too expensive



What They Need To Be



- Useful return more benefit than cost of construction
- Intuitive low adoption barrier
- Fast rapid prototyping
- Flexible modular, able to retool to new uses
- Talkative communicates easily with other technologies
- Cheap design cycle not limited by cost of development
- Open open source, transparent & easily modified
- Adopted Community of developers expanding functionality

Informatics must co-evolve with scientific innovation



How to Translate Bases to Clinical Outcomes?



- GWAS Approach
 - 100s to 1000s samples
 - Well characterized phenotypes
 - Sequence genes/genomes in extreme
 - phenotypes
 - Sequence deep to discover variants
 - Rare loss / gain of function variants
 - Small (additive) effects of many variants expected
 - Annotation
 - Validation of putative SNPs on larger cohorts
 - Identify drug targets
 - Patient stratification, Dx



BUSINESS TECHNOLOGY

www.pnas.org/cgi/doi/10.1073/pnas.0812824106

Overview of NGS Workflow





Informatics Solutions





Informatics Infrastructure





Future Directions







The Details



WikiLIMS

- A system that is quickly constructed and easily revised
- A system that can handle many data types
- A modular and transparent system
- A Web-based system



Simplified tracking information



← Older edit					
Line 5:	Line 5:				
Library type=pool	Library type=pool				
ILab member assigned to prepare library=Bioteam	ILab member assigned to prepare library=Bioteam				
	+ ISample receipt confirmed=Yes				
	+ ISample QA confirmed=No				
»	}}				

- Email and RSS notifications for every step in workflow
- Wiki Revision Control explains *who* did *what* and *when*
- Lab Managers can revert and undo tasks



Monitoring Sequencing Runs– Calendar view



fizer

Intelligent Annotation: Leveraging Web Services





Combining raw data and reporting



stry: FLX _2009_10										
_2009_10										
	_27_02_28_	10_SSF-FLX-cluster_	fullProce	essingAmp	olicons					
iew Files									Name	Last modified Size Description
Lane 😱	Project		Copies	Bead .	Total	Total	Total	Con		
Number	Phase 🌩 Name	Library Name 👙	per 🌩 Bead	Count	Raw 🌩 Wells	Keypass Wells	Keypass %	Aver Len	1 ATG 454Reads fna	- 27-Oct-2009-01:24_3.0M
									1 ATG 454Reads qual	27-Oct-2009 01:24 7 7M
1	UCSF01.01	UCSF01.01_SFLP011	0.5	750,000	444,624	437,127	98.3	22	1. TCA.454Reads.fna	27-Oct-2009 01:24 74M
2	UCSF01.01	UCSF01.01_SFLP012	0.5	750,000	482,013	475,751	98.7	22	1.TCA.454Reads.qual	27-Oct-2009 01:24 190M
Total					926,637	912,878	98.5	22	2.ATG.454Reads.fna	27-Oct-2009 01:24 2.3M
									2.ATG.454Reads.qual	27-Oct-2009 01:24 6.0M
									2.TCA.454Reads.fna	27-Oct-2009 01:24 72M
									2.TCA.454Reads.qual	27-Oct-2009 01:24 185M
									454BaseCallerMetrics.csv	27-Oct-2009 01:24 6.4K
									454BaseCallerMetrics.txt	27-Oct-2009 01:24 20K
									454DataProcessingDir.xml	27-Oct-2009 01:24 351
									2 454QualityFilterMetrics.csv	27-Oct-2009 01:24 954
									454QualityFilterMetrics.txt	27-Oct-2009 01:24 2.1K
ory: 454Ru	ns								454RuntimeMetricsAll.csv	27-Oct-2009 01:24 8.3K
									154RuntimeMetricsAll.txt	27-Oct-2009 01:24 17K
									DM_2009_10_27_14_37_03-F4PLOWH01/	27-Oct-2009 12:57 -
									DM_2009_10_27_14_39_14-F4PLOWH02/	27-Oct-2009 12:57 -
									P 2009 10 27 12 39 37 mapping target set 1-F4P	LOWH01/ 27-Oct-2009 09:59 -
									P_2009_10_27_12_49_07_mapping_target_set_1-F4P	LOWH02/ 27-Oct-2009 11:05 -
									SigProcRunBackupComplete	27-Oct-2009 01:25 0
									dataRunParams.xml	27-Oct-2009 01:24 6.0K
									gsRunProcessor.log	27-Oct-2009 01:24 9.2K

BUSINESS TECHNOLOGY

....

Simple Integrations



page discussion view source history RIN06 02 001 Date: 2010-01-29 Chemistry: Titanium D_2010_01_30_01_35_18_SSF-FLX-cluster_fullProcessing ᠿ View Files Data retreived from Controls with database via 1 line Copies Total Control Control Control Project Bead. Total Raw Total 98% Control Lane 🔺 Library Phase per 🖕 Keypass Average Keypass Pass Number Name Count Wells Keypass % accuracy Mixed % Name Wells Filter % Bead Length Wells of code over 200 bp % 995,379 443.6 0.8 1 983,282 98.8 8,427 93.0 94.8 2 440.7 1,033,895 1,023,075 99.0 8,417 90.5 91.8 0.8 2,029,274 page discussion view source history Total View source for RIN06 02 001 You do not have permission to edit this page, for the following reason: The action you have requested is limited to users in the group: Users. You can view and copy the source of this page: {{454Run|id=33|run dt=0000-00-00} ٠ == Primer Design == [[Image:primer.design.titanium.png]] == 6Mer MIDs == <PRE> +-----| name | sequence | -----+ | PFMID1 | ACGTGT - I | PFMID2 | ACTCTC | ** PFMID3 | AGATCG | | PFMID4 | AGTGAC | ** ----

Getting Started





Project Overview



•

	page	discussion	edit history	move	watch	A Gie	зоау тутак	my preferences	my water	hist my con	unbullons		
all the	Targets												
	Home	Targets	Phage Display	Hybridoma	PEAPS	Vectors F	Protocols	nformatics	Antibo	ody Atlas	Users		
vigation	Target	Name	Uniprot	Structure	Antibodies	Functional	Humanize	d Epitope K	(nown	In-vivo e	fficacy		
Main Page	944		🕼 HUMAN 🖉	yes	yes	yes	yes	yes	Ves		yes		
Community portal Current events	5he	•	HUMAN 🗗	yes	yes	yes	yes	no		ye	s		
Recent changes	-	P47	HUMAN 🖉	yes	yes	yes	no	som	е	ye	s		
Random page Help	107		HUMAN 🗗	yes	no	no	no	no		n	0		
rch	TRUE	100	HUMAN @	yes	yes	yes	yes	yes	•	ye	s 🗟		
	-	Nå 🗌	HUMAN 🗗	yes	yes	?	?	?		2)		
Go Search	Chanter .		a strategy and a	-	yes	?	no	no		n)		
		8	HUMAN @	yes	yes	?	?	?		?)		
Vhat links here	-	r	HUMAN @	yes	yes	?	?	?		?)		
Related changes		k –	HUMAN P	yes	yes	?	?	?		2)		
Active p	/ roject	S				•		1					
						Decisio	ns and	d curre	nt st	atus			

Super Simple Editing



R	Rich Editor							
	page discussion Editing Targe	edit history	move w	atch	2	, Gilesday my talk	my preferences my wa	tchlist my contributions
	Disable rich editor [Ope	n Rich editor in new (1882) ● ● ○ ○ ● ■ <i>I</i> ⊻ ▲	v window] Ma & i I I III Ø ≈∈ ×₂ ײ I IΞ	<> □ 🔐 H	= [[1] (\$) Σ \$ \$ []] [2]	s (c)		
avigation	Target Name	Uniprot	Structure	Antibodies	Functional	Humanized	Epitope Known	In-vivo efficacy
Main Page Community portal			yes	yes	yes	yes	yes	yes
Current events			yes	yes	yes	yes	no	yes
Recent changes			yes	yes	yes	no no	some	yes
Random page			yes	no	no	no	no	no
Help			ves	ves	ves	ves	ves	ves
arch			ves	ves	?	?	?	?
			-	ves	?	no	no	no
Go Search			ves	ves	?	?	?	?
lbox			ves	ves	?	?	?	2
What links here			ves	ves	?	?	?	2
Related changes			,00	ves	?	no	no	2
	1 14		,	Regula	≠ ar HTM	l table		•

Individual Target Page





"Small things, loosely joined, written fast"





- Bring together diverse systems with loose interfaces
- Design and leverage **APIs** wherever possible
- Avoid all-in-one solutions
- Scale components independently



WikiLIMS@Pfizer



.





Data Analysis



Infrastructure as a Service (laaS) for NGS Data





Balancing Trade-offs...



- Physical computing
- Virtual computing
- Availability
- Performance
- Investment
- Cost savings





Amazon S3 Bulk Data Ingest/Export



- How do I forklift 200TB of data into the cloud?
 - Fast network connectivity is expensive
 - Ship physical disks to Amazon to load into S3
- What this means:
 - S3 storage is cheap, deep, and geographically replicated
 - Primary analysis data moved into remote utility storage
 - Data would rarely come back
 - Need to reprocess or analyze?
 - Spin up "cloud" servers and re-analyze in-situ





- Bootstrapping new resources can happen in **minutes** instead of **months**
- Shift capital expenditure to operating expenditure (CapEx/OpEx)
- Primary analysis done on-site or commercial sequencing facility; data moved to remote utility storage
- Easier to distribute and collaborate across multiple sites
- Data would rarely need to move back
- Re-analyze data 'in the cloud' on virtual clusters
- Drop common analysis pipelines into a shared service tier



CloudSobriety



- Security concerns should not be overlooked
 - Regulatory compliance
- Amazon may be the only provider of bulk import/export services
 - What if Amazon goes out of business?
- Life Science tends to be I/O bound; I/O performance in the cloud is still lacking
- IT staff needs to understand 'the cloud' and learn to use it
- Critical to quantify your own internal operating costs



Growth Strategy



- Organic mirroring Pfizerpedia
- Providing information carrots
- Restricted to Rinat groups
- Open source?



Advantages



- Flexibility
- Cost
- Complexity
- Adoption





- Jacob Glanville
- Adam Kraut (BioTeam)
- Steve Pitts
- Andrea Rossi





- That's it
- Giles Day giles.mr.day@pfizer.com
- Adam Kraut <u>kraut@biotam.net</u>

