Data life cycle management in a large scale project

Open Science and Reproducibility Series Workshop III - Data management and Open Data, 22.05.2017

Mark Ibberson and Robin Liechti

Vital-IT, Department of Systems Biology and Medicine, SIB Swiss Institute of Bioinformatics, Lausanne



Vital-IT & SwissProt

Strategic Goal

To provide key competencies & research support to the national life science community



Integrated Coordinated Resource



www.vital-it.ch

Contents

- Introduction why is good data life cycle management important
- An example from a recent project
- Conclusion and take-home messages

Past, Present & Future

The Past

GENOMICS 28, 131-139 (1995)

Genomic Organization of the Human T-Cell Receptor Variable α (TCRAV) Gene Cluster

MARK R. IBBERSON, JOHN P. COPIER,¹ AND ALEX K. SO²

Rheumatology Unit, Hammersmith Hospital, Royal Postgraduate Medical School, Du Cane Road, London W12 ONN, England

Received November 29, 1994; accepted May 11, 1995

We are moving towards a brighter future



- Knowledge in pdfs
- Knowledge resources
- Resource interoperability

The Present

0⁰⁰1¹ 0 0¹-10⁰



0 0 1 d19 9 1 1

Focus is still on publishing

00



0 0 1 19 9 1 1

Scientists are also journalists ..

DATA-DRIVEN JOURNALISM = PROCESS **STORY** VISUALIZE FILTER rising value to public DATA Mirko Lorenz, 2010 Attribution 2.0 Generic (CC BY 2.0)

To publish we need a story ..

Some of what we do gets captured

0 0 1 19 9 1 1

To maximize the value of research we need to be able to interact with data at all levels

Current issues

- Often only raw data is available
 - Loss of time for reanalysis
 - Expertise is often not available
- Methods often unclear or incomplete
 - Lack of reproducibility
- Knowledge is "lost"

How can we do this better?

FAIR

Mark D. Wilkinson et al. The FAIR Guiding Principles for scientific data management and stewardship. Sci Data. 2016 Mar 15;3:160018.

Data is the new oil and we are refining it

An example

IMIDIA is an Innovative Medicines Initiative (IMI) joint public/private partnership for type 2 diabetes research

Integration and Interoperability at the heart of IMIDIA

Working with FAIR principles

Working with FAIR principles

Specific templates used for each data type

SiB			Deimid	ia			Logged in as Mark Ibberson		
Swiss Institute of Bioinformatics			Work Package	e 2A change					
	News	Upload	Experiments	Mice	Analysis	Users			
	Experiment	t: test 🕨					×		
	Select a ty	pe of data							
	2.								
			blood sam	blood sampling					
			fat samp	oling					
			gene expre	ession					
			insulin secreti	on in vitro	K	Accoss to			
				ion in vivo		ALLESS II			
			insulin tolera	ince test	←	different			
			islets isol	lation		tomplate			
			islets lipid	omics	K	template	.5		
			oral glucose to	erance test					
			organ	15					
			pancreas mo	rphology					
			plasma lipio	domics					
			urine sam	pling					

Templates downloadable from a central database

1. Experim	ent: test 🕨				~
2. oral glud	cose tolerance test				~
B. Provide	the list of filenames to up	load			
This step is	used to associate one or multiple dat	afiles to one or multiple	mice		
	1. Download:	1			
		XLS			
		thi	s Microsoft Excel templa	te	
	2. Open the downle	baded file in Excel.			
	The file has two	sheets. Each row of	the first sheet contains the	he following columns:	
	1. Names of a	d SELF if the oral of	upioad ucose tolerance test mea	surements are contained	in the second sheet of
	this excel fil	e.		· · · · · · · · · · · · · · · · · · ·	
	2. Mouse ID o	the word ANNOT if	the uploaded file is not li	nked to a mouse.	
	3. Day of mea	surement. One of D2	2, D10, D30 or D90		
	4. Date of data 5. Comment	Collection (1111-W	IM-DD)		
	(→show sample	file)			
	The second sheet	might contain additio	nal sample, extract or m	easurement annotations	Users are free to modify
	its format following	→ these instruction	s	case. short annotations.	source are not to modify
	3. If one file is not I	inked to a mouse, w	rite ANNOT in the MOUS	E_ID column	
	4. Upload this anno	tated excel file using	g the form below.		
	Data files are up	loaded in the next s	tep		
		Data type oral glu	ucose tolerance test		
	Annotate	d Excel file Choose	File No file chosen		
			g upload		

We worked with the lab biologists to adapt their own templates and standardize the input data and metadata

01⁰¹¹ 0 0¹10⁰

Unique cross-site sample identifiers

- Enabled tracking of samples between different centers
- Not restricted to any particular project
- Short enough to fit on the side of a sample tube!

Curation work is essential but is often underestimated

Egyptian scribe with Papyrus

Biocurator with notebook

Working with FAIR principles

Physiological and molecular data were collected

Similar data were generated in human and mouse

A web portal captured the data

Working with FAIR principles

Using the Vital-IT distributed platform

Vital-IT : A "cloud" HPC

- > 8000 cores
- > 15 PB near-line/archive data
- GPU/CPU based computing

The infrastructure is centrally managed from Lausanne with a distributed support

Working with FAIR principles

Account validation

19 0 1 1

Data ownership rights managed automatically

	SiB Swiss Institute of Bioinformatics	đ		Work Package	ICI 2A change			Logged in as Mark Ibberson		
		News	Upload	xperiments	Mice	Analysis	Users			
			E	Beta Cell Mass An	alysis					
		Details	Owner Group Work Package Request Sownload Create date Experiment type Description	Tracy Gorman AstraZeneca WP2A download upon request 2012-01-19 WP2A: main study						
	ſ	Files		0100000						
		2 results 2HF49_raw_data.xis	Fie p	Data Type Filter : ancreas morphology	Owner Tracy Gorman	Size U	pload date Actions 8.06.2012 0			
		IMIDIA_EXPERIMENT_	Reque	sts an	-		8.05.2017 🛔			
		Ļ	neque			K				
			Gran	ts acc	cess					
		L								
г ч			Ĩ							
Email										
										(i
			Ń							Č.
	$ \square $									
			L							
								Dr	ningt m	- ombor
	Da	ata "o	wner"					PI	oject fr	lemper

Working with FAIR principles

Access to detailed, harmonised sample information

SIB Swiss Institute of				Dei	midia				Logged in as Mark ibberson
Bioinformatics	Work Package 2B change								
	News Upload Experiments Patients Analysis Users								
	Browse patien								
	(
	Filters:								
				+	add filter				
	Number of potion								
	Number of patien	15. 044 - NU	inder of itssues. 1136 - Nun	iber of samples (i	icidaling / excluding se	cuons). 44597 5	SUI	mmary view	
	location	patients		tissues	-	samples		last update	
	Dresden * 26	8 patients	268 Partially pancreatic patien	it 572 tissues	40 islets	2876 samples	7 GTA 2.5% block	08.05.2017	
					264 pancreatic tissue		151 GTA 4% block		
			Patie	nt & :	sample	2	IS) I block block		
						-	block sections		
			meta	data	from s	ever	GTA 0.02% block		
			_		-		ma		
	Q		Europ	bean	sites		sample	05 11 0010	
	Hannover • 12	patients	1 Organ conor	1 100000	Ti pancieatto tiasee	11 Jampios	ample	01.10.2013	
			11 Partially pancreatic patient						
	Eli Lilly 🔻 44	patients	44 Organ donor	43 tissues	43 islets	69 samples	30 DNA sample	20.11.2013	
	Pisa 🔻 20	8 patients	208 Organ donor	501 tissues	141 blood 152 islets	1494 samples	208 DNA sample 77 GTA 2.5% block	04.10.2016	
					208 pancreatic tissue		150 Islets (IS)		
							201 LCM block 591 PFA block		
							141 Plasma		
							126 RNA sample		

🖶 download...

Intuitive web-based tools to explore the data

01011 01-10⁰

Parallel plots to visualise data across multiple conditions

- Gene profiles -

Trait profiles (click to show) -

Filters (click to hide) absolute ratios	Adj.p-value cutoff 0.05	p-value cutoff 0.01
select pathway:	Pathway, GO-CC, GO-BP or	r GO-MF
user gene list: 🛨		
trait filter:select :	selected trait filters	
DBA2J:D30 between 1.9 and 3.2		

Intuitive web-based tools to explore the data

1. Module vs Trait Heatn	Interactive heatmaps for intuitive data mining
- Filters	
Search gene:	show reset my favorites standard: fuzzy modules: ratios exprs
Filter on pathway:	p-value cutoff: 0.05 adj. p-value cutoff: 1
grey60 (121) maroon (67) palevioletred3 (70) salmon2 (51) coral1 (64) darkmagenta (5782) mediumpurple3 (155) coral3 (552) floralwhite (133) white (247) violet (90) brown2 (54) antiquewhite2 (40) blub2 (974) purple (188) lightslateblue (120) plum (129) darkseagreen3 (83) black (69) black (69) black (65) darkolivegreen3 (159) yellow (237) firebrck4 (152) darkulvegreen (90) darkel (116) deeppink (22) blueviolet (25)	

Mouse and Human data integration was possible

Working with FAIR principles

IMIDIA sustainability and interoperability

Goal : Transform IMIDIA data and results into a standardized format

RDF= Resource Description Framework CDISC= Clinical Data Interchange Standards Consortium

What is RDF ?

Susie is_child_of John Smith Web page has_author Robin Ethanol is_oxidized_to acetaldehyde Protein kinase regulates enzyme activity

Other RDF sources are "plugged in" to IMIDIA data

Federated queries can be used to run queries across several SPARQL end-points

Working with FAIR principles

SIB people involved in IMIDIA

Mark Ibberson Data analysis and coordination

Robin Liechti Database, web interface and analysis tools

Lou Gotz Web interface development

Jerven Bolleman (SwissProt) RDF

Leonore Wigger Statistical models

Frédéric Burdet Data management Bioinformatics, RNASeq analysis

Dmitry Kuznetsov Database architect, RDF

Diana Marek CDISC term mapping

Nicolas Guex Algorithm development and data analysis

Anne Niknejad Lipid, gene and pathway annotations

Ioannis Xenarios Work package leader, director of Vital-IT and SwissProt

Roberto Fabbretti Vital-IT infrastructure management

Brian Stevenson

Alan Bridge (SwissProt) Lipid annotations

Lucilla Aimo (SwissProt)

Lipid annotations

Some notes

+ SIB

- FAIRification was only possible because of protocol harmonisation, data centralisation and curation right from the start of the project
- Additional funding was set aside to enable the work
- A lot of people with diverse skills were needed
- The process took time and we shouldn't underestimate costs
 - ~4PM/year for data curation and model design/encoding
 - This is sustained for 7 years and not by a single PhD/Postdoc

Some outputs from the project

- Candidate lipid biomarker for T2D
 - Capture, storage, sharing + **discovery**
- FAIRified data for mouse and human
 - Interoperable with Uniprot, SwissLipids etc.
 - CDISC (clinical standards) compatible
- SIB (Vital-IT) is now Data Coordination Center for 2 new IMI projects RHAPSODY and BEAt-DKD
 - Aligning multiple clinical cohorts to CDISC
 - Federated database infrastructure and analysis

Cell Reports

Plasma Dihydroceramides Are Diabetes Susceptibility Biomarker Candidates in Mice and Humans

Graphical Abstract

Leonore Wigger, Céline Cruciani-Guglielmacci, Anthony Nicolas, ..., Christophe Magnan, Mark Ibberson, Bernard Thorens

Correspondence mark.ibberson@sib.swiss (M.I.), bernard.thorens@unil.ch (B.T.)

In Brief

Authors

Wigger et al. find that several sphingolipids in mouse plasma correlate with glucost tolerance and insulin secretion. Quantitative analysis of these and closely related lipids in human plasma from two cohorts reveal that dihydroceramides are significantly elevated in individuals progressing to diabetes, up to 9 years before disease onset.

Take-home messages

- We are having to deal with more and more diverse data
- Main challenge is making sure knowledge is not lost
 - How to make research data sustainable and reusable
 - Resources need to be put aside for this
 - This can impact future funding allocation (e.g. H2020)
- FAIR is a good concept to work towards
 - Need to think and plan the necessary steps early on (ideally before the project starts)
- FAIR does not necessarily mean Open Access
 - Many datasets (e.g. Clinical) need to remain restricted access
 - These datasets still need to be interoperable with other private and public resources

Thank you!

0 0 1 19 9 1 1

0⁰⁰1¹ 0 0² 10⁰

Backup slides

01⁰11 0 0¹10⁰

WP2A – mouse data

WP2A – essential classes

WP2A – classes and properties

25 classes, 26 properties, 30'488'338 triples

- imidiaWP2A:DataStructures (80985)
 - imidiaWP2A:geneExpressionCount (341)
 - imidiaWP2A:massMeasurement (2398)
 - imidiaWP2A:substanceAmount (78246)
 - imidiaWP2A:Diets (2)
 - imidiaWP2A:Measurements (7088)
 - imidiaWP2A:MeasurementType (13)
 - imidiaWP2A:Mice (1536)
 - imidiaWP2A:mouseAgeValues (7)
 - imidiaWP2A:MouseCohorts (56)
- imidiaWP2A:MouseLines
- imidiaWP2A:Stimuli (744)
- imidiaWP2A:Substances (165)
- imidiaWP2A:Tissues (7)
- imidiaWP2A:Traits (20)

- imidiaWP2A:ageOfTheSubjectedMouse
- imidiaWP2A:calculatedAreaUnderCurve
- imidiaWP2A:calculatedMeanValue
- imidiaWP2A:dateOfDayZero
- imidiaWP2A:delayUnit
- imidiaWP2A:geneticLineage
- imidiaWP2A:internalMouseID
- imidiaWP2A:measuredAfterDelay
- imidiaWP2A:measuredOnDate
- imidiaWP2A:measuredOnTissue
- imidiaWP2A:measuredResults
- imidiaWP2A:measuredSubstance
- imidiaWP2A:measurementTrait

- imidiaWP2A:measurementType
- imidiaWP2A:normalizedCountForGene
- imidiaWP2A:sampleConcentration
- imidiaWP2A:sampleConcentrationUnit
- imidiaWP2A:sampleMass
- imidiaWP2A:sampleMassUnit
- imidiaWP2A:sampleVolume
- imidiaWP2A:sampleVolumeUnit
- imidiaWP2A:subjectedMouseCohort
 - imidiaWP2A:subjectedMouse
- imidiaWP2A:subjectTreatedBeforeExperimentBy
- imidiaWP2A:treatedBy
- imidiaWP2A:usedSubstance

WP2B – patient data

WP2B – essential classes

WP2B – classes and properties

29 classes, 148 properties, 13'150'154 triples

- imidiaWP2B:AffymetrixProbes (54675)
- ImidiaWP2B:DataStructures (25540)
 - imidiaWP2B:AFFYdata (24818)
 - imidiaWP2B:geneExpressionCount
 - imidiaWP2B:massMeasurement
 - imidiaWP2B:substanceAmount (722)
 - imidiaWP2B:Measurements (4251)
 - imidiaWP2B:MeasurementType (22)
 - imidiaWP2B:Patient (494)
 - imidiaWP2B:Samples (3957)
- imidiaWP2B:SampleTypes (3960)
- imidiaWP2B:Specimens (1007)
- imidiaWP2B:Stimuli (1441)
- imidiaWP2B:Substances
 - imidiaWP2B:TissueBlocks (72)
 - imidiaWP2B:Tissues (7)

- imidiaWP2B:affymetrixProbe
- imidiaWP2B:Bioanalyzer_data
- imidiaWP2B:calculatedAreaUnderCurve
- imidiaWP2B:calculatedMeanValue
- imidiaWP2B:CDISCequivalent
- imidiaWP2B:cDNA_yield
- imidiaWP2B:Chip_name
- imidiaWP2B:delayUnit
- imidiaWP2B:FlowCell_number
- imidiaWP2B:Library_concentration
- imidiaWP2B:Library_date
- imidiaWP2B:Library_fragment_size
- imidiaWP2B:Library_local_ID
- imidiaWP2B:Library_type
- imidiaWP2B:Library_volume

