

A FAIR Data Sharing Framework for Large-Scale Human Cancer Proteogenomics

Islam M^{1,2}, Christiansen J³, Mahboob S⁴, Valova V⁴, Baker M⁴, Capes-Davis D⁴, Hains P⁴, Balleine R^{1,4}, Zhong Q^{1,4}, Reddel R^{1,4}, Robinson P^{1,4}, **Tully B⁴**

Brett Tully

28 Nov 2018

@brett_tully

btully@cmri.org.au

Big-Data Approach to Clinical Decision Making

Delivering **molecular data** to cancer **clinicians**,
in a **clinically-relevant time** frame,
to maximise the accuracy of **treatment decisions**

Complex Project; Many Moving Parts



Co-Directors



Roger Reddel

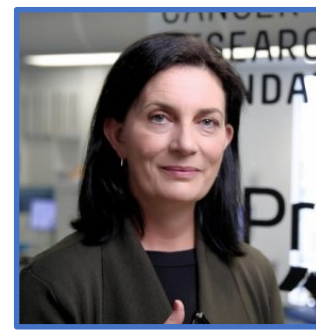


Phil Robinson



Peter Hains

Cancer Proteomics



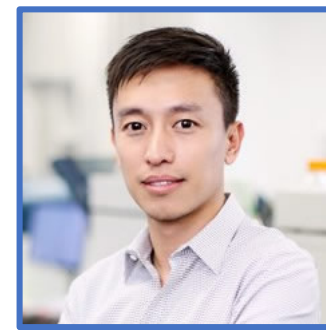
Rosemary Balleine

Cancer Pathology



Brett Tully

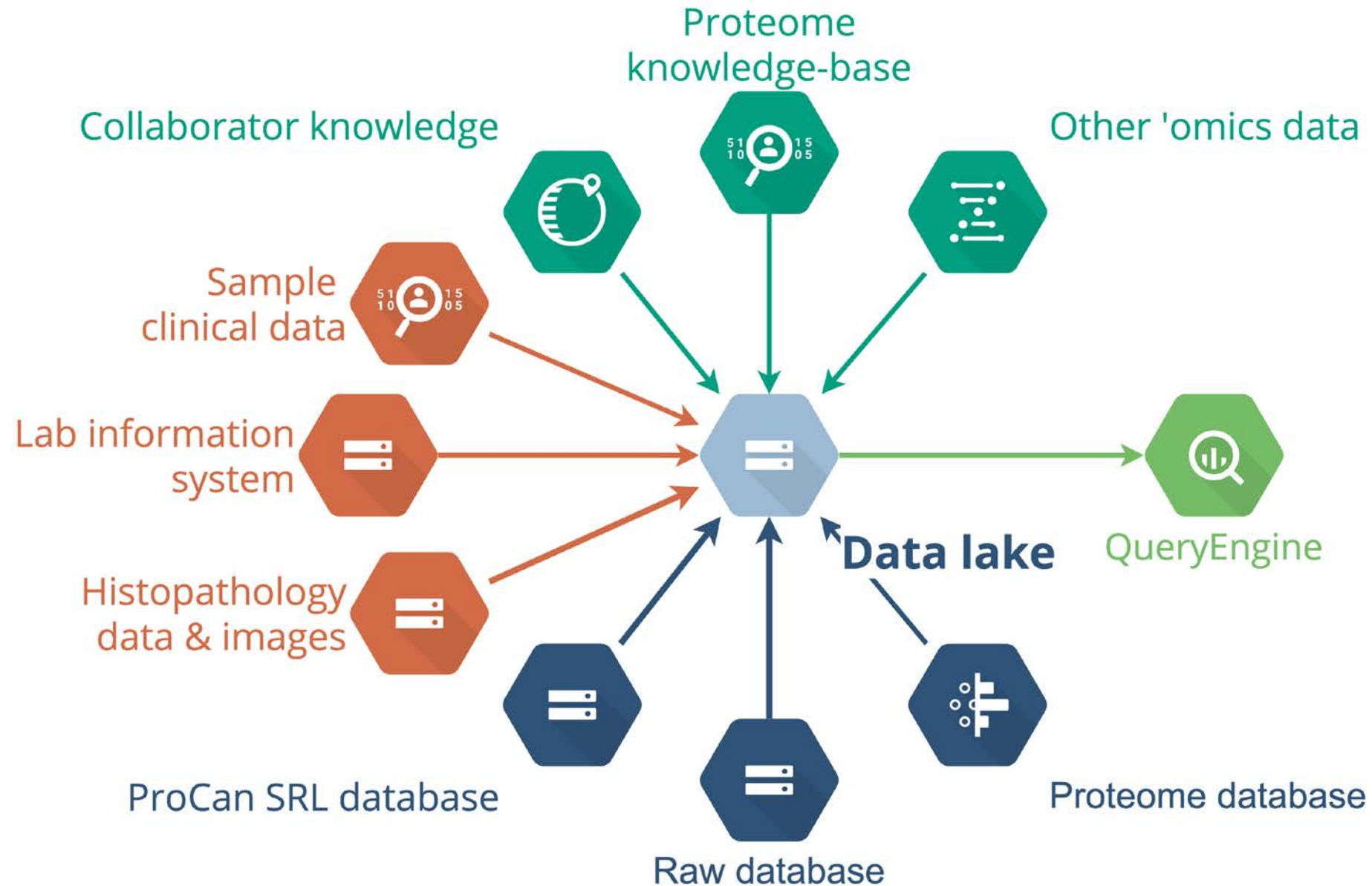
Software Engineering



Qing Zhong

Cancer Data Science

ProCan Data Lake Aggregates Many Sources



FAIR Data Sharing in ProCan

Findable

Accelerates scientific discovery

Accessible

Interoperable

Enhances integrity, transparency,
and reproducibility

Reusable

Findable

by both humans and machines

- Discoverable, well-defined metadata
- Persistent unique URLs, or Document Identifiers (DOIs)
- Machine-readable metadata

ProCan Challenges

- Unique IDs and machine readable metadata are easy to create
- Human discoverability is much more difficult → context dependent

Accessible

using standard protocols

- Data retrievable by their unique identifier
- Open, free, and implementable protocol
- Can be subject to constraints: ethical, privacy, security, commercial

ProCan Challenges

- Integration of many domains: pathology, LIMS, multi-omics, analytics
- Integration of 100's collaborators: each with different agreements
- Sustainable funding: on-going costs for long-term storage & access

Interoperable

with other systems and data resources

- Industry/community standard formats & vocabularies
- Where possible, data accessible in open-formats
- Minimal intervention required to combine with 3rd party data

ProCan Challenges

- Proteomics data largely produced in proprietary vendor formats
- Pan-cancer = cross-discipline vocabularies = complex ontology

Reusable

and reproducible via richly described metadata

- Clear and accessible usage license
- Fully described provenance to community-defined standard
- Completeness meeting community-defined expectations

ProCan Challenges

- Existing public repositories; context-specific, and non-overlapping
- Not all data generated internally; dependent on 3rd party processes

Proposed FAIR Shared Responsibility Framework

Data Custodian (DC)

Sample Provider/Collaborator
ProCan

Data Management
Access Management
Publication

Interoperability & Reusability
Data Quality
(Meta)Data standardisation

Hosting Institutes &
Data Steward (ProCan)

Hosting Institute

Children's Medical Research Institute (ProCan)
Collaborator's Institute
International Repositories

Fit-for-purpose Infrastructure
Authentication
Data Storage
Compute
Retention
Discovery Services
Data Submission
Transfer Protocols

Proposed Risk-based Access, Sharing and Governance Model

Data	Non-Sensitive					Sensitive
Governance Access	Low Open	Low Registered	Moderate Registered	High Controlled	Highest Controlled	User Least Trusted Most Trusted

Proposed Risk-based Access, Sharing and Governance Model

Data	Non-Sensitive					Sensitive
Governance	Low	Low	Moderate	High	Highest	User
Access	Open	Registered	Registered	Controlled	Controlled	Least Trusted
			QA & QC Data			Most Trusted
			Published research output			
			Analysed data + minimum metadata			

Proposed Risk-based Access, Sharing and Governance Model

Data	Non-Sensitive					Sensitive
Governance Access	Low Open	Low Registered	Moderate Registered	High Controlled	Highest Controlled	User
			<p>Lightly aggregated data De-identified analysed data Non-identifiable clinical data</p>			Least Trusted
			<p>Moderately aggregated data De-identified analysed data Non-identifiable clinical data</p>			
			<p>Highly aggregated data De-identified analysed data Non-identifiable clinical & demographic data</p>			Most Trusted

Proposed Risk-based Access, Sharing and Governance Model

Data	Non-Sensitive					Sensitive
Governance Access	Low Open	Low Registered	Moderate Registered	High Controlled	Highest Controlled	User
			<p>Lightly aggregated data De-identified analysed data Re-identifiable clinical data for this dataset only</p>			Least Trusted
			<p>Moderately aggregated data De-identified analysed data Re-identifiable minimum clinical data for analysis</p>			
			<p>Sample verification and quality information De-identified tissue sample and related metadata Re-identifiable minimum clinical & demographic data for analysis</p>			Most Trusted

Proposed Risk-based Access, Sharing and Governance Model

Data	Non-Sensitive					Sensitive	User
Governance	Low	Low	Moderate	High	Highest		
Access	Open	Registered	Registered	Controlled	Controlled		
			<p>Lightly aggregated data</p> <p>Genetics & multi-omics data</p> <p>Re-identifiable clinical data</p>				Most Trusted
			<p>Moderately aggregated data</p> <p>Genomics & multi-omics data</p> <p>Re-identifiable clinical data</p>				
			<p>Highly aggregated genomics & multi-omics data</p> <p>Commercial and pharmaceutical data</p> <p>Re-identifiable clinical & demographic data</p>				

Proposed Risk-based Access, Sharing and Governance Model

Data	Non-Sensitive					Sensitive
Governance Access	Low Open	Low Registered	Moderate Registered	High Controlled	Highest Controlled	User
	<p>Clinical decision support system – non-personal aggregated data Re-identifiable aggregated clinical data</p>					Least Trusted
	<p>Clinical decision support system Re-identifiable clinical & demographic data</p>					
	<p>Clinical decision support system Clinical & demographic data</p>					Most Trusted

Acknowledgements

Co-Directors: Prof. Roger Reddel & Prof. Phil Robinson

Pathology: Prof. Rosemary Balleine & team

Proteomics: Dr Peter Hains & team

Software Engineering: Dr Brett Tully & team

Data Science: Dr Qing Zhong & team

***Intersect:** Dr Mohammad Islam, Dr Jeff Christiansen



THE UNIVERSITY OF SYDNEY