

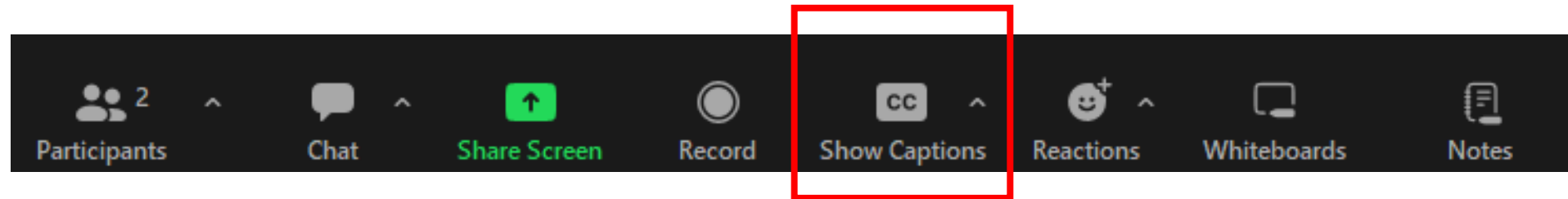
Data Deposit for Health Sciences Research

January 21, 2025



Turning 'live captions' on and off

- On your meeting controls, click on **“Show Captions”**



Land Acknowledgement

We wish to acknowledge this land on which the University of Toronto operates.

For thousands of years it has been the traditional land of the Huron-Wendat, the Seneca, and the Mississaugas of the Credit.

Today, this meeting place is still the home to many Indigenous people from across Turtle Island and we are grateful to have the opportunity to work on this land.

Housekeeping

- This webinar is being recorded and transcribed
- A link to the recording and presenter slides will be sent to all participants after the session
- Please put questions into the chat, we will hold all questions until the end of the presentations

Purpose & Agenda

Bring together the University of Toronto tri-campus and TAHSN health sciences research community for facilitated conversations about research data management.

Learning Objectives:

- Preparing data for deposit
- Considerations for selecting a repository and managing access
- Finding and using secondary data sets

1. **Data deposit basics**
2. **Panelist presentations**
3. **Panel discussion and Q&A**

Data Deposit – The Basics

What

Transferring research data to a data repository for secure storage and preservation

Researchers control data accessibility

Includes:

- Collected data
- Accompanying documentation
- Source code
- Software
- Metadata
- Other supplementary materials

Why

Purpose: Ensure that data are securely preserved and accessible to researchers after project completion

Can support reuse, validation, replication and links with other research

Requirements & considerations

- Funders (e.g., [Tri-Agency Research Data Management Policy](#))
- Publication
- Disciplinary norms
- Ethical, cultural, legal, and commercial requirements

How

Factors

- [FAIR Principles](#)
- Indigenous data sovereignty
- Data sensitivity & confidentiality

Types of repositories

- Disciplinary
- Multidisciplinary/generalist
- Institutional

Resources:

- [University of Toronto Libraries – Data Repositories](#)

Panelists



Dr. Rachel Harding

Assistant Professor, Department of
Pharmacology and Toxicology &
Principal Investigator, Structural
Genomics Consortium



Conrad Pow

Senior Lead for Digital Health at
Diabetes Action Canada



Dr. Daniel Roth

Associate Professor, Department of
Paediatrics & Clinician-Scientist,
Division of Paediatric Medicine at
SickKids



Dr. Michael M. Hoffman

Associate Professor, Department of
Medical Biophysics & Senior
Scientist, Princess Margaret Cancer
Centre

The SGC is a global public-private partnership focused on open drug discovery





4,000+
DEPOSITED
STRUCTURES



4,450+
PLASMIDS
DISTRIBUTED

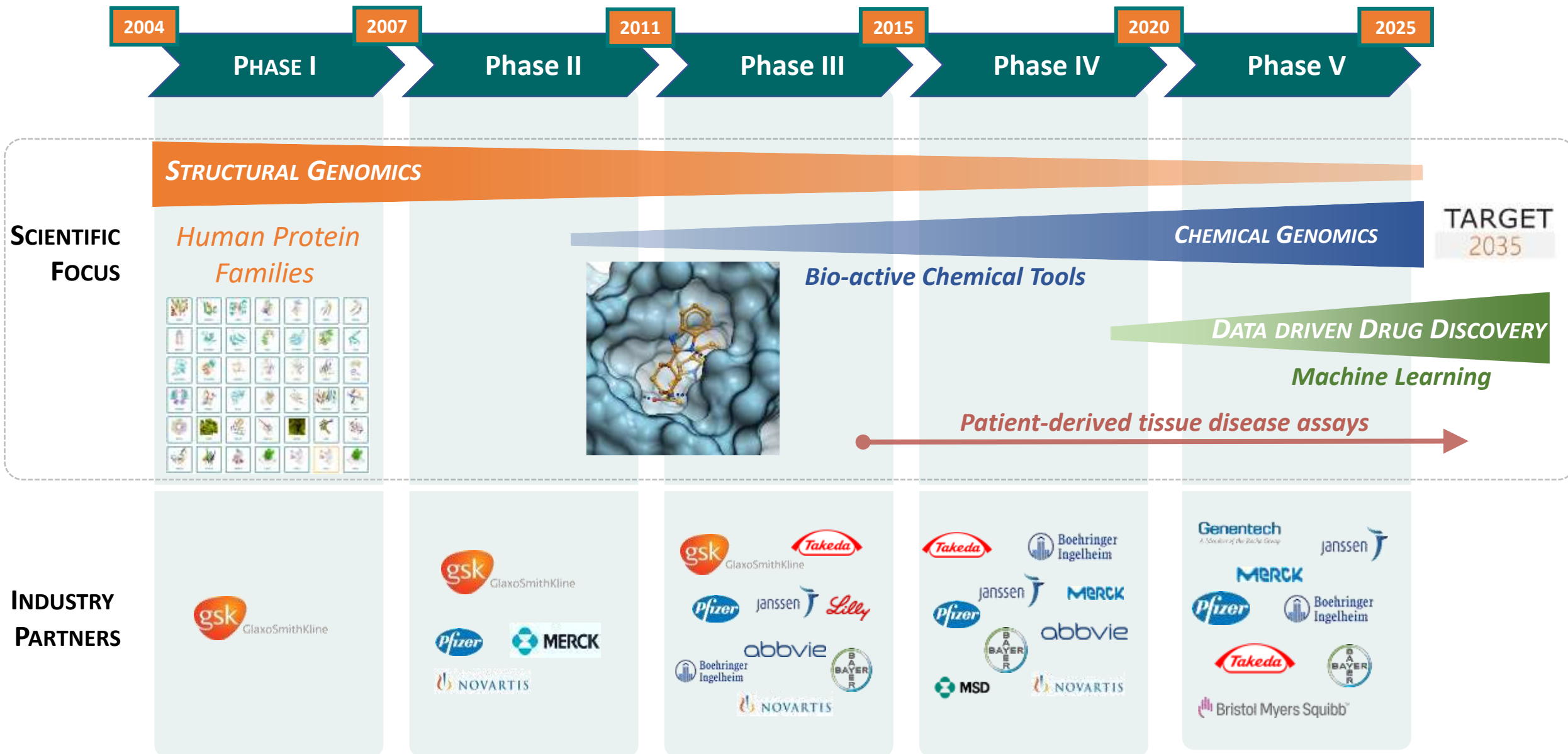


1,500+
DETAILED PURIFICATION
PROTOCOLS



1,100+
PEER REVIEWED
PUBLICATIONS OF
STRUCTURES

Evolving science and partners to address pressing global needs



Chemical probe = a drug-like small molecule that selectively modulates the activity of a specific protein in cells

Probe Criteria

- Potent:** $KD/IC_{50} < 100 \text{ nM}$ *in vitro*
- Selective:** $> 30x$ over related proteins
- Cell activity:** at $1 \mu\text{M}$ or less
- Low/No off-target cellular toxicity**
- Ideally, a negative control molecule**

HIGHLY ENABLING
RESEARCH TOOLS

Bibliometric analysis shows that chemical probes were the most impactful reagent/tool to enable researchers to work on new targets



- Used to interrogate protein function in cells
- Chemical counterpart to genetic knock-out/knockdown methods
- Enables development of novel disease target hypotheses
- The first step in a drug discovery program
- All shared in public domain

www.theSGC.org/chemical-probes/

SGC CHEMICAL PROBES BY THE NUMBERS



DISCOVERED

200+

Novel chemical probes developed in collaboration with industry and academic partners



DISTRIBUTED

50,000+

Samples of chemical probes distributed globally by SGC and trusted vendors



CITATIONS

13,000+

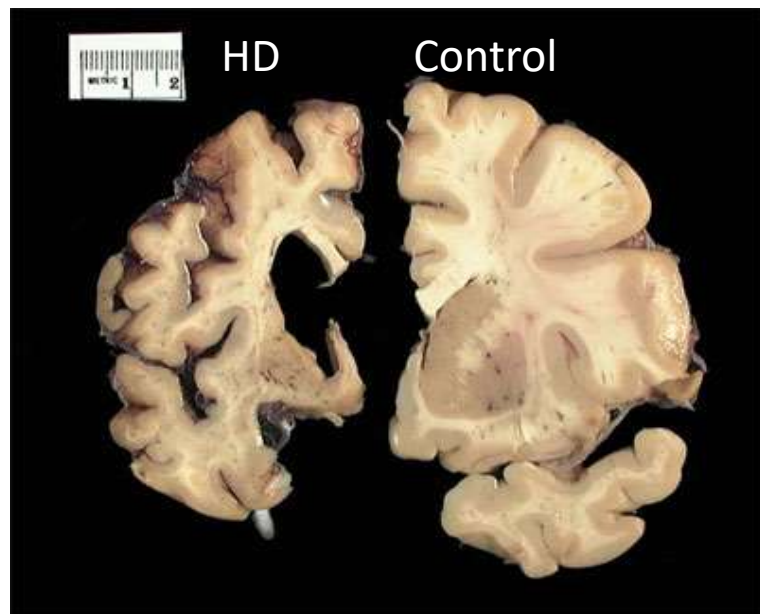
SGC chemical probes used by scientists around the world



CLINICAL TRIALS

85+

Clinical trials and late-stage preclinical programs based on therapeutic hypotheses generated with SGC chemical probes

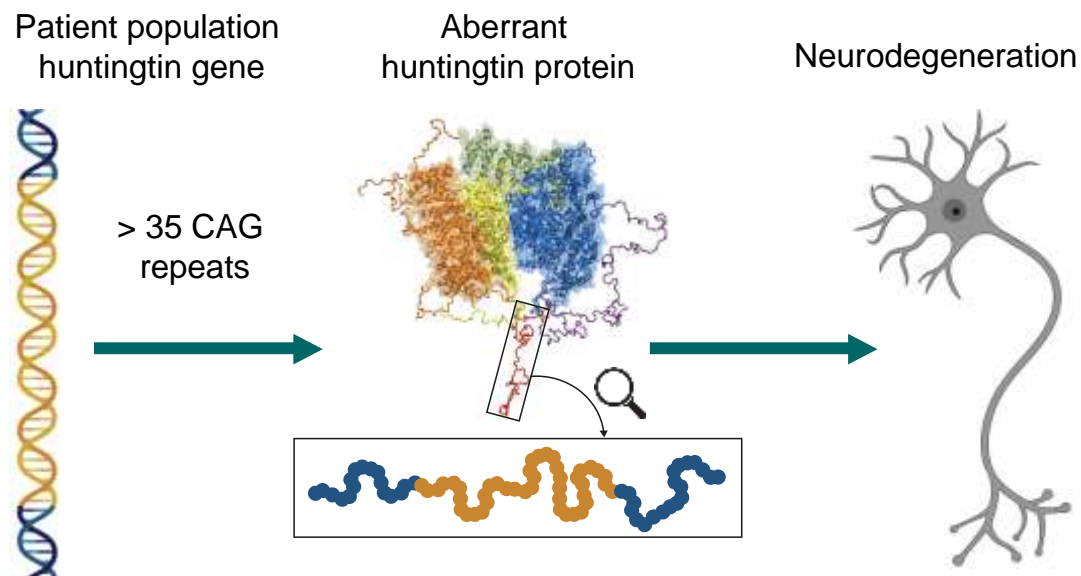
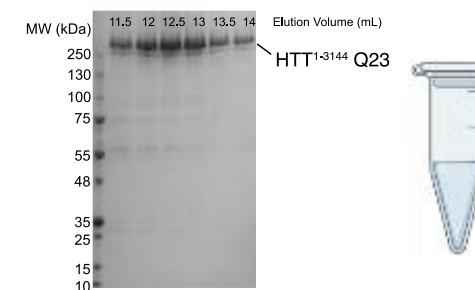


Outstanding Questions:

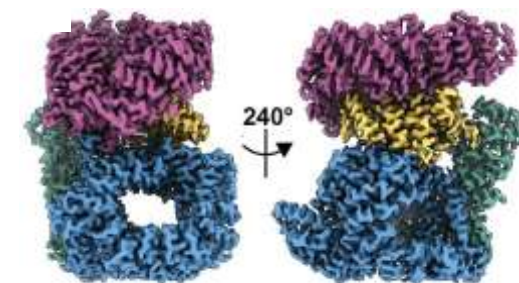
Why is there a mutation threshold for disease?
 How does this (relatively) small change to the huntingtin protein initiate neurodegeneration?

New Approaches to Seek Answers:

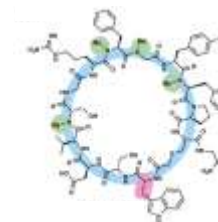
Robust protein biochemistry
Samples shared freely with community



Structure-function analysis
Data deposited prior to publication



Novel modalities for chemical biology
Develop unencumbered tools





Diabetes Research Connect

Conrad Pow
Sr. Lead, Digital Health
Diabetes Action Canada

Diabetes Action Canada

- Pan-Canadian Diabetes Research Network with **120 researchers**, and over **400 Patient Partners**
- Funded by the **Canadian Institutes for Health Research (CIHR)**, funding partners (i.e. Diabetes Canada) and through **philanthropic support**
- Strengthening healthcare systems through multi-disciplinary collaboration among experts in primary and specialist care, **health data management**, analytics and policy
- Advancing Patient-oriented research by fostering collaborations between **people living with diabetes** and research teams dedicated to preventing complications
- Primary focus on addressing the needs of **equity denied communities**

Diabetes Research Connect



- **Secure Virtual Platform**
- **Largest National Primary Care Data Repository**
- **Over 1,200 primary care providers**
- **Over 1.8 million patient data**
- **Approval from DAC's Digital Health Governing Committee**
- **All projects co-designed with Patient Partners**

Data Collection

Data is collected through an **REB approved, systematic process** involving the extraction of **de-identified** patient information from electronic medical records (EMRs) used by primary care providers across Canada.

Data Collection Process:

1. Data Extraction
2. Data De-identification
3. Data Standardization
4. Data Storage
5. Data Quality Assurance

Data Sources



- De-identified National EMR data
- Agreements from 8 Provinces and 1 Territory

EMR Vendors (and products):

Telus: (Kinlogix, Med Access, Medesync, Nightingale, Practice Solutions, Telin, Wolf)

QHR: (Accuro, Jonoke, Healthscreen, xwave)

OSCAR

Da Vinci

Healthquest

InputHealth

IntraHealth

Purkinje

P&P

- Inclusion of Ontario Community Health Centres (BIRT)

Data

Data Extraction:

- CPCSSN does **not** extract all patient data
- Focuses mainly on **structured data**
- Does **not** extract notes or PDFs

Data De-identification:

- No identifiable information is extracted, and only **de-identified data** is processed

Extracted Data Fields:

- Includes billing, health conditions, reasons for visits, lab results, family history, medications ICD-9, referrals, vital stats...

Extraction Schedule:

- Data extractions occur biannually (June 30 and December 31)

Research Environment

Researchers are provided with a suite of analytical and development tools to conduct their studies. Software provided includes:

- SAS
- R/R Studio
- Stata
- Anaconda
- Jupyter Notebook
- Sublime Text



Researchers may request additional software applications be installed in their environments

*Researchers must provide their own licenses for software not already present in the SRE.

** Software that utilizes concurrent licensing (i.e. SPSS) cannot be used due to disabled internet access.

Role of DRC in Health Research



Impact of government-funded insulin pump programs on insulin pump use in Canada: a cross-sectional study design using the National Diabetes Repository. Weisman, A. et al.

Increased Accessibility: Government funding has made insulin pumps more accessible, confirming the positive impact of financial support on technology uptake.

Persistent Disparities: Despite funding, disparities persist, with lower-income individuals less likely to use insulin pumps, indicating that financial support alone is insufficient.

Non-Financial Barriers: Other factors, including healthcare provider biases and limited access to specialized care, contribute to the underutilization of insulin pumps among certain populations.

Need for Comprehensive Strategies: To achieve equitable access, it's crucial to address both financial and non-financial barriers through comprehensive approaches that include education, support, and systemic changes.

Role of DRC in Health Research



Achievement of treatment targets among patients with type 2 diabetes in 2015 and 2020 in Canadian primary care. Lau, D. et al.

Medication Usage: The study found suboptimal use of statins and ACE inhibitors or ARBs, especially among women, suggesting potential gaps in adherence to clinical guidelines or prescribing practices.

Trends Over Time: Between 2015 and 2020, there was an increase in HbA1c target achievement, stability in LDL-C target rates, but a decline in blood pressure target achievement and the use of recommended medications, indicating areas needing quality improvement.

Implications for Practice: The findings underscore the necessity for targeted interventions to enhance the comprehensive management of type 2 diabetes, with particular attention to gender disparities and the declining trends in blood pressure control and medication usage.

This study emphasizes the critical role of primary care in managing type 2 diabetes and the ongoing need for quality improvement initiatives to ensure patients achieve comprehensive treatment targets.

Thank You!

Questions?



Image: DALL-E



Preparing and depositing health research data for public access: SEPSiS project experience

Daniel Roth

January 21, 2025



icddr,b

SickKids®

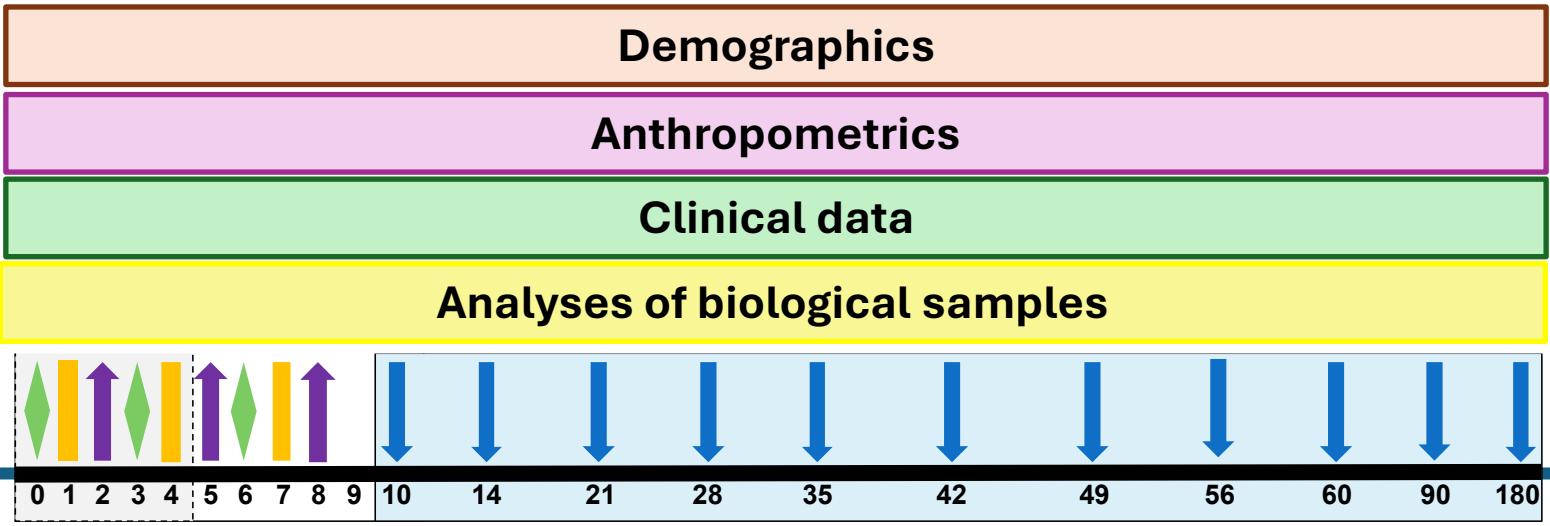
SEPSiS Project in Dhaka, Bangladesh



November 2020



November 2022



0 – 4 days of age

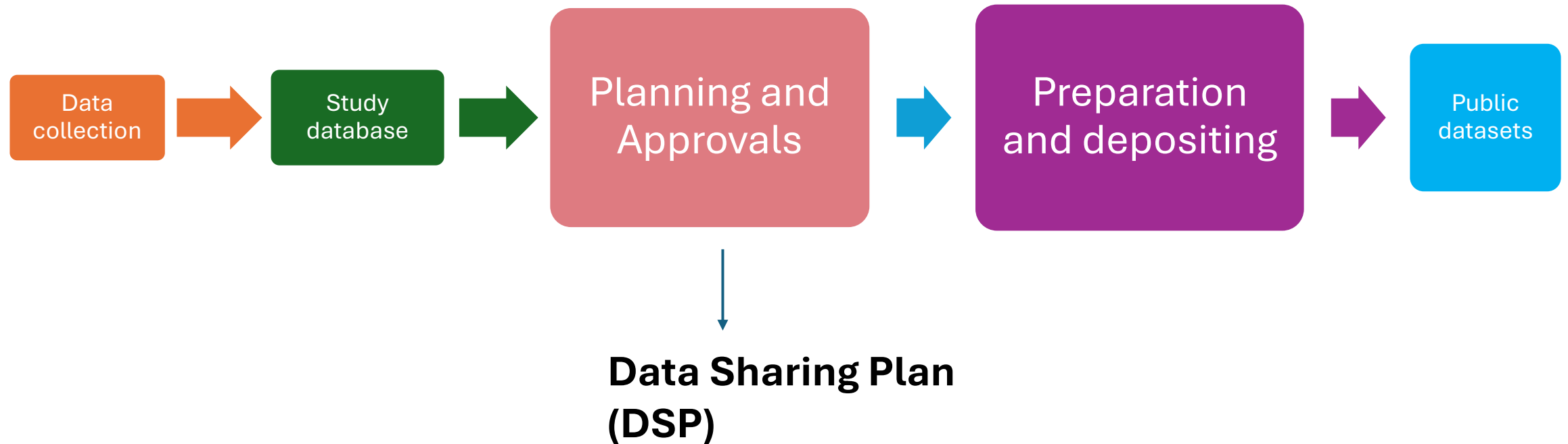
6 months of age

N = 2,458 infants

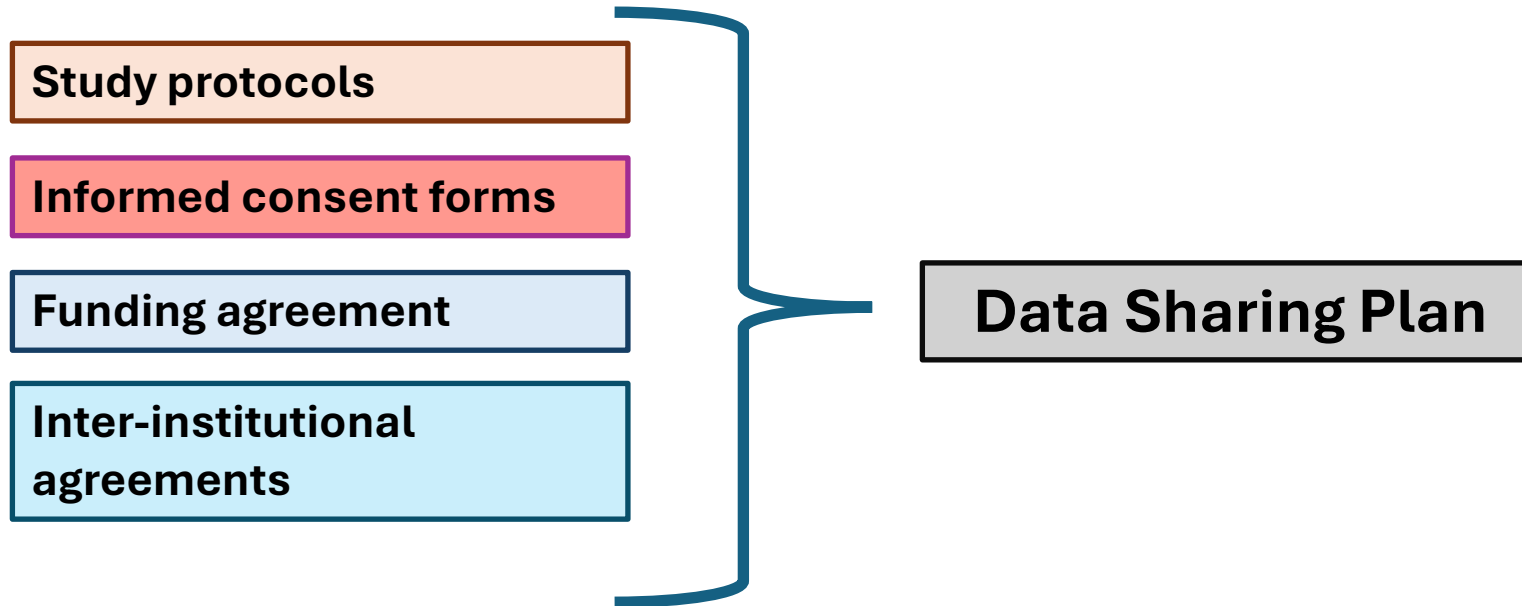
Research Data Sharing – initial assumption



Research Data Sharing – reality



Data Sharing Plan (DSP)



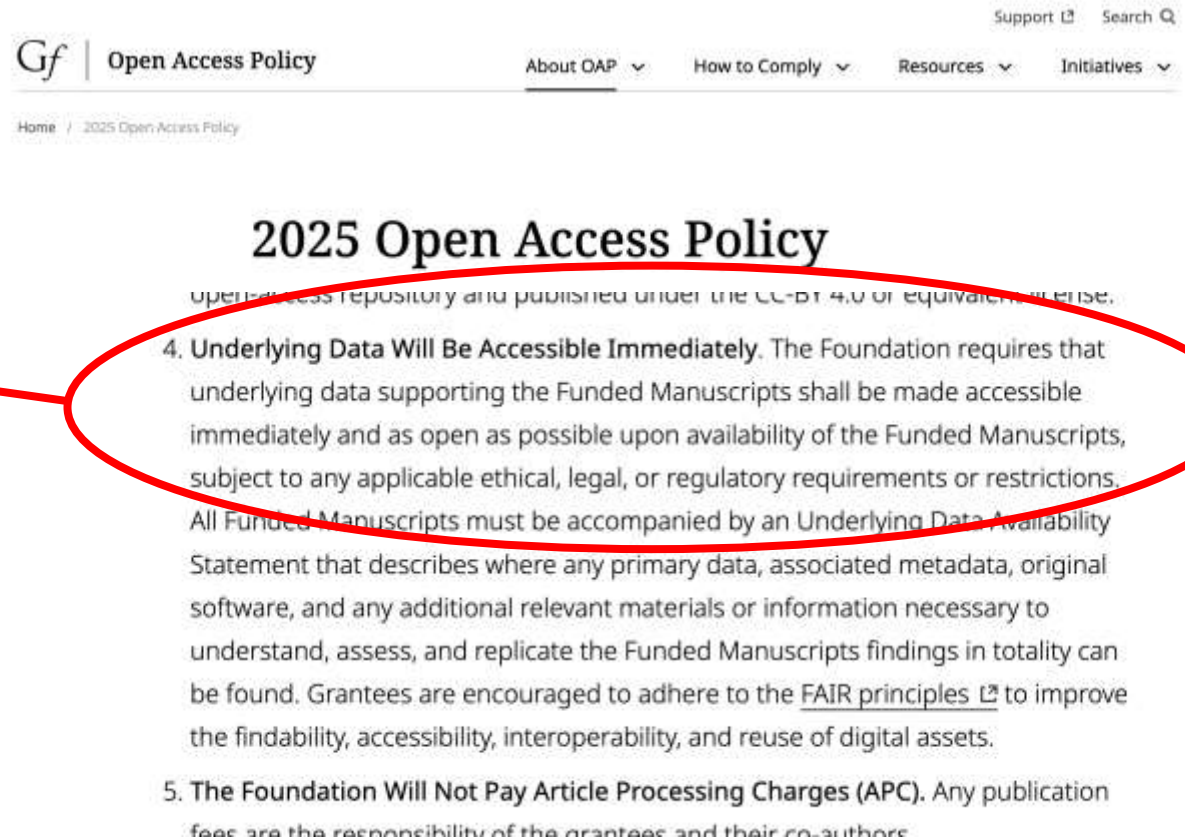
Data Sharing Plan (DSP)

Study protocols



Informed consent forms





Funding agreement

Inter-institutional agreements



The screenshot shows the '2025 Open Access Policy' page from the Gates Foundation. The page header includes the logo 'Gf | Open Access Policy' and navigation links for 'About OAP', 'How to Comply', 'Resources', and 'Initiatives'. The main heading is '2025 Open Access Policy'. Below the heading, there is a list of policy points. Point 4, 'Underlying Data Will Be Accessible Immediately', is circled in red. A blue bracket on the left side of the slide groups the four items in the list (Study protocols, Informed consent forms, Funding agreement, and Inter-institutional agreements) and points to the circled text in the screenshot.


Support  Search 

Gf | Open Access Policy About OAP  How to Comply  Resources  Initiatives 

Home / 2025 Open Access Policy

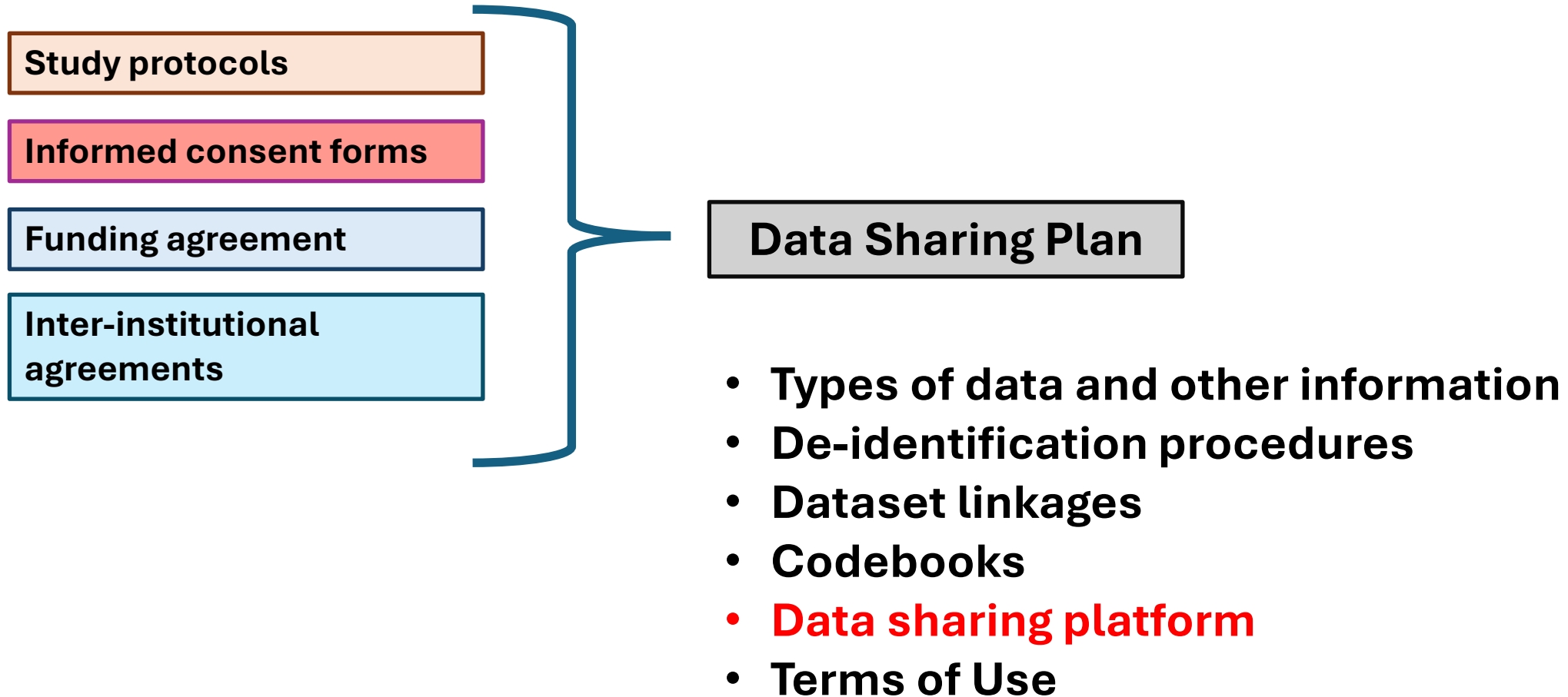
2025 Open Access Policy

open access repository and published under the CC-BY 4.0 or equivalent license.

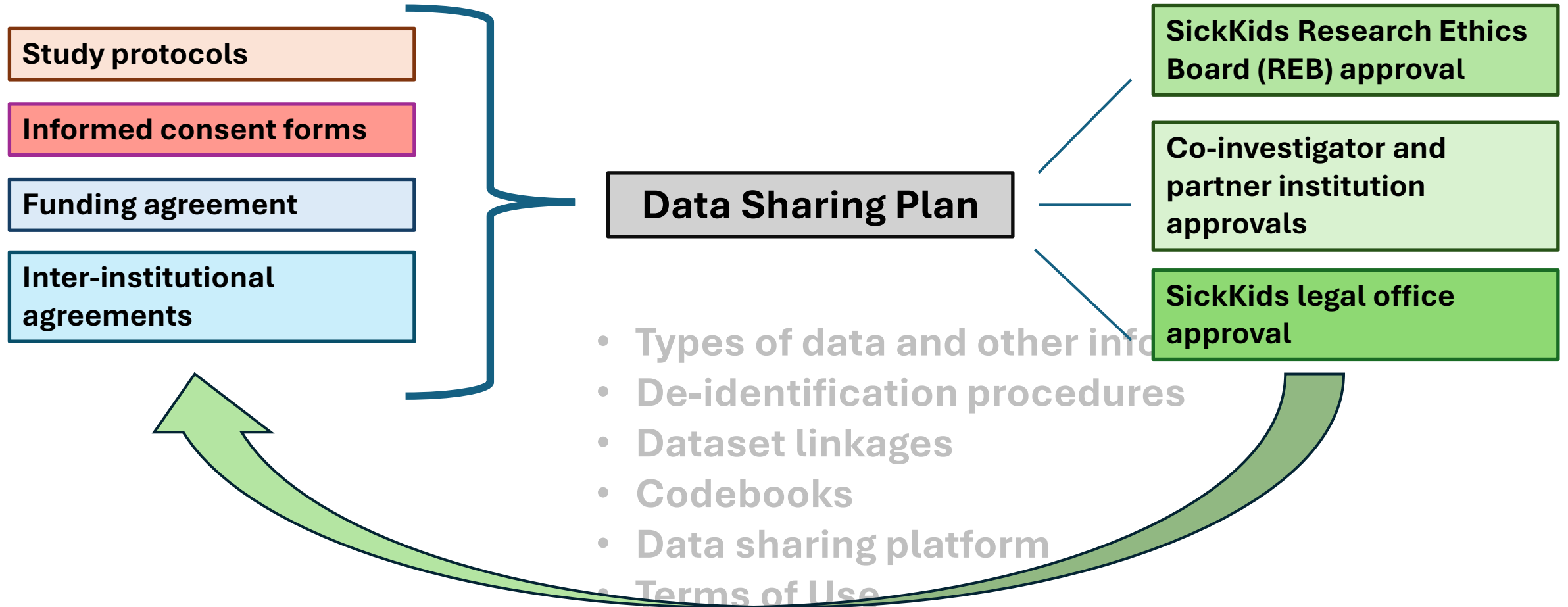
4. **Underlying Data Will Be Accessible Immediately.** The Foundation requires that underlying data supporting the Funded Manuscripts shall be made accessible immediately and as open as possible upon availability of the Funded Manuscripts, subject to any applicable ethical, legal, or regulatory requirements or restrictions. All Funded Manuscripts must be accompanied by an Underlying Data Availability Statement that describes where any primary data, associated metadata, original software, and any additional relevant materials or information necessary to understand, assess, and replicate the Funded Manuscripts findings in totality can be found. Grantees are encouraged to adhere to the [FAIR principles](#)  to improve the findability, accessibility, interoperability, and reuse of digital assets.
5. **The Foundation Will Not Pay Article Processing Charges (APC).** Any publication fees are the responsibility of the grantees and their co-authors.

<https://openaccess.gatesfoundation.org/open-access-policy/>

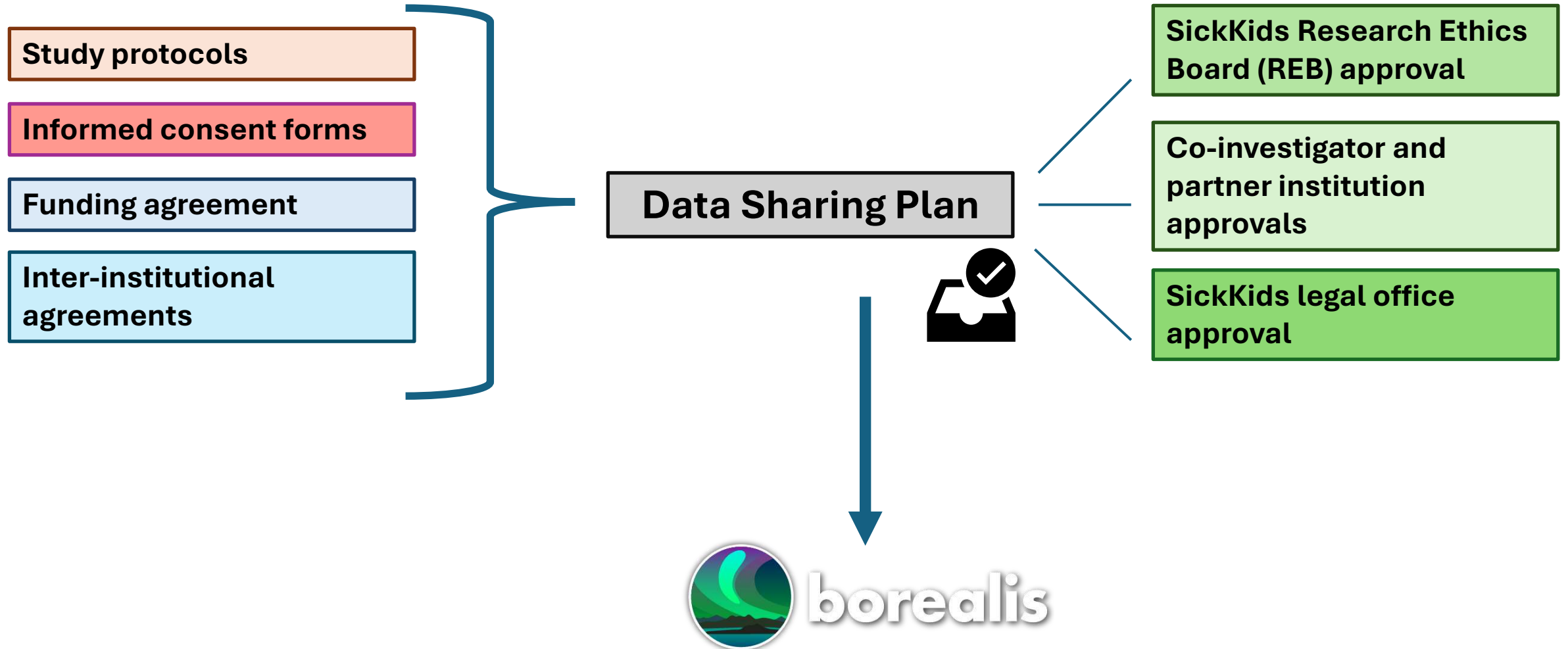
Data Sharing Plan

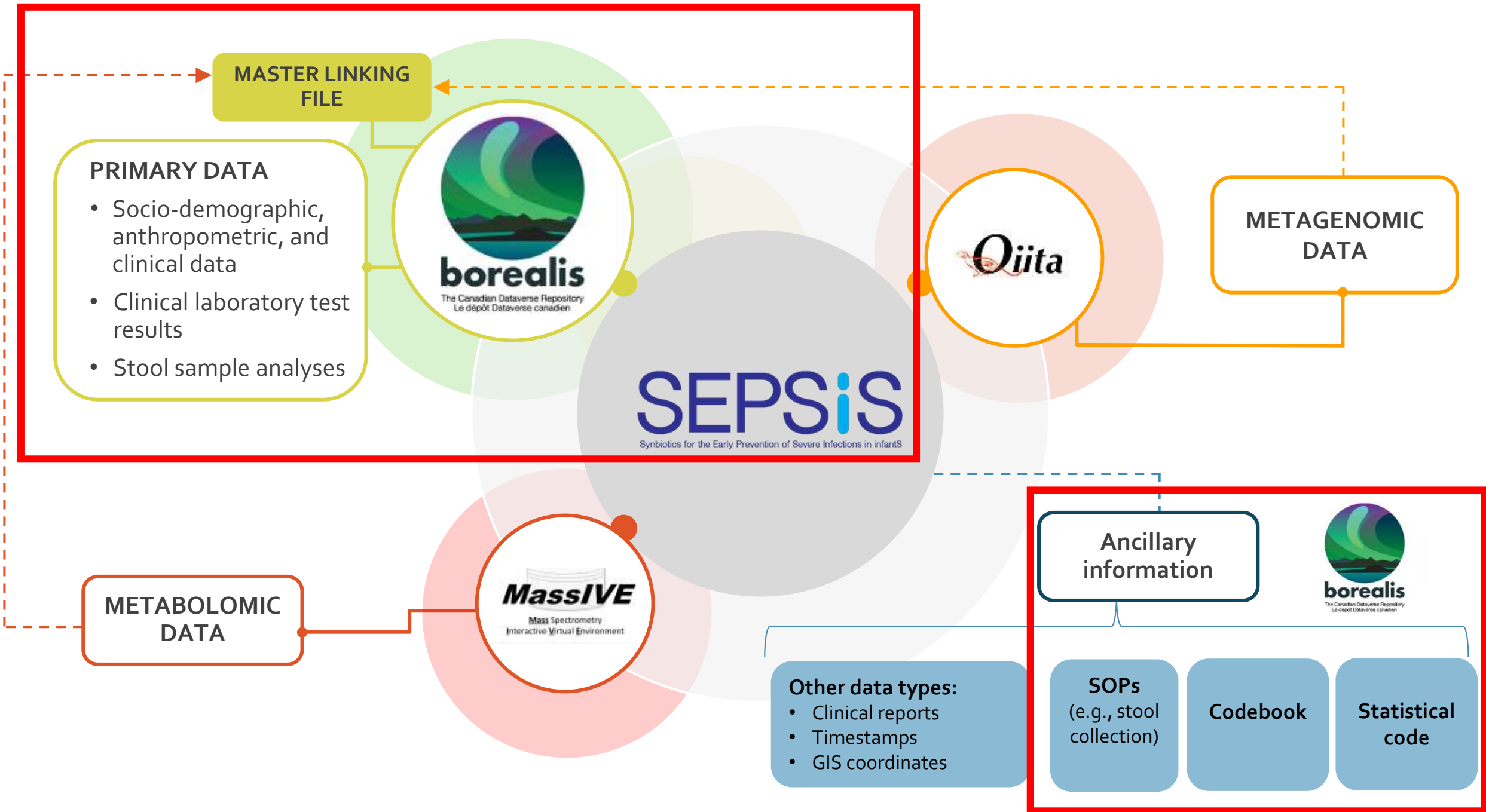


Data Sharing Plan



Data Sharing Plan



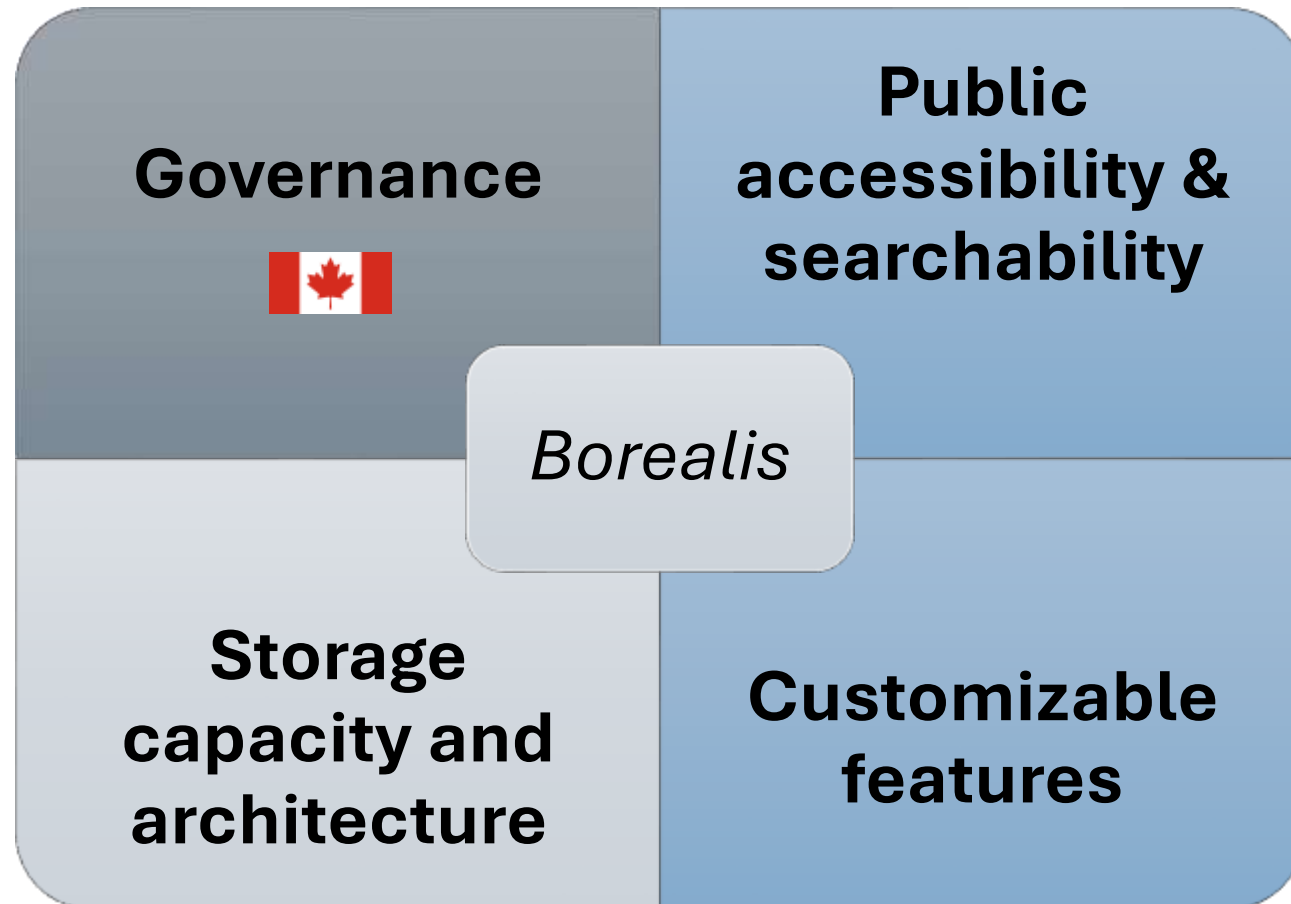


Borealis – Canadian Dataverse Repository



borealis

<https://borealisdata.ca>



- a. Custom Terms of Use**
- b. Guestbook** - user data (i.e., name, email, institution, and position)



<https://borealisdata.ca>

Dataverse:

Sub-dataverses:

https://borealisdata.ca/dataverse/SEPSiS_Project

The screenshot shows the Borealis SEPSiS project page. At the top, the Borealis logo and 'SEPSiS' title are visible, along with navigation links for Search, User Guide, Support, English, and Daniel Roth. The main content area includes a 'SEPSiS Project' section with a description: 'By accessing or using the SEPSiS study dataset, you attest that you have read, understood, and agree to abide by our Custom Dataset Terms.' Below this is a breadcrumb trail: 'Borealis > University of Toronto Dataverse > Roth Lab Dataverse >'. There are buttons for 'Contact', 'Share', and 'Edit'. The 'SEPSiS Project' section is followed by a description of the dataverse and two study protocols: 1) Severe infections and the intestinal microbiome in young infants in Dhaka, Bangladesh; 2) Safety, tolerability and effects on the microbiome of neonatal administration of Lactiplantibacillus plantarum ATCC 202195. Below the protocols are three sub-dataverse cards: 'SEPSiS - Observational cohort study manuscripts', 'SEPSiS - L. plantarum phase II trial manuscripts', and 'SEPSiS - Protocols and supporting materials'. At the bottom, there is a search bar with 'Advanced Search' and 'Add Data' buttons. The search results show 1 to 3 of 3 results, including the three sub-dataverses listed above.



<https://borealisdata.ca>

Dataset:

Files:

<https://borealisdata.ca/dataset.xhtml?persistentId=doi:10.5683/SP3/MOSXFC>

borealis
SEPSiS
Resource for the Early Prevention of Severe Infections in Infants

Search ▾ User Guide Support English ▾ Log In

SEPSiS - Observational cohort study manuscripts

By accessing or using the SEPSiS study dataset, you attest that you have read, understood, and agree to abide by our Custom Dataset Terms.

(The Hospital for Sick Children (SickKids))

Borealis > University of Toronto Dataverse > Roth Lab Dataverse > SEPSiS Project > SEPSiS - Observational cohort study manuscripts >

Severe infection among young infants in Dhaka, Bangladesh: effect of case definition on incidence estimates

Version 6.0

Fung, Alastair; Heasley, Cole; Pell, Lisa G.; Bassani, Diego G.; Shah, Prakesh S.; Morris, Shaun K.; Hamer, Davidson H.; Islam, Mohammad Shahidul; Al Mahmud, Abdullah; Pullenayegum, Eleanor; Saha, Samir K.; Haque, Rashidul; Hossain, Md Iqbal; Chen, Chun-Yuan; Emdin, Abby; O'Callaghan, Karen M.; Loutel, Miranda G.; Sultana, Shamima; Billah, S.M. Masum; Gaffar, S.M. Abdul; Karim, Enamul; Sayed, Sharika; Yeasmin, Sultana; Hoque, Md. Mahbubul; Ahmed, Tahmeed; Sarker, Shafiqul A.; Roth, Daniel E., 2024, "Severe infection among young infants in Dhaka, Bangladesh: effect of case definition on incidence estimates", <https://doi.org/10.5683/SP3/MOSXFC>, Borealis, V6, UNF:6:E7kgr1SyEpO2q+HoWvx7Dg== [fileUNF]

Cite Dataset ▾ Learn about Data Citation Standards

Description Manuscript data, code, and related materials from the SEPSiS manuscript entitled "Severe infection among young infants in Dhaka, Bangladesh: effect of case definition on incidence estimates" [Submitted for publication]

Subject Medicine, Health and Life Sciences

Keyword Bacterial Infections, Incidence, Infant, Newborn, Resource-Limited Settings, Sepsis

License/Data Use Agreement Custom Dataset Terms

Files Metadata Terms Versions

Search this dataset...

Filter by
File Type: All ▾ Access: All ▾ File Tag: All ▾

Group by Folder
Group by Tag

Sort ▾

1 to 3 of 3 Files

Download ▾

README_SEPSiS_Obs_SI.pdf
Adobe PDF · 115.2 KB
Published Sep 24, 2024
9 Downloads
MD5: 46e...07e

Customizable features



<https://borealisdata.ca>

Dataset Terms

License/Data Use Agreement Our Community Name is **SEPSIS**. Good science practices expect you to properly cite your data. Please use the data citation shown on the dataset page.

Custom Dataset Terms — the following Custom Dataset Terms have been defined for this dataset.

Terms of Use By accessing or using the Synbiotics for the Early Prevention of Severe Infections in Infants (SEPSIS) study dataset, you attest that you have read, understood, and agree to abide by the following terms of use. If you do not or are not authorized by your employer or affiliated institution to agree to abide by these terms, please do not access the dataset through Borealis, and instead contact the SEPSIS investigators directly through the "contact owner" function.

Name

Email *

Institution *

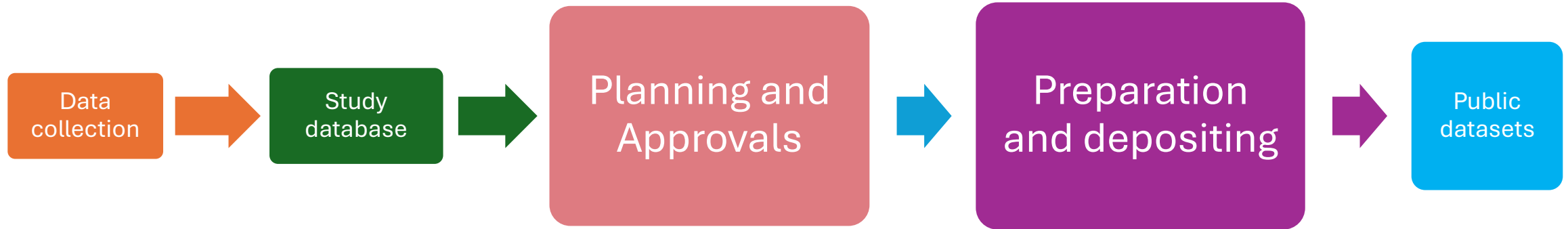
Position *

Additional Questions If the "Institution" field above auto-populates to "OTHER", please enter the name of your institution in the text box below.

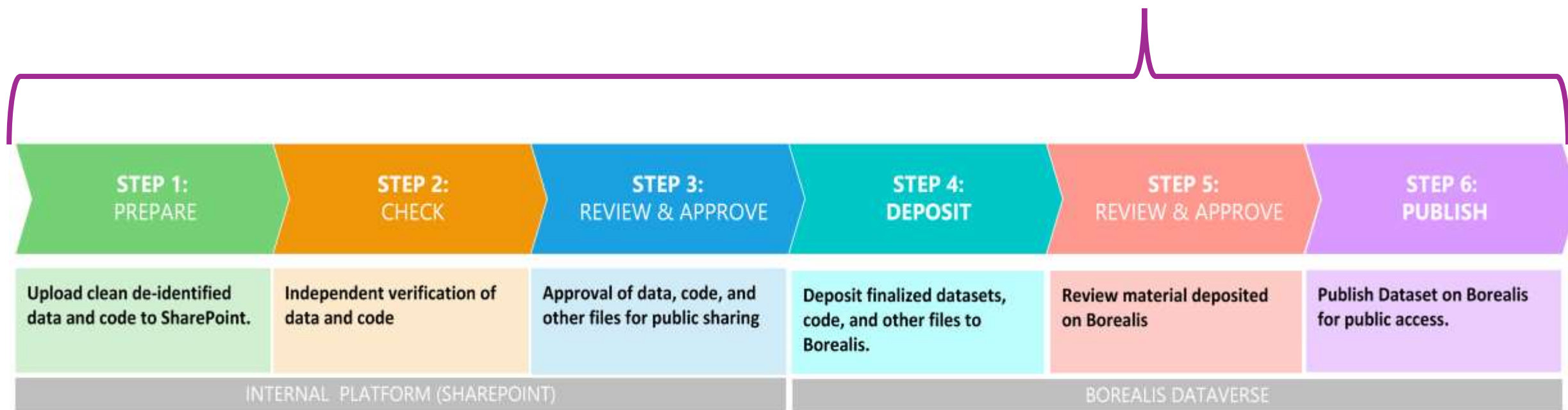
Custom Terms of Use

Guestbook entry

Step-wise process for public sharing of SEPSiS data on Borealis



Step-wise process for public sharing of SEPSiS data on Borealis



Key messages and challenges

- **Plan in advance.**
 - Needs to be addressed in protocol and consent forms.
- Requires dedicated time, personnel and expertise.
- No off-the-shelf roadmap for depositing primary data from REB-approved human subjects research
 - Engage your collaborating partners.
 - Consider requirements of funding agencies.
- Institutional approval processes
 - Institutional norms that emphasize privacy/confidentiality
- Linkages across different platforms.
- Consider the user experience (e.g., formatting files in a way to alleviate need to repeatedly sign guestbook)



Acknowledgements

- SickKids SEPSiS team
 - Veselina Stefanova
- SEPSiS project partner organizations
- Gates Foundation
- Borealis



Sharing research data: why, when, and how?

Michael M. Hoffman

Princess Margaret Cancer Centre
University Health Network

Department of Medical Biophysics
University of Toronto

Vector Institute

<https://hoffmanlab.org/>

Bluesky: [@michaelhoffman.bsky.social](https://bsky.app/profile/@michaelhoffman.bsky.social)

Disclosures

Competing interests:

- Patent application, licensed to Adela

Full disclosure of potential competing interests:

<https://github.com/michaelmhoffman/disclosure/>

Why share data?

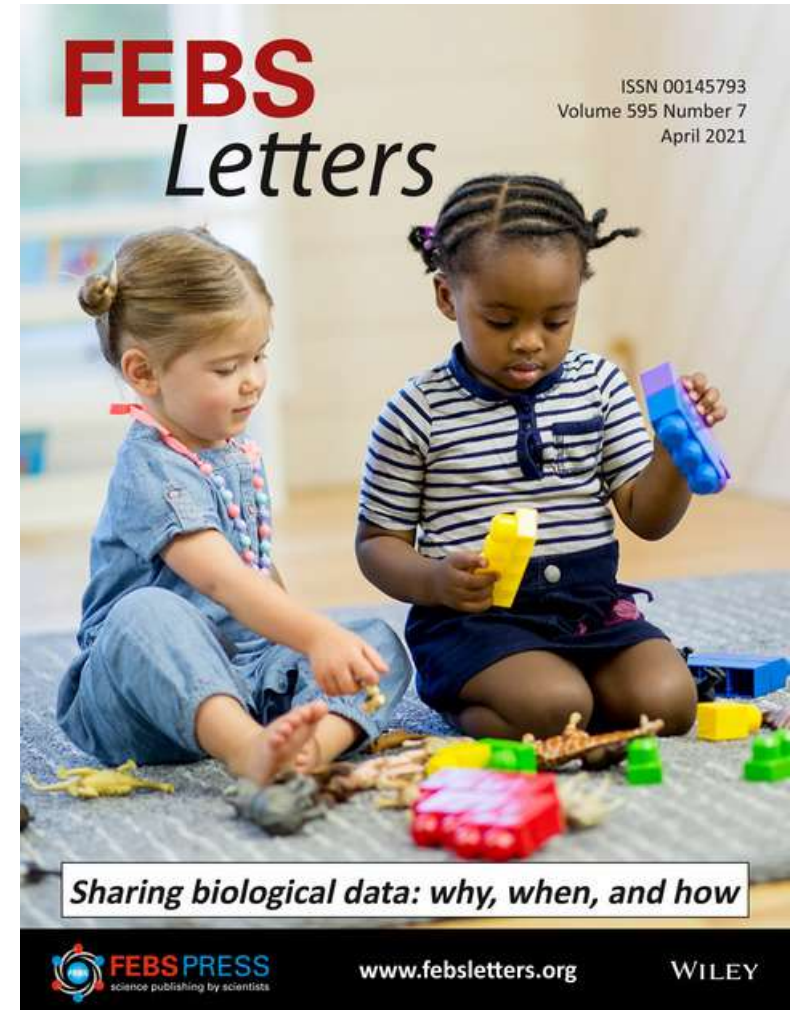
- Further progress of research and science
- Maximize the benefits provided by investment of funds, work, participation
- Best way to ensure that you can get the data later!

Wh data is needed?

- Can never reproduce without the data
- Raw data required
- Pre-existing data?
 - Include information and code required to download and pre-process data

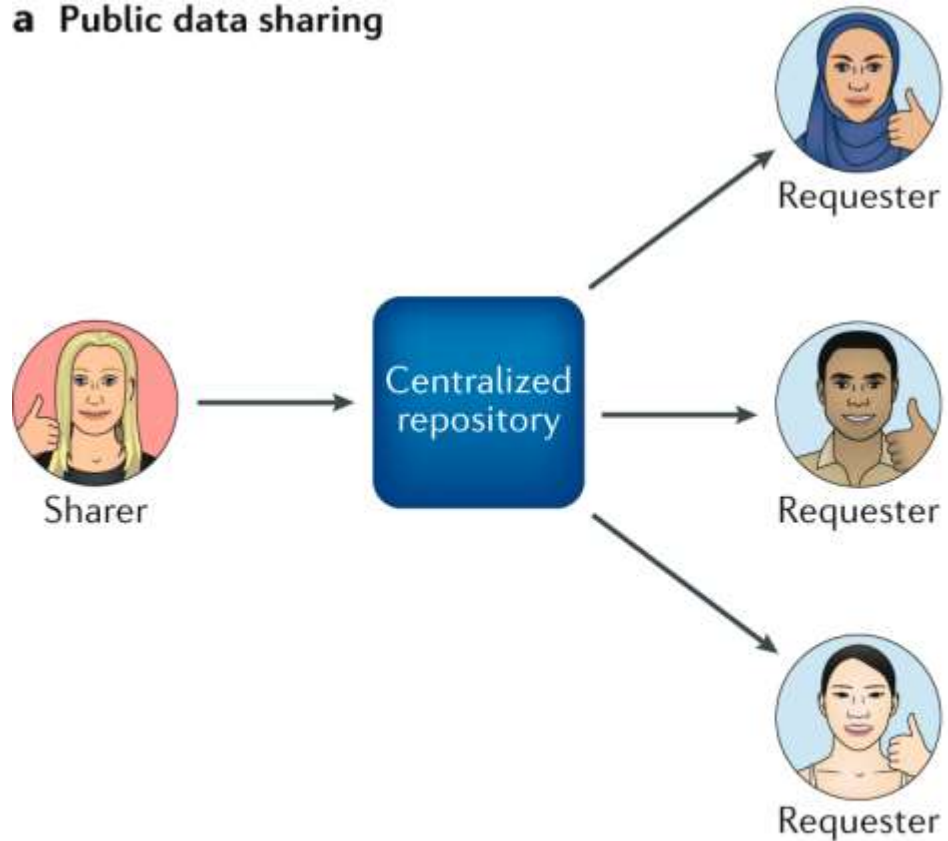
Where does the data go?

- Specialist repositories, when possible
 - Gene expression → Gene Expression Omnibus
 - Microscopy → BioImage Archive
- Generalist repositories, otherwise
 - ≤50 GB: Zenodo
 - >50 GB: Dryad
- Never under 100% author control
 - Not lab web site!
 - Not GitHub!!
 - Not an S3 bucket!!!



Sharing data

a Public data sharing



```

{
  "average_age": "56 years",
  "eye_colors": {
    "Brown": "2",
    "Green": "4",
    "Blue": "2",
  },
  "sex": {
    "male": "4",
    "female": "4",
  },
  "zipcode": {
    "59448": "3",
    "43001": "5"
  }
}

```



a Summarized participant data

```

{
  "average_age": "60 years",
  "eye_colors": {
    "Brown": "4",
    "Green": "8",
    "Blue": "4",
  },
  "sex": {
    "male": "8",
    "female": "8",
  },
  "zipcode": {
    "59448": "6",
    "43001": "10"
  }
}

```



b Summarized participant data with noise

```

{
  "participant_id": "1",
  "age": "34 years",
  "eye_color": "Brown",
  "sex": "male",
  "zipcode": "59448"
},
{
  "participant_id": "2",
  "age": "44 years",
  "eye_color": "Blue",
  "sex": "female",
  "zipcode": "43001"
}
.....

```

Participant data



c Redacted participant data



```

{
  "participant_id": "1",
  "age": "34 years",
  "eye_color": "Brown",
  "sex": "male",
  "zipcode": "59448" "state": "Montana"
},
{
  "participant_id": "2",
  "age": "44 years",
  "eye_color": "Blue",
  "sex": "female",
  "zipcode": "43001" "state": "Ohio"
}
.....

```

d Synthetic participant data



```

{
  "participant_id": "1",
  "age": "34 years",
  "eye_color": "Blue",
  "sex": "female",
  "zipcode": "59448"
},
{
  "participant_id": "2",
  "age": "44 years",
  "eye_color": "Green",
  "sex": "female",
  "zipcode": "43001"
}
.....

```




Reproducibility standards for machine learning in the life sciences

To make machine-learning analyses in the life sciences more computationally reproducible, we propose standards based on data, model and code publication, programming best practices and workflow automation. By meeting these standards, the community of researchers applying machine-learning methods in the life sciences can ensure that their analyses are worthy of trust.

Benjamin J. Heil, Michael M. Hoffman, Florian Markowitz, Su-In Lee, Casey S. Greene and Stephanie C. Hicks

The field of machine learning has grown tremendously within the past ten years. In the life sciences, machine-learning models are rapidly being adopted because they are well suited to cope with the scale and complexity of biological data. However, there are drawbacks to using such models. For example, machine-learning models can be harder to interpret than simpler models, and this opacity can obscure learned biases. If we are going to use such models in the life sciences, we will need to trust them. Ultimately all science requires trust —no scientist can reproduce the results from every paper they read. The question, then, is how to ensure that machine-learning analyses in the life sciences can be trusted.

One attempt at creating trustworthy analyses with machine-learning models revolves around reporting analysis details such as hyperparameter values, model architectures and data-splitting procedures. Unfortunately, such reporting requirements are insufficient to make analyses trustworthy. Documenting implementation details without making data, models and code publicly available and usable by other scientists does little to help future scientists attempting the same analyses and less to uncover biases. Authors can only report on biases they already know about, and without the data, models and code, other scientists will be unable to discover issues post hoc.

For machine-learning models in the life sciences to become trusted, scientists must prioritize computational reproducibility. That is to say that third parties should be able to obtain the same results as the original authors by using their published data, models and code. By doing so, researchers can ensure the accuracy of reported results and detect biases in the models.

Analyses and models that are reproducible by third parties can be examined in depth and, ultimately, become worthy of trust. To that end, we believe the

Table 1 | Proposed reproducibility standards

	Bronze	Silver	Gold
Data published and downloadable	x	x	x
Models published and downloadable	x	x	x
Source code published and downloadable	x	x	x
Dependencies set up in a single command		x	x
Key analysis details recorded		x	x
Analysis components set to deterministic		x	x
Entire analysis reproducible with a single command			x

life sciences community should adopt norms and standards that underlie reproducible machine learning research.

The menu

While many regard the computational reproducibility of a work as a binary property, we prefer to think of it on a sliding scale that reflects the time needed to reproduce. Published works fall somewhere on this scale, which is bookended by 'sovereign', for a completely irreproducible work, and 'zero', for a work where one can automatically repeat the entire analysis with a single keystroke. As in many cases it is difficult to impose a single standard that divides work into 'reproducible' and 'irreproducible', we instead propose a menu of three standards with varying degrees of rigor for computational reproducibility (Table 1).

1. **Bronze standard.** The authors make the data, models and code used in the analysis publicly available. The bronze standard is the minimal standard for reproducibility. Without data, models and code, it is not possible to reproduce a work.
2. **Silver standard.** In addition to meeting the bronze standard: (1) the

dependencies of the analysis can be downloaded and installed in a single command; (2) key details for reproducing the work are documented, including the order in which to run the analysis scripts, the operating system used and system resource requirements; and (3) all random components in the analysis are set to be deterministic. The silver standard is a midway point between minimal availability and full automation. Works that meet this standard will take much less time to reproduce than ones only meeting the bronze standard.

3. **Gold standard.** The work meets the silver standard, and the authors make the analysis reproducible with a single command. The gold standard for reproducibility is full automation. When a work meets this standard, it will take little to no effort for a scientist to reproduce it.

While reporting has become a recent area of focus¹, excellent reporting can look akin to a nutritional information panel. It describes information about a work, but it is insufficient for reproducing the work. In the best case scenario, it provides a summary of what the researchers who conducted

Bronze standard

The minimum standard:

- **Data** published and downloadable
- **Models** published and downloadable
- **Code** published and downloadable

Silver standard

Bronze standard plus:

- **Dependencies** in single command+key
analysis details recorded+deterministic

Gold standard

- Entire analysis reproducible with a **single command**

**How do we make data sharing
a priority?**

Meaningful assessment improves research

- Promotes value of all scholarly outputs

- Journal articles
- Preprints
- Datasets
- Software
- Protocols
- Research materials
- Well-trained researchers
- Societal outcomes and policy changes

- Focuses on the merits of the work

- Reduces JIF-chasing
- Facilitates Open Science practices
- Improves rigor and reproducibility
- Enhances collaboration



DORA
Declaration
of Research
Assessment

Call to action

- **Your papers and grant apps:** point out your careful data sharing!
- Signal reproducibility practices important
 - **Paper reviews:** compare against standard
 - **Grant reviews:** commitment to standard
 - **Hiring:** ask for previous record



Hoffman Lab

New upload

Records

Requests

Members

Settings

Curation policy

29 results found

Sort by Most viewed ▾

Versions

 View all versions

May 30, 2017 (v2)

Dataset

Open

Mappability of the mouse and human genomes and methylomes with Umap and Bismap

Karimzadeh, Mehran ; Ernst, Carl ; Kundaje, Anshul ; and 1 other

This dataset consists of single-read mappability (Bed files) and multi-read mappability (Wiggle files) of human and mouse genomes and methylomes (bisulfite-converted genome). We provide mappability information for the two most recent assemblies of each organism, and for four different read lengths (24 bp, 36 bp, 50 bp, and 100 bp).

Part of Hoffman Lab

Uploaded on May 30, 2017

1 more versions exist for this record

👁 3476 📄 779

Access status

 Open

29

Resource types

 Software

14

 Dataset

9

 Other

5

 Publication

1

May 30, 2017 (v2)

Software

Open

Umap and Bismap: quantifying genome and methylome mappability

Karimzadeh, Mehran ; Ernst, Carl ; Kundaje, Anshul ; and 1 other

The free Umap software package identifies uniquely mappable regions of any genome. Its Bismap extension identifies mappability of the bisulfite converted genome.

Part of Hoffman Lab

Uploaded on May 30, 2017

1 more versions exist for this record

👁 1464 📄 122

Subjects

 Virtual ChIP-seq

4

October 10, 2018 (3.0.0)

Dataset

Open

Virtual ChIP-seq predictions of binding of 36 transcription factor in Roadmap Epigenomics Project tissues

Thank you!

Michael Hoffman

Bluesky: [@michaelhoffman.bsky.social](https://bsky.app/profile/michaelhoffman.bsky.social)



[Wilson et al. 2021. *FEBS Lett* 595:845](#)

Panel Discussion



Upcoming Event

Learning Together Discussion Group: **The Fundamentals of OCAP®**

Overview:

- Access to the *Fundamentals of OCAP®* online course developed by the *First Nations Information Governance Centre*
- Participation in a 6-week synchronous discussion group
- Dates: Spring 2025
- Registration opening soon! Subscribe to the CRIS newsletter to stay informed.



Centre for Research
& Innovation Support

Upcoming Event



SharePoint for Research

Wed, Feb 26, 2025
10:00 - 11:00AM

Register now: <https://cris.utoronto.ca/events/>



 UNIVERSITY OF TORONTO
Centre for Research & Innovation Support



What is a Data Management Plan and Why Do I Need One?

March 26, 2025 | 1:00pm - 2:00pm

REGISTER NOW: [HTTPS://CRIS.UTORONTO.CA/EVENTS/](https://cris.utoronto.ca/events/)

 UNIVERSITY OF TORONTO
Centre for Research & Innovation Support

 UNIVERSITY OF TORONTO
LIBRARIES

Thank you!

- A link to the recording, presenter slides, and feedback form will be sent out after the session
- Follow-up questions can be addressed to cris@utoronto.ca