

Data standardization

A (very) short introduction

Jean-Baptiste Poline

MNI, Brain Imaging Centre, McGill, Montreal

Credit : S. Urchs

Part I: Motivations - if necessary !

Part II: Standardizing : meaning

Part III: How do we do it ?

Part IV: Conclusion

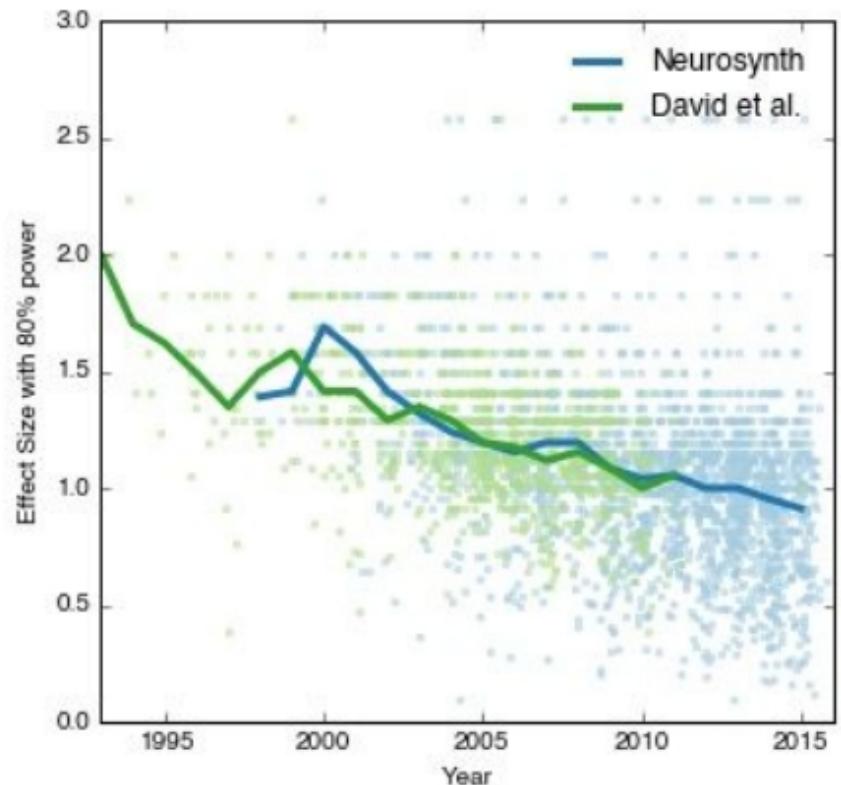
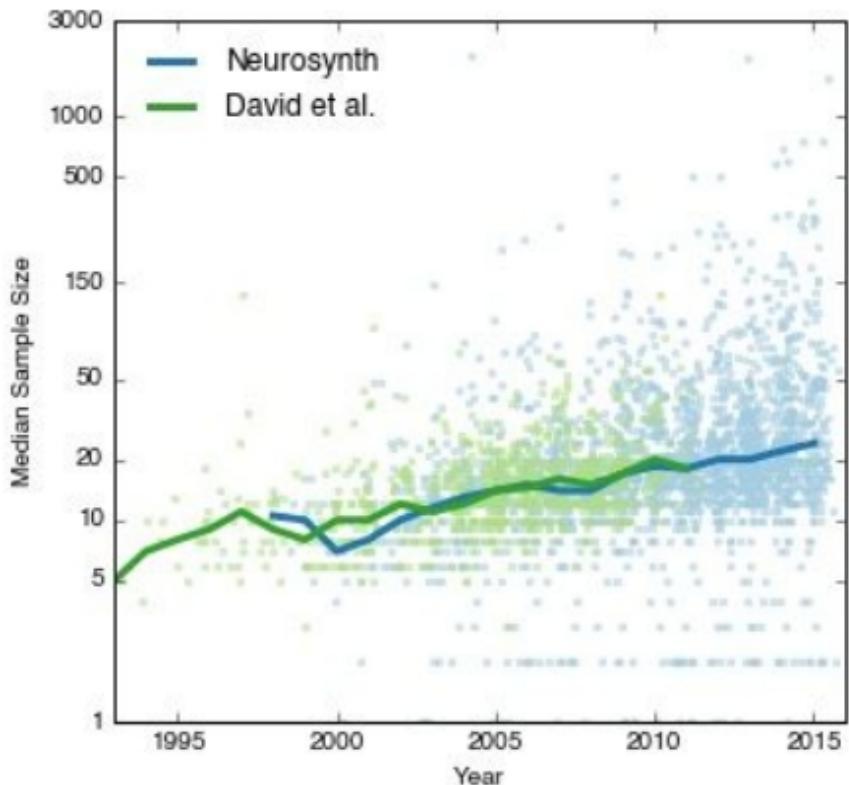
Part I: Motivations - if necessary !

**Part II: What does harmonization
mean ?**

Part III: How do we do it ?

Part IV: Conclusion

Feeling the Future



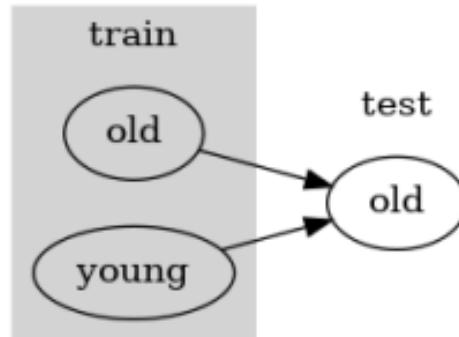
Poldrack et al., PNAS, 2016

Feeling the Future

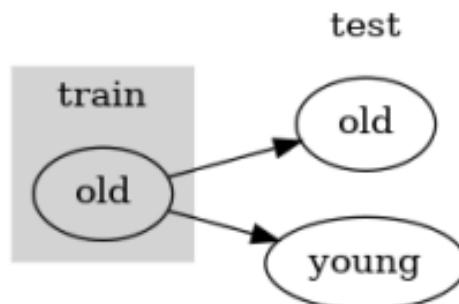
Paradigm	Intersection mask	mask size (vox)	Cohen D			BOLD		
			P10	median	P90	P10	median	P90
MOTOR	Bilateral Precentral Gyrus	12894	0.158	0.628	1.070	0.505	2.707	8.582
	Bilateral Supplementary motor cortex	3418	0.211	0.716	1.197	0.911	4.033	12.510
	Left putamen	1532	0.114	0.513	0.864	0.586	2.388	4.318
	Right putamen	1437	-0.008	0.369	0.749	-0.045	1.696	3.609
WM	Bilateral Middle frontal gyrus	7116	0.101	0.474	0.837	0.130	0.986	2.504
EMOTION	Left amygdala	1133	0.265	0.534	1.065	0.516	1.198	3.379
	Right amygdala	1082	0.308	0.645	1.140	0.581	1.350	3.557
GAMBLING	Left accumbens	455	0.138	0.310	0.461	0.369	0.849	1.440
	Right accumbens	417	0.141	0.332	0.488	0.373	0.981	1.618

With effect size = 0.5 => Power ~ 30%

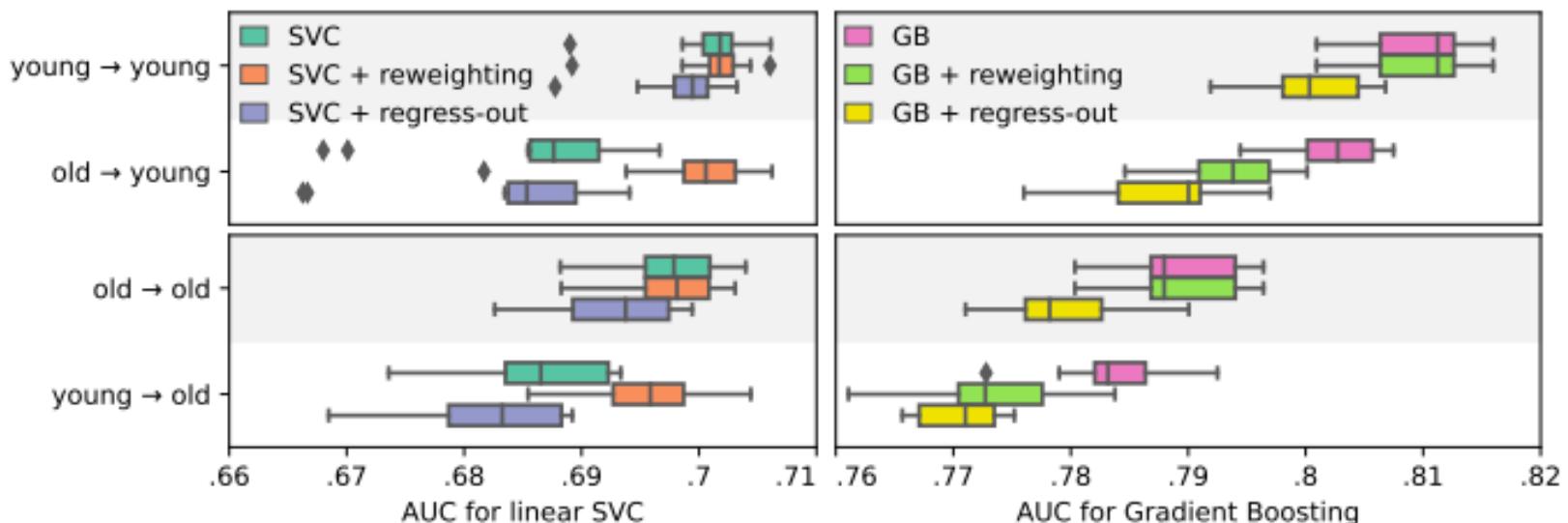
- What do we lose by training on the wrong population?



- How does a model perform when deployed on a new population?

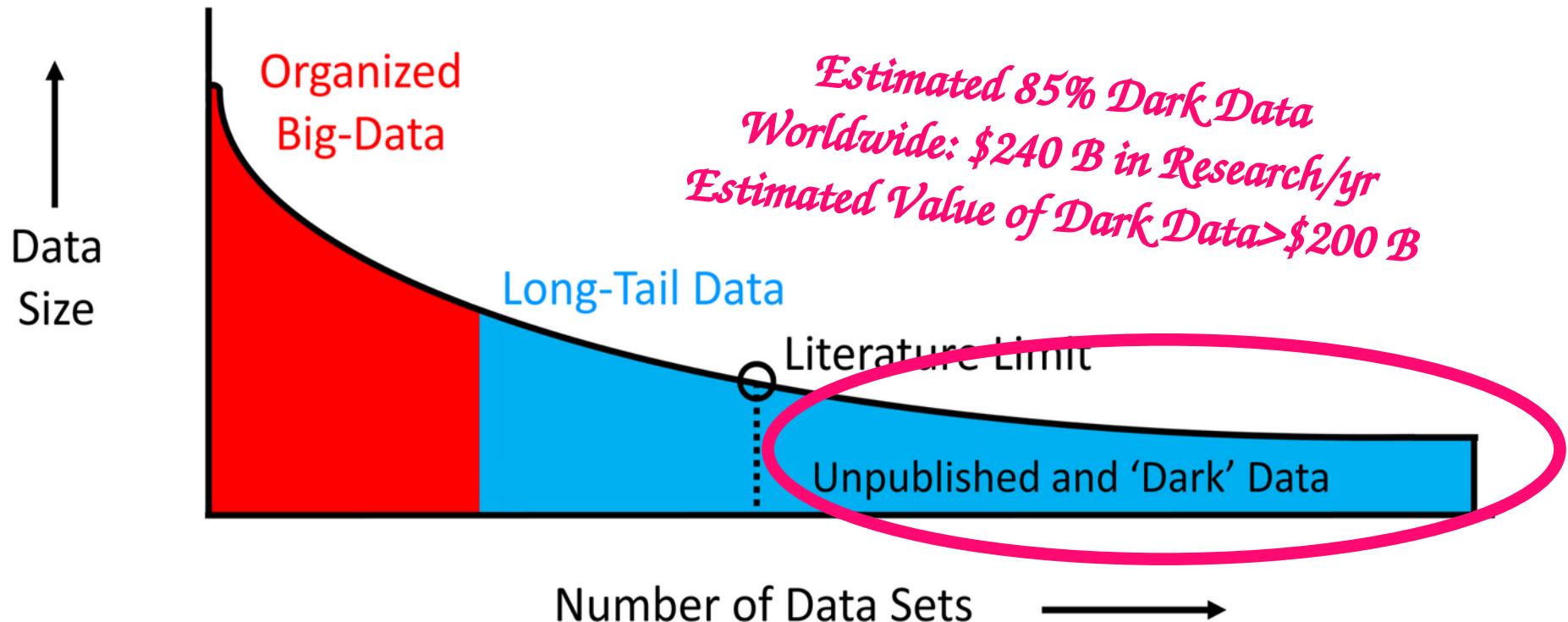


Predicting smoking status in the UKBiobank (10-fold CV, n train = 90K, n test = 9K)



- Dataset shifts threaten generalization
- Simpler models are not always more robust
- "Regressing out nuisance variables" does not help

Jerome Dockes et al., in prep.



McLeod et al., Lancet, 2014

Ferguson et al., 2014, *Nature Neuroscience*

- A less rare case than usually thought !
- No license
- Database not containing what it describes
- Wrong QC – QC unreliable
- Headers of files are not correct (cf the Left/Right issue)
- Provenance of data is lost
- **SAM1 SAM2 SAM3:**
<https://www.youtube.com/watch?v=N2zK3sAtr-4>

From HCP:

"With the releases of FreeSurfer 7.X, there have been some regressions in surface placement performance when running FreeSurfer inside the HCP Pipelines.

At this time, I would recommend sticking with FreeSurfer 6.0 while we get these issues sorted out."

Part I: Motivations - if necessary !

**Part II: What does standardisation
mean - and provides ?**

Part III: How do we do it ?

Part IV: Conclusion

- **Formats and units**

- Information is exactly the same
- Nifti / dicoms
- Age in years or months, Fahrenheit or Celcius

- **Easily “mappable” measures, versions**

- I measure motor control with UDPRS III in my population but I adapted the test to my clinical population

- **Related meaning**

- I measure the concept of anxiety with TestA or TestB

Standardization benefits



individual researcher



open dataset consortium



clinical research institute / clinical consortium
Douglas Research Center / Enigma consortium

“I need to expand my training sample to be more diverse (and larger)”

“Where do I find data of PD patients that have completed the MOCA?”

Standardization benefits



individual researcher



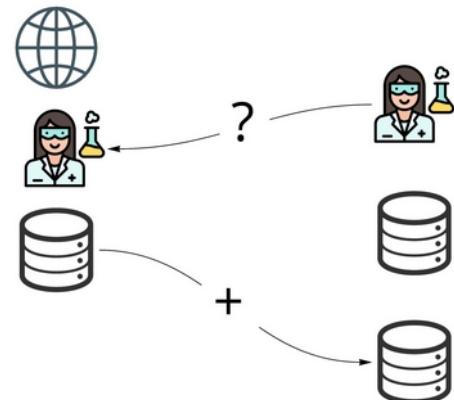
open dataset consortium



clinical research institute / clinical consortium
Douglas Research Center / Enigma consortium

"I need to expand my training sample to be more diverse (and larger)"

"Where do I find data of PD patients that have completed the MOCA?"



Standardization benefits



individual researcher



open dataset consortium



clinical research institute / clinical consortium
Douglas Research Center / Enigma consortium

“I need to expand my training sample to be more diverse (and larger)”

“Where do I find data of PD patients that have completed the MOCA?”

“What data do our PIs have?”

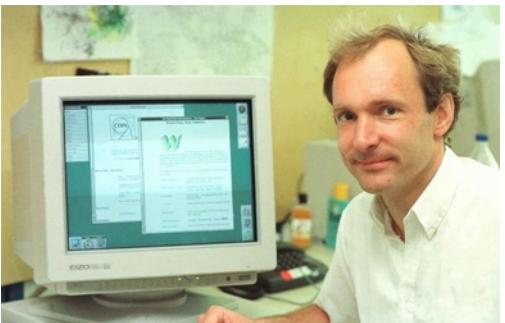
And

“How can we advertise our data without sharing them?”

schema.org

- What if we had google power on datasets ?
- The briefest introduction ever to schema.org

Documents



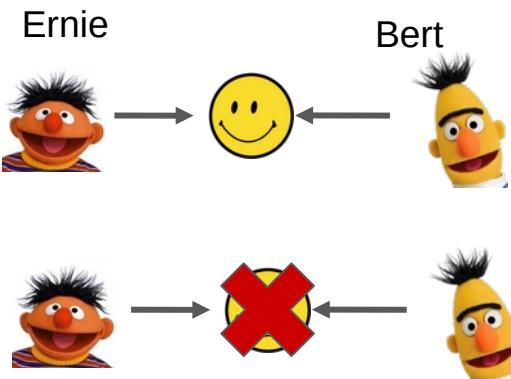
World Wide Web (WWW)

- Physical location is abstracted away
- You just care about documents
- Hyperlinks



Tim Berners-Lee, CERN
March 1989, May 1990

Services



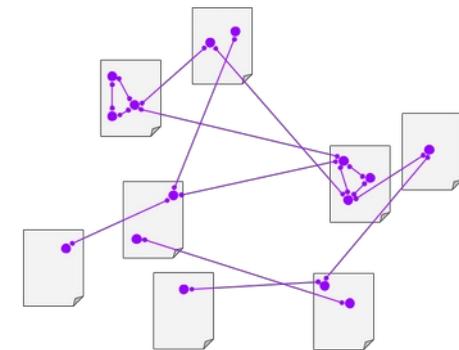
“List all the friends of Ernie!”



Create information

Meanings

“Link facts, not documents”



“Make the meaning of facts machine readable”



Berners-Lee, T., Hendler, J. Publishing on the semantic web. Nature 410, 1023–1024 (2001).

Link data : 0.101

RDF triples form (linked) graphs

<Bob> <is a> <person>.

<Bob> <is a friend of> <Alice>.

<Bob> <is born on> <the 4th of July 1990>.

<Bob> <is interested in> <the Mona Lisa>.

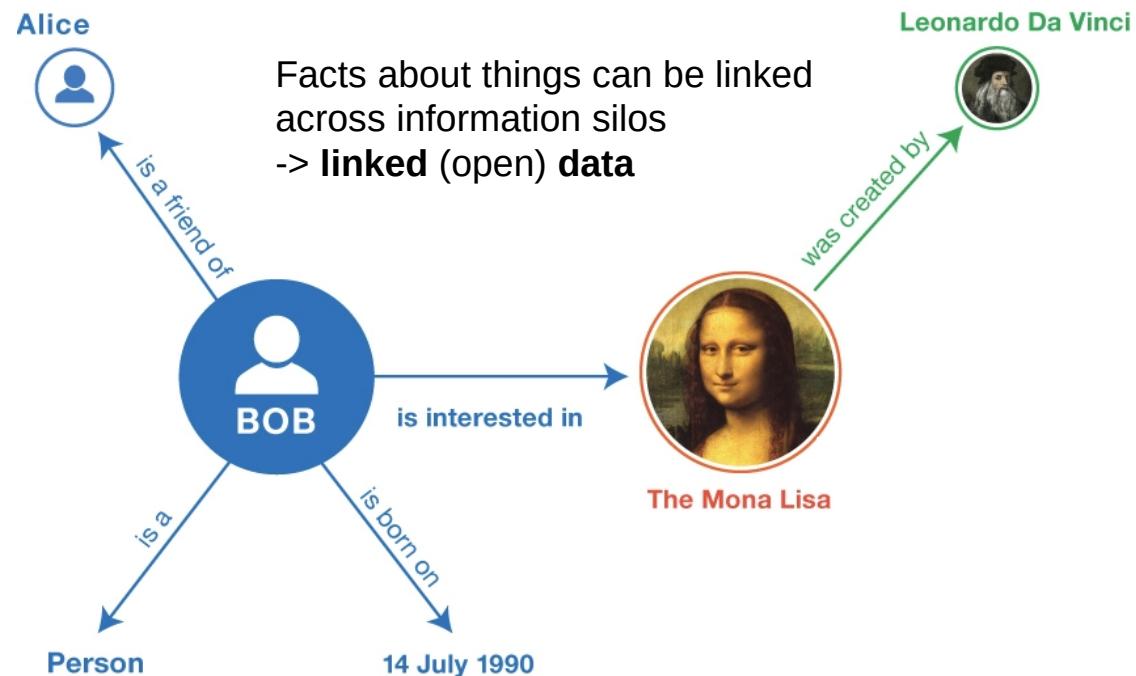
<the Mona Lisa> <was created by> <Leonardo da Vinci>.

vocabularies? -> many!

Theory



Practice



<https://www.w3.org/TR/rdf11-primer/>

OMOP origin stories: the 2004 VIOXX scandal



FDA conducts regulatory action to withdraw approved painkiller because of **increased stroke risk**

Merck was alleged to have **known and suppressed** this information!

- > lawsuits
- > congress hearings



nature

[Explore content](#) [About the journal](#) [Publish with us](#)

[nature](#) > [news](#) > [article](#)

[Published: 13 November 2007](#)

Merck settles Vioxx lawsuits for \$4.85 billion

[Meredith Wadman](#)

[Nature](#) (2007) | [Cite this article](#)

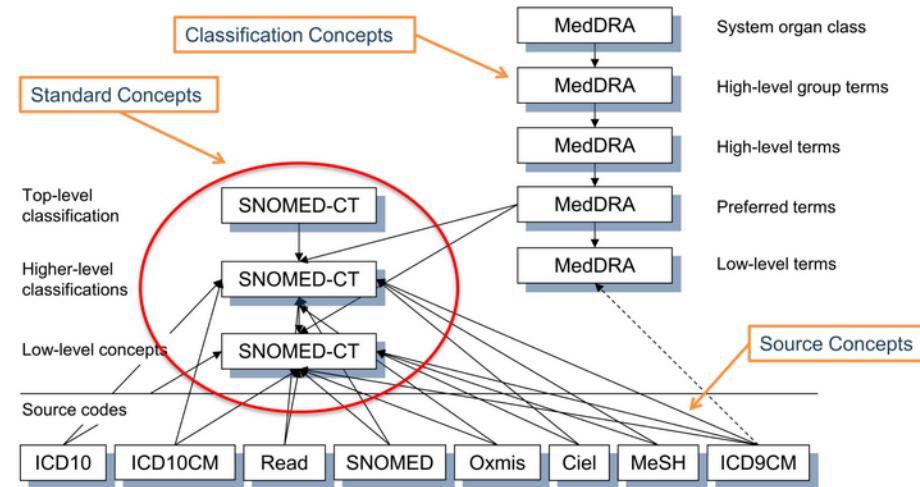
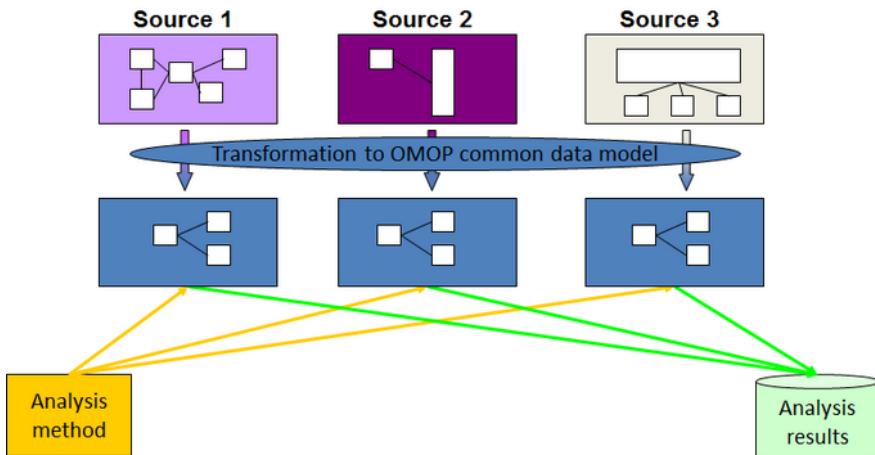
3893 Accesses | 3 Citations | 129 Altmetric | [Metrics](#)

But drug firm maintains it was not at fault over arthritis drug.

Three years after it pulled its blockbuster painkiller Vioxx from the market, Merck has agreed to pay \$4.85 billion to settle nearly 27,000 lawsuits that claim the arthritis drug caused heart attacks and strokes.

No system to integrate observational data

OMAP



* Transformation to a common data model

* generally only one vocabulary per domain (**SNOMED** for diagnosis, **RxNorm** for drugs ...)



PARTICIPANTS

126 661 070 people observed for at least 365 days before 1 January 2017, 2018, or 2019 from 13 databases.

Part I: Motivations - if necessary !

**Part II: What does harmonization
mean - and provides ?**

Part III: How do we do it ?

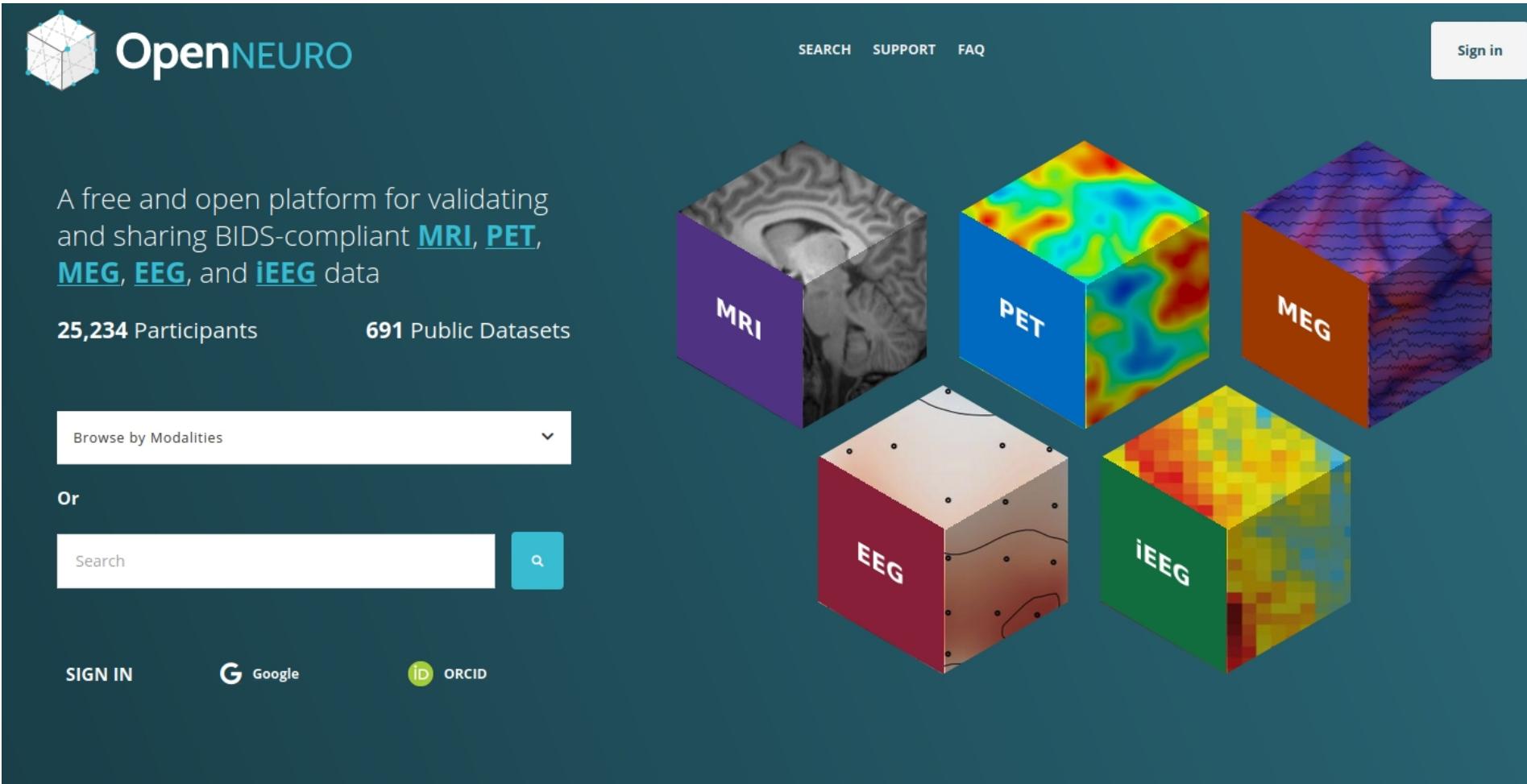
Part IV: Conclusion

Use the simple tools first



- File name **matters**
- Format **matters**
 - Excel or CSV ? Why?
 - Docx or markdown?
 - Proprietary ?
 - License ?
 - Copyright ? (quarto)
 - Supported ?
 - ...
- Where we store **matters**
 - Platform governance
 - Security / ...
 - No free lunch: When will I pay ?
 - ...

- BIDS
- NIDM terms
- BIDS->NIDM
- Show annotation and search tools
- Why should I put my data in BIDS ?
 - And how much will I suffer in the process ?



The screenshot shows the OpenNeuro website homepage. At the top left is the "OpenNEURO" logo with a 3D cube icon. At the top right are links for "SEARCH", "SUPPORT", and "FAQ", and a "Sign in" button. Below the header, a main text area says: "A free and open platform for validating and sharing BIDS-compliant [MRI](#), [PET](#), [MEG](#), [EEG](#), and [iEEG](#) data". To the left, there are statistics: "25,234 Participants" and "691 Public Datasets". Below these are two search/filter sections: "Browse by Modalities" (with a dropdown arrow) and a "Search" bar with a magnifying glass icon. At the bottom, there are links for "SIGN IN", "Google", and "ORCID". On the right side, five 3D cubes represent different data modalities: MRI (purple, showing brain slices), PET (blue, showing metabolic activity), MEG (orange, showing magnetic field patterns), EEG (red, showing electrode placement), and iEEG (green, showing spatial activity patterns).

BIDS is the “worst standard to the exclusion of all others”

- *Participant.tsv* variables
 - Variables **column names** are not standardized
 - a much wider space than imaging !
 - Variable **values** (i.e. data element, cells) are not standardized
- *Sidecar json* files
 - Terms are not harmonized
 - Content not standardized
- Consequence:
 - You cannot easily recreate a dataset from several datasets / search across participants

- 📁 dicomdir/
 - 📁 1208200617178_22/
 - 📄 1208200617178_22_8973.dcm
 - 📄 1208200617178_22_8943.dcm
 - 📄 1208200617178_22_2973.dcm
 - 📄 1208200617178_22_8923.dcm
 - 📄 1208200617178_22_4473.dcm
 - 📄 1208200617178_22_8783.dcm
 - 📄 1208200617178_22_7328.dcm
 - 📄 1208200617178_22_9264.dcm
 - 📄 1208200617178_22_9967.dcm
 - 📄 1208200617178_22_3894.dcm
 - 📄 1208200617178_22_3899.dcm
 - 📁 1208200617178_23/
 - 📁 1208200617178_24/
 - 📁 1208200617178_25/



- 📁 my_dataset/
 - 📄 participants.ttl → participants.tsv + participants.json
 - 📁 sub-01/
 - 📁 anat/
 - 📄 sub-01_T1w.nii.gz
 - 📁 func/
 - 📄 sub-01_task-rest_bold.nii.gz
 - 📄 sub-01_task-rest_bold.json
 - 📁 dwi/
 - 📄 sub-01_dwi.nii.gz
 - 📄 sub-01_dwi.json
 - 📄 sub-01_dwi.bval
 - 📄 sub-01_dwi.bvec
 - 📁 sub-02/
 - 📁 sub-03/
 - 📁 sub-04/

Example of sidecar json

```
{  
    "MeasurementToolMetadata": {  
        "Description": "Adult ADHD Clinical Diagnostic Scale V1.2",  
        "TermURL": "http://www.cognitiveatlas.org/task/id/trm_5586ff878155d"  
    },  
    "adhd_b": {  
        "Description": "B. CHILDHOOD ONSET OF ADHD (PRIOR TO AGE 7)",  
        "Levels": {  
            "1": "YES",  
            "2": "NO"  
        }  
    },  
    "adhd_c_dx": {  
        "Description": "As child met A, B, C, D, E and F diagnostic criteria",  
        "Levels": {  
            "1": "YES",  
            "2": "NO"  
        }  
    }  
}
```

ADHD200 - Brown

-  sub-0026001
-  sub-0026002
-  sub-0026004
-  sub-0026005

Participants.tsv file:

participant_id	gender	age	handedness	verbal_iq	performance_iq	full4_iq	qc_rest_1	qc_anatomical_1
26001	Male	16.92	Right	133	104	120	Pass	Pass
26002	Male	15.68	Right	106	106	107	Pass	Pass
26004	Female	14.99	Right	119	123	125	Pass	Pass
26005	Female	15.16	Right	116	131	126	Pass	Pass
26009	Male	16.91	Left	113	81	97	Pass	Pass
26014	Female	16.21	Right	101	102	102	Pass	Pass
26015	Female	15.2	Right	127	98	113	Pass	Pass
26016	Male	16.07	Right	120	96	109	Pass	Pass
26017	Female	14.56	Right	95	87	89	Pass	Pass
26022	Male	17.83	Right	105	111	109	Pass	Pass
26024	Female	17.77	Right	89	83	85	Pass	Pass
26027	Female	11.28	Right	108	103	106	Pass	Pass
26030	Female	14.51	Right	121	119	123	Pass	Pass
26039	Female	14.19	Right	125	117	124	Pass	Pass
26040	Female	13.67	Right	111	93	102	Pass	Pass
26041	Female	13.68	Right	129	120	128	Pass	Questionable
26042	Female	13.82	Right	106	110	109	Pass	Pass

Brown Site - ADHD200

- ▶  sub-0026042
- ▶  sub-0026043
- ▶  sub-0026044
- ▶  sub-0026045
- ▶  sub-0026050
- ▶  sub-0026052
- ▶  sub-0026053
- ▶  sub-0026054
- ▶  sub-0026055
- ▶  sub-0026057
-  T1w.json
-  task-rest_bold.json
-  nidm.ttl
-  dataset_description.json
-  participants.json
-  participants.tsv

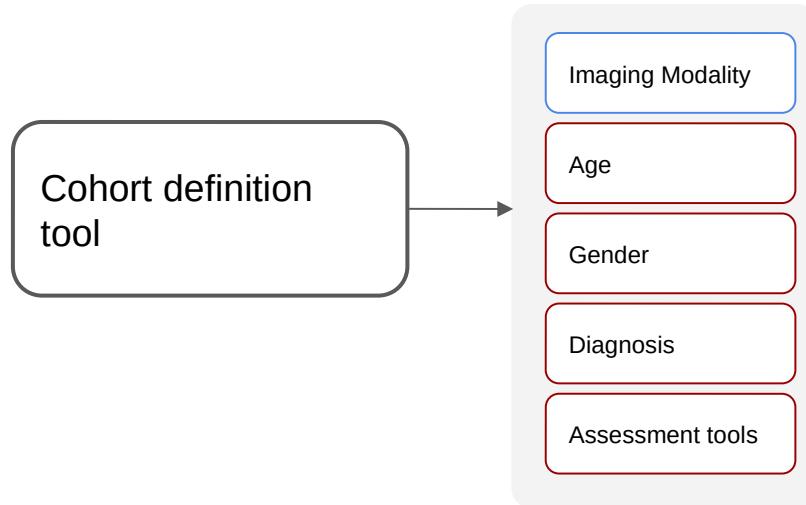
```

"handedness": {
  "minValue": "NA",
  "maxValue": "NA",
  "hasUnit": "NA",
  "label": "handedness",
  "description": "Simple handedness determined using the Edinburgh Handedness Inventory",
  "source_variable": "handedness",
  "valueType": "http://www.w3.org/2001/XMLSchema#complexType",
  "associatedWith": "NIDM",
  "levels": {
    "Left": "0",
    "Right": "1",
    "Ambidextrous": "3"
  },
  "isAbout": {
    "url": "http://uri.interlex.org/base/ilx_0104886",
    "label": "Handedness assessment"
  }
},
"verbal_iq": {
  "hasUnit": "",
  "minValue": "0",
  "maxValue": "200",
  "label": "verbal_iq",
  "description": "verbal IQ standard score measured with WISC-IV, WASI, WISCC-R, two-subt",
  "source_variable": "verbal_iq",
  "valueType": "http://www.w3.org/2001/XMLSchema#integer",
  "associatedWith": "NIDM",
  "isAbout": {
    "url": "http://uri.interlex.org/base/ilx_0739359",
    "label": "verbal intelligence quotient"
  }
}

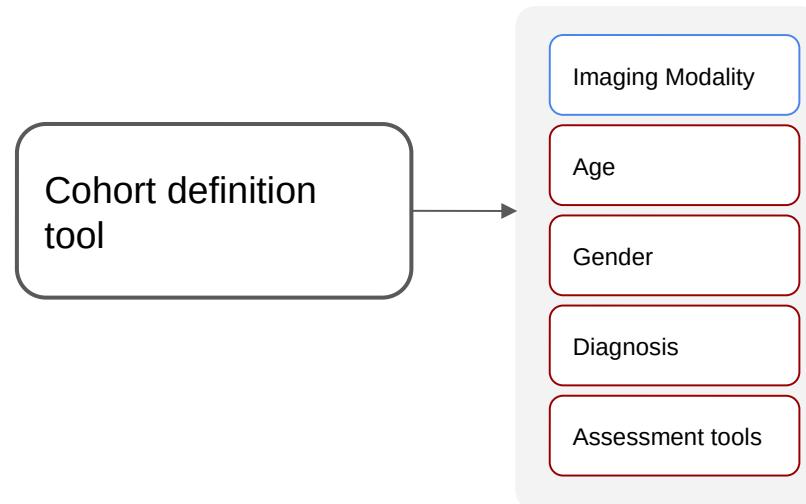
```



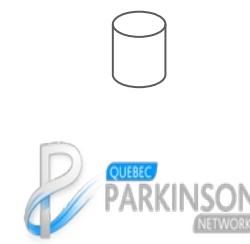
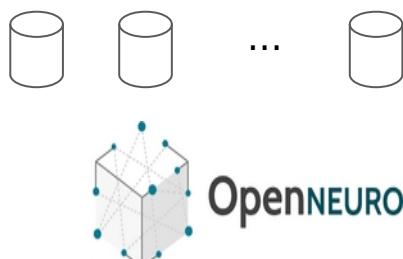
tools and services



tools and services



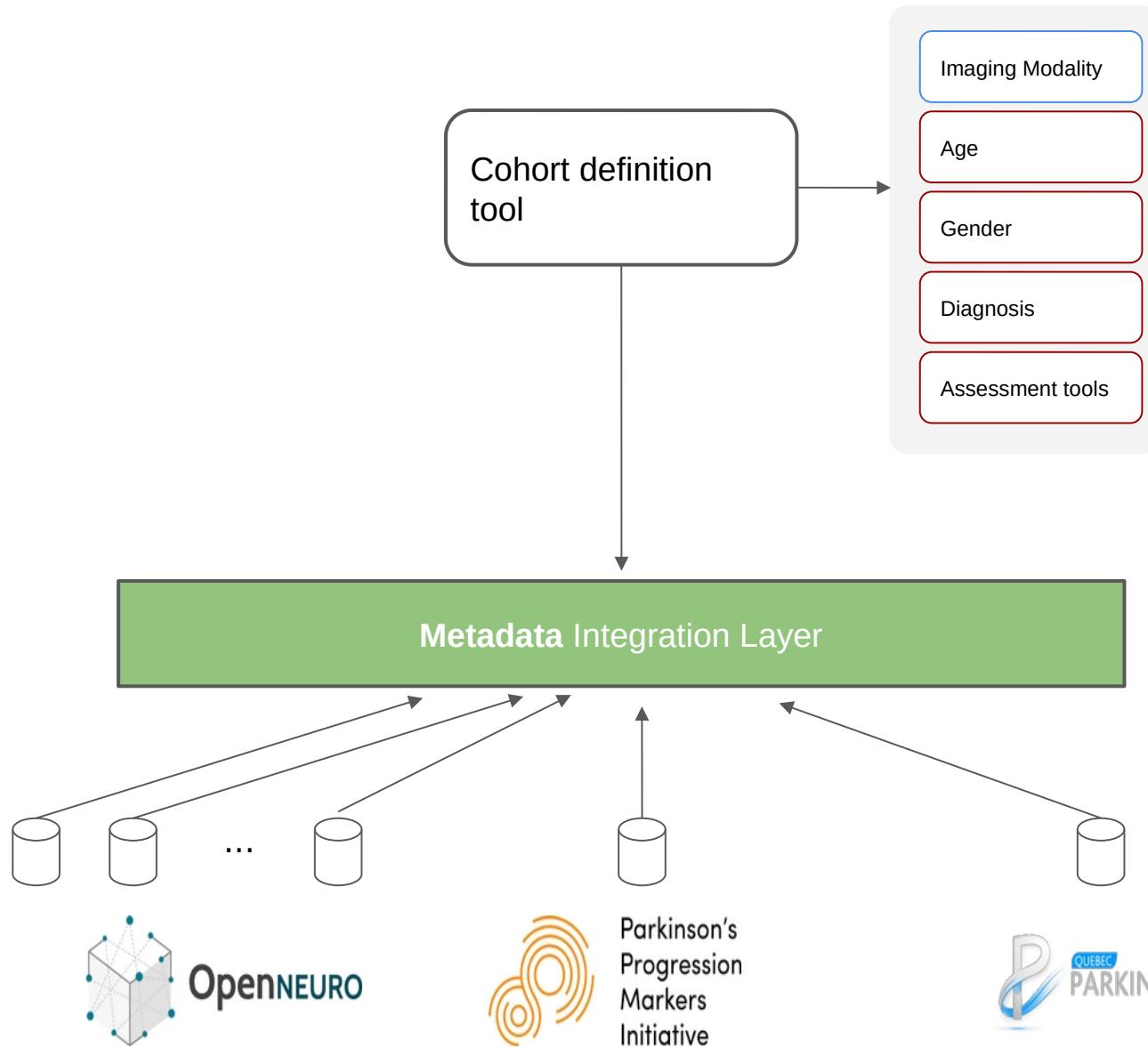
distributed data repositories



tools and services

integrated metadata

distributed data repositories



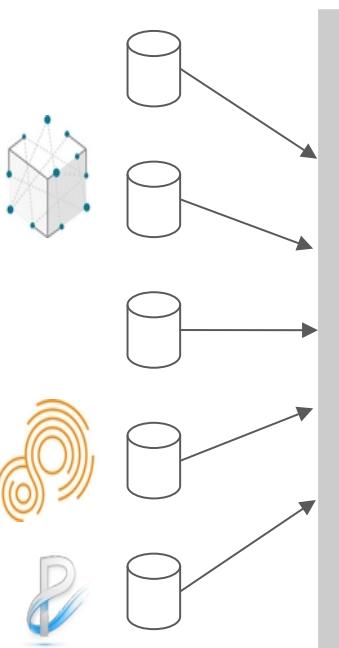
Metadata integration workflow

imaging metadata

```
sub-control01/  
anat/  
  sub-control01_T1w.nii.gz
```

observational metadata

age	gender	group
44	f	depr_no_treatment
21	f	depr_cbt
28	m	depr no treatment



local
metadata

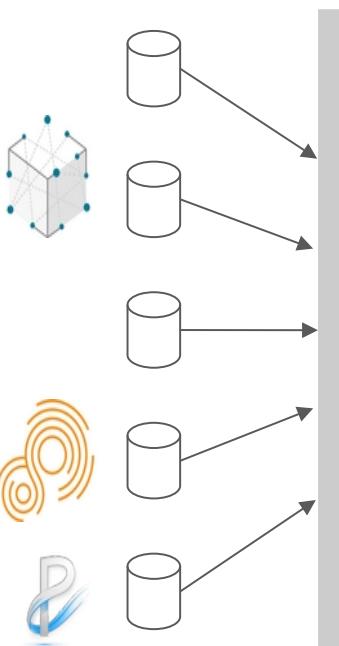
Metadata integration workflow

imaging metadata

```
sub-control01/
  anat/
    sub-control01_T1w.nii.gz
```

observational metadata

age	gender	group
44	f	depr_no_treatment
21	f	depr_cbt
28	m	depr no treatment



local
metadata

1: Integration



Imaging Modality

Age

"019Y"

Gender

"M"

Diagnosis

"PD"

Assessment tools

"MCTOT"

gathering the metadata

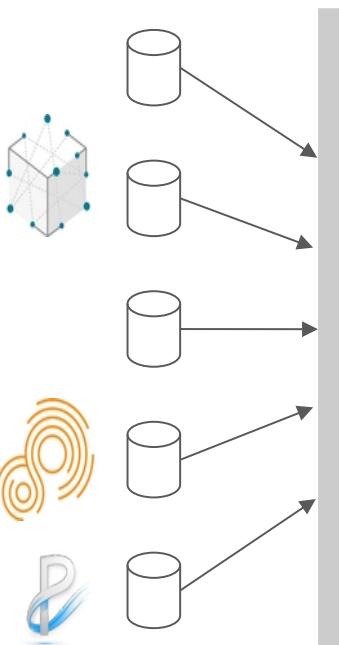
Metadata integration workflow

imaging metadata

```
sub-control01/
  anat/
    sub-control01_T1w.nii.gz
```

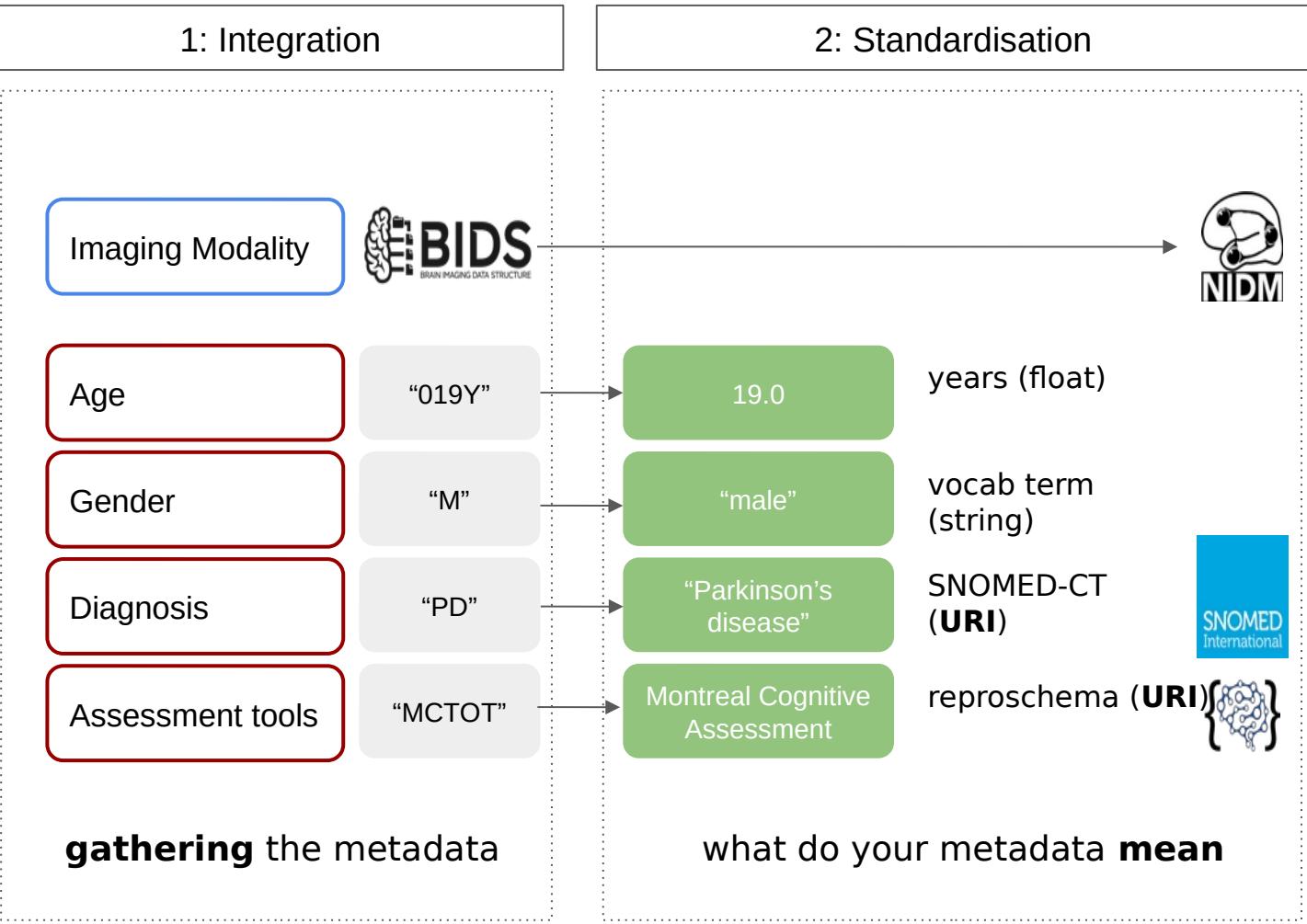
observational metadata

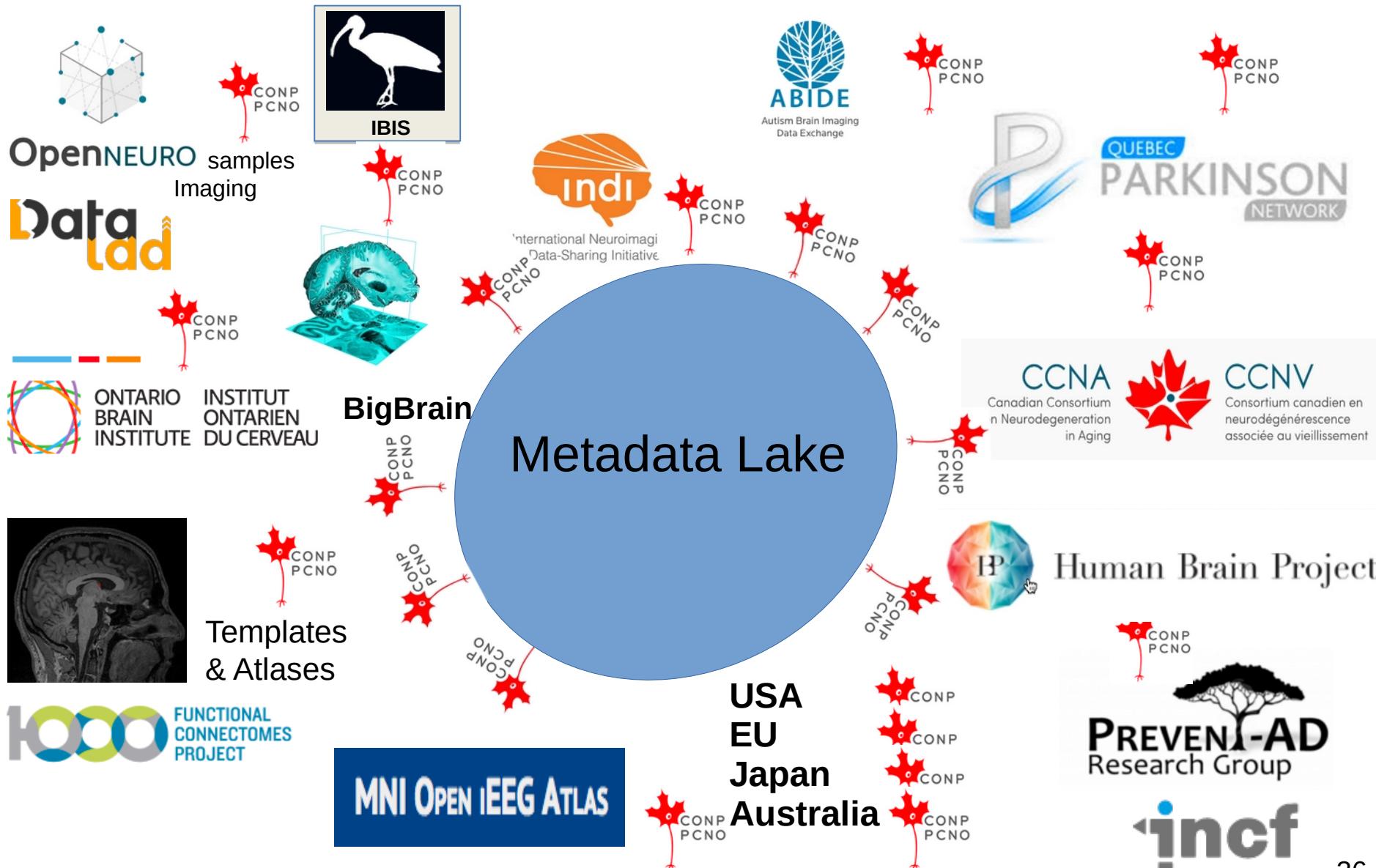
age	gender	group
44f		depr_no_treatment
21f		depr_cbt
28m		depr no treatment



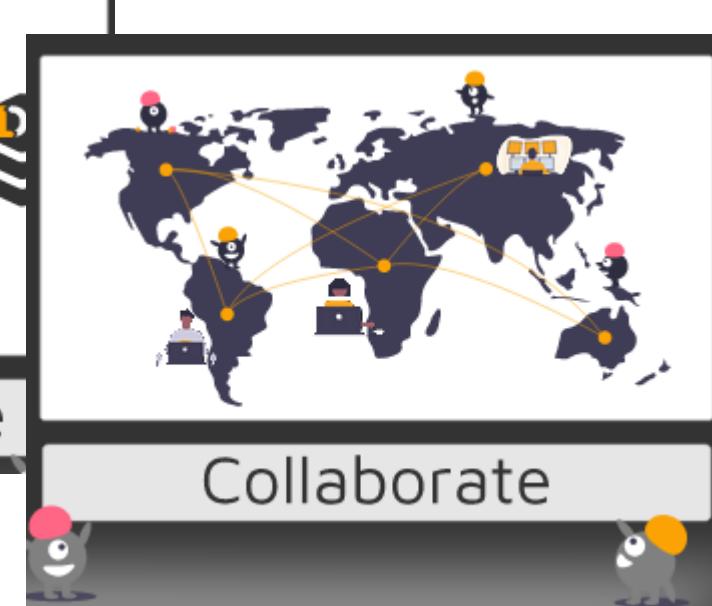
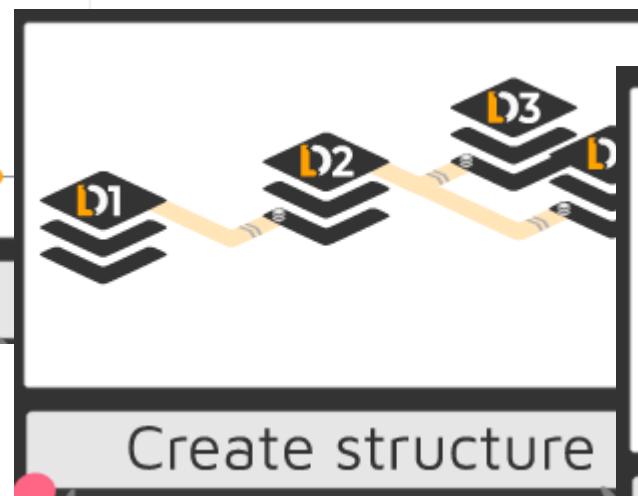
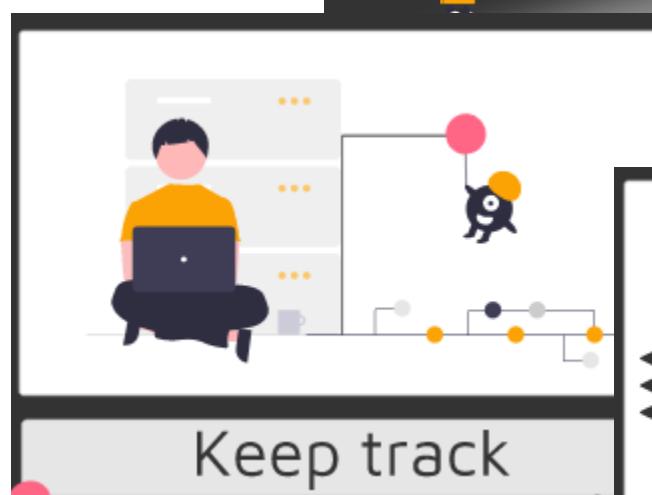
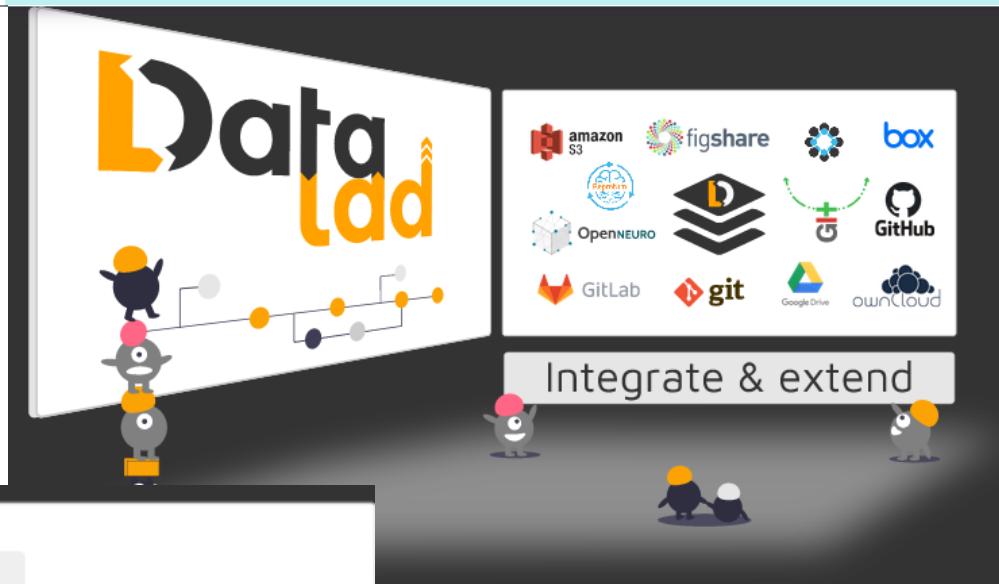
local
metadata

gathering the metadata





The mighty DataLad





Satrajit Ghosh



Seb. Urchs



Jean-Baptiste Poline



J. Armoza



Dorota Jarecka



Yaroslav Halchenko



Jeffrey Grethe



Sanu Abraham



Nazek Queder



Troy Sincomb



Karl Helmer



Tom Gillespie



David Kennedy



Theo VanErp



National Institute
of Mental Health



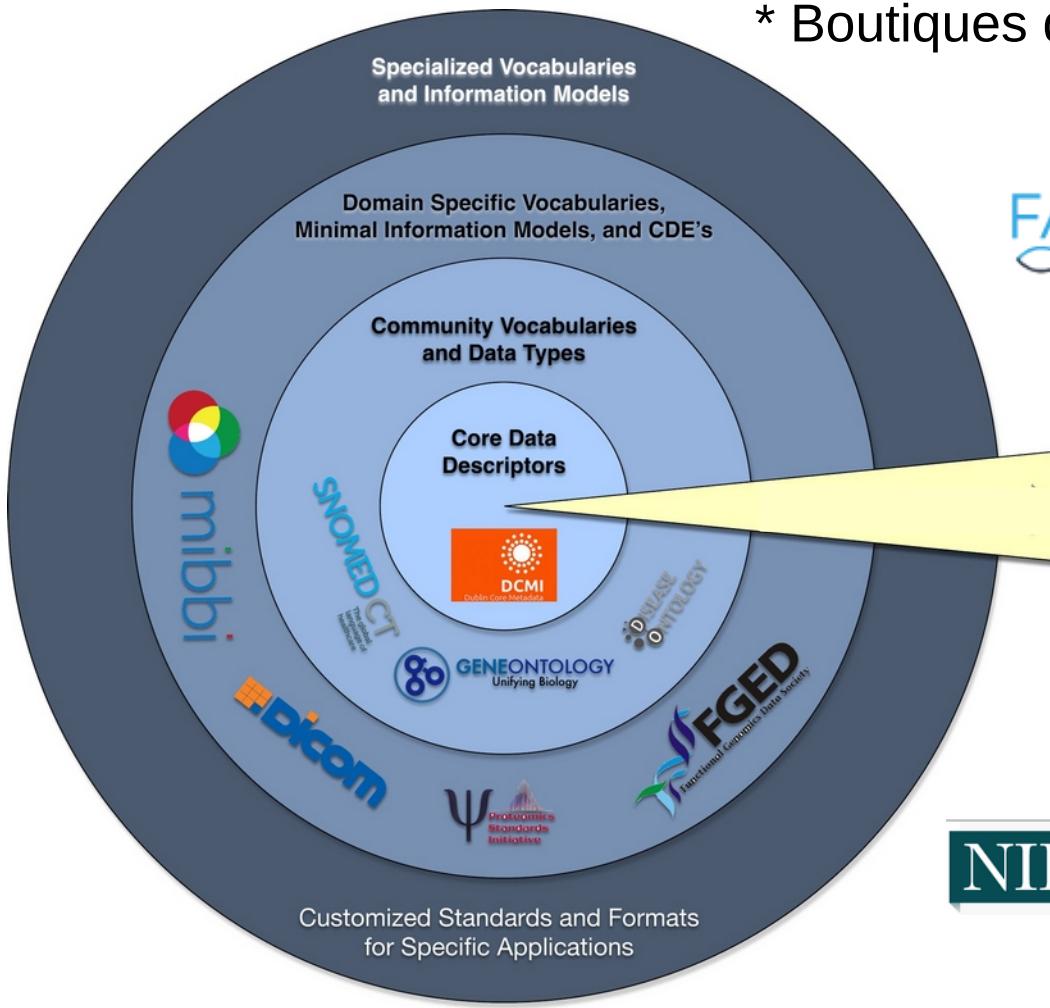
NIPYPE



1. Poor statistical procedures
2. Issues in data and software
3. A cultural issue: Publication practices and research incentives

Neuroinformatics platforms may partly bring answers to these three gaps

- * Reuse existing SDT (DATS, schema.org)
- * Develop metadata SDT (INCF, FDE)
- * Develop Landing page STD (Monii)
- * Boutiques descriptors STD



FAIRsharing.org
 standards, databases, policies

FORCE11
 The Future of Research Communications and e-Scholarship

in**ncf**



RDA
 RESEARCH DATA ALLIANCE



NIF

NEURODATA
 WITHOUT BORDERS

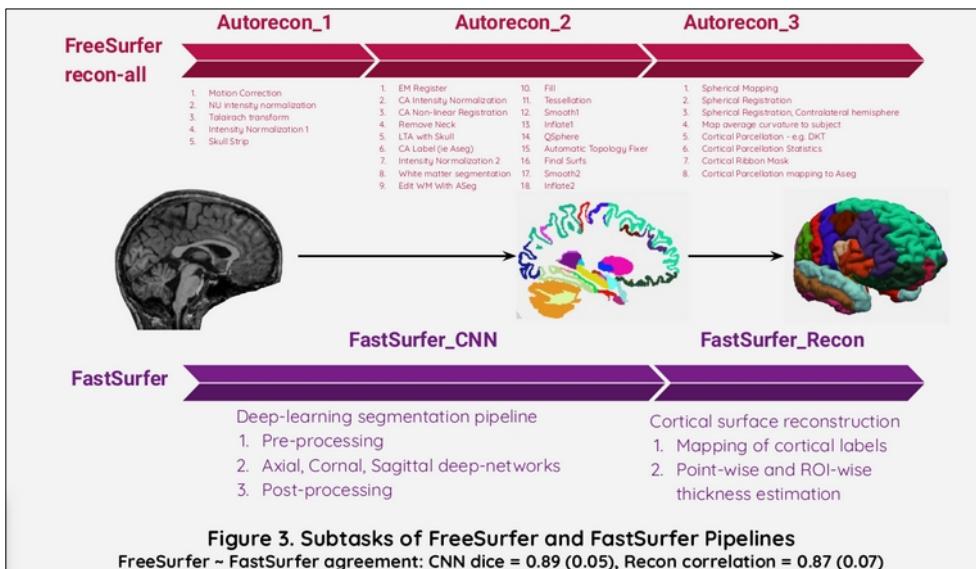
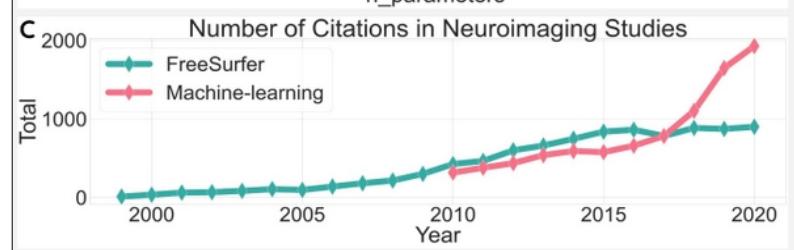
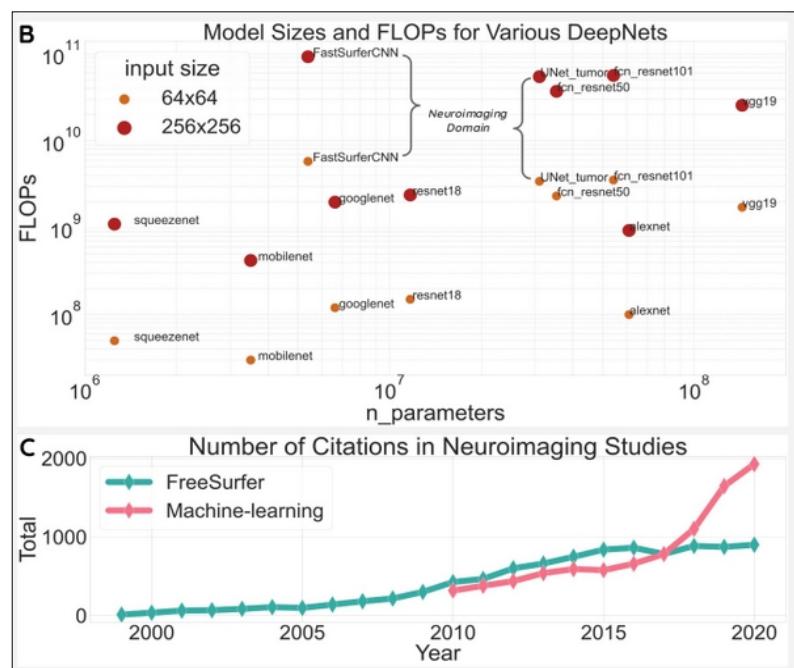
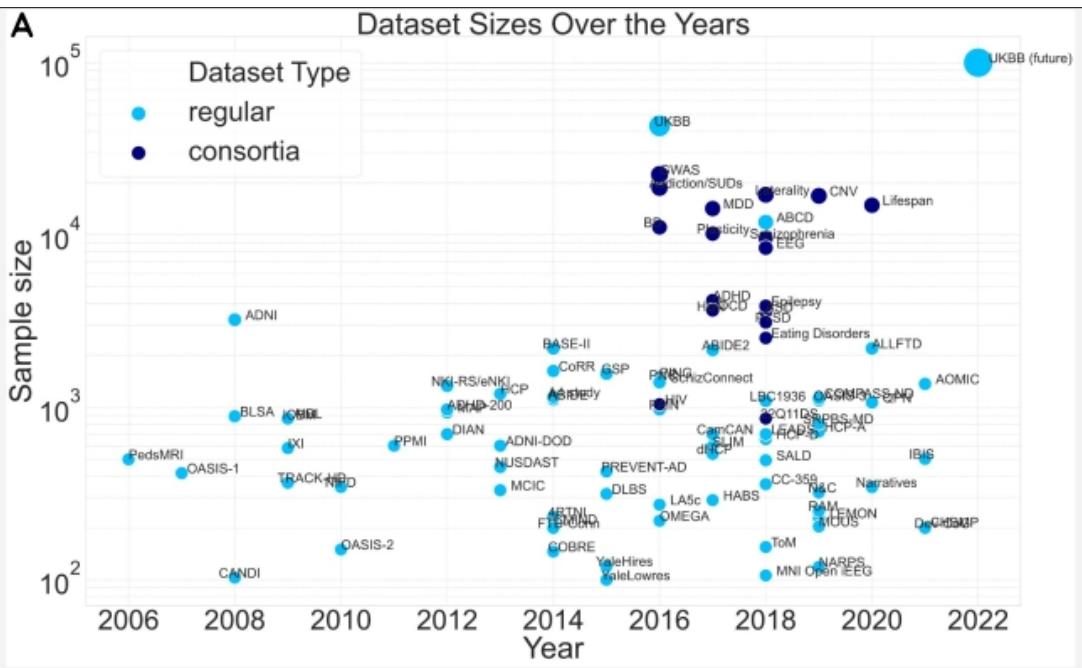
Part I: Reproducibility: background

Part II : Gap analysis

Part III : NeuroInformatics platforms

Part IV : Conclusion

Is this sustainable ?



Nikhil Bhagwat
et al., 2021

- Sustainability : We need to recalibrate our research objectives at individual and infrastructural levels with computing “cost” in mind
- Promote journals not handled by large companies when you can
 - When asked for review or editing
 - When publishing if possible
- Promote / use / contribute to tools, standards and best practices that are
 - Community developed and governed
 - Helping the larger vision of data and code sharing
- Choose your research project wisely
 - Will it be in extreme competition ?
 - Will it help move the scientific community in a better place ?
- Is it my job ?

Thank you

- Lab@McGill: <https://neurodatascience.github.io/>
- **McGill** colleagues: S. Brown, T. Glatard, G. Kiar, A. Evans, C. Greenwood, A. DeGuise and others
- **ReproNim** colleagues: D. Kennedy, D. Keator, S. Ghosh, M. Martone, J. Grethe, M. Hanke, Y. Halchenko
- **Berkeley** colleagues: S. Van der Walt, M. Brett, J. Millman, Dan Lurie, M. D'Esposito, et al
- **Pasteur** colleagues: G. Dumas, R. Toro, T. Bourgeron, and others
- **Paris** colleagues: B. Thirion, G. Varoquaux, V. Frouin, et al
- **Funders:** McGill HBHL, HBHL NeuroHub, NIH, NIMH

neurodatascience.github.io



Jean-Baptiste Poline
Principal Investigator



Kendra Oudyk
PhD student



Ting Zhang
PhD student



Gaël Varoquaux
Visiting Professor (Research Director,
Inria)



Elizabeth DuPre
PhD student



Jacob Sanz-Robinson
MSc student



Peer Herholz
Postdoctoral researcher



Qing Wang (Vincent)
Postdoctoral researcher



Jérôme Dockès
Postdoctoral researcher



Alexandre Hutton
Research Assistant



Adam Trefonides
Research Software



**Kate
Kim**

Data Analysis projects: Methods development Infrastructures and Tools

- QW
- JD
- PH
- KK

- KO
- ED
- TZ
- GV

- AH
- JSR
- AT
- JB



Fondation
Brain Canada
Foundation



HEALTHY BRAINS
HEALTHY LIVES

