# Perspectives from a domain specific data repository:
## *The National Sleep Research Resource*

Susan Redline, M.D.,M.P.H.

Professor of Sleep Medicine

Harvard Medical School

Brigham and Women's Hospital

Beth Israel Deaconess Medical Center
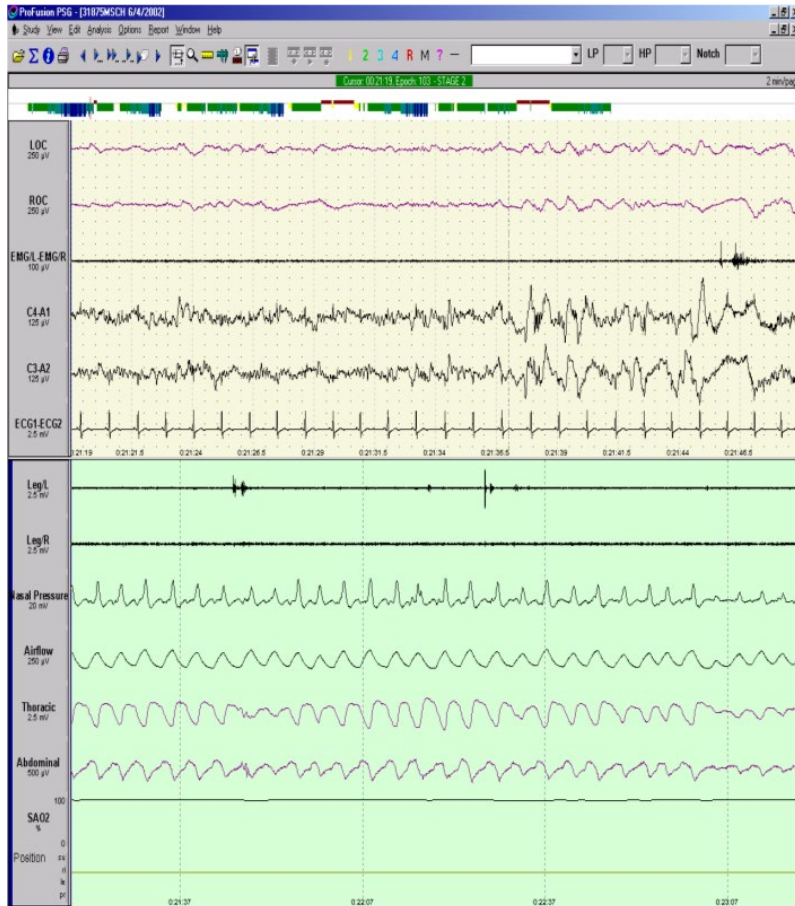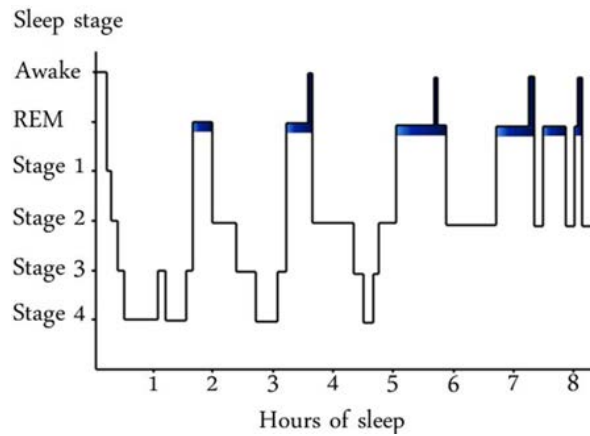
sredline@bwh.harvard.edu

# Outline

- Data specific domain: potential and challenges
  - Sleep and Circadian Data

- Goals and organization of the National Sleep Research Resource

- The user community
  - Defining their needs
  - Measuring Impact

- Challenges

# Reservoir of sleep data

- **2,800** accredited sleep labs in the U.S.

- **845,569 sleep studies** were performed in 2014
  - Increasing per year

- **~400 MB/study– 340 TB/yr**

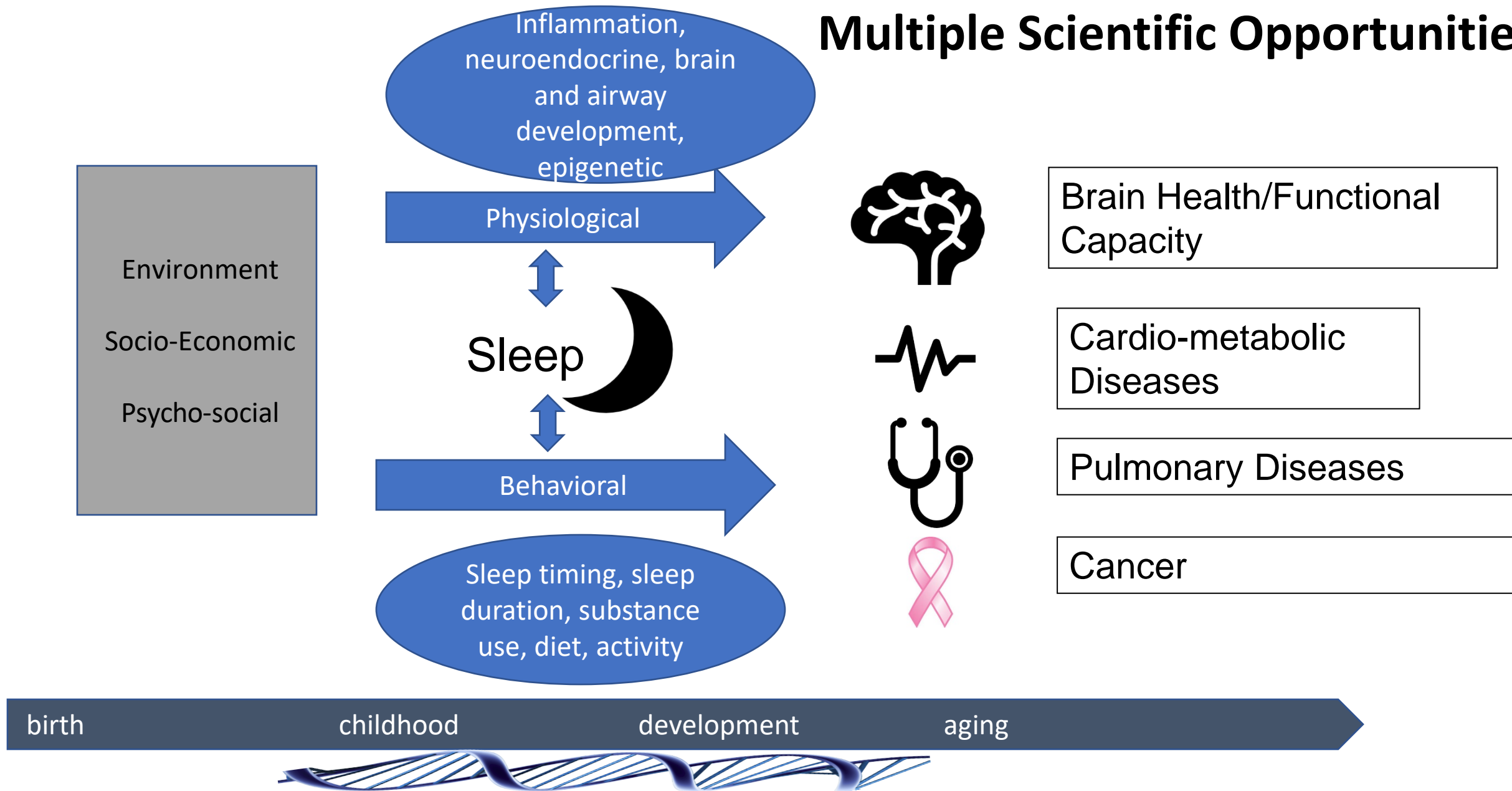- **NIH and industry-funded research sleep studies**
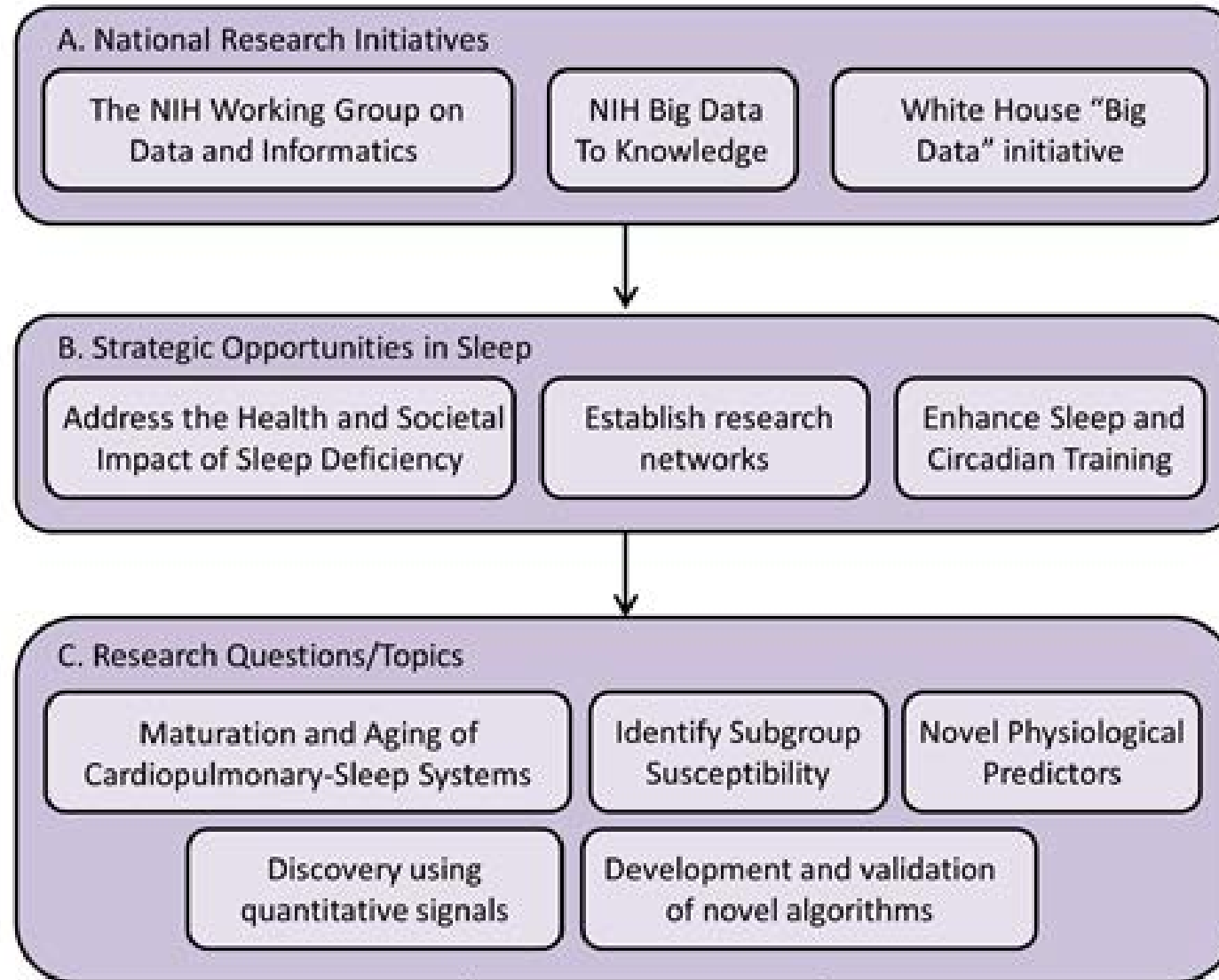
# Untapped Data Signatures

- Hours of physiological signals
- Cross-talk between physiological systems
- Temporal and dynamic features
- Days of physiological/behavioral signals
- Multiple ways to annotate

Multiple Scientific Opportunities

Inflammation, neuroendocrine, brain and airway development, epigenetic

Physiological

Sleep

Behavioral

Sleep timing, sleep duration, substance use, diet, activity

Environment

Socio-Economic

Psycho-social

Brain Health/Functional Capacity

Cardio-metabolic Diseases

Pulmonary Diseases

Cancer

birth          childhood          development          aging

**Critical timepoints; Cumulative risk models**

# Sleep / Circadian Big Data Opportunities



A. National Research Initiatives

- The NIH Working Group on Data and Informatics
- NIH Big Data To Knowledge
- White House "Big Data" initiative

B. Strategic Opportunities in Sleep

- Address the Health and Societal Impact of Sleep Deficiency
- Establish research networks
- Enhance Sleep and Circadian Training

C. Research Questions/Topics

- Maturation and Aging of Cardiopulmonary-Sleep Systems
- Identify Subgroup Susceptibility
- Novel Physiological Predictors
- Discovery using quantitative signals
- Development and validation of novel algorithms

# Challenges in sleep data analysis



- **Data sets heterogeneous, some poorly annotated and difficult to harmonize**
  - Different collection protocols, lack of standardized montages, variable scoring
  - Lack of accepted sleep ontologies/variable vocabularies
    - Summary data and raw signals



- **Limited data types**
  - Focus on summary data
    - Untapped potential of advanced signal processing/machine learning

- **Few "open" sources of well-defined signals, linked covariates, and analysis tools**

# Gaps in data access and appropriate tools

- **Many web portals have a..**
  - Limited ability to query and visualize data
  - Limited ability to directly access data
  - Limited ability to access tools for visualizing and processing data
- **Large data analytics**
  - High data storage/egress costs
- **Access/download procedures**
  - Concerns over privacy/security
- **"Sandboxes" needed for**
  - Collaboration, promote documentation (transparency/reproducibility)
- **Barriers to users unfamiliar with dataset or dependencies on others**

# National Sleep Research Resource: sleepdata.org (2014-)

**Provide users web-based tools to assist with preliminary exploration of data within and across data sets and identify subsets of data most useful using clearly mapped terms**

**Community resource to deposit and access "raw" or complex primary data (physiological signals), including processed physiological signals**

**Provide users access to a hub of tools for processing physiological signals as well as a resource to support communications among sleep researchers**

**Partner with and link to other resources, such as BioLINCC and dbGAP (BioData Catalyst)**

# National Sleep Research Resource

Free research data and tools.

## What interests you?

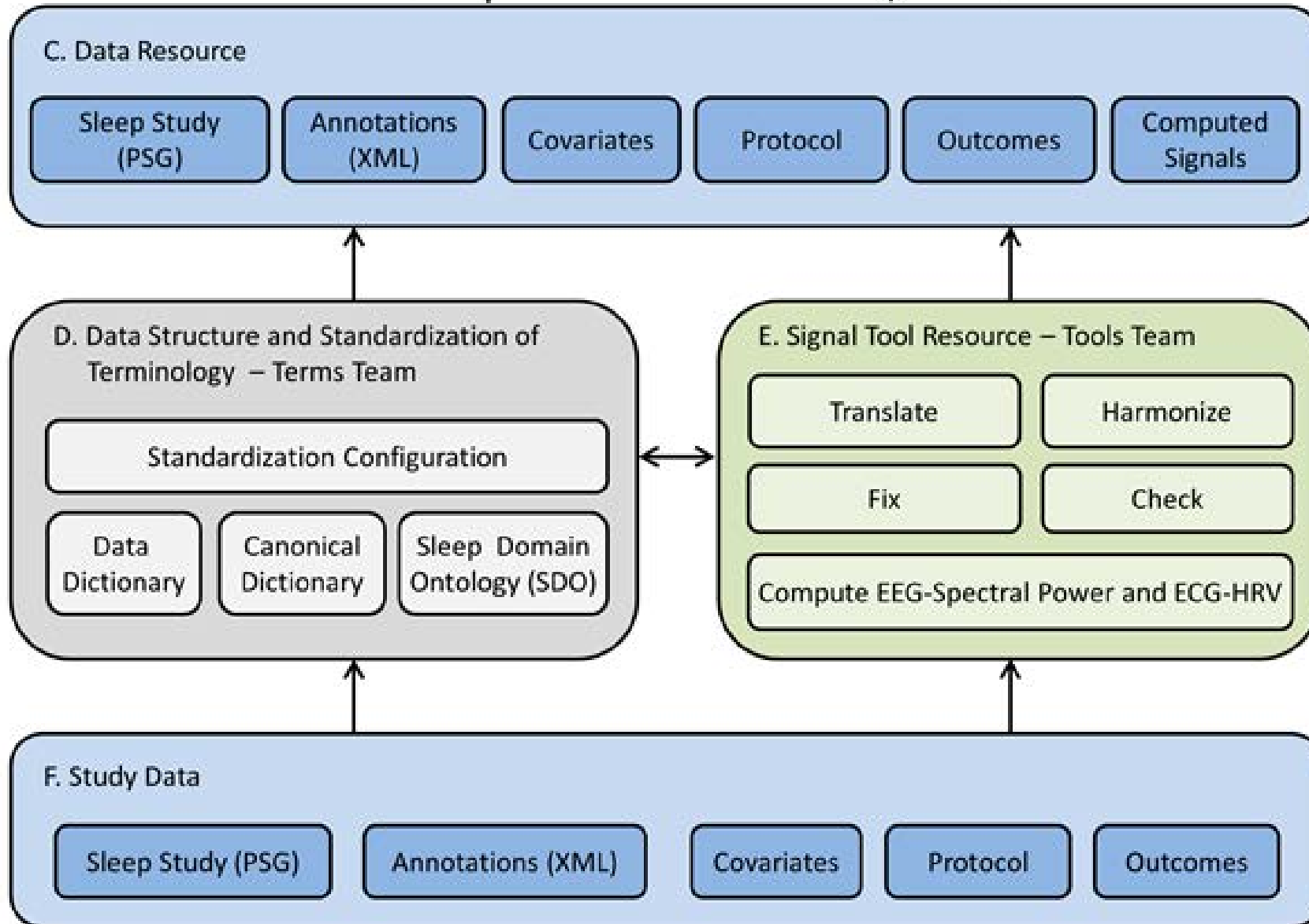| POLYSOMNOGRAPHY | ACTIGRAPHY | DATASETS | SHARING DATA |
|---|---|---|---|

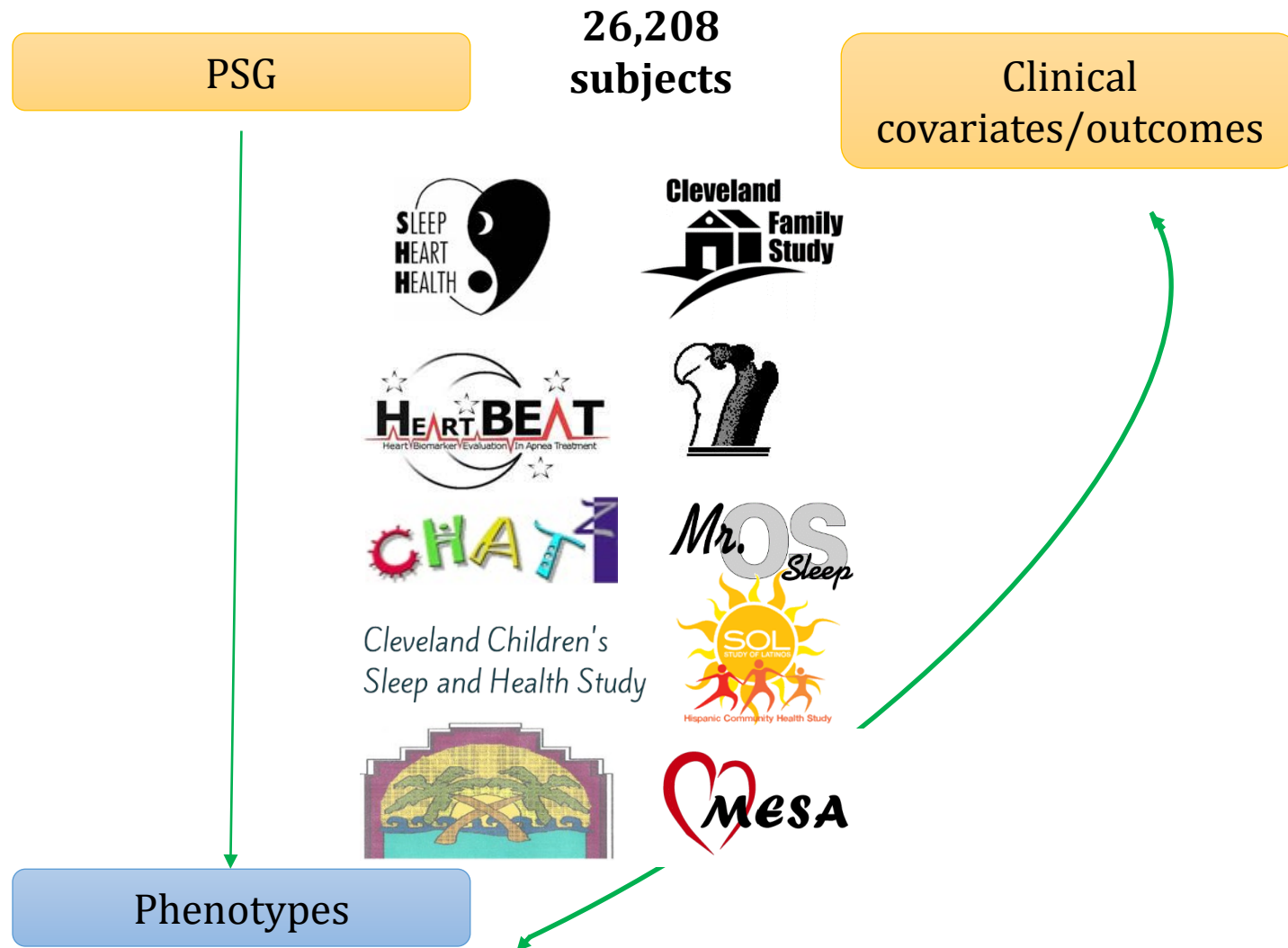Not sure? View our most popular datasets.



## Visualizing NSRR data

To paraphrase the adage, a picture is worth a thousand numbers. In order to investigate some basic properties of NSRR datasets, here we generate a number of whole-dataset visualizations. To make sense of these images, we'll employ a remarkably complex computational pattern recognition and dimension reduction framework, a.k.a. the human visual system. Keep reading ▸

# National Sleep Research Resource: Sleepdata.org

# Data Integration

# Available Data

**31,580** EDFs from 27,151 subjects

**19,235** PSGs with EEG or ECG spectral analysis results

**4,064** actigraphy files

**5,324** terms annotated to structured definitions

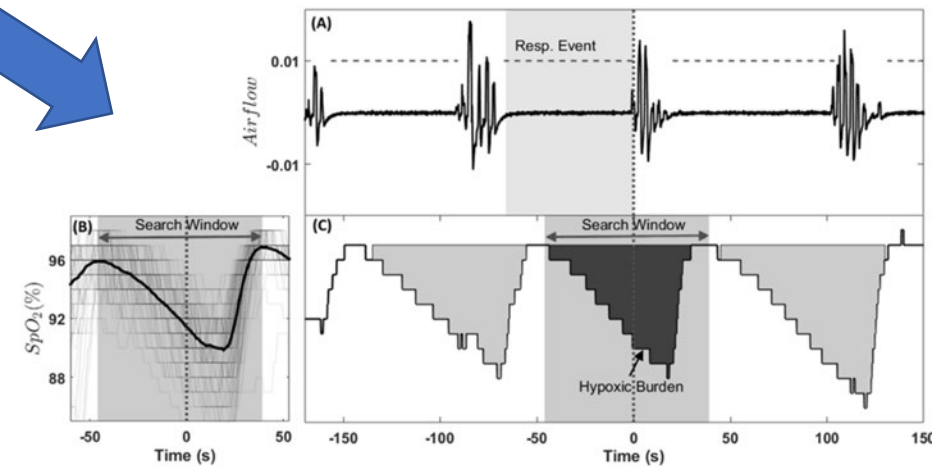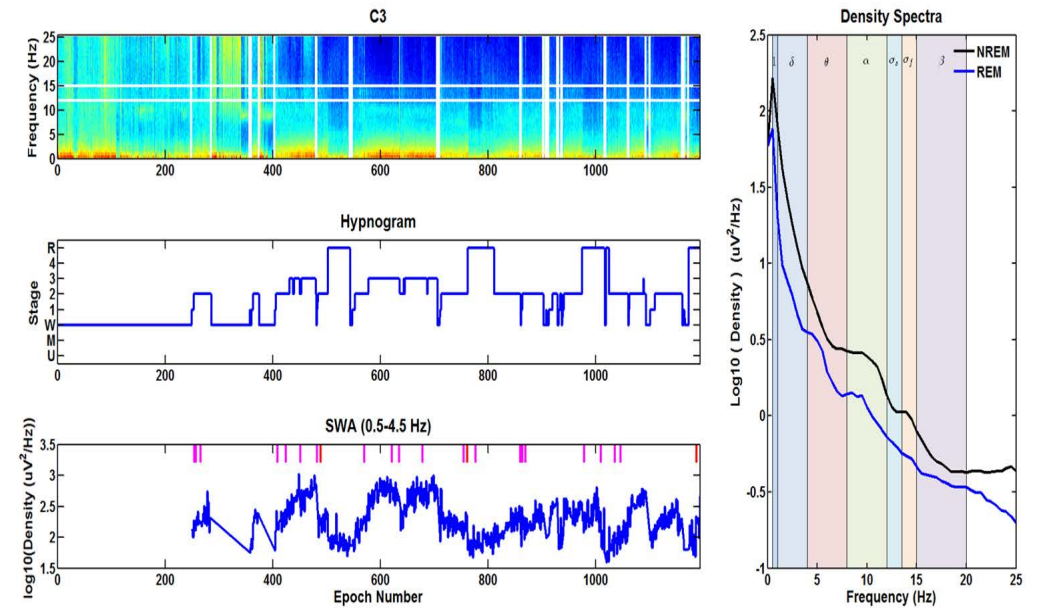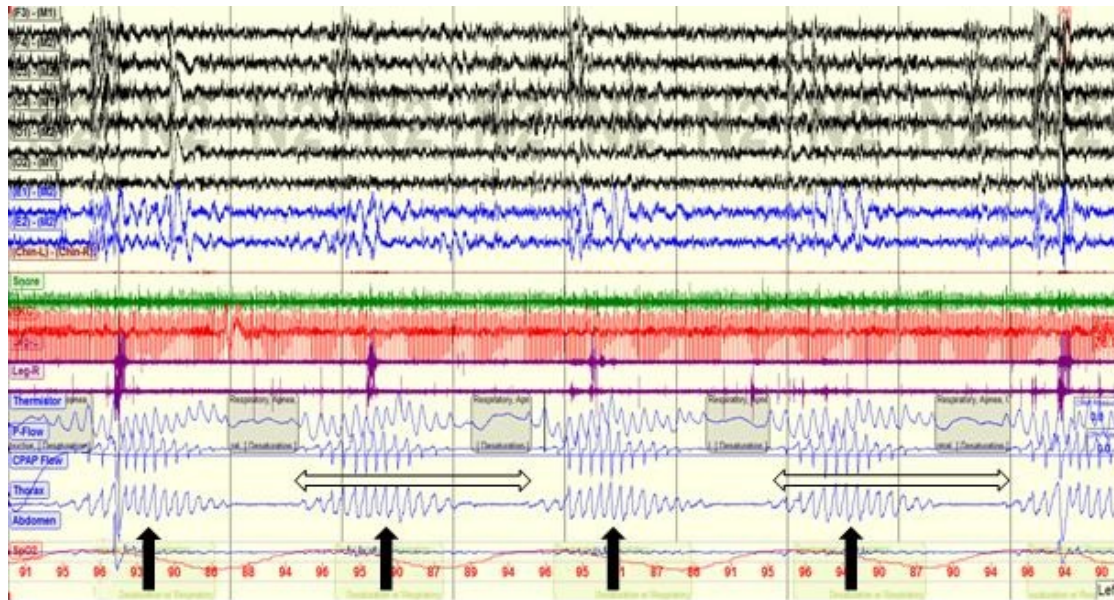**4,681** with provenance attributes

# Quantitative Signal Analysis



Figure 1: Example of hypoxic burden calculation for an individual respiratory event (Resp. Event), Panel A

*Search across 1,000s of variables...*



Variables logically grouped by type

# National Sleep Research Resource

Free research data and tools.

NSRR   About   Datasets   Tools   Forum   Blog                🔍 Search   Sign in

## What interests you?

| POLYSOMNOGRAPHY | ACTIGRAPHY | DATASETS | SHARING DATA |

Not sure? View our most popular datasets.

## Tools for the analysis of sleep data

The NSRR is revamping its Tools pages and needs your help! Have you developed a tool for the analysis of sleep data that you'd like others to know about and use? Do you have some tricks and tips for using existing packages that you'd like to share? What about a write-up listing your favorite tools, explaining how you use them and what's good about them? Or perhaps you'd like to share some data analytic problems that aren't met by existing tools? If so, we'd love for you to submit a gues. Keep reading ▸

By shaunpurcell on *September 13, 2019* in Tools                                💬 0

# *Easy, but not uncontrolled, access to data…*



Interface for DAUA and IRB approval required for data access

# Share Your Data on the NSRR

The National Sleep Research Resource (NSRR) is an NHLBI-funded resource designed to host and share data from major sleep cohort studies and clinical trials. All shared study data must be de-identified using the HIPAA Safe Harbor method and must adhere to the data sharing language stated in the participant informed consent. Records and files from participants who did not consent to data sharing must be redacted before submitting to the NSRR.

The NSRR creates a unique space to share and link covariate data, complex physiological data, and quantitative signal (e.g. EEG, ECG) processing results. The NSRR team will guide you through the process of preparing and uploading your datasets to the NSRR.

Uploading data to the NSRR satisfies requirements of the NIH Data Sharing Policy. For future grants, please consider including data sharing language that mentions the NSRR.

---

**What you will do:**

- Compile documentation (e.g. manuals, questionnaires) about your data
- Prepare final datasets with data dictionaries and descriptions
- Remove all identifiers from dataset and raw data files
- Upload files to NSRR through Secure File Transfer Protocol (SFTP)

**What we will do:**

- Assist you during each step of the submission process
- Review uploaded data to ensure all identifiers have been removed
- Establish an institutional data use agreement (if needed)
- Ensure that only the users you want to access your data receive access
- Create a repository for your dataset to organize documentation and data files

# Assessing Impact

- **User Base**
  - Register/Access data

- **Products**
  - Use/publish data
  - Contribute data
  - Discoveries/new tools
  - Support new grants
  - Training

- **Engagement**
  - Interactive user community- collaboration, blogs, etc

# Assessing Impact: Access Data

- **6,041 registered users**
  - **1793 approved DUAs**

- Over **13 million files** downloaded, over **321 TB** of data
  - **2 TB** data per week

- **Ease of access**
  - Time interval from access to approval
    - User-friendly on-line DUA

# Assessing Impact: Use Data/Publish Results

- Publications
  - Epidemiological associations
  - Discovery/replication
    - New signals/Associations
  - Machine learning
  - Algorithm development/validation
- Tracking difficult
  - DUA: Cite grant / resource
  - NEED: Datasets as "citable" object
  - Track "Impact Factor" of resource

# Assessing Impact: Grants

- Training grants (NIH, AHA, AASM), R21s, RO1s

# Assessing Impact: Training

- Multiple levels
  - High school– post graduate

  - For example,
    - \> 100 Georgia Tech students: capstone project
    - OSHU Data Wrangling courses/workbooks
    - Basis for Harvard ML course
    - Resource for a biostatistics book/course

# Assessing Impact: Contribute Data

- **New contributors**
  - Individuals
  - New cohorts: ~15 new cohorts identified
  - NIH (NIMH; E.S.P)

- Incentives for data sharing

- Reducing "friction"
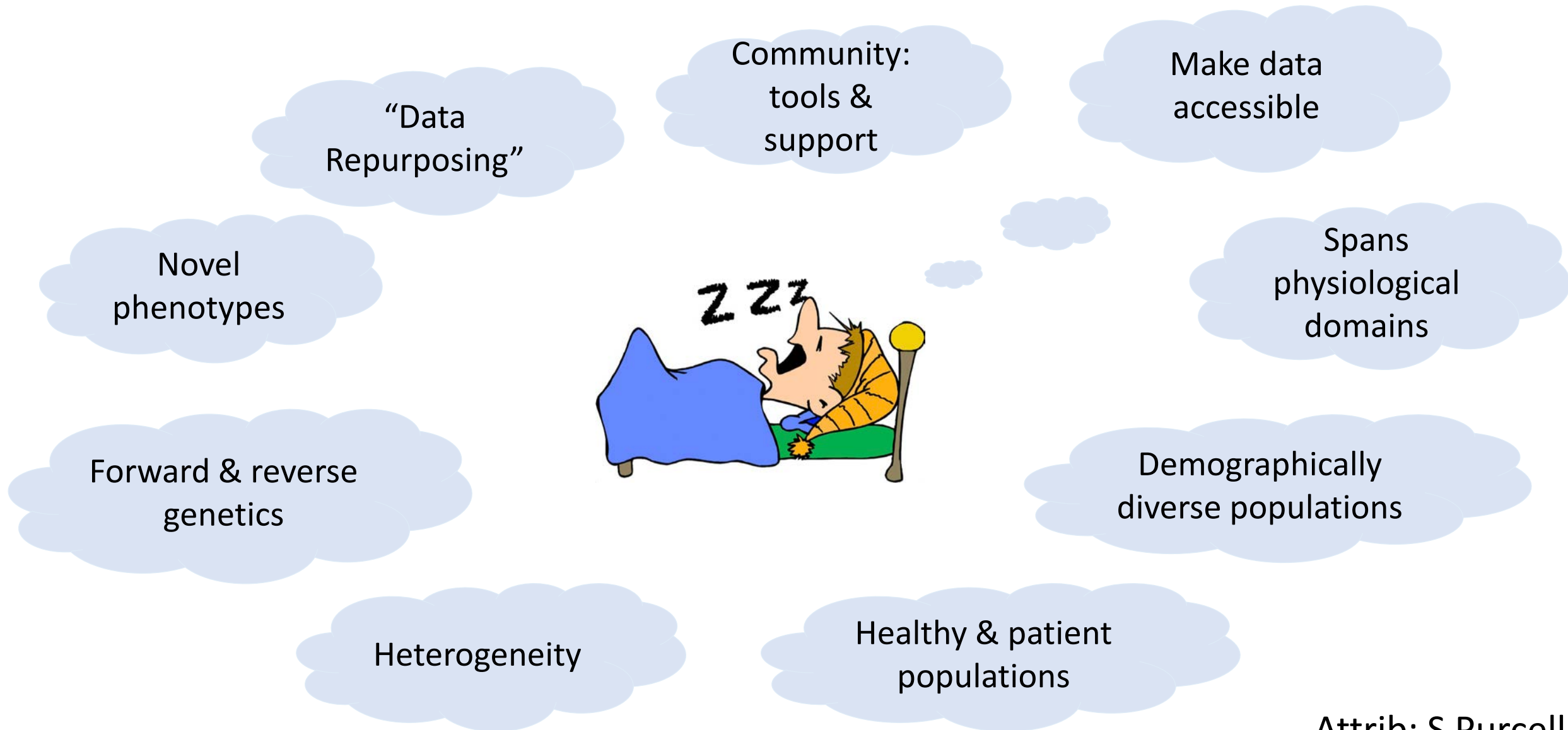  - Regulatory
  - Data structure/documentation

# Assessing Impact: Contribute CDEs

- > 5,000 variables mapped to standards
  - ICSD-3, NIH CDE
- >4000 variables mapped to provenance data
  - Bioportal CMS (wiki)

# Summary: Challenges/Needs

- Systematize citation process/Orchid registrations
  - Data resource impact factor?
- Link NIH funded grants for secondary data to sources?
- Trainee impact
  - Inventories of courses/books/trainee grants
- Data and tool contributions
  - Publish/highlight attributions

# Acknowledgements

- Brigham and Women's Hospital
  - Dennis Dean
  - Matthew Kim
  - Sara Mariani
  - Daniel Mobley
  - Remo Mueller
  - **Shaun Purcell**
  - Michael Rueschman
  - Ying Zhang


- NHLBI

- University of Kentucky
  - GQ Zhang
  - Satya Sahoo
  - Licong Cui


- Beth Israel Deaconess Medical Center
  - Ary Goldberger
  - Madalena Costa