

# Derived Study Elements

**Building an Analytic Dataset in N3C**

**N3C LOGIC LIAISONS:**

Johanna Loomba, ME

Andrea Zhou, ME

# Learning Objectives

What does it mean to build an analytic dataset for observational research?

What are important considerations I should take into account?

What are some available N3C tools that can help me?

# Background

# Observational Research: An N3C Study Workflow Perspective

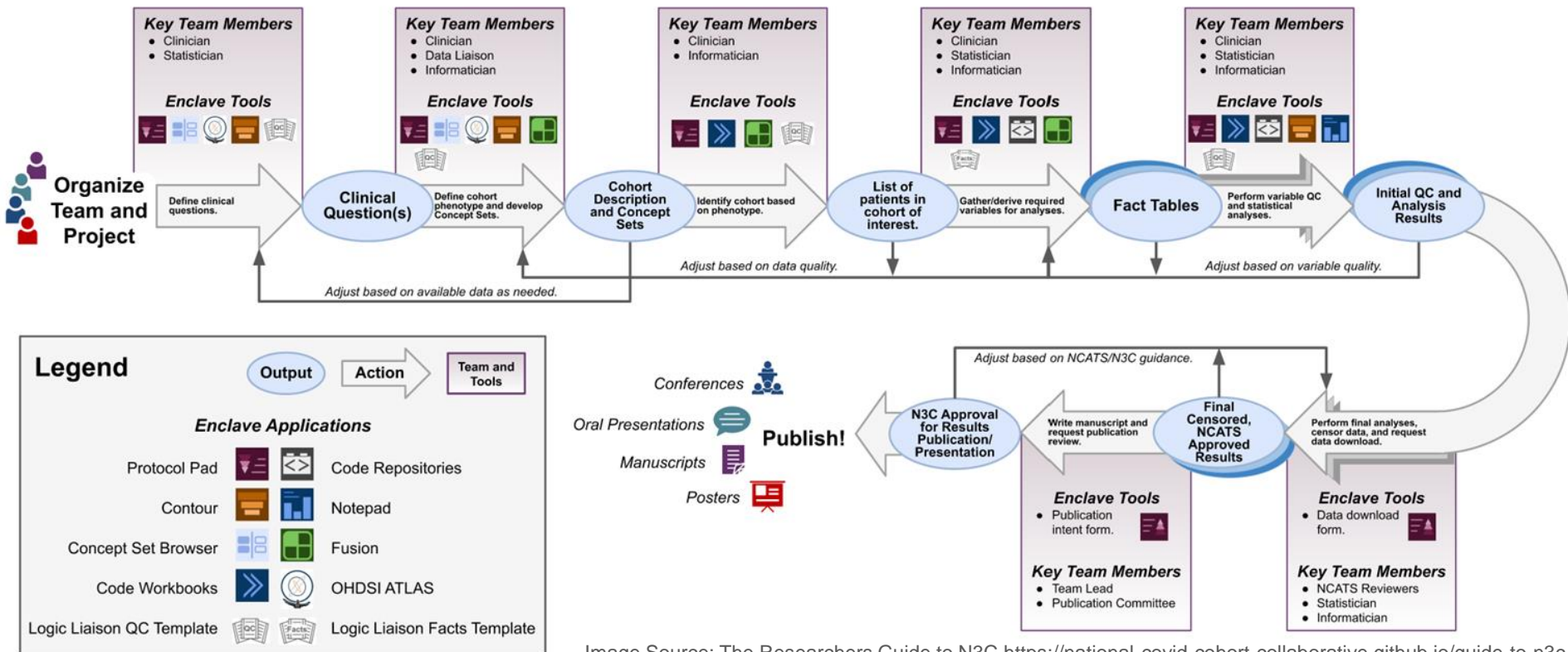
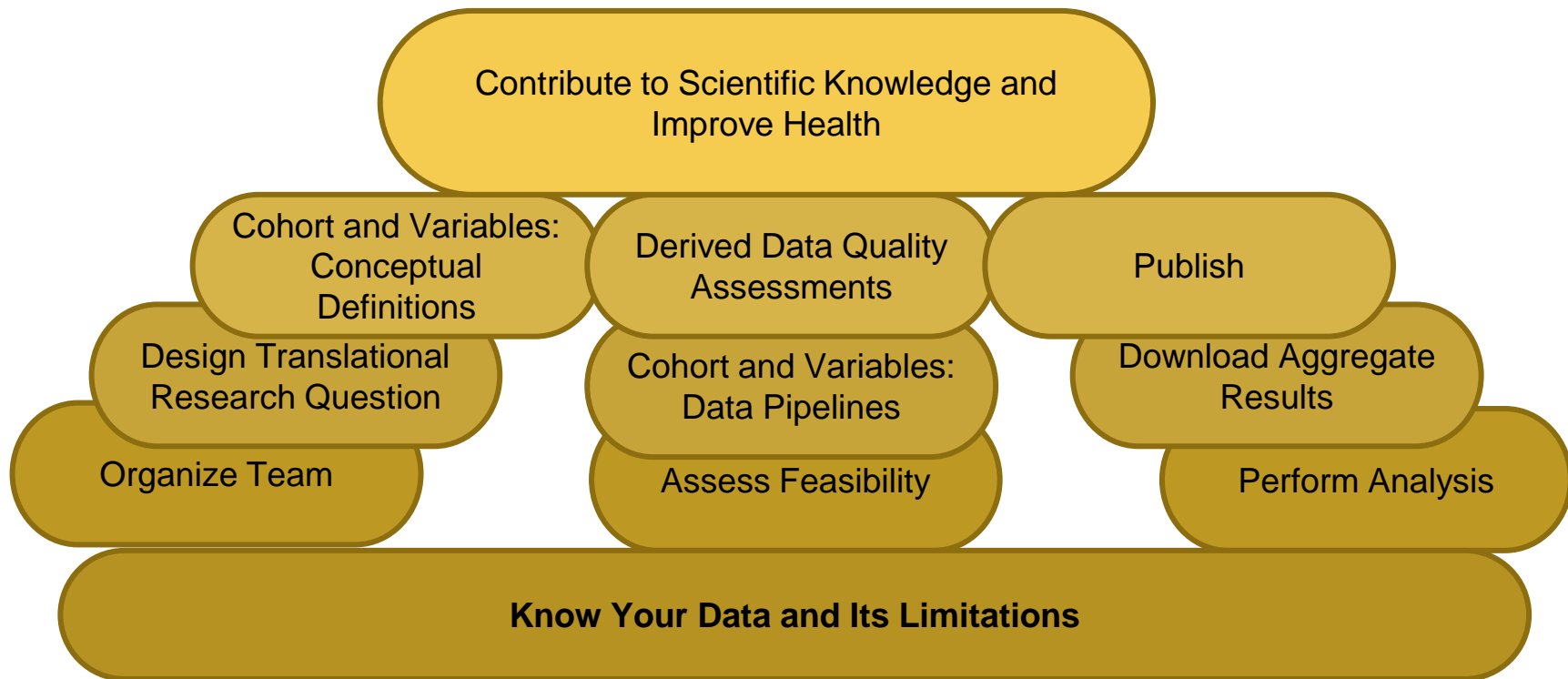


Image Source: The Researchers Guide to N3C <https://national-covid-cohort-collaborative.github.io/guide-to-n3c-v1/chapters/tools.html>

# Observational Research: Foundations Matter!



# Harmonized Clinical Records to Analytic Results

person

drug\_exposure

condition\_occurrence

measurement

device\_exposure

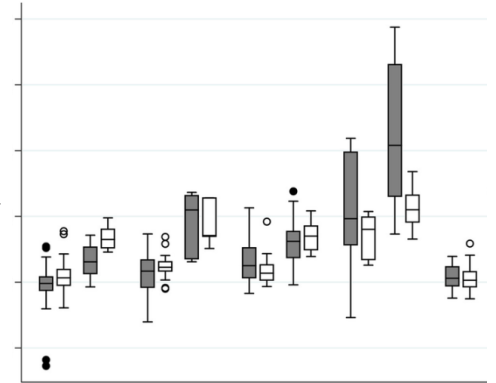
procedure\_occurrence

observation

visit\_occurrence

death

	v1	v2	...	vX
Patient 1	0	0	...	1
Patient 2	0	1	...	0
...	...	...	...	...
Patient n	1	0	...	1



# Creating Cohorts and Variables

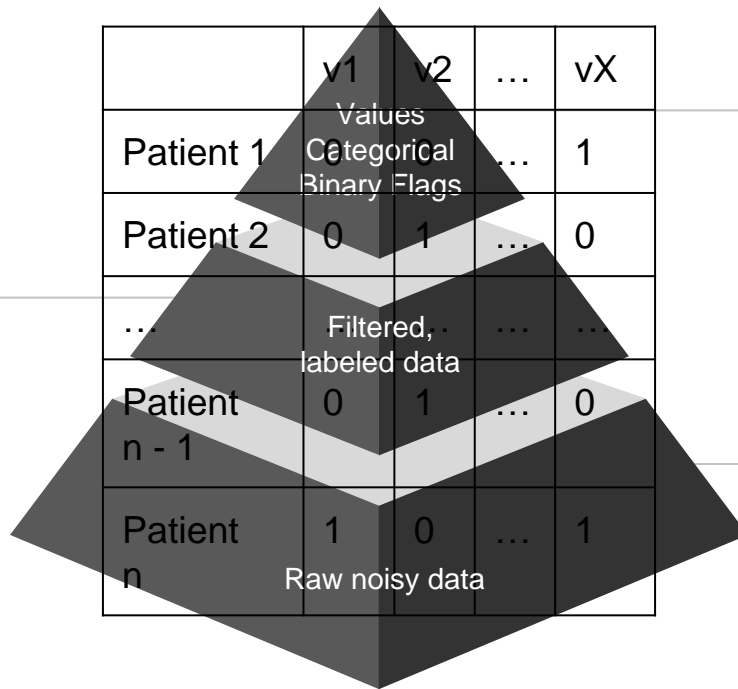
## Apply Concept Sets

Use to collapse raw data.

Leverage community built concept sets.

Evaluate specificity and sensitivity in light of you team's specific aims.

2



## Apply additional logic

Apply relevant logical operators and time relationships between concepts

3

Example: Patient severity at the time of COVID is defined by logical relationships between six concept sets and applies two temporal tests.

## Source Data

Find and evaluate relevant tables and fields in your source data.

1

Note limitations.

# Important Considerations



# Observational Health Data Limitations: Heterogeneity

Geographic Effects

Clinical Site Features

Electronic Health Record Systems

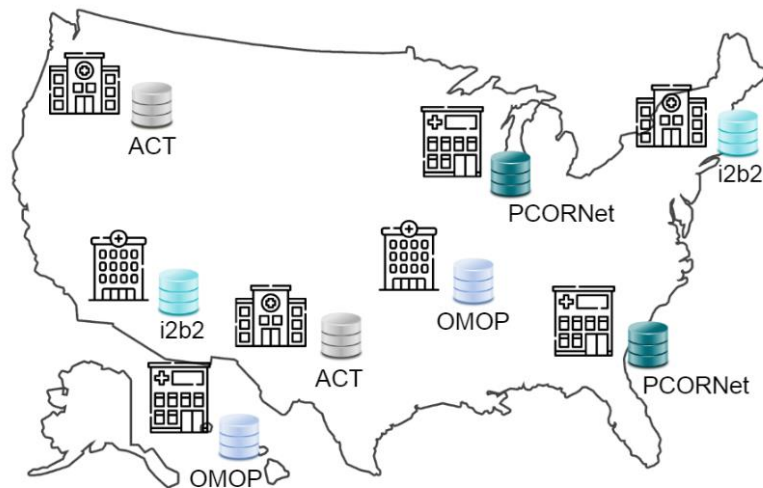
Clinical Data Models, Vocabularies, and Versions

Clinical Workflows, Care Practices, and Charting

Patient Access to Care, Trust, Transparency, Compliance, and Context

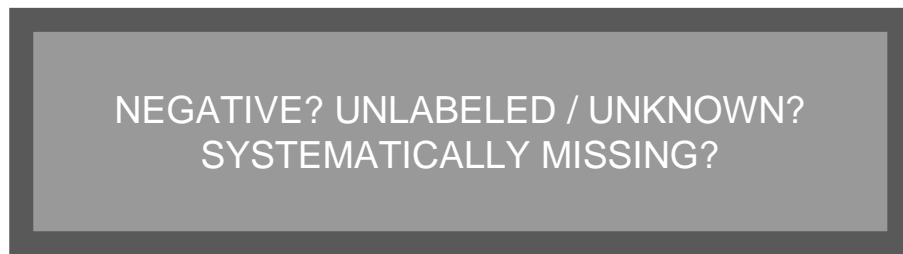
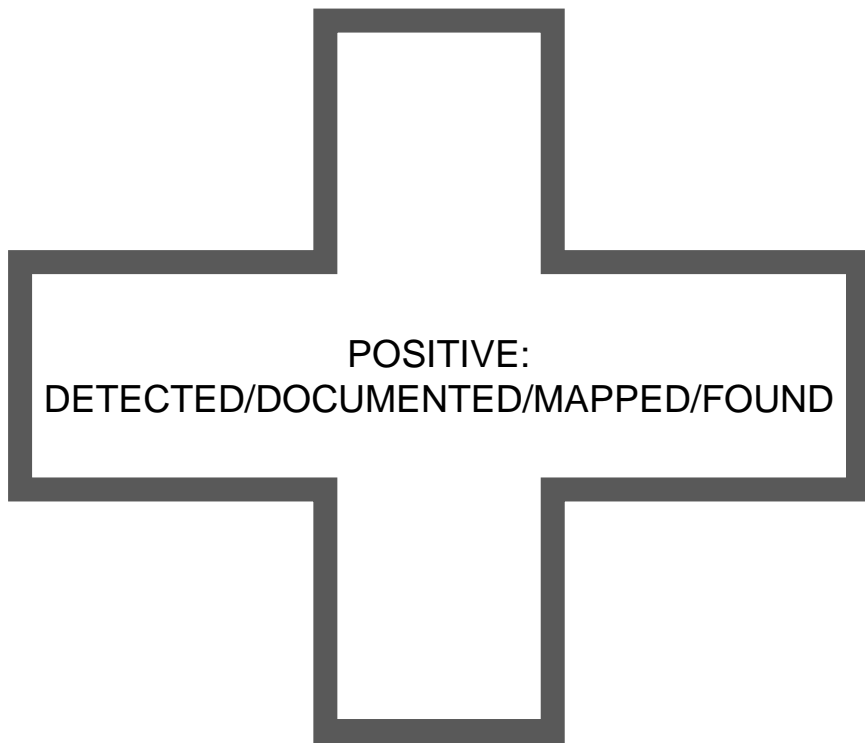
Privacy Preserving Record Linkage (PPRL): Data Augmentation Limited to certain patient populations and sites (i.e. Medicare, Death, Cross-Site Patient Linkage)

Etc.....

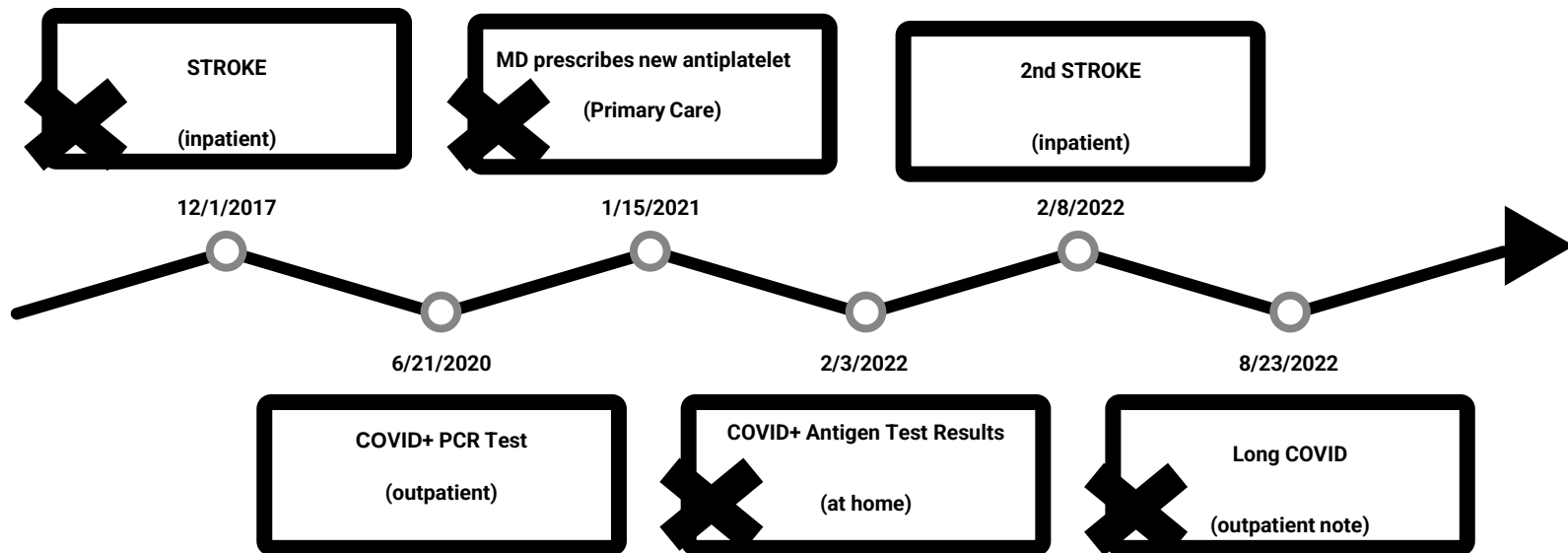


Relevant N3C Tables: Manifest, OMOP, PPRL Data

# Observational Data Limitations: Positive & Unlabeled Data



# Observational Health Data: Time-Related Limitations

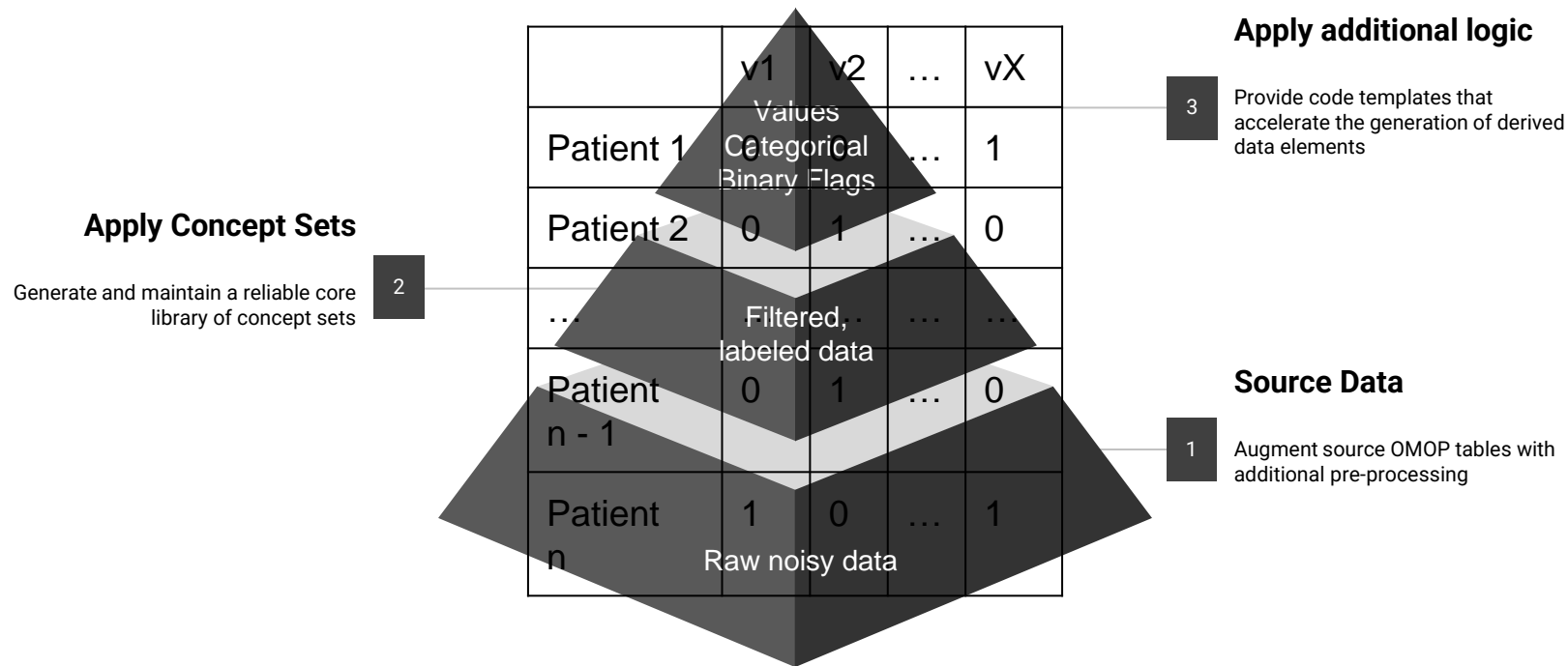


BEWARE of DATE SHIFTING, DATA LOSS, UNRELIABLE START/STOP DATES (chronic issues, patient lost to follow-up), EFFECT OF OBSERVATION PERIODS (data collection periods, data freshness), PHENOTYPE (patient selection effects)

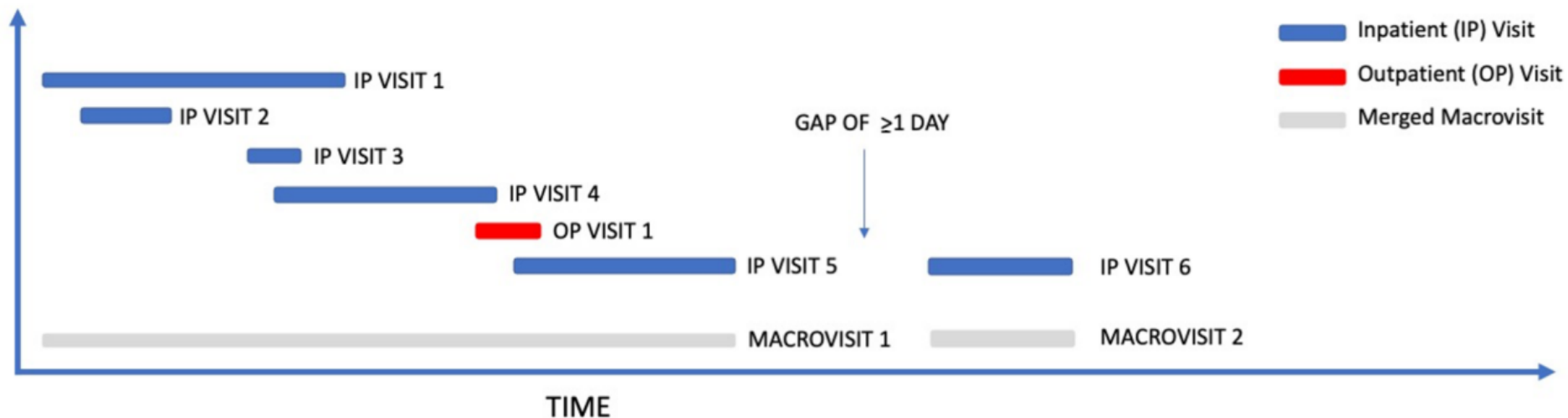
Relevant N3C Tables: Manifest, OMOP, PPRL Data

# Community Tools

# Eliminating Redundant Efforts: Data and Logic Liaison Tools



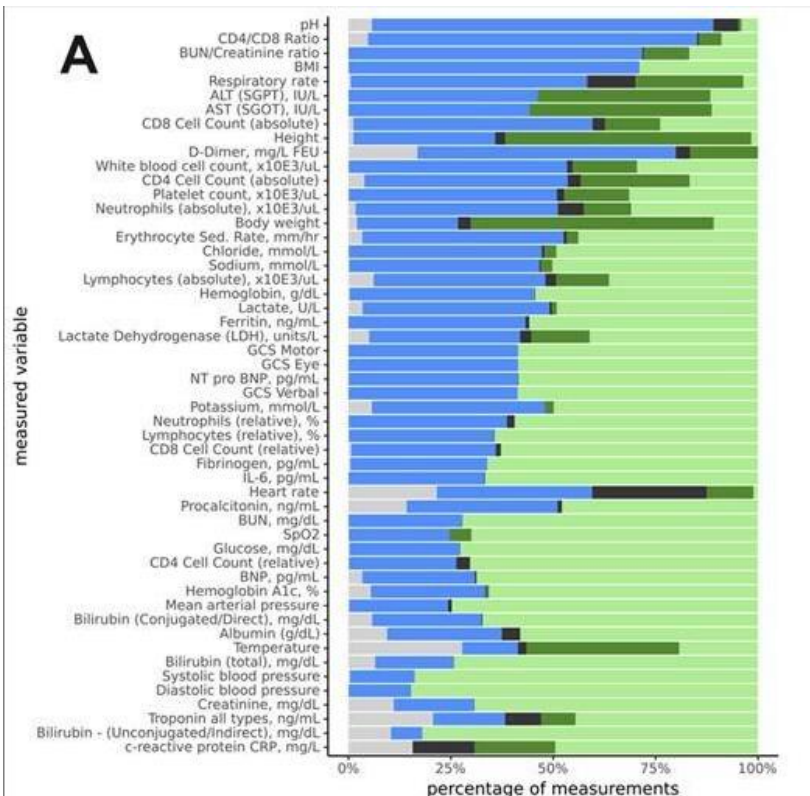
# Data Cleaning: N3C Macrovisits



Clinical encounter heterogeneity and methods for resolving in networked EHR data: A study from N3C and RECOVER programs. Peter Leese, Adit Anand, Andrew Girvin, Amin Manna, Saaya Patel, Yun Jae Yoo, Rachel Wong, Melissa Haendel, Christopher G Chute, Tellen Bennett, Janos Hajagos, Emily Pfaff, Richard MoffittmedRxiv 2022.10.14.22281106; doi: <https://doi.org/10.1101/2022.10.14.22281106>

N3C Microvisit to Macrovisit Map Table:  
macrovisit\_id, macrovisit\_start\_date, macrovisit\_end\_date

# Data Cleaning: N3C Measurement Harmonization



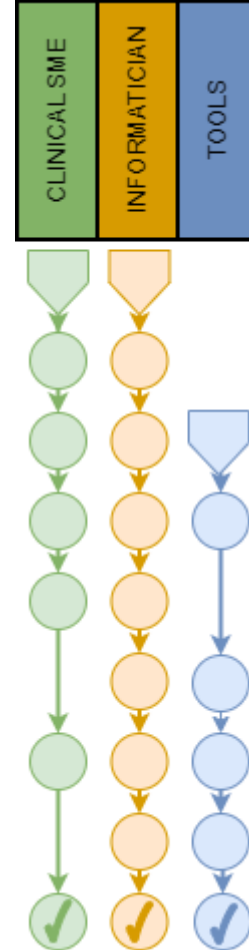
	Harmonized Unit Present in Source
	Unit Mapped and Value Converted
	Unit Absent but Inferred and Converted
	Unit Missing (Unable to infer)
	Unit Mapped to Null (Nonsensical Value)

Bradwell KR, Wooldridge JT, Amor B, et al. Harmonizing units and values of quantitative data elements in a very large nationally pooled electronic health record (EHR) dataset. *J Am Med Inform Assoc.* 2022;29(7):1172-1182. doi:10.1093/jamia/ocac054

**N3C Measurement Dataset Fields:**  
 harmonized\_unit\_concept\_id, harmonized\_value\_as\_number

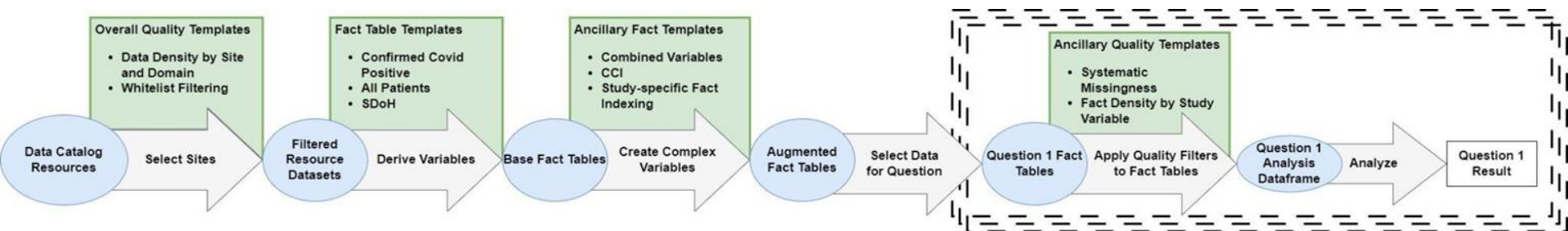
# N3C Recommended Concept Sets

1. **Assemble Team:** N3C Clinical Domain Team Leads (or delegates), N3C Data or Logic Liaison (informaticians and analysts)
2. **Articulate Goals:** Define the scope, intentions, and limitations
3. **Explore:** Use OHDSI tools to explore the CDM and identify candidate codes
4. **Compare:** Use both authoritative and community-built concept sets
5. **Review:** Present to a broader team of clinical experts
6. **Collapse:** Reduce the intensional concept set expression as parsimonious as possible, retaining all the approved concepts collected in prior steps.
7. **Document:** Intention, Limitations, Provenance, and expert Reviews
8. **Present:** Present for final vetting at the Data Liaison informaticists' meeting.
9. **Publish:** Marked as N3C Recommended and published to Zenodo



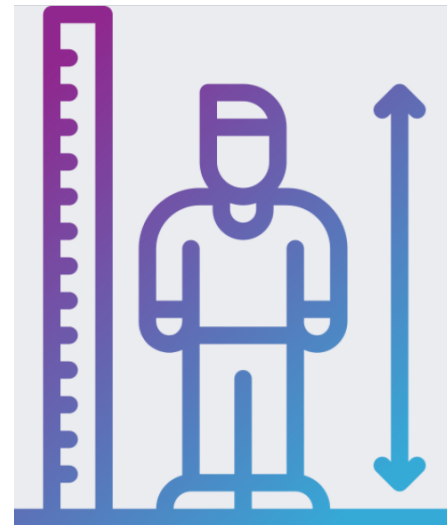


# Eliminating Redundant Work: Logic Liaison Templates

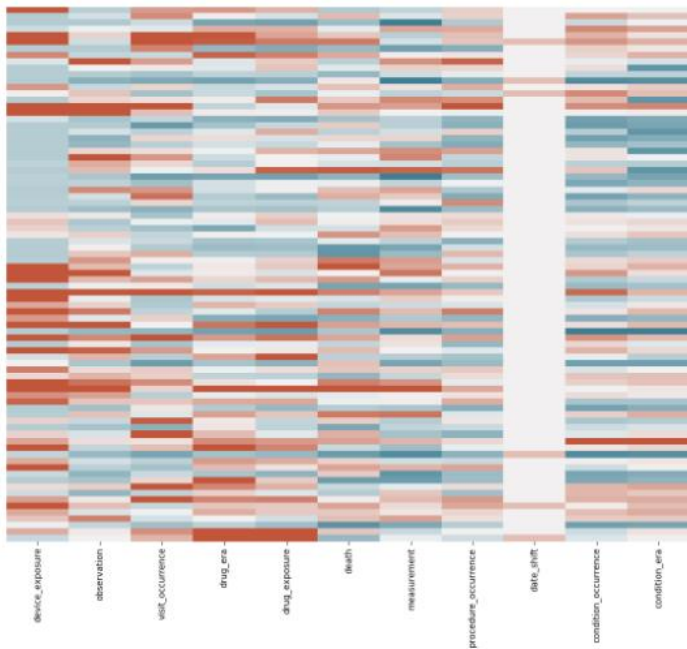


# Logic Liaison Tables: Minimal Data Cleaning/Imputation

- BMI/Obesity:
  - Unreasonable height/weight/BMI thresholds are applied and can be configured by the user
  - BMI computed when not reported and obesity imputed using  $BMI > 30$
- Unreasonable date detection:
  - DOB
  - Visits
  - Death



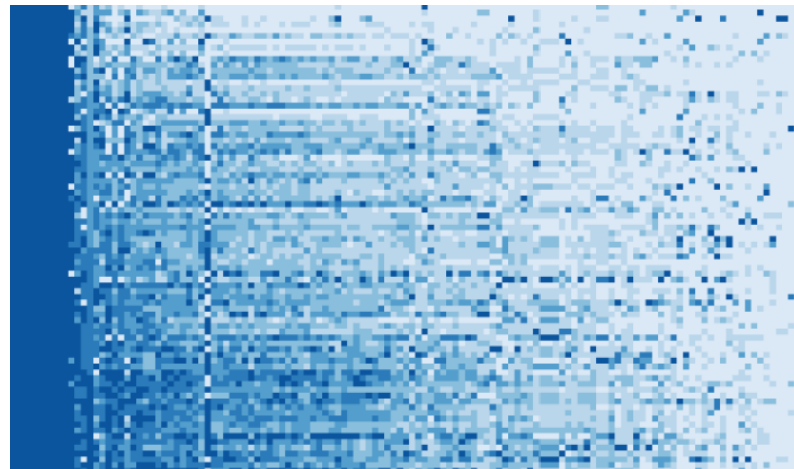
# Logic Liaison Quality Templates: Domain Density and Data Partner “Whitelist Filtering”



- How rich is each site's data by OMOP Domain?
- Which sites meet customizable data quality standards

# Logic Liaison Quality Templates: Fact Density and Systemic Missingness

Quality must be assessed again  
AFTER creating your derived variables  
and PRIOR to analysis



MALIGNANTCANCER\_post\_covid\_indicator  
MALIGNANTCANCER\_before\_or\_day\_of\_covid\_indicator  
Long\_COVID\_diagnosis\_post\_covid\_indicator  
Long\_COVID\_clinic\_visit\_post\_covid\_indicator  
LL\_IMV\_during\_weak\_covid\_hospitalization\_indicator  
LL\_IMV\_during\_strong\_covid\_hospitalization\_indicator  
LL\_ECMO\_during\_weak\_covid\_hospitalization\_indicator  
LL\_ECMO\_during\_strong\_covid\_hospitalization\_indicator

N3C Logic Liaison Quality Templates: Fact Density by Site Visualization

# Key Takeaways



Work with a breadth of domain experts

Don't recreate the wheel

Interrogate your raw data and derived  
dataframes

Leverage metadata

Understand and disclose limitations

Carefully derive features

Test your assumptions

Patience and humility

