# OHDSI community efforts on COVID-19 disease natural history: Status update and look forward to 'life after COVID'

Patrick Ryan, PhD

Janssen Research and Development
Columbia University Irving Medical Center

on behalf of OHDSI community:
CHARYBDIS study leads: Anthony Sena, Kristin Kostka,
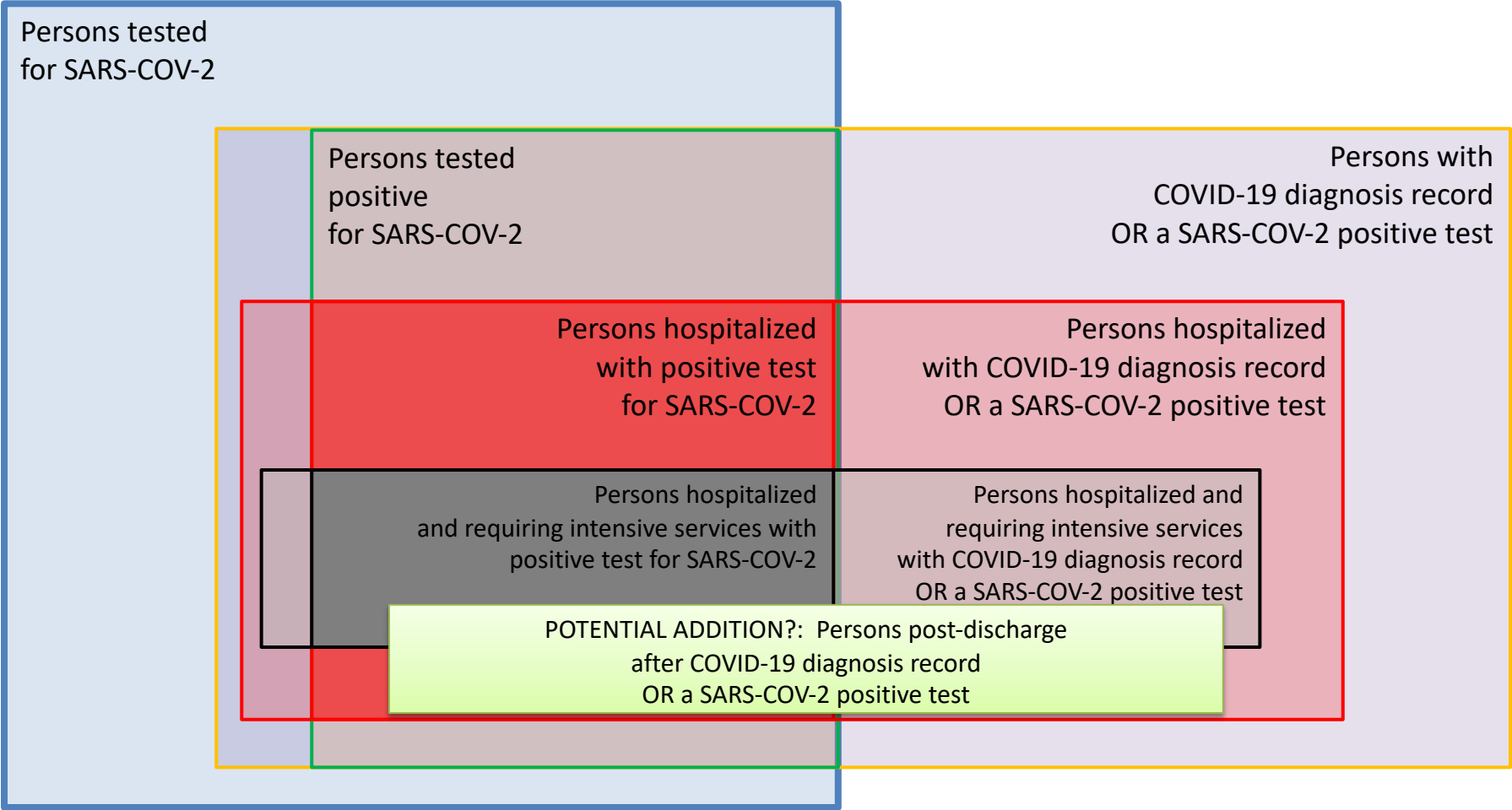Talita Duarte-Salles, Albert Prats-Uribe

# Open collaboration requires FULL transparency in every step of the research process

- Protocol and analysis source code freely available and directly downloadable: https://github.com/ohdsi-studies/Covid19CharacterizationCharybdis

- Phenotype definitions are both human-readable and computer-executable using ATLAS against any OMOP CDM:
https://atlas.ohdsi.org/

- All analysis results will be available for public exploration through interactive R shiny application:
http://data.ohdsi.org/Covid19CharacterizationCHARYBDIS/

- The study is a living evidence repository: any data partners can execute analysis and share aggregate results at any point, including updates as data accumulate
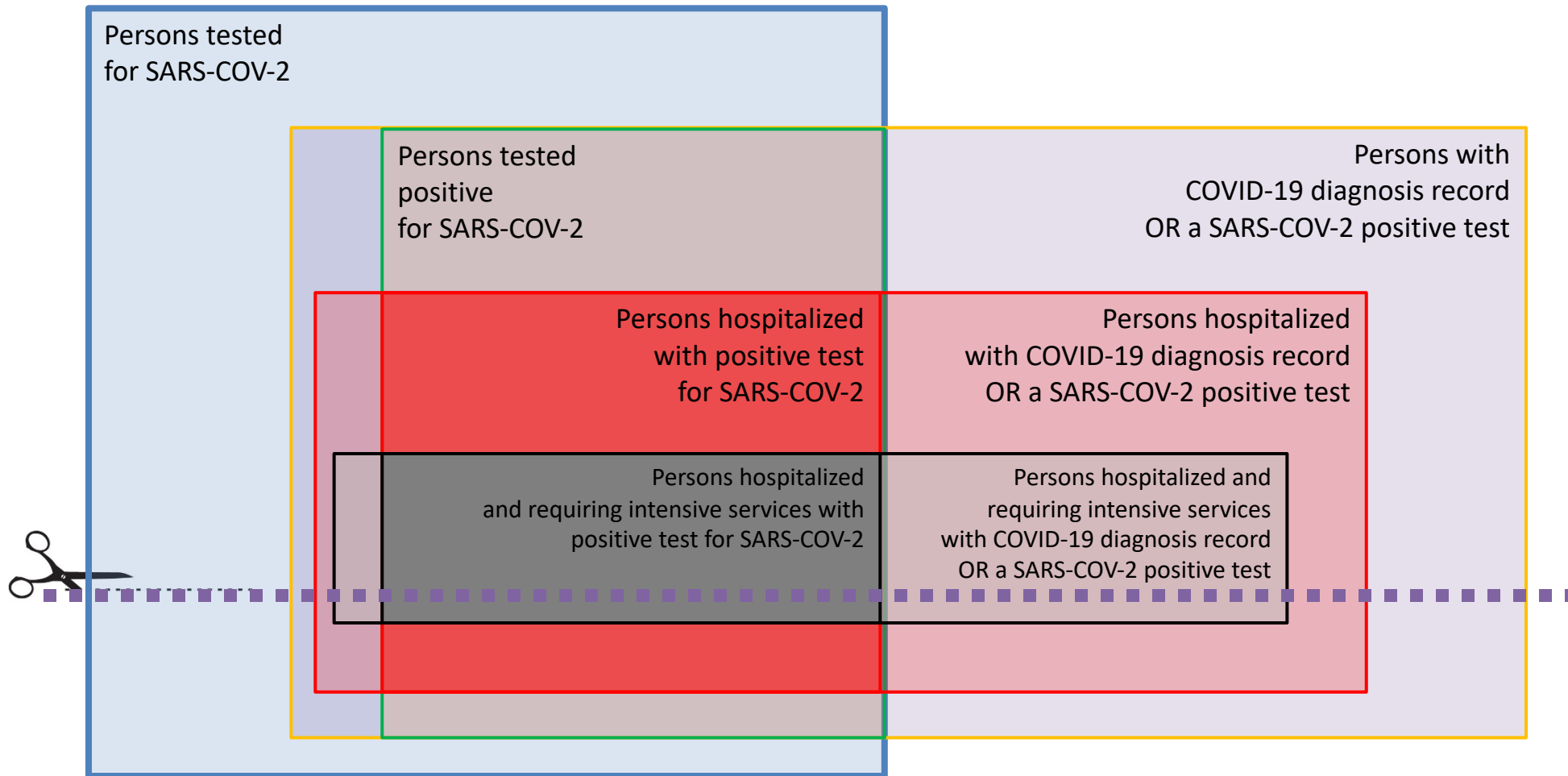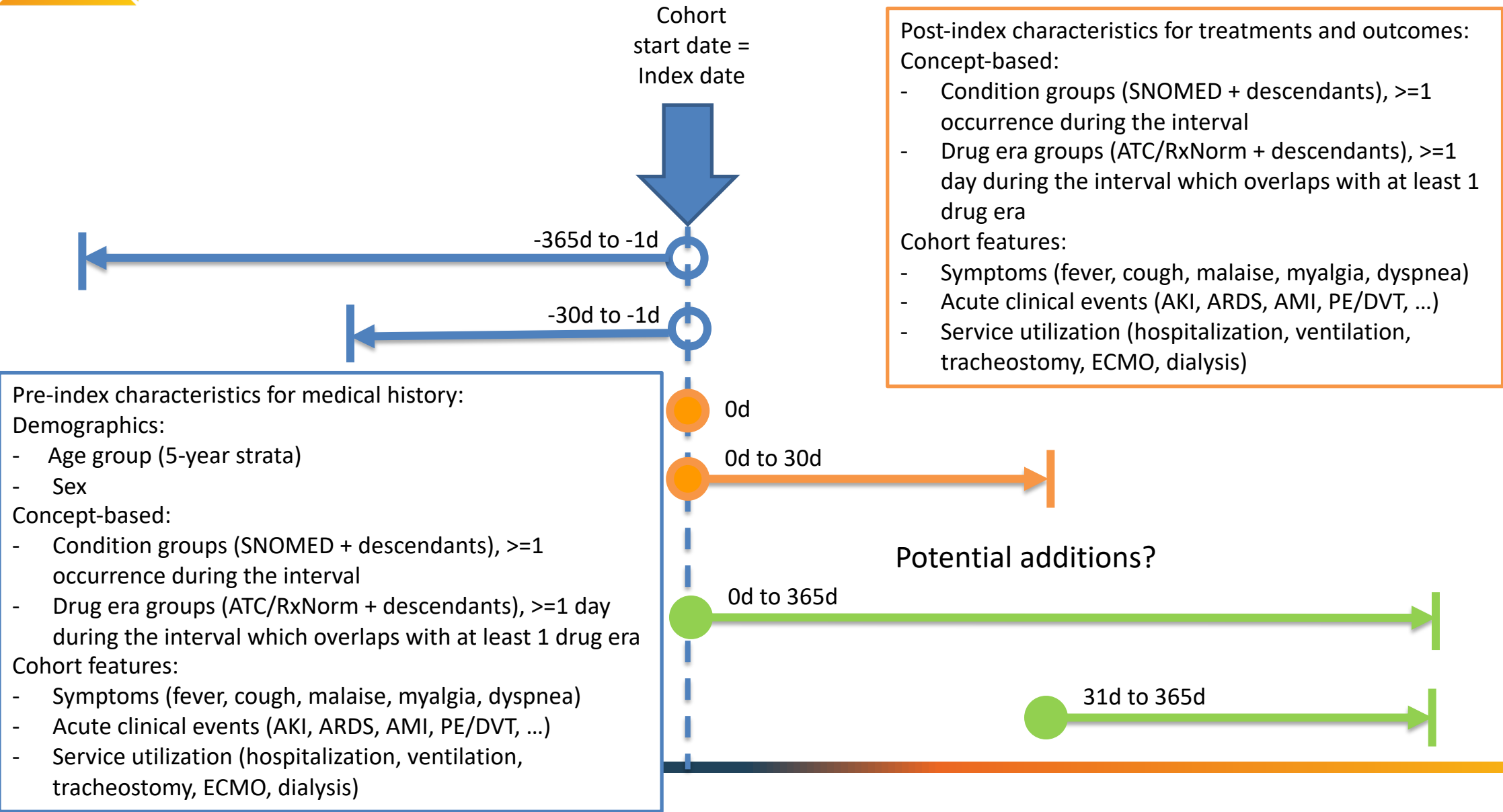
**Join the Journey!**

# CHARYBDIS target cohorts



Persons tested for SARS-COV-2

Persons tested positive for SARS-COV-2

Persons with COVID-19 diagnosis record OR a SARS-COV-2 positive test

Persons hospitalized with positive test for SARS-COV-2

Persons hospitalized with COVID-19 diagnosis record OR a SARS-COV-2 positive test

Persons hospitalized and requiring intensive services with positive test for SARS-COV-2

Persons hospitalized and requiring intensive services with COVID-19 diagnosis record OR a SARS-COV-2 positive test

POTENTIAL ADDITION?: Persons post-discharge after COVID-19 diagnosis record OR a SARS-COV-2 positive test

# CHARYBDIS subgroup cohorts

Persons tested for SARS-COV-2

Persons tested positive for SARS-COV-2

Persons with COVID-19 diagnosis record OR a SARS-COV-2 positive test

Persons hospitalized with positive test for SARS-COV-2

Persons hospitalized with COVID-19 diagnosis record OR a SARS-COV-2 positive test

Persons hospitalized and requiring intensive services with positive test for SARS-COV-2

Persons hospitalized and requiring intensive services with COVID-19 diagnosis record OR a SARS-COV-2 positive test

Stratification cohorts:
- Age:  <18, >65
- Gender:  Female/Male
- Race:  Black/White
- Index month
- Hypertension
- Type 2 Diabetes
- Heart disease
- Obesity
- Asthma
- COPD
- Chronic kidney disease
- End stage renal disease
- Cancer
- Autoimmune conditions
- Dementia
- HIV
- Pregnant women
- **Follow-up time: >=30d**

# CHARYBDIS Time windows

Cohort start date = Index date

-365d to -1d

-30d to -1d

0d

0d to 30d

Potential additions?

0d to 365d

31d to 365d

Post-index characteristics for treatments and outcomes:
Concept-based:
- Condition groups (SNOMED + descendants), >=1 occurrence during the interval
- Drug era groups (ATC/RxNorm + descendants), >=1 day during the interval which overlaps with at least 1 drug era

Cohort features:
- Symptoms (fever, cough, malaise, myalgia, dyspnea)
- Acute clinical events (AKI, ARDS, AMI, PE/DVT, …)
- Service utilization (hospitalization, ventilation, tracheostomy, ECMO, dialysis)

Pre-index characteristics for medical history:
Demographics:
- Age group (5-year strata)
- Sex

Concept-based:
- Condition groups (SNOMED + descendants), >=1 occurrence during the interval
- Drug era groups (ATC/RxNorm + descendants), >=1 day during the interval which overlaps with at least 1 drug era

Cohort features:
- Symptoms (fever, cough, malaise, myalgia, dyspnea)
- Acute clinical events (AKI, ARDS, AMI, PE/DVT, …)
- Service utilization (hospitalization, ventilation, tracheostomy, ECMO, dialysis)

# Data partners contributing to CHARYBDIS thusfar

| Database name | Geography | Data type |
|---|---|---|
| Premier | US (National) | Hospital billing |
| Optum EHR | US (National) | Electronic health records |
| Iqvia Open Claims | US (National) | Administrative claims |
| VINCI (VA) | US (National) | Electronic health records |
| STARR (Stanford) | US (CA) | Electronic health records |
| TRDW (Tufts) | US (MA) | Electronic health records |
| CUIMC (Columbia) | US (NY) | Electronic health records |
| SIDIAP | Spain | Electronic health records |
| SIDIAP-H | Spain | EHR-hospital linkage |
| HM Hospitales | Spain | Hospital billing |
| ICPI | Netherlands | Electronic health records |
| CPRD | UK | Electronic health records |
| HIRA | South Korea | Administrative claims |
| DCMC | South Korea | Electronic health records |

All databases standardized to OMOP CDM v5.3

# Live demo of CHARYBDIS

**CHARYBDIS** ☰

Show 25 entries        Search: [              ]

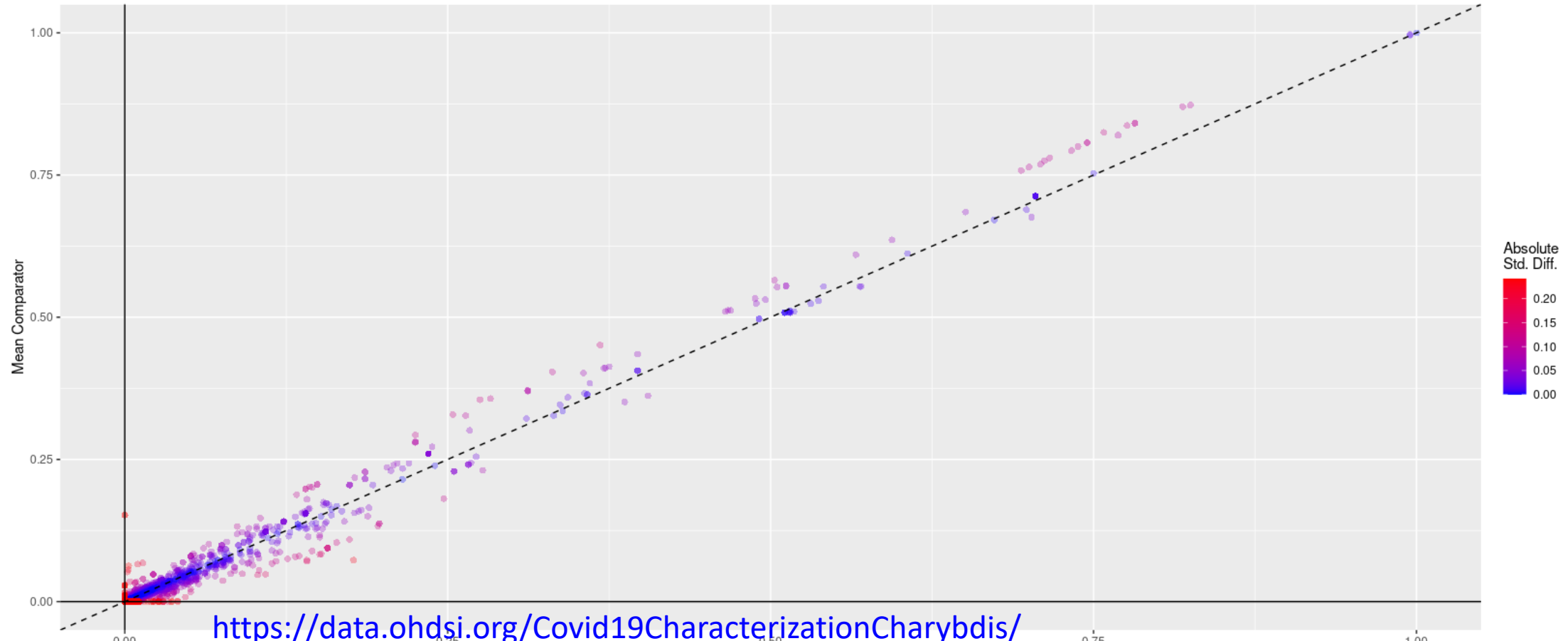| Cohort | Strata | CDM_Premier_COVID_v1240 Subjects | HIRA Subjects | optum_ehr_covid_v1239 Subjects | SIDIAP Subjects | STARR-OMOP Subjects | TRDW Subjects | IQVIA_OpenClaims Subjects | CUIMC Subjects |
|---|---|---|---|---|---|---|---|---|---|
| Persons tested for SARS-CoV-2 with no required prior observation | All | 219,230 | 230,268 | 411,580 | 150,187 | 56,881 | 6,950 | 783214 | 22094 |
| Persons tested for SARS-CoV-2 with at least 365d prior observation | All | 1,289 | 230,268 | 355,014 | 148,468 | 39,877 | 3,719 | 739518 | 18053 |
| Persons with a COVID-19 diagnosis or a SARS-CoV-2 positive test with no required prior observation | All | 66,132 | 7,603 | 45,508 | 124,221 | 4,788 | 1,250 | 493949 | 10437 |
| Persons tested with a COVID-19 diagnosis record or a SARS-CoV-2 positive test with no required prior observation | All | 21,503 | 6,013 | 43,386 | 42,325 | 4,095 | 1,097 | 74793 | 7998 |
| Persons tested positive for SARS-CoV-2 with no required prior observation | All | | | 42,909 | 37,975 | 1,880 | 1,035 | | 6959 |
| Persons with a COVID-19 diagnosis or a SARS-CoV-2 positive test with at least 365d prior observation | All | 194 | 7,603 | 37,880 | 122,058 | 3,328 | 664 | 466191 | 8519 |
| Persons tested with a COVID-19 diagnosis record or a SARS-CoV-2 positive test with at least 365d prior observation | All | 63 | 6,013 | 36,048 | 41,916 | 2,741 | 574 | 70301 | 6497 |
| Persons tested positive for SARS-CoV-2 with at least 365d prior observation | All | | | 35,624 | 37,604 | 902 | 520 | | 5625 |
| Persons hospitalized with a COVID-19 diagnosis record or a SARS-CoV-2 positive test with no required prior observation | All | 36,019 | 7,599 | 13,283 | 18,364 | 744 | 326 | 139971 | 3439 |
| Persons hospitalized with a SARS-CoV-2 positive test with no required prior observation | All | | | 12,451 | 13,644 | 128 | 232 | | 3075 |
| Persons hospitalized with a COVID-19 diagnosis record or a SARS-CoV-2 positive test with at least 365d prior observation | All | 132 | 7,599 | 10,534 | 18,197 | 615 | 186 | 133091 | 2600 |
| Persons hospitalized with a SARS-CoV-2 positive test with at least 365d prior observation | All | | | 9,841 | 13,520 | 86 | 140 | | 2344 |
| Persons hospitalized and requiring intensive services with a COVID-19 diagnosis record or a SARS-CoV-2 positive test with no required prior observation | All | 8,373 | 130 | 1,719 | | 62 | 102 | 15184 | 86 |
| Persons hospitalized and requiring intensive services with a SARS-CoV-2 positive test with no required prior observation | All | | | 1,611 | | 19 | 73 | | 58 |
| Persons hospitalized and requiring intensive services with a COVID-19 diagnosis record or a SARS-CoV-2 positive test with at least 365d prior observation | All | 28 | 130 | 1,345 | | 46 | 40 | 14633 | 56 |
| Persons hospitalized and requiring intensive services with a SARS-CoV-2 positive test with at least 365d prior observation | All | | | 1,253 | | 12 | 31 | | 40 |

Showing 1 to 16 of 16 entries      Previous | 1 | Next

https://data.ohdsi.org/Covid19CharacterizationCharybdis/

# Live demo of CHARYBDIS



https://data.ohdsi.org/Covid19CharacterizationCharybdis/

# Live demo of CHARYBDIS

# Live demo of CHARYBDIS



https://data.ohdsi.org/Covid19CharacterizationCharybdis/

# Live demo of CHARYBDIS cohort diagnostics



https://data.ohdsi.org/Covid19CharacterizationCharybdisDiagStrata/

# Live demo of CHARYBDIS cohort diagnostics



https://data.ohdsi.org/Covid19CharacterizationCharybdisDiagStrata/

# Using Twitter to characterize the COVID disease natural history and 'life after COVID'

Juan M. Banda

www.panacealab.org

Georgia State University

# Preface: Twitter is gaining attention for health-related research since 2009



Number of publications per year

Results of PubMed Query for Twitter and Health

# Benefits of using Twitter:

1) Good population representation

2) Everybody can post and have an account

3) Anonymity = unfiltered opinions

4) Data is freely available*

5) Tons of data generated each day (hundreds of millions of tweets get posted every day)

6) Easy filtering (hashtag usage, people mentions)



Distribution of global users over age groups



Number of users per country (millions)

https://www.statista.com/statistics/242606/number-of-active-twitter-users-in-selected-countries/

# Traditional disadvantages of using Twitter:

- Messy data (plenty of misspellings, shorthand, emojis, etc.)
  - There are at least 25 different ways people misspell hydroxychloroquine

- Attribution is an issue – are people just mentioning something or did it happen to them?

- Freely available data is only a 1% sample of whole set

- Collection is hard and needs to be ongoing for days/weeks before getting considerable mass

# The COVID opportunity – highly focused data

## The dataset:

- 490+ Million Tweets

- ONLY COVID related chatter is included

Longitudinal – January 27ᵗʰ to today… and growing



Dataset: https://doi.org/10.5281/zenodo.3723939

Pre-print: https://arxiv.org/abs/2004.03688

Recent additions: https://github.com/thepanacealab/covid19_twitter

# Current work: Drug characterization

- Methods to deal with misspellings and noisiness of data:

Table 2. Drug ingredient mentions found

| Drug Ingredient | Frequency |
|---|---|
| hydroxychloroquine | 204,879 |
| remdesivir | 72,841 |
| chloroquine | 49,915 |
| oxygen | 37,961 |
| vitamin D | 25,445 |
| dexamethasone | 25,142 |
| zinc | 24,843 |
| azithromycin | 16,079 |
| ibuprofen | 8,469 |
| ivermectin | 6,390 |



Figure 1. Timeline of Tweets with potential drug treatment mentions.

- Charybdis-like characterization over countries (work with Dani Prieto-Alhambra – University of Oxford)

* https://openreview.net/forum?id=qlGPXs9FWa

# Current Work: Symptom/condition detection

- Self-reported symptoms on Twitter vs EHR lists *
  - Can we find related symptoms both found on EHR's (Callahan, A., Steinberg, E., Fries, J.A. et al. Estimating the efficacy of symptom-based screening for COVID-19. npj Digit. Med. 3, 95 (2020). https://doi.org/10.1038/s41746-020-0300-0) but on Twitter?

| Term | Frequency |
|------|-----------|
| pneumonia | 110124 |
| infection | 71882 |
| influenza | 36390 |
| cough | 35753 |
| anxiety | 34658 |
| pain | 12773 |
| depression | 12189 |
| asthma | 8307 |

* https://github.com/thepanacealab/covid19_biohackathon/tree/master/user_symptoms

# What does this lead to?

- Since we can find symptoms and drugs, we can also find people that had COVID and their symptoms after infection!

  - On-going work with Dani Prieto-Alhambra and others
    - Incorporates methods shown before + manual review by clinicians

JAMA doi:10.1001/jama.2020.12603

Some very preliminary findings:

| | |
|---|---|
| fatigue | 789 |
| shortness of breath=dyspnea | 701 |
| chest pain | 687 |
| palpitations | 674 |
| anxiety | 212 |
| post-exertional malaise | 36 |
| Tired = fatigue | 36 |
| muscle pain = myalgia | 35 |

# The gory details:

- Technical stuff:
  - "Building tools and frameworks for large-scale social media mining: Creating data infrastructure for COVID-19 research" **dair.ai meetup 7/22:** https://www.meetup.com/dair-ai/events/271690722/

- Extended version of today's short talk:
  - "Leveraging the OHDSI vocabulary to characterize the COVID-19 epidemic using Twitter data and NLP" **OHDSI community call 7/21**: https://www.ohdsi.org/web/wiki/doku.php?id=projects:ohdsi_community