28 Days, 28 Phenotypes

OHDSI

Phenotype
Phebruary

forums.ohdsi.org

Join The Conversation!

# OHDSI Phenotype Phebruary: lessons learned

**2022 OHDSI Symposium**

**14 October 2022 Bethesda MD**

**Azza Shoaibi, Joel Swerdel, Allan Wu, Gowtham Rao, Adam Black, Evan Minty, Asieh Golozar, Rupa Makadia, Jill Hardin, Erica Voss, Tiffany J. Callahan, Juan Banda, Anna Ostropolets, Claudia Pulgarin, Marcela Rivera, David Vizcaya, Patrick Ryan**

# Phenotype development and evaluation is yet to become a completed chapter in the book of OHDSI

- Phenotypes are the foundational elements in almost every real-world analysis.

- The reliability of the generated evidence depends on the validity of the phenotypes.

- Yet, the science of phenotype development and evaluation is relatively immature.

- We have built best practices and end to end process, tools and packages for characterization, estimation and prediction.

- But for phenotyping- **"addressed in the limitation section."**

**" Phenotype Phebruary":** I realized that becoming a master of karate was not about learning 4,000 moves but about doing just a handful of moves 4,000 times." — Chet Holmes

- We collectively started a discussion on 28 phenotypes over 28 days

- Followed 5 step process:

| clinical description | prior work | cohort definition | evaluate | Discuss & Share |

### 1. Join the conversation

- Discussions will be held on forums.ohdsi.org
- Each day will be a new thread
- Explore the definitions and review the results provided
- Reply with your thoughts, reflections, insights and question

### 2. Evaluate the cohort definitions in your data

- Execute cohort definitions and CohortDiagnostics in your CDM
- Share insights you learn from your data on the forums
- Share results to compile across the network on data.ohdsi.org



**Phenotype Phebruary • Daily Threads & What We Learned**

"Phenotype Phebruary" was a community-wide initiative to both develop and evaluate phenotypes for health outcomes that could be investigated by the community. Patrick Ryan introduced this initiative in both a video presentation and a forum post, and each of the conversations around the "28 phenotypes for 28 days" are being held within the OHDSI forums.

This page will provide direct links to each forum post, which is where conversations around each specific phenotype should be held. The video on the right includes "phun phacts" shared about each phenotype during our weekly community calls.

#### Daily Phenotype Phebruary Links

*(future dates are subject to change)*

Feb. 1 • Type 2 Diabetes Mellitus
Feb. 2 • Type 1 Diabetes Mellitus
Feb. 3 • Atrial Fibrillation
Feb. 4 • Multiple Myeloma
Feb. 5 • Alzheimer's Disease
Feb. 6 • Hemorrhagic Events
Feb. 7 • Neutropenia
Feb. 8 • Kidney Stones
Feb. 9 • Delirium
Feb. 10 • Systemic Lupus Erythematosus
Feb. 11 • Suicide Attempts
Feb. 12 • Parkinson's Disease and Parkinsonism
Feb. 13 • Attention Deficit Hyperactivity Disorder
Feb. 14 • Hypertension *(Video Description)*
Feb. 15 • Acute Myocardial Infarction
Feb. 16 • Heart Failure
Feb. 17 • Cardiomyopathy
Feb. 18 • Multiple Sclerosis
Feb. 19 • Triple Negative Breast Cancer
Feb. 20 • Pulmonary Hypertension
Feb. 21 • Prostate Cancer
Feb. 22 • HIV
Feb. 23 • Hidradenitis Suppurativa
Feb. 24 • Anaphylaxis
Feb. 25 • Depression
Feb. 26 • Non-Small-Cell Lung Cancer
Feb. 27 • Drug-Induced Liver Injury
Feb. 28 • Severe Visual Impairment And Blindness
Bonus • Acute Kidney Injury

# 15 phenotypes were developed, evaluated and discussed and we learned few things

**Thematic learnings:** clinical description, phenotype development and phenotype evaluation. The themes identified belonged to **5 different types of lessons**: tips, strategies, debatable topics, challenges, and opportunities.

# Lessons learned: Phenotype development

## Tips/strategies

Evaluate **all** types of measurement error

Use patient profile to get a sense of validity. Identify disqualifying patterns.

Explore in CD:
temporal stability, expected trends, patients composition, index event misclassification

Estimate measurement error & quantify trade-offs by PheValuator or APHRODITE

## Challenges

Subjective

Time-consuming and complex

Lack an approach to evaluate exit criteria & washout periods for reoccurring events



**What objective diagnostic criteria can we apply to determine fitness-for-use'?**

# Lessons learned: Phenotype evaluation

## Tips/strategies

Model the clinical idea not the analytical use case

Code selection is a clinical choice. The material consequences should be empirically investigated

Notions like "primary position" need to be standardized and evaluated

Differentiate between situations where a data is not "fit-for use" and situations where logic is not "fit" for the data

## Opportunities

Systematically assess multiple look back periods and recommend one.

Combine a heuristic-based approach (APHRODITE) with a rule-based approach (Atlas)

Develop a PubMed search strategy to find published/evaluated phenotypes

OHDSI/Aphrodite
[in development]

Contributors 3    Issue 1    Stars 31    Forks 10

PHOEBE
About
Initial Concept
PHenotype Observed Entity Baseline Endorsements (PHOEBE)

**Do we customize phenotypes to specific data sources/analytical use ?**

# The choices are NOT:

**1: Code list**
**2: Code list with chart reviews**



**Required Inputs:**

Clinical description

Prior phenotype knowledge from literature or other sources

**Expected Outputs:**

{Cohort definition(s)}

Cohort diagnostics

PheValuator results (when possible)

Evaluation report

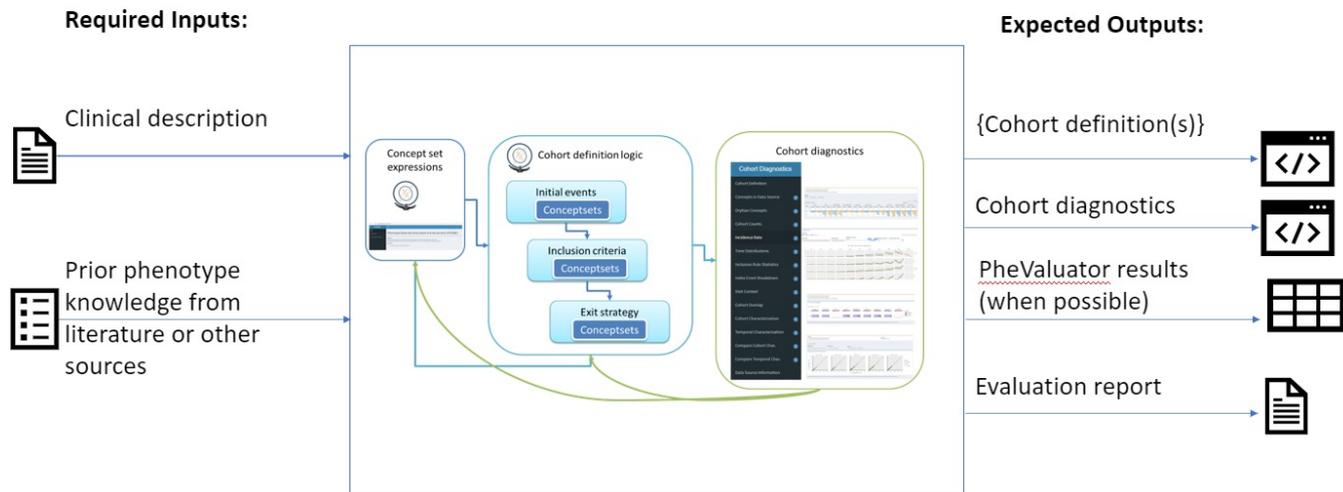- Phenotyping is complex, multidimensional and requires exchange of knowledge, learnings and insights across collaborators from different background and expertise

- Large scale characterization (e.g.CD), Diagnostic predictive models (e.g., PheValuator) and structured review of patient's profile are potentially effective and novel strategies for phenotype evaluation.

- We are getting closer to a standardized process. **But** further collaboration is needed to formalize a **scalable and reproducible** processes and establish **empirically-driven objective diagnostics**