# ICPE 2021 Symposium
# New frontiers in computable phenotyping for medical product safety evaluation

xx

**Presented at ICPE 2021 All Access**

Sentinel

# Improving Outcome Ascertainment by Applying Natural Language Processing and Machine Learning to Electronic Health Record Data:
# Identifying Anaphylaxis

David S. Carrell, PhD

Kaiser Permanente Washington Health Research Institute

# Overview

1. Motivation & objectives
2. Study design
    1. Study cohort
    2. Natural language processing (NLP)
    3. Structured data
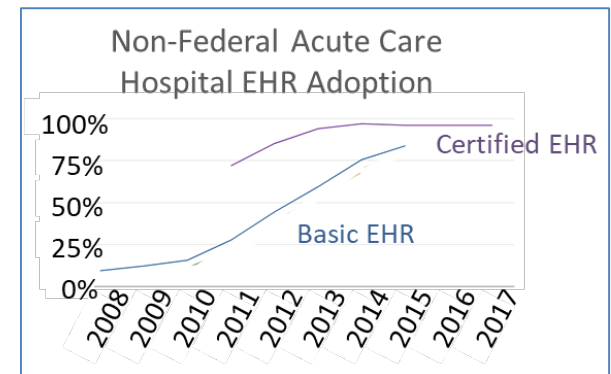    4. Machine learned-models
3. Results and implications

# **Motivation**: Improving ARIA sufficiency

Existing algorithms …

- Rely on structured data (Dx, Px, Rx, demographics, …)
- Have good sensitivity
- Lack positive predictive value
  - <2/3 are true cases (Walsh et al. 2013)

A challenging outcome to model

- Rare (limited training data)
- "Rule-out" coding/mis-diagnosis
- Complex diagnosis
  - Ball et al. 2018: NLP of chart notes may help

EHR data = opportunity?



Non-Federal Acute Care Hospital EHR Adoption

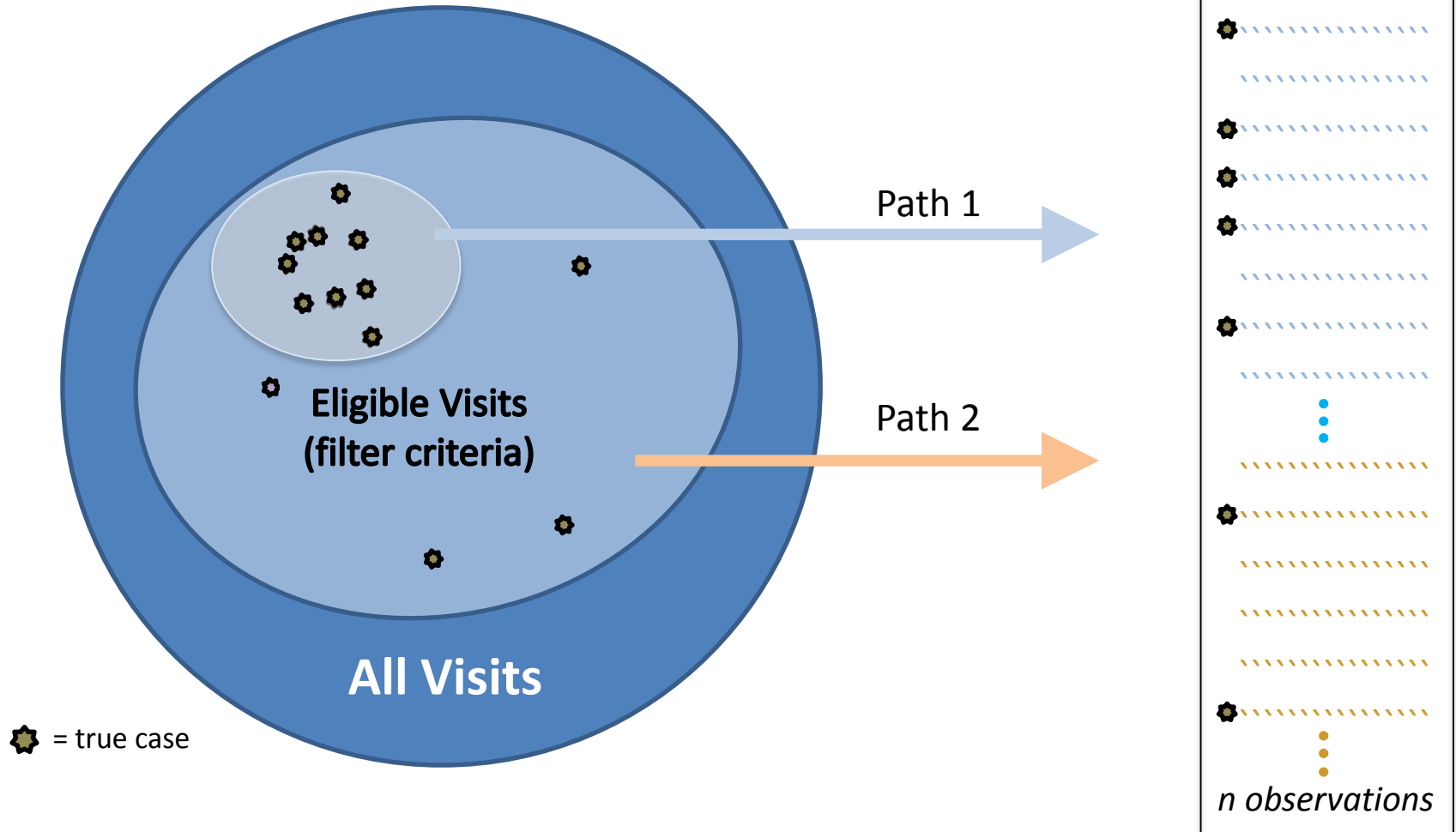# **Objective**: Improve outcome identification

- Use **NLP**-extracted data to enrich covariates
  - Are clinical diagnostic criteria documented?
    - Organ system involvement (e.g., skin, respiratory, BP)
    - Clinical course (e.g., rapid onset)
  - Telltale utilization
    - Treatments (e.g., *multiple* epinephrine administrations)
    - Hospital admission "for observation"
  - Are competing explanations described?
- Use **machine learning** to better model "signal" in a rich set of covariates

# **Design:** population, outcomes, covariates

- Study period: 10/2015 – 12/2018
- Population: Age ≥1-year
  - Kaiser Permanente Washington (KPWA)
  - Kaiser Permanente Northwest (KPNW)
- Eligibility
  - Anaphylaxis diagnosis (ED/inpatient or outpatient)
  - ≥12 months prior enrollment (*w/o anaphylaxis diagnosis*)
- Gold standard outcomes (clinician review)
- Covariates (manually engineered)
  - Structured: Demographics, Dx, Px, Rx, encounters
  - NLP-derived: Symptoms, clinical criteria, …

# Stratified Random Sampling

Goal is to sample enough cases, while ensuring the analytic dataset faithfully represents the source population

# **Design:** Gold standard creation

- KPWA:
  - Dual blind manual review by clinicians
  - Decisions recorded on spreadsheet
- KPNW
  - Dual blind manual review by non-clinician abstractors following a written protocol
  - Decisions, supporting documentation in REDCap
  - Difficult cases → clinician review

# **Design:** Manual covariate curation

- Clinicians & informaticists reviewed/discussed charts

Nose: No rhinorrhea.
Mouth: Mild swelling
Neck: Nontender, supple, no lymphadenopathy
Lymphatic: No lymphadenopathy noted.
Cardiovascular: Normal heart rate, normal rhythm, no murmurs, no rubs, no gallops. Intact distal pulses, no tenderness, no cyanosis, no clubbing.
Respiratory: Normal breath sounds, no respiratory distress, no wheezing, no chest tenderness. No severe stridor, severe wheezing
Abdomen: Bowel sounds are present. Abdomen is soft, no tenderness, no masses, no rebound or guarding. No organomegaly. No hernia.
GU: No CVA tenderness. Bladder is nontender and not distended.
Skin: Erythema noted about the face and minimally to the hands
Back: No tenderness
Musculoskeletal: No tenderness to palpation or major deformities noted. No back or cervical spine tenderness. No edema.

Pt after her CTA ABdomen she develop allergic /anaphylactic reaction in ED with nausea/vomting and tachycardia and hypotensive and she became hypoxic ,even so she had many ct with contrast without any reactions

She received multiple rounds of epinephrine , benadryl ,decadron ,pepcid

SHE FEEL MUCH BETTER NOW except some dizziness when she walk

- Curated structured and NLP covariates we judged *clinically relevant and feasible*

- *We did <u>not</u> use gold standard labels to curate covariates (due to small sample size)*

# **Design:** Structured covariates

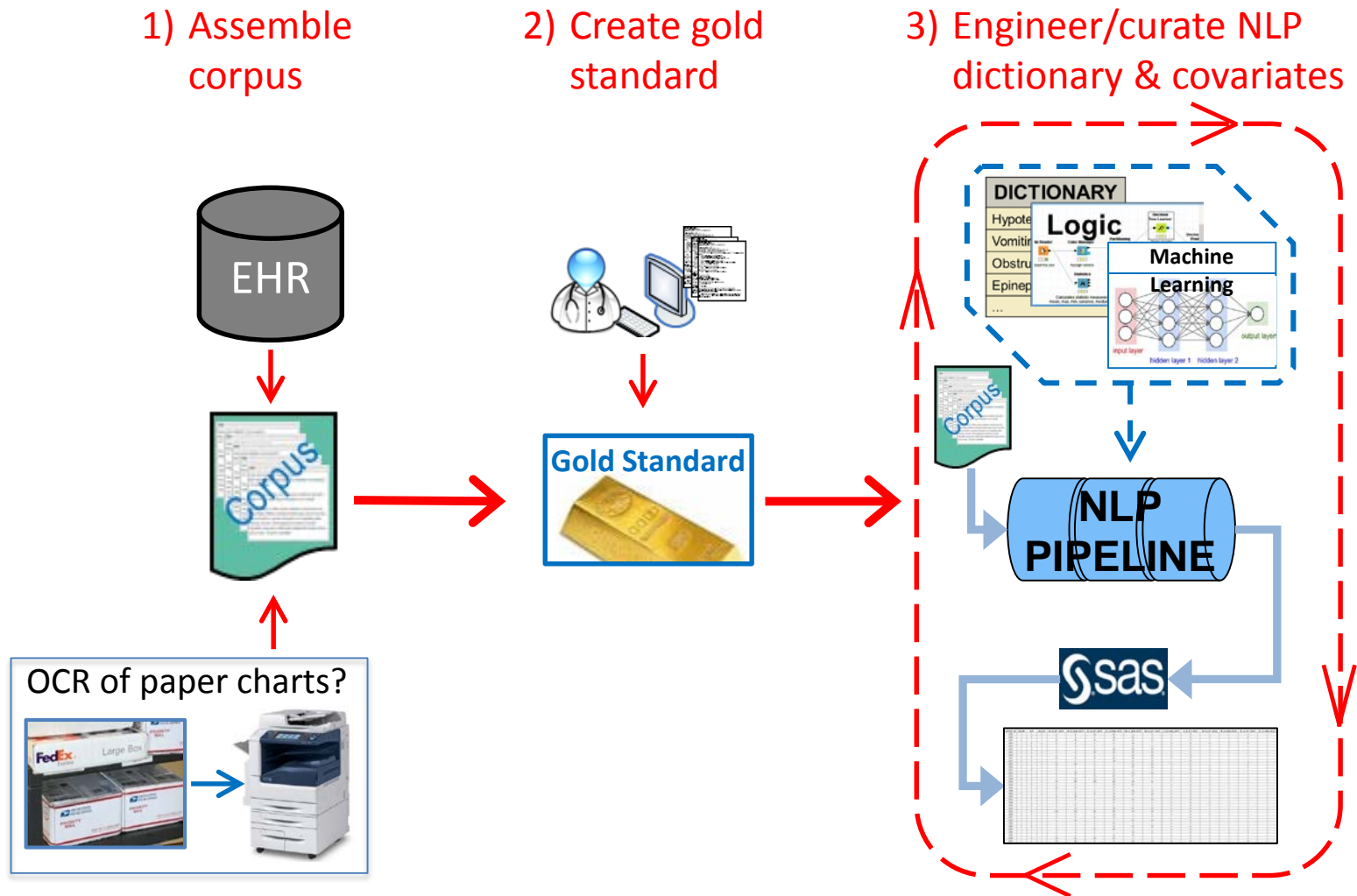Manually curated from the Sentinel common data model

| Anaphylaxis Structured Covariates | |
|---|---|
| **Category** | **Count** |
| Demographics (age, sex, race, enrollment history) | 6 |
| Care setting (ED, IP, outpatient) | 6 |
| History of allergic reaction/anaphylaxis | 4 |
| Exposures (e.g., imaging dye, immunotherapy) | 3 |
| Treatment (e.g., epinephrine, steroids, intubation, CPR) | 10 |
| Competing diagnoses (asthma, COPD, angioedema, infection) | 11 |
| Other (summer event, labs, immunology follow-up) | 3 |
| TOTAL: | 43 |

# **Design:** Covariate curation – NLP-derived

## **NLP definitions**

- **NLP** – **Converts** information in **unstructured clinical text** to **structured data** using methods from computer science, artificial intelligence, and computational linguistics

- *Manual* **NLP** – Human curation of NLP dictionaries and NLP-derived covariates guided by domain-specific clinical knowledge, informatics expertise, and "gold standard" data

- *Automated* **NLP** – (semi)automated engineering of NLP dictionaries and covariates using "silver standard" data and data-driven approaches to algorithm development

# **Design:** Covariate curation – NLP process



1) Assemble corpus

2) Create gold standard

3) Engineer/curate NLP dictionary & covariates

EHR

OCR of paper charts?

Gold Standard

DICTIONARY
Logic
Machine Learning

Corpus

NLP PIPELINE

sas

# **Design:** Manual NLP process – dictionary

- 843 terms

  >50% "skin/mucosal"

- Concepts per chart:
  Median:      128
  Min:            9
  Max:        2,092

| ID | CUI | TEXT | SOURCE | SOURCETYPE |
|----|-----|------|--------|------------|
| 3001 | GI001 | abd pain | GI | ABDOPAIN |
| 6001 | SM001 | abdomen with erythema | GI | ABDOPAIN |
| 3002 | GI002 | abdominal pain and shock | GI | ABDOPAIN |
| 2001 | BP001 | acute hypotensive | BPREDUCED | HYPOTENSION |
| 5001 | RC001 | acute hypoxic | RESPCOMP | HYPOXIA |
| 5002 | RC002 | acute respiratory failure | RESPCOMP | RESPFAIL |
| 5003 | RC003 | acute upper airway obstruction | RESPCOMP | AIRWAY |
| 4001 | OT001 | admission diagnosis | OTHER | DIAGNOSIS |
| 4002 | OT002 | admitting diagnosis | OTHER | DIAGNOSIS |
| 5004 | RC004 | airway narrowing | RESPCOMP | AIRWAY CONSTRICTION |
| 5005 | RC005 | airway obstruction | RESPCOMP | AIRWAY CONSTRICTION |
| 6002 | SM002 | airway itch | SKINMUC | AIRWAY |
| 6003 | SM003 | airway remains swolen | SKINMUC | ORALSWELL |
| 6004 | SM004 | airway remains swollen | SKINMUC | AIRWAY |
| 4003 | OT003 | alergic reacton | OTHER | ALLERGREACT |
| 6005 | SM005 | all skin appears red | SKINMUC | RASH |
| 4004 | OT004 | allergic reaction | OTHER | ALLERGREACT |
| 4005 | OT005 | allergic reacton | OTHER | ALLERGREACT |
| 4006 | OT006 | allergic to | OTHER | HYPO |
| 4007 | OT007 | allergies | OTHER | HYPO |
| 4008 | OT008 | allergy comment | OTHER | HYPO |
| 2002 | BP002 | almost passed out | BPREDUCED | SYNCOPE |
| 5006 | RC006 | altered mentation | RESPCOMP | ALTERED MENTATION |
| 1001 | AN001 | anaphalytic shock | ANAPH | ANAPH SHOCK |
| 1002 | AN002 | anaphylactic shock | ANAPH | ANAPH SHOCK |
| 1003 | AN003 | anaphylaxis allergic shock | ANAPH | ANAPH SHOCK |
| 4009 | OT009 | anaphylaxis | OTHER | ANAPH |
| 2003 | BP003 | and hypotensive | BPREDUCED | HYPOTENSION |
| 2004 | BP004 | and passed out | BPREDUCED | SYNCOPE |
| 2005 | BP005 | and shock | BPREDUCED | SHOCK |
| 6006 | SM006 | angioedema | SKINMUC | ANGIOEDEMA |
| 1004 | AN004 | anhylactic shock | ANAPH | ANAPH SHOCK |

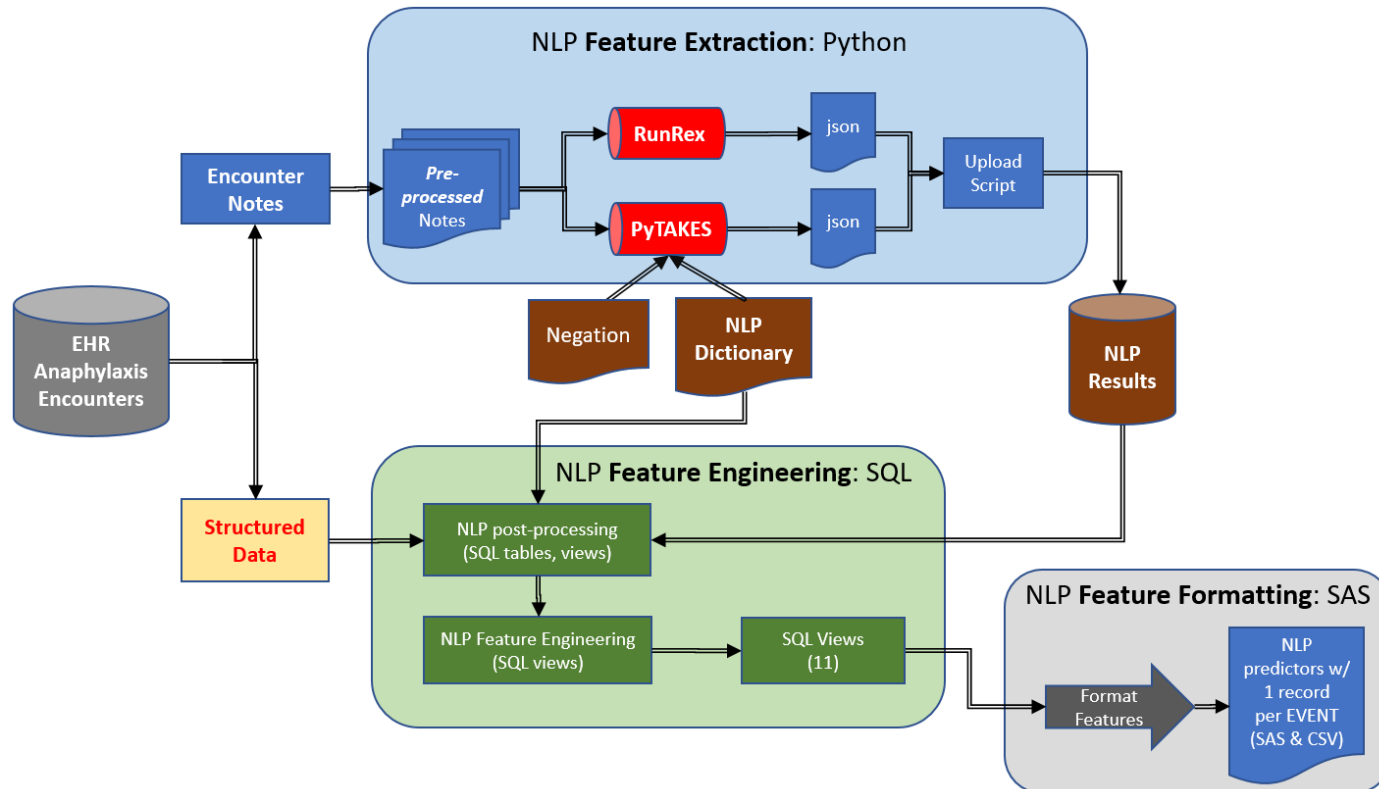# **Design:** Manual NLP process – dictionary

## Anaphylaxis concepts in the NLP dictionary (N terms)

- BRADYCARDIA (13)
- CARDIACARRHYTH (8)
- CARDIOCOLLAPSE (2)
- COLLAPSE (2)
- END ORGAN (2)
- HYPOTENSION (77)
- PALPITATIONS (3)
- SHOCK (3)
- SYNCOPE (30)
- TACHYCARDIA (9)
- **ABDOPAIN (3)**
- **VOMIT (1)**
- AIRWAY (4)
- AIRWAY CONSTRICTION (4)
- ALTERED MENTATION (1)
- APHONIA (3)
- BREATH (6)
- BRONCHOSPASM (1)
- CHEST DISCOMFORT (2)
- CHEST TIGHTNESS (9)

- COARSE BREATH SOUND (4)
- DYSPHONIA (1)
- DYSPNEA (55)
- HOARSENESS (7)
- HYPOXEMIA (6)
- HYPOXIA (3)
- IMPENDING DOOM (2)
- INTUBATION (6)
- LARYNGEAL OEDEMA (1)
- RESP COMPROMISE (3)
- RESP DISTRESS (2)
- RESPFAIL (1)
- RONCHI (2)
- STRIDOR (3)
- TACHYPNEA (5)
- THROAT CLOSURE (14)
- THROAT TIGHTNESS (34)
- TIGHTNESS BREATHING (1)
- VOICE QUALITY (1)
- WHEEZE (8)

- ANGIOEDEMA (102)
- DIFFICULTY SWALLOWING (14)
- DYSPHAGIA (1)
- EDEMA (4)
- ERYTHEMA (42)
- EYE SWELLING (33)
- FACIAL SWELLING (20)
- FLUSH (38)
- HIVES (68)
- ITCHING (14)
- ITCHY SOFT TISSUE (15)
- METALLIC TASTE (1)
- MOUTH (1)
- MOUTHSWELL (4)
- ORALSWELL (4)
- PRURITUS (15)
- RASH (7)
- REACTION (1)
- SOFT TISSUE SWELLING (4)
- SWELLING (31)

- THROAT (4)
- TINGLING (1)
- TINGLY SOFT TISSUE (14)
- URTICARIA (24)
- ALLERGREACT (5)
- ANAPH (5)
- COMPLAINT (12)
- DIAGNOSIS (8)
- DIFFERENTIAL (1)
- HYPO (6)
- IMPRESSION (1)

● **REDUCED BLOOD PRESSURE** ● **GASTROINTESTINAL** ● **RESPIRATORY COMPROMISE** ● **SKIN/MUCOSAL**
● **OTHER**

# **Design:** Transportable NLP system

- Developed & applied at KPWA
- Transported to KPNW via GitHub
  - NLP system (Python), SQL queries, SAS code, documentation
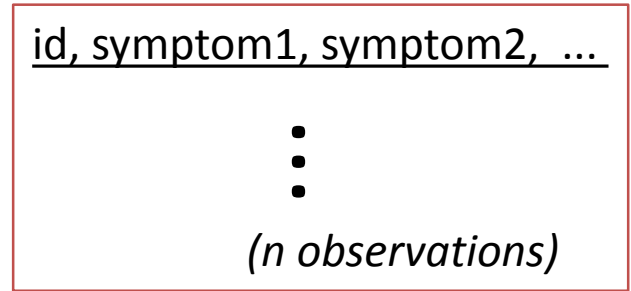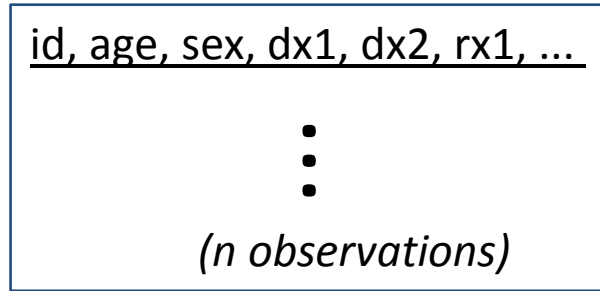
# Design: NLP covariates

- 116 NLP covariates engineered for use in modeling (selected from >450 candidates):

| Anaphylaxis NLP Covariates | |
|---|---|
| **Category** | **Count** |
| Symptoms (skin/mucosal, respiratory compromise, reduced BP) | 10 |
| Anaphylaxis concepts (e.g., wheezing, epinephrine, …) | 66 |
| Diagnostic criteria (e.g., skin/mucosal + [resp. comp. *or* ↓BP]) | 30 |
| Explicit diagnoses of anaphylaxis | 5 |
| "Special features" (e.g., admitted to hospital for observation) | 5 |
| TOTAL: | 116 |

# Model Development

Structured Data in Sentinel CDM + labs     EHR Text-based (NLP) covariates

1. Collect Data

id, age, sex, dx1, dx2, rx1, ...

⋮

*(n observations)*

id, symptom1, symptom2, ...

⋮

*(n observations)*

2. Prescreen Covariates

3. Develop Model

1

2

4. Obtain Predictions, Classifications

| 0.92 | CASE |
|------|------|
| 0.01 | CONTROL |
| 0.84 | CASE |

⋮

| 0.97 | CASE |
|------|------|
| 0.02 | CONTROL |
| 0.63 | CONTROL |

⋮

# What's in the box?

- Logistic regression
- Elastic net
- Bayesian Additive Regression Trees
- Neural network
- Boosted Trees

Super Learner
(a weighted combination)

$$\beta_0 + \beta_1 * age + \beta_2 * ICD10 + \dots$$

Boosted Regression Tree is a hierarchical and supervised machine learning method that combines weak learners (binary splits) to strong prediction rules that allow a flexible partition of the feature space.

input layer

hidden layer 1    hidden layer 2    output layer

$x_5 < c$          $x_5 \geq c$

$\mu_3$

$x_2 < d$    $x_2 \geq d$          $\Leftrightarrow$

$\mu_1$    $\mu_2$

$x_5$

$\mu_3$

$c$

$\mu_1$    $\mu_2$

$d$    $x_2$

# 75 Models

| Algorithm | R package name | Notes on tuning parameters |
|---|---|---|
| 1. Logistic regression | (base) | |
| 2. Elastic net | glmnet | 10-fold cross validation to select optimal alpha and lambda |
| 3. Gradient boosting | xgboost | Variant 1: maximum tree depth = 2<br>Variant 2: maximum tree depth = 4 |
| 4. Bayesian Additive Regression Trees | dbarts | Variant 1: k = 2 (default),<br>Variant 2: k=1 (reduced regularization prior) |
| 5. Neural network (feed forward) | neuralnet | Variant 1: 1 hidden layer containing 1 node<br>Variant 2: 1 hidden layer containing 3 nodes |
| 6. Super Learner | SuperLearner | |

$$3 \quad \text{x} \quad ( 3 \quad \text{x} \quad 8 \quad + \quad 1) \quad = \quad 75$$

| Datasets | Covariate Selection | Variants of six | SL |
|---|---|---|---|
| structured data | none | prediction | weighted |
| structured+NLP | lasso | algorithms | combination |
| struct+clinicianNLP | clustering | | |

# Results

| Path | KPWA (n=239) | | KPNW (n=277) | |
|---|---|---|---|---|
| | Cases | Controls | Cases | Controls |
| 1 | 106 (65.8%) | 55 (34.2%) | 115 (70.6%) | 48 (29.4%) |
| 2 | 48 (61.5%) | 30 (38.5%) | 65 (57.0%) | 49 (43.0%) |
| all | 154 (64.4%) | 85 (35.6%) | 180 (65.0%) | 97 (35.0%) |

# Results

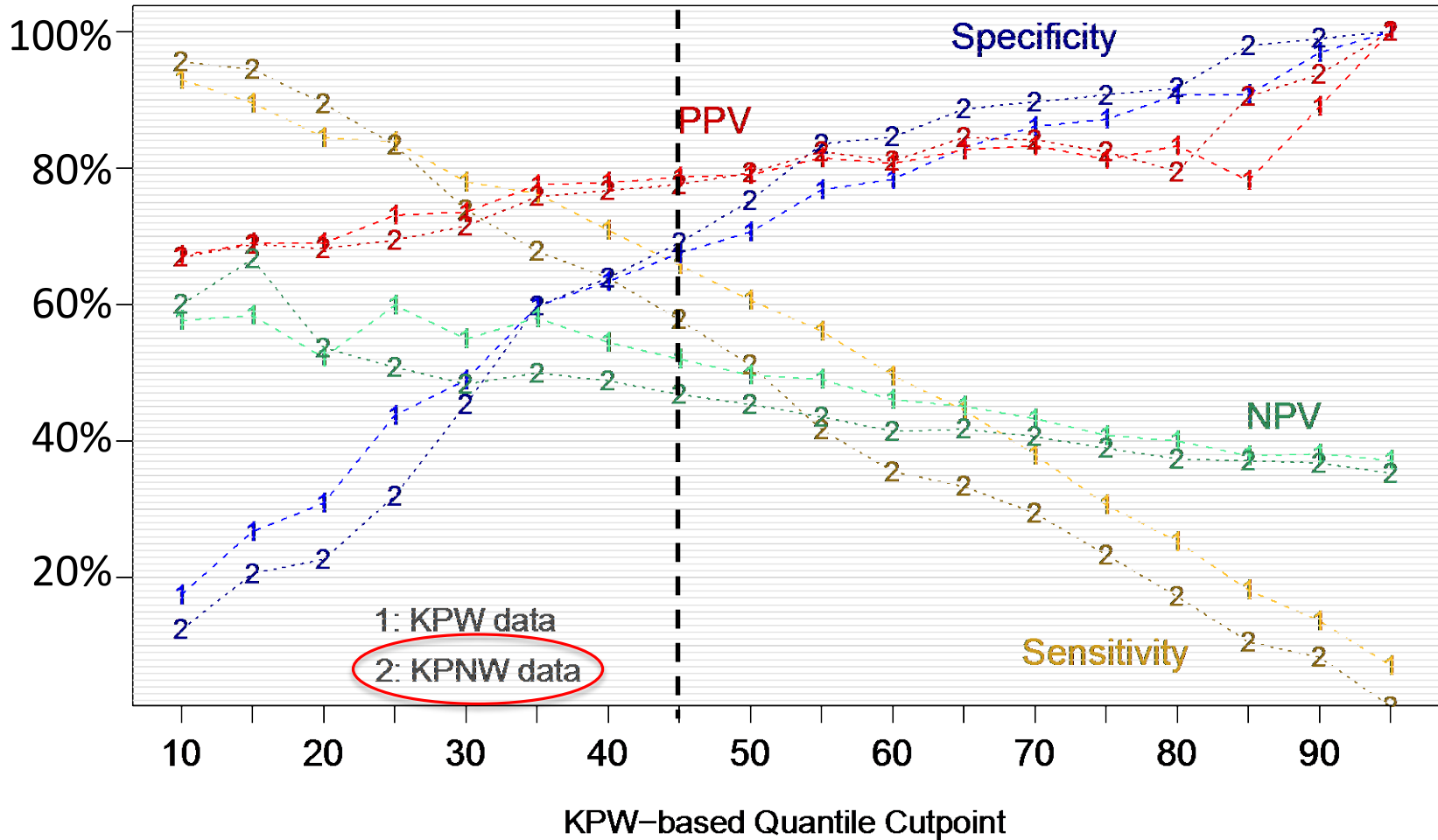Cross-validated AUCs for best models for each KPWA data set

# Results

- Two versions of Bayesian Additive Regression Trees combining structured data with NLP-derived covariates were nearly identical

- BART2-RetainAll generalized best to KP Northwest external validation set
  - cvAUC at KPWA = 0.70, cvAUC at KPNW = 0.67
  - Next step: Choose a prediction risk threshold for classification
    - if risk >=  *threshold*, classify as a case, otherwise a control
    - most interested in high positive predictive value (PPV), high sensitivity (% cases identified)

# **Results:** Performance Metrics

# **Results:** Performance Metrics



23

# Implications

- NLP-derived covariates derived from EHR data improve algorithm performance

- Machine-learning models are well-suited to this type of data

- Next steps:
  - Explore two-stage models (to correct classification errors)
  - Explore modeling all data (KPWA 239 + KPNW 277 = 516)
  - Explore (semi)automated NLP approaches

# Acknowledgements*

*Study team members listed alphabetically*

**FDA**

Adebola Ajao

Robert Ball

Steven Bird

Sara Karami

Yong Ma

Michael Nguyen

Danijela Stojanovic

Mingfeng Zhang

Yueqin Zhao

**Harvard Pilgrim**

Adee Kennedy

Judy Maro

Mayura Shinde

**Kaiser Washington**

Maralyssa Bann

David Carrell

David Cronkite

James Floyd

Monica Fujii

Vina Graham

Kara Haugen

Ron Johnson

Jennifer Nelson

Mary Shea

Jing Zhou

**Kaiser Northwest**

Andrew Felcher

Brian Hazlehurst

Denis Nyongesa

Daniel Sapp

Matthew Slaughter

**Putnam Data Science**

Susan Gruber

**Vanderbilt University**

Cosmin (Adi) Bejan

**HealthCore**

Kevin Haynes

# *Thank You!*

# *Questions & Discussion*

David Carrell – david.s.carrell@kp.org

# Extra Slides

| Priorities | Goals | Initiatives | Outputs |
|---|---|---|---|
| Establishing data infrastructure | Establishing a Sentinel electronic health record (EHR) network requires determining where to source and how to structure the data, as well as implementation of robust governance, harmonization, and quality assurance (QA) processes. | • Horizon scan of EHR databases<br>• Adding unstructured data to the Sentinel common data model<br>• Assessment and validation of source data mappings to improve the reliability and reproducibility of real-world data sources<br>• Harmonizing EHRs from heterogenous systems<br>• Developing and integrating approaches to identifying date and cause of death<br>• FHIR implementation preparedness | • EHR data partners<br>• Set of necessary EHR data elements<br>• EHR common data model<br>• Data governance process<br>• Data harmonization and QA strategy<br>• Data quality metrics<br>• Sentinel death index<br>• FHIR strategy |
| | Frameworks and tools are needed for extracting critical information from EHR data to enable and enhance EHR-based computable phenotyping and to support EHR-based descriptive, inferential, and detection queries in Sentinel. | • Extending machine learning methods development in Sentinel: follow-up analyses for anaphylaxis algorithm and formalization of a general phenotyping algorithm<br>• Scalable automated natural language processing- (NLP-) assisted chart abstraction<br>• Advancing scalable NLP approaches for unstructured EHR data<br>• Improving probabilistic phenotyping of incident outcomes through enhanced ascertainment with NLP | • Computable phenotyping framework<br>• NLP tools for cohort identification, exposure assessment, covariate ascertainment, and outcome identification<br>• Chart review automation approaches<br>• Automated feature extraction tool to improve confounding control in EHR data<br>• NLP-assisted chart abstraction tool |
| | Developing, evaluating, and implementing advanced epidemiologic and statistical methods will enable Sentinel to make best use of EHR data to increase Active Risk Identification and Analysis (ARIA) sufficiency and expand the acceptance and use of real-world data for regulatory decision-making. | • Empirical evaluation of the causal inference effects of utilizing best practices for pharmacoepidemiologic studies<br>• Enhancing causal inference in the Sentinel system: an evaluation of targeted learning and propensity scores<br>• Approaches for handling missing laboratory data<br>• Subset calibration for detecting and correcting for bias<br>• Development of performance metrics and reporting standards<br>• Advancing distributed regression in Sentinel | • Causal inference design and analysis framework<br>• Super learner, target maximum likelihood estimation, complex treatment strategy analysis, missing data, subset calibration, and distributed regression tools<br>• Inferential query performance metrics and reporting standards |
| | Building safety signal detection approaches for specific use cases and in EHR data, in general, will substantially enhance Sentinel's capabilities for ensuring medical product safety but requires special design and analytic methods. | • Evaluation of existing approaches to EHR-based signal detection<br>• Empirical comparison of EHR-based approaches to signal detection in Sentinel<br>• Developing and advancing EHR-based signal detection methods<br>• Advancing methods for safety signal detection for pregnancy and birth outcomes<br>• Developing and evaluating a cancer signal detection tool | • Methodological framework for EHR-based signal detection<br>• General safety signal detection tool for EHR data<br>• Enhanced methods for signal detection for pregnancy and birth outcomes<br>• Tool for cancer safety signal detection |

*slide courtesy of Joshua Gagne*

**Data infrastructure**

- Data partners
- Data elements
- Governance
- Harmonization
- Data quality assurance

**Feature engineering**

- Natural language processing
- Automated feature extraction
- Computable phenotyping

**Causal inference**

- Target trial design
- Advanced, semi-automated analytics
- Subset calibration
- Distributed methods

**Detection analytics**

- Methodological framework
- Statistical methods
- Cancer outcomes
- Pregnancy and birth outcomes

*slide courtesy of Joshua Gagne*

# Variable Importance (struct. + all NLP)

## Top 5 structured:

1. Number of prior years with allergic reaction diagnoses (-)
2. Allergic reaction diagnosis in the prior year (-)
3. Same-day exposure to any imaging procedure (-)
4. Prescription for antihistamines @discharge (-)
5. Prescription for corticosteroids @discharge (-)

## Top 5 NLP-derived:

1. ≥2 affirmative mentions of hypotension
2. Any description of respiratory compromise and reduced BP near a mention of either anaphylaxis as a diagnosis, epinephrine administration, suddenness of onset, or admission for observation
3. ≥2 affirmative mentions of skin/mucosal involvement and either respiratory compromise or reduced blood pressure near anaphylaxis as a diagnosis
4. ≥2 affirmative mentions of wheezing
5. any description of skin/mucosal involvement and reduced blood pressure near a mention of either anaphylaxis as a dx, epinephrine administration, suddenness of onset, or admission for observation

# NLP dictionary: 2. Exploratory query

- Use relational database full-text indexing

- Find Synonyms of "dyspnea"
  - Known: "shortness of breath" and "trouble breathing"
  - Review notes with breath
    - 208 strings yield 5 new terms

| Before_Term | Term | After_Term |
|---|---|---|
| was closing and wheezing and difficulty | breath | ing. She has some mild reactive airway d |
| and throat swelling. Having difficulty | breath | ing and a hard time swallowing saliva. W |
| rhythm.  RESP: Clear to auscultation. | breath | ing comfortably.   Jerico endorses feel |
| like this before. Feels like she cannot | breath | . Cannot swallow. Has not taken anything |
| omplaint: Allergic Reaction; Edema; and | breath | ing Problems      HISTORY AND PHYSICAL E |
| tightening and it was a little hard to | breath | e so comes here for evaluation where she |
| ing      Swelling around eyes, tears, no | breath | ing problems   • Lovastatin    • Sulfa ( |
| en he began to cry and said he couldn't | breath | . He sent Mom a picture of his face- she |
| the first time.  Pt apparently stopped | breath | ing briefly, was given epinephrine and a |

# NLP dictionary: 3. Synonyms

UMLS: Unified Medical Language System – Metathesaurus
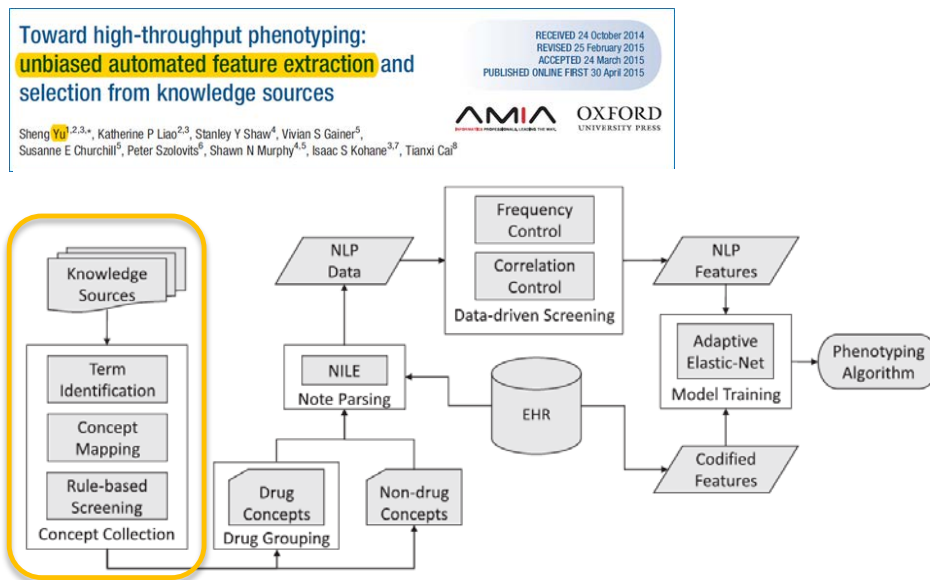


"Dyspnea"

"breathing difficulties"

"DIB"

"difficulty in breathing"

…

# NLP dictionary: Clinical knowledge sources

- 1st step in Yu and colleagues 2015 JAMIA paper "AFEP"



- Important terms will appear in ≥3 clinical knowledge base articles

Yu et al. Toward high-throughput phenotyping: unbiased automated feature extraction and selection from knowledge sources JAMIA 2015;22:993–1000.

# NLP dictionary: Clinical knowledge sources



5 clinical knowledge base articles on the topic anaphylaxis

(+ UpToDate)

367 unique SNOMED terms

90 terms appear in ≥3 sources

# NLP dictionary: Clinical knowledge sources

**90 terms** in the Standard Nomenclature of Medicine, Clinical Terms (SNOMED CT) appeared in at least 3 anaphylaxis knowledge base articles on anaphylaxis.

| Appearing in 5-6 articles | | Appearing in 4 articles | Appearing in 3 articles | |
|---|---|---|---|---|
| Allergens | Blood | Angioedema | Air | Lung |
| Anaphylaxis | Cells[1] | Anxiety | Albuterol | Muscle |
| Diagnosis[1] | Dizziness | Atopy | Antigens | omalizumab |
| Diarrhea | Dyspnea | Basophils | Arteries | Ovum |
| Disease[1] | Exercise | Coughing | Asphyxia | Oxygen |
| Epinephrine | Heart | Edema | Autopsy | Panic |
| Hypersensitivity | Histamine | Esthesia | Chest | Proteins |
| Shock | Hypotension | Flushing | Complication[1] | receptor |
| Skin | Injection | Glucagon | Confusion | Redness |
| Urticaria | Latex | Hoarseness | Congestion | Seizures |
| Venoms | Nausea | Mastocytosis | Extravasation | Services[1] |
| Vomiting | Obstruction | Nose | Eye | Source[1] |
| Wheezing | Pain | Opioids | Gold[2] | Uterus |
| Abdomen | Palpitations | Rhinorrhea | Headache | Vaccines |
| Antibiotics | Pruritus | Stridor | Immunoglobulins | Vancomycin |
| Antibodies | Swelling | Tachycardia | Immunotherapy | Vasodilation |
| Antihistamines | Syncope | Tryptase | Lactams | Veins |
| Aspirin | Tongue | | Larynx | |
| Asthma | | | Lightheadedness | |
| 37 terms ( 13 in 6 and 24 in 5) | | 17 terms | 36 terms | |

[1] Terms unlikely to be useful for distinguishing anaphylaxis cases from non-cases.
[2] "Gold" is an author name appearing in 3 bibliographies (N Engl J Med 2008; 358:28).

# NLP: Feature engineering (manual)

| Diagnostic criteria for anaphylaxis (Sampson/NIAID 2006) | | |
|---|---|---|
| **Sampson Criterion** | **Clinical criteria** | **NLP Features** |
| #1 | Skin/mucosal involvement (SM), *plus either:*<br>    Respiratory compromise (RC) *or*<br>    Reduced blood pressure (BP) | SM+RC<br>SM+BP |
| #2 | Exposure to a likely allergen *for that patient*[1] *plus any 2:*<br>    Skin/mucosal involvement (SM) *or*<br>    Respiratory compromise (RC) *or*<br>    Reduced blood pressure (BP) *or*<br>    Gastrointestinal symptoms (GI) | ~~SM+RC~~[2]<br>~~SM+BP~~[2]<br>SM+GI<br>RC+BP<br>RC+GI<br>BP+GI |
| #3 | Exposure to a known allergen *for that patient*[1] *plus:*<br>    Reduced blood pressure (BP) | None[3] |
| 1. Allergen exposure not operationalized because too difficult to do accurately via NLP.<br>2. This combination not included in criterion #2 because already in criterion #1.<br>3. Not operationalized because w/o allergen exposure reduced BP is non-specific. | | |

Sampson HA, Muñoz-Furlong A, Campbell RL, et al. Second symposium on the definition and management of anaphylaxis: summary report – second national Institute of allergy and infectious disease/food allergy and anaphylaxis network symposium. J Allergy Clin Immunol. 2006;117:391–397.