# Evaluating Use of Methods for Adverse Event Under Surveillance (EUMAEUS)

# Why EUMAEUS?

1) The rapid rollout of COVID-19 vaccines makes it increasingly critical to perform large-scale evaluations of vaccine safety using real-world evidence.

2) Estimate the comparative performance (bias, precision, timeliness) of the case-control, cohort, historical rate, and self-controlled methods for vaccine safety.

# Literature Review

Lana Lai

on behalf of the EUMAEUS task force

# Types of Study Designs

| Study Design | Description | Advantages | Disadvantages | Clinical Applications |
|---|---|---|---|---|
| Case-control | • Comparison of cases vs. non-cases from the same source population from the same time-period | • Uses small data sample from entire cohort, cost efficient<br>• Use matching to control for time-varying confounders | • Confounding by indication<br>• Selection bias<br>• Misclassification of exposure | • Autism spectrum disorders & various vaccines<br>• Inflammatory bowel disease (IBD) & MMR vaccine<br>• Guillain-Barré syndrome (GBS) & H1N1 vaccine |

# Types of Study Designs

| Study Design | Description | Advantages | Disadvantages | Clinical Applications |
|---|---|---|---|---|
| Cohort | • Comparison of incidence rate ratio of adverse events between vaccinated vs. unvaccinated population | • Easy to implement – large amount of data available<br>• Use matching / stratification to control for potential confounders | • Confounding by indication<br>• Misclassification of exposure | • Intussusception & rotavirus vaccine<br>• Autism spectrum disorders & various vaccines |
| Historical Rate (Comparator) Cohort | • Comparison between observed incidence of adverse events vs. expected incidence based on historical data | • Greater statistical power to detect rare adverse events<br>• Improved timeliness in detecting potential safety signals | • Temporal confounders (e.g. seasonality, changing trends in detection of adverse events & variation in diagnostic criteria over time) | • Pediatric vaccines<br>• Tdap vaccine<br>• HPV vaccine<br>• H1N1 vaccine |

# Types of Study Designs

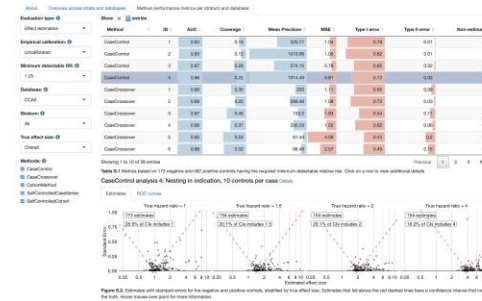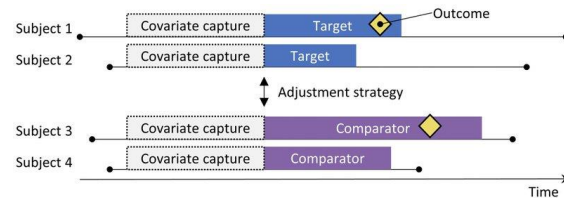| Study Design | Description | Advantages | Disadvantages | Clinical Applications |
|---|---|---|---|---|
| Self-Controlled Case Series (SCCS) / Self-Controlled Risk Interval (SCRI) | • Comparison between incidence rates in exposed time periods vs. incidence rates of self-matched unexposed time periods<br>• SCCS: Cases only<br>• SCRI: Vaccinated population only | • Adjust for time-invariant confounders<br>• SCCS: Multiple occurrences of independent events within an individual can be assessed<br>• SCRI: Less susceptible to misclassification of exposure | • Time-varying confounding (e.g. age, seasonality)<br>• Reverse causality bias | • Guillain-Barré syndrome (GBS) & H1N1 vaccine<br>• Autism spectrum disorders & various vaccines<br>• Seizures & various vaccines |

# Overview of the
# EUMAEUS Experiment Design

Marc Suchard

on behalf of the EUMAEUS task force

**OHDSI**

Observational Health Data Sciences and Informatics

# EUMAUES is an empirical benchmark study

Builds on our prior work evaluation of comparative (drug) effectiveness and safety methods published in *Harvard Data Science Review*



To systematically evaluate the

**performance of methods**

to reliably

**identify vaccine safety signals**

in

**real-world settings**

## How Confident Are We About Observational Findings in Health Care: A Benchmark Study

Martijn J. Schuemie, M. Soledad Cepede, Marc A. Suchard, Jianxiao Yang, Yuxi Tian Alejandro Schuler, Patrick B. Ryan, David Madigan[1], George Hripcsak

[1]Professor of Statistics, Columbia University

# Vaccine safety surveillance methods

Reduce systematically to **four** components:

- Construction of a *counterfactual* ("expected count" without vaccination)

- A *time-at-risk* when safety events can occur

- The *test-statistic* to estimate, and

- A *decision rule* to classify signals from non-signals

# Counterfactual construction

- ## Case-control
  - How often are patients with events vaccinated?

- ## Contemporary non-user comparator cohort method
  - How often do events occur to similar unvaccinated patients?
  - Some variants: **anchoring** (or *not*) on healthcare visit; **matching** (or *not*) on age + sex

- ## Historical rates
  - How often did events occur to other patients in the past?
  - Some variants: **anchoring**; **stratifying** (or *not*) on age + sex

- ## Self-control case series
  - How often did/do events occur in the same patients at different times?

Note: 17 total variations drawn from the literature
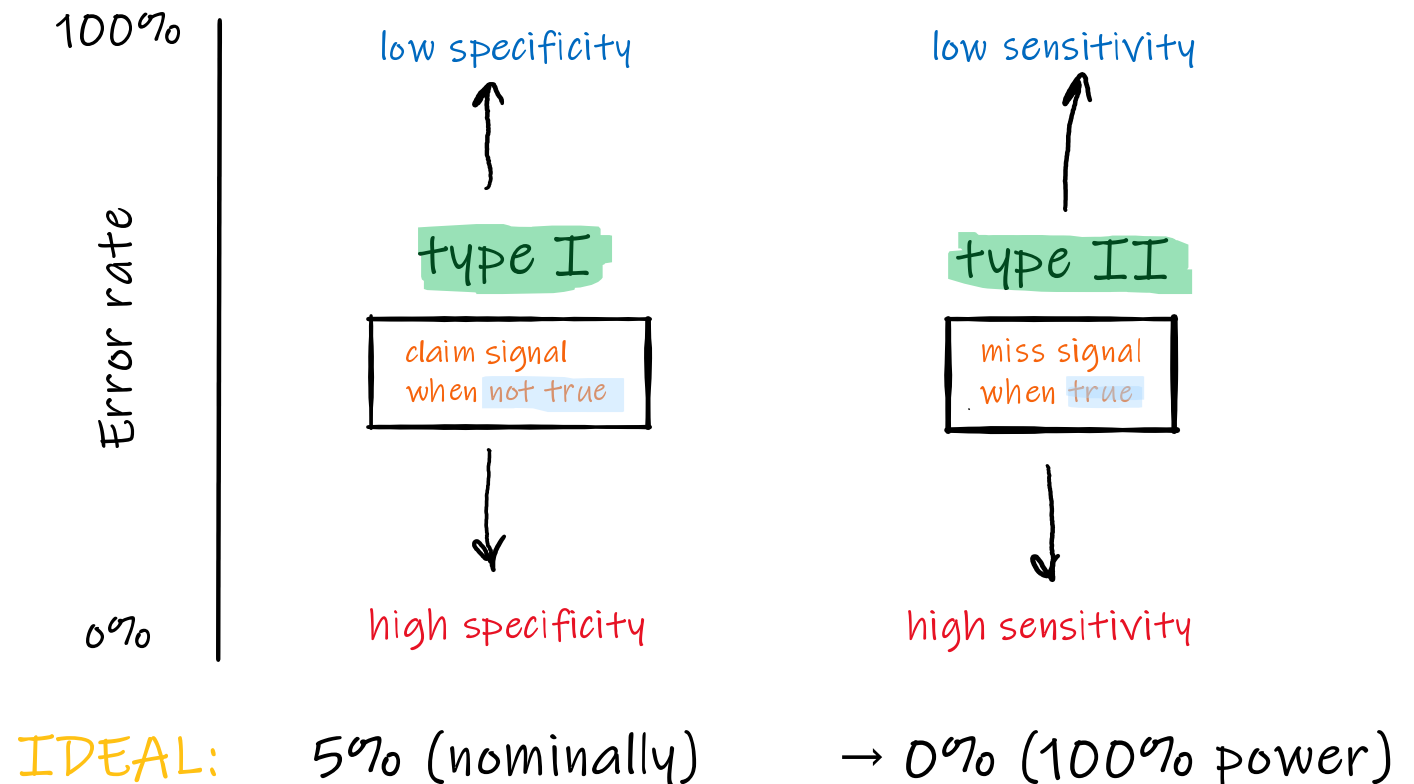
# Time at risk and test-statistics

- A ***time-at-risk*** when safety events can occur:
  - 0-1 days, ***1-28 days*** and 1-42 days after vaccination
  - Dose definition (first, second, both)

- The ***test-statistic*** to estimate:
  - Effect-size estimation (incidence rate ratio, hazard ratio or odds ratio)
  - Log-likelihood ratio (common in vaccine surveillance, allows for corrections for multiple testing over time via ***MaxSPRT***)
  - With and without ***empirical calibration*** (to control for systematic error)

# Method performance metrics

- A **decision rule** to classify signals from non-signals
  - Bias / variance (particularly of the **residual systematic error**)
  - Type 1 error rate
  - Type 2 error rate
  - **Timeliness** to achieve power



100%

Error rate

0%

low specificity

type I

claim signal when not true

high specificity

low sensitivity

type II

miss signal when true

high sensitivity

IDEAL:   5% (nominally)        → 0% (100% power)

# Real-world evidence with 117M estimates

**Exposures of interest:**

- H1N1pdm (`09-`10)
- Seasonal influenza (Fluvirin, `17-`18)
- Seasonal influenza (Fluzone, `17-`18)
- Seasonal influenza (all, `17-`18)
- Zoster (2018, 2 doses)
- HPV (2018, 2 doses)

**Data sources:**

- CCAE
- MDCR
- MDCD
- Optum EHR

**Negative control outcomes (93):**

- Not related to any of these vaccines
- Similar prevalence and %-inpatient diagnoses (severity) to AESI
- Clinical expert review

**Positive control outcomes:**

- Imputed from negative controls
- Known effect sizes (1.5, 2, 4 x)

***Open Science***: pre-specified and registered protocol, open-source analytic code, public access to all results

- https://ohdsi-studies.github.io/Eumaeus/Protocol.html
- https://github.com/ohdsi-studies/Eumaeus/
- https://data.ohdsi.org/Eumaeus/

# Prelude to the results

- Which methods are *least bias* in the real-world?
  - Effect of counterfactual anchoring
  - Effect of confounding adjustment
- What is the *trade-off* to achieve, say, 50% power?
- Should we *combine multiple designs* (signal generation / evaluation) to improve performance?
- Is *sequential testing ($\alpha$-spending) correction* a panacea?
- Do *2$^{nd}$ doses* influence method choice?

# Bias, precision and timeliness of historical rate comparison methods

### Xintong Li

on behalf of the EUMAEUS task force

Recall the advantages of historical comparator design:
- Greater statistical power
- Improved timeliness

Especially useful at early stage after vaccine introduction

Historical comparator is from:
- literature
- **within same database / population (best-case scenario)**
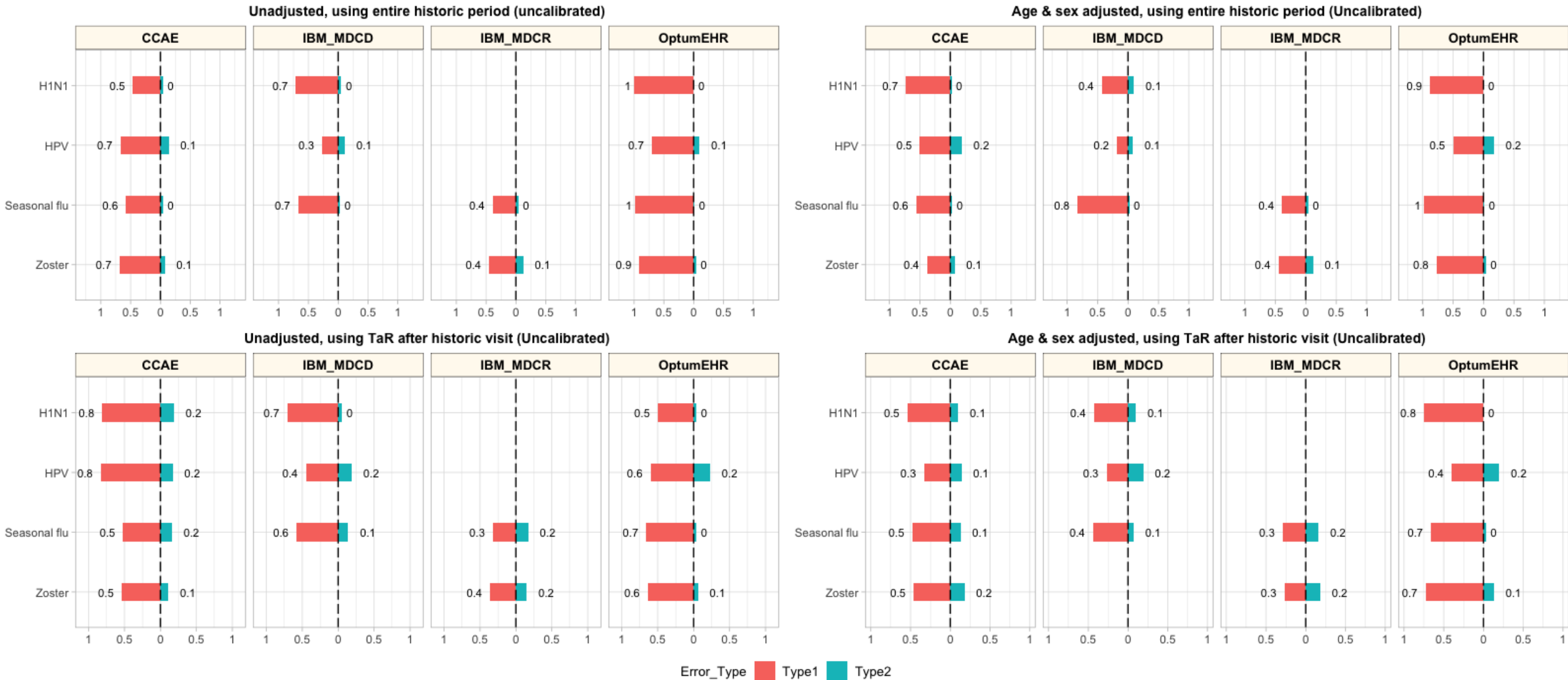- others

# Choice of design

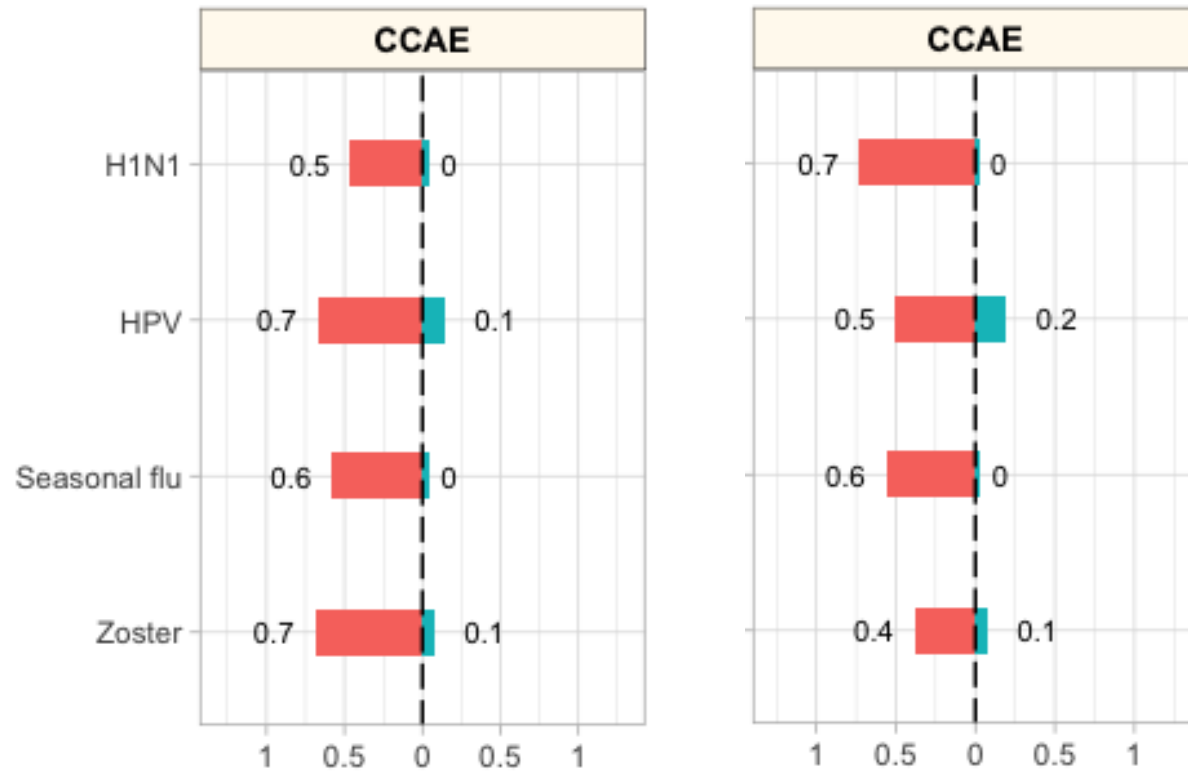| Population | Time-at-risk | Calibration |
|---|---|---|
| Unadjusted | Entire year | Yes |
| Adjusted for age and sex | Relative to outpatient visit | No |

# Historical comparison in general: Sensitive but not specific

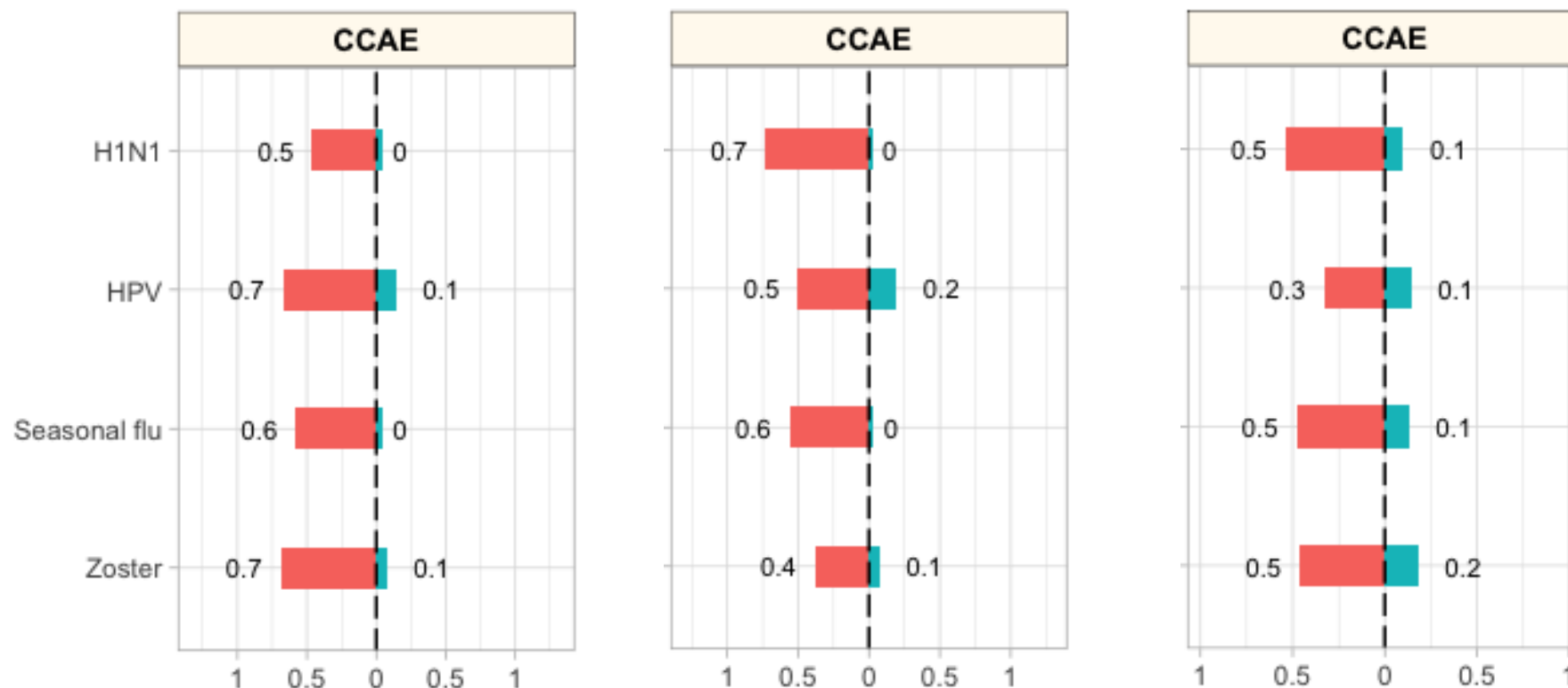# Sensitive but not specific



Unadjusted,
entire historical period

Age and sex adjusted,
entire historical period

Adjust for age and sex reduced type 1 error.

# Sensitive but not specific



Unadjusted,
entire historical period

Age and sex adjusted,
entire historical period

Age and sex adjusted,
Time-at-risk after
historic visit

After adjusting for age and sex, anchoring on visit further reduce type 1 error.

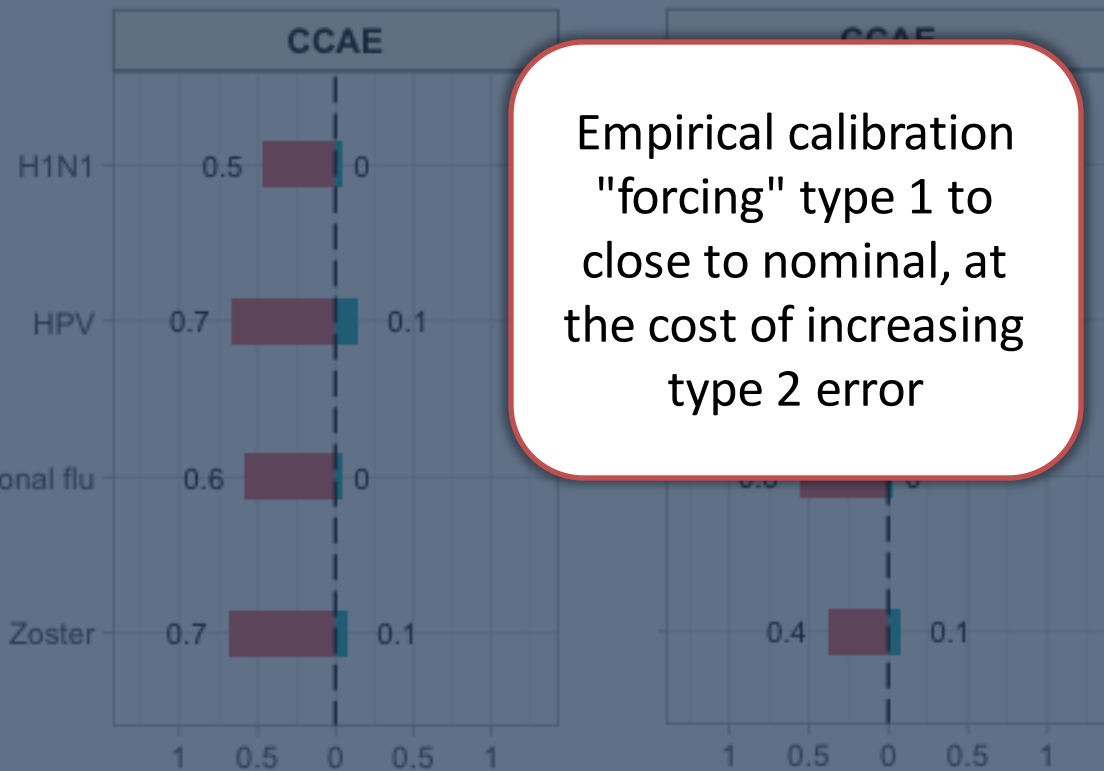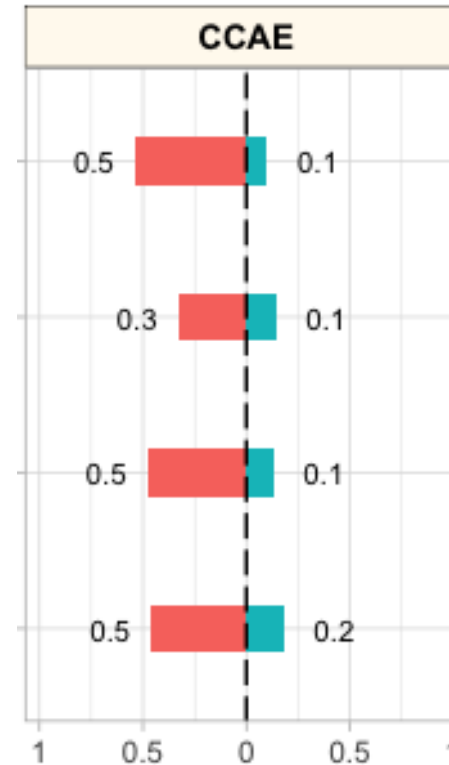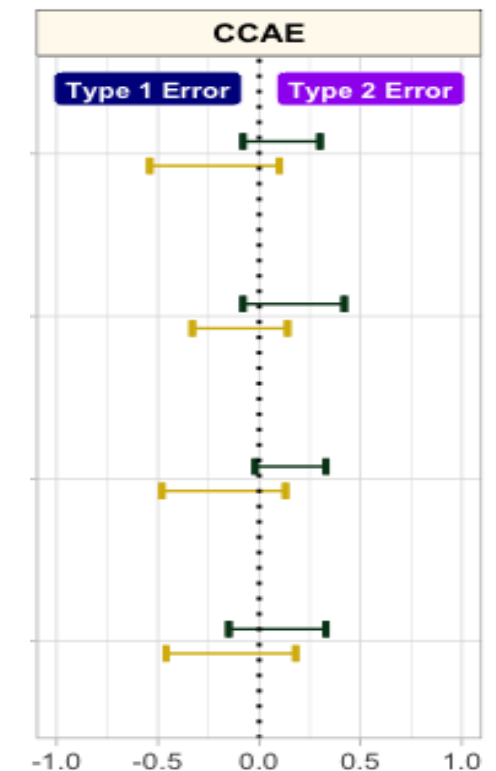Empirical calibration "forcing" type 1 to close to nominal, at the cost of increasing type 2 error

Unadjusted, entire historical period

Age and sex adjusted, entire historical period

Age and sex adjusted, Time-at-risk after historic visit

Age and sex adjusted, Time-at-risk after historic visit

21

# Empirical calibration: reduce type 1, increase type 2

# Higher and faster uptake, earlier detection



**H1N1 vaccination**

**Seasonal flu vaccination (All)**

**Zoster vaccination (Shingrix)**

**HPV vaccination (Gardasil 9)**

Database
CCAE

Analysis
Adjusted for age and sex, no anchoring

Calibration
No

True effect size: ● True effect size = 1 ● True effect size = 2 ● True effect size = 1.5 ● True effect size = 4   p < 0.05  ○ FALSE ● TRUE

# Conclusion

- Sensitive but not specific: overestimate risks

- Age-sex adjustment reduce false positive

- Anchoring on visit reduce false positive

- Empirical calibration: forced type 1 error back to normal, at the cost of increasing type 2 error.

- For vaccine with high uptake speed: can detect earlier, stabilized estimation.

# Combining Methods in a Safety Surveillance System

Faaizah Arshad

on behalf of the EUMAEUS task force

# Introduction

| Sensitive method | → | Specific method | → | High sensitivity & specificity |

- ## HIV testing
  - Two part test: 1) highly sensitive (few false negatives); 2) highly specific (eliminate false positives)
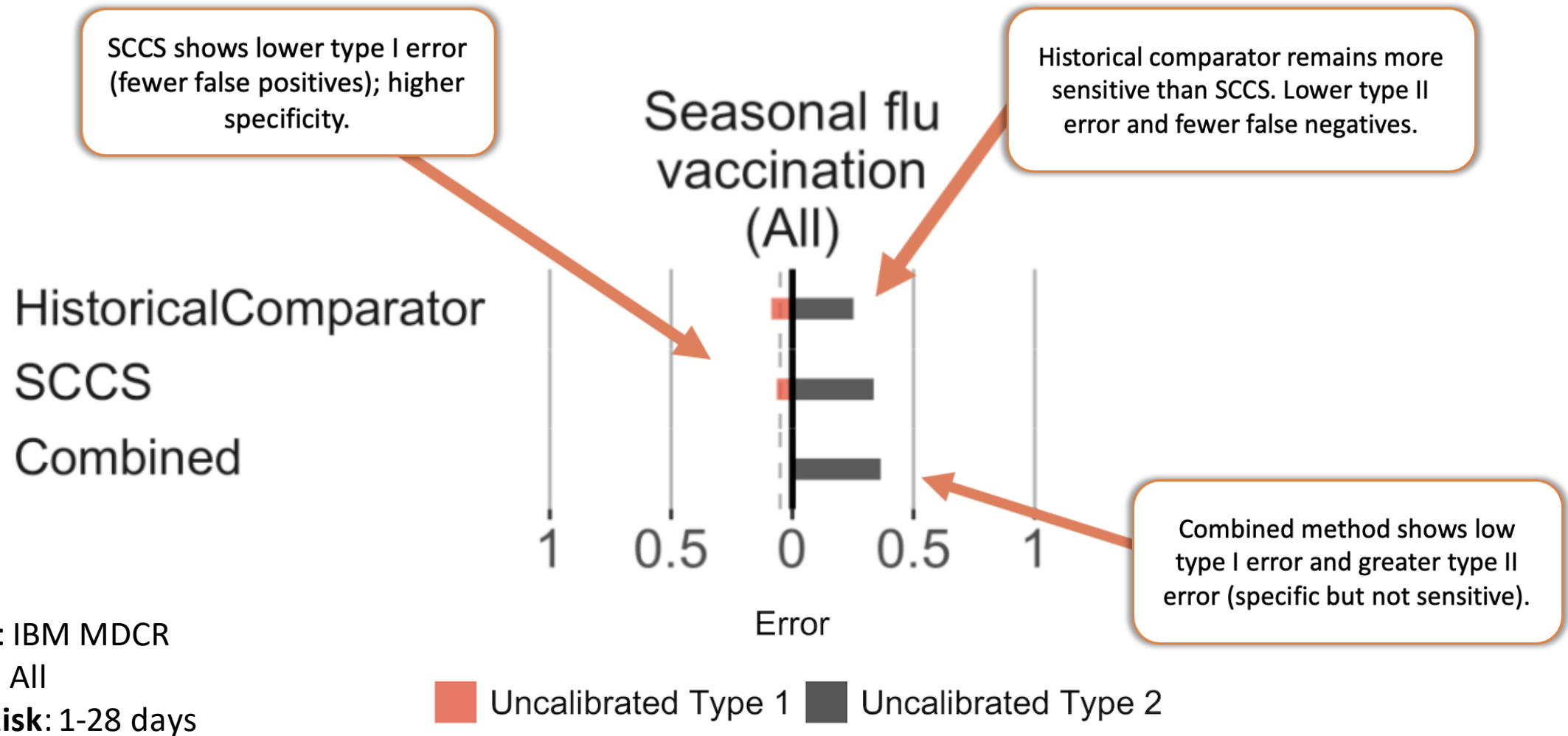
# Methods

- We hypothesized that sequentially combining methods might be desirable for population-level COVID-19 vaccine safety surveillance.

- Method 1: historical comparator (sensitive / cheap)

- Method 2: self-controlled case series (specific)

- Combined: Method 1 → Method 2

# Uncalibrated type I and II errors for all outcomes



SCCS shows lower type I error (fewer false positives); higher specificity.

Historical comparator remains more sensitive than SCCS. Lower type II error and fewer false negatives.

Combined method shows low type I error and greater type II error (specific but not sensitive).

Seasonal flu vaccination (All)

HistoricalComparator
SCCS
Combined

1    0.5    0    0.5    1

Error

**Database**: IBM MDCR
**Outcome**: All
**Time-at-Risk**: 1-28 days
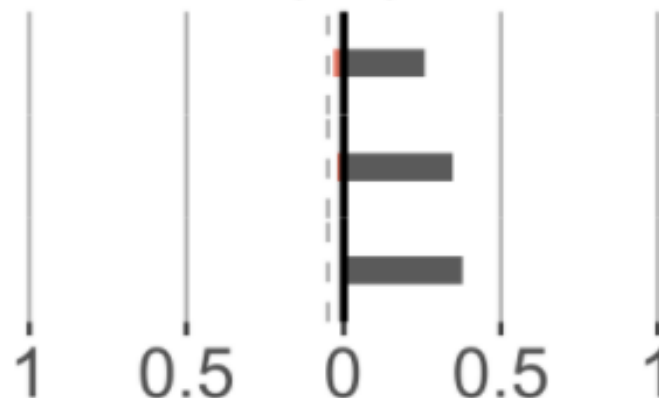
Uncalibrated Type 1    Uncalibrated Type 2

# Calibrated type I and II errors for all outcomes

Calibration tries to fix the type I error rate (closer to nominal); most noticeable for historical comparator.

After calibration, historical comparator still most sensitive. Reduced type I error for historical comparator and SCCS.



Seasonal flu vaccination (All)

HistoricalComparator
SCCS
Combined

1   0.5   0   0.5   1

Error

Calibrated Type 1        Calibrated Type 2

**Database**: IBM MDCR
**Outcome**: All
**Time-at-Risk**: 1-28 days

# Conclusion

- Reject hypothesis.

- Sequentially combining sensitive and specific methods does not improve performance over using a single method.

- Future vaccine monitoring should consider the sequence of methods used to ensure accurate signal detection.

# Estimation for Two-Dose Vaccines

## Ty Stanford

on behalf of the EUMAEUS task force

Aim:

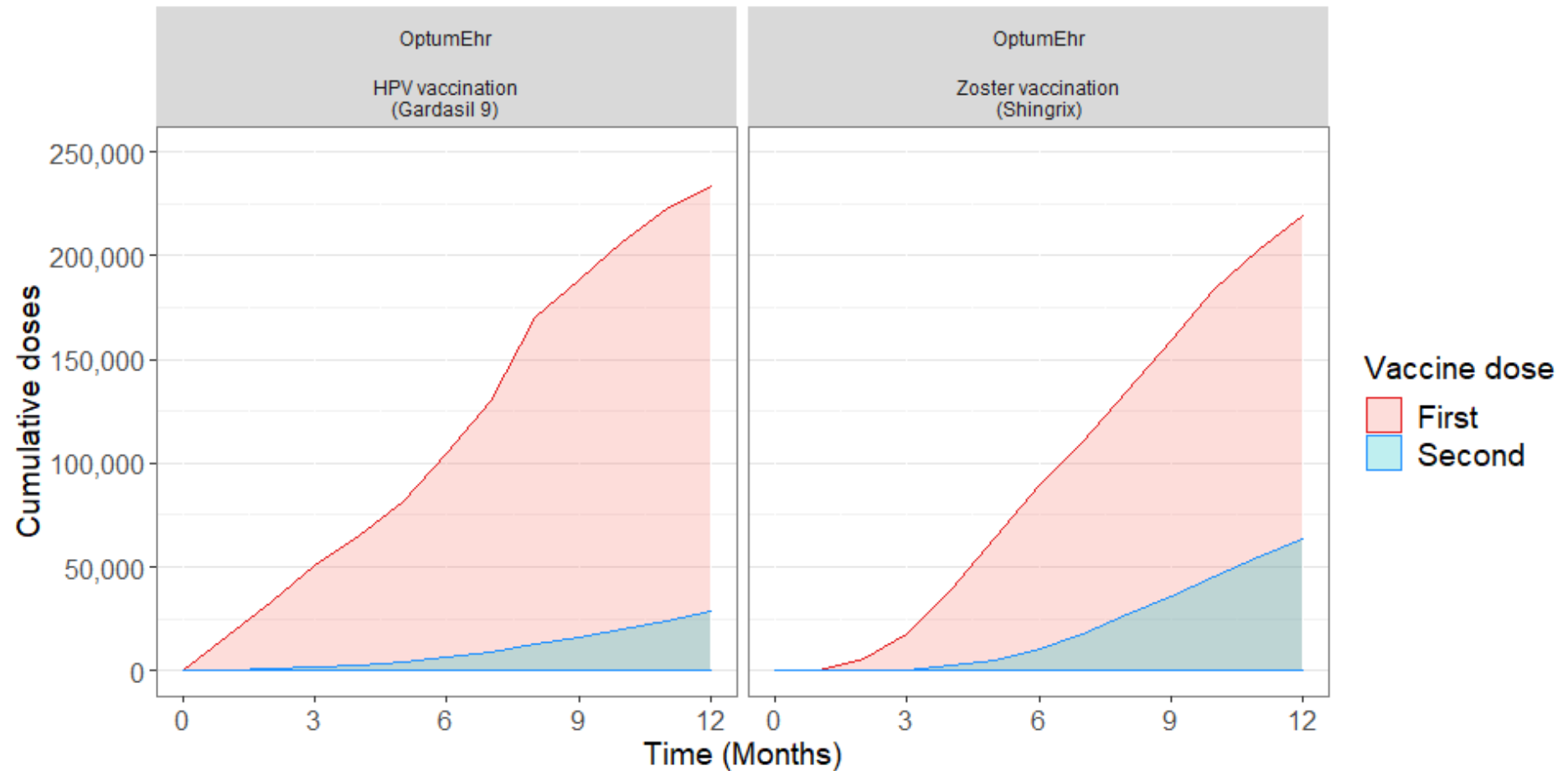- Does the inclusion of data from the 2nd dose, among vaccines with 2 doses, reduce type II error?

Data:

- This limits EUMAEUS data to

  (CCAE, Optum EHR) x (HPV vaccine, Zoster vaccine) combinations

# Dose accumulation



| Database | Dose | HPV vaccination (Gardasil 9) | Zoster vaccination (Shingrix) |
|---|---|---|---|
| Optum EHR | First | 233985 | 219665 |
| | Second | 28336 | 63464 |

# To calibrate or not to calibrate?

Uncalibrated

Calibrated

Dose

1st only

# Adding 2nd dose: Cohort Design

Dose
(1st only) vs (1st & 2nd)

### Uncalibrated

### Calibrated

# Adding 2nd dose: SCCS

Dose

(1st only) vs (1st & 2nd)



Uncalibrated

Calibrated

# Conclusion

- Inclusion of the 2$^{nd}$ dose can increase the power
  - marginally in this case, likely as a result of a marginal increase in sample size
- The most important factor is *empirical calibration*
  - more data doesn't magically negate issues with specific designs
- Future work to understand the issues better:
  - Larger proportion of 1$^{st}$ doses to also have 2$^{nd}$ doses (with differing rates)
  - Underlying signals (positive controls) to have varying effects after each dose

# Comparison of performance across methods

Martijn Schuemie

on behalf of the EUMAEUS task force

# Same data & question, different methods: different results

| Analysis choices | Effect size (95% CI) |
| --- | --- |
| Matched case-control | 3.33 (0.64-15.44) |
| PS-weighted cohort method | 4.42 (1.83-10.42) |
| Historical comparator | 3.50 (1.92-5.87) |
| SCCS | 1.07 (0.59-1.81) |



Estimated effect size

**Exposure**
H1N1pdm vaccinations

**Outcome**
Contusion of toe

**Database**
Optum EHR

39

# Comparing on type 1 and type 2 error



Case-control and historical comparator tend to generate many false positives

# Empirical calibration: restoring type 1



H1N1 vaccination

Type 1 | Type 2

Matched case-control

PS-weighted cohort method

Historical comparator

SCCS

Error

Legend:
- Uncalibrated Type 1
- Uncalibrated Type 2
- Calibrated Type 1
- Calibrated Type 2

Calibration makes methods comparable.

After calibration (fixed type 1), SCCS has lowest type 2 error

# Adjusting for systematic error and sequential testing

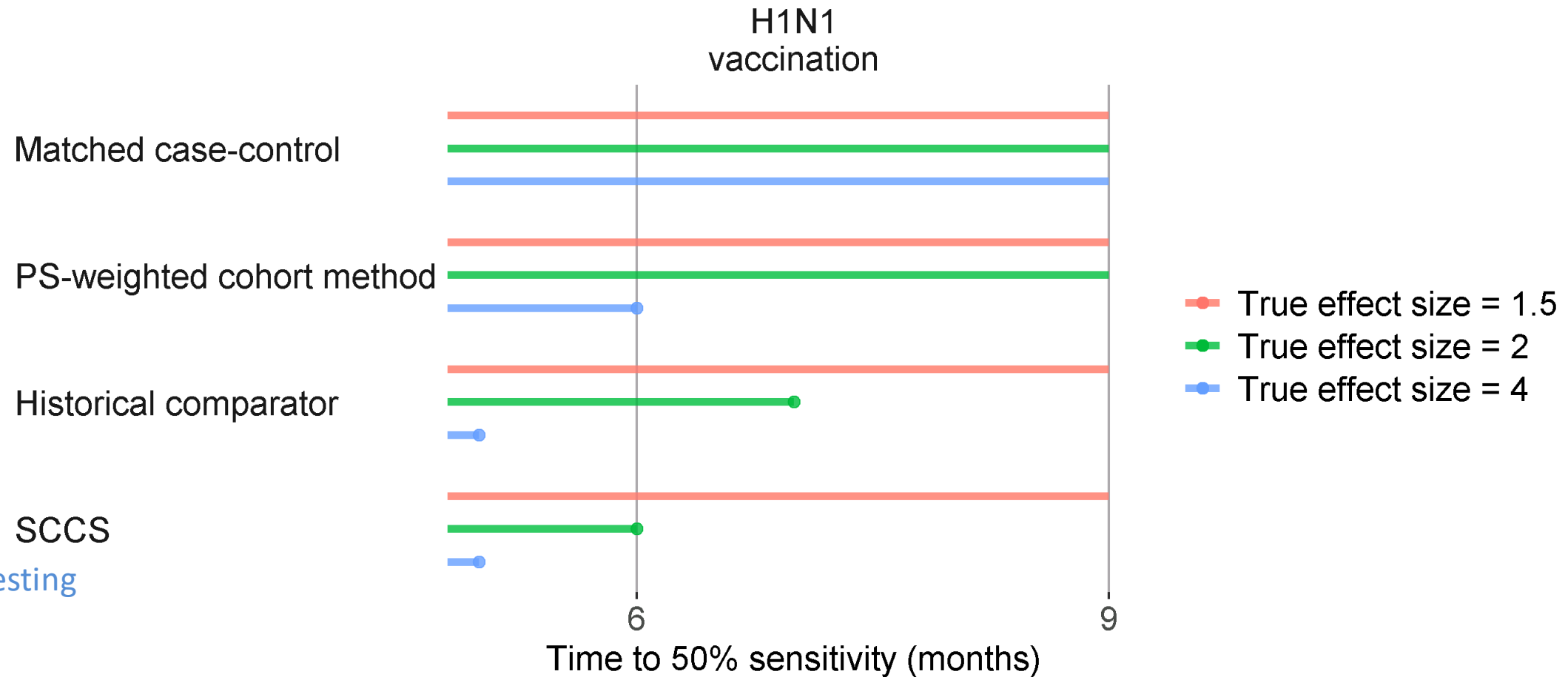| | Type 1 error | |
| :--- | :---: | :---: |
| | **Historical comparator** | **SCCS** |
| **Uncalibrated, no adjustment for sequential testing** | 28.0% | 4.3% |
| **Uncalibrated, MaxSPRT** | 18.3% | 2.2% |
| **Calibrated, no adjustment for sequential testing** | 10.8% | 5.4% |
| **Calibrated, MaxSPRT** | 6.5% | 4.3% |

**New!**

**Exposure**
H1N1pdm vaccinations

**Outcome**
All negative controls

**Database**
Optum EHR

Adjusting for **systematic error**
has bigger impact than
adjusting for **sequential testing**

# Time to 50% sensitivity (after calibration)

H1N1
vaccination

Matched case-control

PS-weighted cohort method

Historical comparator

SCCS

True effect size = 1.5
True effect size = 2
True effect size = 4

6                    9

Time to 50% sensitivity (months)

Adj. for sequential testing
MaxSPRT

Adj. for systematic error
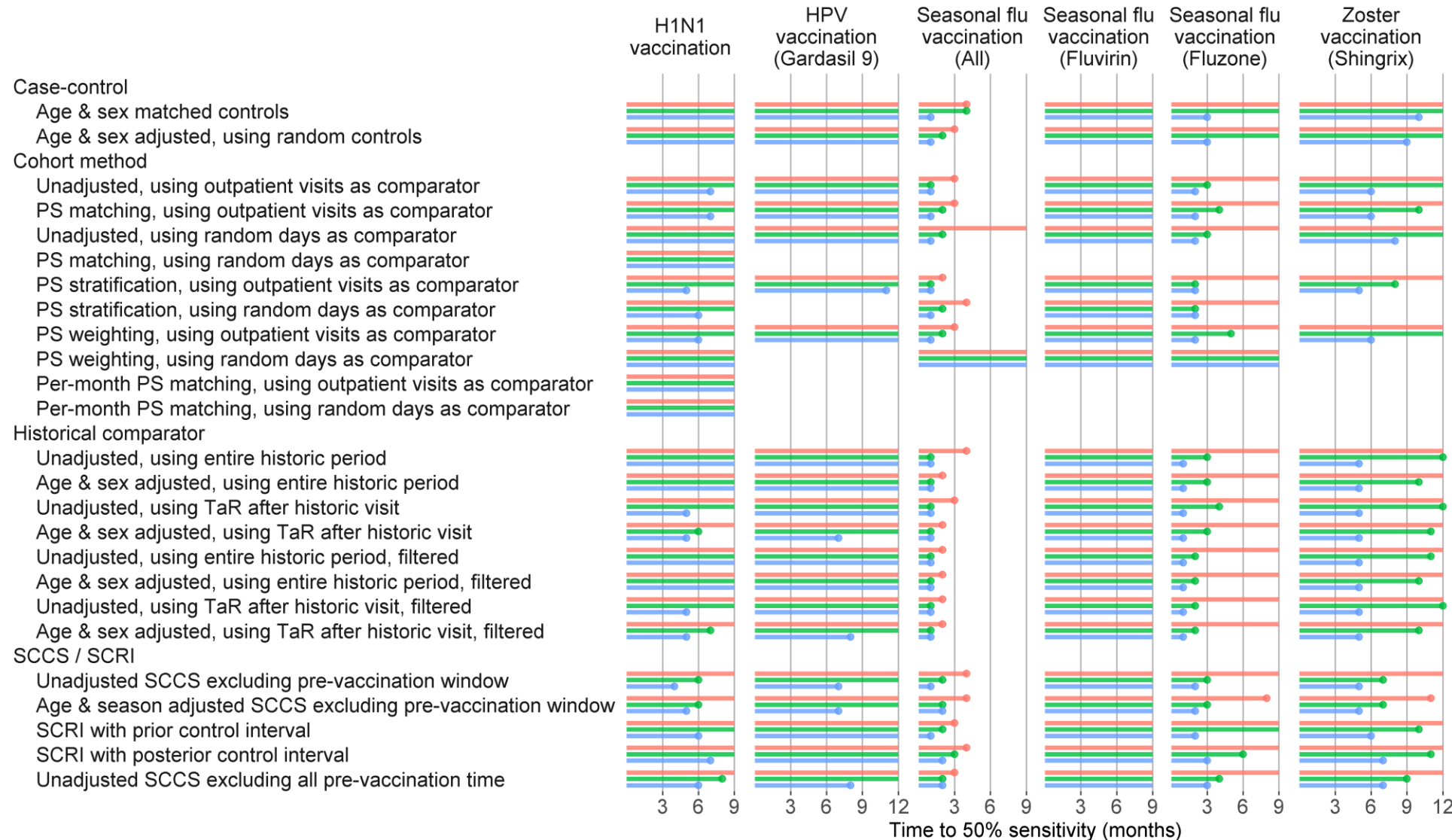Empirical calibration

Database
Optum EHR

SCCS has shortest time to detection

All methods struggle to achieve 50% sensitivity
for small effects

43

# More or less consistent across methods / outcomes /databases

Adj. for sequential testing
MaxSPRT

Adj. for systematic error
Empirical calibration

Database
Optum EHR

# Conclusions

- Many methods show large systematic error / type 1 error

- Empirical calibration can restore type 1 error to nominal, at the cost of increasing type 2 error
  (depending on magnitude of systematic error)

- Empirical calibration often has bigger impact than adjusting for sequential testing
  (should do both)

- After calibration and adj. for sequential testing SCCS seems overall best
  (shortest time to detection)

- No method achieves high sensitivity for small true effect sizes
  (on these data)

# Recommendations for a safety surveillance system

## Martijn Schuemie

on behalf of the EUMAEUS task force

# Recommendations

- Many methods (e.g. case-control & historical comparator) have positive bias, causing many false positives (high type 1 error)
  - Include negative controls and use empirical calibration
  - Include self-controlled designs
  - Always use confounding adjustment
  - Carefully consider anchoring of counterfactual
- Detecting more than half of true adverse effects may require accepting more false positives (e.g. using calibrated $p < 0.10$)
- Combining multiple designs likely doesn't improve performance
  - Do not distinguish between 'signal generation' and 'signal evaluation'
- Second dose often underpowered to contribute to evidence