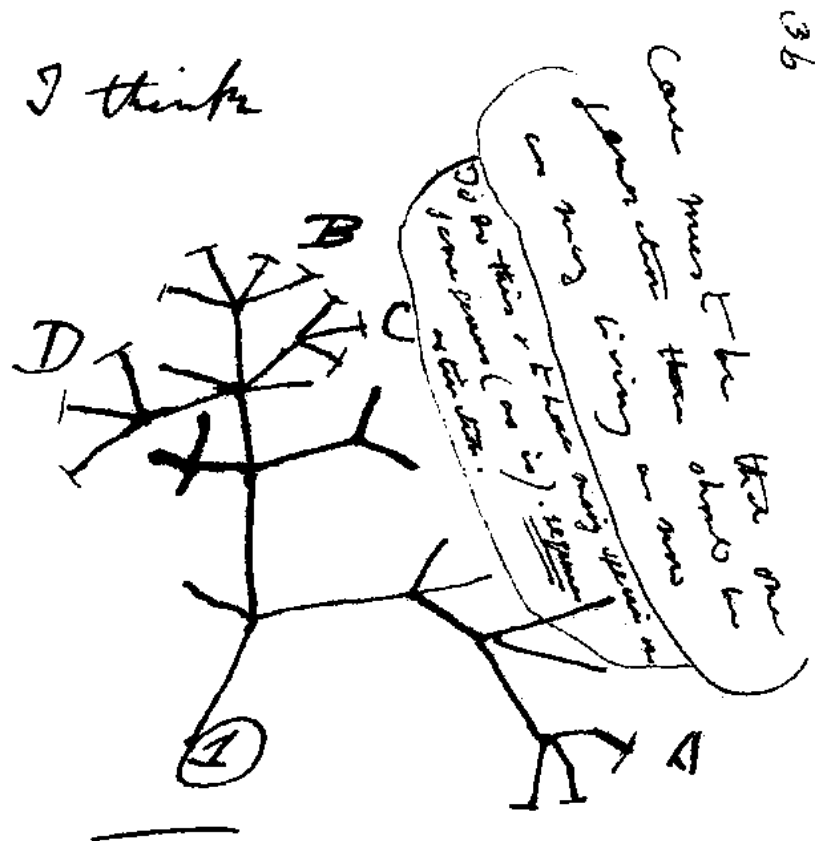# Biomedical Discovery through Data Mining and Data Science

November 14th, 2016
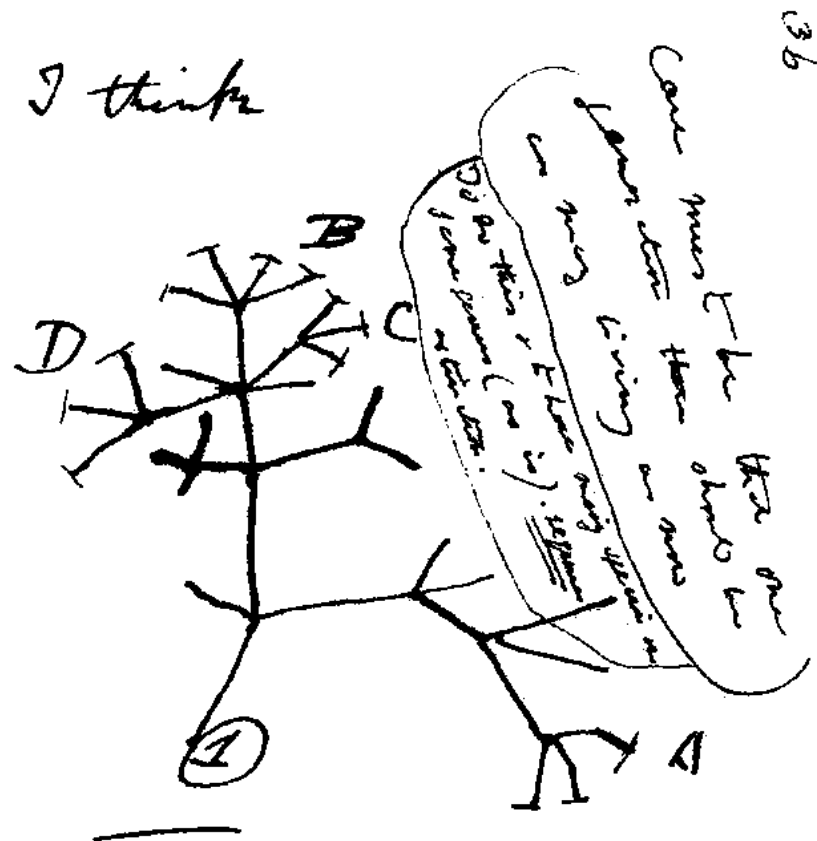
Nicholas P. Tatonetti, PhD
Columbia University

# Observation is the starting point of biological discovery

# Observation is the starting point of biological discovery



- Charles Darwin observed relationship between geography and phenotype

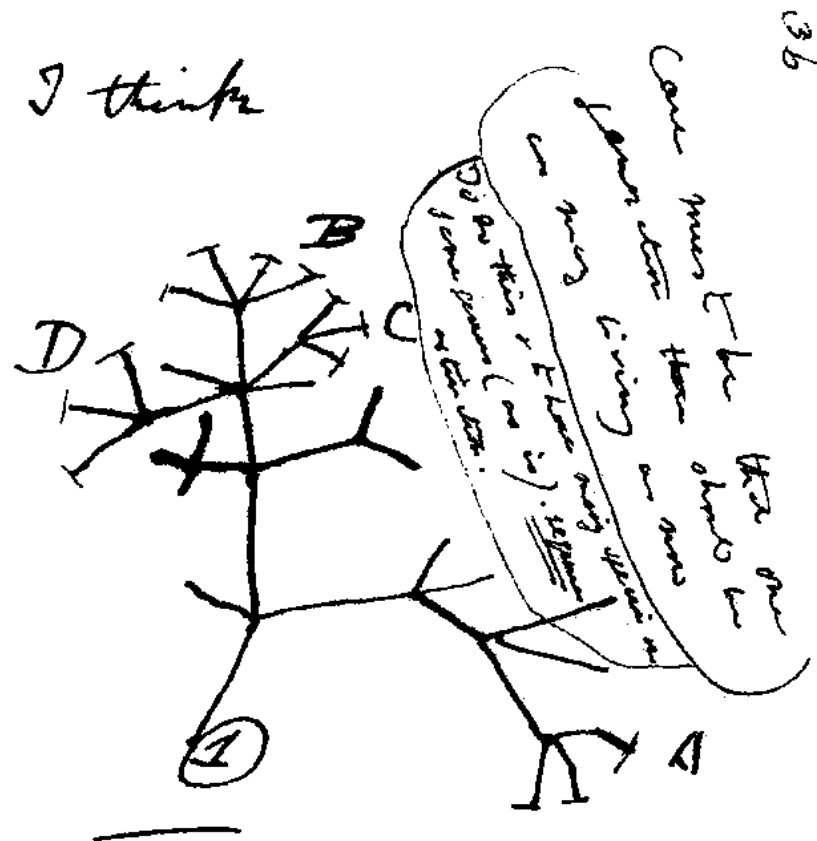# Observation is the starting point of biological discovery

- Charles Darwin observed relationship between geography and phenotype

- William McBride & Widukind Lenz observed association between thalidamide use and birth defects

# The tools of observation are advancing

# The tools of observation are advancing

- Human senses

# The tools of observation are advancing

- Human senses

  - sight, touch, hearing, smell, taste

# The tools of observation are advancing

- Human senses

  - sight, touch, hearing, smell, taste

- Mechanical augmentation

# The tools of observation are advancing

- Human senses

    - sight, touch, hearing, smell, taste

- Mechanical augmentation

    - binoculars, telescopes, microscopes, microphones

# The tools of observation are advancing

- Human senses

  - sight, touch, hearing, smell, taste

- Mechanical augmentation

  - binoculars, telescopes, microscopes, microphones

- Chemical and Biological augmentations

# The tools of observation are advancing

- Human senses

  - sight, touch, hearing, smell, taste

- Mechanical augmentation

  - binoculars, telescopes, microscopes, microphones

- Chemical and Biological augmentations

  - chemical screening, microarrays, high throughput sequencing technology

# The tools of observation are advancing
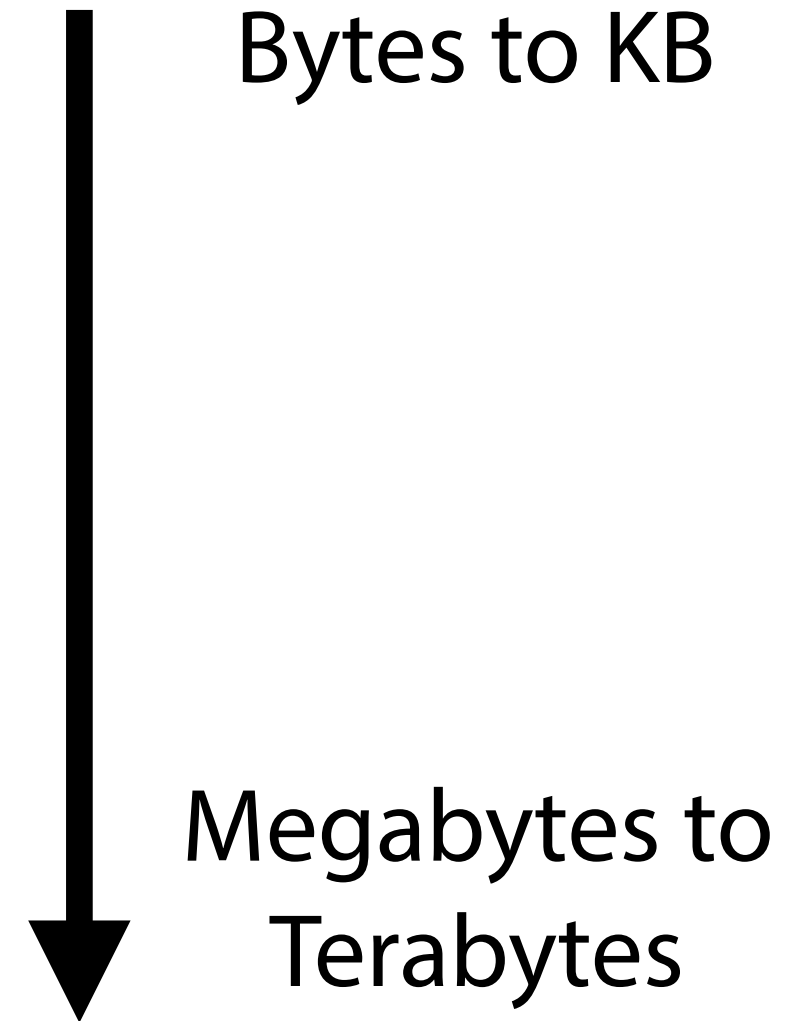
- Human senses

    - sight, touch, hearing, smell, taste

- Mechanical augmentation

    - binoculars, telescopes, microscopes, microphones

- Chemical and Biological augmentations

    - chemical screening, microarrays, high throughput sequencing technology
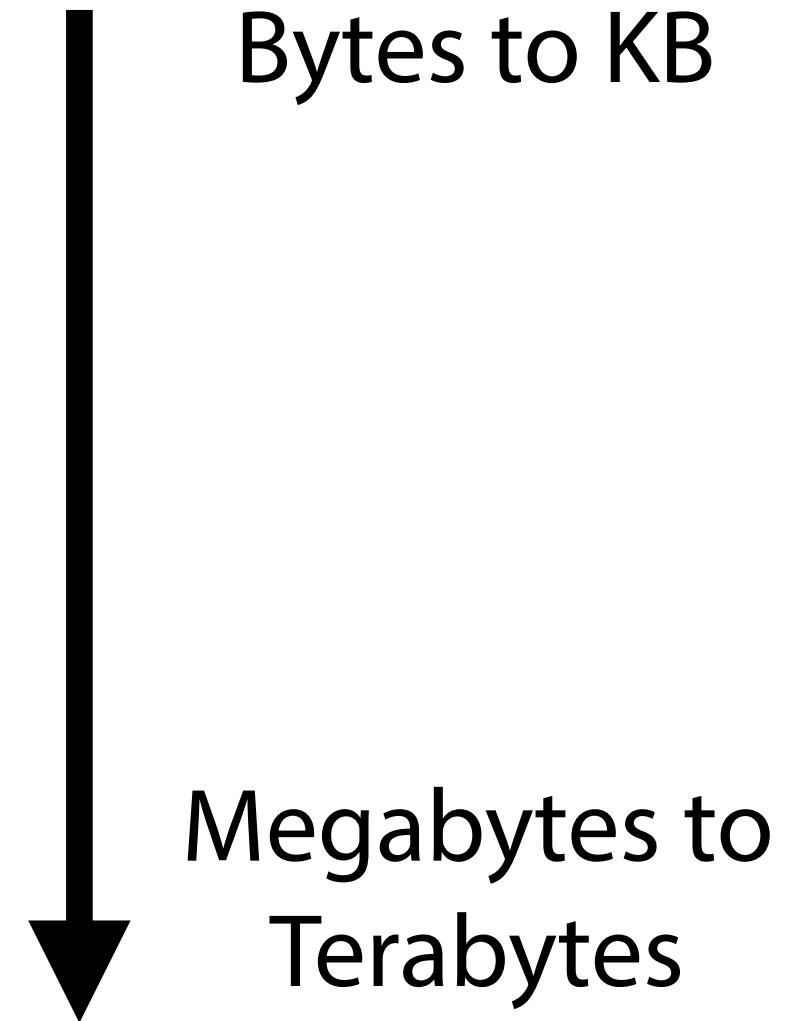
Bytes to KB

Megabytes to Terabytes

# The tools of observation are advancing

- Human senses

  - sight, touch, hearing, smell, taste

- Mechanical augmentation

  - binoculars, telescopes, microscopes, microphones

- Chemical and Biological augmentations

  - chemical screening, microarrays, high throughput sequencing technology
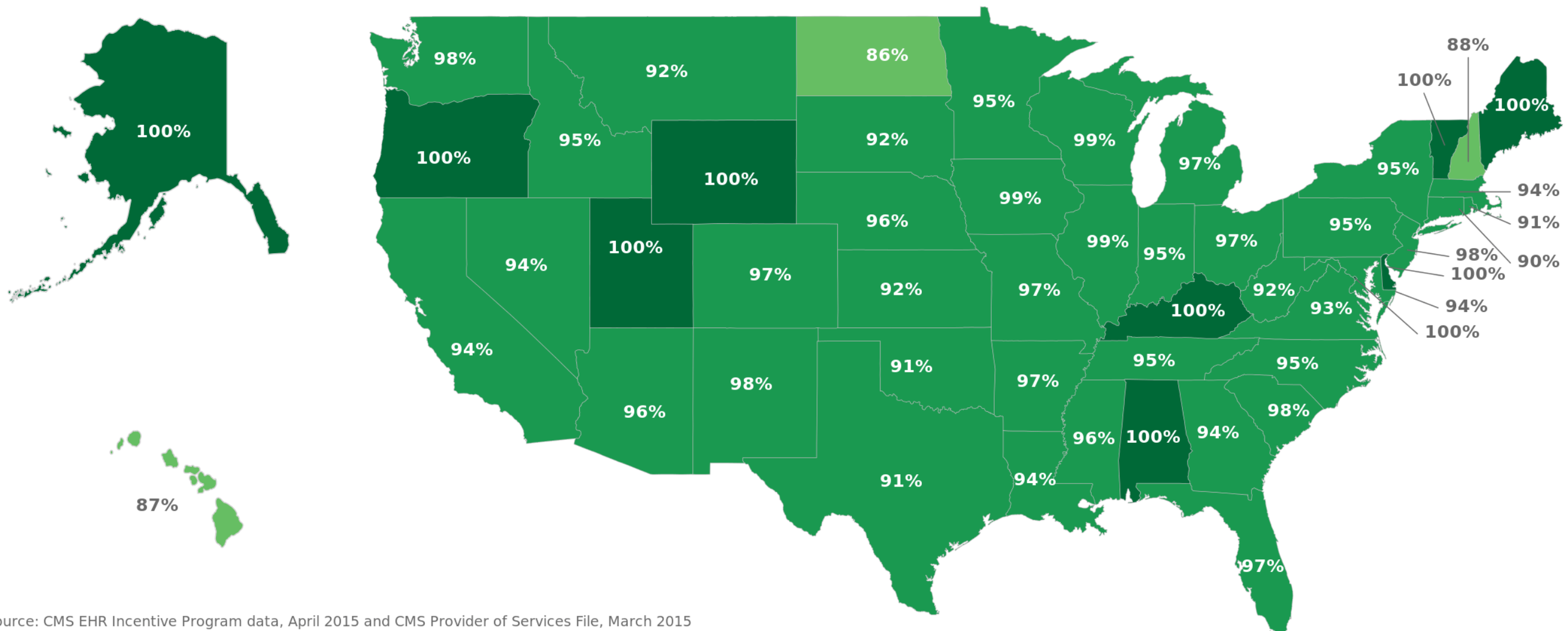
- What's next?

Bytes to KB

Megabytes to Terabytes

# Your doctor is observing you like never before

## >99% of Hospitals have Electronic Health Records

# Every drug order is an experiment.

# Observation analysis in a *peta*byte world

# Observation analysis in a *peta*byte world

- Darwin, McBride, and Lenz were working with *kilo*bytes of data

# Observation analysis in a *peta*byte world

- Darwin, McBride, and Lenz were working with *kilo*bytes of data

- Today's scientists are observing *tera*bytes and *peta*bytes of data

# Observation analysis in a *peta*byte world

- Darwin, McBride, and Lenz were working with *kilo*bytes of data

- Today's scientists are observing *tera*bytes and *peta*bytes of data

- The human mind simply cannot make sense of that much information
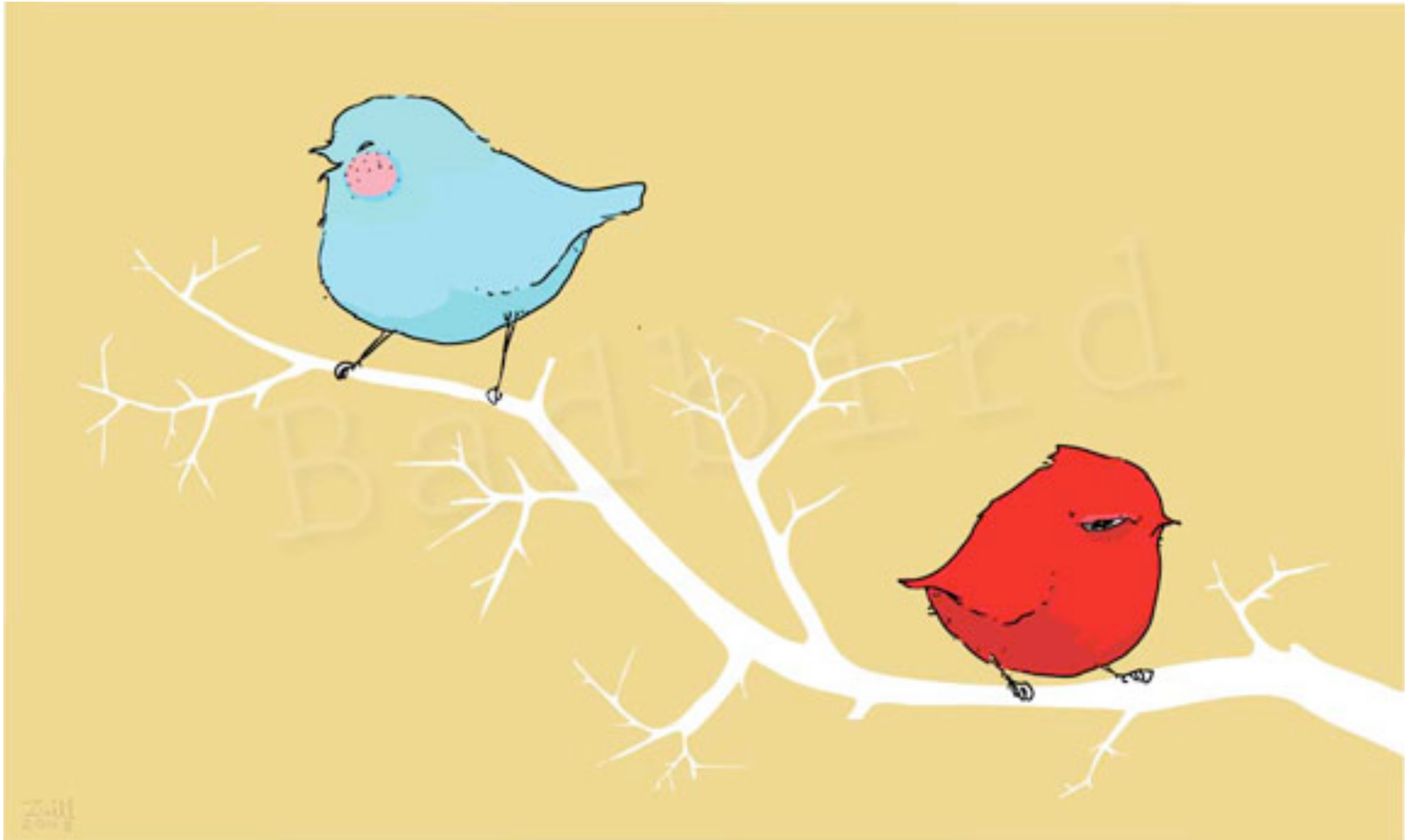
# Observation analysis in a *peta*byte world

- Darwin, McBride, and Lenz were working with *kilo*bytes of data

- Today's scientists are observing *tera*bytes and *peta*bytes of data

- The human mind simply cannot make sense of that much information

- Data mining is about making the tools of data analysis ("hypothesis generation") catch up to the tools of observation

But, there's a problem…

# Bias confounds observations

Let's focus on just one example...

Let's focus on just one example...

**Drug-Drug Interactions**

# Drug-drug interactions (DDIs)

# Drug-drug interactions (DDIs)

- DDIs can occur when a patient takes 2 or more drugs

# Drug-drug interactions (DDIs)

- DDIs can occur when a patient takes 2 or more drugs

- DDIs cause unexpected side effects

# Drug-drug interactions (DDIs)

- DDIs can occur when a patient takes 2 or more drugs

- DDIs cause <span style="color:#a02020">unexpected side effects</span>

  - 10-30% of adverse drug events are attributed to DDIs

# Drug-drug interactions (DDIs)

- DDIs can occur when a patient takes 2 or more drugs

- DDIs cause <span style="color:darkred">unexpected side effects</span>

  - 10-30% of adverse drug events are attributed to DDIs

- Understanding of DDIs may lead to better outcomes
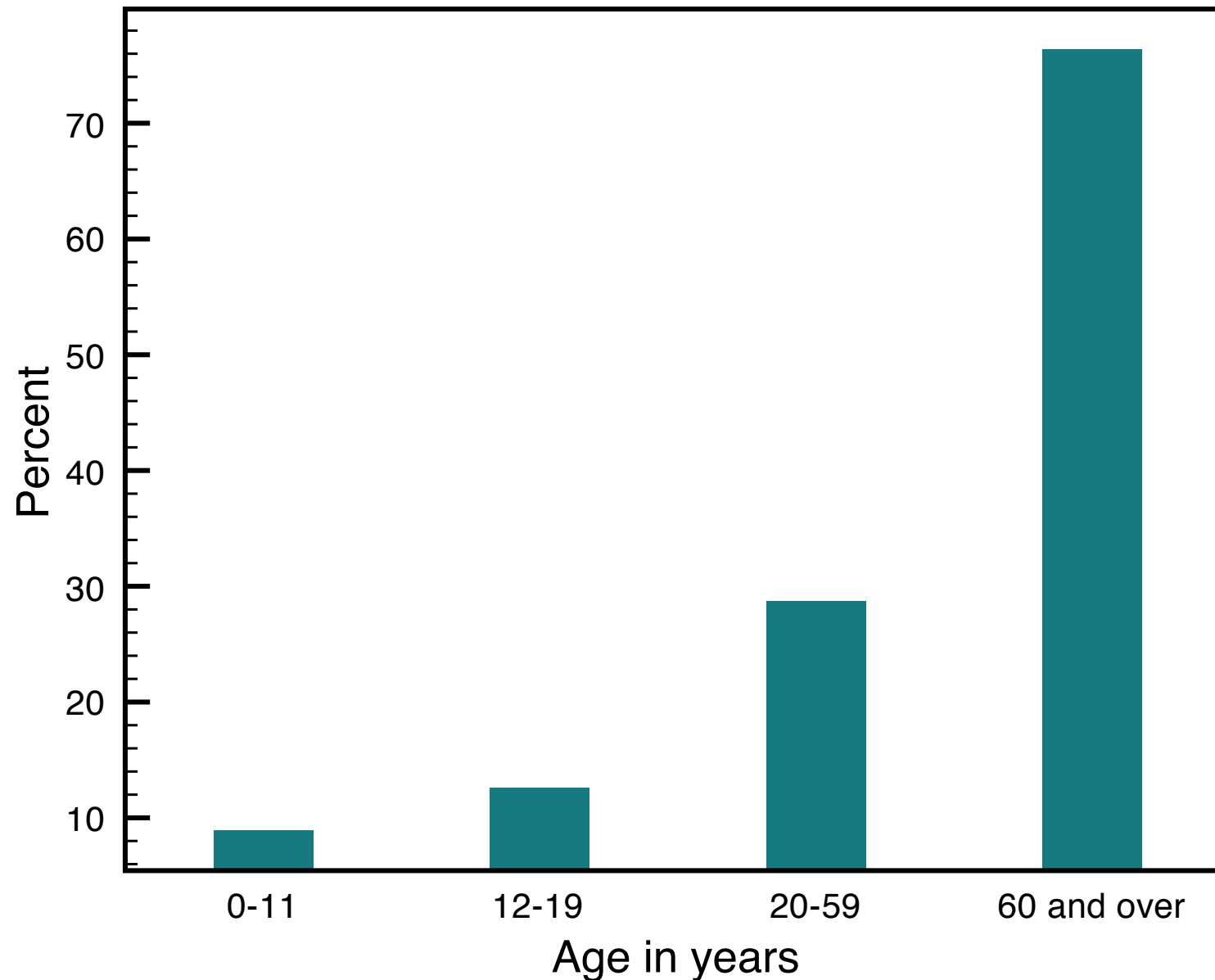
# Drug-drug interactions (DDIs)

- DDIs can occur when a patient takes 2 or more drugs

- DDIs cause <span style="color:darkred">unexpected side effects</span>

  - 10-30% of adverse drug events are attributed to DDIs

- Understanding of DDIs may lead to better outcomes

  - precaution in prescription

# Drug-drug interactions (DDIs)

- DDIs can occur when a patient takes 2 or more drugs

- DDIs cause <span style="color:#AA2222">unexpected side effects</span>

  - 10-30% of adverse drug events are attributed to DDIs

- Understanding of DDIs may lead to better outcomes

  - precaution in prescription

  - synergistic therapies

# Polypharmacy increases with age

**Percent of people on two or more drugs by age**
**United States 2007-2008**



SOURCE: CDC/NCHS, National Health and Nutrition Examination Survey

76% of older Americans used two or more prescription drugs

# More needs to be done to understand and identify drug-drug interactions

# More needs to be done to understand and identify drug-drug interactions

- Clinical trials do not typically investigate drug-drug interactions

# More needs to be done to understand and identify drug-drug interactions

- Clinical trials do not typically investigate <span style="color:red">drug-drug interactions</span>

- **Observational studies** are the only systematic way to detect drug-drug interactions

# Large population databases enable DDI discovery

- Contain clinical data on millions of patients over many years

- Currently being used to establish single drug adverse events (pharmacovigilance)

- Eg. **Spontaneous Adverse Event Reporting Systems**

  - Collect adverse event reports for a patient (a snapshot in time)

  - Maintained by WHO > FDA > Health Canada

# Observational data are messy

**Drugs**

METFORMIN

ROSIGLITAZONE

PRAVASTATIN

TACROLIMUS

PREDNISOLONE

**Adverse Events**

ACUTE RESP. DISTRESS
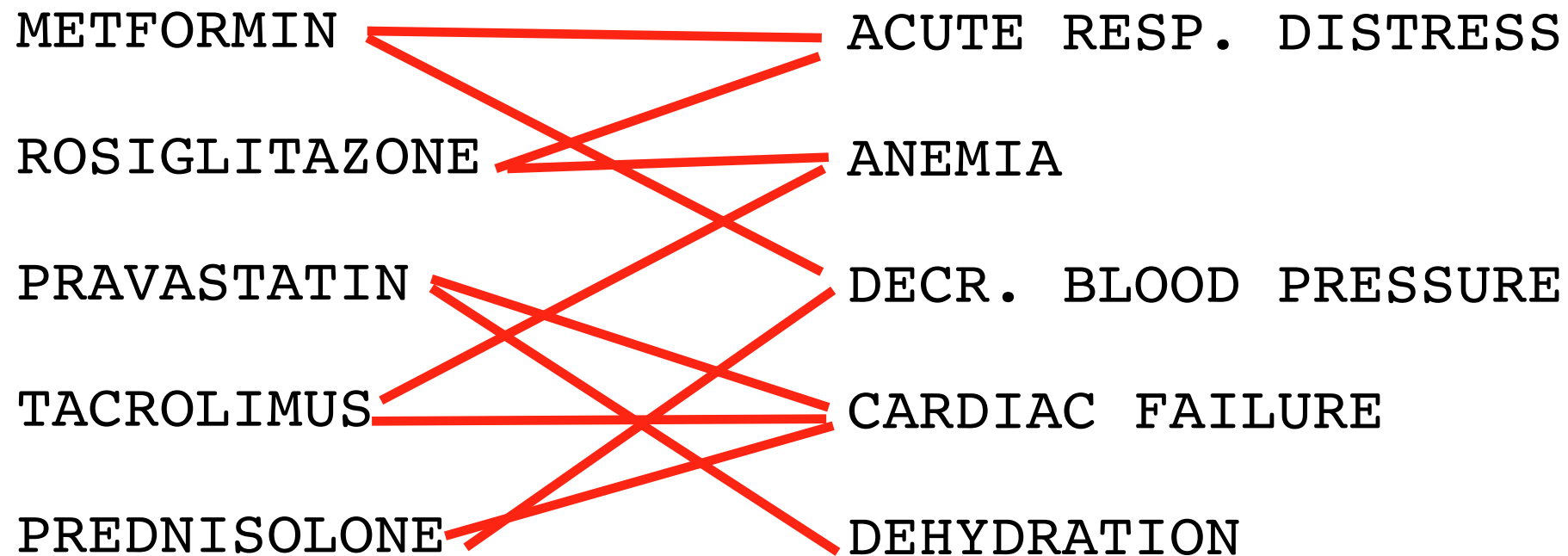
ANEMIA

DECR. BLOOD PRESSURE

CARDIAC FAILURE

DEHYDRATION

# Observational data are messy

- Many drugs, many adverse events

| Drugs | Adverse Events |
|---|---|
| METFORMIN | ACUTE RESP. DISTRESS |
| ROSIGLITAZONE | ANEMIA |
| PRAVASTATIN | DECR. BLOOD PRESSURE |
| TACROLIMUS | CARDIAC FAILURE |
| PREDNISOLONE | DEHYDRATION |

# Observational data are messy

- Many drugs, many adverse events

  - what causes what?

**Drugs**                              **Adverse Events**
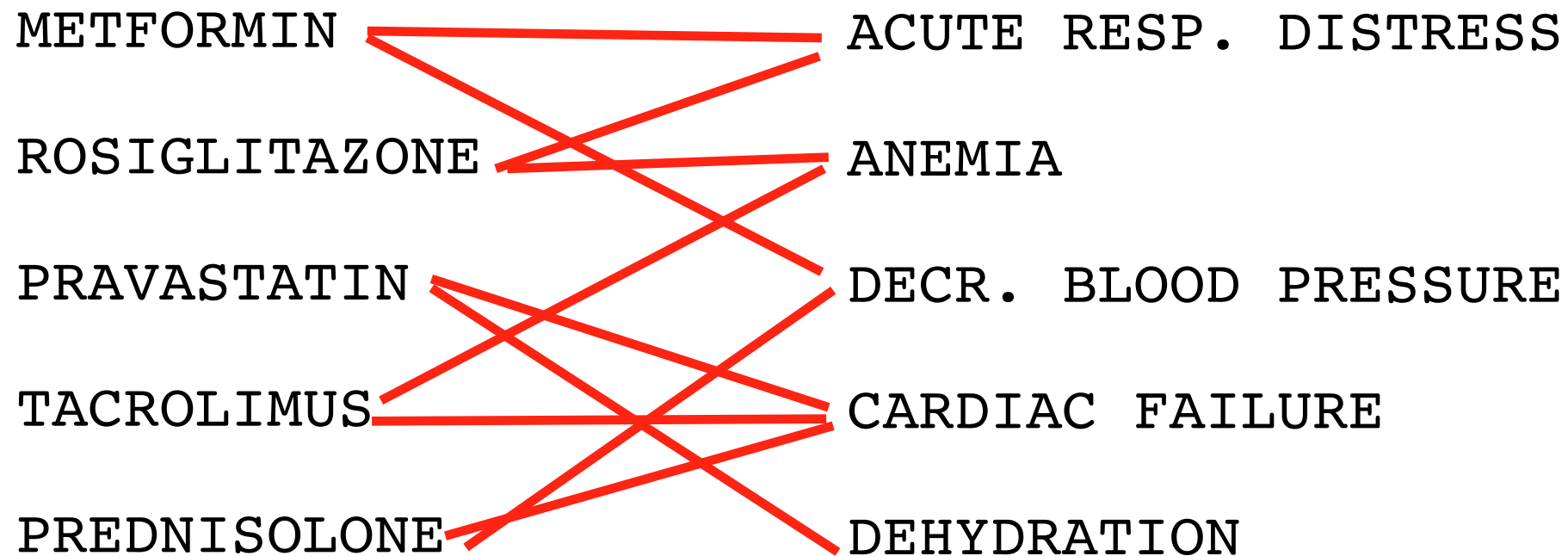
METFORMIN                              ACUTE RESP. DISTRESS

ROSIGLITAZONE                          ANEMIA

PRAVASTATIN                            DECR. BLOOD PRESSURE

TACROLIMUS                             CARDIAC FAILURE

PREDNISOLONE                           DEHYDRATION

# Observational data are messy

- Many drugs, many adverse events

  - what causes what?

| Drugs | Adverse Events |
|---|---|
| METFORMIN | ACUTE RESP. DISTRESS |
| ROSIGLITAZONE | ANEMIA |
| PRAVASTATIN | DECR. BLOOD PRESSURE |
| TACROLIMUS | CARDIAC FAILURE |
| PREDNISOLONE | DEHYDRATION |

most of these red lines are false - which are true?

# Observational data are confounded

- Spontaneous reporting systems are observational data sets (unknown biases)

- noise from concomitant drug use (***co-Rx effect***)

  - drugs co-prescribed with Vioxx more likely to be associated with heart attacks

- noise from indications (***indication-effect***)

  - drugs given to diabetics more likely to be associated with hyperglycemia

# SCRUB
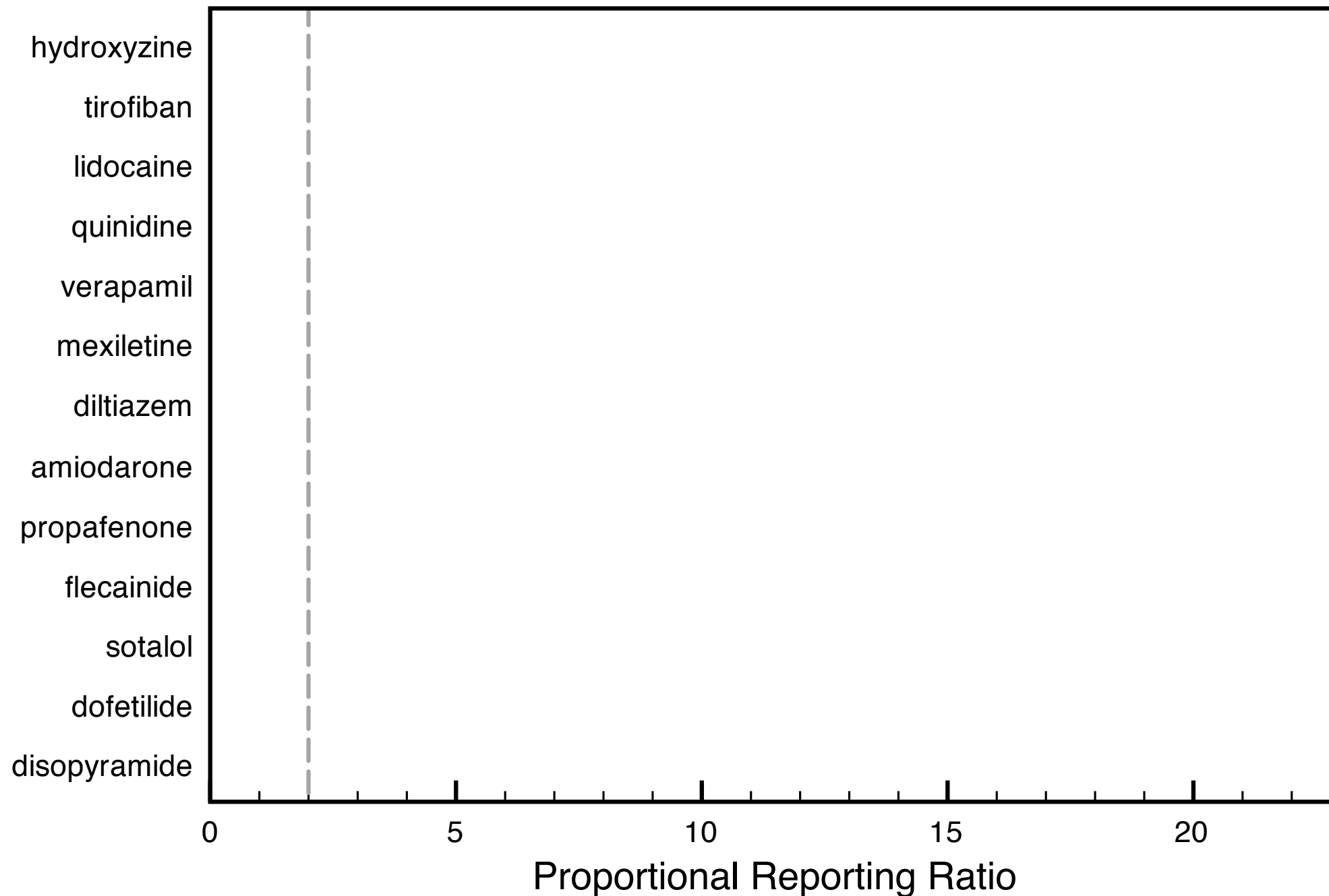## Statistical CorRection of Uncharacterized Bias

- *Implicitly* corrects for confounding of both observed and missing variables

- Assumes some combination of the **drugs** and **indications** describes the patient covariates
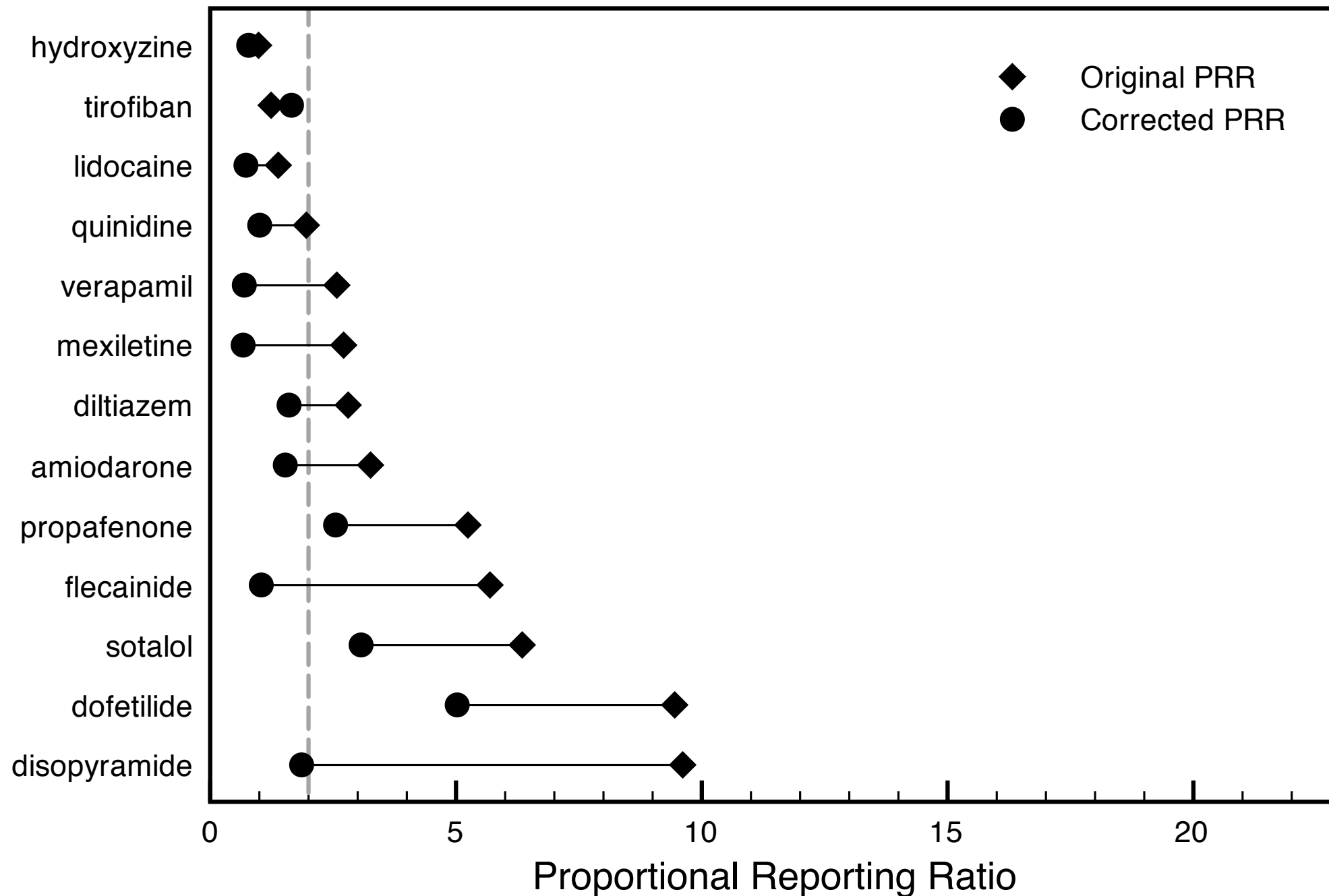
- Only works on very large data sets

N. Tatonetti et al., *Science Translational Medicine* (2012)

# Method corrects for **indication biases**

**Anti-arrhythmics and Arrhythmia**



Proportional Reporting Ratio plot with y-axis labels (top to bottom): hydroxyzine, tirofiban, lidocaine, quinidine, verapamil, mexiletine, diltiazem, amiodarone, propafenone, flecainide, sotalol, dofetilide, disopyramide. X-axis: Proportional Reporting Ratio (0, 5, 10, 15, 20).

# Method corrects for **indication biases**



Anti-arrhythmics and Arrhythmia

# Method corrects for **indication biases**



Anti-arrhythmics and Arrhythmia

Implicit correction of age differences in exposed vs non-exposed

# Bias, corrected. Missing data?

**If there are no observations
then no associations can be found.**

# Diseases can be identified by the side effects they elicit

# Diseases can be identified by the side effects they elicit

Diabetes

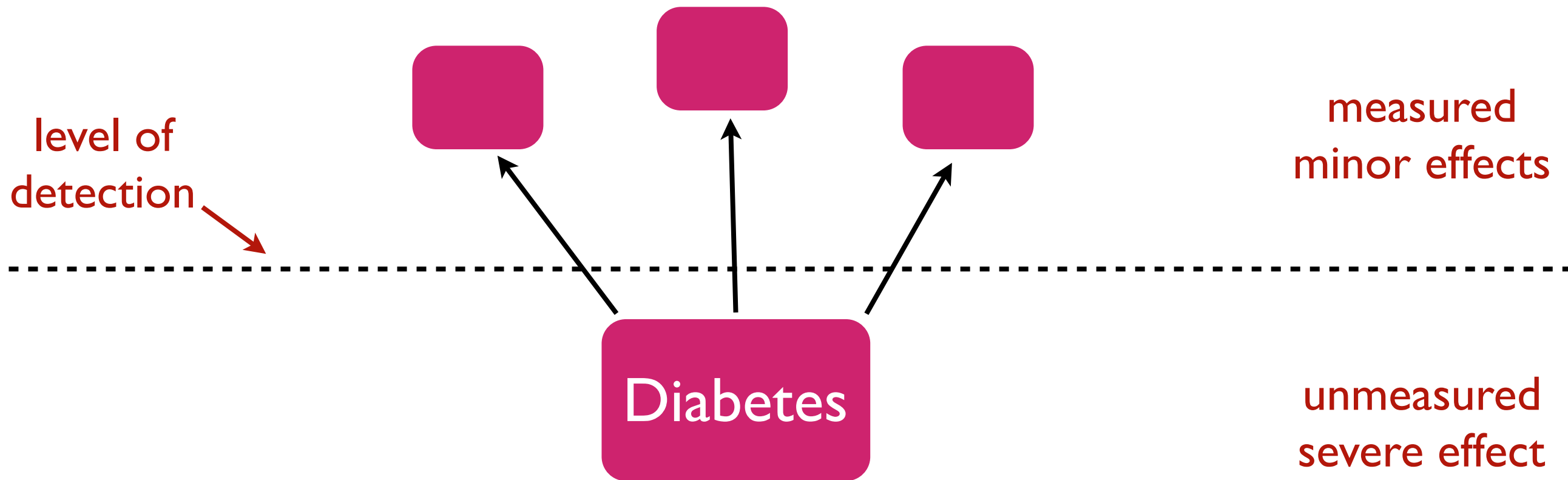# Diseases can be identified by the side effects they elicit

level of
detection

Diabetes

# Diseases can be identified by the side effects they elicit
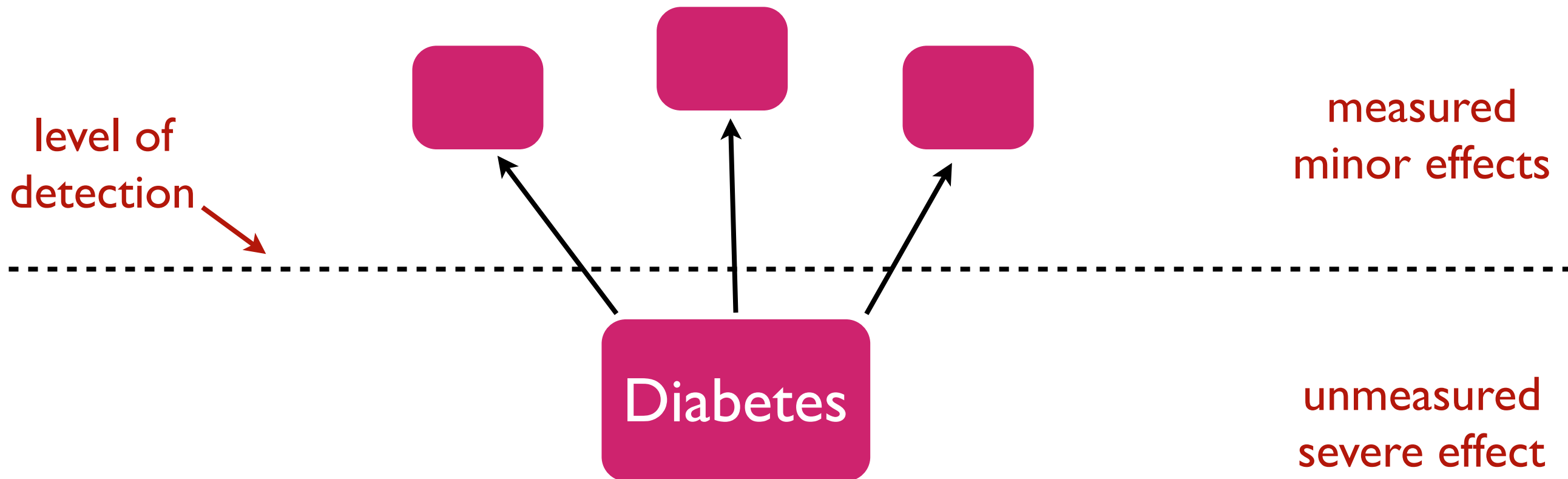
level of detection

Diabetes

unmeasured severe effect

21

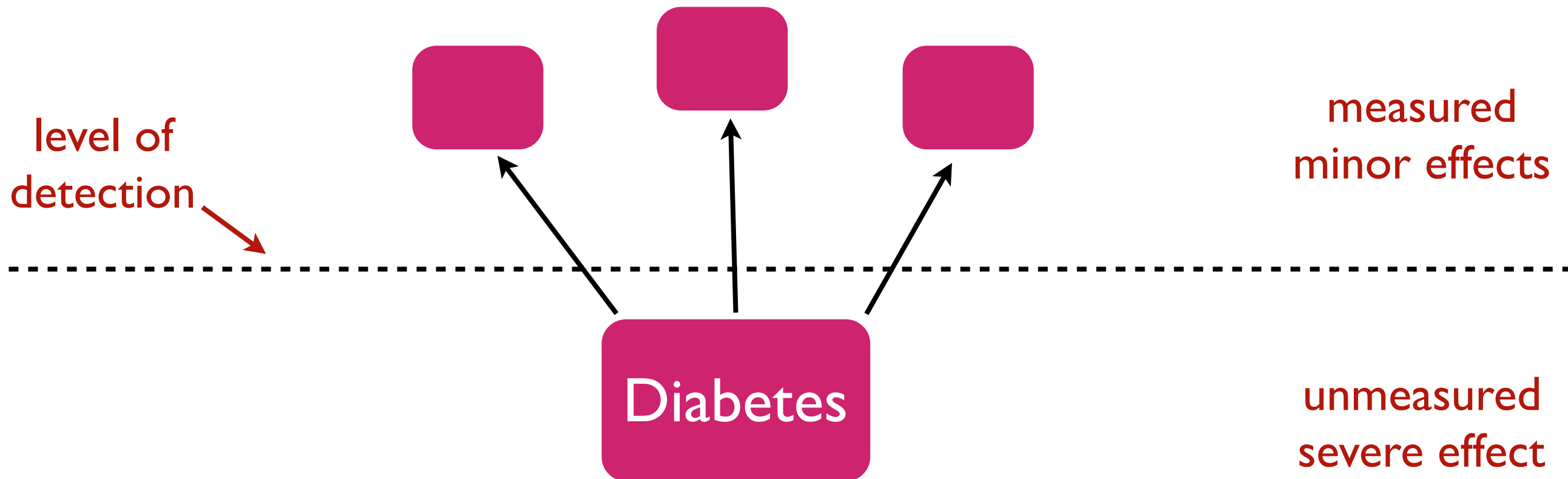# Diseases can be identified by the side effects they elicit



level of detection

measured minor effects

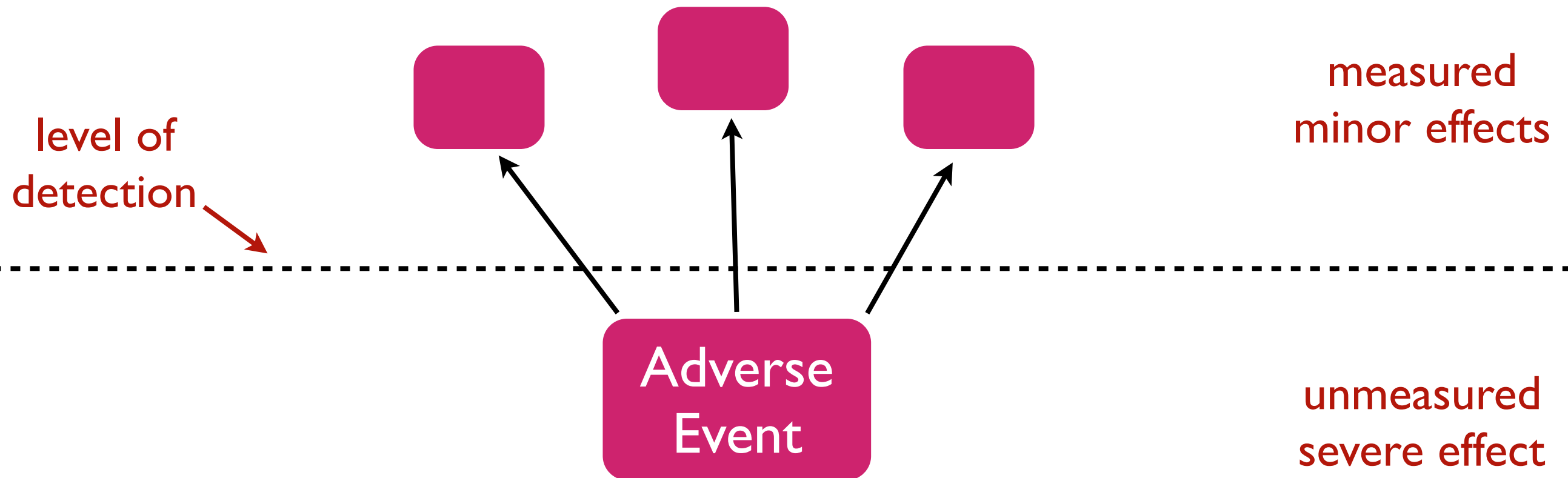Diabetes

unmeasured severe effect

21

# Diseases can be identified by the side effects they elicit

- physicians use observable side effects to form hypothesis about the underlying disease



level of detection

measured minor effects

Diabetes

unmeasured severe effect

# Diseases can be identified by the side effects they elicit

- physicians use observable side effects to form hypothesis about the underlying disease

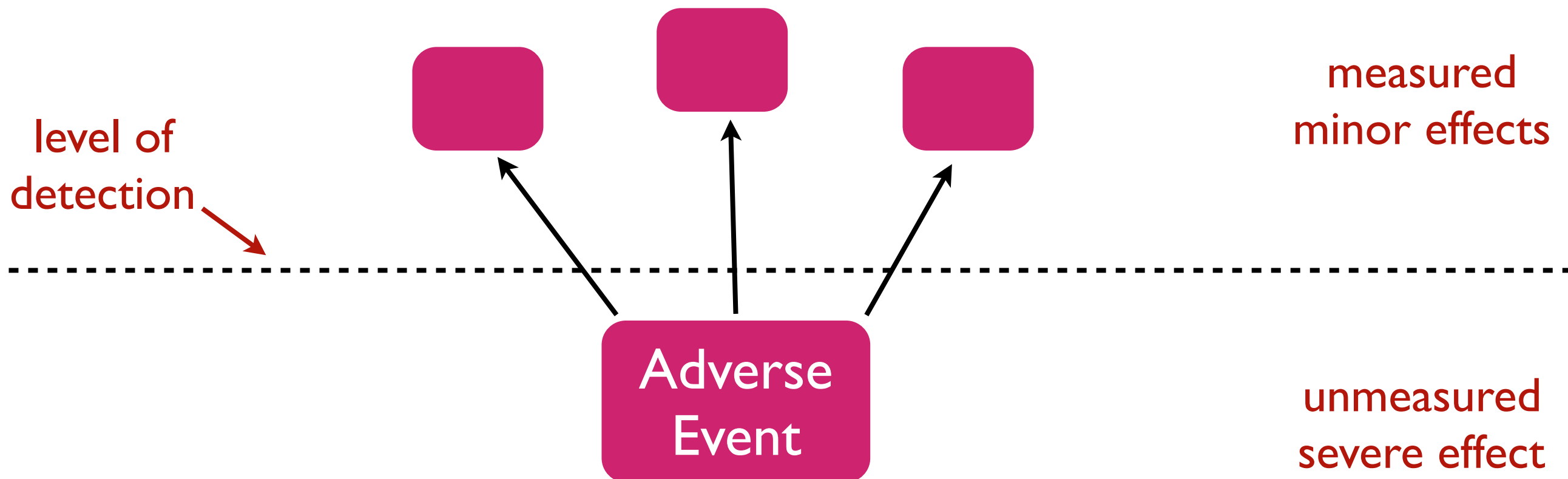- e.g. you can't *see* diabetes, but you can *measure* blood glucose

level of detection

measured minor effects

Diabetes

unmeasured severe effect

# Severe ADE's can be identified by the presence of more minor (and more common) side effects



level of
detection

measured
minor effects
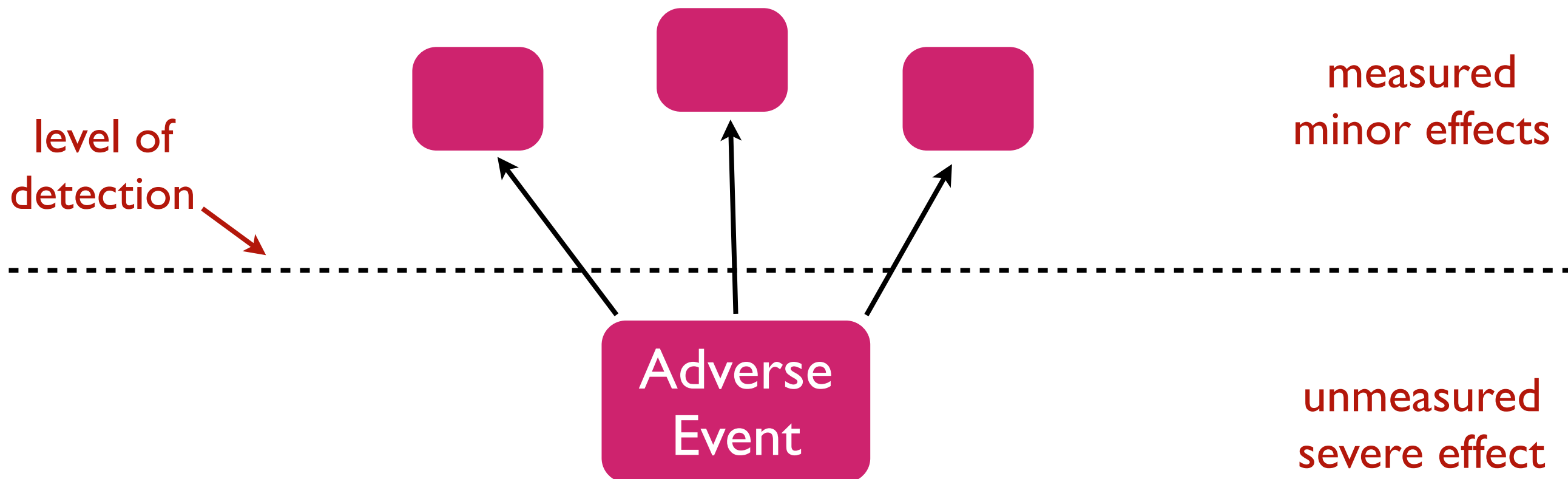
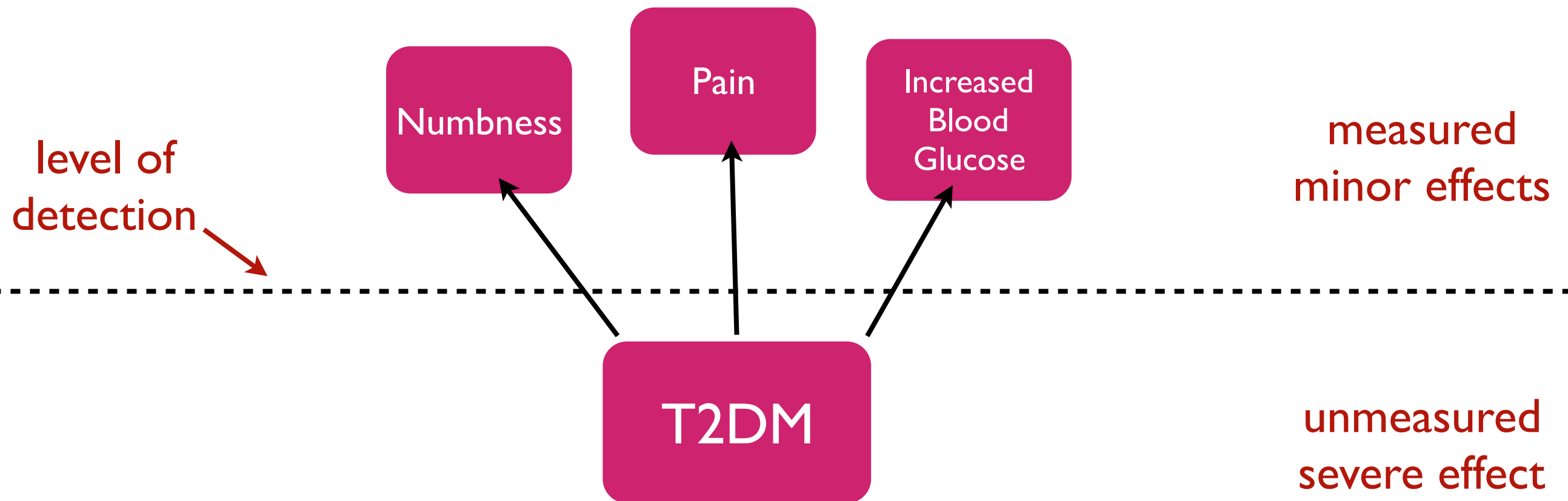Adverse
Event

unmeasured
severe effect

# Severe ADE's can be identified by the presence of more minor (and more common) side effects

- First, identify the common side effects that are harbingers for the underlying severe AE



level of detection

measured minor effects

**Adverse Event**

unmeasured severe effect

# Severe ADE's can be identified by the presence of more minor (and more common) side effects

- First, identify the common side effects that are harbingers for the underlying severe AE

- Then, combine these side effects together to form an "effect profile" for an adverse event



level of detection

measured minor effects

Adverse Event

unmeasured severe effect

# Severe ADEs can be identified by the presence of more minor (and more common) side effects
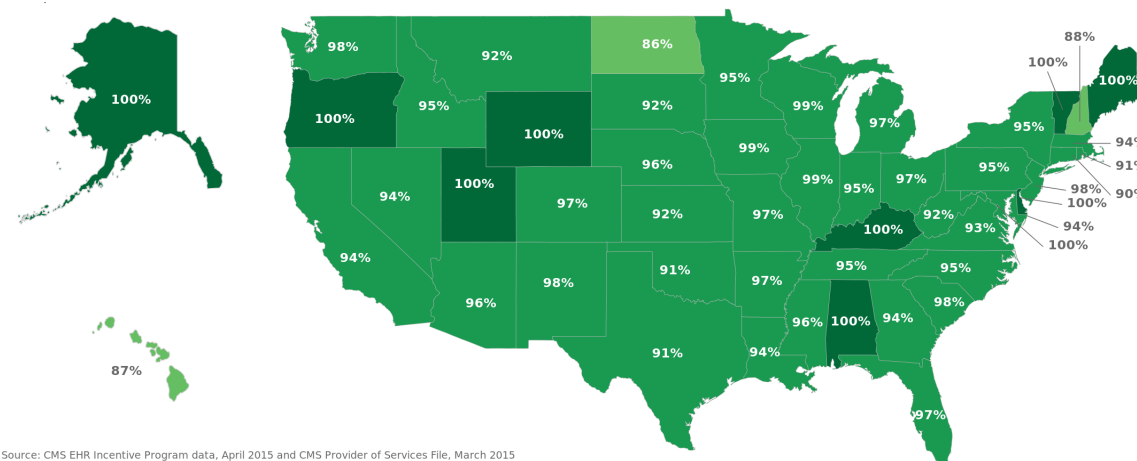
# DDI prediction validation

**Table S3** Novel drug-drug interaction predictions for diabetes related adverse events.

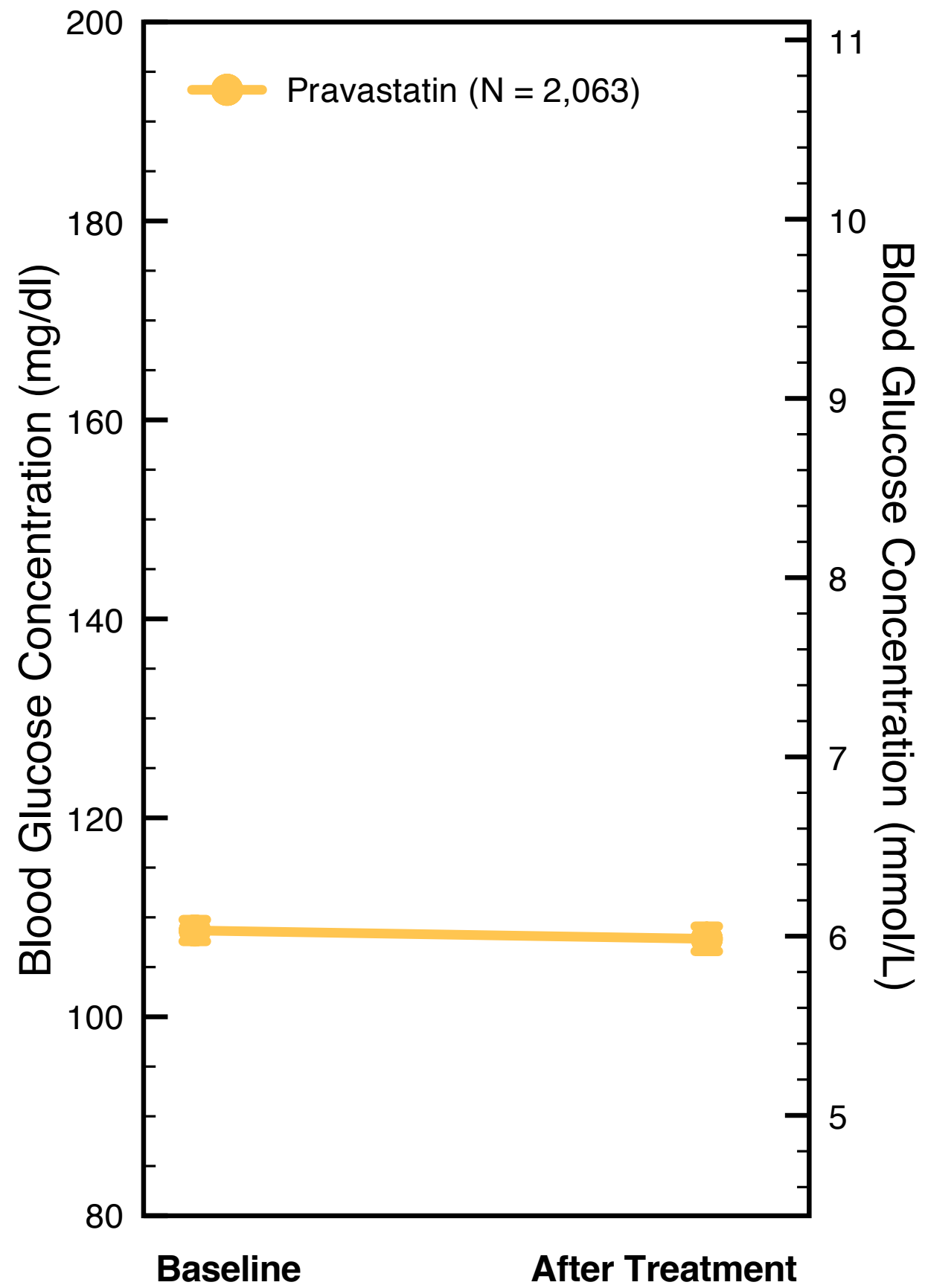| Rank | Drug A | Drug B | Score | Minimum Randomization Rank | Known DDI exists |
|------|--------|--------|-------|----------------------------|------------------|
| 38 | PAROXETINE HCL | PRAVASTATIN SODIUM | 11.3518960149 | 62 | |
| 72 | DIOVAN HCT | HYDROCHLOROTHIAZIDE | 7.1786599539 | 89 | |
| 94 | CRESTOR | PREVACID | 4.7923771645 | 148 | |
| 107 | DESFERAL | EXJADE | 3.97220625 | 129 | |
| 159 | COUMADIN | VESICARE | 0.8928376683 | 169 | |
| 160 | DEXAMETHASONE | THALIDOMIDE | 0.8928376683 | 168 | CRITICAL |
| 170 | FOSAMAX | VOLTAREN | 0.5033125 | 1138 | |
| 175 | ALIMTA | DEXAMETHASONE | 0.2442375 | 197 | |

- Focus on top hit from diabetes classifier

- paroxetine = depression drug, pravastatin = cholesterol drug

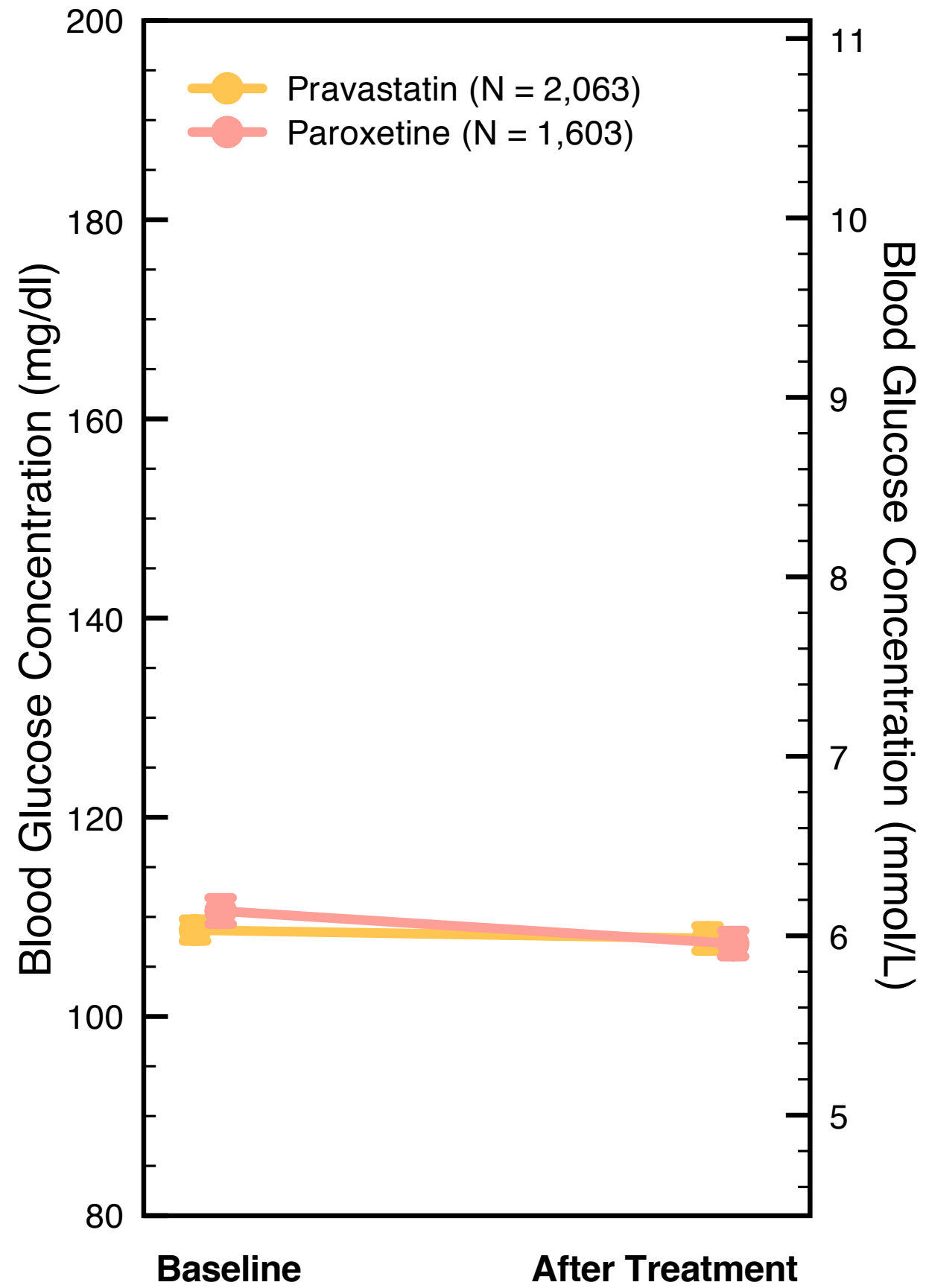- Popular drugs, est. ~1,000,000 patients on this combination!

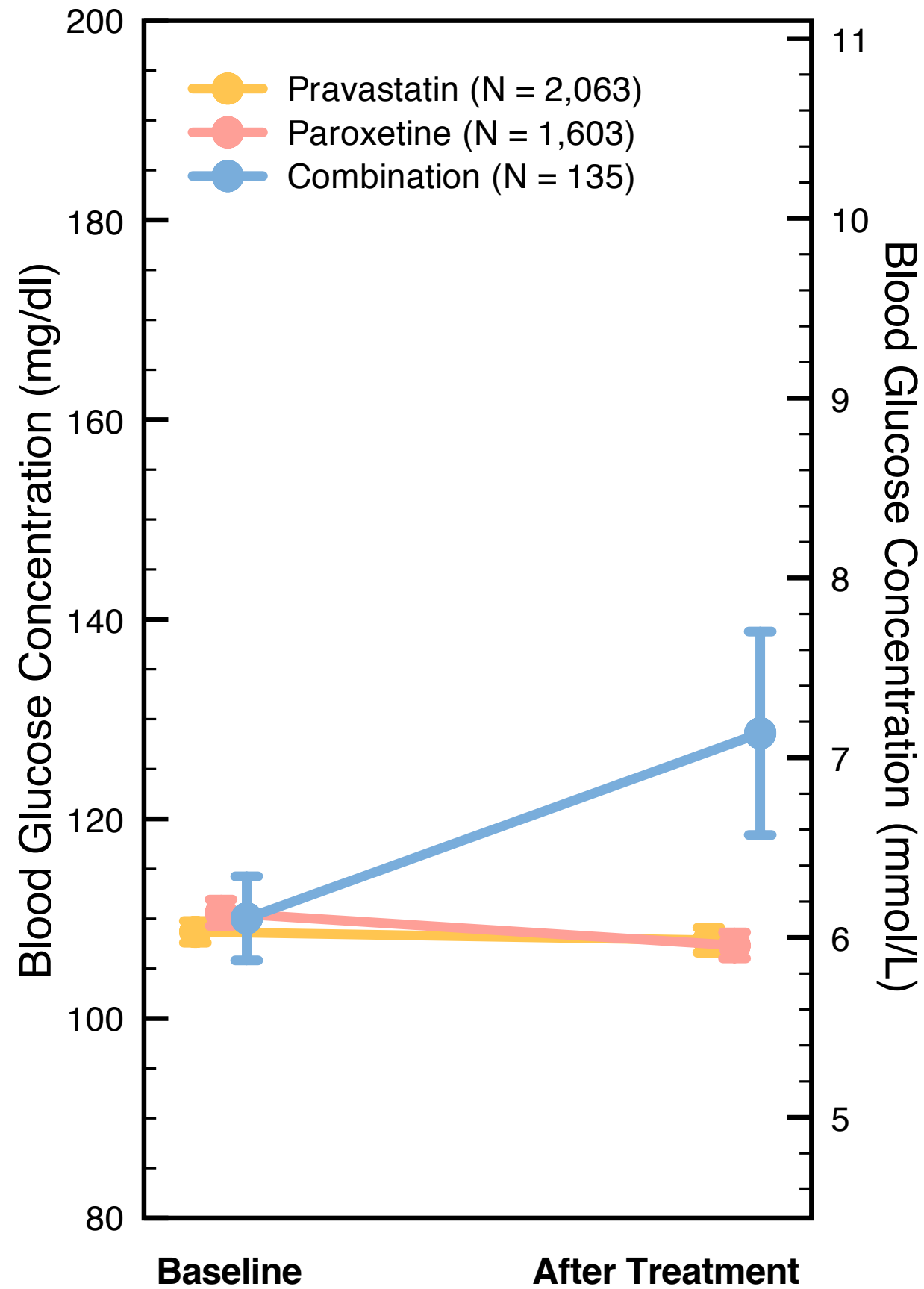# Analyzed blood glucose values for patients on either or both of these drugs

## To the electronic health records…

Blood Glucose Concentration (mg/dl)

Blood Glucose Concentration (mmol/L)

Pravastatin (N = 2,063)

Baseline

After Treatment

Tatonetti, et al. *Clinical Pharmacology & Therapeutics* (2011)

Tatonetti, et al. *Clinical Pharmacology & Therapeutics* (2011)

Tatonetti, et al. *Clinical Pharmacology & Therapeutics* (2011)

# no diabetics



Pravastatin (N = 2,063)
Paroxetine (N = 1,603)
Combination (N = 135)

Blood Glucose Concentration (mg/dl)

Blood Glucose Concentration (mmol/L)

Baseline    After Treatment

## no diabetics

Pravastatin (N = 2,063)
Paroxetine (N = 1,603)
Combination (N = 135)

Blood Glucose Concentration (mg/dl)
Blood Glucose Concentration (mmol/L)

Baseline        After Treatment

## including diabetics

Pravastatin
Paroxetine
Combination (N=177)

Blood Glucose Concentration (mg/dl)

Baseline        After Treatment

Tatonetti, et al. *Clinical Pharmacology & Therapeutics* (2011)

no diabetics

Pravastatin (N = 2,063)
Paroxetine (N = 1,603)
Combination (N = 135)

Blood Glucose Concentration (mg/dl)
Blood Glucose Concentration (mmol/L)

Baseline    After Treatment

including diabetics

Pravastatin
Paroxetine
Combination (N=177)

Blood Glucose Concentration (mg/dl)

~60mg/dl increase

Baseline    After Treatment

Tatonetti, et al. *Clinical Pharmacology & Therapeutics* (2011)

# Informatics methods have taken us far, skeptics remain

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model

  - 10 control mice on normal diet (Ctl Ctl)

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model

  - 10 control mice on normal diet (Ctl Ctl)
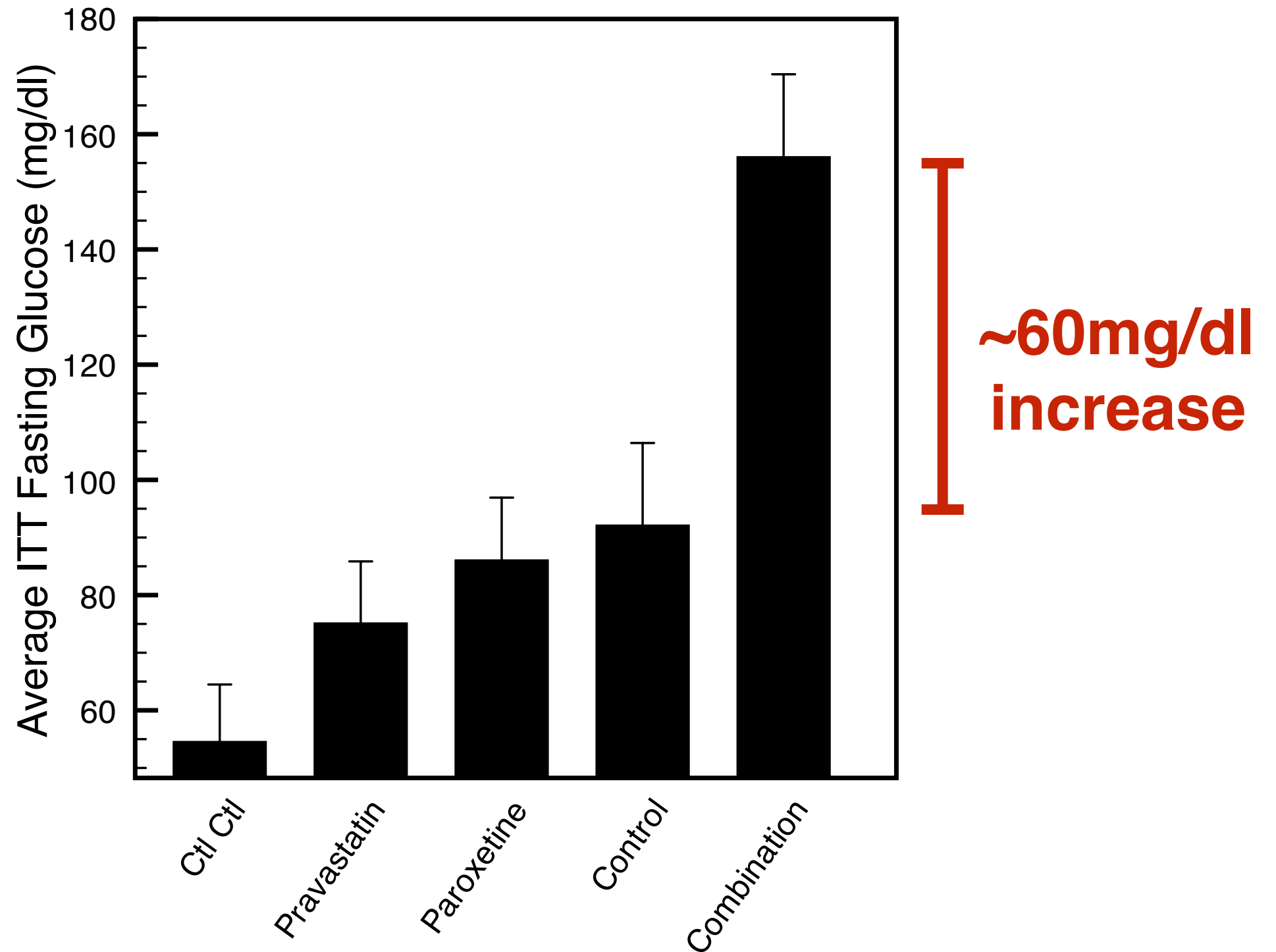
  - 10 control mice on high fat diet (HFD)

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model

  - 10 control mice on normal diet (Ctl Ctl)

  - 10 control mice on high fat diet (HFD)

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model

  - 10 control mice on normal diet (Ctl Ctl)

  - 10 control mice on high fat diet (HFD)

**Simulating Pre-Diabetics**

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model
  - 10 control mice on normal diet (Ctl Ctl)
  - 10 control mice on high fat diet (HFD)

**Simulating Pre-Diabetics**

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model

  - 10 control mice on normal diet (Ctl Ctl)

  - 10 control mice on high fat diet (HFD)

**Simulating Pre-Diabetics**

# Informatics methods have taken us far, skeptics remain

- Insulin Resistant Mouse Model

  - 10 control mice on normal diet (Ctl Ctl)

  - 10 control mice on high fat diet (HFD)

  - 10 mice on pravastatin + HFD

  - 10 mice on paroxetine + HFD

  - 10 mice on combination + HFD

# Summary of fasting glucose levels

# Replication is vital to science

- In biology we would never trust a result that hasn't been replicated

- Why should **algorithms** be any different?

# Drug-drug interactions and acquired Long QT Syndrome (LQTS)

- Long QT syndrome (LQTS): congenital or drug-induced change in electrical activity of the heart that can lead to potentially fatal arrhythmia: *torsades de pointes* (TdP)

- 13 genes associated with congenital LQTS

- Drug-induced LQTS usually caused by blocking the hERG channel (*KCNH2*)



From Berger et al., *Science Signaling* (2010)

# Identify acquired LQTS drug-drug interactions using Latent Signal Detection



level of detection

bradycardia

AFib

tachycardia

LQTS

measured minor effects

unmeasured severe effect

*Lorberbaum, et al. Drug Safety (2016)*

# Latent Signal Detection of acquired LQTS

## Top Prediction:
## Ceftriaxone + Lansoprazole

- Ceftriaxone — common in-patient cephalosporin antibiotic

- Lansoprazole — proton-pump inhibitor used to treat GERD, one of the most commonly taken drugs in the world

- In the EHR: Patients on the combination have QT intervals 11ms longer, on average and are **1.5X as likely to have a QT interval > 500ms**

|  | White | Black/African American | Other, including Hispanic | Asian |
|---|---|---|---|---|
| Females | 11.1 ± 3.1 ms** (N=220) | -1.3 ± 7.4 ms (N=91) | 6.0 ± 4.9 ms (N=78) | 13.2 ± 4.8 ms (N=4) |
| Males | 15.1 ± 4.1 ms** (N=164) | 0.7 ± 7.2 ms (N=53) | 10.5 ± 6.6 ms (N=46) | 8.3 ± 12.5 ms (N=4) |

** p < 0.01, one sample Student's T test

*Lorberbaum, et al. Drug Safety (2016)*
*Lorberbaum, et al. JACC (In press)*

- Predicted QT-DDI: **ceftriaxone** (cephalosporin antibiotic) and **lansoprazole** (proton pump inhibitor)

- Neither drug alone has any evidence of QT prolongation/ hERG block

- Predicted QT-DDI: **ceftriaxone** (cephalosporin antibiotic) and **lansoprazole** (proton pump inhibitor)
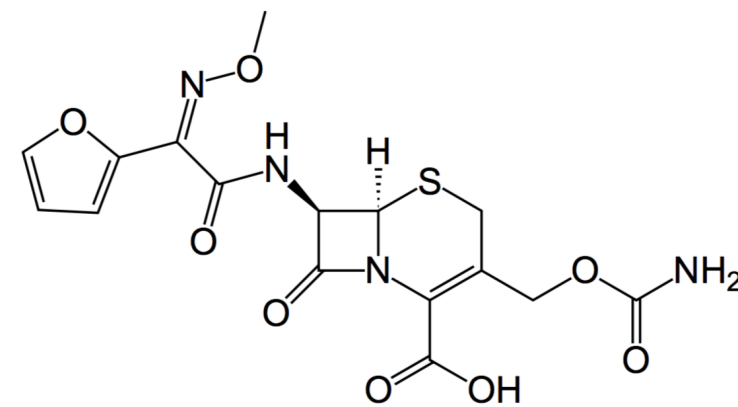
- Predicted QT-DDI: **ceftriaxone** (cephalosporin antibiotic) and **lansoprazole** (proton pump inhibitor)

- <u>Negative control</u>: lansoprazole + **cefuroxime** (another cephalosporin) – no evidence in FAERS of an interaction

- Predicted QT-DDI: **ceftriaxone** (cephalosporin antibiotic) and **lansoprazole** (proton pump inhibitor)

- Negative control: lansoprazole + **cefuroxime** (another cephalosporin) – no evidence in FAERS of an interaction



Ceftriaxone



Cefuroxime

# FAERS

Ceftriaxone+
Lansoprazole

*Lorberbaum, et al. In Revision*

**FAERS**

Ceftriaxone+ Lansoprazole

*Lorberbaum, et al. In Revision*

FAERS
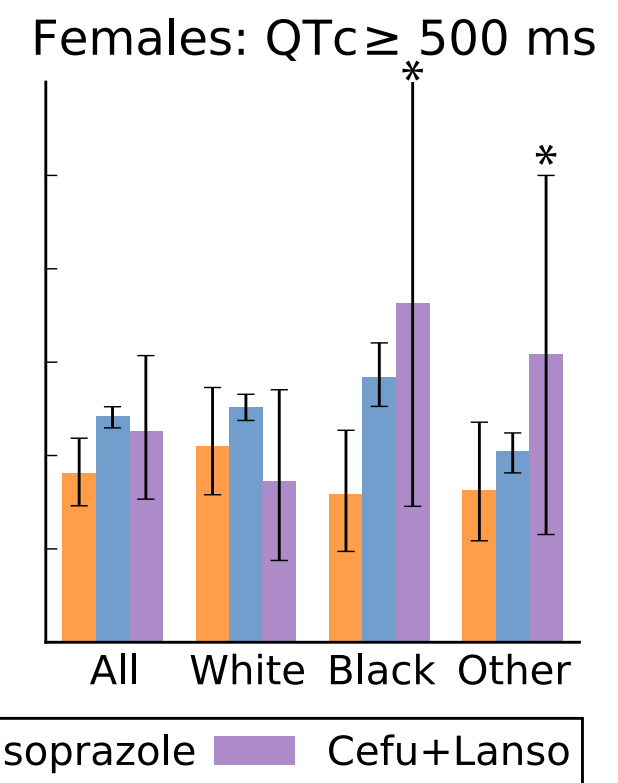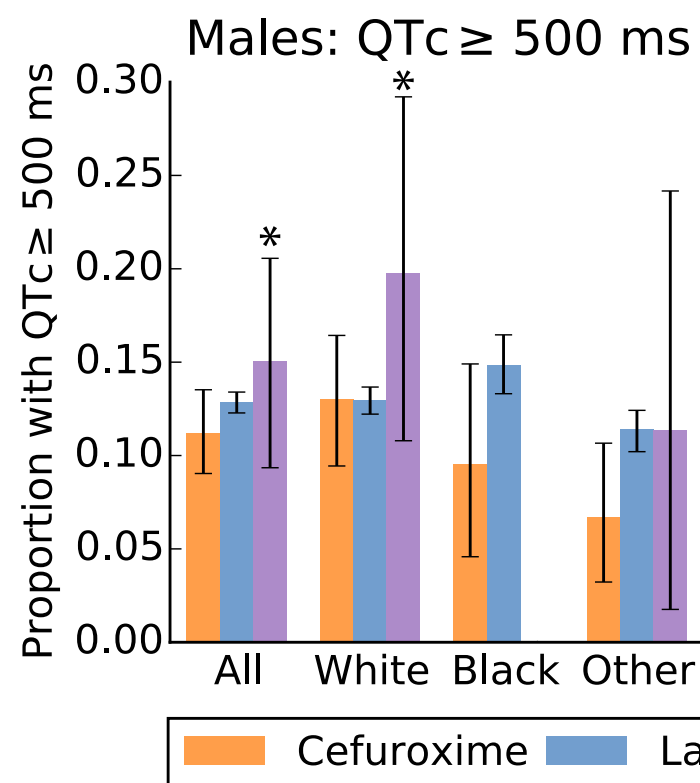
Ceftriaxone+Lansoprazole

Cefuroxime+Lansoprazole

Lorberbaum, et al. In Revision

# Electronic Health Records



Ceftriaxone+Lansoprazole

**Ceft+Lanso ΔQTc**

**Males: QTc ≥ 500 ms**

**Females: QTc ≥ 500 ms**

● Males    ● Females

■ Ceftriaxone  ■ Lansoprazole  ■ Ceft+Lanso

Cefuroxime+Lansoprazole

**Cefu+Lanso ΔQTc**

**Males: QTc ≥ 500 ms**

**Females: QTc ≥ 500 ms**

● Males    ● Females

■ Cefuroxime  ■ Lansoprazole  ■ Cefu+Lanso

*Lorberbaum, et al. In Revision*

# Electronic Health Records



Lorberbaum, et al. In Revision

# What is the mechanism?

# MADSS

## Modular Assembly of Drug Safety Subnetworks

- Use network analysis to build AE neighborhoods: a subset of the interactome surrounding AE "seed" proteins

- Score each protein on connectivity to seeds using:
  - Mean first passage time
  - Betweenness centrality
  - Shared neighbors
  - Inverse shortest path

- Overarching hypothesis: drugs targeting proteins within an AE neighborhood more likely to be involved in mediating that AE



- ● Protein
- ● Seed protein
- — Interaction
- ■ Adverse event (AE)
- ▲ Drug known to cause AE
- △ Drug predicted to cause AE

Lorberbaum, et al. *Clin. Pharmacol. Ther.* (2015)
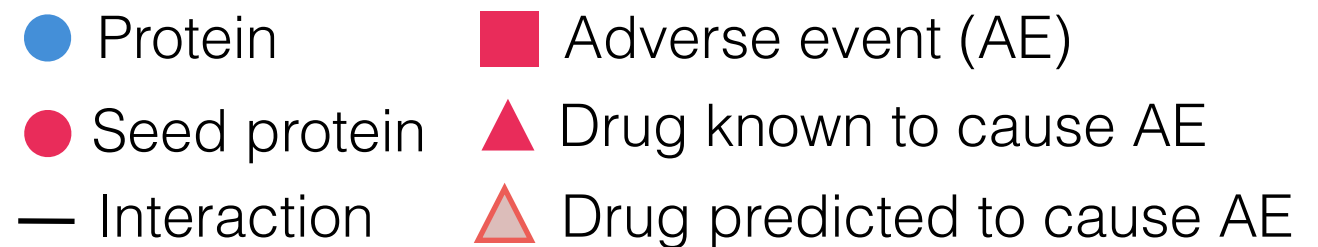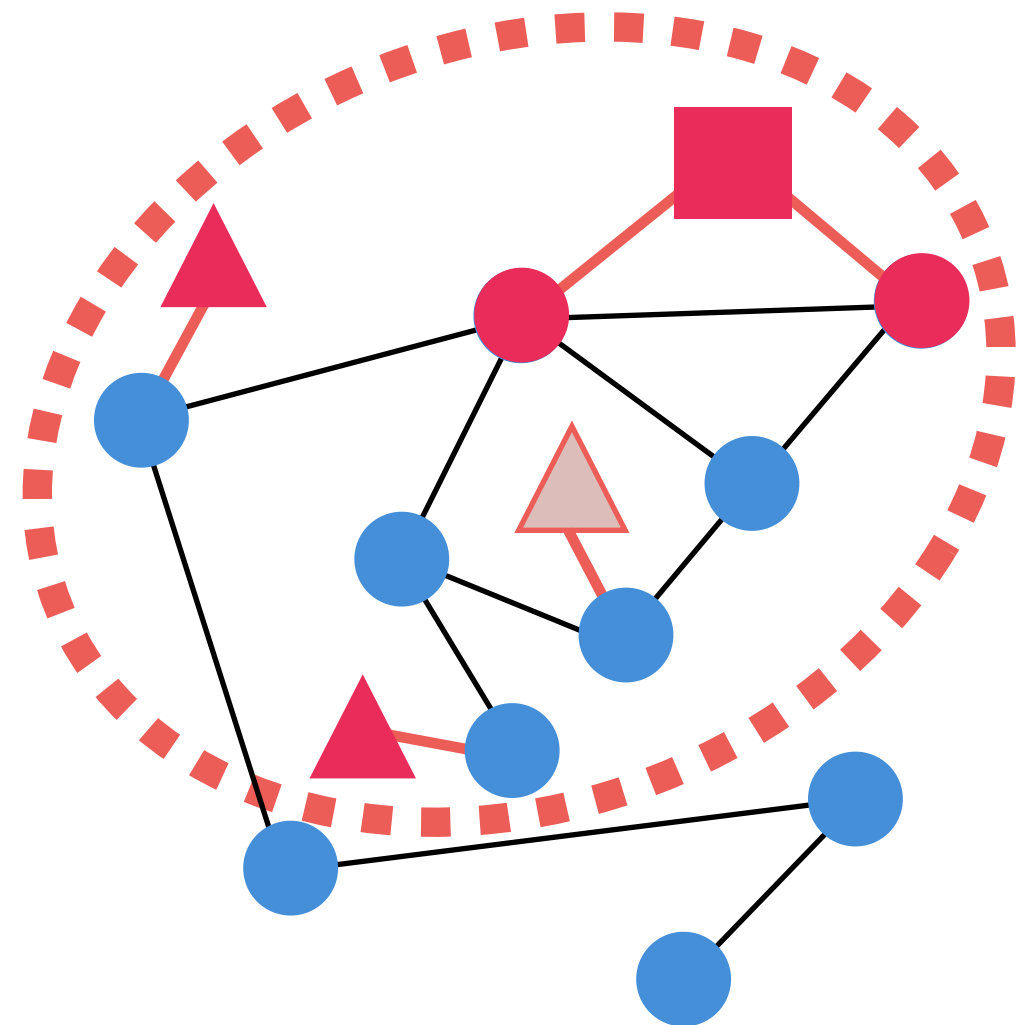
# MADSS

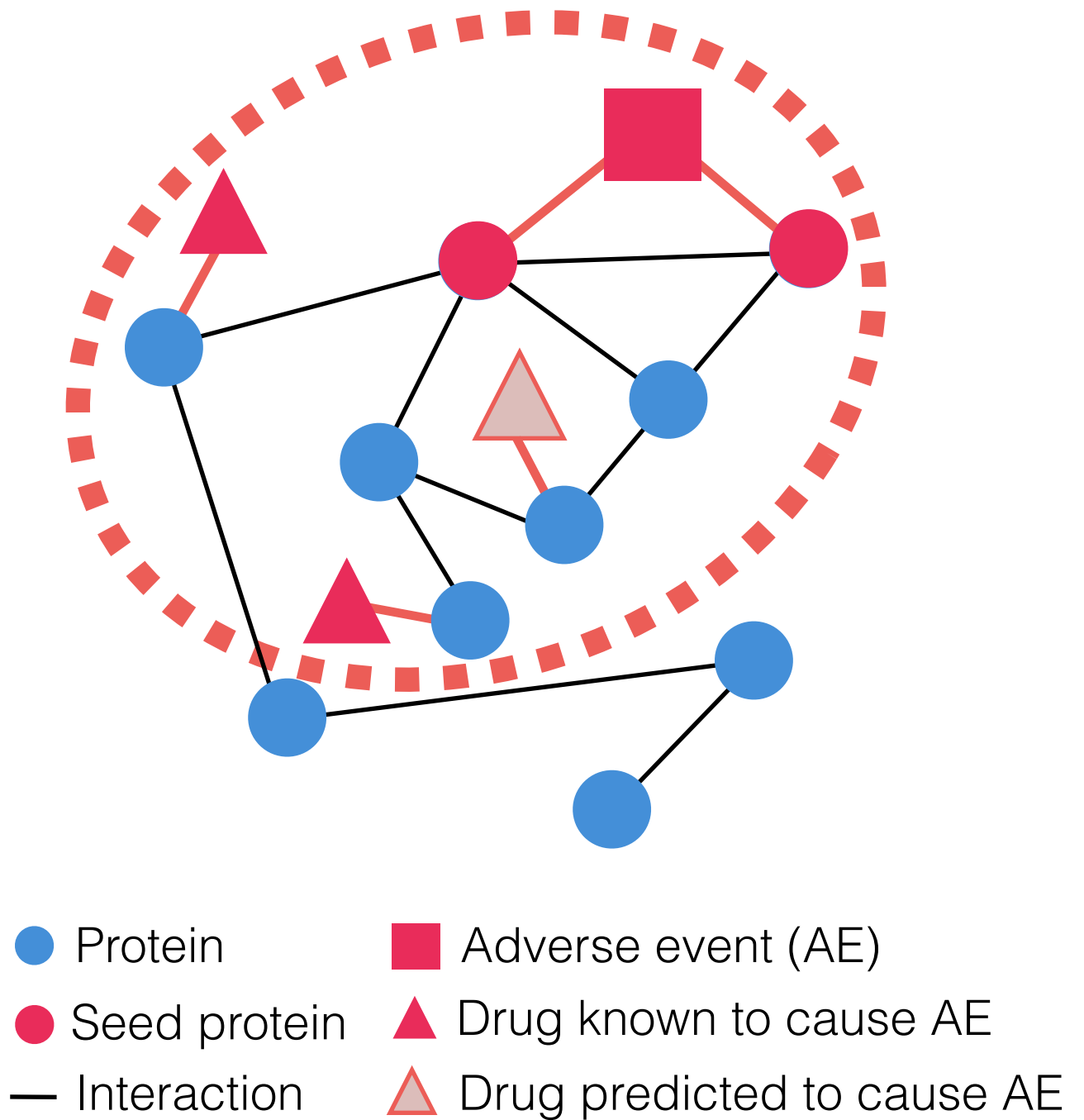## Modular Assembly of Drug Safety Subnetworks

- Use network analysis to build AE neighborhoods: a subset of the interactome surrounding AE "seed" proteins

- Score each protein on connectivity to seeds using:
  - Mean first passage time
  - Betweenness centrality
  - Shared neighbors
  - Inverse shortest path

- Overarching hypothesis: drugs targeting proteins within an AE neighborhood more likely to be involved in mediating that AE

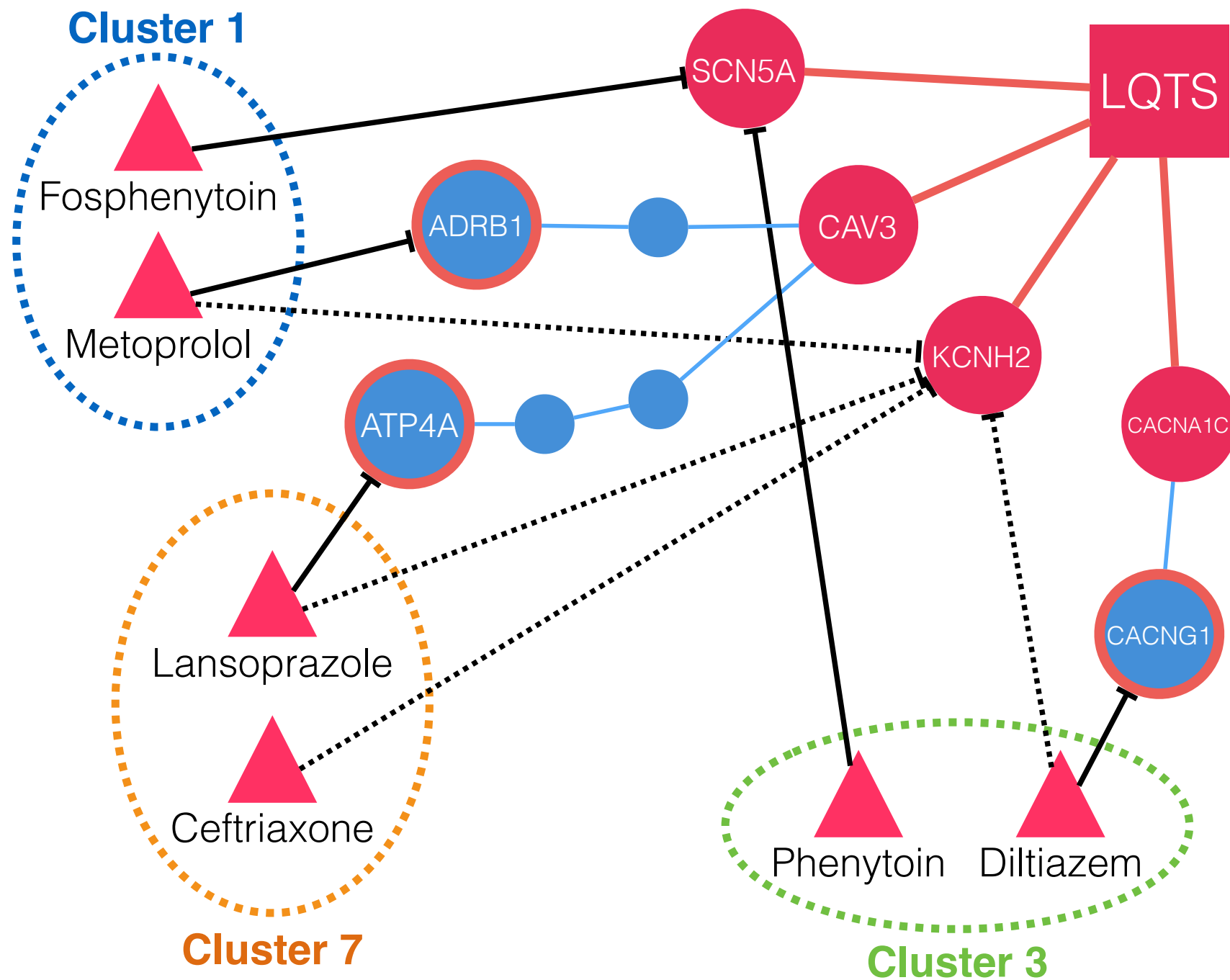- **Ran MADSS using 13 LQTS genes as seeds**



- 🔵 Protein
- 🔴 Seed protein
- — Interaction
- 🟥 Adverse event (AE)
- 🔺 Drug known to cause AE
- 🔺 Drug predicted to cause AE

Lorberbaum, et al. *Clin. Pharmacol. Ther.* (2015)

Putative mechanisms of QT-DDIs

# Automated Patch Clamp

- Collaboration with Rocky Kass (CUMC Pharmacology Dept.)

- Take HEK293 cells over-expressing the hERG channel

- Perform a single-cell patch clamp experiment

  - control

  - ceftriaxone alone

  - lansoprazole alone

  - combination of ceftriaxone and lansoprazole



*Nanion Patchliner*



Voltage protocol: step to +40mV followed by a return to -40mV

*Lorberbaum, et al. JACC (In press)*

# Ceftriaxone+Lansoprazole



Legend panel 1:
- Vehicle Control
- + Ceftriaxone 0.1μM
- + Ceftriaxone 10μM
- + Ceftriaxone 100μM

Legend panel 2:
- Vehicle Control
- Lansoprazole 1μM
- + Ceftriaxone 0.1μM
- + Ceftriaxone 10μM
- + Ceftriaxone 100μM

Legend panel 3:
- Vehicle Control
- Lansoprazole 10μM
- + Ceftriaxone 0.1μM
- + Ceftriaxone 10μM
- + Ceftriaxone 100μM

*Lorberbaum, et al. JACC (In press)*

# Ceftriaxone+Lansoprazole

Vehicle Control
+ Ceftriaxone 0.1μM
+ Ceftriaxone 10μM
+ Ceftriaxone 100μM

Vehicle Control
Lansoprazole 1μM
+ Ceftriaxone 0.1μM
+ Ceftriaxone 10μM
+ Ceftriaxone 100μM

Vehicle Control
Lansoprazole 10μM
+ Ceftriaxone 0.1μM
+ Ceftriaxone 10μM
+ Ceftriaxone 100μM

# Cefuroxime+Lansoprazole

Vehicle Control
+ Cefuroxime 0.1μM
+ Cefuroxime 10μM
+ Cefuroxime 100μM

Vehicle Control
Lansoprazole 1μM
+ Cefuroxime 0.1μM
+ Cefuroxime 10μM
+ Cefuroxime 100μM

*Lorberbaum, et al. JACC (In press)*

# Ceftriaxone+Lansoprazole

Vehicle Control
+ Ceftriaxone 0.1µM
+ Ceftriaxone 10µM
+ Ceftriaxone 100µM

Vehicle Control
Lansoprazole 1µM
+ Ceftriaxone 0.1µM
+ Ceftriaxone 10µM
+ Ceftriaxone 100µM

Vehicle Control
Lansoprazole 10µM
+ Ceftriaxone 0.1µM
+ Ceftriaxone 10µM
+ Ceftriaxone 100µM

## Ceft+Lanso effect on hERG current

Change from Control

- Ceftriaxone + 10µM Lansoprazole
- Ceftriaxone + 1µM Lansoprazole
- Ceftriaxone alone

Ceftriaxone Concentration (µM)

# Cefuroxime+Lansoprazole

Vehicle Control
+ Cefuroxime 0.1µM
+ Cefuroxime 10µM
+ Cefuroxime 100µM

Vehicle Control
Lansoprazole 1µM
+ Cefuroxime 0.1µM
+ Cefuroxime 10µM
+ Cefuroxime 100µM

## Cefu+Lanso effect on hERG current

Change from Control

- Cefuroxime + 1µM Lansoprazole
- Cefuroxime alone

Cefuroxime Concentration (µM)

*Lorberbaum, et al. JACC (In press)*

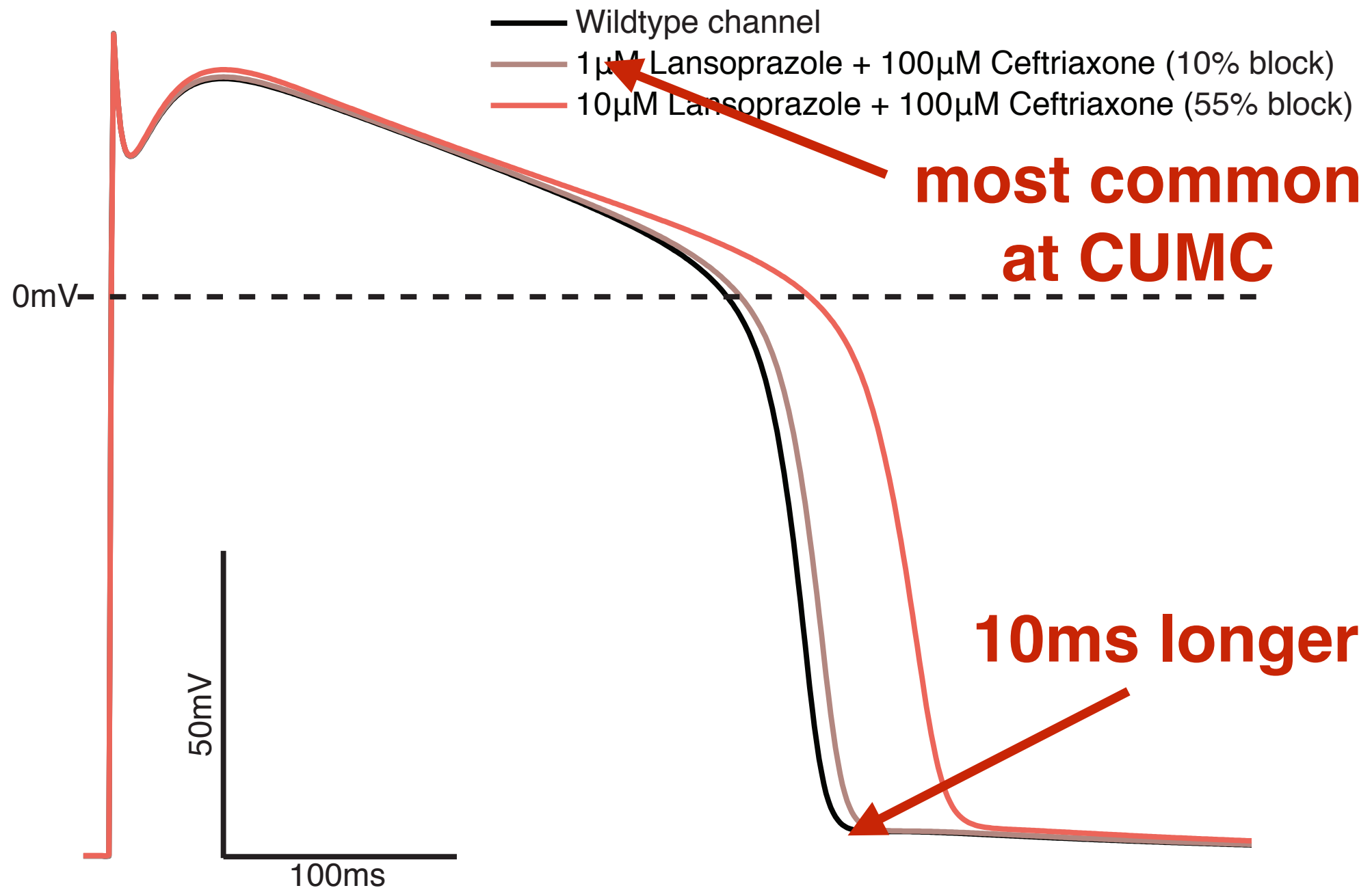# A Computational Model of the Human Left-Ventricular Epicardial Myocyte

Vivek Iyer, Reza Mazhari, and Raimond L. Winslow
The Center for Cardiovascular Bioinformatics and Modeling and the Whitaker Biomedical Engineering Institute, The Johns Hopkins University School of Medicine and Whiting School of Engineering, Baltimore, Maryland

# Computational model of human ventricular myocyte



Legend:
- Wildtype channel
- 1µM Lansoprazole + 100µM Ceftriaxone (10% block)
- 10µM Lansoprazole + 100µM Ceftriaxone (55% block)

0mV

50mV

100ms

*Lorberbaum, et al. JACC (In press)*

# Computational model of human ventricular myocyte



Lorberbaum, et al. JACC (In press)

# Data mining clinical information

- Drug-drug interactions can be discovered using observational data

  - **paroxetine/pravastatin**

  - **ceftriaxone/lansoprazole**

- EHR data accurately predict prospective experiments

# Thank you

tatonettilab.org
nick.tatonetti@columbia.edu
**@nicktatonetti**

## Current Lab Members

Rami Vanguri, PhD
Kayla Quinnies, PhD
Alexandra Jacunski
**Tal Lorberbaum**
Mary Boland
Joseph Romano
Yun Hao
Phyllis Thangaraj
Alexandre Yahi
Fernanda Polubriaginof, MD

**Tal Lorberbaum**
PhD Candidate in Cellular Physiology and Biophysics
Computational biology, systems pharmacology, protein structure modeling

## Collaborators

David Goldstein, PhD
Krzysztof Kiryluk, MD, MS
David Vawdrey, PhD
Robert Kass, PhD
Kevin Sampson, PhD
Brent Stockwell, PhD
George Hripcsak, MD, MS
Ziad Ali, MD, DPhil
Ray Woosley, MD, PhD (Credible Meds)
Konrad Karczewski, PhD (Broad/MGH)
Joel Dudley, PhD (Mount Sinai)
Li Li, PhD (Mount Sinai)
Patrick Ryan, PhD (OHDSI)
Russ Altman (Stanford)
Issac Kohane (HMS)
Shawn Murphy (HMS)

## Funding

**COLUMBIA UNIVERSITY MEDICAL CENTER**
*Discover. Educate. Care. Lead.*