

# #CovidComplete

An Inside Look at Covid-19 Forecasting

By Steve McConnell

# — Copyright

© Copyright Steve McConnell.  
All Rights Reserved.

These materials may be freely copied and distributed for personal, educational, and non-commercial purposes so long as the following attribution is included:

*These materials are © Copyright Steve McConnell.  
Additional information is available at [stevemcconnell.com/covid](http://stevemcconnell.com/covid).*



WORLD

THE WORLD IS  
TEMPORARILY CLOSED



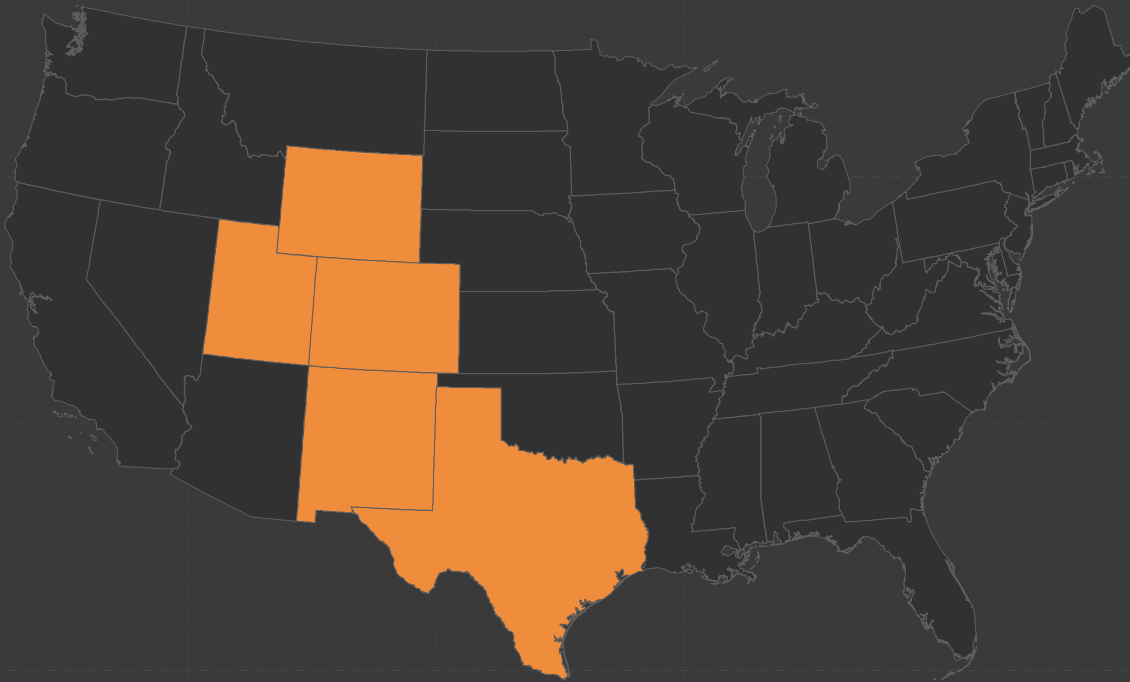
# Goals for this Talk

- ❑ Diffuse emotion and politics; show why you should ignore IHME's forecasts
- ❑ Demystify forecasting (vs. speculation)
- ❑ Explore Covid-19 data issues; correct common data misreporting issues
- ❑ Provide background on the CDC's forecasting process
- ❑ Explain what's possible and what is not possible in Covid-19 forecasting
- ❑ Demonstrate state of the art of current Covid-19 forecasts

# — Why Create State Forecasts?

- ▣ Guides near-term activity planning and event planning (1-4 weeks)
- ▣ Anticipates regional surges
- ▣ Identify areas that will need additional supplies, equipment, staff
- ▣ Identify areas that can spare supplies, equipment, staff
- ▣ Provides data-based feedback on policy decisions (as long as the forecast is based on data rather than the assumptions being tested)

# — Why create National Forecasts?



Guides near-term activity and event planning (1-4 weeks)



Provides data-based feedback on policy decisions



Provides important context for state-level decisions (how much capacity is expected to be available nationally?)



Avoid panic (or create it)



# My personal role in Covid-19 Forecasting

Microsoft

2

Second Edition

# CODE COMPLETE



A practical handbook of software construction

**Steve McConnell**

Two-time winner of the *Software Development Magazine* Jolt Award



For most of my career I have focused on understanding the data analytics of software development, including quality, productivity, and estimation. The techniques I've learned from working with noisy data, bad data, uncertainty, and forecasting all apply to COVID-19.



# SOFTWARE ESTIMATION



*Demystifying the Black Art*

Steve McConnell

Two-time winner of *Software Development* magazine's Jolt Award



Experts tend to use simple estimation strategies, even when their level of expertise in the subject being estimated is high — *Steve McConnell, Software Estimation*

# Coronavirus Disease 2019 (COVID-19)



- Your Health
- Community, Work & School
- Healthcare Workers & Labs
- Health Depts
- Cases & Data
- More

## Cases, Data & Surveillance

US Cases & Deaths +

Cases & Deaths by County

Testing Data in the US +

Hospitalizations & +

### CASES, DATA & SURVEILLANCE

# COVID-19 Forecasts: Deaths

Updated Sept. 30, 2020 [Print](#)



Observed and forecasted new and total reported COVID-19 deaths as of September 28, 2020.

On This Page

## Interpretation of Forecasts of

Special Populations Data +

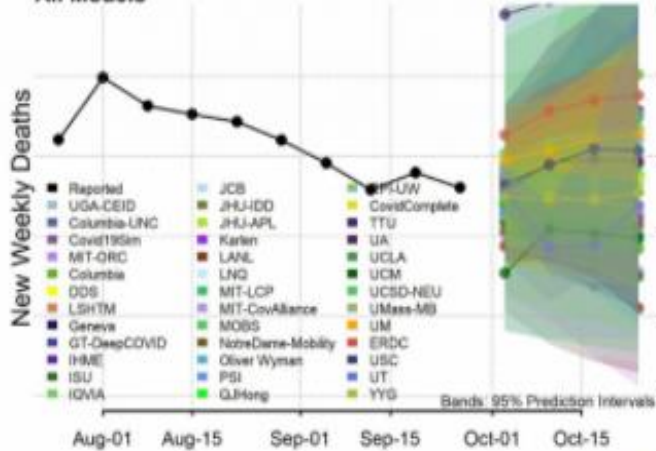
COVIDView Weekly Summary +

Sequencing for SARS-CoV-2 (SPHERES)

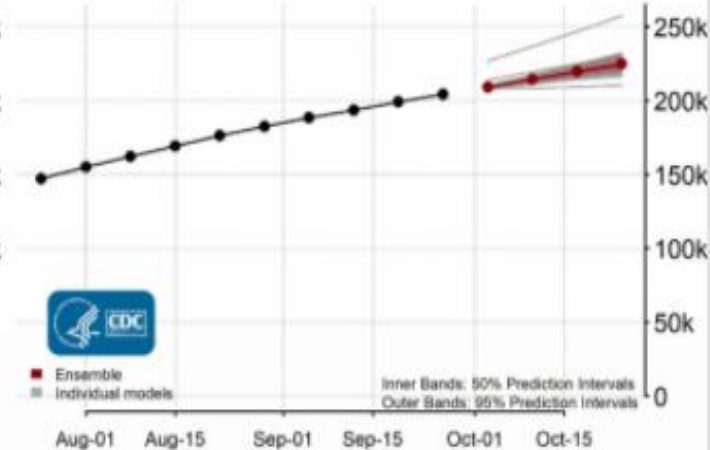
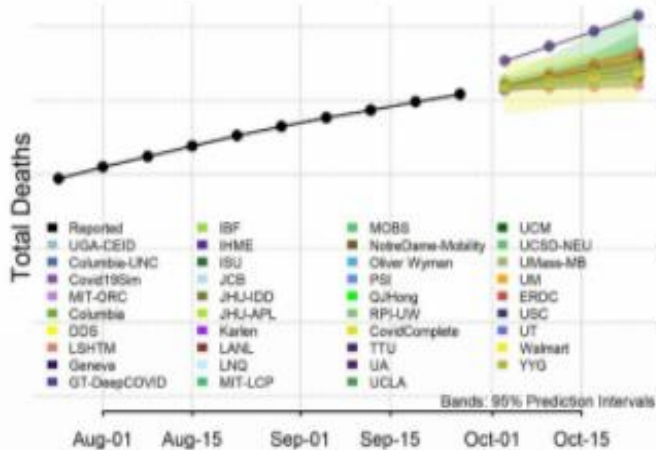
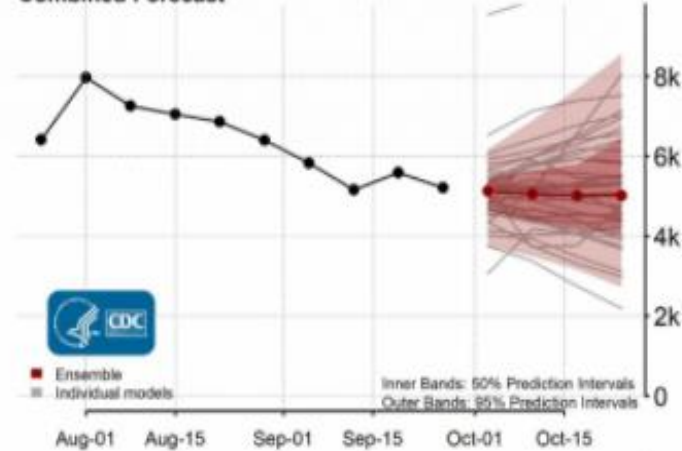
Epidemiology for COVID-19 +

# National Forecast

## National Forecast All Models



## Combined Forecast

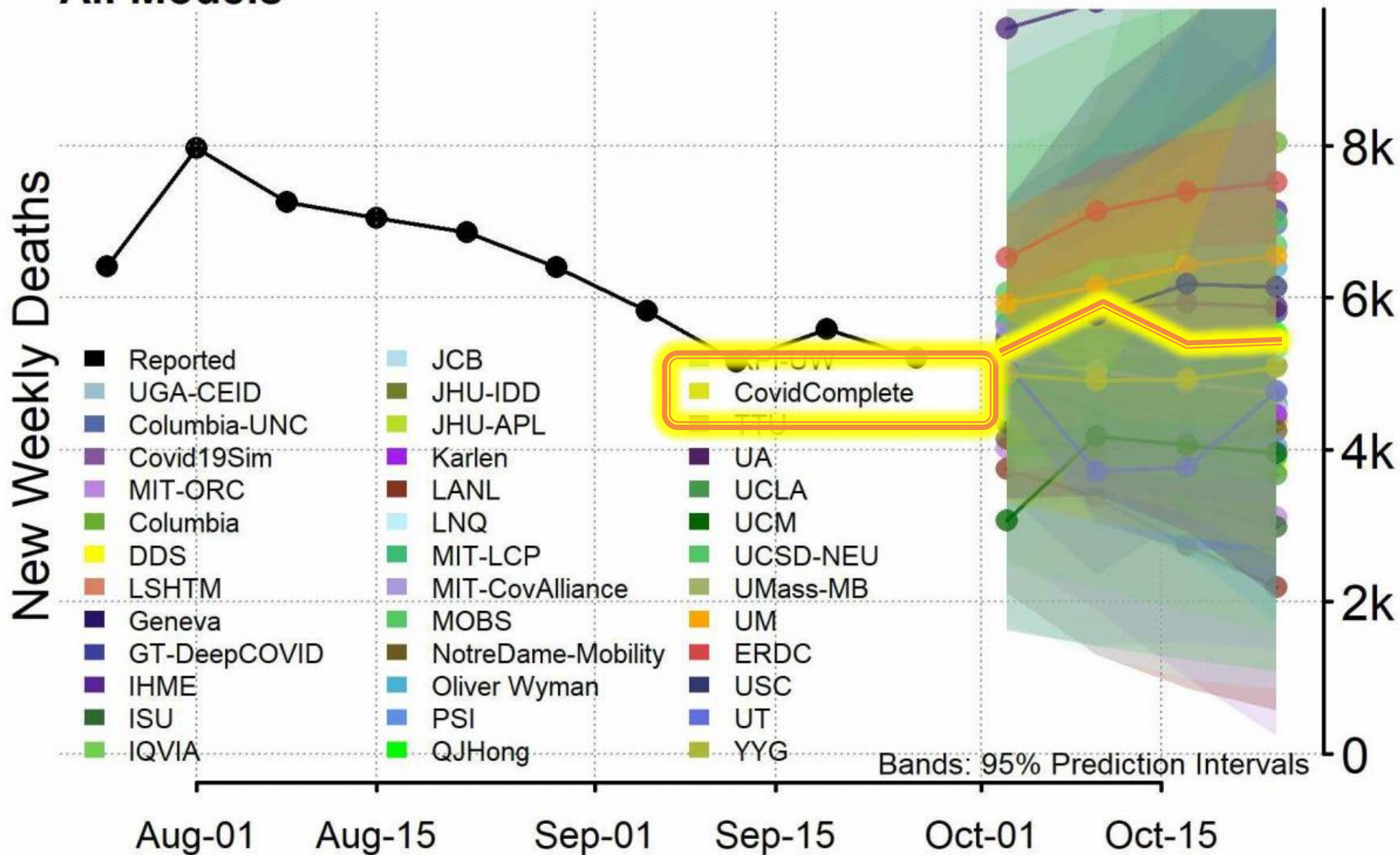


Get Email Updates

To receive email updates about COVID-19, enter your email address:

# National Forecast

All Models





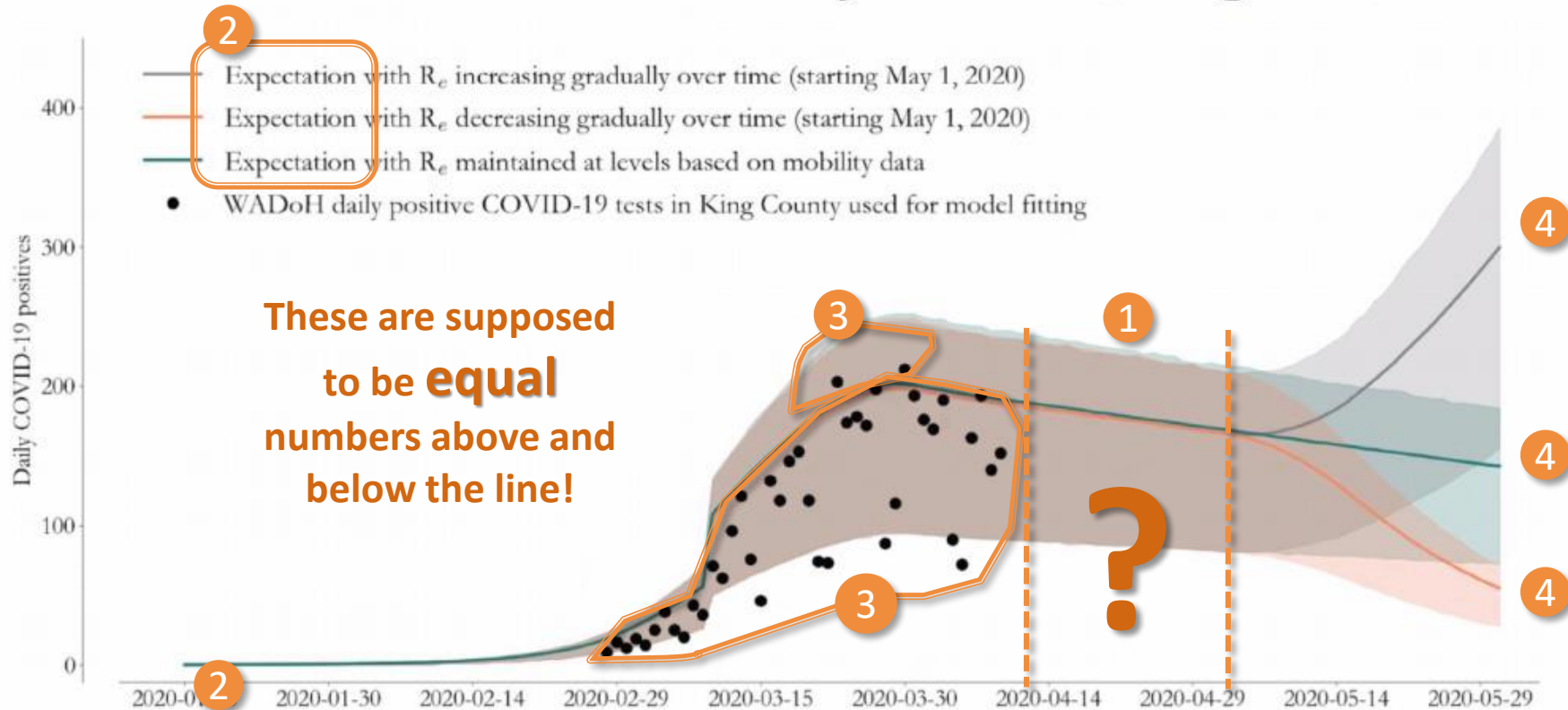
Why I got Involved



I don't like playing games with forecasts during a global pandemic

# — Playing games with data – earlier

## COVID-19 Case Projections (King Co.)



Annotation and critique by Steve McConnell.

Copyright © 2019 Intellectual Ventures Management, LLC (IVM). All rights reserved.

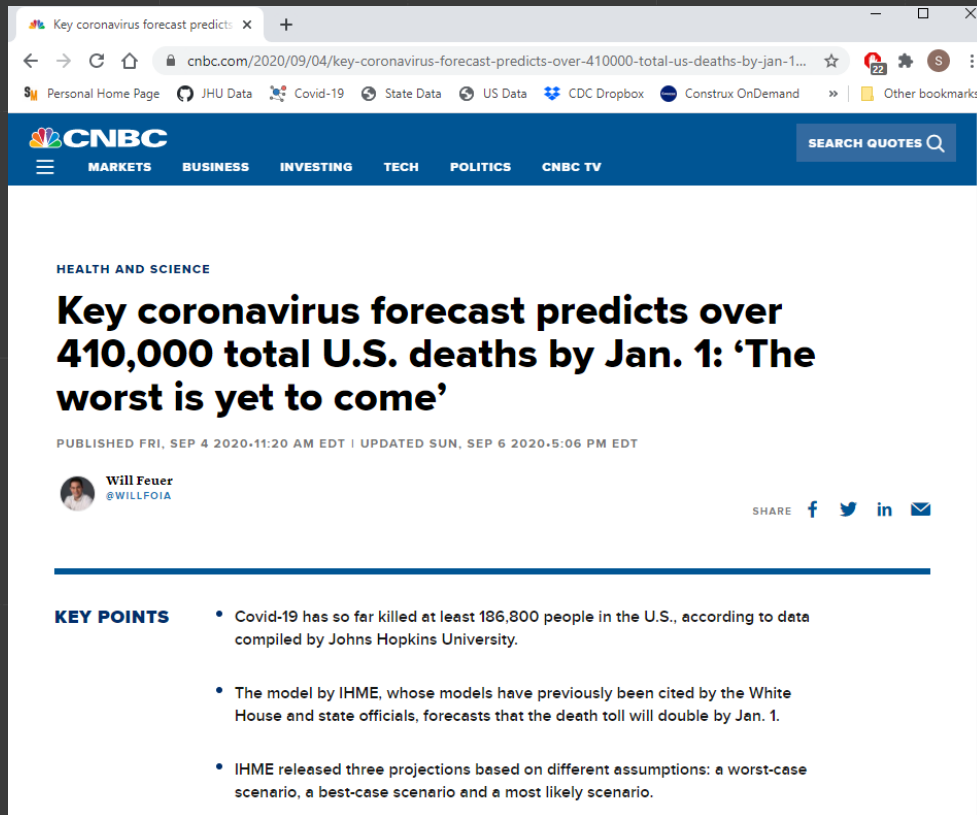


# Playing games with Data, Recently

The screenshot shows a web browser window with the URL [cnbc.com/2020/09/04/key-coronavirus-forecast-predicts-over-410000-total-us-deaths-by-jan-1...](https://www.cnbc.com/2020/09/04/key-coronavirus-forecast-predicts-over-410000-total-us-deaths-by-jan-1...). The page features the CNBC logo and navigation menu. The main headline is "Key coronavirus forecast predicts over 410,000 total U.S. deaths by Jan. 1: 'The worst is yet to come'", with "410,000 total U.S. deaths" underlined in orange. The article is categorized under "HEALTH AND SCIENCE" and was published on Friday, September 4, 2020, at 11:20 AM EDT, with an update on Sunday, September 6, 2020, at 5:06 PM EDT. The author is Will Feuer (@WILLFOIA). Social sharing icons for Facebook, Twitter, LinkedIn, and Email are visible.

scenario, a best-case scenario and a most likely scenario. The most likely estimates that Covid-19 will kill 410,450 people in the U.S. by Jan. 1.

- The model by IHME, whose models have previously been cited by the White House and state officials, forecasts that the death toll will double by Jan. 1.
- IHME released three projections based on different assumptions: a worst-case scenario, a best-case scenario and a most likely scenario.



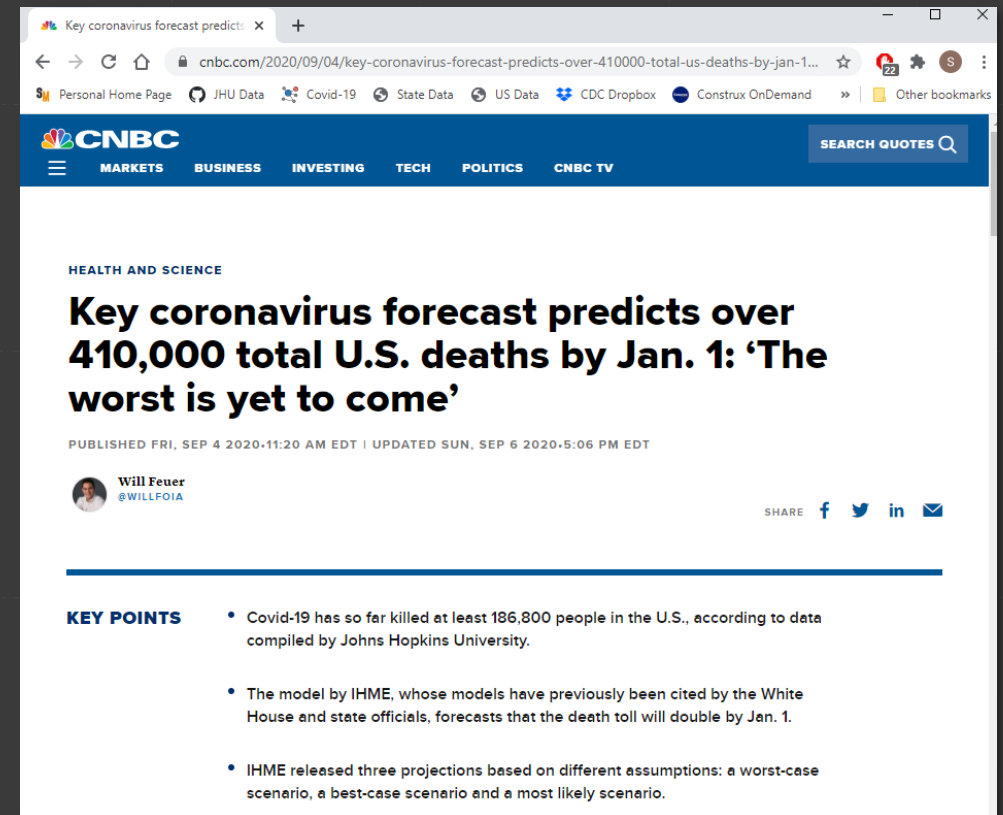
# What would need to happen for this to be true?

- The death trend would need to more than double, immediately, despite trending down 10%/week for the past few weeks, along with leading indicators also trending down
- Need to incur more deaths in four months than we've had since the beginning of the pandemic



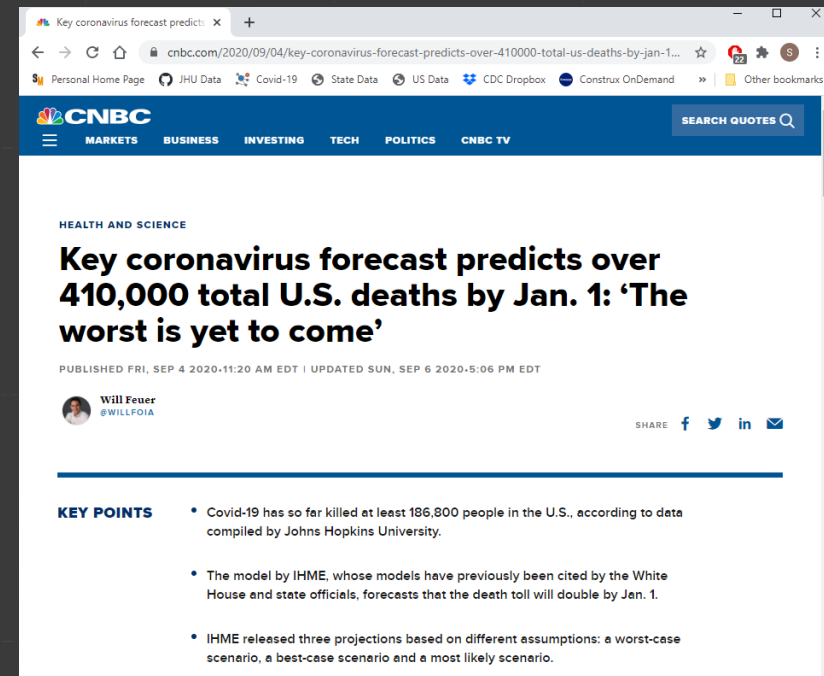
# What would need to happen for this to be true?

- Need to average 1800 deaths per day for 115 days straight, which was 100 more days than we'd had that many deaths so far
- No one can take any corrective action to reduce deaths even after deaths skyrocket, and they must continue not to take any action for **4 months**

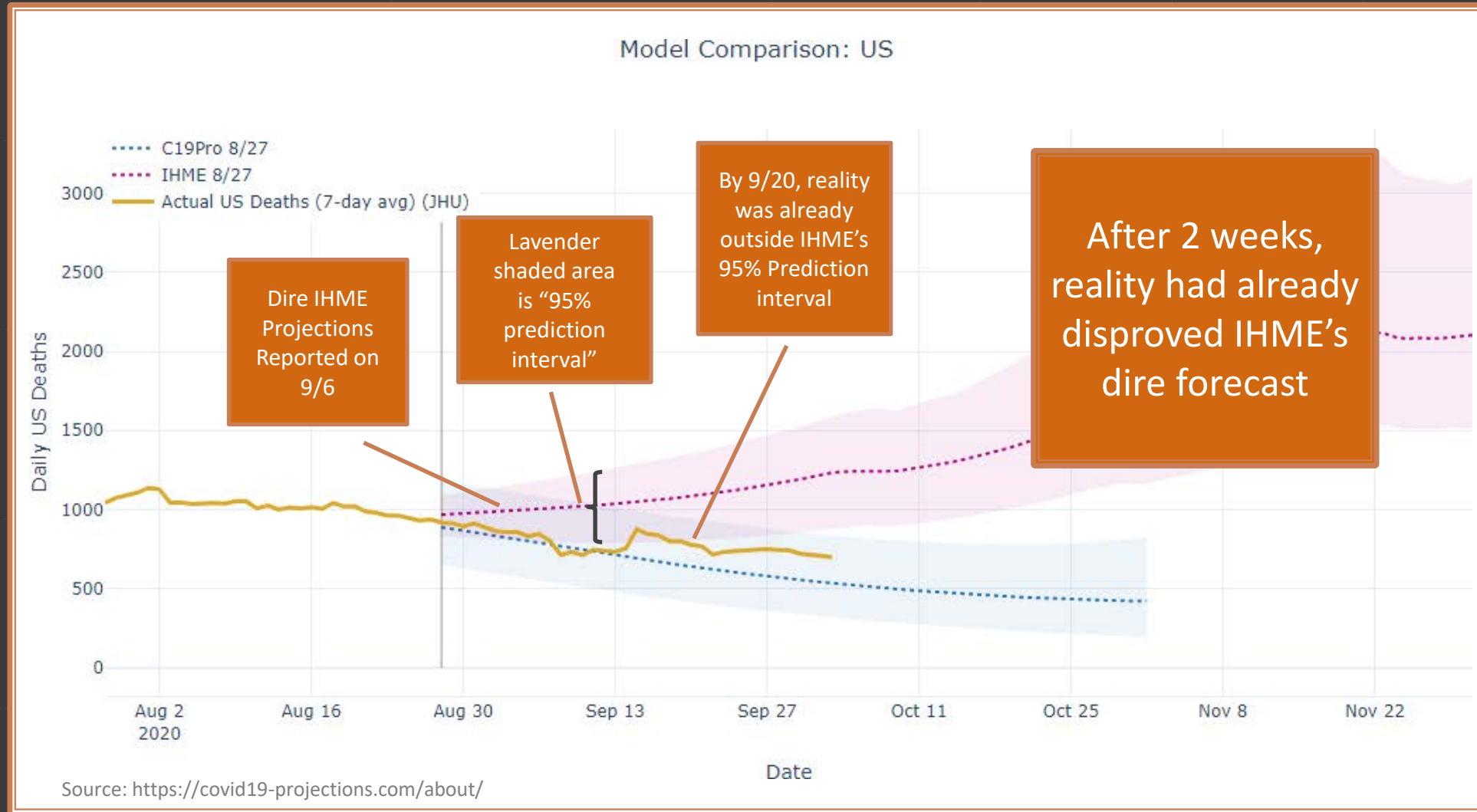


# — My Observations

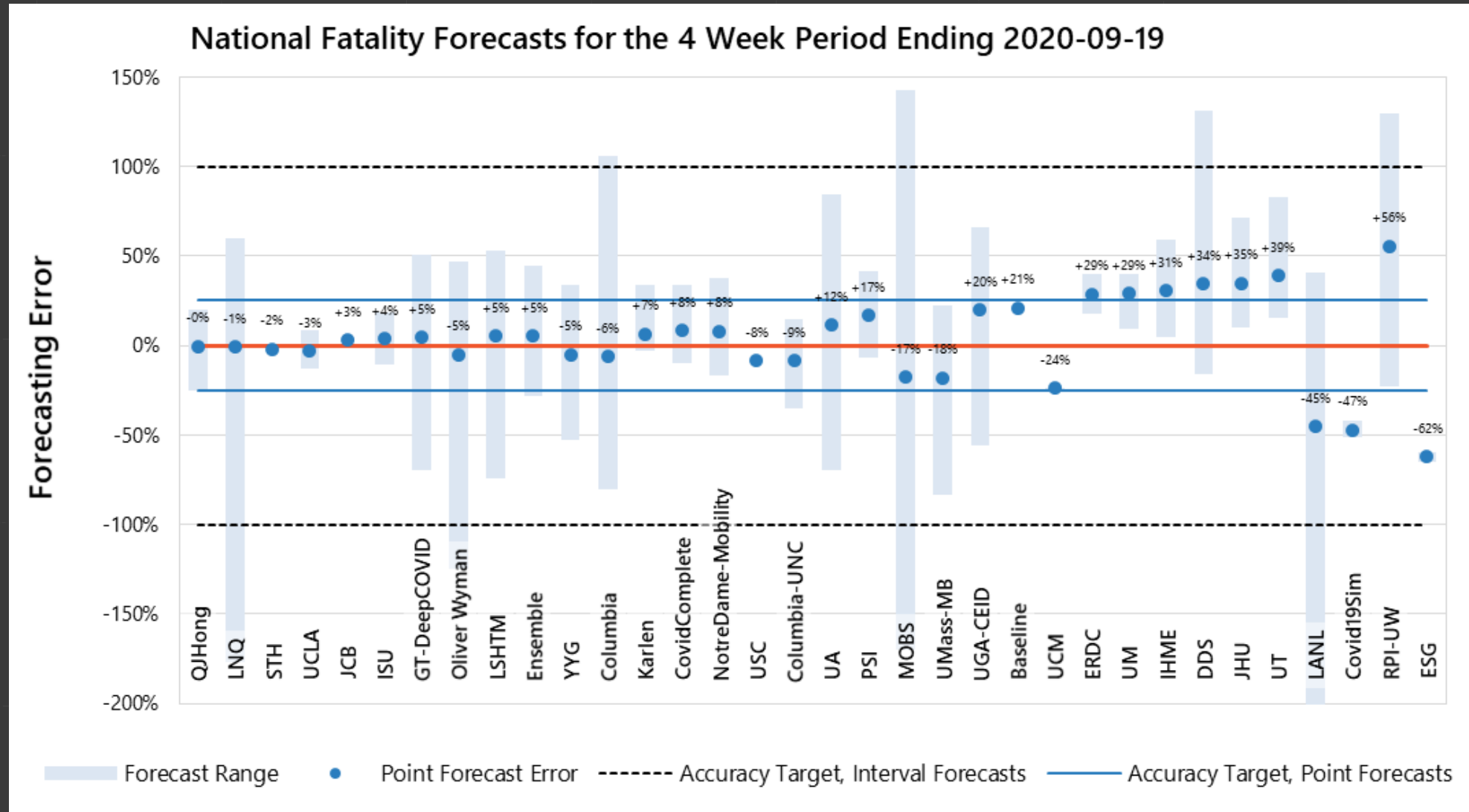
- ❑ This forecast is not going to happen
- ❑ This forecast is completely unfounded
- ❑ This forecast is in no way “most likely”



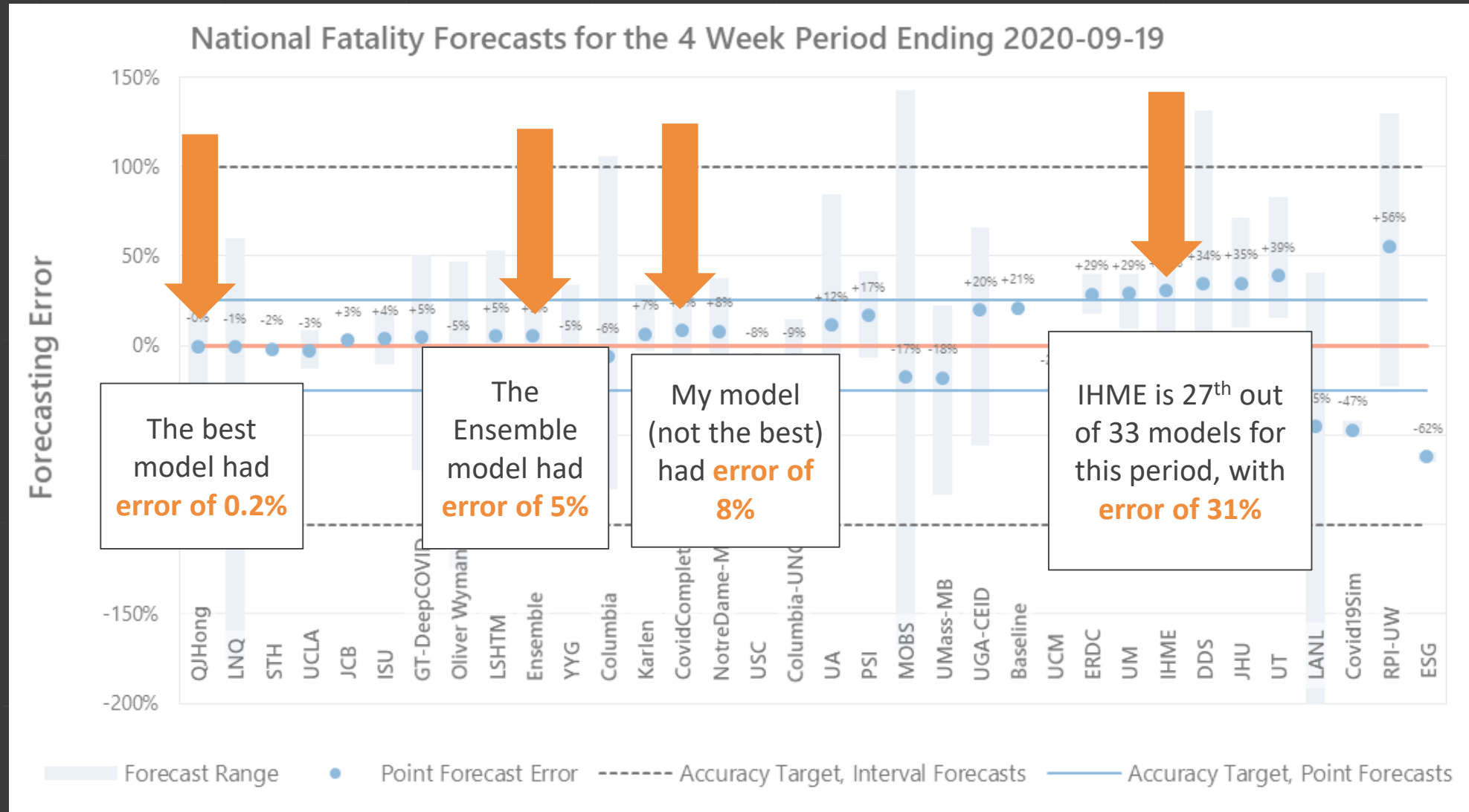
# — This “Forecast” is completely unfounded



# How Accurate were Other Forecast for the same period?



# How Accurate were Other Forecast for the same period?



www.stevemcconnell.com



Why I got Involved

I just wanted to know what is going on with  
the pandemic, without all the spin



- The Basics: What is Forecasting?

Forecasting is not the same as speculation, but we are seeing them treated interchangeably





Coronavirus Latest news U.S. map World map FAQ Vaccine tracker Coronavirus Living Extraordinary People

PostEverything • Perspective

# Scientists want to predict covid-19's long-term trajectory. Here's why they can't.

Our research suggests that forecasts are unreliable further than a few weeks out



Dining outdoors in Manhattan on Monday. (Jeenah Moon/Reuters)

By **Nicholas Reich** and **Caitlin Rivers**

September 15, 2020 at 11:22 a.m. PDT

## What is possible with forecasting

- ✓ 1-4 week time horizons
- ✓ Accurate national forecasts
- ✓ Pretty accurate state forecasts, especially for more active states
- Anything further out than that is speculation (even if it's dressed up with graphs and numbers, which most of it is)

# Forecasting vs. Speculation

- **Forecasting**
  - Inputs are known
  - Outputs (forecasts) are calculated based on known relationships
- **Speculation**
  - Inputs are guesses
  - Outputs are guesses piled on top of other guesses



# Knowns

---

- Number of people who have already tested positive
- Approximate relation between positive tests and cases
- Approximate fatality rates by age and comorbidity status
- Progression of the virus in individuals

# Unknowns

---

- Specific policies that specific states will implement (or relax) in the future
- Dates those policies will be implemented or relaxed
- Effectiveness of the policies, with good compliance
- Effectiveness of the policies, with whatever compliance we actually get
- Availability of vaccines
- Effectiveness of vaccines, when and if available

The Basis of Forecasting  
*“What We Know”*



# We Know Much More About Covid-19 Now Than We Did on 1/22/20



Fatality is highly age-related (~50% of deaths are age 80+)



Fatality is highly co-morbidity related (~90% of fatalities involve at least one co-morbidity)



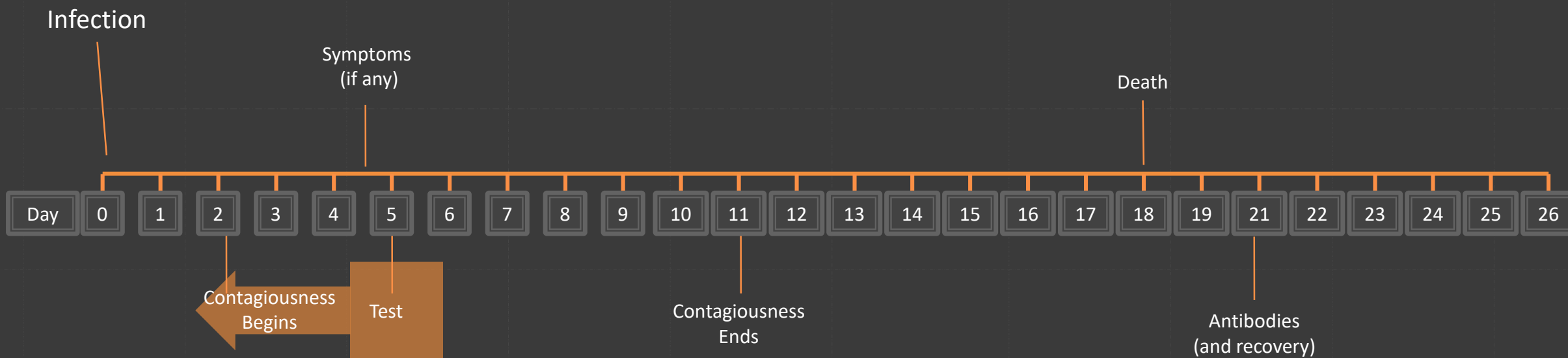
A lot more people are infected than have tested positive

- Early in the pandemic: 10-20x cases per positive test
- Now—3-5x cases per positive test

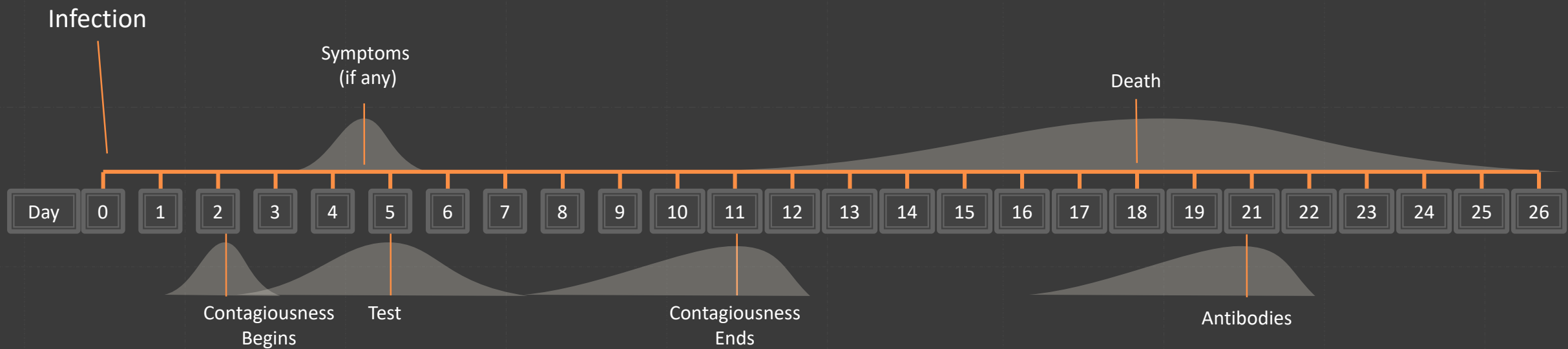


We know the timing of the course of the disease

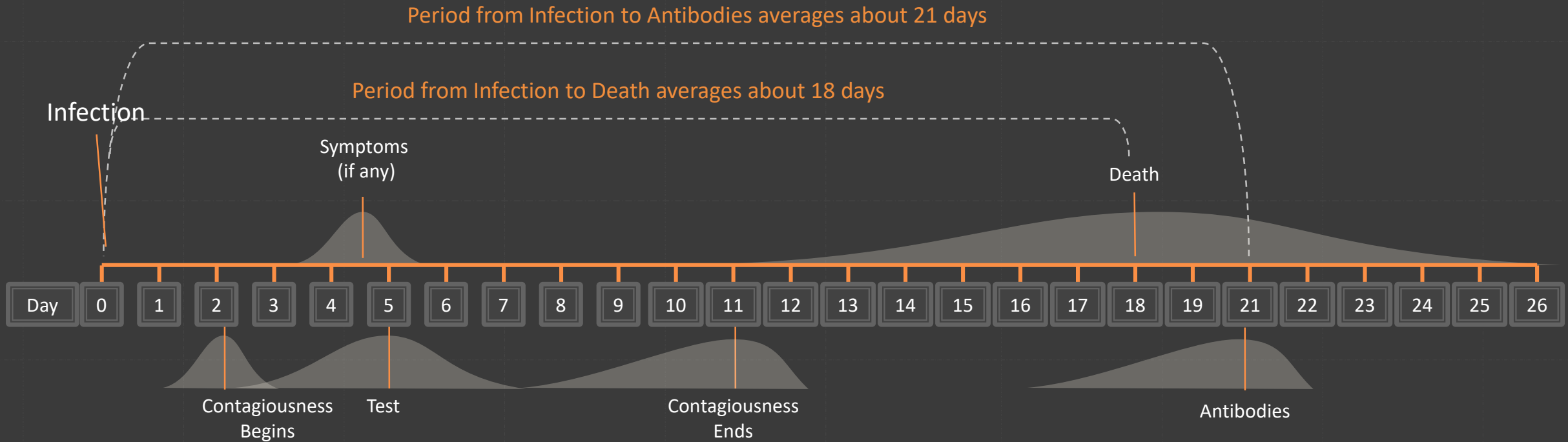
# — We have a pretty good view of the timing of the disease



— Of course there's variability in all of this



# — This allows us to draw some inferences

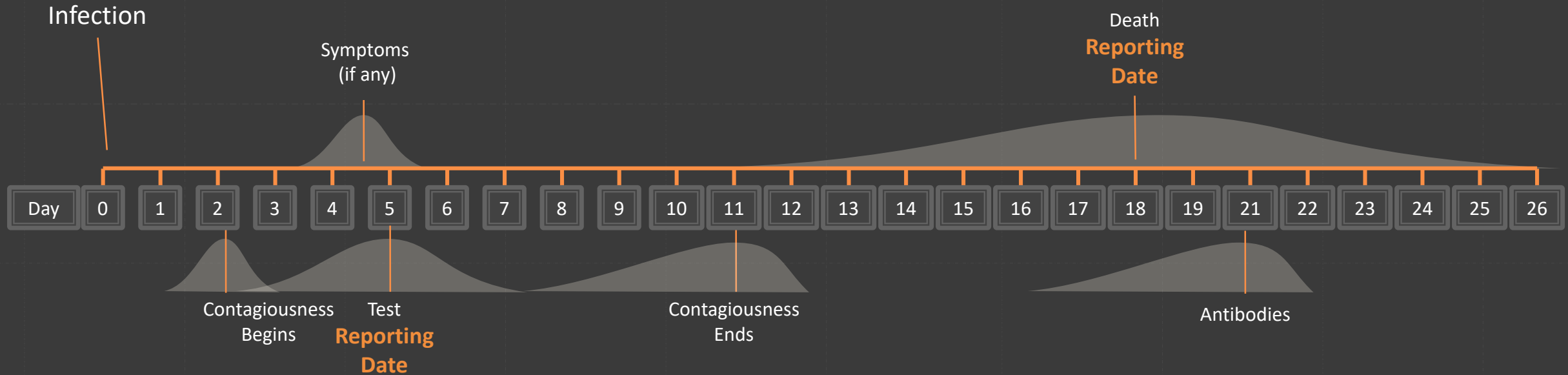




# — This allows us to draw some inferences



# This is more about reporting dates, more than actual dates



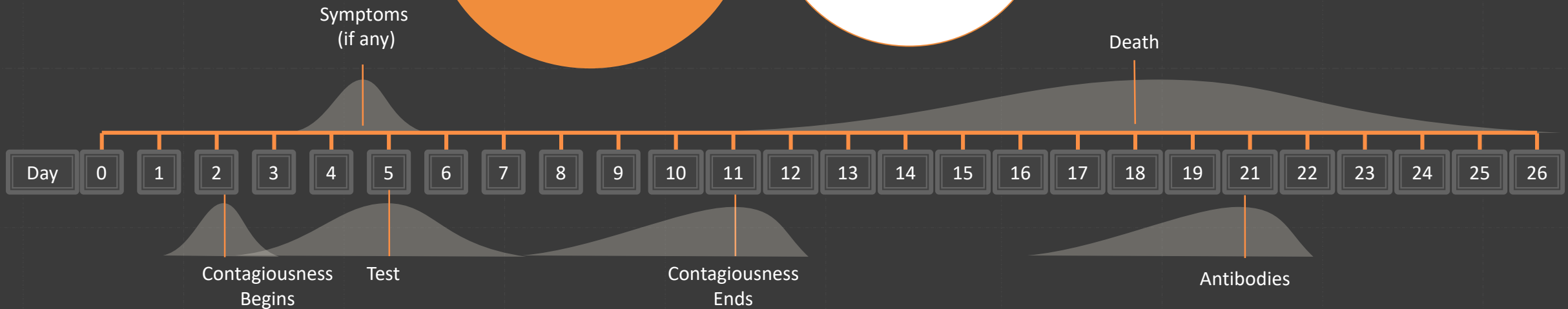


What Does This Have to do With  
Forecasting?

# Timing forms the basis of forecasting (for my method)

We can do this with terrific accuracy at the national level, and with fair accuracy at the state level

And because trends tend to continue, we can project forward another 14 days



We can take this number, which we know on date X

We can apply what we've learned about the ratio of positive tests to deaths

We can forecast this number, **14 days ahead**

— Forecasting  
Begins with Data

# Current State of Data, Part 1



# Current State of Data, Part 2





# Problems with Cumulative Data

Issues with cumulative data, examples

- Bing/covid graph
- JHU graphic
- What is cumulative data good for?

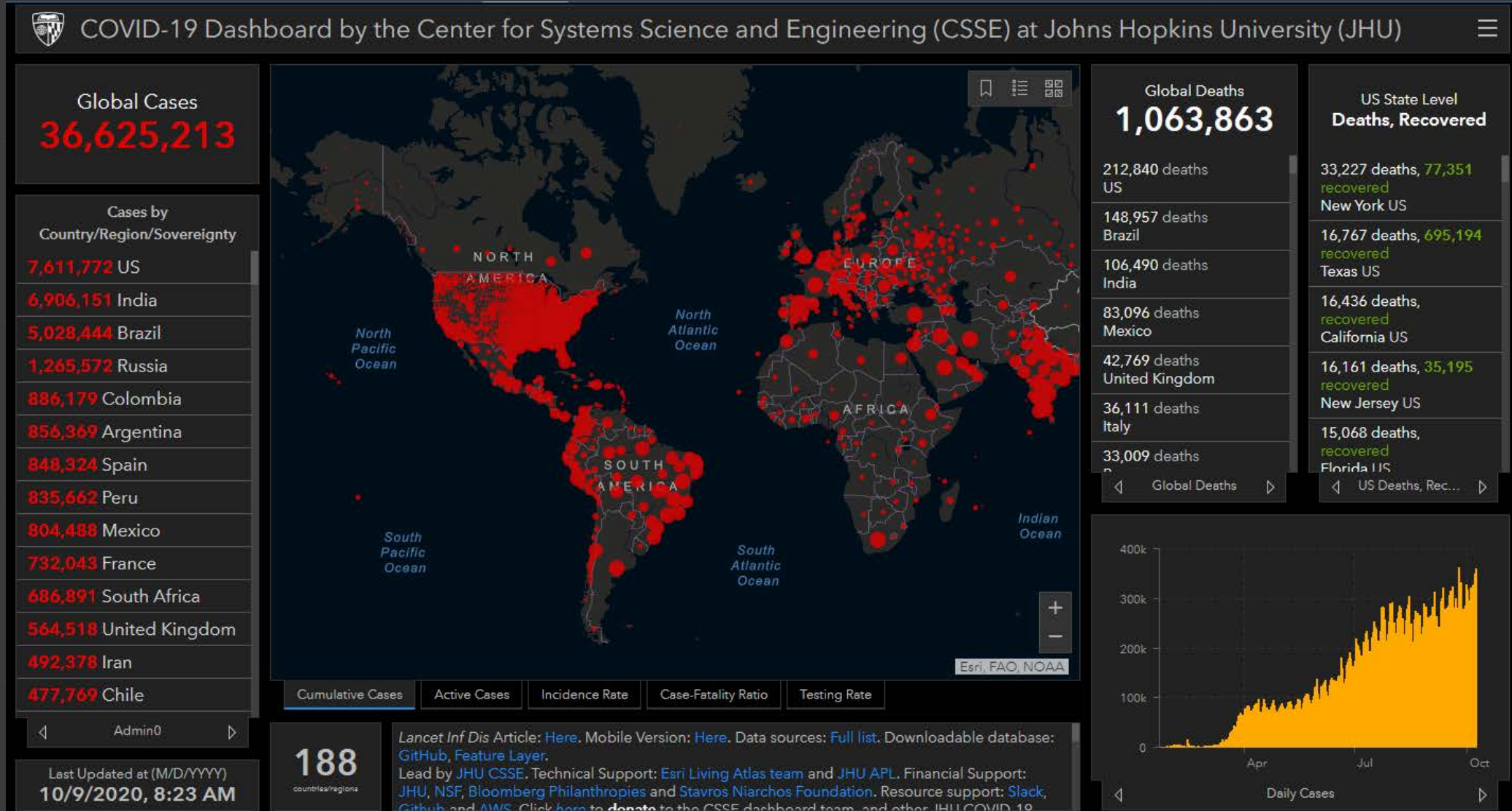


# — Typical Cumulative Graph

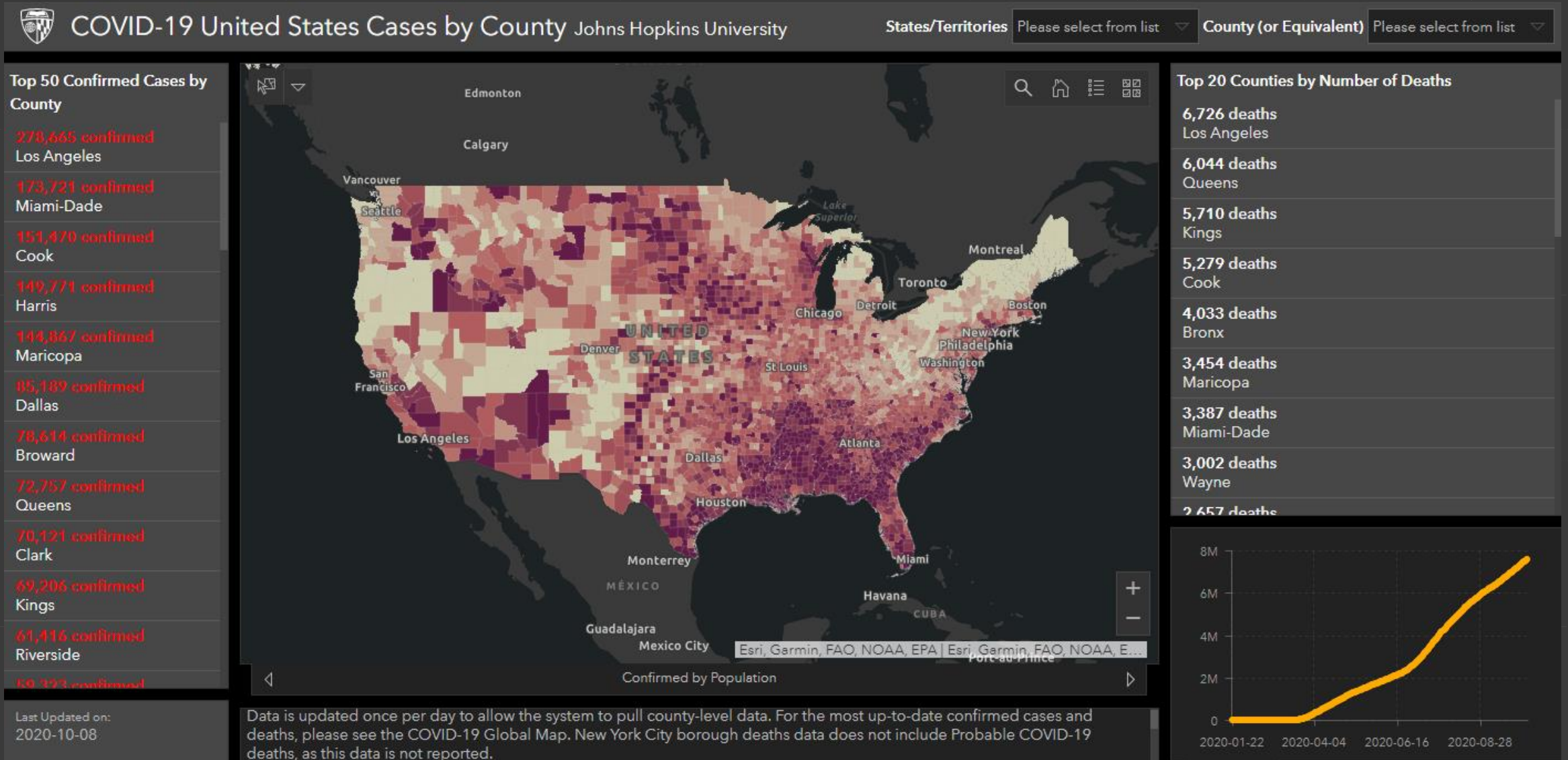


Source: <https://www.bing.com/covid/local/unitedstates?vert=graph>

# Johns Hopkins Coronavirus Dashboard



# Johns Hopkins Coronavirus Dashboard – Newer Graphic



Source: <https://coronavirus.jhu.edu/us-map>

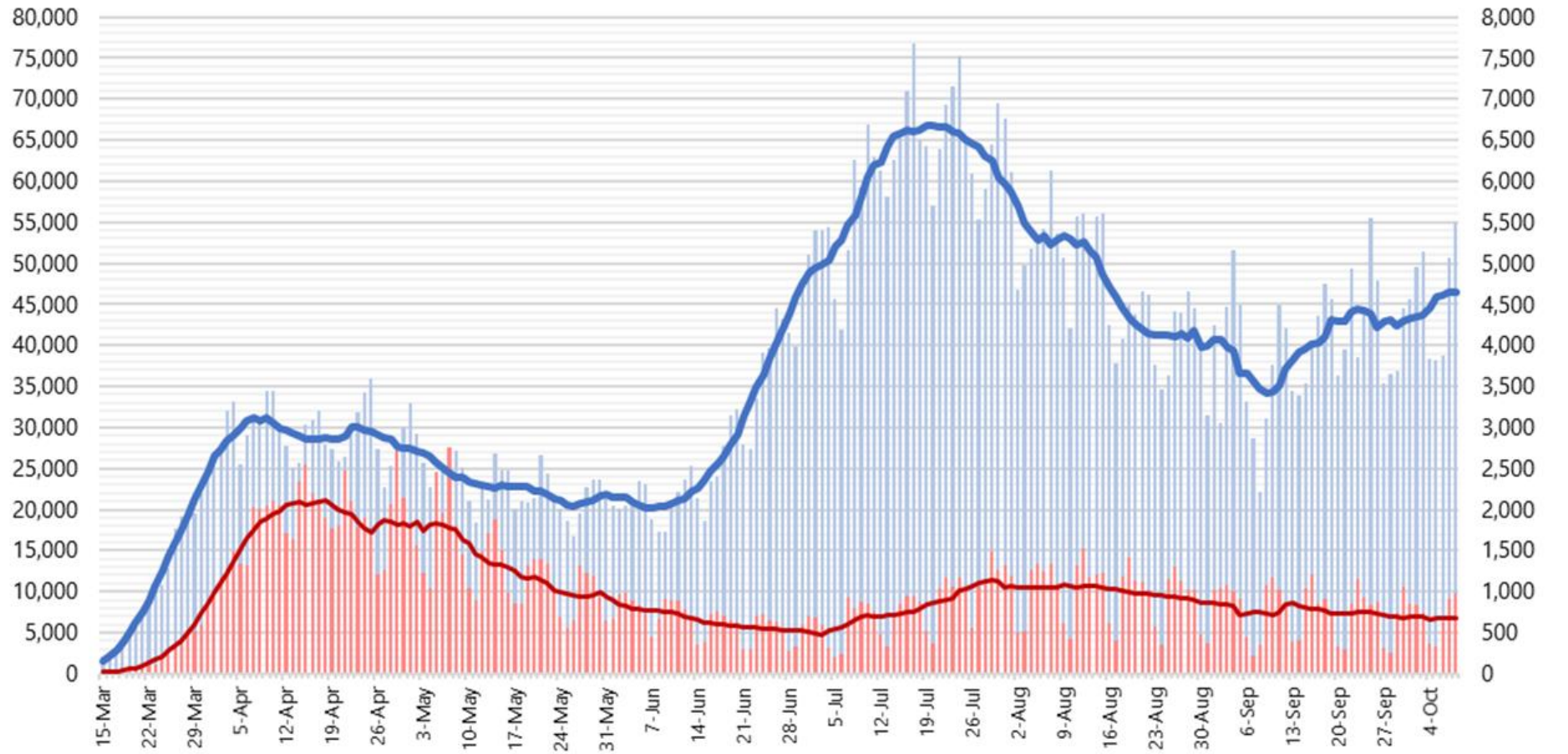


# Incremental Data

- ▣ Incremental data is more meaningful than cumulative data
- ▣ But it also opens the door to a whole new set of data quality issues

# Daily Graphic

US Daily Positive Tests and Deaths as of 10/8/20

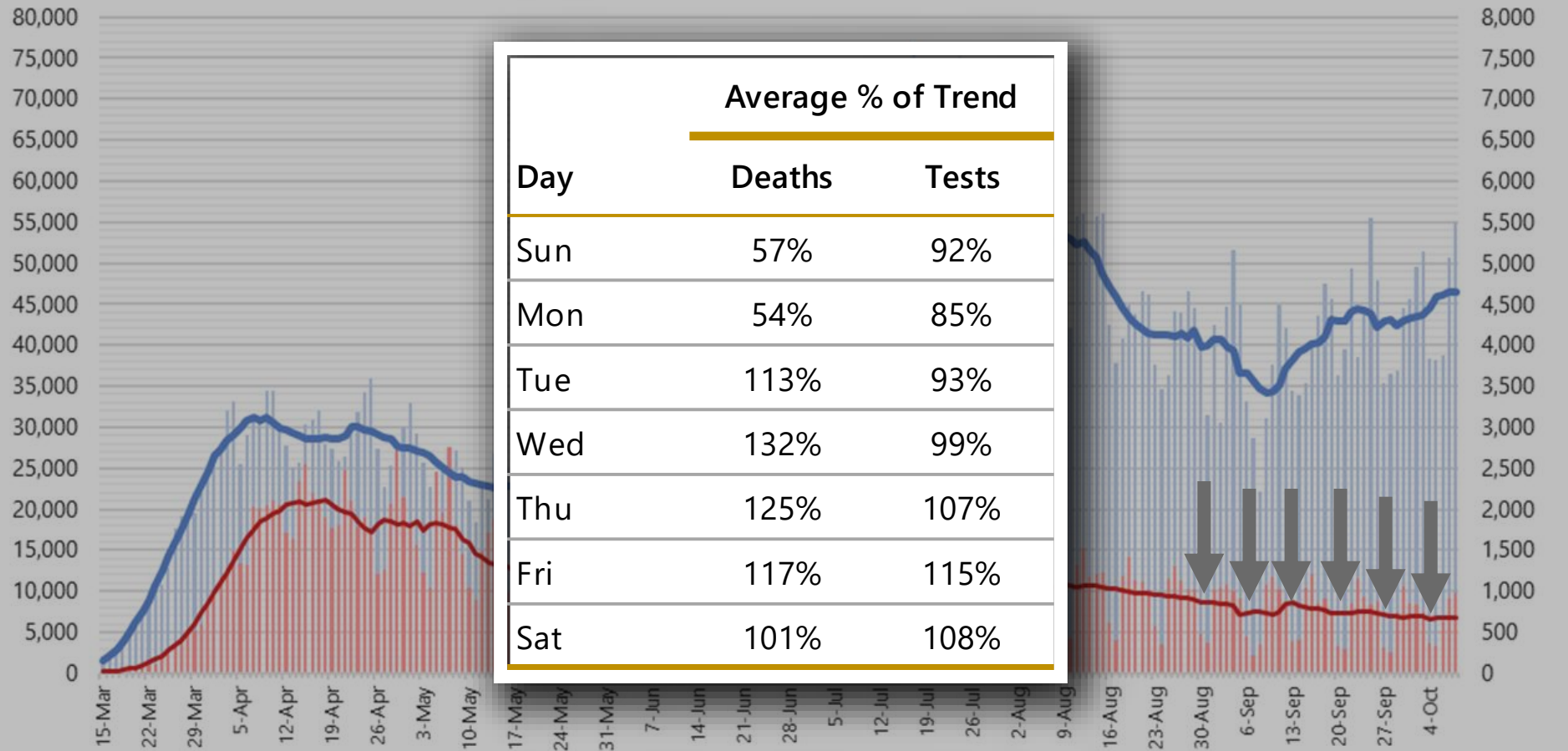


Source: [stevemccconnell.com/covid](http://stevemccconnell.com/covid)

Positive Tests Deaths 7-Day Average Positive Tests 7-Day Average Deaths

# — Issue #1 with Daily Data – Sundays!

US Daily Positive Tests and Deaths as of 10/8/20

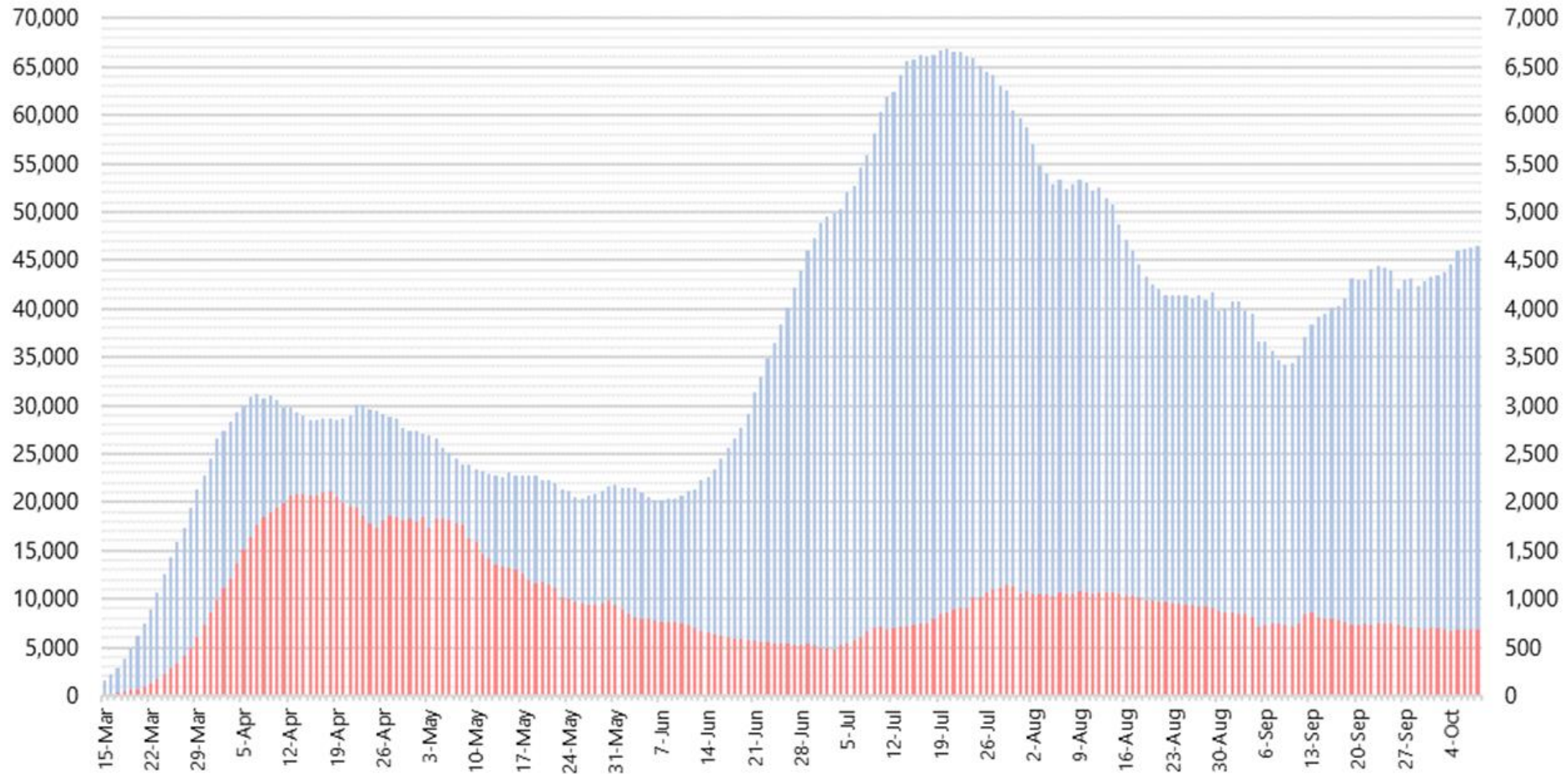


Source: [stevemccconnell.com/covid](http://stevemccconnell.com/covid)

■ Positive Tests 
 ■ Deaths 
 — 7-Day Average Positive Tests 
 — 7-Day Average Deaths

# Issue #1 with Daily Data – Sundays!

US Smoothed Daily Positive Tests and Deaths as of 10/8/20

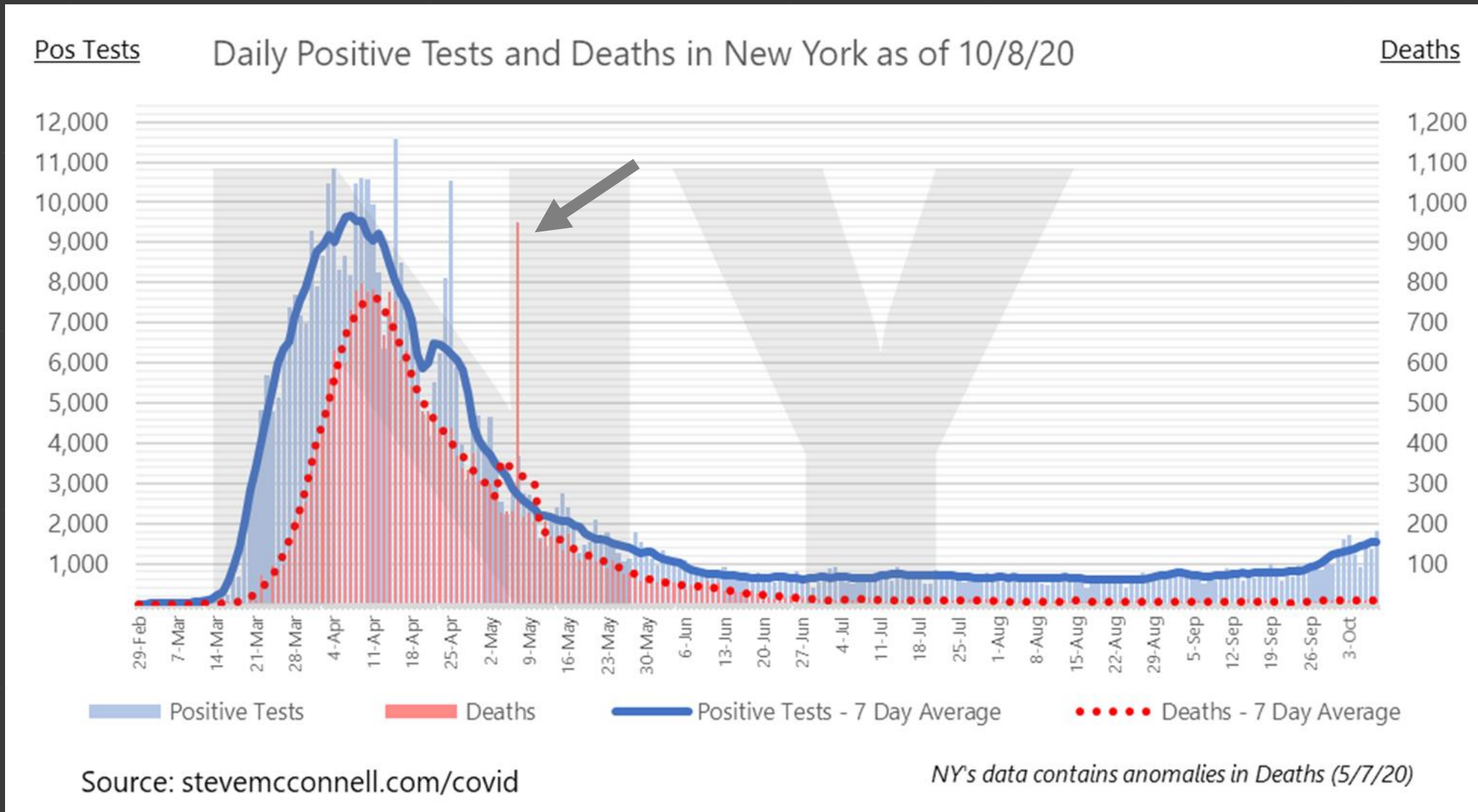


Source: [stevemccconnell.com/covid](http://stevemccconnell.com/covid)

■ Positive Tests

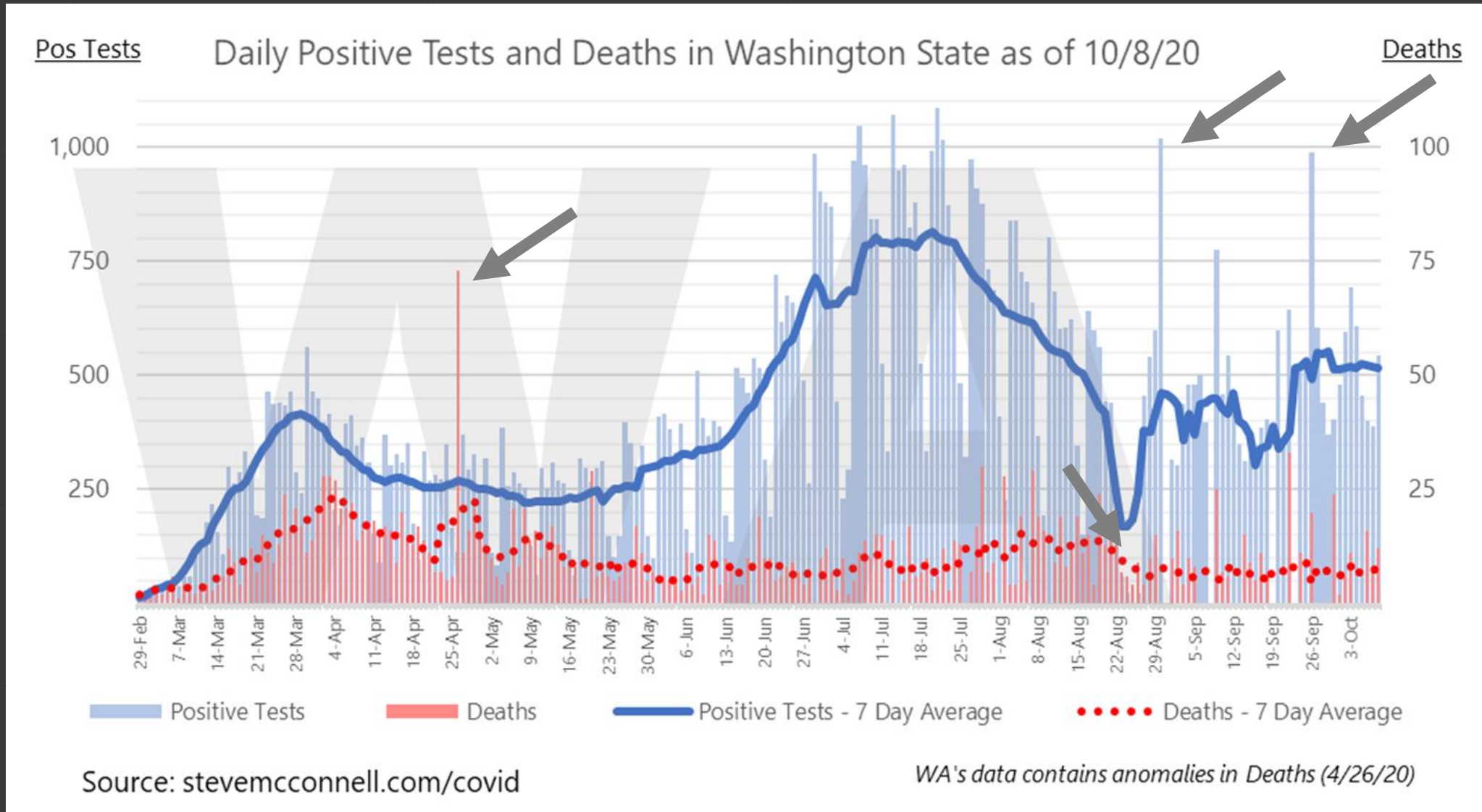
■ Deaths

# Issue #2 with Daily Data – Spikes in State Level Data

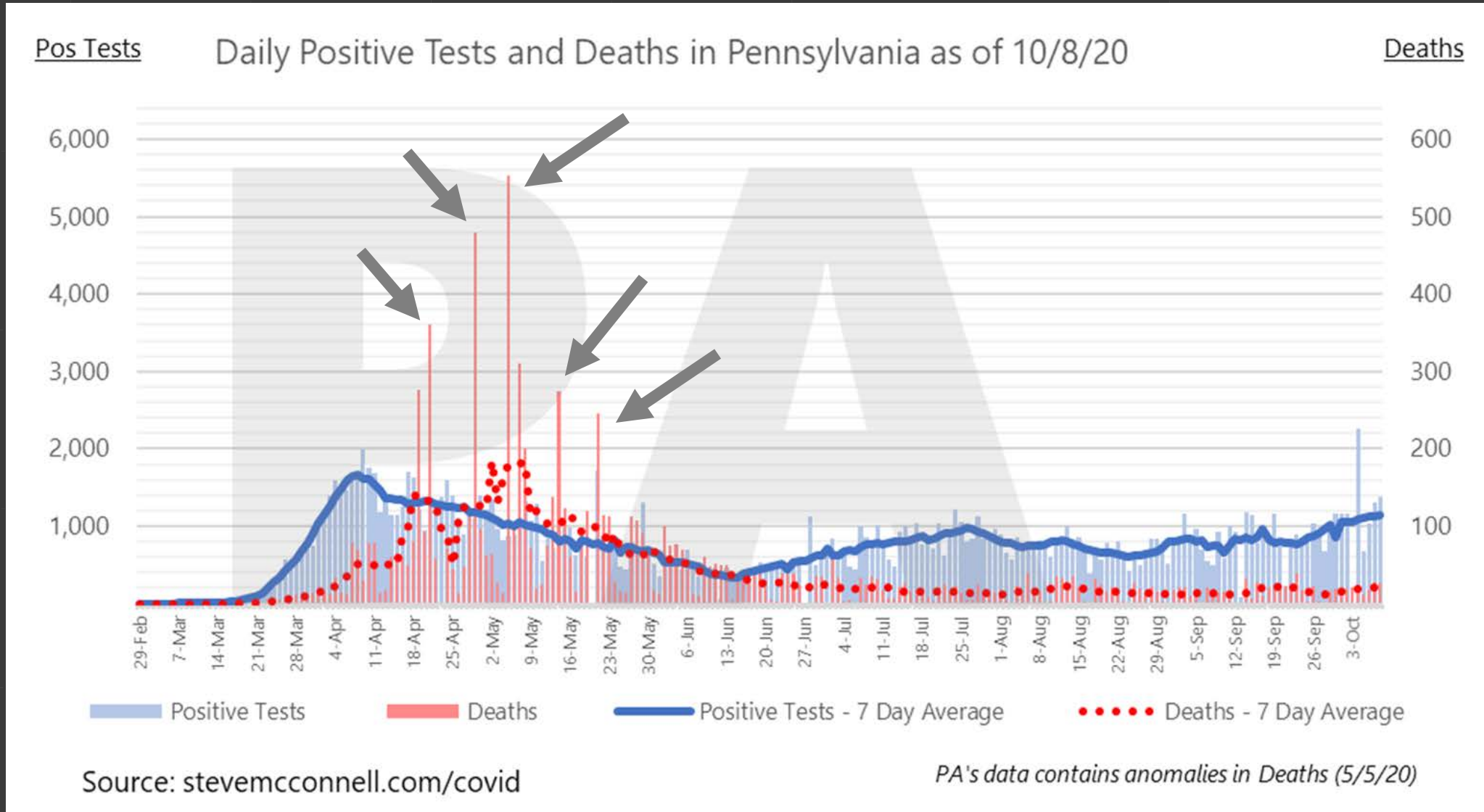




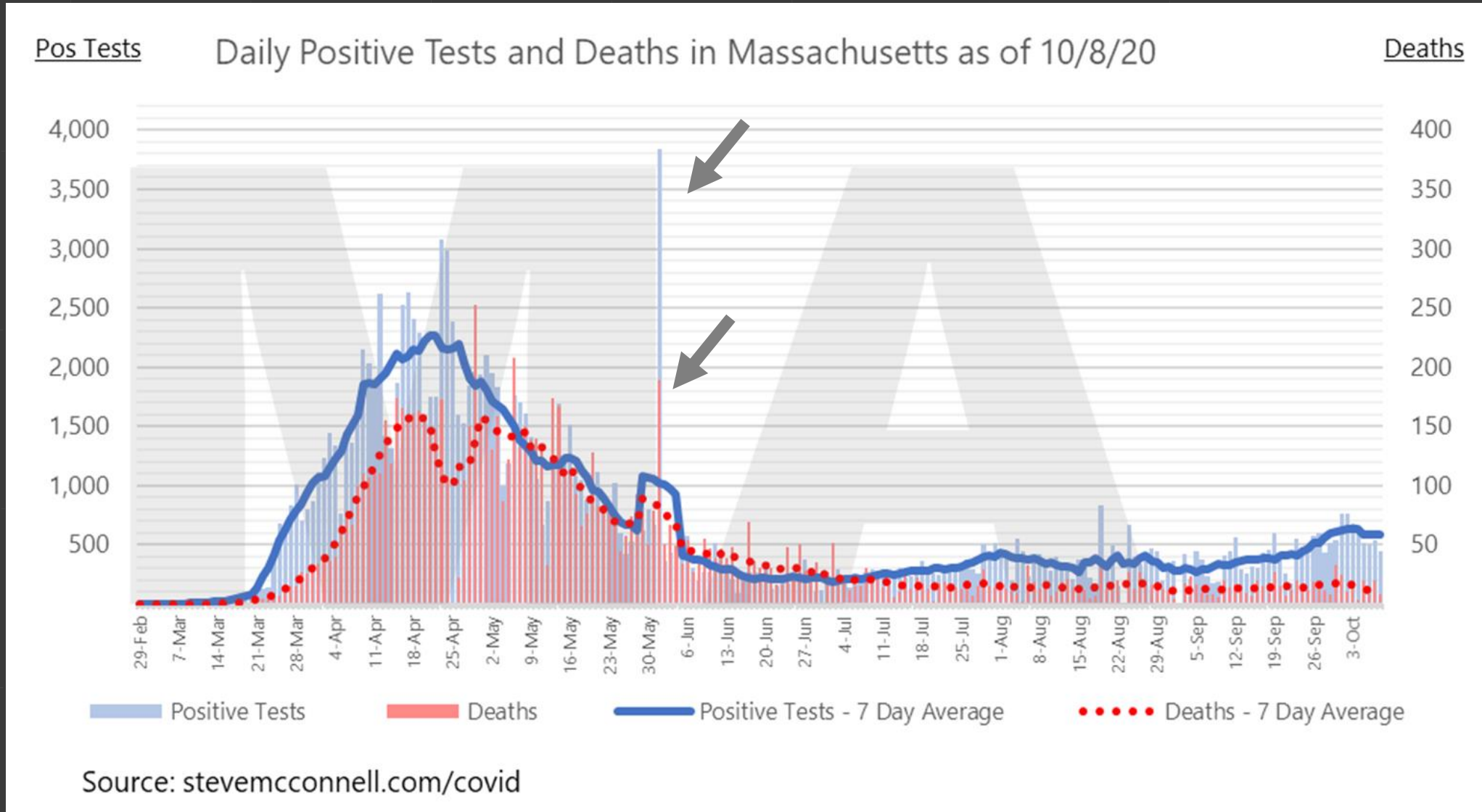
# Issue #2 with Daily Data – State Level Data



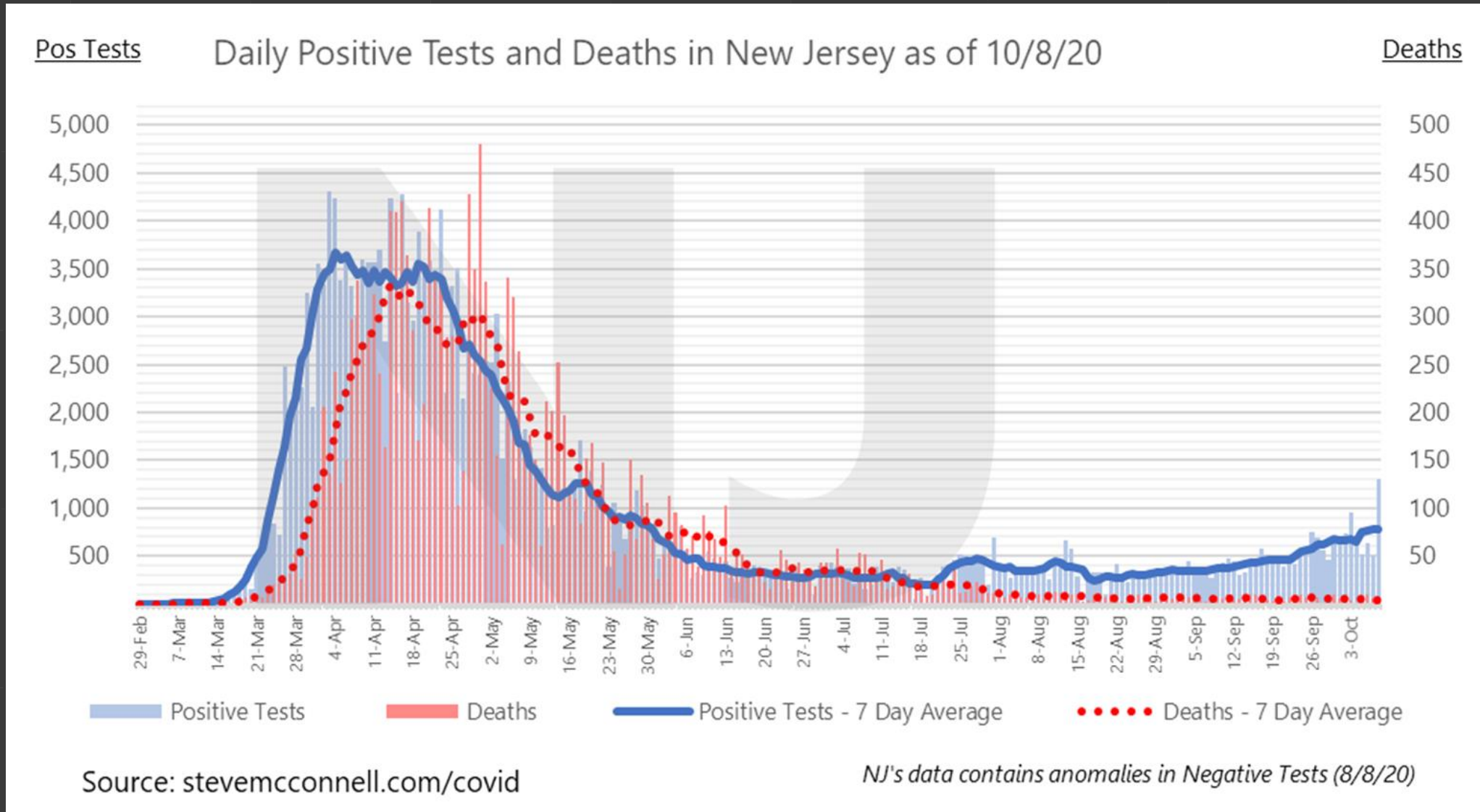
# Issue #2 with Daily Data – State Level Data



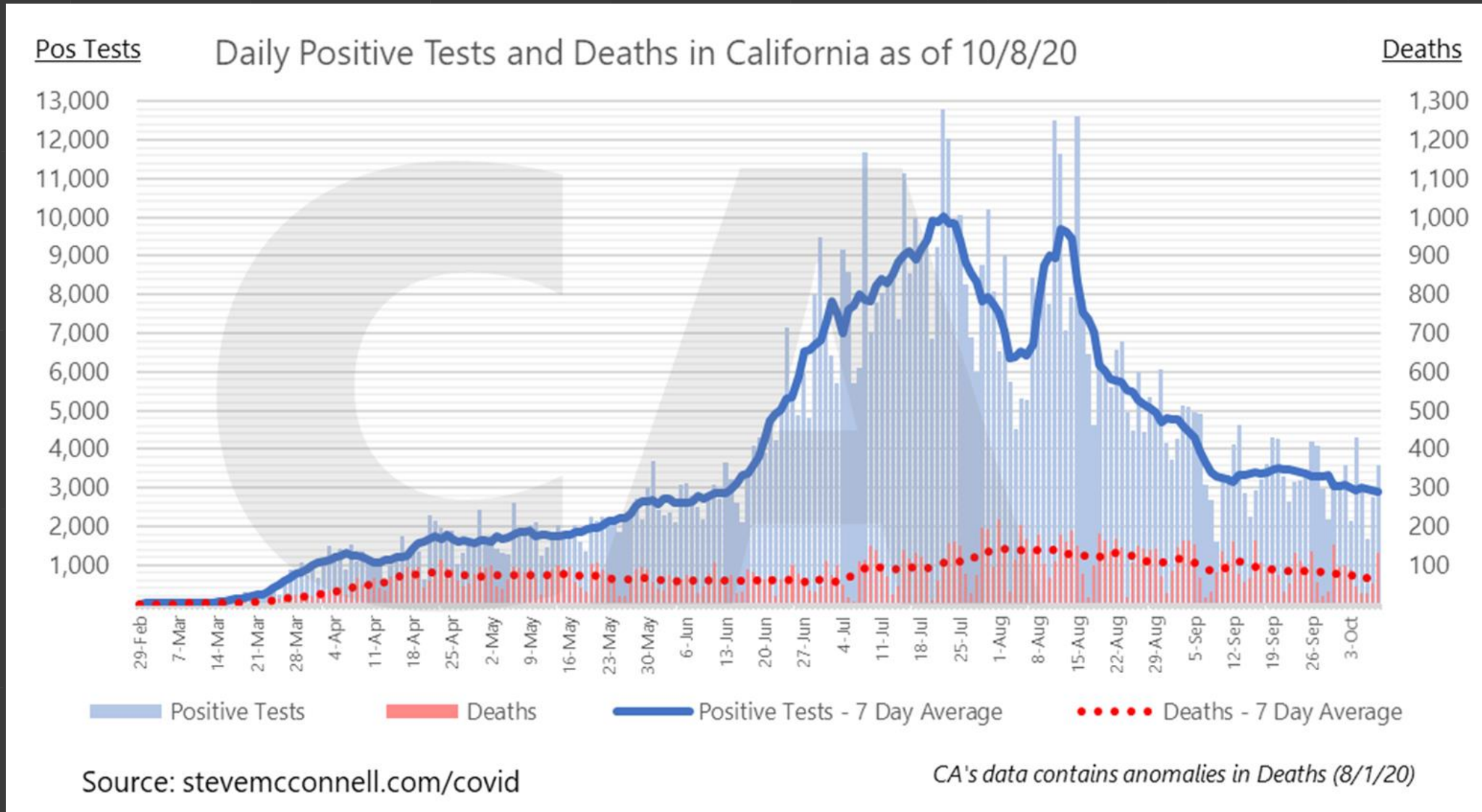
# Issue #2 with Daily Data – State Level Data



# Issue #2 with Daily Data – State Level Data



# Issue #2 with Daily Data – State Level Data



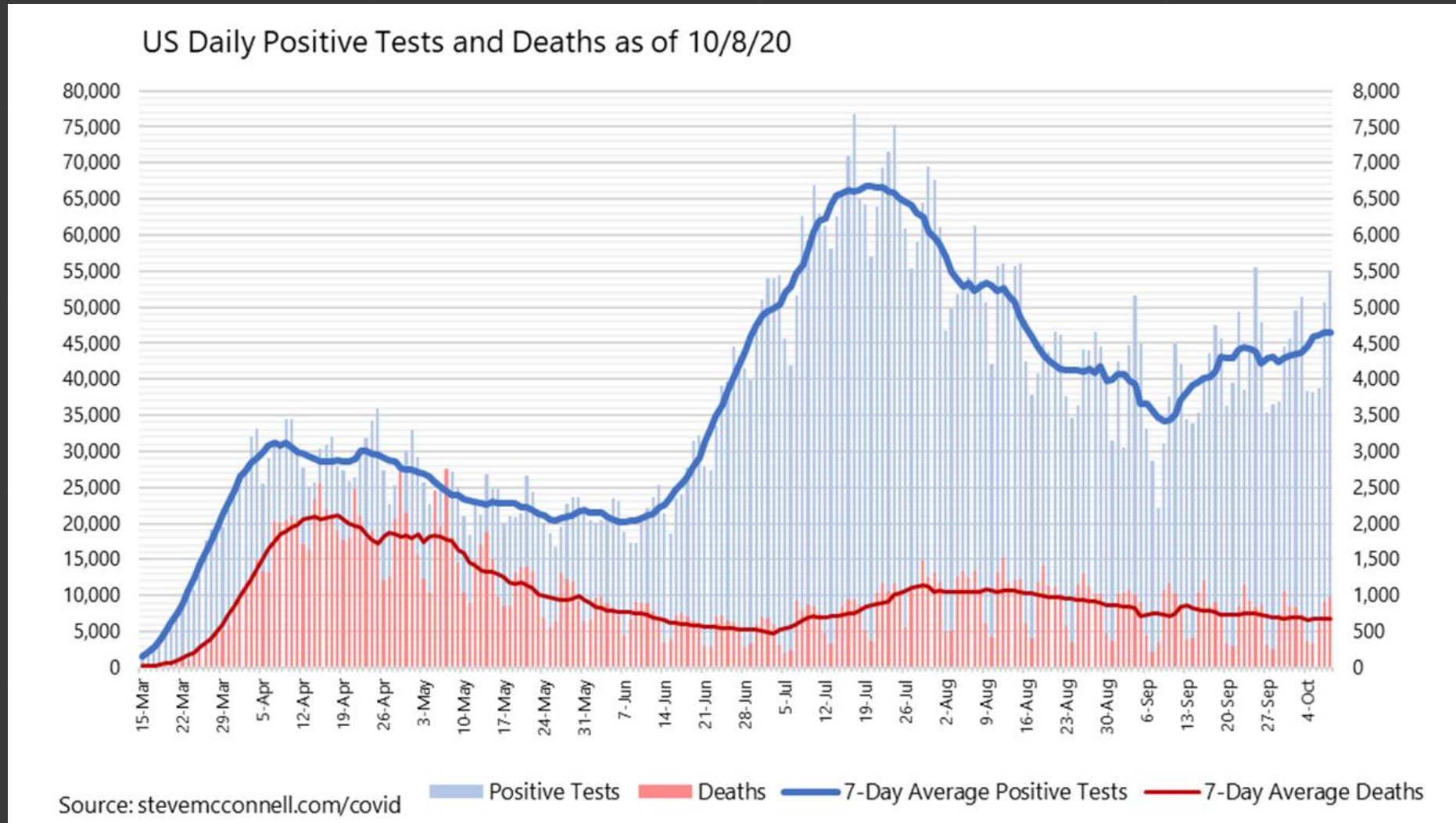
# — State Level Data Smoothing

- ❑ Due to non-regular reporting frequencies, smoothed data is actually more accurate than non-smoothed data
- ❑ The smoothing period must be a multiple of 7 days
- ❑ By the same reasoning, weekly data is more accurate than daily data
- ❑ Smoothing helps, but it doesn't change the fact that there are still weird spikes in the data

# — Covid-19 Data Issues

- ▣ Issues with daily data
  - Part 1: Sundays (the weekly cycle)
  - Part 2: State-level data spikes and corrections
- ▣ Alternatives
  - Running averages
  - Deltas

# — Data Smoothing – Running Averages



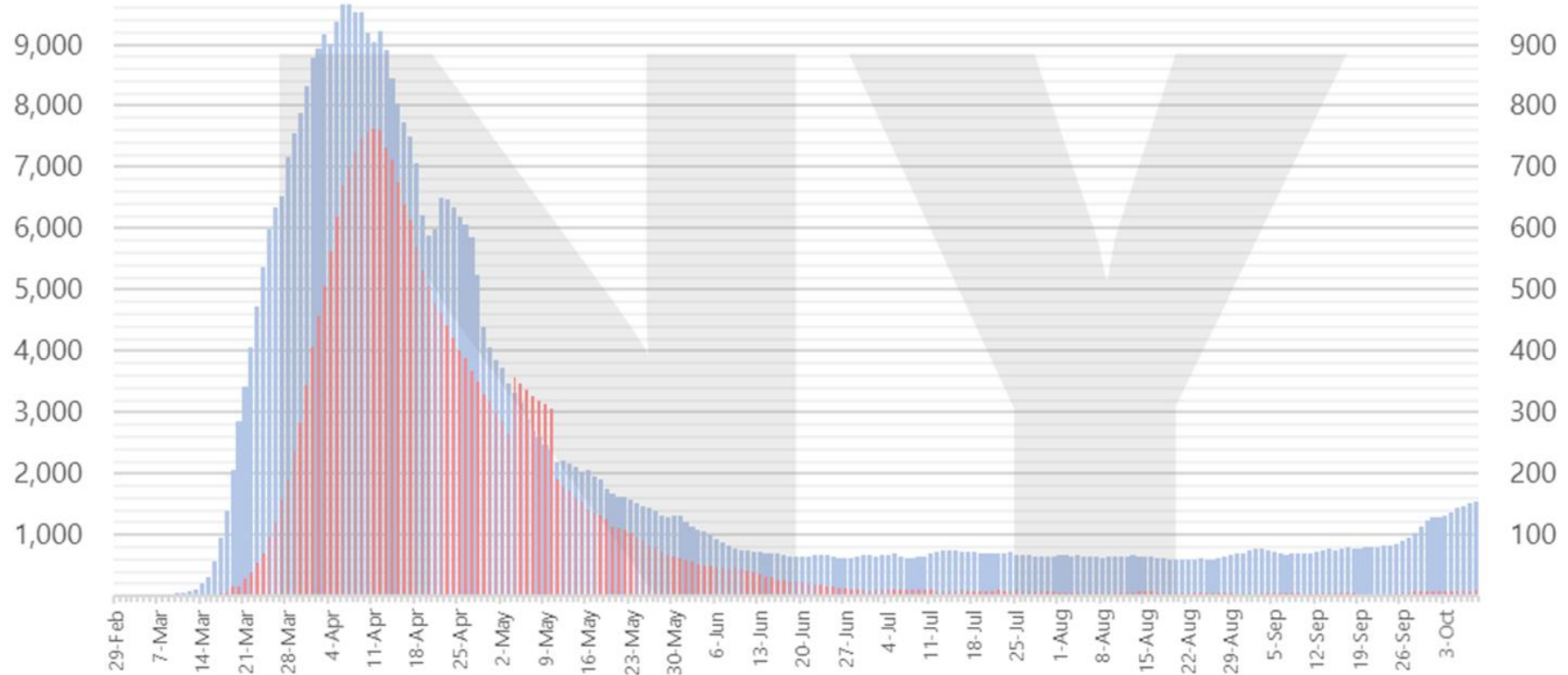


# — Data Smoothing – Running Averages

Pos Tests

Daily Positive Tests and Deaths in New York - Smoothed as of 10/8/20

Deaths



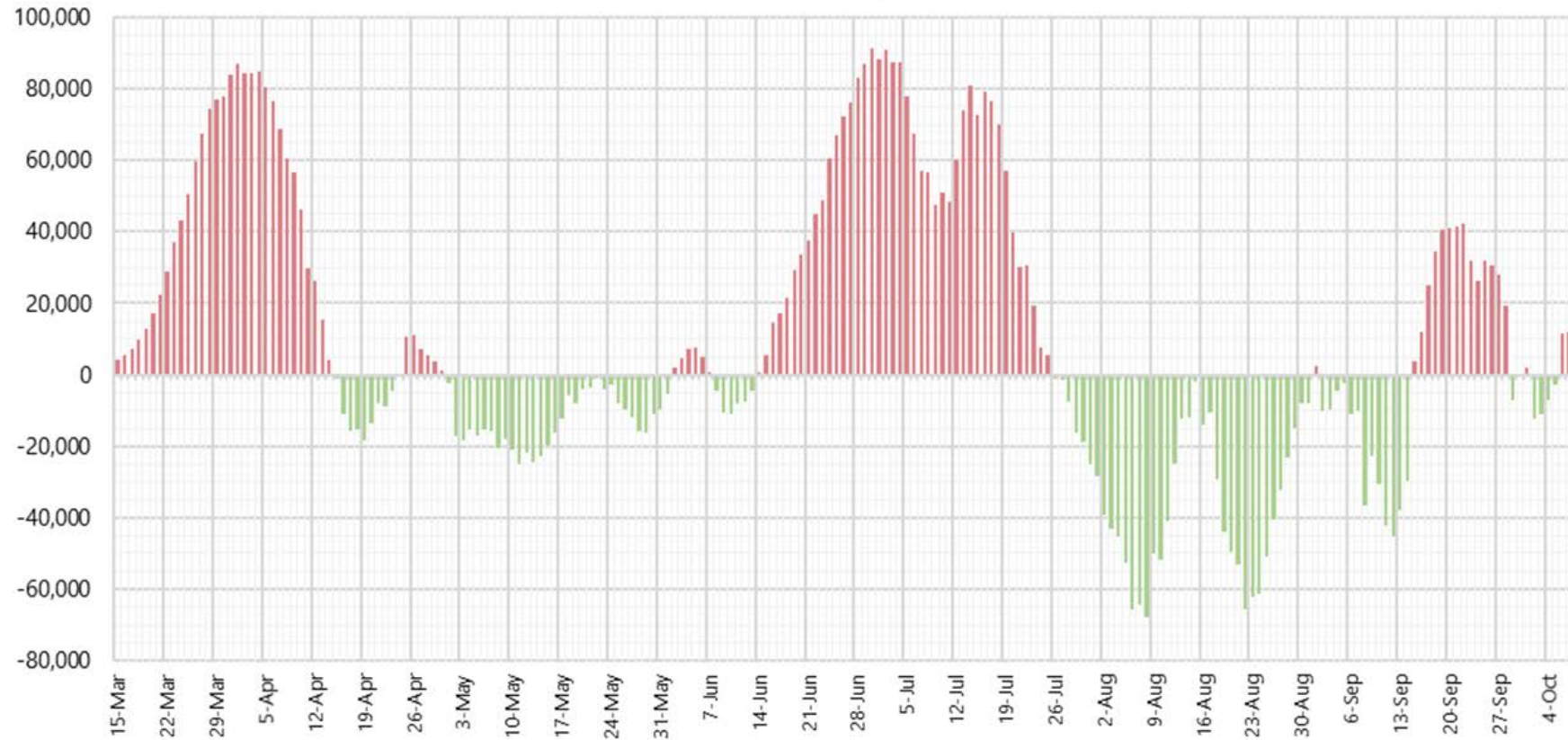
Source: [stevemccconnell.com/covid](http://stevemccconnell.com/covid)

■ Positive Tests

■ Deaths

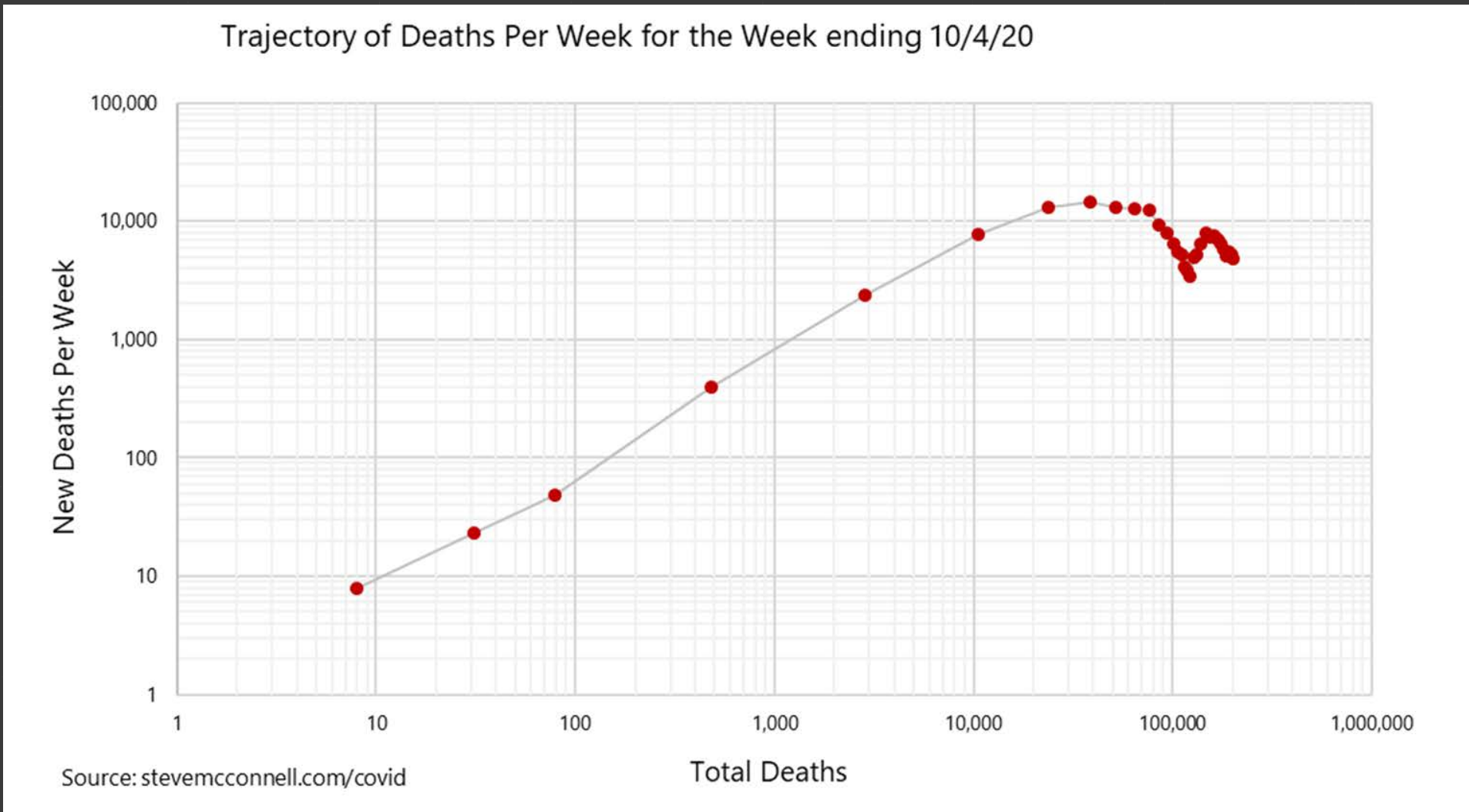
# — Data Smoothing – Deltas (linear)

Delta in Consecutive 7-Day Periods' Positive Tests for the US as of 10/8/20  
(i.e., Tests Previous 1-7 days vs. previous 8-14 days)

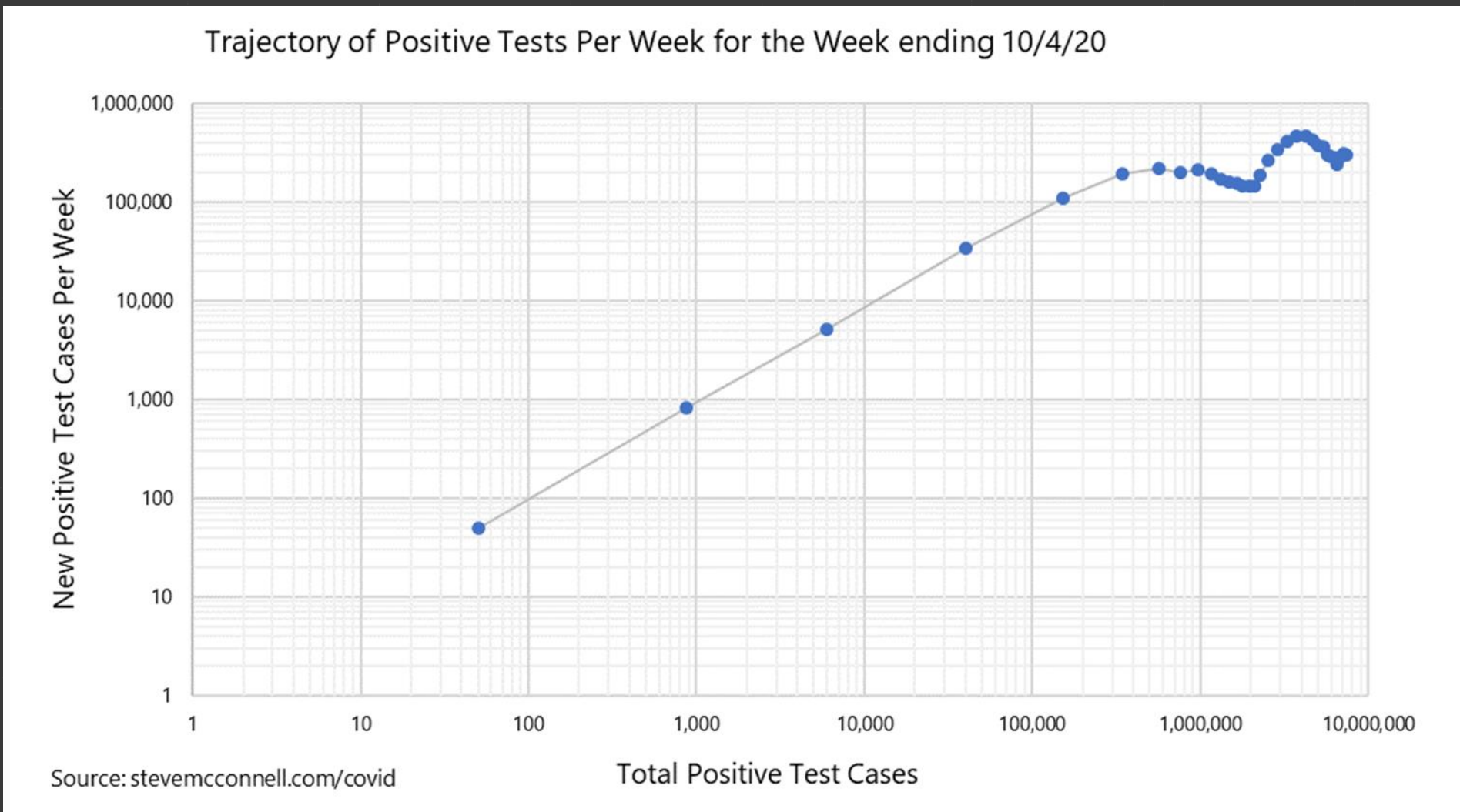


Source: [stevemccconnell.com/covid](http://stevemccconnell.com/covid)

# — Data Smoothing – Deltas (log)



# — Data Smoothing – Deltas (log)



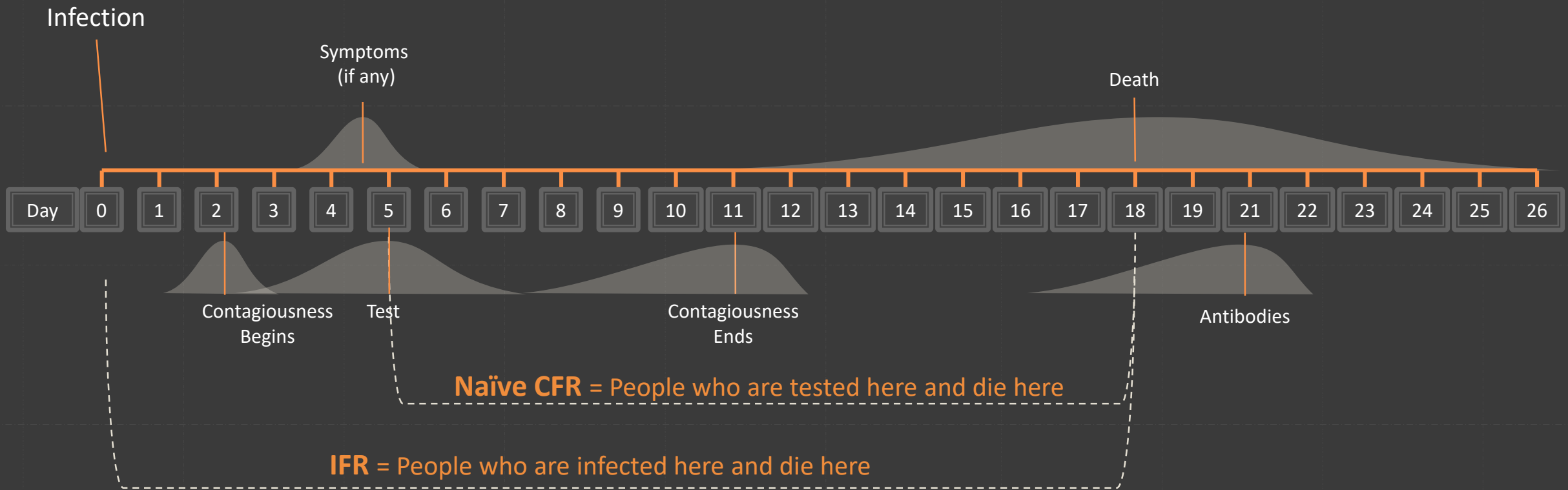
- Specific Data Relationships

# Calculating CFR/IFR

# — CFR vs. IFR - Terminology

- ❑ “Case” – people who are diagnosed/tested as positive
- ❑ “Infected” – people who get the virus whether they’re tested or not
- ❑ “CFR” = Case Fatality Rate
  - This is number of fatalities divided by number of “cases” (positive tests)
- ❑ “Junk CFR” = Today’s fatalities divided by today’s positive tests
- ❑ “Naïve CFR” = Today’s fatalities divided by the number of positive tests at some sensible point in the past (e.g., 14 days)
- ❑ “IFR” = Infection Fatality Rate – percentage of people infected who die, whether they’ve been tested or not
- ❑ BOTH Naïve CFR and IFR are useful for certain purposes

# — IFR vs. CFR





# Usefulness of IFR

---

Useful for assessing risk

- To an overall population
- To specific groups within the population

# Usefulness of CFR

---

Useful for estimation purposes

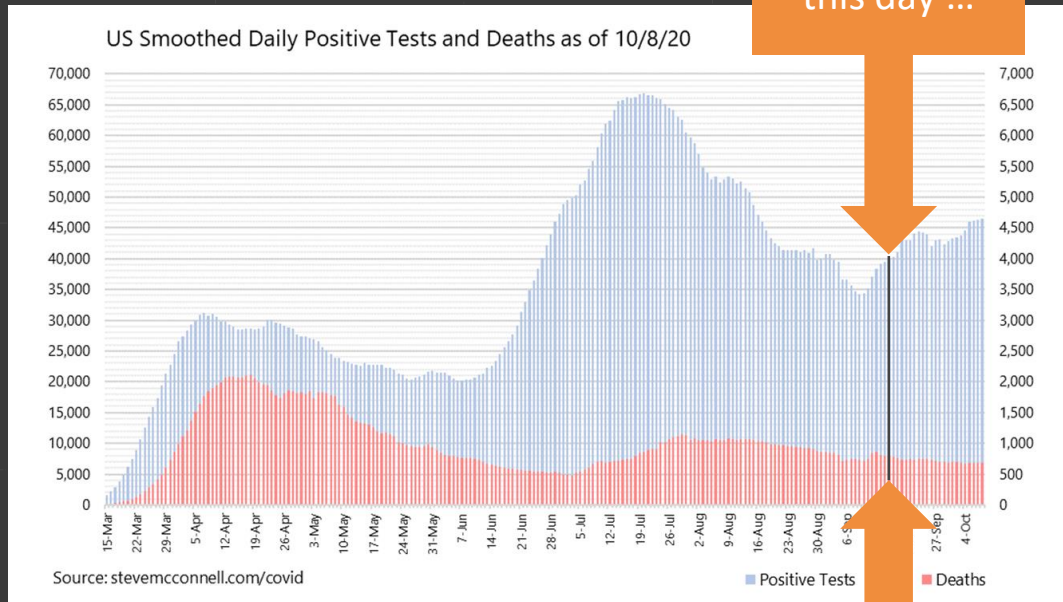
- IF you get all your facts straight (data assumptions)
- Calculations are usually presented incorrectly in ways that undermine clear understanding of the data

Data  
Correction #1

Lining up time periods correctly  
for Naïve CFR

## Correcting “Junk CFR” (line up time periods correctly)

- People do not get infected, test positive, and die on the same day
- But commonly reported CFRs effectively assume that’s the way it happens



The people infected on this day ...

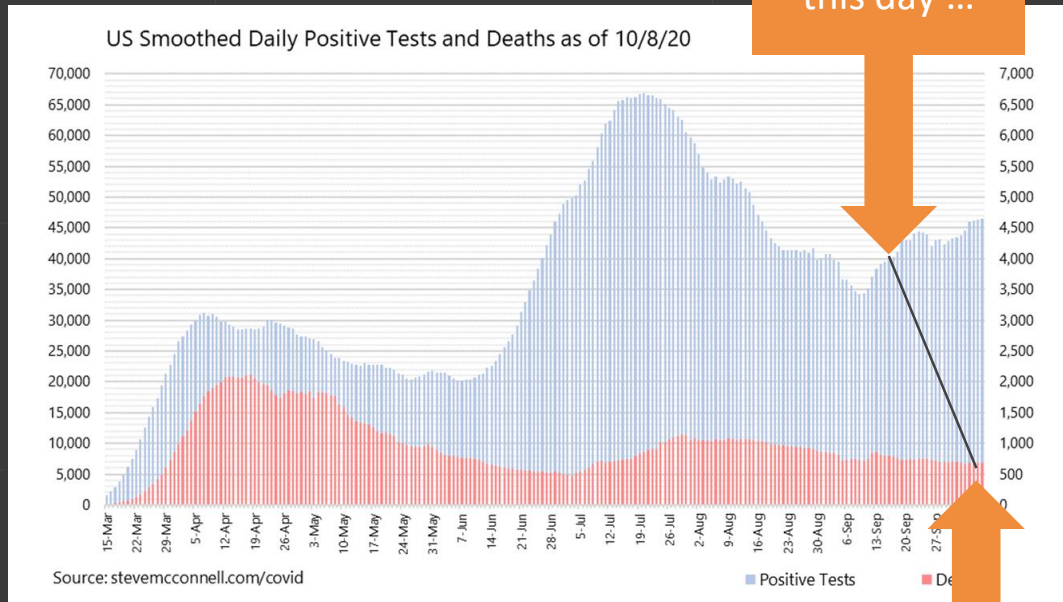
... do not die on the same day

## Correcting “Junk CFR” (line up time periods correctly)

- People do not get infected, test positive, and die on the same day
- But commonly reported CFRs effectively assume that’s the way it happens

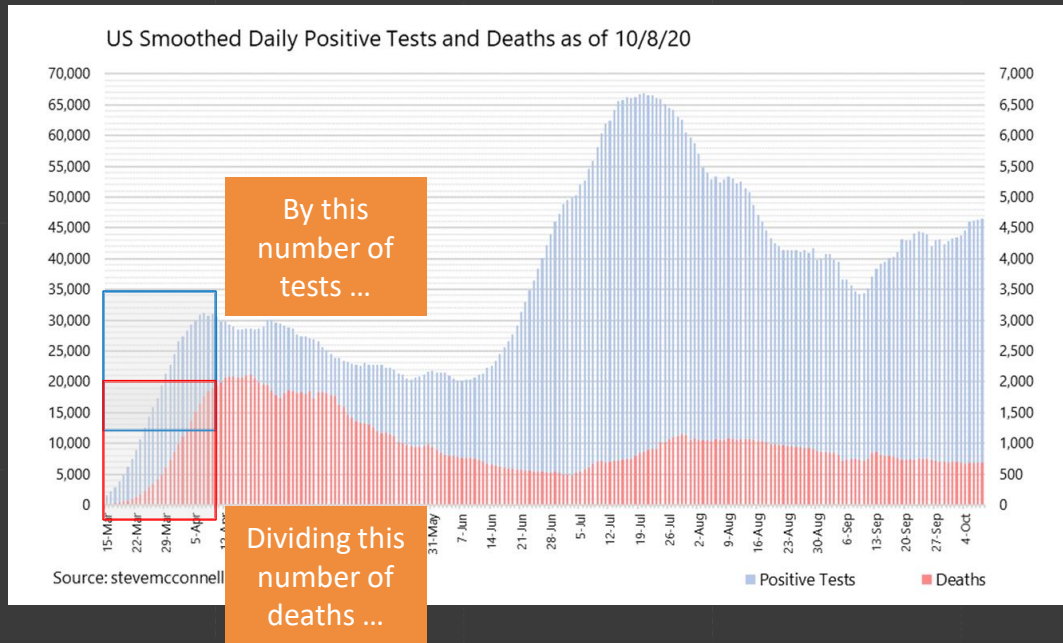
The people infected on this day ...

... die, on average, about 14 days later



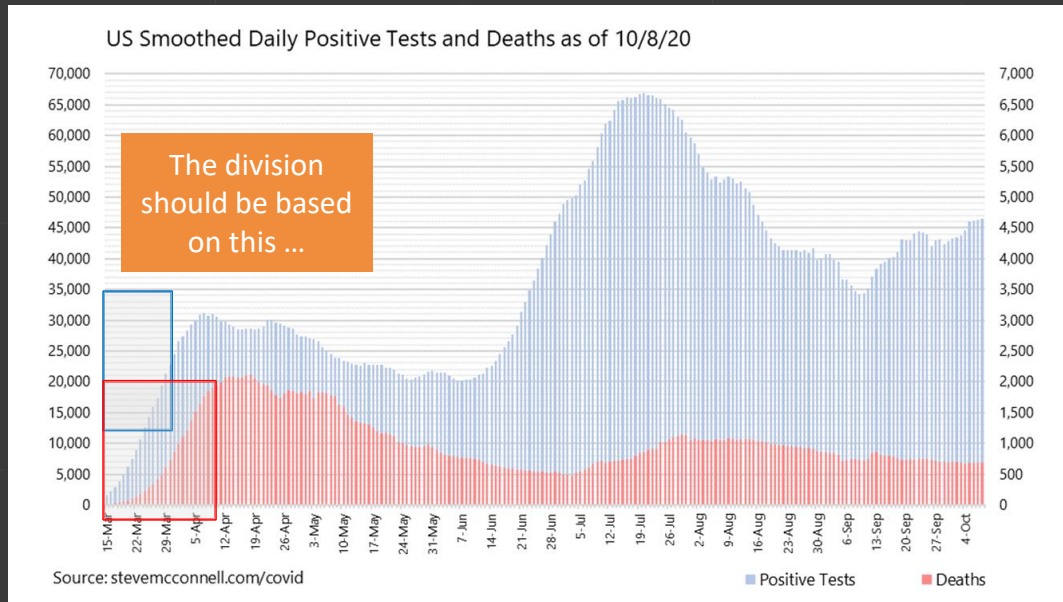
## Correcting “Junk CFR” (line up time periods correctly)

- This mistake gives rise to dramatically inaccurate calculations of fatality rates, especially in the early days of the pandemic



## Correcting “Junk CFR” (line up time periods correctly)

- This mistake gives rise to dramatically inaccurate calculations of fatality rates, especially in the early days of the pandemic
- This correction results in a **higher** CFR than was reported by the media at the time (but that doesn't imply the IFR was higher)



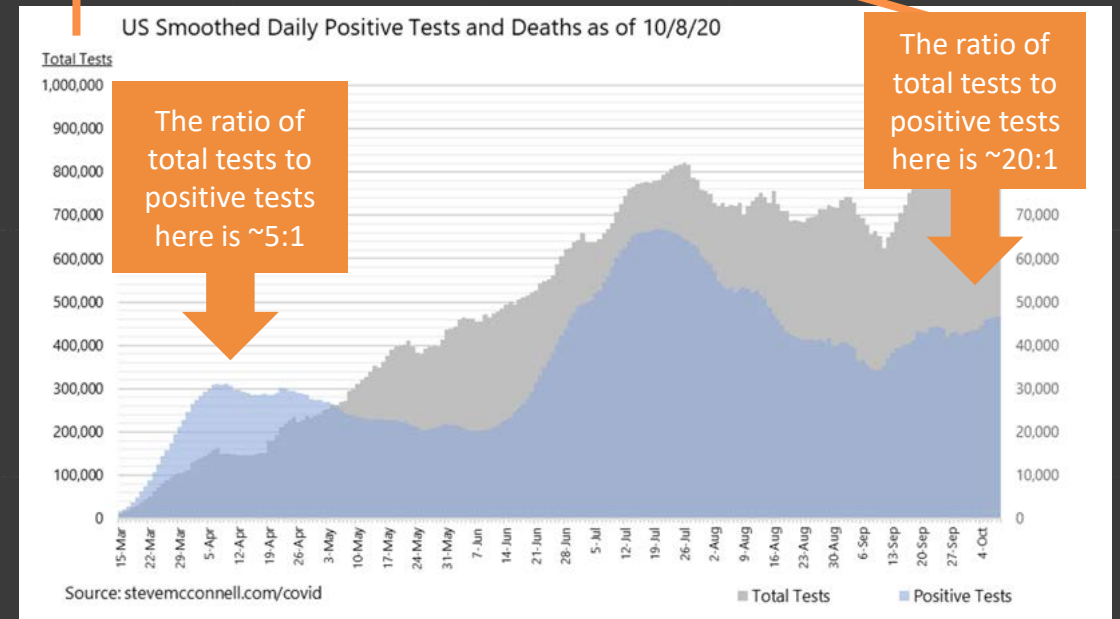
Data  
Correction #2

Tests <> Cases

# Relationship between positive tests and cases has changed

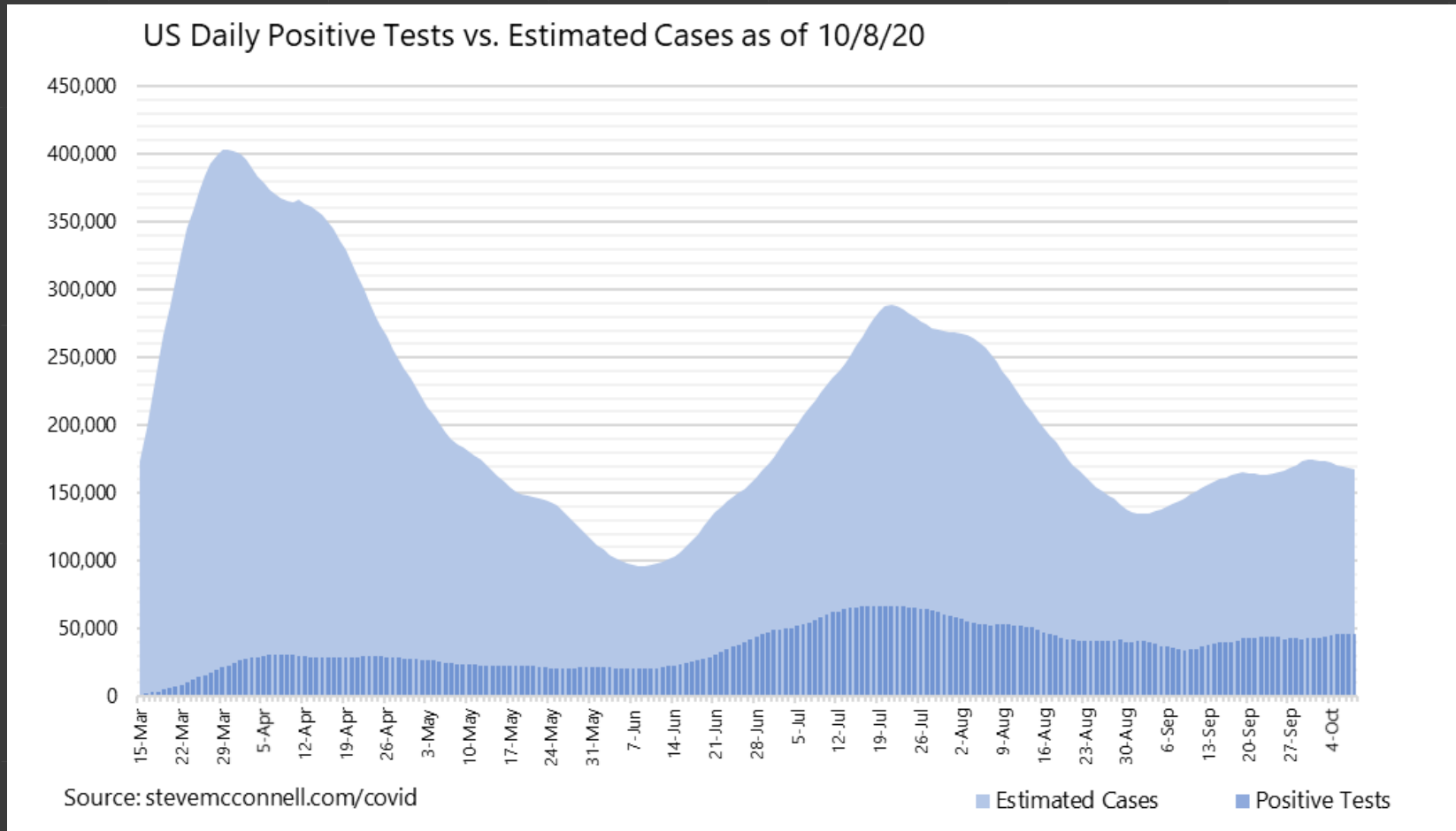
- The level of testing has steadily increased over the course of the pandemic
- Early in the pandemic, there were **10-20** cases per positive test
- More recently, there are **3-5** cases per positive test

Note: the total tests axis is 10x the positive tests axis

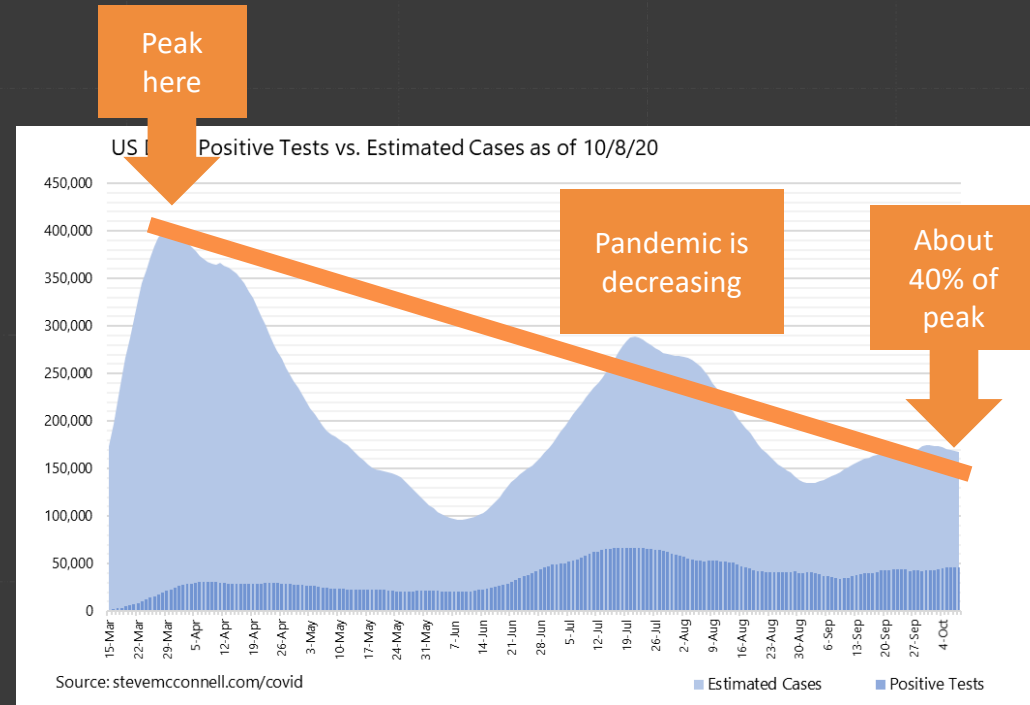
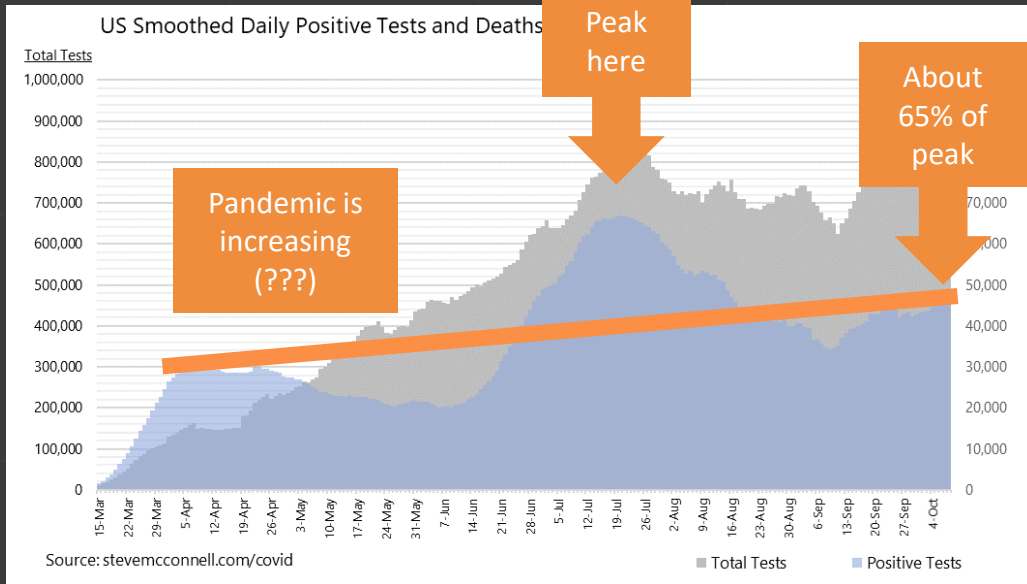




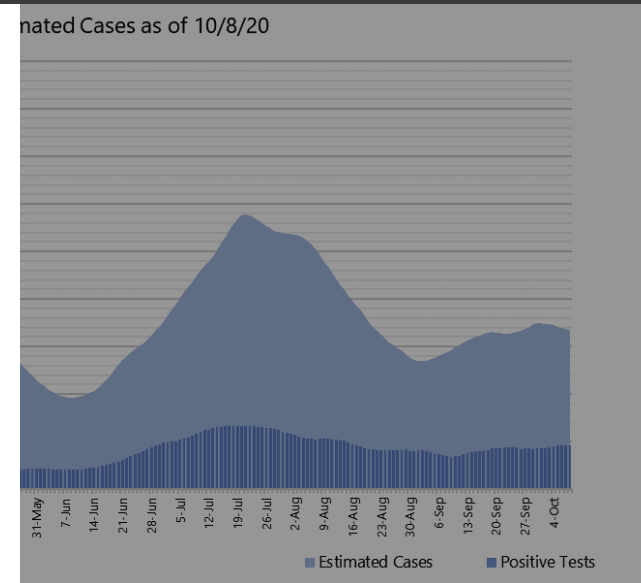
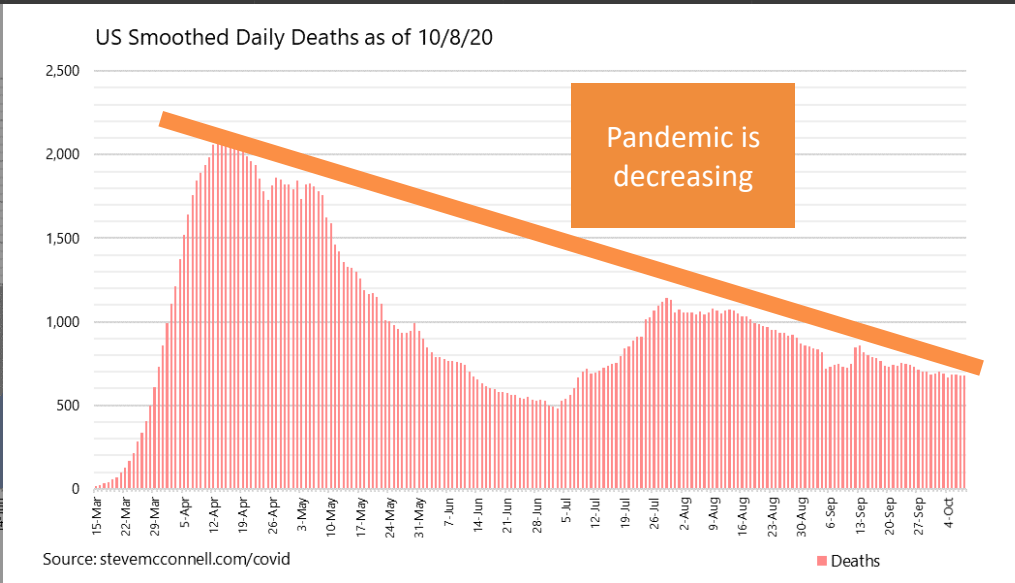
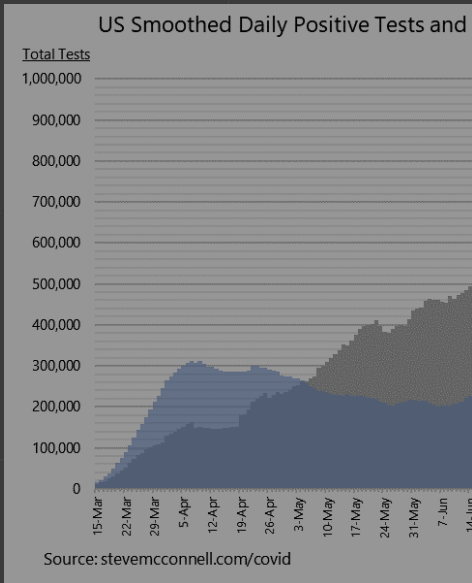
# Positive Tests vs. Estimated Cases Over Time



# — More on Positive Tests vs. Cases



# — The Estimated Cases Trend Aligns with Death Trend (The Positive Tests Trend Does not)



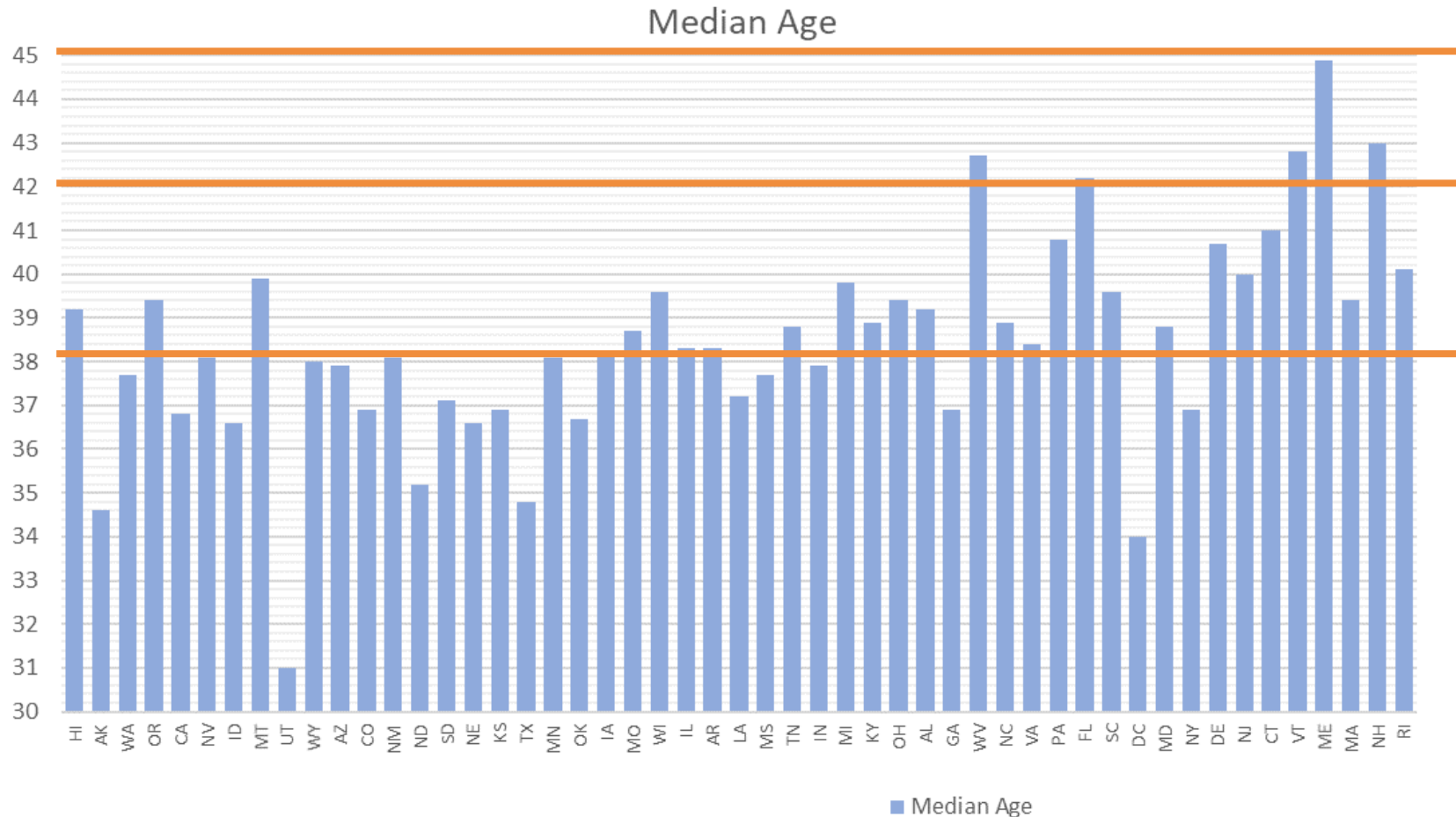
# — Calculating Meaningful CFR

- ❑ Get our time periods for tests and deaths lined up correctly
- ❑ Get clear about the idea that tests  $\leftrightarrow$  infections
- ❑ Calculate a correct number for infections

Data  
Correction #3

Account for Age

# — Remaining CFR Issue: CFR is Age-Based



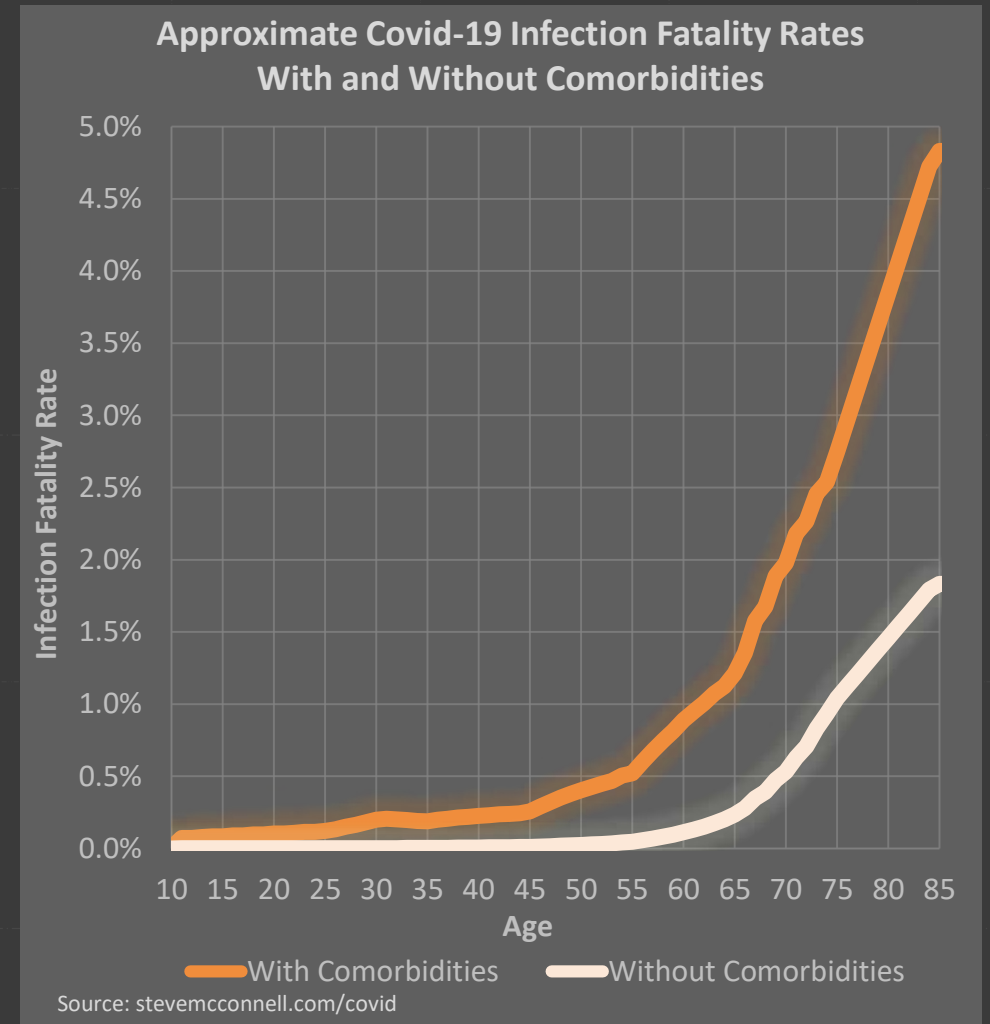
Italy's average is 45.1

Spain's average is 42.3

Average age in the US is 38.2

# IFRs and Co-Morbidity

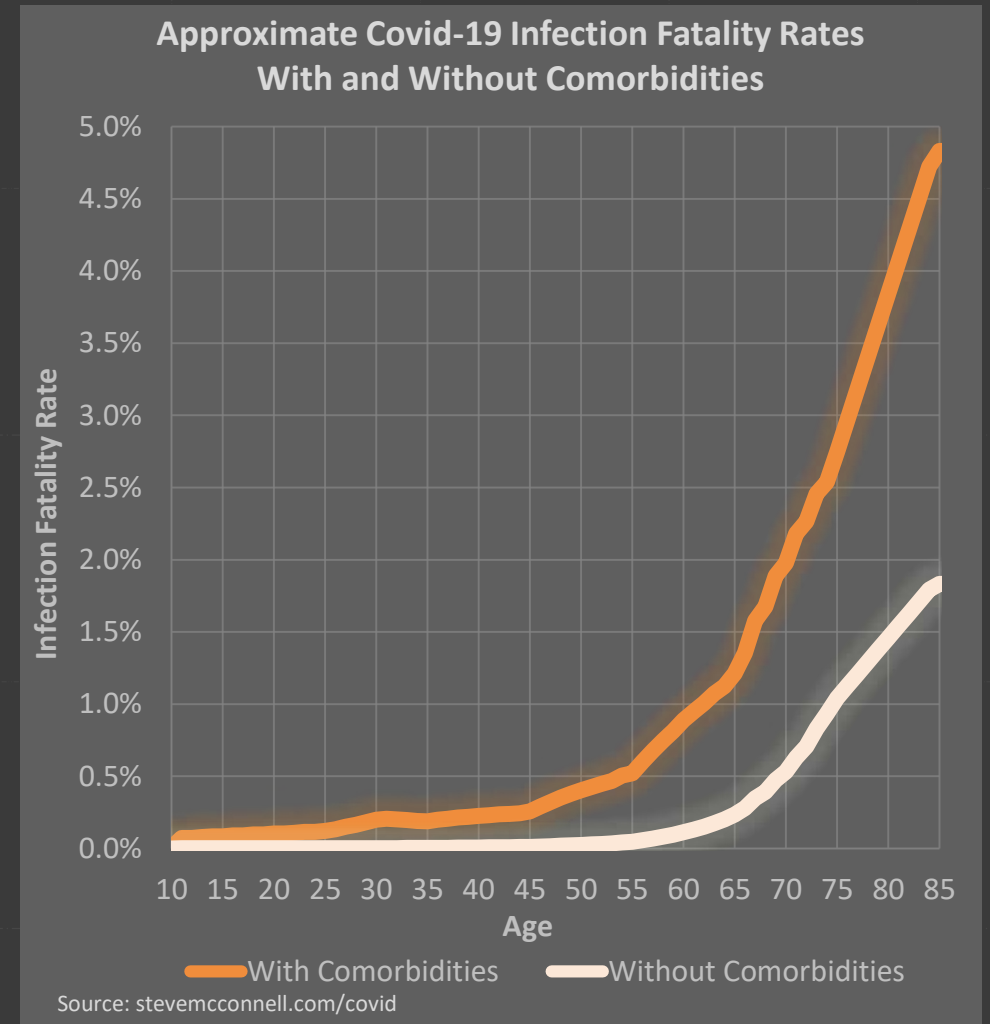
- 90% of deaths have involved at least 1 co-morbidity
- Risk of death with a co-morbidity is HIGHER than published; risk without is LOWER
- At younger ages this makes a HUGE difference in risk



# Most serious comorbidities, according to CDC

- Serious heart conditions, such as heart failure, coronary artery disease, or cardiomyopathies
- Cancer
- Chronic kidney disease
- COPD
- Obesity (BMI > 30)
- Sickle cell disease
- Solid organ transplantation
- Type 2 diabetes mellitus

CDC list of risk levels of specific comorbidities and level of evidence for each:  
<https://www.cdc.gov/coronavirus/2019-ncov/need-extra-precautions/evidence-table.html>





# Summary of “Tools” we can use for Forecasting

- Corrections to raw data reporting (weekly cycle issues)
- Corrections to timeline (lag from positive tests to deaths)
- Understanding of what is needed to estimate infections (not tests)
- Potential to adjust for the key factor of age

— CDC

Covid-19 Forecasting

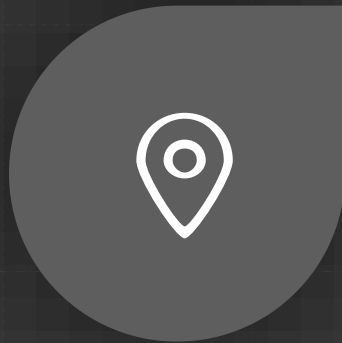
[CDC Forecasting Website](#)

# — CDC Covid-19 Forecast Process

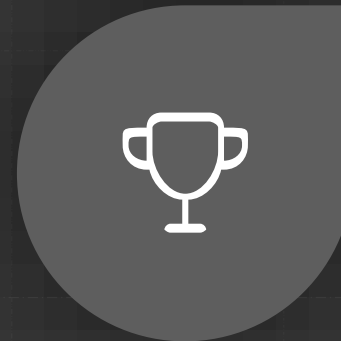
CDC work is overseen by a team at University of Massachusetts Amherst



Modeling groups submit forecasts to the CDC



Specific organization is Reich Lab based in the Department of Biostatistics and Epidemiology



Forecasts that meet certain criteria are combined into the “Ensemble” model, which is the forecast model of record for the CDC

# — Modeling submissions to the CDC

Forecast are submitted weekly, on Mondays

Forecasts are submitted via github

All forecast data is public and can be accessed by anyone

Forecasts are submitted in csv files with specific formatting requirements

Submitting a forecast to the Ensemble model requires generation of about 10,000 forecast records / week / team to participate

# CDC Covid-19 Forecasts

## Types

“Point” forecast

“Prediction interval”  
(commonly known  
as 95% confidence  
ranges)

## Quantities

“Cases”  
(positive tests)

Deaths

## Regions

51+ states and  
territories

US National

## Time Period

Week 1

Week 2

Week 3

Week 4

# — Forecast Teams

## Universities

Carnegie Mellon  
Columbia  
Columbia/UNC  
Georgia Tech  
Harvard  
Iowa State  
Johns Hopkins  
London School of  
Hygiene and Tropical  
Medicine  
MIT  
Northeastern  
Notre Dame

## Universities

RPI  
Texas Tech  
UCLA  
USC  
UW/IHME  
University of AZ  
U Cal Merced  
UCSD  
University of Geneva  
University of Georgia  
U Mass Amherst  
U Michigan  
U Texas Austin

## Research Labs

Los Alamos National Lab  
US Army Research and  
Development Center  
Walmart Labs Data  
Science Team  
Covid19 Simulator  
Consortium (MassGen,  
Harvard, Georgia Tech,  
Boston Medical Center)

## Individuals and Firms

Discrete Dynamical  
Systems  
Institute for Business  
Forecasting  
IQVIA  
John Burant  
Karlen Working Group  
LockNQuay  
Oliver Wyman  
Predictive Science Inc  
Qi-Jun Hong  
Robert Walraven  
Steve McConnell  
Youyang Gu

# Forecast Teams and Coordination

The better forecast teams and the worst are not what you'd expect them to be, considering some of the institutions involved

There's a weekly forecast call Tuesdays at noon Pacific time to review the week's forecasts

The university teams can be quite academic ("once you know how the sausage is made, you'll never eat sausage again")

Some well-known, highly respected institutions are submitting **terrible** forecasts

# Groups use Various Methods

- Massive data sets, e.g., 500 million records
- Esoteric data, e.g., use of mobility data from cell phone records
- Machine learning
- AI
- Bayesian analysis, Monte Carlo simulations, etc.
- Pre-existing infectious diseases models

Model descriptions are available from the [CDC website](#) (many just link to github)





One of the most enduring and useful conclusions from research on forecasting is that simple methods are generally as accurate as complex methods.

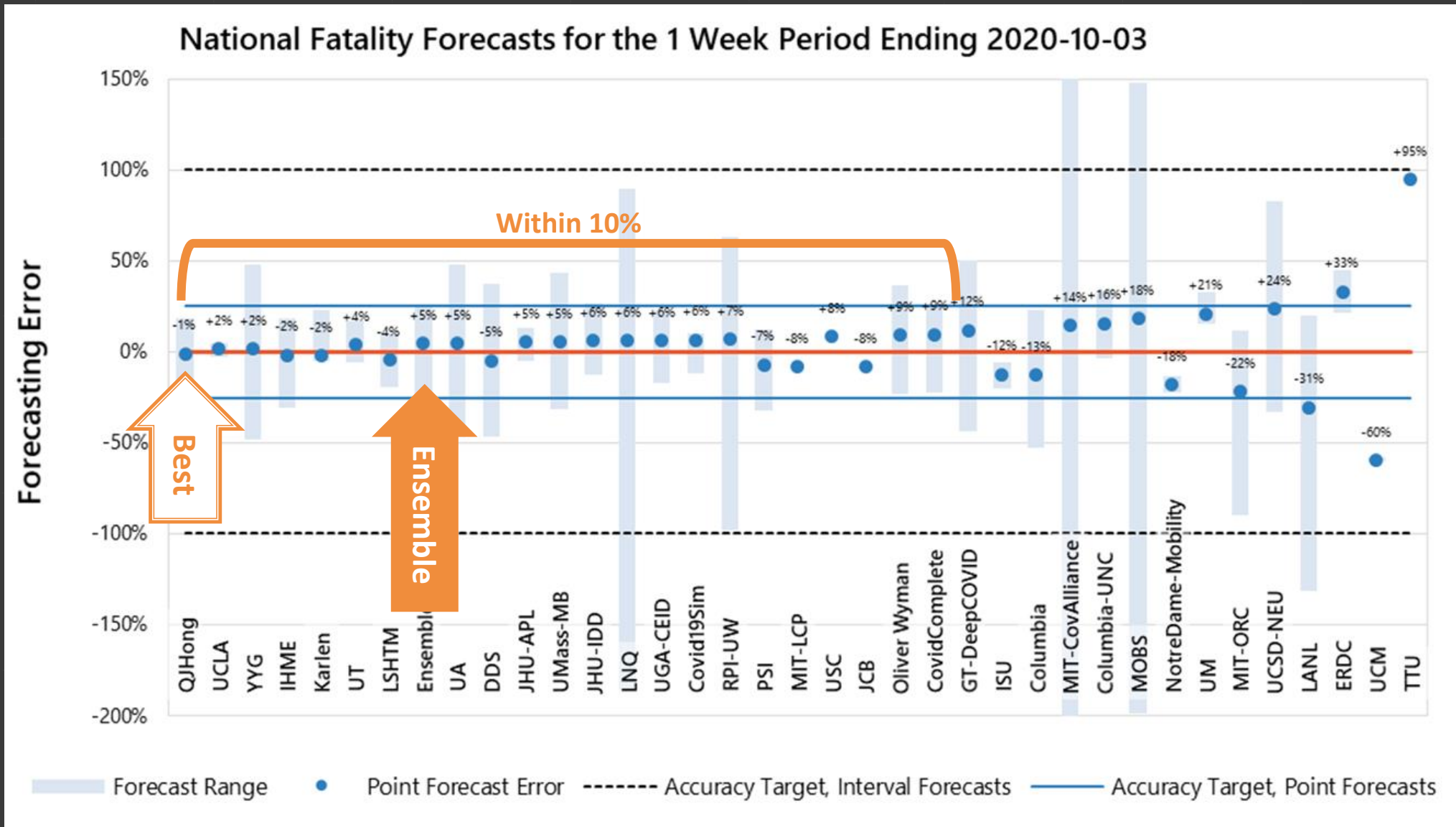
— *J. Scott Armstrong,*  
*Principles of Forecasting*



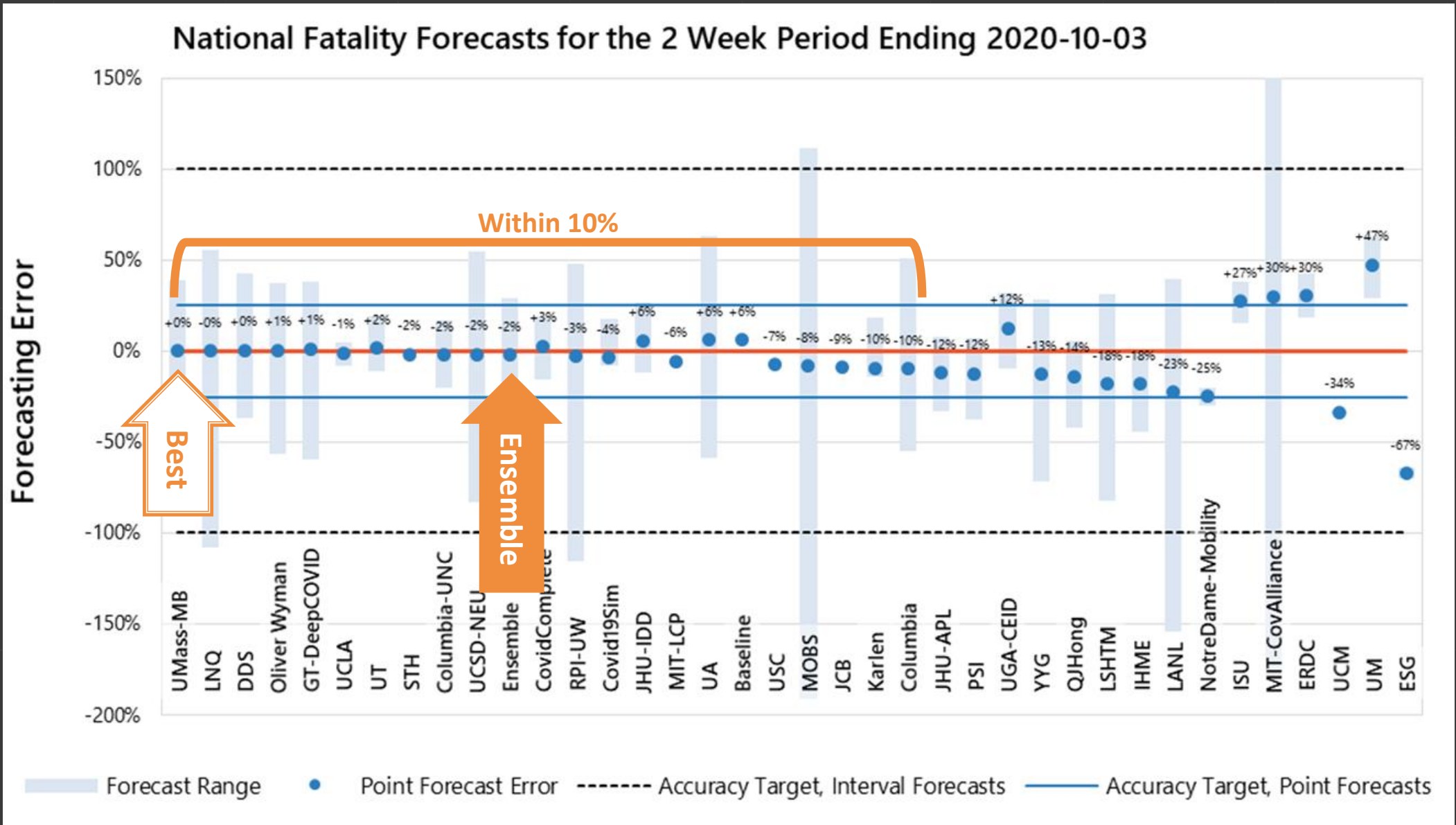


# What's Possible with Forecasting Nationally?

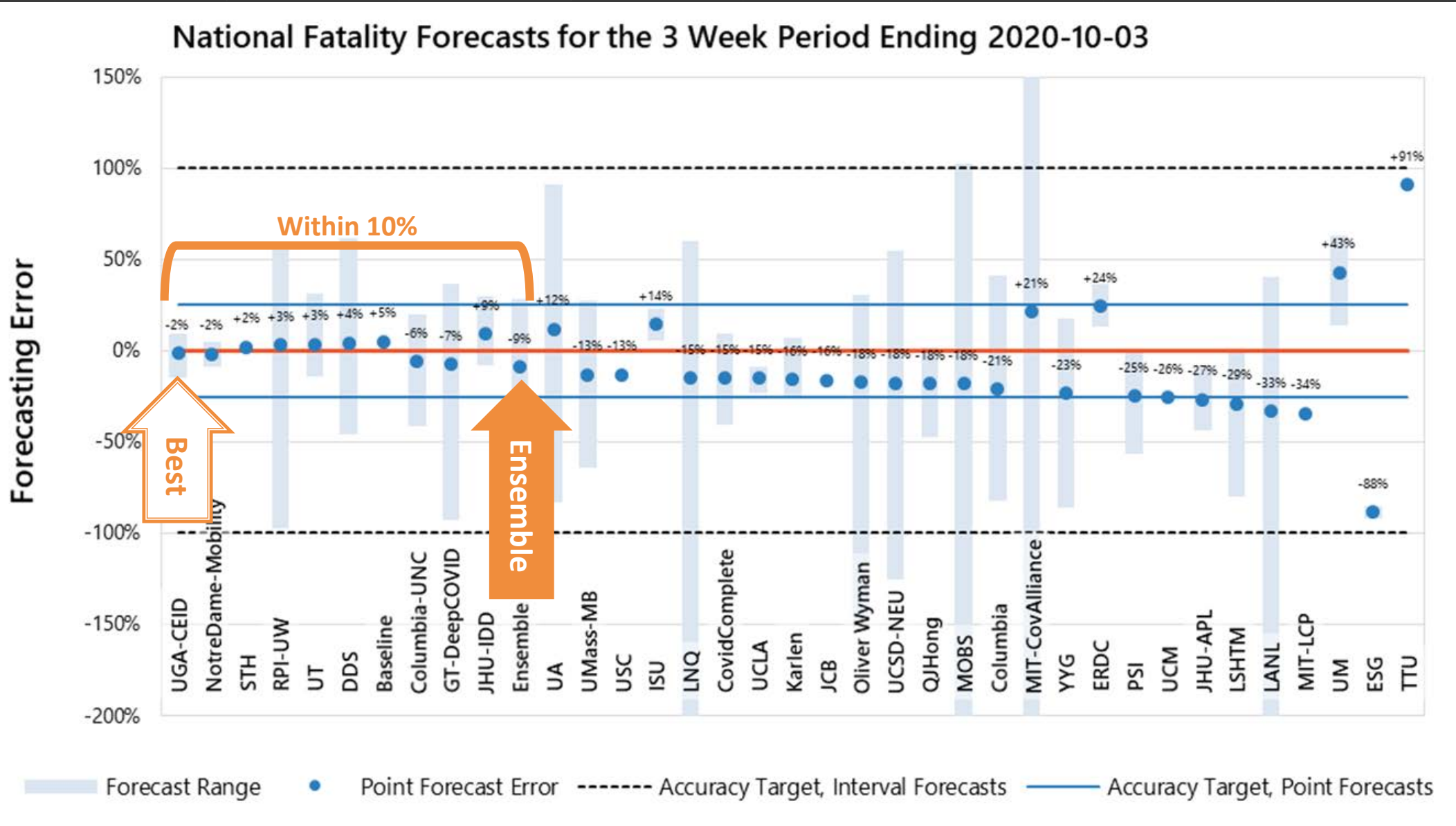
# Recent National Forecasts



# Recent National Forecasts

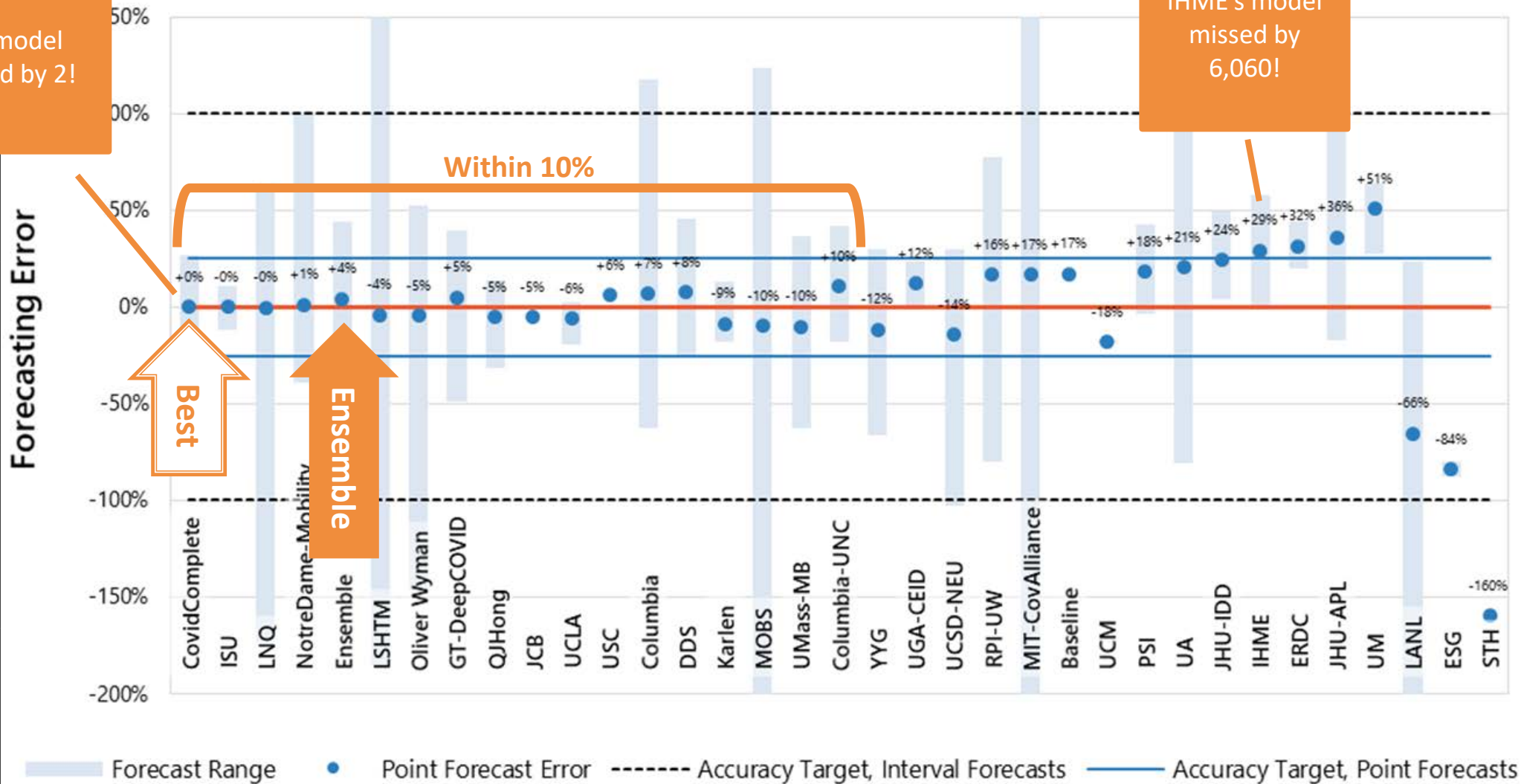


# Recent National Forecasts

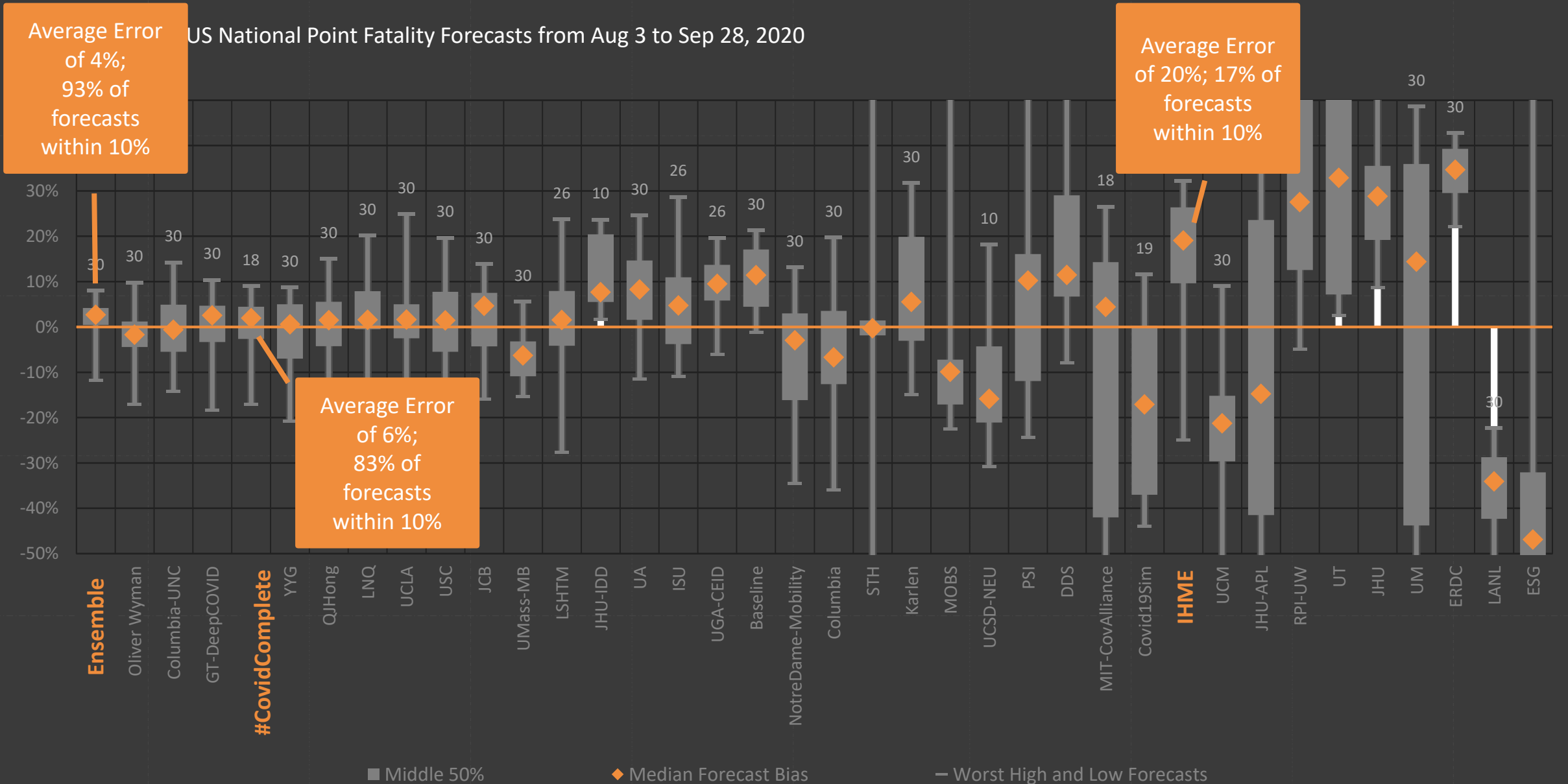


# Recent National Forecasts

National Fatality Forecasts for the 4 Week Period Ending 2020-



# Accuracy over time for 1, 2, 3, and 4-week forecasts

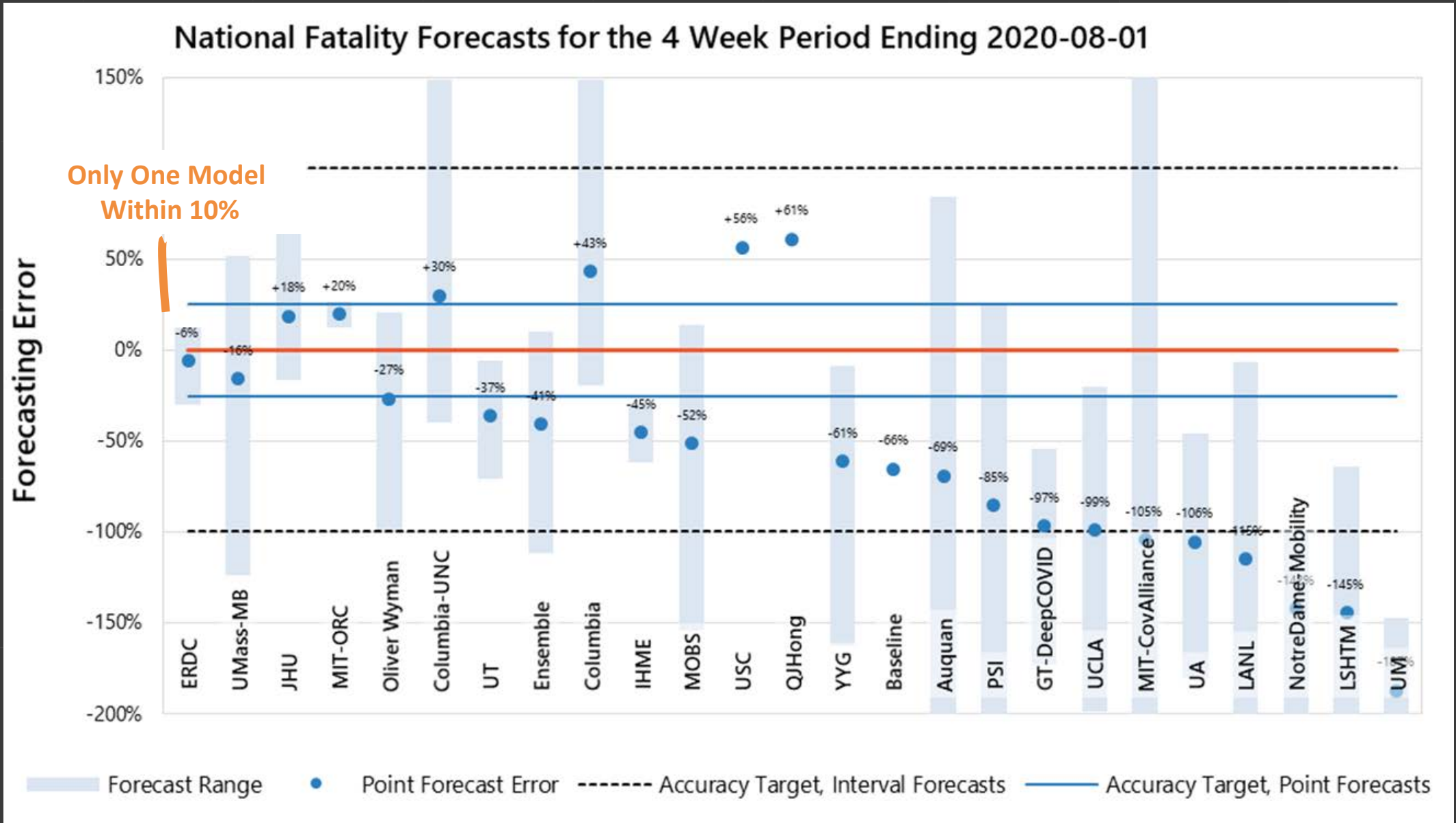


Covid-19

National Forecasts Have  
Improved Markedly  
Since July

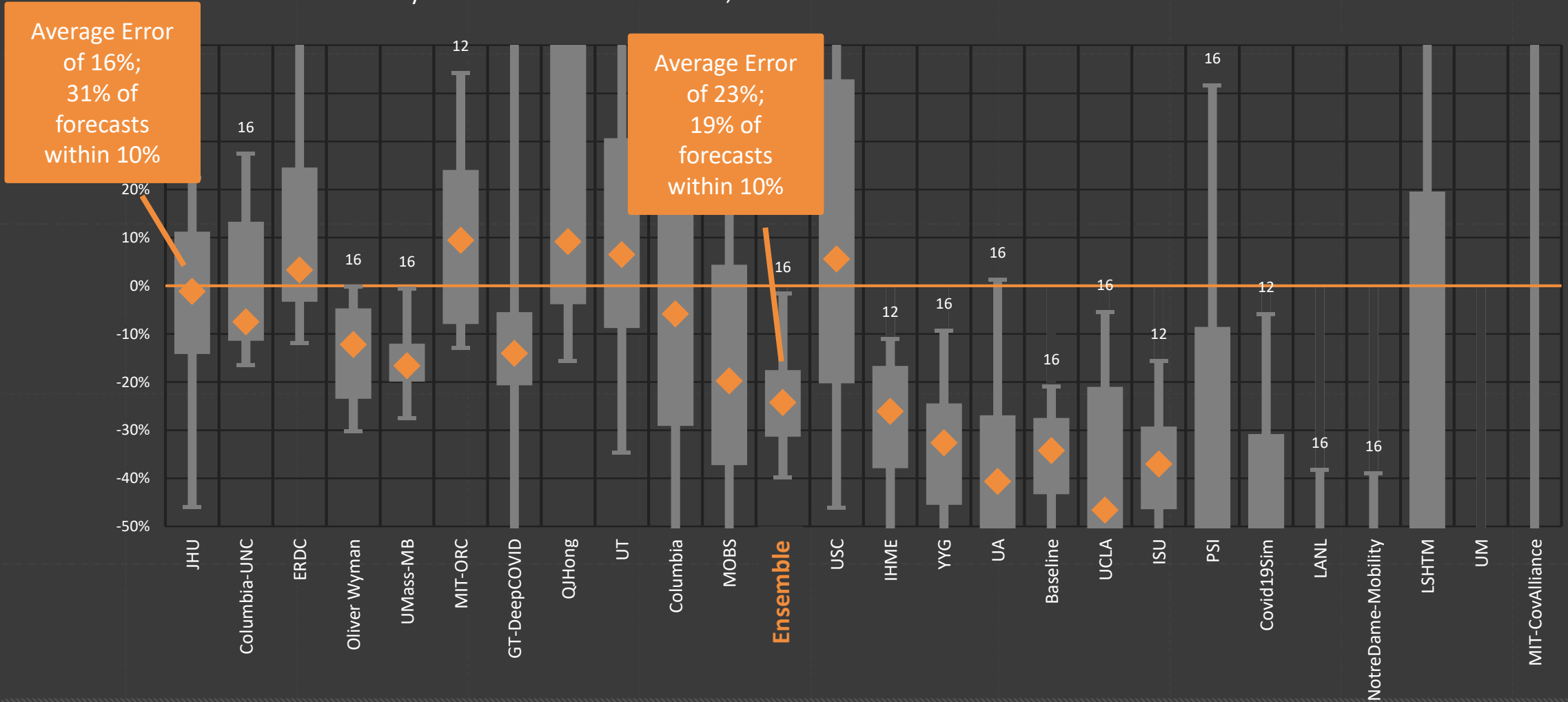


# Earlier National Forecasts

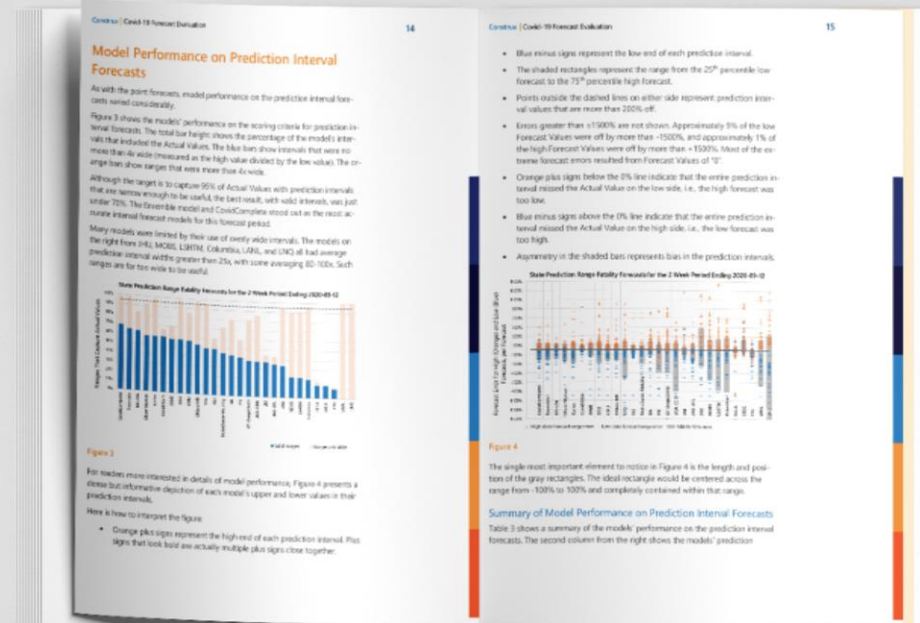
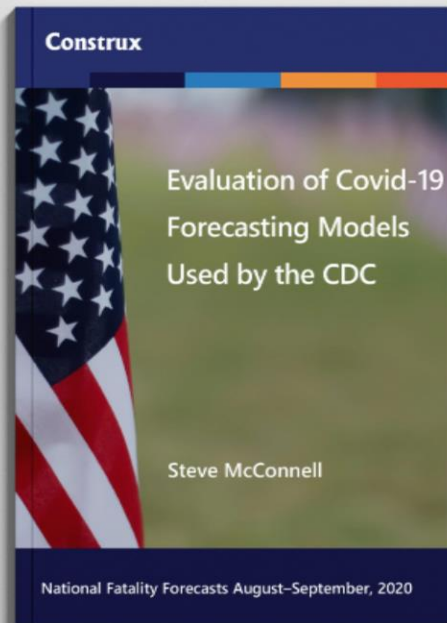


# Longer-term history of earlier national forecasts

US National Point Fatality Forecasts from Jul 6 to Jul 27, 2020



# — More National Forecast Evaluations



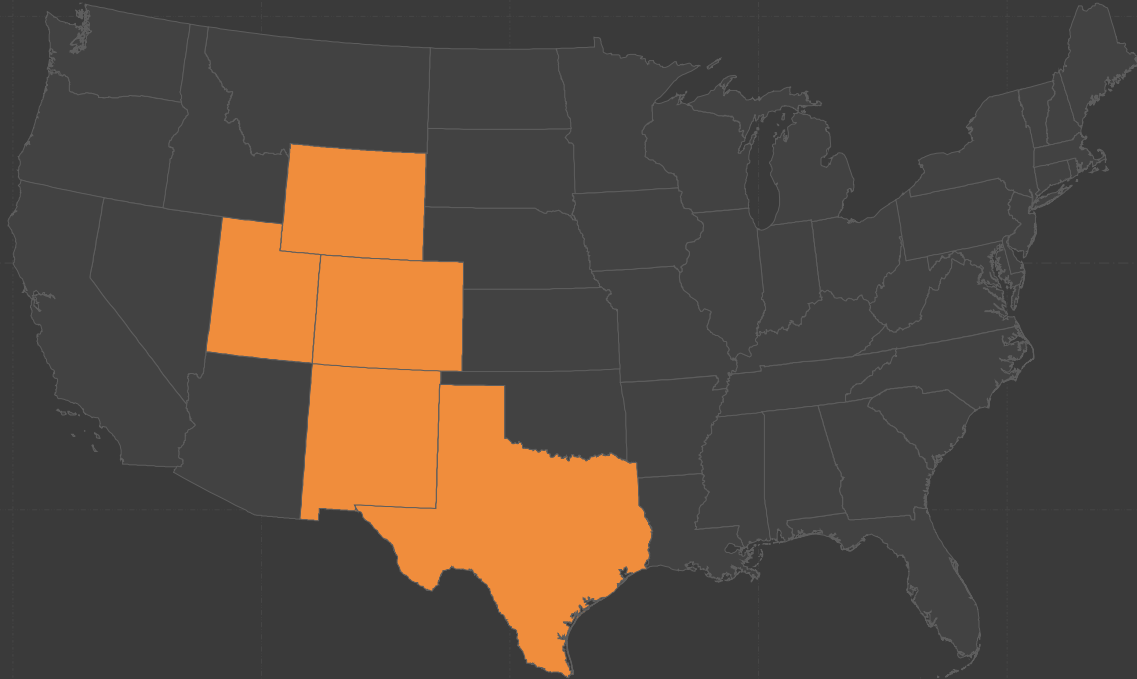
[www.stevemccconnell.com/covid](http://www.stevemccconnell.com/covid)



# What's Possible with Forecasting at the State Level?

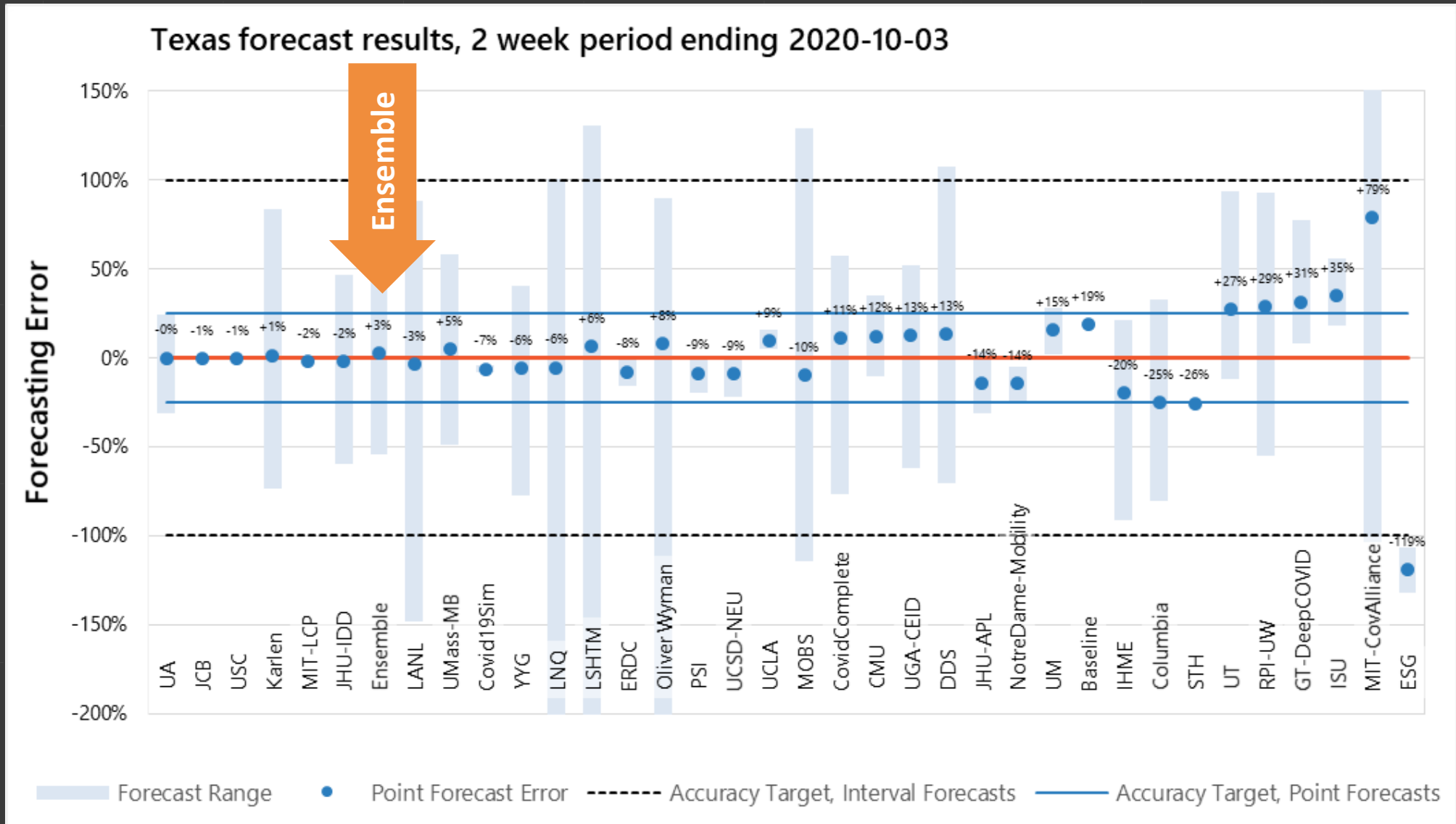
---

# State Level Forecasts

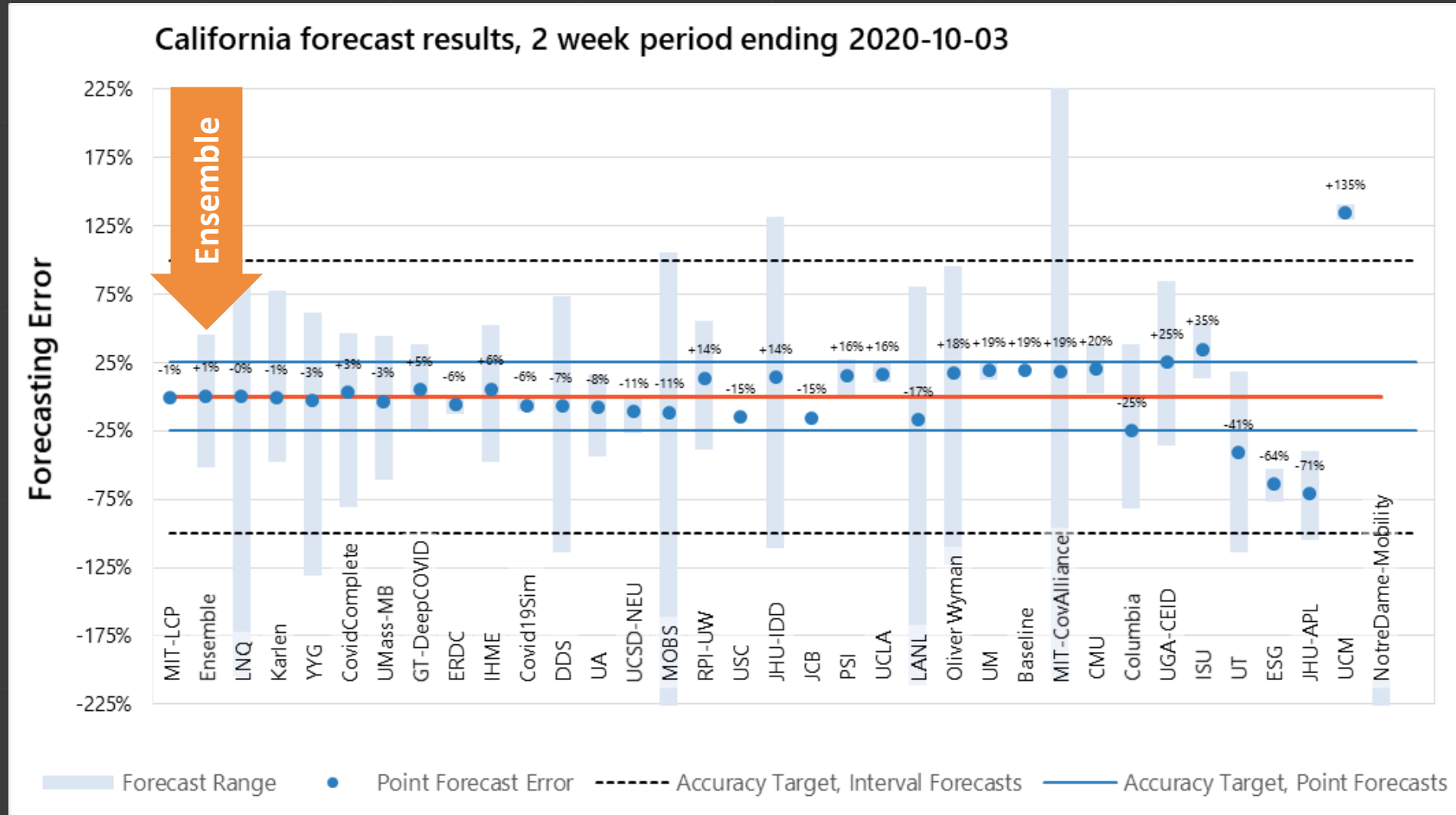


- Numbers of total tests, positive tests, cases, and deaths are smaller
- Data is rougher
- Models tend to do better on states with bigger numbers
- Accurate forecasting is arguably more important than national forecasts

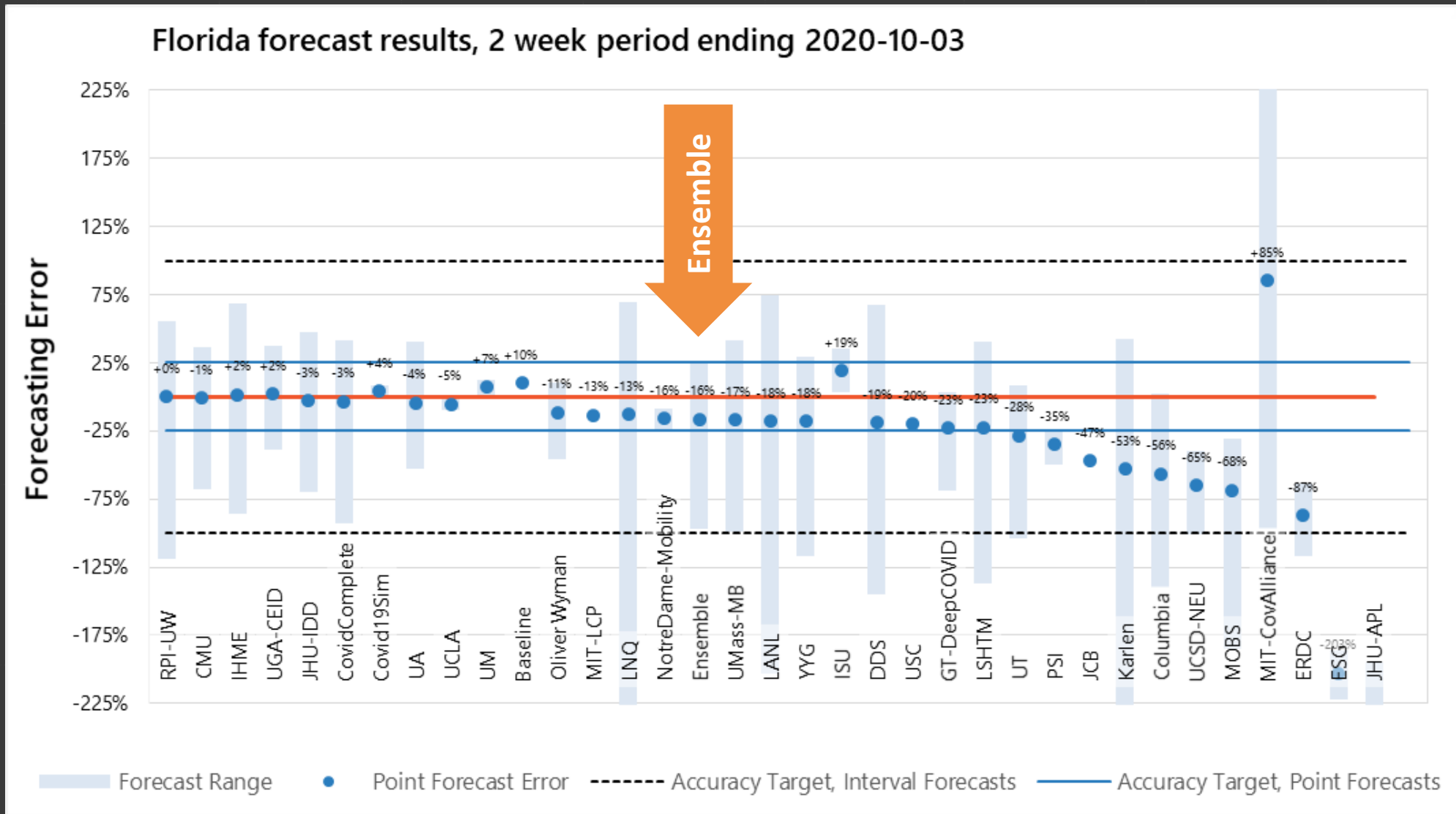
# State Forecast Examples



# State Forecast Examples

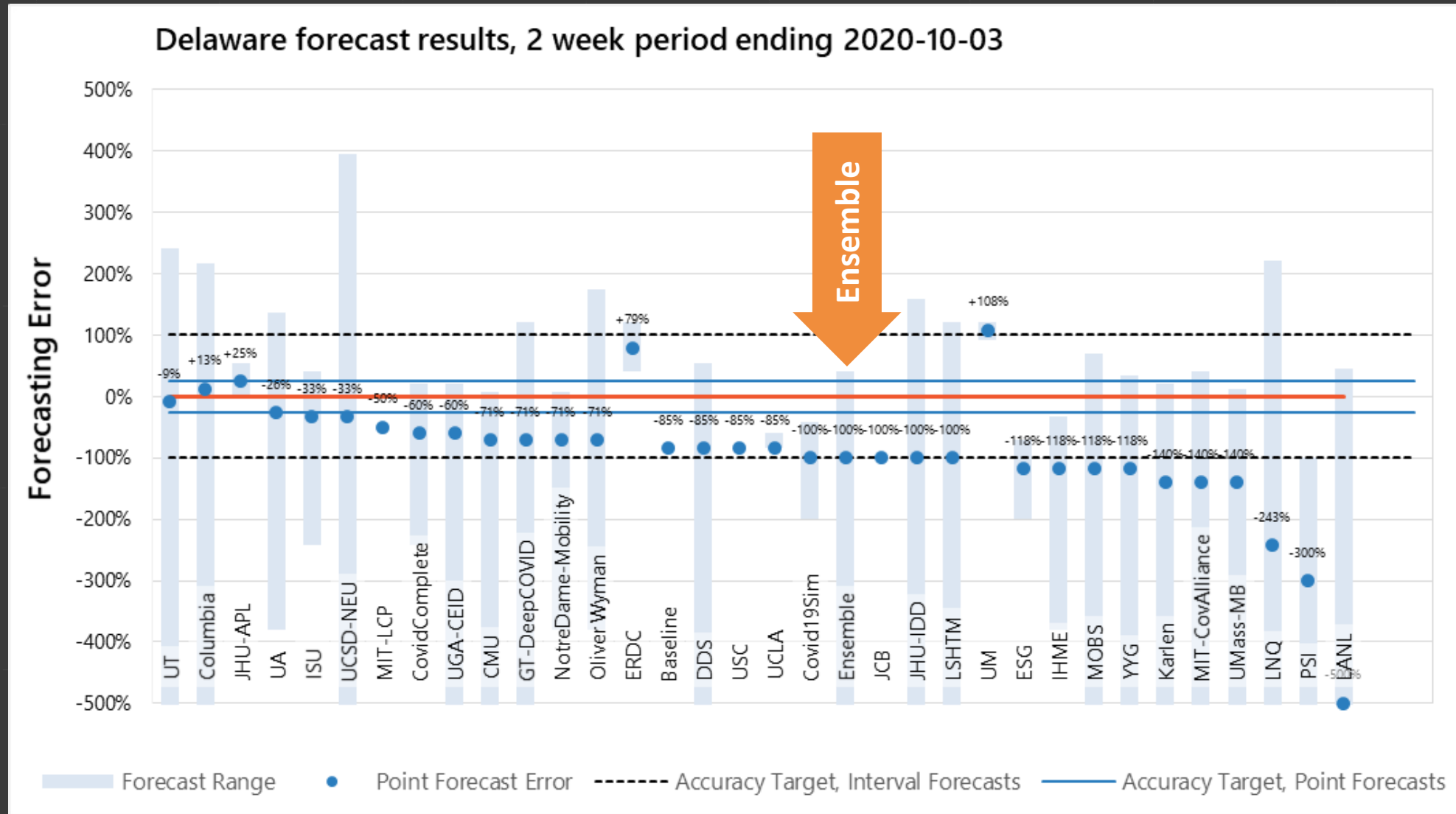


# State Forecast Examples

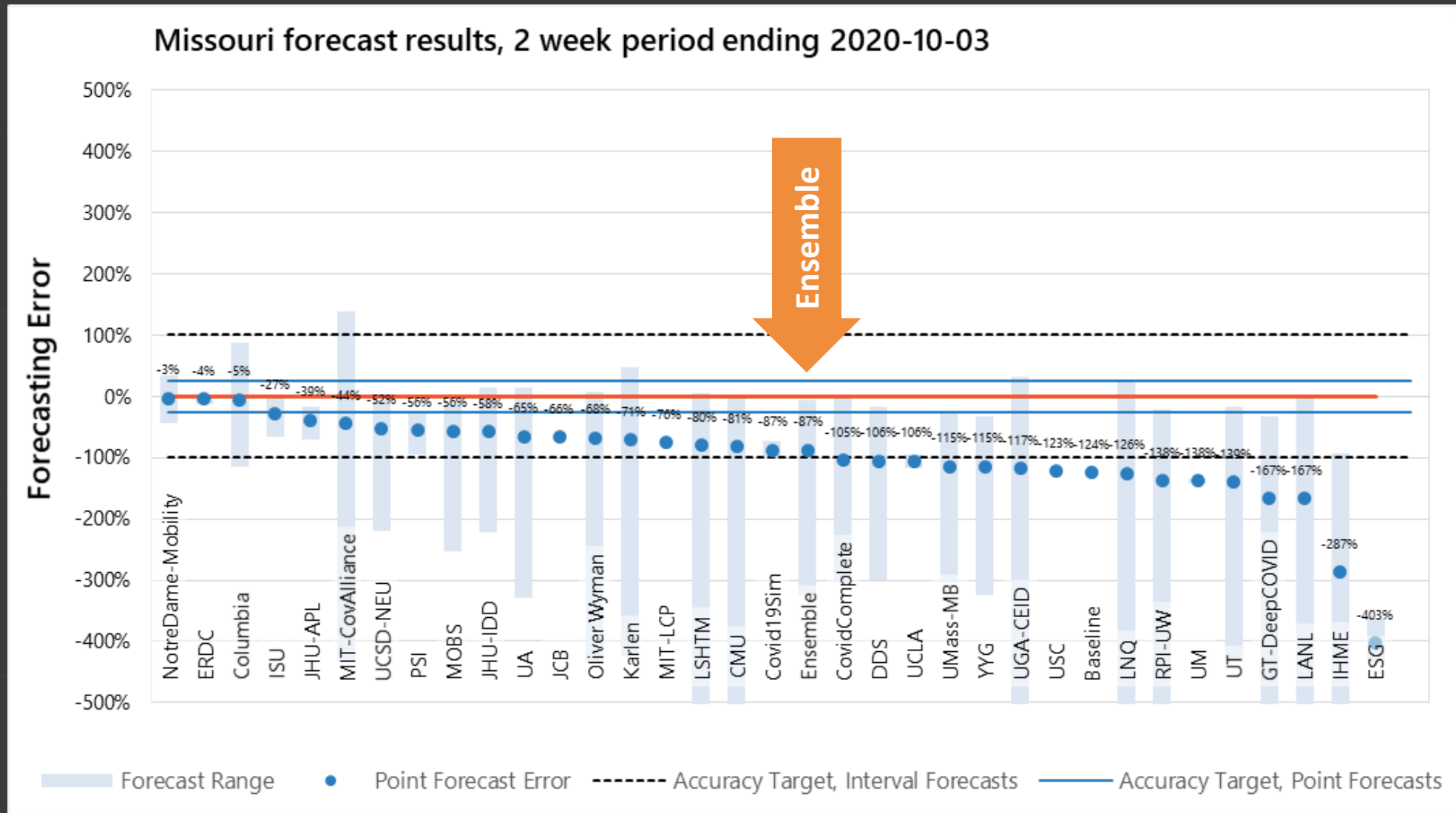




# State Forecast Examples

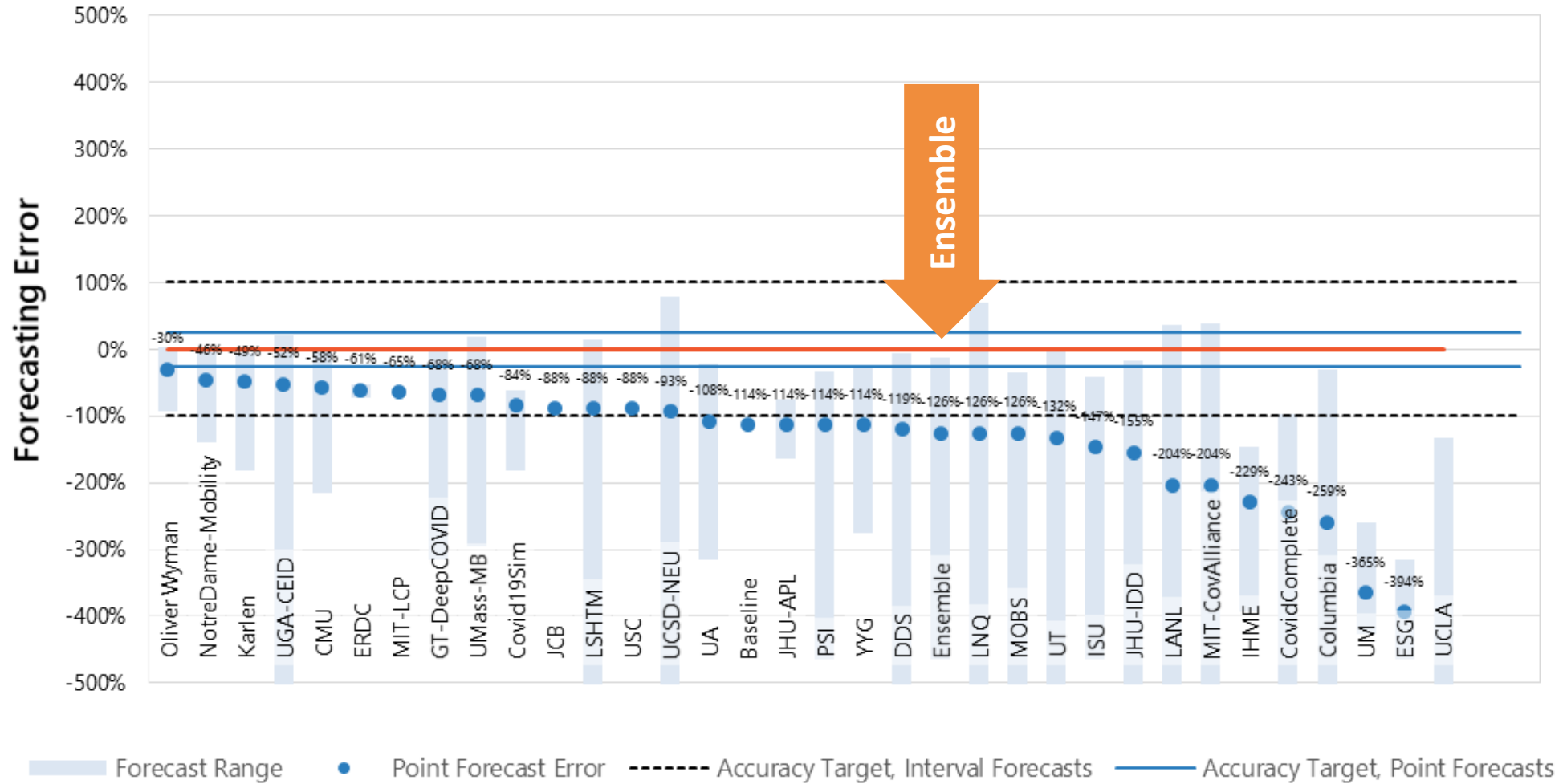


# State Forecast Examples



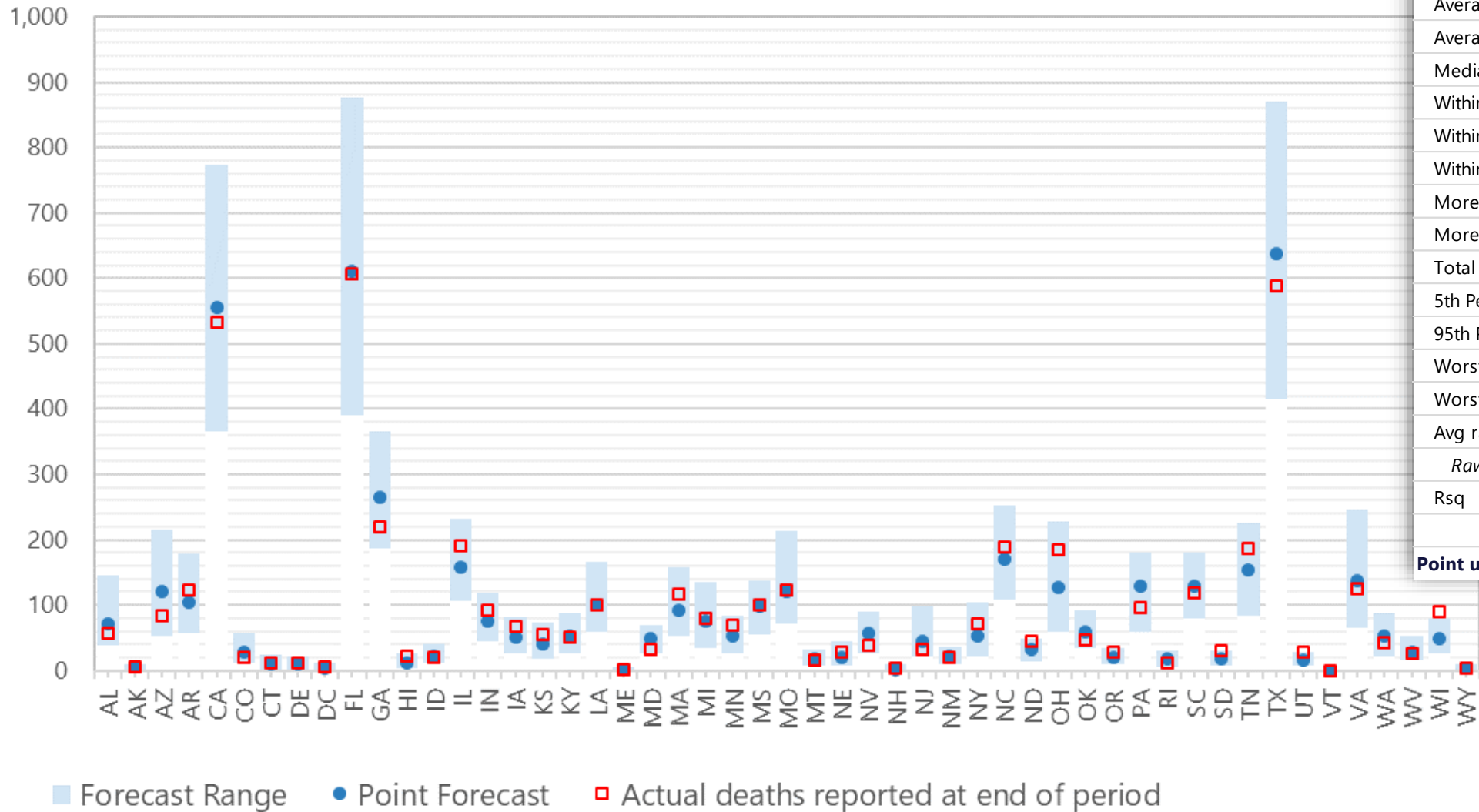
# State Forecast Examples

North Dakota forecast results, 2 week period ending 2020-10-03



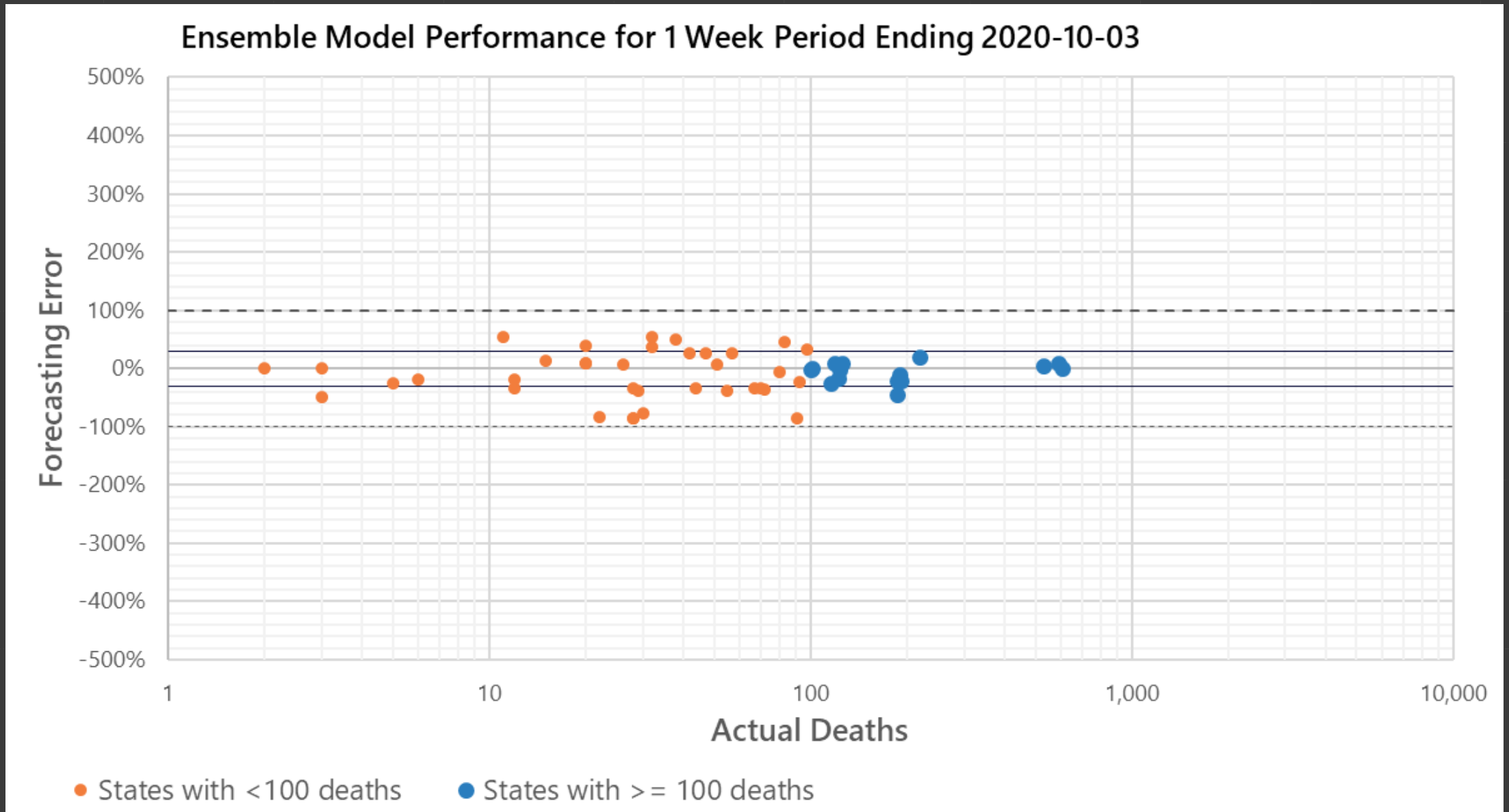
# State Forecasting at its Best: 1 Week Horizon

Ensemble Model Performance for 1 Week Period Ending 2020-10-03

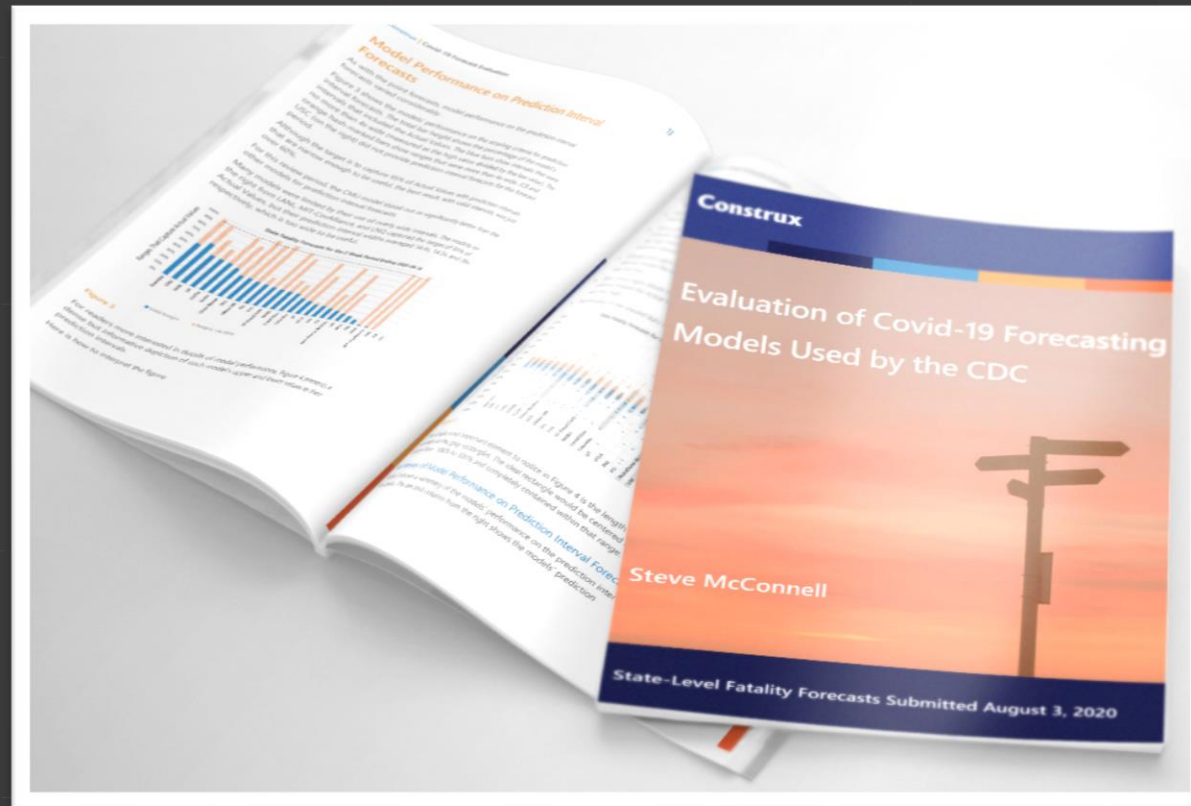


Point Forecasts	
Median error magnitude:	25%
Average error magnitude	27%
Average bias:	-8%
Median error bias:	-1%
Within 10%:	31%
Within 25% of actual:	53%
Within 50% of actual:	88%
More than 100% high	0%
More than 100% low	0%
Total Missed by more than 100%	0%
5th Percentile:	-80%
95th Percentile:	+47%
Worst overestimate:	+55%
Worst underestimate:	-87%
Avg raw error	14
Raw error as %	15%
Rsq	0.98
<b>Point usefulness</b>	
	<b>71%</b>

# — State Forecasting at its Best: 1 Week Horizon



# — More State Forecast Evaluations



[www.stevemcconnell.com/covid](http://www.stevemcconnell.com/covid)

# — More Overall Forecast Evaluations



[www.stevemccconnell.com/covid](http://www.stevemccconnell.com/covid)

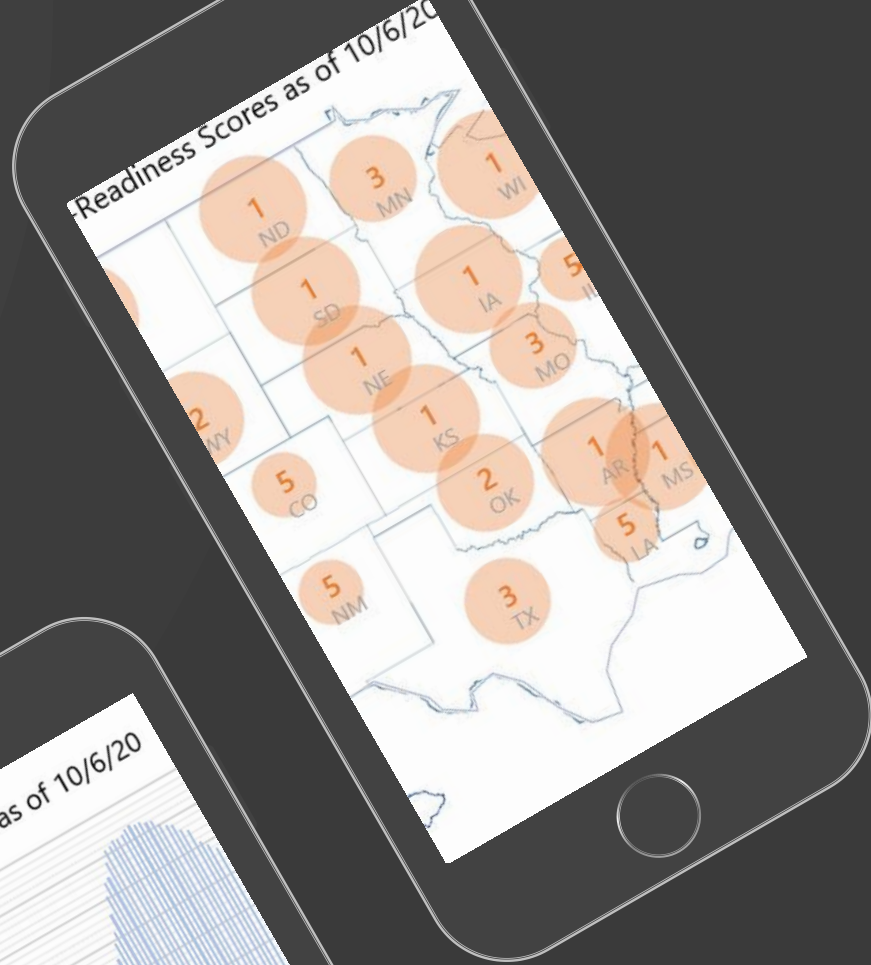


# Summary



# Summary

- Accurate forecasts are possible on 1-4 week horizons
- National forecasts were not accurate into early summer, but have become quite accurate in late summer and fall
- State forecasts are not as accurate as national forecasts; forecasts for bigger states tend to be more accurate, but not always
- The most accurate forecasts come from private individuals as often as they come from universities; the most prestigious universities are not producing accurate forecasts
- The CDC's Ensemble model is consistently among the most accurate forecast models for nearly all periods and types of forecasts



# Covid-19 Spin-Free Data Center

[stevemcconnell.com/covid](http://stevemcconnell.com/covid)

# — Useful Links

- SteveM's Covid-19 Data Center  
<https://www.stevemccconnell.com/covid>
- SteveM's article on detailed age and comorbidity risk  
<https://medium.com/@stevemcc/the-one-graph-you-need-to-understand-covid-19-comorbidities-5c7c8b64254f>
- CDC list of specific comorbidity risks  
<https://www.cdc.gov/coronavirus/2019-ncov/need-extra-precautions/evidence-table.html>
- CDC forecast hub  
<https://viz.covid19forecasthub.org/>
- Washington state data dashboard  
<https://coronavirus.wa.gov/what-you-need-know/covid-19-risk-assessment-dashboard>
- Tableau international comparison data  
[https://public.tableau.com/profile/jonas.nart#!/vizhome/COVID19\\_15844962693420/COVID19-TrendTracker](https://public.tableau.com/profile/jonas.nart#!/vizhome/COVID19_15844962693420/COVID19-TrendTracker)

# — Discussion