

Predicted Incrementality by Experimentation (PIE)

Florian Zettelmeyer

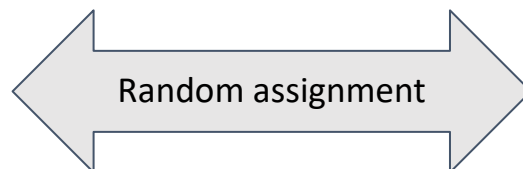
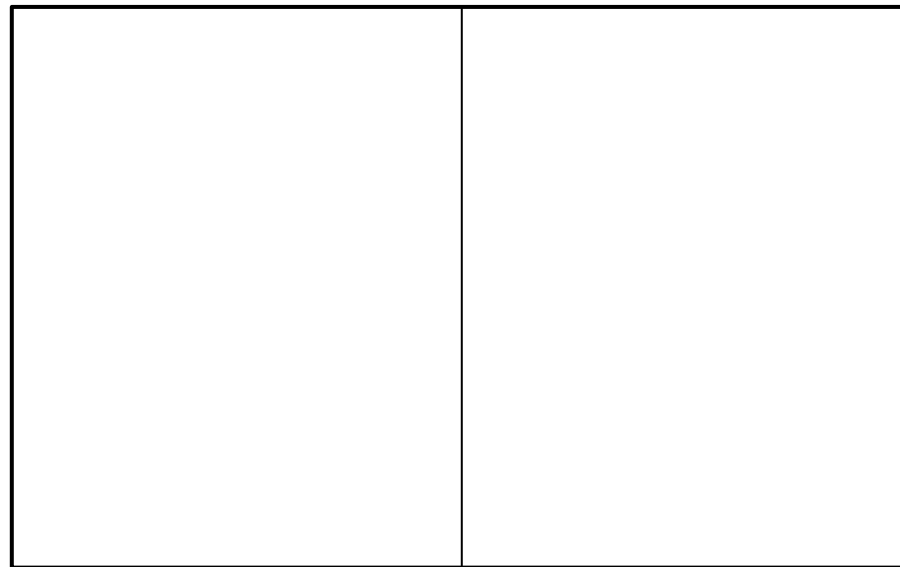
Kellogg School of Management, Northwestern University

Joint work with Brett Gordon (Northwestern) and Robert Moakler (Meta)

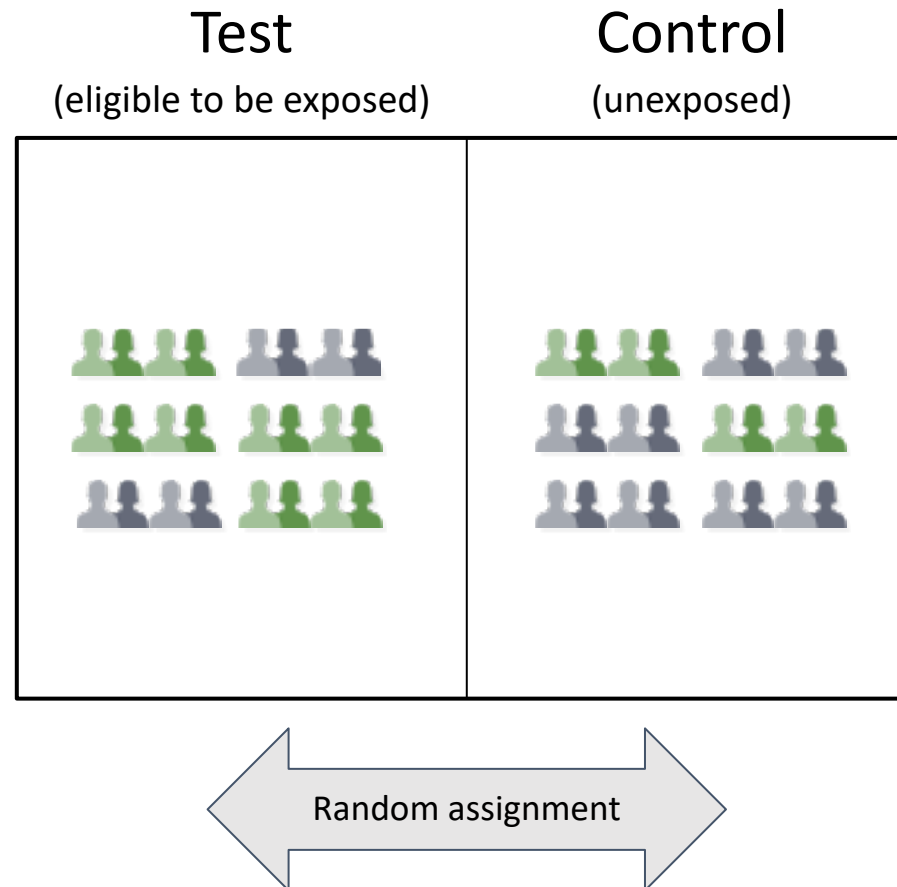
MSI 2024 Summit

Randomized controlled trials (RCTs) are recognized as the “gold standard” to measure incremental (causal) effects

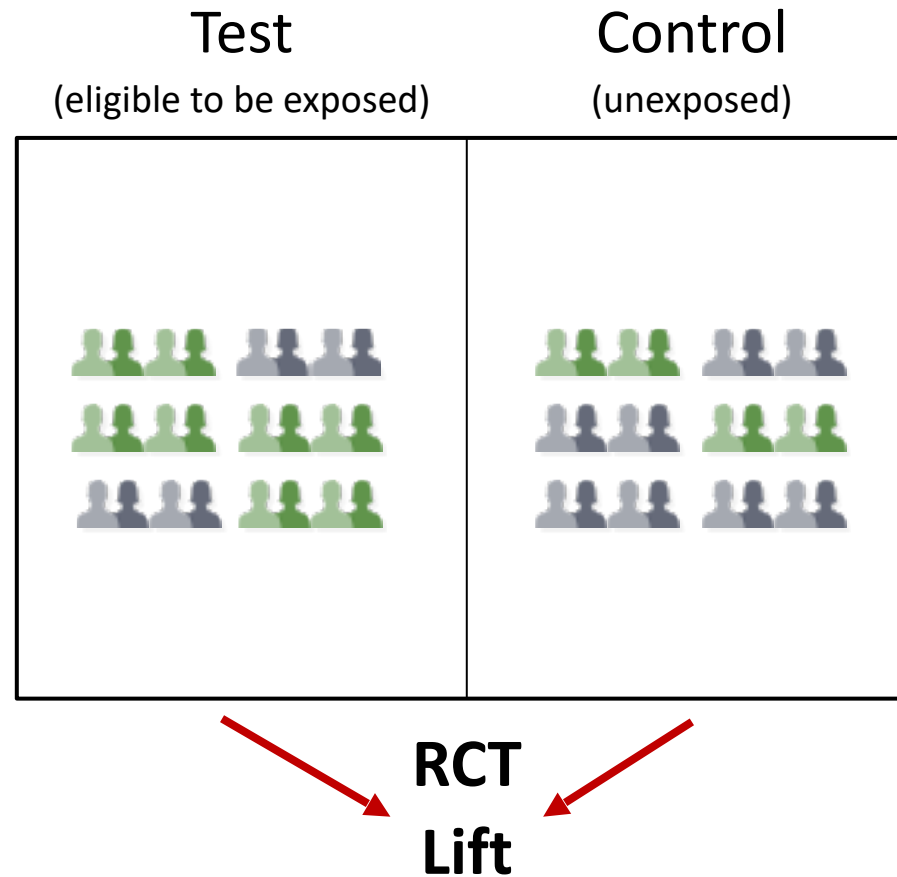
Target Audience



Randomized controlled trials (RCTs) are recognized as the “gold standard” to measure incremental (causal) effects



Randomized controlled trials (RCTs) are recognized as the “gold standard” to measure incremental (causal) effects

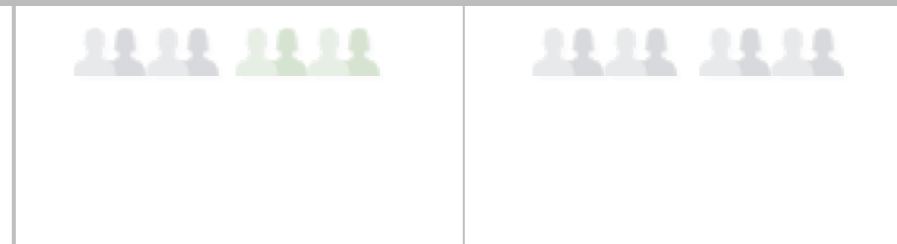


Randomized controlled trials (RCTs) are recognized as the “gold standard” to measure incremental effects

Test

Control

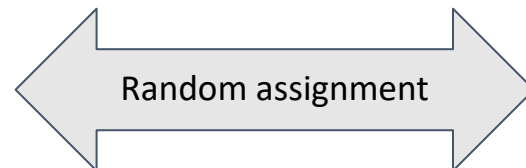
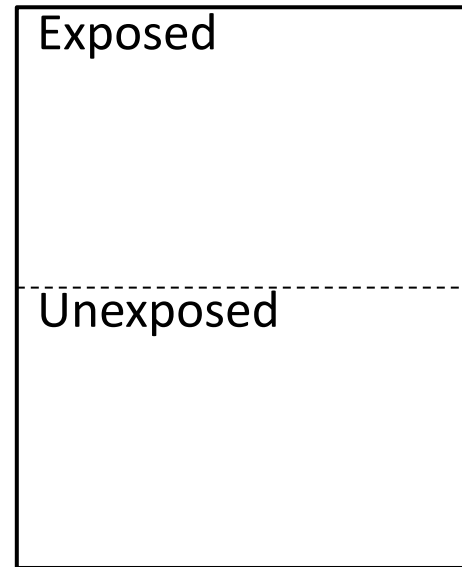
But what if you can't run an RCT?



RCT
Lift

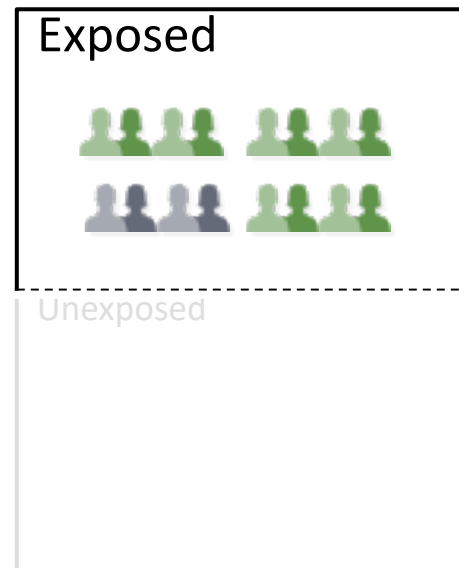
Imagine we don't have an RCT...

Target Audience



Proxy metrics are a common approach

Target Audience



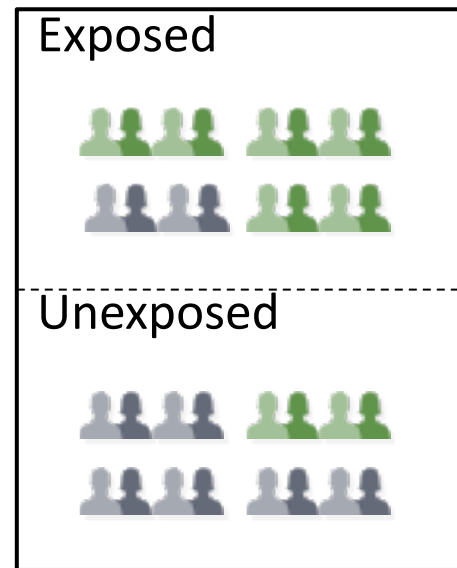
Last Click (LC):

- Start with an outcome (e.g., purchase)
- Attribution window (e.g., 7 days)
- “Attribute” purchase to ad that was clicked last in attribution window
- No counterfactual (control group)

Problem: Can over or underestimate incrementality

Alternatively, we could compare outcomes between people **who saw** versus **did not see** the ad campaign

Target Audience



Problem: Suffers from “selection bias” into ad exposure

“Undoing” the selection induced by ad-targeting algorithms using causal inference approaches has been unsuccessful

EVIDENCE

- *Gordon, Zettelmeyer, Chapsky, Bhargava (2019), Marketing Science*
 - Compare RCTs with (observational) program evaluation approaches
 - 15 studies, hand-selected -> cannot come close replicating RCT results
- *Gordon, Moakler, Zettelmeyer (2023), Marketing Science*
 - 1673 RCTs, representative
 - SPSM, Double/Debiased ML + Deep Learning for propensity score
 - From 30 to (nearly all ~ 5000) logged features at FB
 - Equally depressing ...

RCT Lift vs. Lift from SPSM & DML

Funnel Level of Outcome	Median Lift		
	RCT	SPSM	DML
Upper	29%		
Mid	18%		
Lower	5%		

The **RCT lift** estimates and ...

- **SPSM** are statistically different in $1482 / 1673 = 89\%$ of the RCTs
- **DML** are statistically different in $1258 / 1673 = 75\%$ of the RCTs

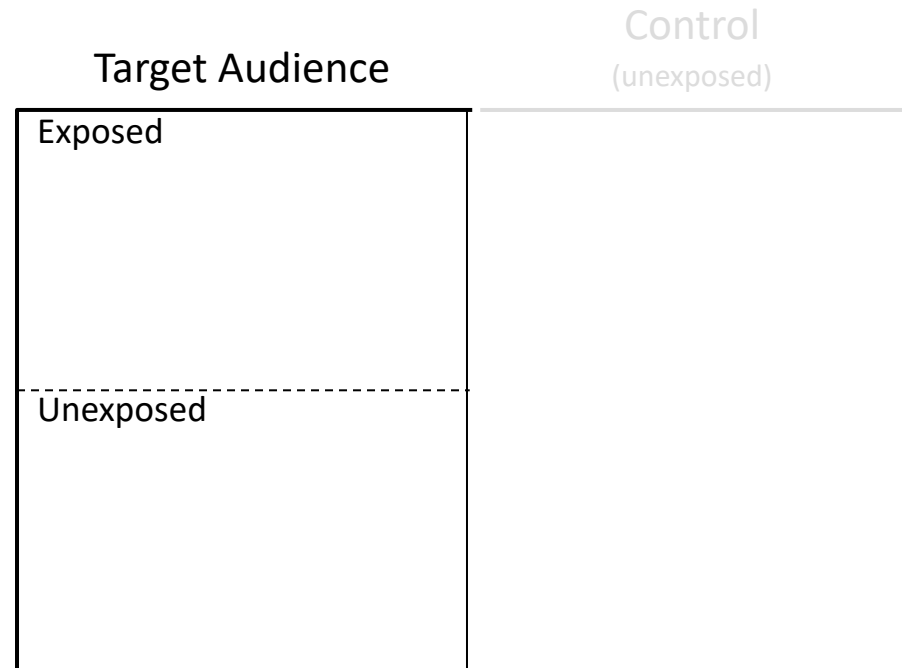
So, what should advertisers do?

Don't have the data for observational methods

And can't run RCTs all the time

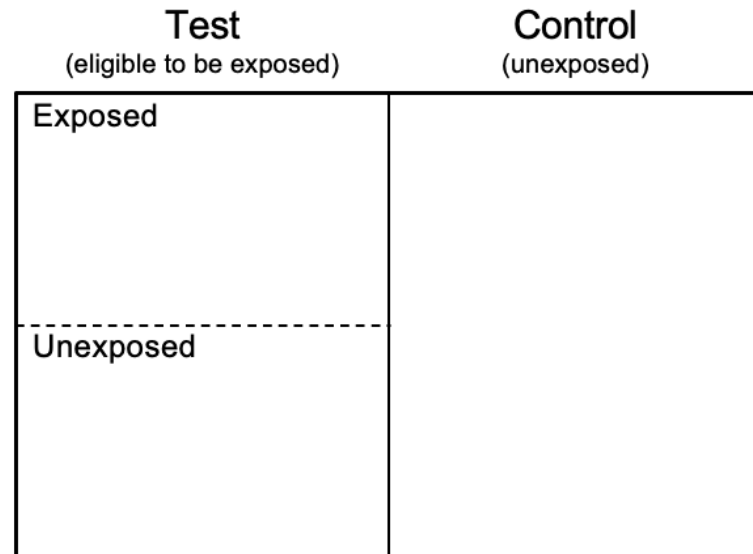
We have tried to estimate the causal effect of advertising without RCTs by controlling for **user-level selection bias**

“Traditional” Causal Inference methods



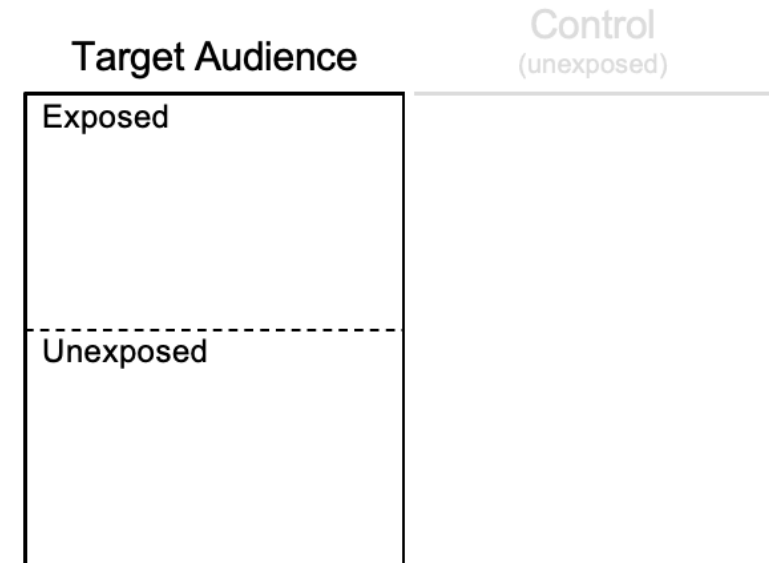
We often have RCTs for a **subset** of advertising campaigns...

Ad campaigns as RCTs



RCT Lift

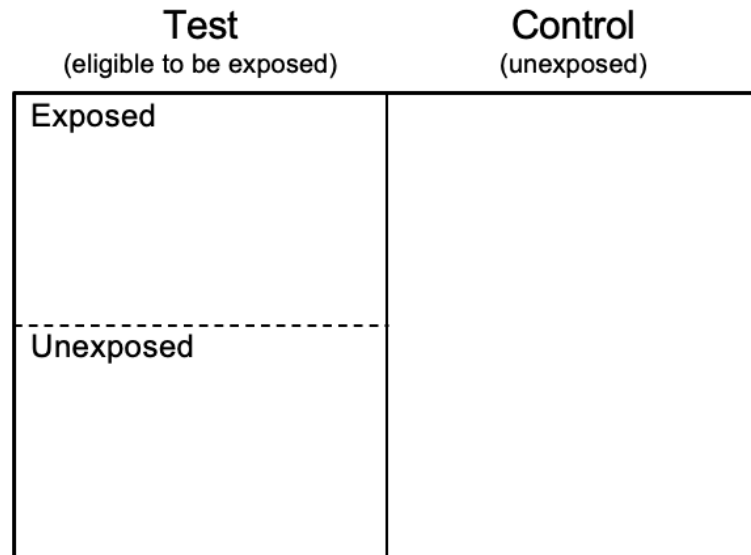
Ad campaigns not as RCTs



RCT Lift?

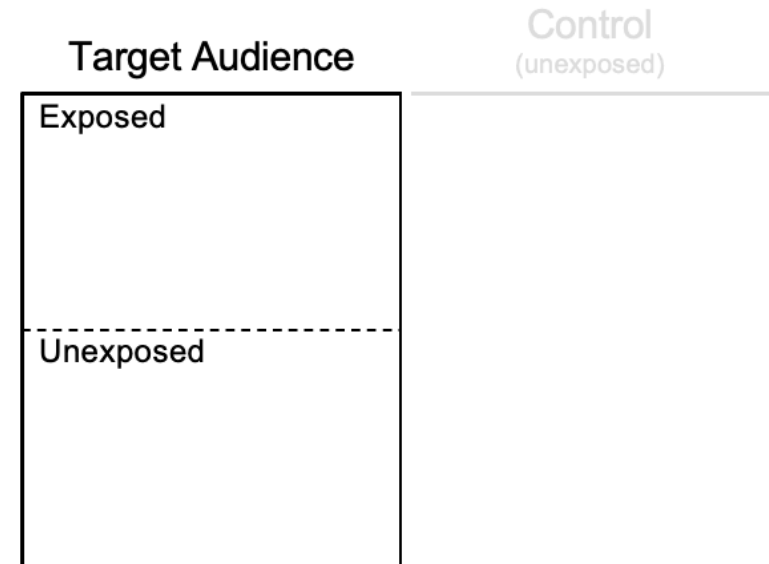
Predictive Incrementality by Experimentation (PIE)

Ad campaigns as RCTs



RCT Lift

Ad campaigns not as RCTs



RCT Lift



Using a database of RCTs, how well could we predict a **new campaign's RCT Lift** if it was **not run as an RCT?**

- Unit of observation is an RCT campaign, not a user
- Shift to predictive models, instead of those from causal inference
- Use our RCT dataset to assess the performance of PIE models

PIE – a first cut

Incremental Conversion in $RCT_r = \theta(\# \text{ of } LC^w \text{ Conversions in } RCT_r) + \text{Error}_r$

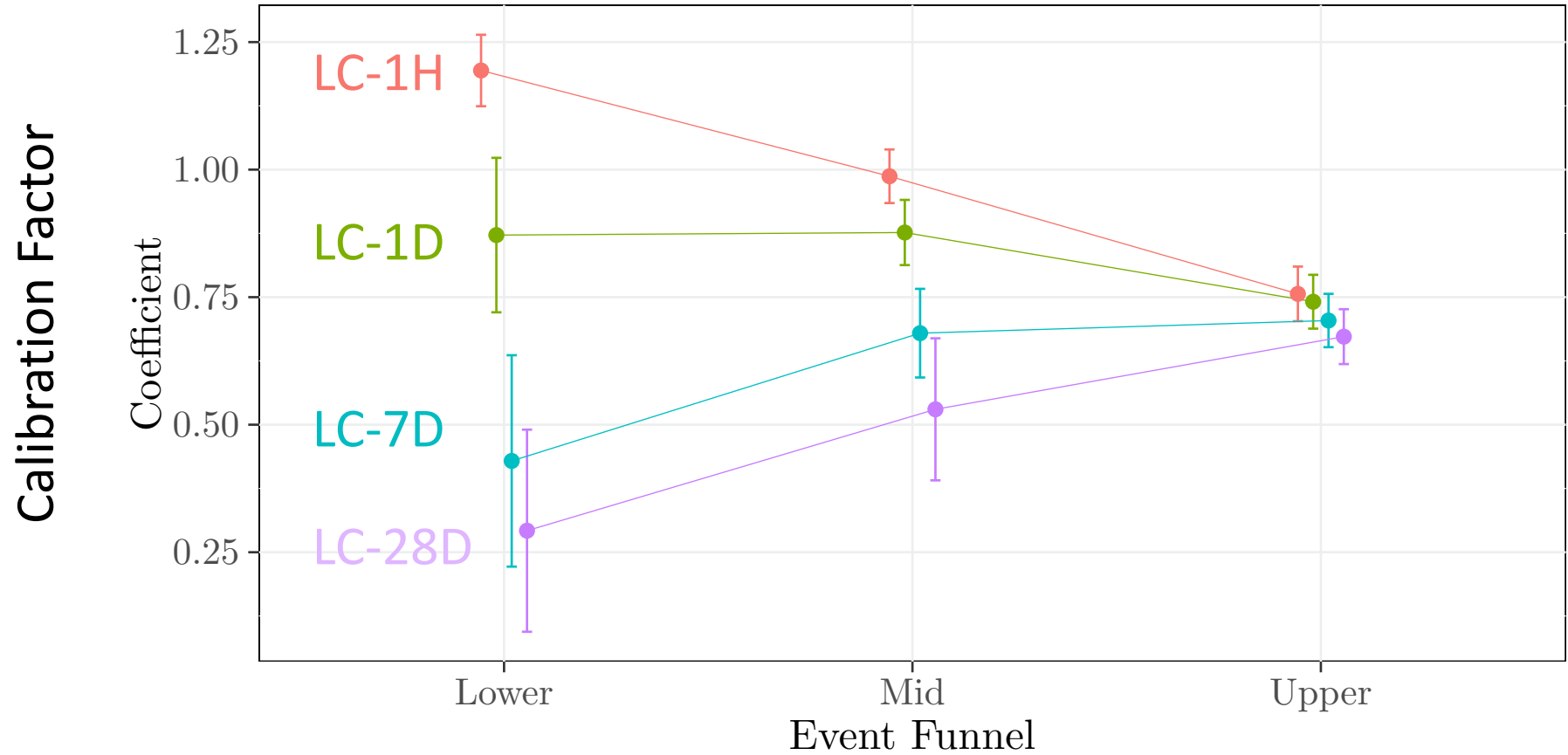


“Calibration Factor”

We estimate this model separately by:

- Funnel levels (Lower, Mid, Upper)
- Last click attribution windows $w \in \{1 \text{ hour}, 1 \text{ day}, 7 \text{ days}, 28 \text{ days}\}$

Calibration factor model by funnel and attribution window



CI based on robust SEs

To generalize this approach, we normalize incremental conversions by ad spend and add more features

RCT Incremental Conversions
per Dollar (ICPD)

$$\text{ICPD}_r = f(X_r^{\text{pre}}, X_r^{\text{post}}; \theta) + \varepsilon_r,$$

Advertiser-campaign characteristics
known **before** campaign was run

Campaign: targeting criteria, bidding
params, optimization goal, budget, etc.

Advertiser: vertical, experimentation
experience, etc.

Proxy metrics known **after** campaign was run

Last click conversion counts by {1H, 1D, 7D, 28D}

If available, other post-campaign metrics could
be used (e.g., view-through conversion counts)

Key: None of the features rely on the RCT control group

We try a variety of models and two feature sets

$$f(X_r^{\text{pre}}, X_r^{\text{post}}; \theta)$$

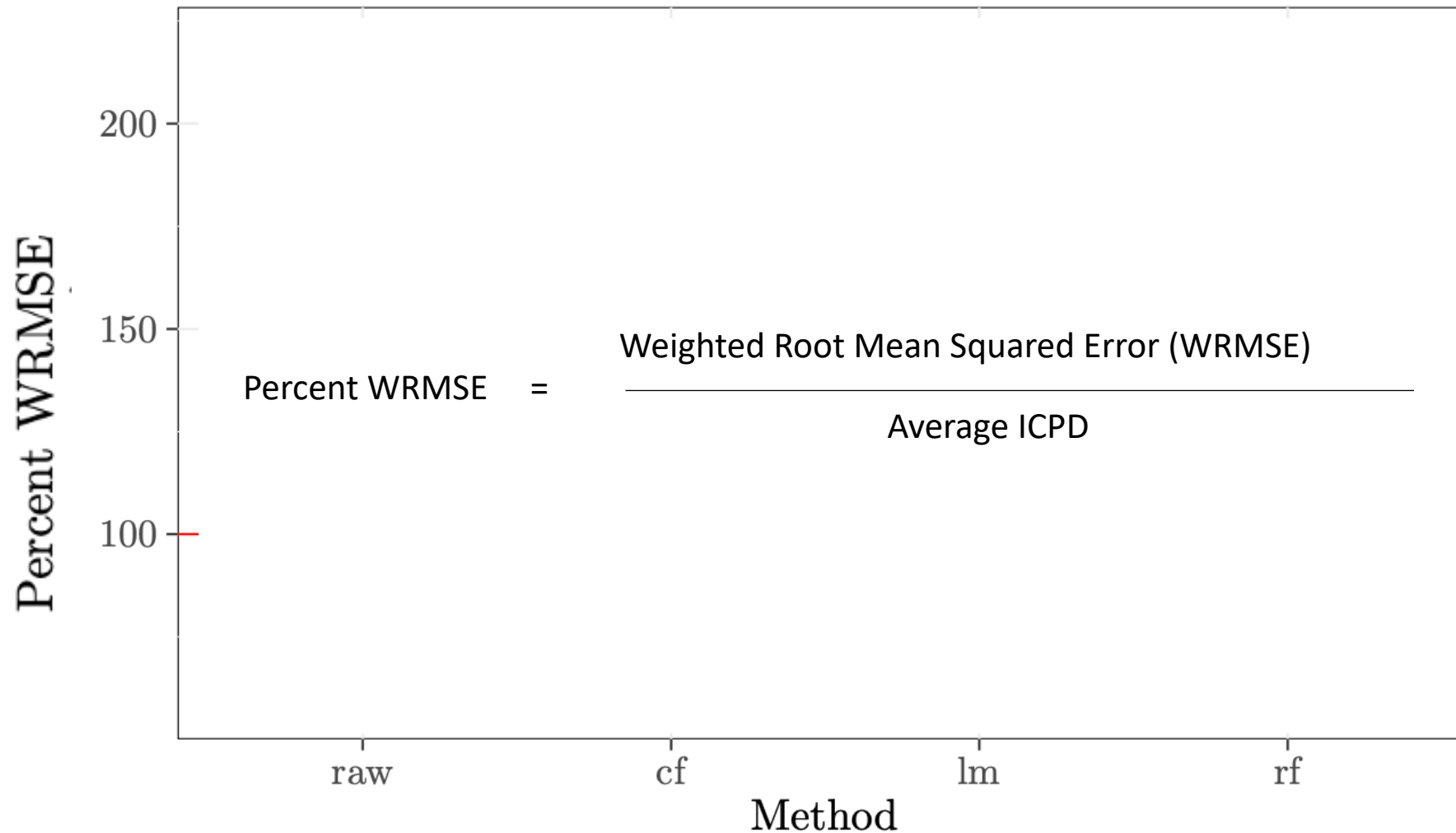
- Models for

- **raw**: how well does each last click metric perform by itself?
- **cf**: calibration factor model (the “first cut”)
- **lm**: linear regression
- **rf**: random forest

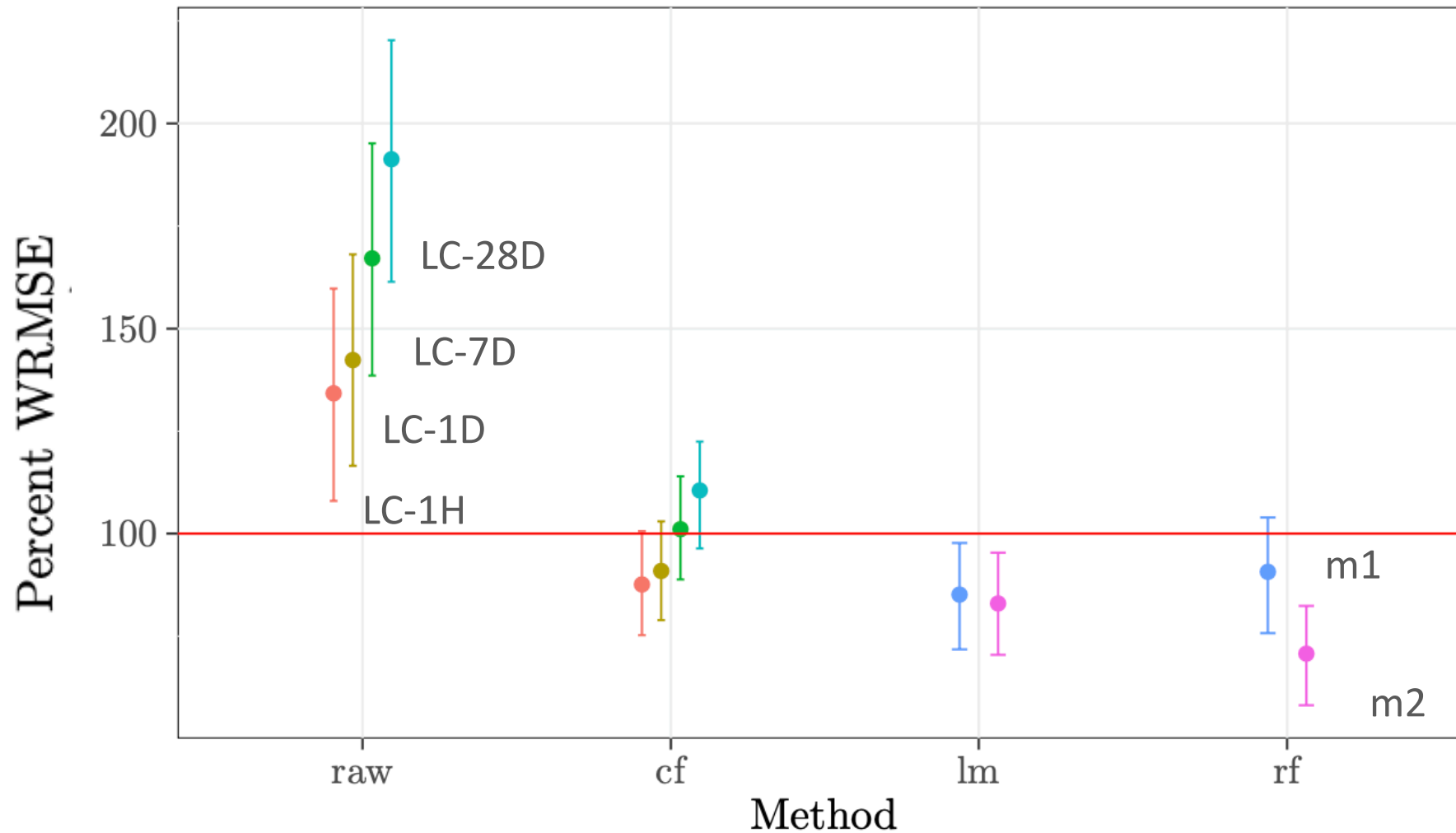
- Features used in **lm** and **rf**

- **m1**: $X_r^{\text{post}} = \{\text{LC-1h, LC-1d, LC-7d, LC-28d}\}$ $X_r^{\text{pre}} = \{\}$
- **m2**: $X_r^{\text{post}} = \{\text{LC-1h, LC-1d, LC-7d, LC-28d}\}$ $X_r^{\text{pre}} = \{\text{everything}\}$

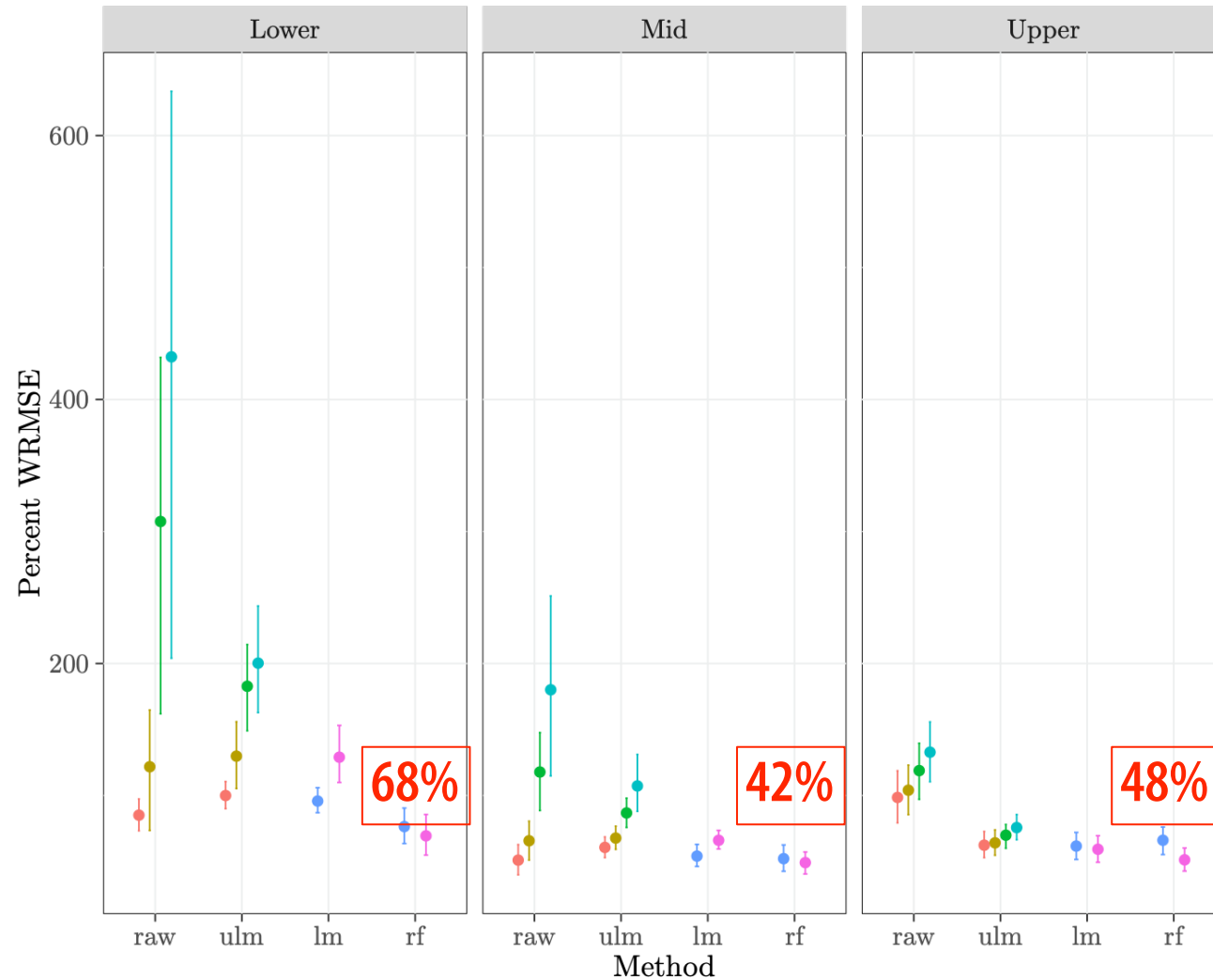
Assess the models using **Percent WRMSE** based on Leave-One-Out Predictions



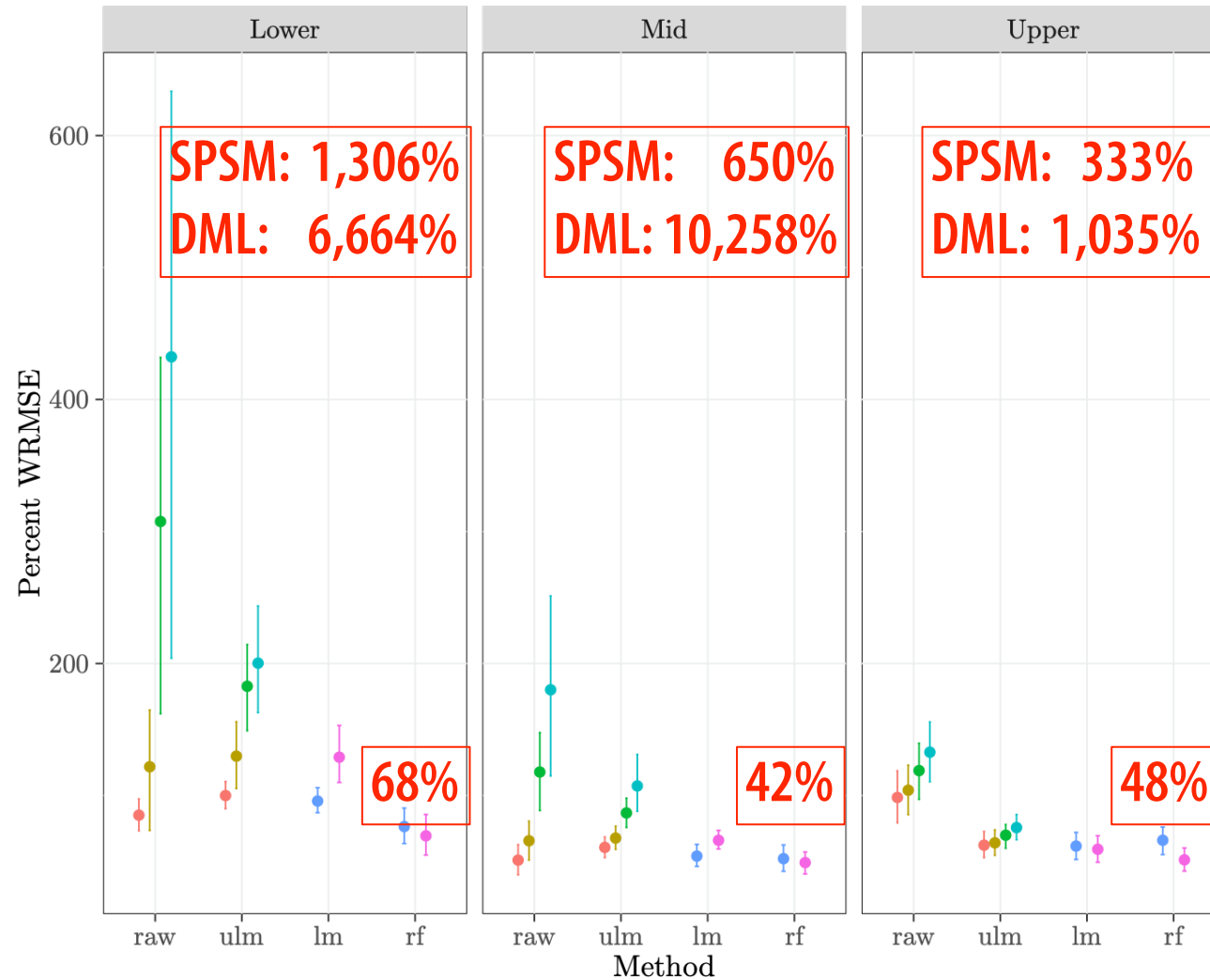
A random forest (rf) with all features performs best



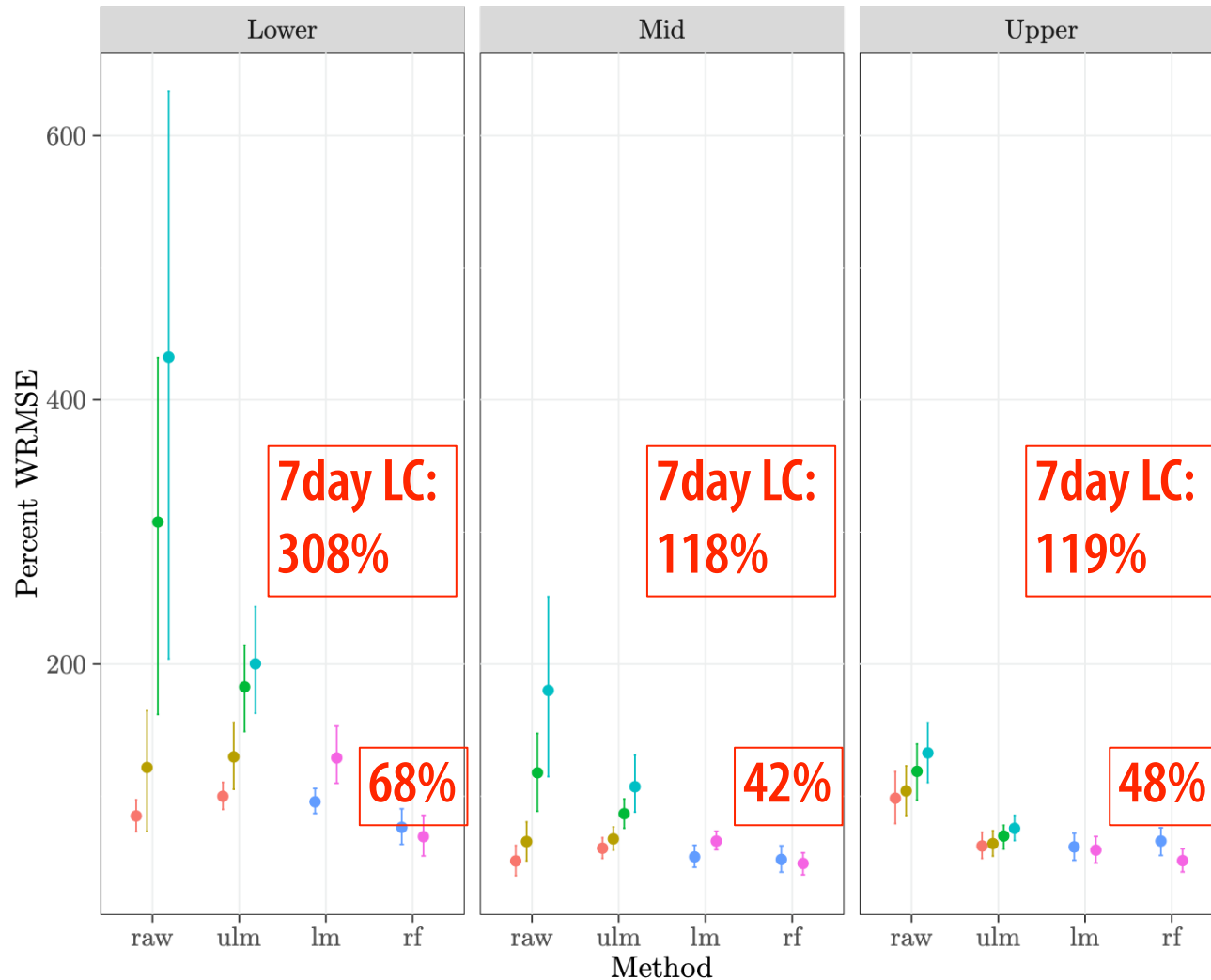
The percent WRMSE for the best specification is **around 50%**



PIE does much better than SPSM or DML



PIE also does **substantially better** than the industry-standard 7-day Last-Click Attributed Conversions



Recap: What exactly is PIE?

PROBLEM - Advertisers can't rely on observational data, nor can they always run RCTs ...
...but they still need to measure advertising effects

Recap: What exactly is PIE?

PROBLEM - Advertisers can't rely on observational data, nor can they always run RCTs ...
...but they still need to measure advertising effects

KEY IDEA - Advertisers can still run RCTs for a subset of campaigns ...
- Use RCTs to predict ad effects for new campaign's that were not run as an RCT

Recap: What exactly is PIE?

PROBLEM - Advertisers can't rely on observational data, nor can they always run RCTs ...
...but they still need to measure advertising effects

KEY IDEA - Advertisers can still run RCTs for a subset of campaigns ...
- Use RCTs to predict ad effects for new campaign's that were not run as an RCT

STANDARD - Estimate ad effects using campaign and user characteristics before campaign ran

PREDICTION - Uses "pre-determined" features: e.g. FinTech ads have higher lift than CPG ads

Recap: What exactly is PIE?

PROBLEM - Advertisers can't rely on observational data, nor can they always run RCTs ...
...but they still need to measure advertising effects

KEY IDEA - Advertisers can still run RCTs for a subset of campaigns ...
- Use RCTs to predict ad effects for new campaign's that were not run as an RCT

STANDARD PREDICTION - Estimate ad effects using campaign and user characteristics before campaign ran
- Uses "pre-determined" features: e.g. FinTech ads have higher lift than CPG ads

PIE - Predict ad effects using performance features after the campaign starts
- Uses "post-determined" features: e.g. clicks, last-click conversion, page views, ...
- Anything in the treatment group that is correlated with causal ad effects
- Move from causal inference to a prediction problem

PIE will work when post-determined features are predictive and the relationship is stable

WHEN PIE IS LIKELY TO WORK

- **PREDICTIVE:** RCTs need to measure causal effect (not too noisy)
Post-determine features need to contain some causal signal
(empirical question)
- **STABLE:** We need the nature of this relationship to be stable over time
(i.e., no concept shift)

So, is PIE useful for practice?

HOW SHOULD WE MEASURE CAUSAL ADVERTISING EFFECTS?

- Attribution models are biased
- Causal inference models don't work
- RCTs are only viable option ... but are infeasible at scale
- PIE makes RCTs scalable
- In our testing PIE has smaller confidence intervals than raw RCTs

Thank you!