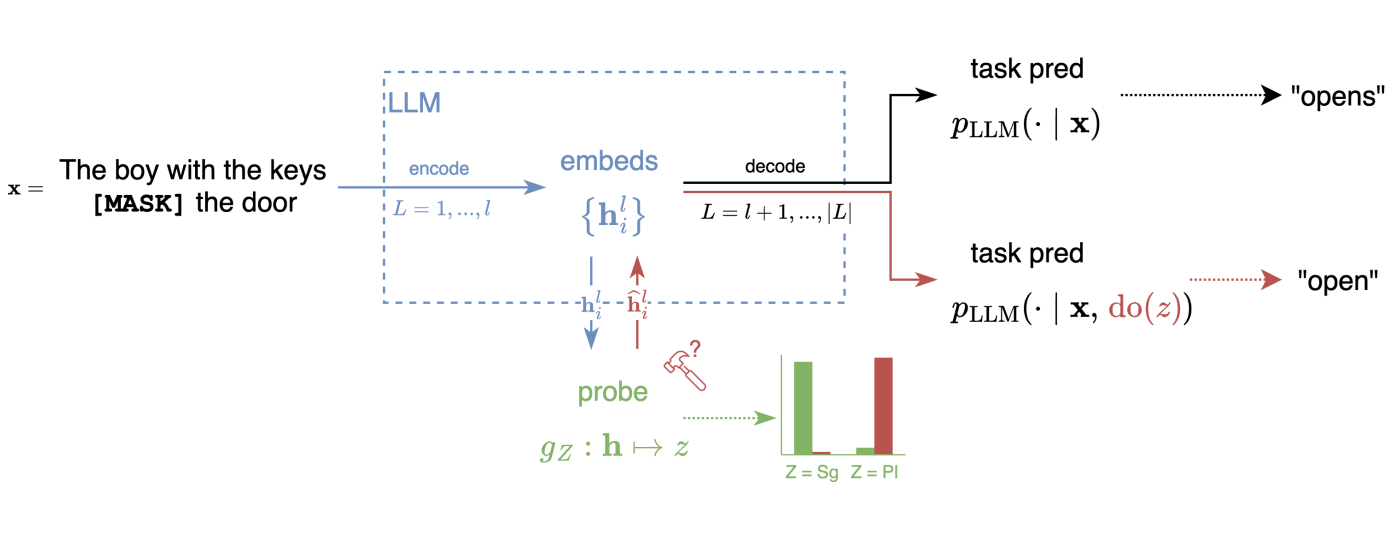# Measuring the Reliability of Causal Probing Methods:
## Tradeoffs, Limitations, and the Plight of Nullifying Interventions
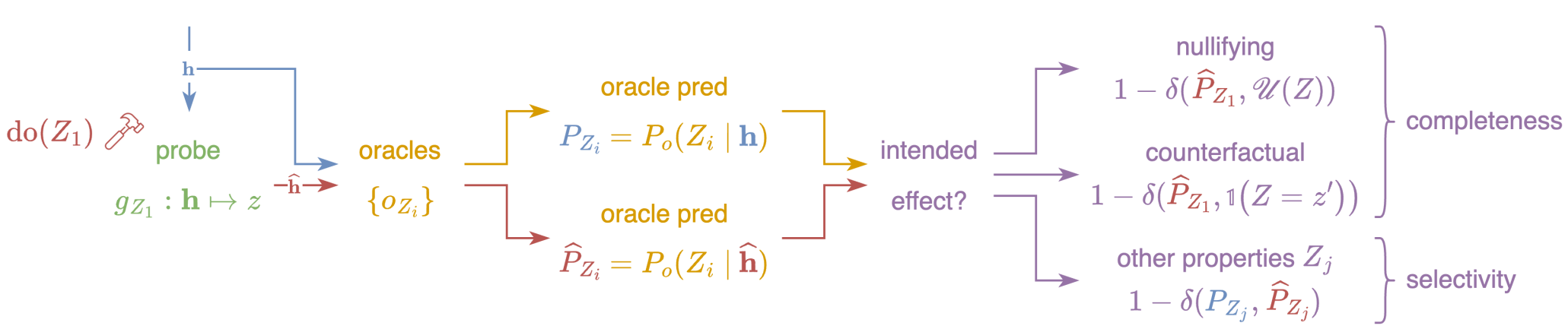
Marc E. Canby*, Adam Davies*, Chirag Rastogi, Julia Hockenmaier
*Equal contribution      University of Illinois Urbana-Champaign
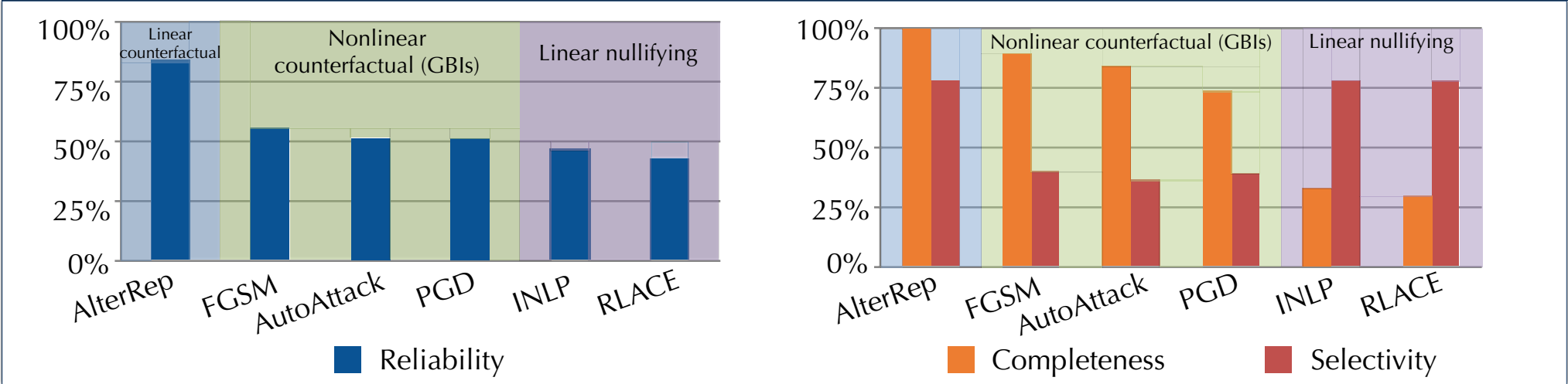
# How *reliable* are interp methods?



## *Completeness:* How much is targeted property damaged?
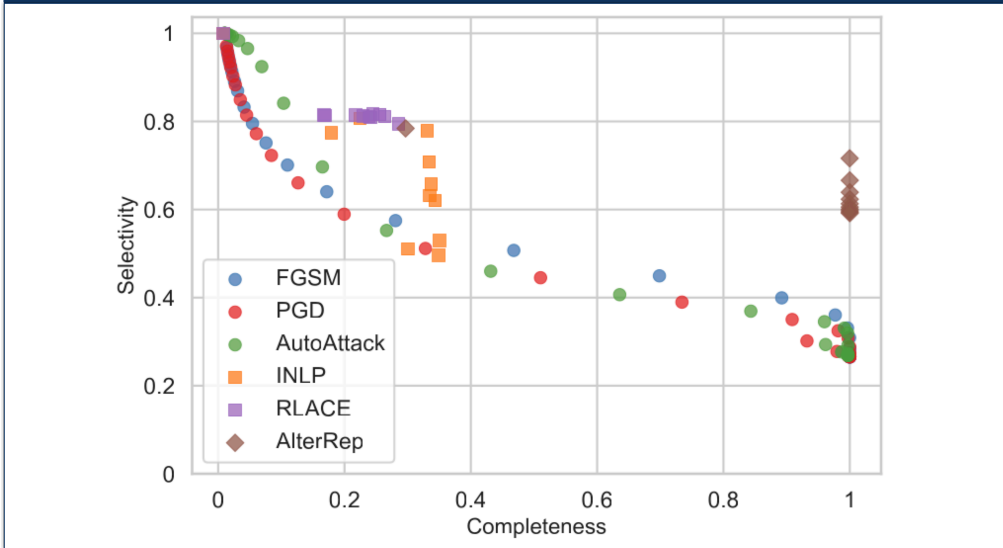## *Selectivity:* How little collateral damage is done?



# Experiments: Subject-verb agreement

`The boy with the keys [MASK] the door` → `open` vs. `opens`

$Z_1 = $ # subject ("`boy`")       $Z_j = $ # distractor ("`keys`")

How reliable are interventions on `[MASK]` vector in final layer of BERT?

## How reliable are causal probing interventions? Nullifying methods are not reliable!



Reliability



Completeness        Selectivity

## Completeness & selectivity are a tradeoff!



## Greater reliability → Greater Δ Task Acc.