

Text as Causal Mediators

Research Design for Causal Estimates of Differential Treatment of Social Groups via Language Aspects

We propose a research design for (observational) causal estimates of the effects of social group signals on speakers' responses via language as causal mediators.

Identification assumptions

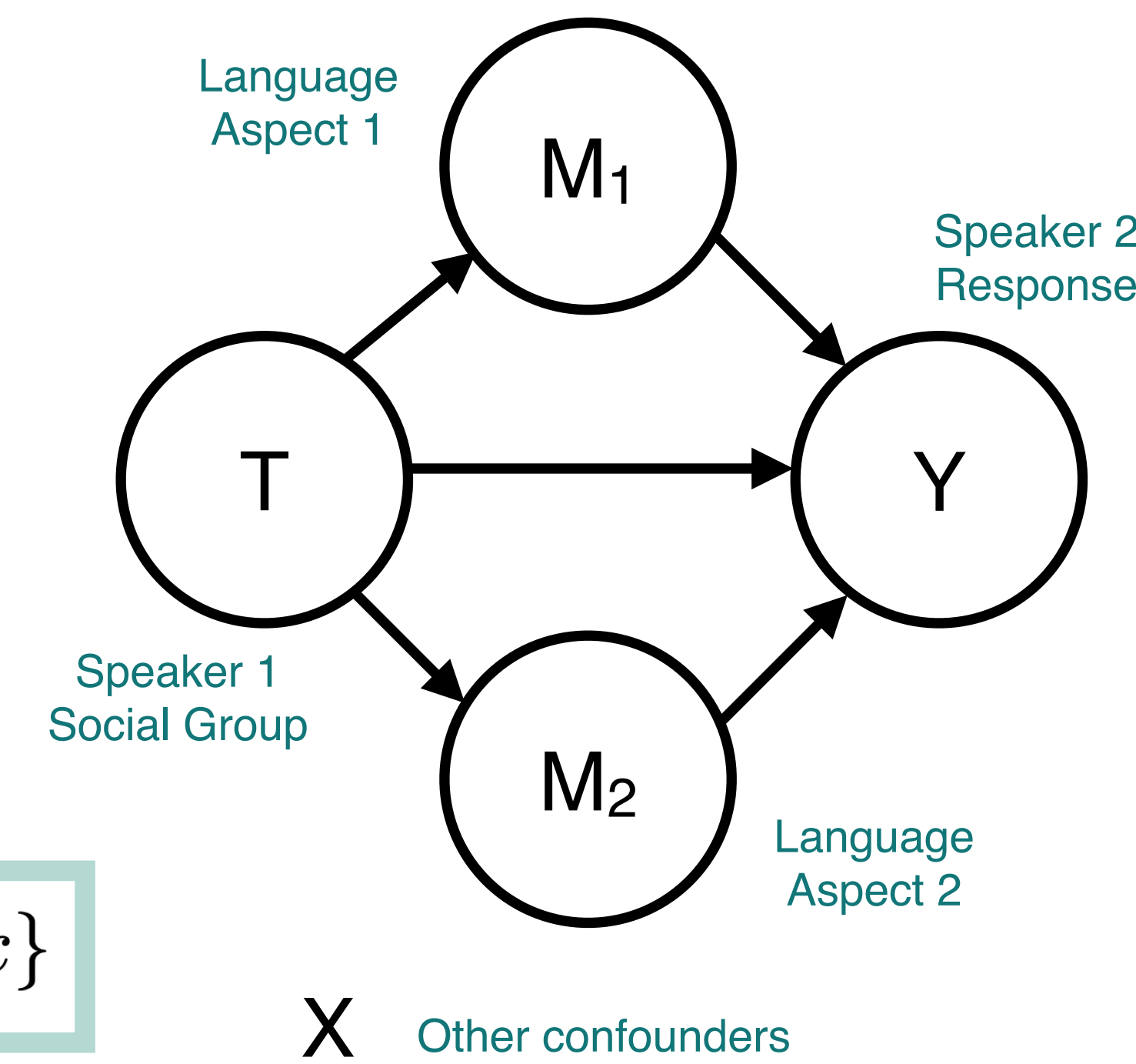
- Sequential ignorability (Imai et al, 2010)

$$\{Y_i(t', m), M_i(t)\} \perp\!\!\!\perp T_i \mid X_i = x$$

$$Y_i(t', m) \perp\!\!\!\perp M_i(t) \mid \{T_i = t, X_i = x\}$$
- Mediator Independence

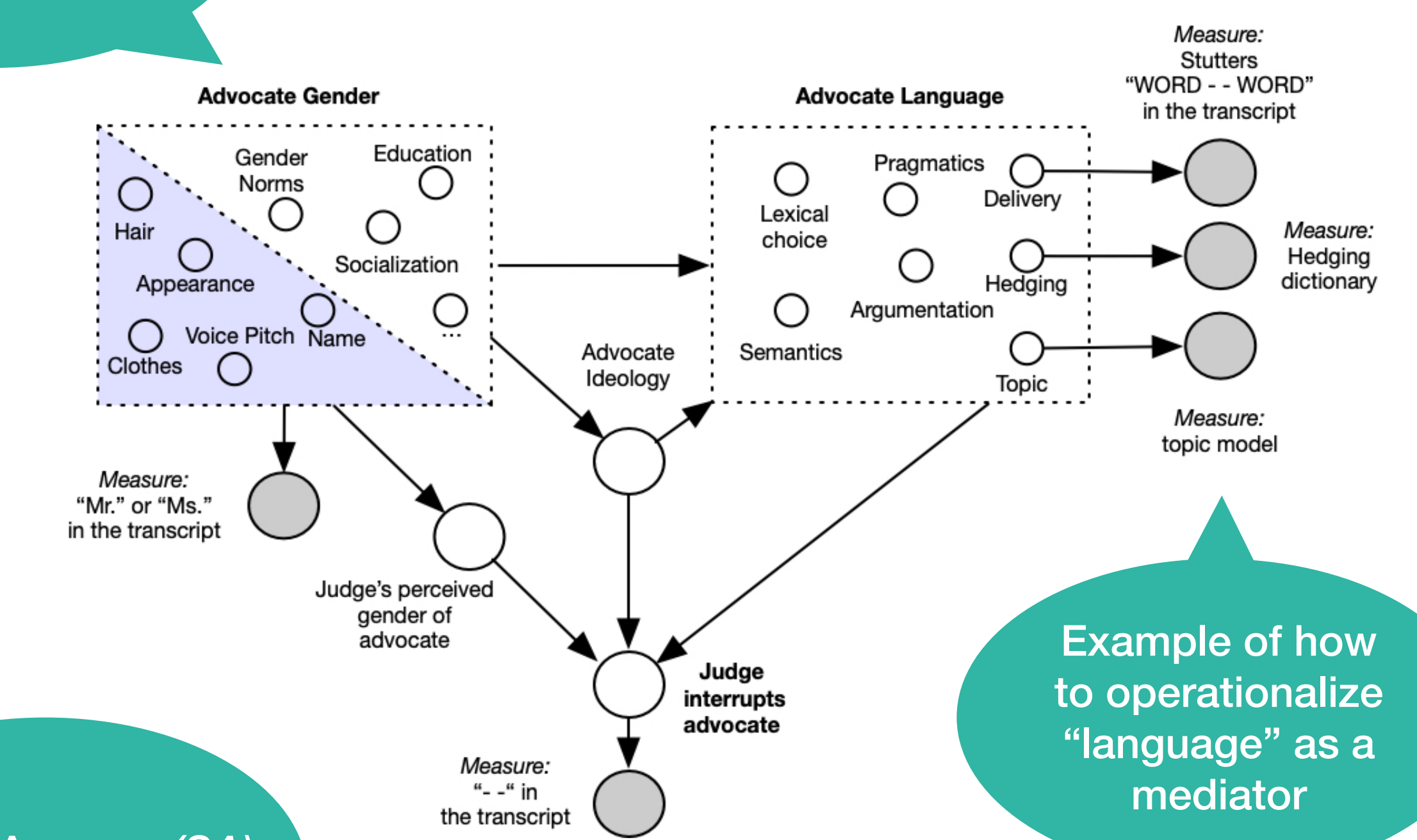
$$\forall j, j' : M_i^j(t) \perp\!\!\!\perp M_i^{j'}(t) \mid \{T_i = t, X_i = x\}$$

Simplified causal graph



Addresses "social group" as a causal treatment

Expanded causal graph



Example of how to operationalize "language" as a mediator

Estimation

$$SA-NDE^j = \frac{1}{N} \sum_{i=1}^N \sum_{x \in \mathcal{X}} \sum_{m \in \mathcal{M}^j} \left(\hat{f}^j(Y|M_i^j = m, T_i = 1, X_i = x) - \hat{f}^j(Y|M_i^j = m, T_i = 0, X_i = x) \right) \hat{g}^j(m|T_i = 0, X_i = x)$$

$$SA-NIE^j = \frac{1}{N} \sum_{i=1}^N \sum_{x \in \mathcal{X}} \sum_{m \in \mathcal{M}^j} \hat{f}^j(Y|M_i^j = m, T_i = 0, X_i = x) \left(\hat{g}^j(m|T_i = 1, X_i = x) - \hat{g}^j(m|T_i = 0, X_i = x) \right)$$

- Fit \hat{f} , \hat{g} on train set. Then at inference time, apply the fitted models to real confounders from the test/inference set and "counterfactual" treatment and mediator values.
- We will also need models of mediators given raw text. An open question is whether to infer this separately or jointly with \hat{f} and \hat{g} .

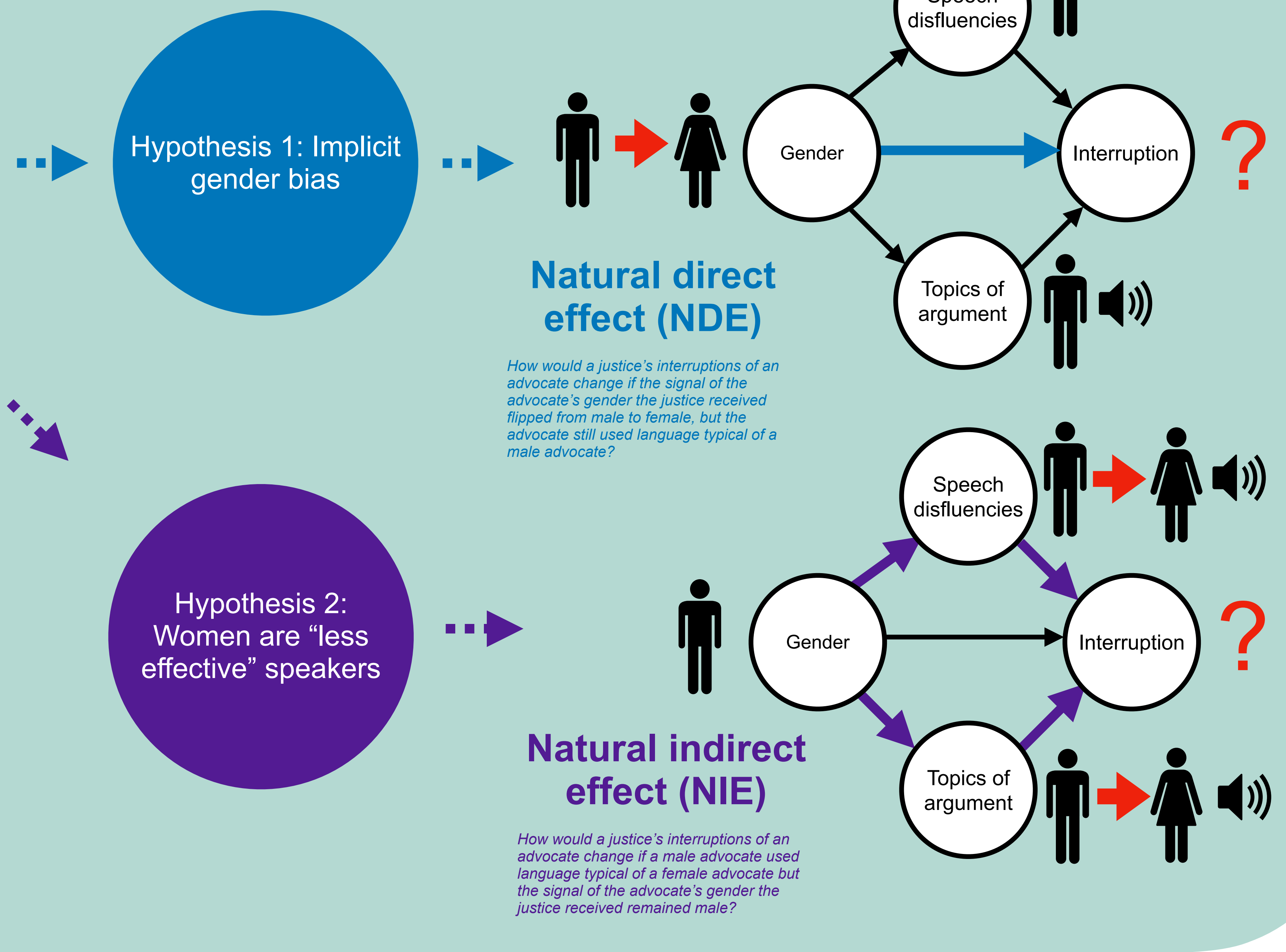
Future Directions

- Empirical estimates from real data
- Address causal dependence between temporal utterances
- Analyze between-judge and between-court temporal estimates

Theoretical case study from U.S. Supreme Court oral arguments



Q: Why do some justices interrupt female advocates more than male advocates?



Example: Lozano v. Montoya Alvarez (2013)

Ann O'Connell Adams (advocate): Well— Interruption

Antonin Scalia (justice): I mean, it seems to me it just makes that article impossible to apply consistently country to country. Hedging Speech Disfluencies

Ann O'Connell Adams (advocate): No, I don't think so. And—and, the other signatories have—have almost all, I mean I think the Hong Kong court does say that it doesn't have discretion, but it said in that case nevertheless it would, even if it had discretion, it wouldn't order the children returned. But the other courts of signatory countries that have interpreted Article 12 have all found a discretion, whether it be in Article 12 or in Article 8. And if I— Interruption

Antonin Scalia (justice): Have they exercised it? Have they exercised it, that discretion which they say is there?



Katherine A. Keith
 Then: University of Massachusetts Amherst, College of Information and Computer Sciences
 Now: Postdoctoral researcher, AI2
 Next: Assistant Professor, Williams College
kkeith@cs.umass.edu



Doug Rice
 University of Massachusetts Amherst, Department of Political Science
drice@legal.umass.edu



Brendan O'Connor
 University of Massachusetts Amherst, College of Information and Computer Sciences
brenocon@cs.umass.edu

Paper & References:



Causal Inference + NLP Workshop, EMNLP 2021