

# Generative AI Powered Causal Inference

Kosuke Imai

Harvard University

Oxford Computational Political Science Group

University of Oxford

May 20, 2025

Joint work with Kentaro Nakamura (Harvard Kennedy School)

# Motivation

- Generative AI is transforming medicine, education, marketing, etc.
- Can methodologists get some help from generative AI too?

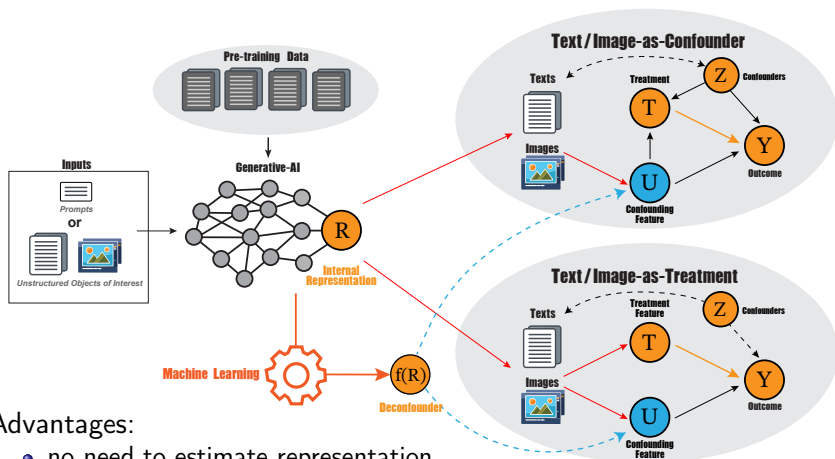
How can a methodologist use LLM to help improve research? Give me a short answer.



A methodologist can use large language models (LLMs) to enhance research by automating literature reviews, generating hypotheses, designing study frameworks, and analyzing data. LLMs can also assist in drafting and refining research papers, coding qualitative data, and even suggesting relevant methodologies or statistical approaches, making the research process faster and more rigorous.

# Using Generative AI to Improve Causal Inference

- Generative AI Powered Causal Inference (GPI)
- GenAI-assisted causal inference with unstructured data



- Advantages:
  - no need to estimate representation
  - avoid functional form assumptions
  - better empirical performance

# Generative AI: Definition and Assumption

- Deep generative model:

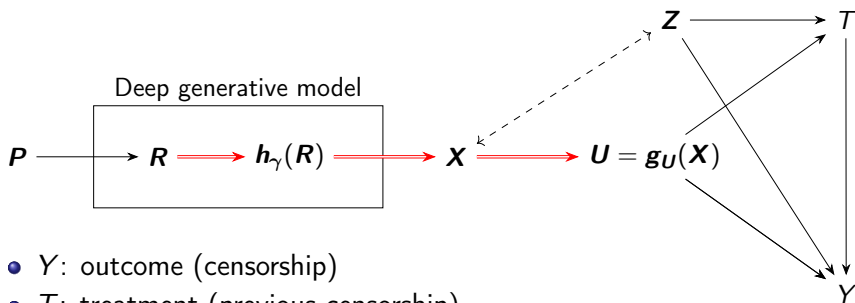
$$\mathbb{P}(\mathbf{X}_i \mid \mathbf{h}_\gamma(\mathbf{R}_i)),$$
$$\mathbb{P}(\mathbf{R}_i \mid \mathbf{P}_i).$$

- $\mathbf{P}_i$ : prompt
  - $\mathbf{X}_i$ : unstructured generated
  - $\mathbf{R}_i$ : hidden states or internal representations
  - $\mathbf{h}_\gamma(\mathbf{R}_i)$ : deterministic function from hidden states to the last layer
- 
- **Deterministic decoding**:  $\mathbb{P}(\mathbf{X}_i \mid \mathbf{h}_\gamma(\mathbf{R}_i))$  is degenerate
  - Use of open-source GenAI for replicability

# Text as Confounder: Chinese Censorship (Roberts *et al.* 2020)

- Do Chinese social media users who had their post censored become more likely to be censored for later posts or self-censor themselves?
  - Treatment: whether or not a post was censored
  - Outcomes: censorship during four weeks after a censored post
    - 1 number of posts
    - 2 proportion of censored posts
    - 3 proportion of missing posts
  - structural confounders: lagged outcomes, date of the post (dummies)
  - text-as-confounder: contents of posts
- Original analysis: Matching (CEM) with topic proportions (STM) and propensity score (inverse regression)
- Our reanalysis:
  - Text reuse with Llama 3
  - Apply the proposed method:
    - 1 entire sample (4155 users; 75324 Weibo posts)
    - 2 matched sample (628 users; 879 posts)

# Assumptions



- $Y$ : outcome (censorship)
- $T$ : treatment (previous censorship)
- $\mathbf{Z}$ : observed structured confounding variables
- $\mathbf{X}$ : unstructured confounding object
- $\mathbf{U} = g_U(\mathbf{X})$ : unknown and unstructured confounding variables
- **Strong latent ignorability:**

$$\{Y_i(t)\}_{t \in \mathcal{T}} \perp\!\!\!\perp T_i \mid \mathbf{Z}_i = \mathbf{z}, \mathbf{U}_i = \mathbf{u}, \quad \text{for all } \mathbf{z} \in \mathcal{Z}, \mathbf{u} \in \mathcal{U}$$
$$\mathbb{P}(T_i = t \mid \mathbf{Z}_i = \mathbf{z}, \mathbf{U}_i = \mathbf{u}) > 0 \quad \text{for all } t \in \mathcal{T}, \mathbf{z} \in \mathcal{Z}, \mathbf{u} \in \mathcal{U}$$

# Identification

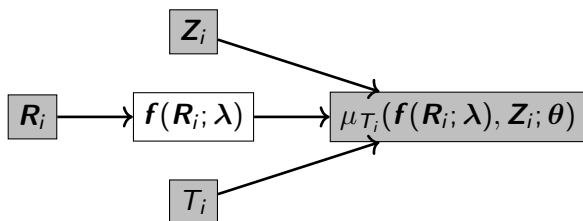
- There exists a **deconfounder**  $\mathbf{f} : \mathbb{R}^r \mapsto \mathbb{R}^q$  with  $q \leq r$  that satisfies:

$$Y_i \perp\!\!\!\perp \mathbf{R}_i \mid T_i, \mathbf{Z}_i, \mathbf{f}(\mathbf{R}_i)$$

- Adjusting for the deconfounder and  $\mathbf{Z}$  identifies the marginal distribution of any potential outcome  $Y(t)$ :

$$\mathbb{P}(Y_i(t) = y) = \int_{\mathbb{R}^q} \int_{\mathcal{Z}} \mathbb{P}(Y_i = y \mid T_i = t, \mathbf{Z}_i, \mathbf{f}(\mathbf{R}_i)) dF(\mathbf{Z}_i) dF(\mathbf{R}_i)$$

# Estimation via Neural Network



- Conditional expectation function:

$$\mu_{T_i}(\mathbf{f}(\mathbf{R}_i), \mathbf{Z}_i) := \mathbb{E}[Y_i(t) \mid \mathbf{f}(\mathbf{R}_i), \mathbf{Z}_i]$$

- Loss function for the outcome model and deconfounder:

$$\{\hat{\lambda}, \hat{\theta}\} = \operatorname{argmin}_{\lambda, \theta} \frac{1}{N} \sum_{i=1}^N \{Y_i - \mu_{T_i}(\mathbf{f}(\mathbf{R}_i; \lambda), \mathbf{Z}_i; \theta)\}^2$$

- Estimate the propensity score using the estimated deconfounder

$$\pi(\mathbf{f}(\mathbf{R}_i, \hat{\lambda}), \mathbf{Z}_i) = \mathbb{P}(T_i = 1 \mid \mathbf{f}(\mathbf{R}_i, \hat{\lambda}), \mathbf{Z}_i)$$



# Double Machine Learning (Chernozhukov et al. 2018)

- Cross-fitting for the binary treatment case:
  - 1 randomly divide the data into  $K$  folds
  - 2 for each  $k = 1, \dots, K$ , use the  $k$ th fold as the test set and the remaining  $k - 1$  folds as the training set
    - 1 randomly split the training set further into two subsets
    - 2 use the first subset to estimate outcome model and deconfounder
    - 3 use the second subset to estimate propensity score given deconfounder
  - 3 Compute the ATE estimator as:

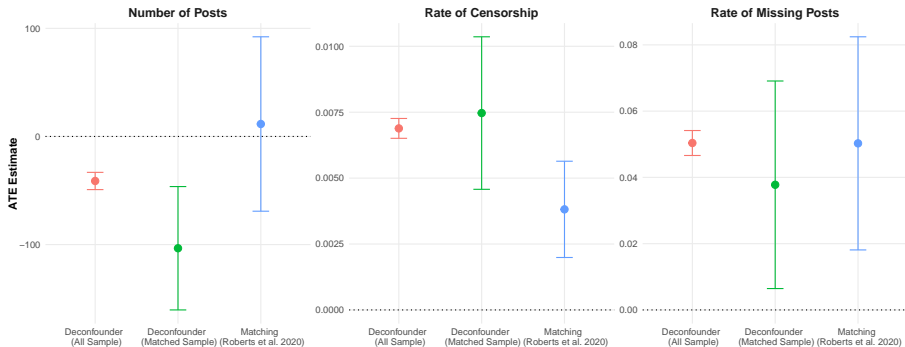
$$\hat{\tau} = \frac{1}{nK} \sum_{k=1}^K \sum_{i: I(i)=k} \hat{\mu}_1^{(-k)}(\hat{\mathbf{f}}^{(-k)}(\mathbf{R}_i), \mathbf{Z}_i) - \hat{\mu}_0^{(-k)}(\hat{\mathbf{f}}^{(-k)}(\mathbf{R}_i), \mathbf{Z}_i) \\ + \frac{T_i \{Y_i - \hat{\mu}_1^{(-k)}(\hat{\mathbf{f}}^{(-k)}(\mathbf{R}_i), \mathbf{Z}_i)\}}{\hat{\pi}^{(-k)}(\hat{\mathbf{f}}^{(-k)}(\mathbf{R}_i), \mathbf{Z}_i)} - \frac{(1 - T_i) \{Y_i - \hat{\mu}_0^{(-k)}(\hat{\mathbf{f}}^{(-k)}(\mathbf{R}_i), \mathbf{Z}_i)\}}{1 - \hat{\pi}^{(-k)}(\hat{\mathbf{f}}^{(-k)}(\mathbf{R}_i), \mathbf{Z}_i)}$$

- Double robustness, asymptotic normality

# Empirical Analysis

- Reproduced all the texts using open-source LLaMa3-8B
- Internal representation: last token of the final layer,  $\dim(\mathbf{R}) = 4080$
- Automated hyperparameter tuning via Optuna (Akiba et al. 2019)
  - $\dim(\mathbf{f}(\mathbf{R})) = 2048$
  - depth of hidden layers = 2
  - size of hidden layers after deconfounder =  $[50, 1]$
- 2-fold cross-fitting:
  - clustered standard errors at the user level
  - truncation of extreme propensity scores (Dorn, 2025)

# Empirical Results



- Our analysis shows higher rates of censorship and self-censorship
- Full sample analysis is much more efficient

## Residual Correlations with Candidate Confounder

- Confounder: proportion of 30 censorship related keywords (Fu et al. 2013)
- Extract the residuals from each method
- Compute Spearman's rank correlation with the confounder and  $p$ -value

Outcome	Proposed method		Matching
	Full sample	Matched sample	Matched sample
Number of posts	0.005 (0.330)	-0.027 (0.476)	0.022 (0.353)
Rate of censorship	-0.003 (0.421)	0.017 (0.647)	0.071 (0.005)
Rate of missing posts	-0.002 (0.653)	-0.025 (0.504)	0.043 (0.062)

# Image as Treatment: Facial Features and Election Results

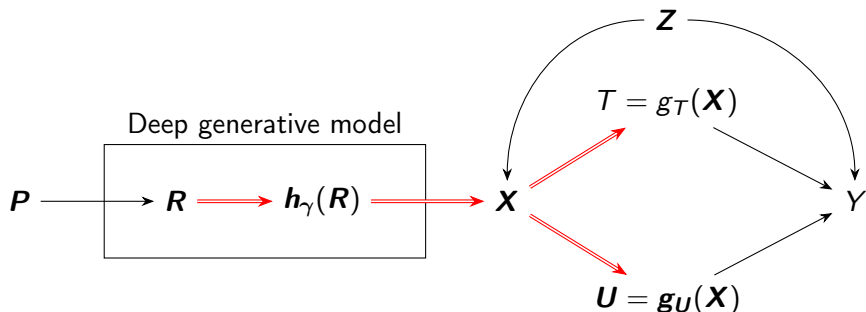
(Lindholm et al. 2024)

- How does the visual appearance of political candidate predict their electoral success?
- Data: 7,080 Danish politicians with candidate photos



- Treatment variables: facial features (continuous scores)
  - 1 attractiveness
  - 2 trustworthiness
  - 3 dominance
- Outcome: Election results (number of votes standardized via z-score)
- Structured confounding variables: age, gender, education
- We wish to adjust other facial confounding features

# Assumptions



- **Separability:**

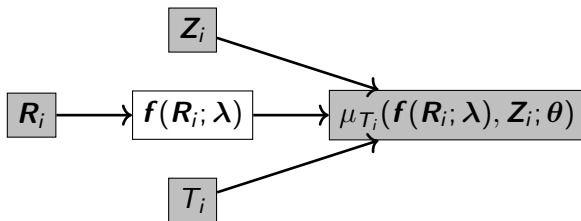
$$Y_i(\mathbf{X}_i) = Y_i(g_T(\mathbf{X}_i), g_U(\mathbf{X}_i)) = Y_i(T_i, U_i)$$

- **Lemma:** separability implies **overlap**

$$\mathbb{P}(T_i = t \mid \mathbf{U}_i = \mathbf{u}, \mathbf{Z}_i = \mathbf{z}) > 0 \quad \text{for all } \mathbf{u} \in \mathcal{U}, \mathbf{z} \in \mathcal{Z}$$

# Identification, Estimation, and Inference

- Identification
  - Existence of (possibly non-unique) deconfounder
  - Adjusting for the deconfounder yields nonparametric identification
- Estimation and inference



- 1 estimate the outcome models and deconfounder via Neural Network
  - 2 estimate the propensity score using the estimated Deconfounder
  - 3 inference via Double Machine Learning
- DragonNet (Shi et al. 2019) jointly estimates the outcome models, propensity score, and deconfounder, leading to the lack of overlap

# Empirical Analysis

- Reproduce all images using **Stable diffusion** (ver. 1.5)
- Original image:

$$\dim(\mathbf{X}) = 304(\text{width}) \times 304(\text{height}) \times 3(\text{RGB}) = 277248$$

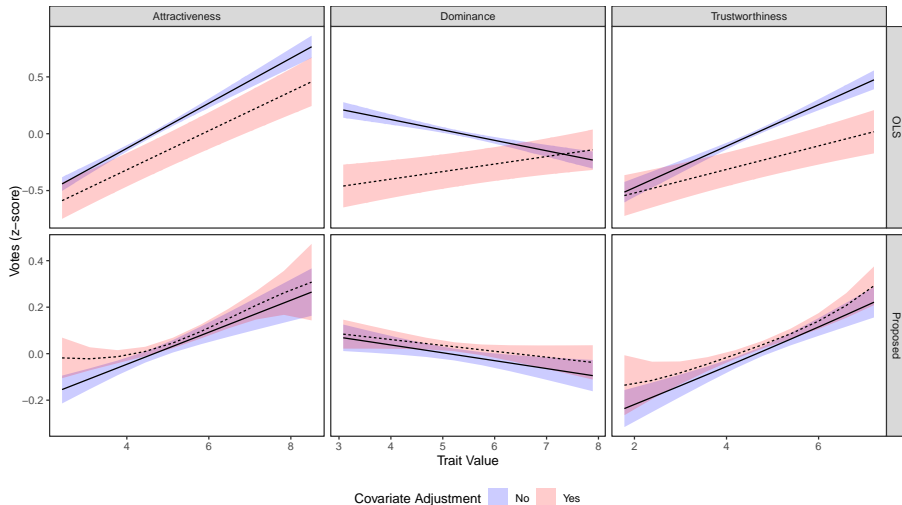
- Internal representation:  $\dim(\mathbf{R}) = 16384$
- Neural network architecture:
  - $\dim(\mathbf{f}(\mathbf{R})) = 1024$
  - depth of hidden layers = 2
  - size of hidden layers after deconfounder = [200, 1]
- Nonparametrically estimate the average effect curve

$$\xi_t := \mathbb{E}[Y_i(t, \mathbf{U}_i)]$$

- Doubly-robust pseudo-outcome approach (Kennedy et al. 2017)



# Empirical Results



- Unlike OLS, the proposed method is not sensitive to the inclusion of structured confounding variables

# Text as Treatment: Persuasion and Rhetoric

(Blumenau and Lauderdale, 2022)

- Which types of political rhetorics are most persuasive?
- Forced choice conjoint experiment with texts
- Total of 336 political arguments
  - 12 policy issues: tuition fees, fracking, etc.
  - 14 rhetorical elements: cost and benefit, morality, etc.
  - for or against
- Outcome: Persuasiveness of arguments
  - one argument is more persuasive than the other
  - equally persuasive

## Example Text Pair

- Policy topic: building a third runway at Heathrow:

### Appeal to authority / For

The Airports Commission, an independent body established to study the issue, have argued that expanding Heathrow is the most effective option to address the UK's aviation capacity challenge

### Appeal to history / Against

History show us that most large infrastructure projects do not lead to significant economic growth, which suggests that the expansion of Heathrow will fail to pay for itself

- Can we adjust for the unstructured confounding features of texts?

# The Structural Model

- The original **Bradley-Terry** type model:

$$\log \left[ \frac{\mathbb{P}(Y_{jj'(i)} \leq k)}{\mathbb{P}(Y_{jj'(i)} > k)} \right] = \delta_k + (\alpha_{P_j} S_j + \beta_{T_j} + \gamma_j) - (\alpha_{P_{j'}} S_{j'} + \beta_{T_{j'}} + \gamma_{j'})$$

where  $i$  indexes respondents,  $j$  indexes arguments,  $P_j$  denotes policy area,  $S_j$  denotes for/against, and  $T_j$  denotes rhetoric

- Our **semiparametric model**:

$$\log \left[ \frac{\mathbb{P}(Y_{j(i),j'(i)} \leq k)}{\mathbb{P}(Y_{j(i),j'(i)} > k)} \right] = \delta_k + \mu(T_j, \mathbf{U}_j) - \mu(T_{j'}, \mathbf{U}_{j'})$$

- Persuasiveness of rhetoric  $T_j = t$

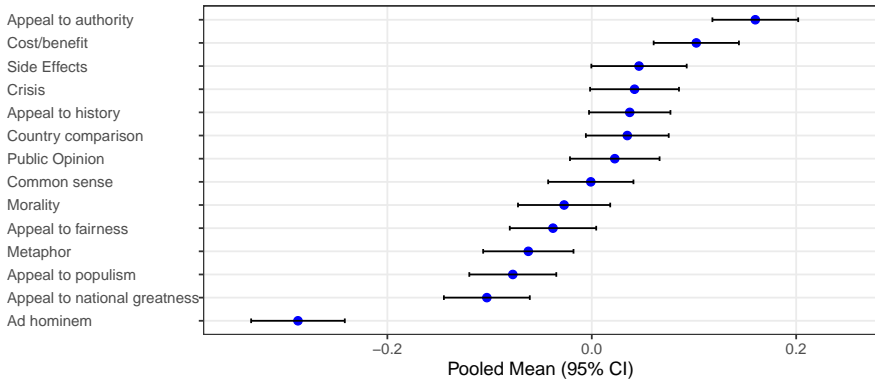
$$\beta(t) := \mathbb{E}[\mu(t, \mathbf{U}_j)]$$

- Estimate  $\beta(t)$  using the deconfounder  $\mathbf{f}(\mathbf{R}_j)$

# Empirical Analysis

- Reproduce all texts using Llama3–8B
- Internal representation: last token of the final layer,  $\dim(\mathbf{R}) = 4096$
- Neural network architecture:
  - $\dim(\mathbf{f}(\mathbf{R})) = 1024$
  - depth of hidden layers = 2
  - size of hidden layers after deconfounder =  $[200, 1]$
- Quantify uncertainty via Monte Carlo dropout (Gal and Ghahramani 2016)

Pooled Mean Estimates with Confidence Intervals



- Stronger effects for ad hominem, appeal to authority, and cost/benefit

# Concluding Remarks

- Generative AI can be used to improve causal inference
  - can generate treatments at scale
  - enables the extraction of true internal representation
  - better causal representation learning
- Open-source software **GPI** (GenAI Powered Inference) is available at <https://gpi-pack.github.io/>
- Further extensions
  - causal inference with multimodal data (e.g., videos)
  - interpretation of estimated deconfounder
  - discovery of treatment concepts
  - policy learning with unstructured treatments