# Causal Discovery with Language Models as Imperfect Experts

Stephanie Long [*1]   Alexandre Piché [*234]   Valentina Zantedeschi [*2]   Tibor Schuster [1]   Alexandre Drouin [24]

**servicenow.**     **Mila**

Causality Discussion Group - 2024

# Scope of the work

- NOT ABOUT causal reasoning of Large Language Models

- ABOUT leveraging information from related tasks
  - by querying an (imperfect) expert
  - via variables' meta-data (e.g., their name or description)
  - to reduce uncertainty in data-based causal discovery methods
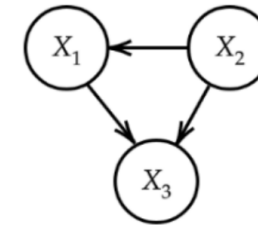
# Causal Discovery

# Markov Equivalence Class



How to reduce uncertainty?

# Causal Discovery with Expert Knowledge



Data-driven Algorithm e.g., PC

$$p(Z \rightarrow Y) = 0.9 \qquad p(Z \leftarrow Y) = 0.1$$
$$p(Z \rightarrow X) = 0.5 \qquad p(Z \leftarrow X) = 0.5$$

Expert orientations and their probability

Bayesian Inference

**Key point: We do NOT assume that the experts are perfect**

# Causal Discovery with Expert Knowledge



$$\min \quad \left| \mathcal{M}^{E,S} \right|$$

Final MEC size

$$\text{such that } p\left( G^\star \in \mathcal{M}^{E,S} \right) \geq 1 - \eta,$$

Probability orientations are correct

Tolerance to error

# Causal Discovery with Expert Knowledge

$$\min \; |\mathcal{M}^{E,S}|$$

Final MEC size

$$\text{such that} \; \boxed{p(G^\star \in \mathcal{M}^{E,S})} \geq 1 - \eta,$$

Probability orientations are correct

Tolerance to error

estimated via Bayesian inference:

P(edges are correctly oriented | we observed such expert orientations)

# Causal Discovery with Expert Knowledge

$$\min \ \left| \mathcal{M}^{E,S} \right|$$

Final MEC size

$$\text{such that } p\left( G^\star \in \mathcal{M}^{E,S} \right) \geq 1-\eta,$$

Probability orientations are correct

Tolerance to error

Hyper-parameter

# Expert Model

Assumption: Expert makes independent decisions

$$p(E_1, \ldots, E_u \mid O_1, \ldots, O_u, \ldots, O_{k+u}) = \prod_i p(E_i \mid O_i)$$

We can factorize the likelihood

True orientations $O_1 - O_2 - \ldots - O_{k+u}$

Expert orientations $E_1 \quad E_2$

# Bayesian Posterior

$$p(O_1, O_2 | E_1, \ldots, E_u) = \frac{p(E_1,\ldots,E_u | O_1,O_2)\, p(O_1,O_2)}{p(E_1,\ldots,E_u)}$$

Likelihood

Prior

posterior

Normalization constant

# Edge orientations are inter-dependent



Posterior cannot be factorized as we do for the likelihood

*Perković, Emilija, Markus Kalisch, and Maloes H. Maathuis. "Interpreting and using CPDAGs with background knowledge." arXiv preprint arXiv:1707.02171 (2017).*

# Edge orientations are inter-dependent



Prior := uniform over graphs in MEC
(we marginalize to get prior and posterior probabilities of subset of edges)

*Perković, Emilija, Markus Kalisch, and Maloes H. Maathuis. "Interpreting and using CPDAGs with background knowledge." arXiv preprint arXiv:1707.02171 (2017).*

# Considered experts

- $\epsilon$-expert: gives wrong orientation with constant probability of error

- LLM: ? we trust their confidence estimate
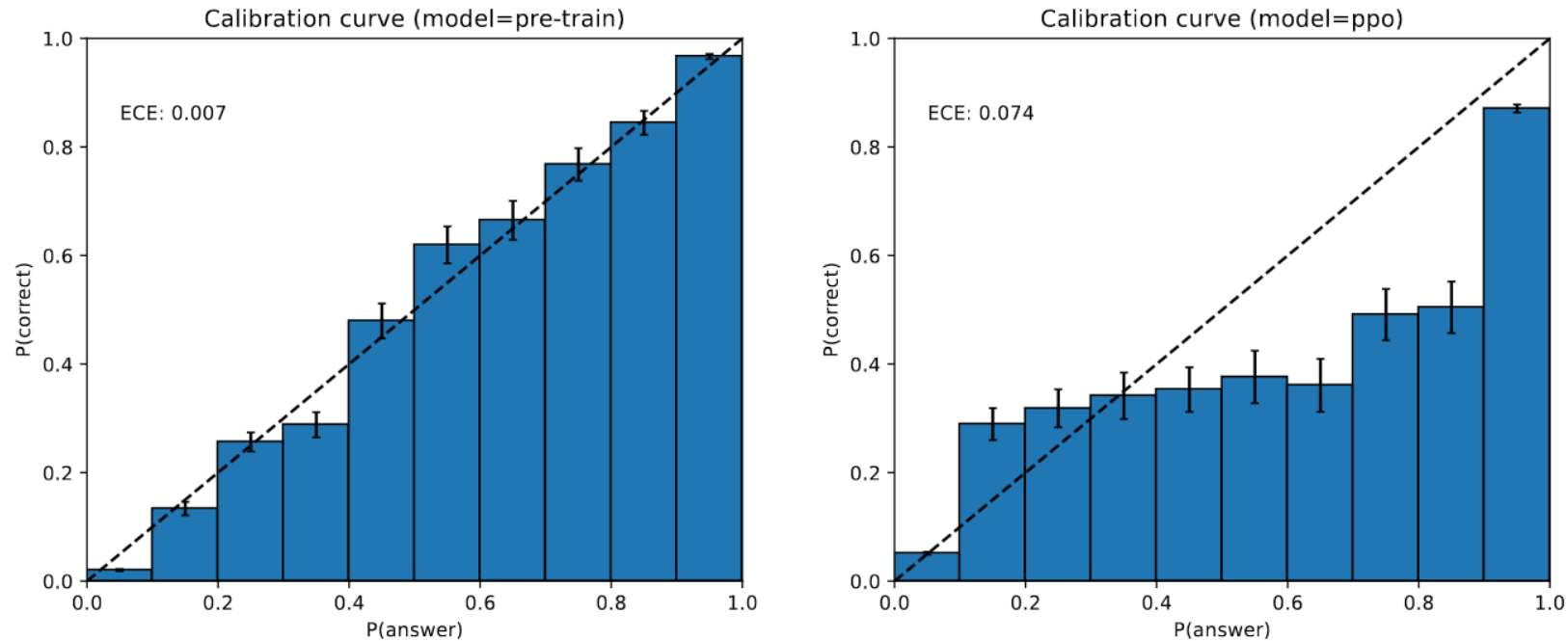
# Are LLMs calibrated?



**Figure 8.** Left: Calibration plot of the pre-trained GPT-4 model on a subset of the MMLU dataset. On the x-axis are bins according to the model's confidence (logprob) in each of the A/B/C/D choices for each question; on the y-axis is the accuracy within each bin. The dotted diagonal line represents perfect calibration. Right: Calibration plot of the post-trained GPT-4 model on the same subset of MMLU. The post-training hurts calibration significantly.

Achiam, Josh, et al. "Gpt-4 technical report." *arXiv preprint arXiv:2303.08774* (2023).
Kadavath, Saurav, et al. "Language models (mostly) know what they know." *arXiv preprint arXiv:2207.05221* (2022).

# Scoring orientations with LLMs

Among these two options which one is the most likely true:
(A) lung cancer causes cigarette smoking
(B) cigarette smoking causes lung cancer'
The answer is:

We compute likelihood of (A) and likelihood of (B)
… and normalize

# Randomizing the prompt

Among these two options which one is the most likely true:
(A) $\{\mu_i\}$ $\{\text{verb}_k\}$ $\{\mu_j\}$
(B) $\{\mu_j\}$ $\{\text{verb}_k\}$ $\{\mu_i\}$
The answer is:

# Greedy Algorithm

$$\min \; \left| \mathcal{M}^{E,S} \right|$$

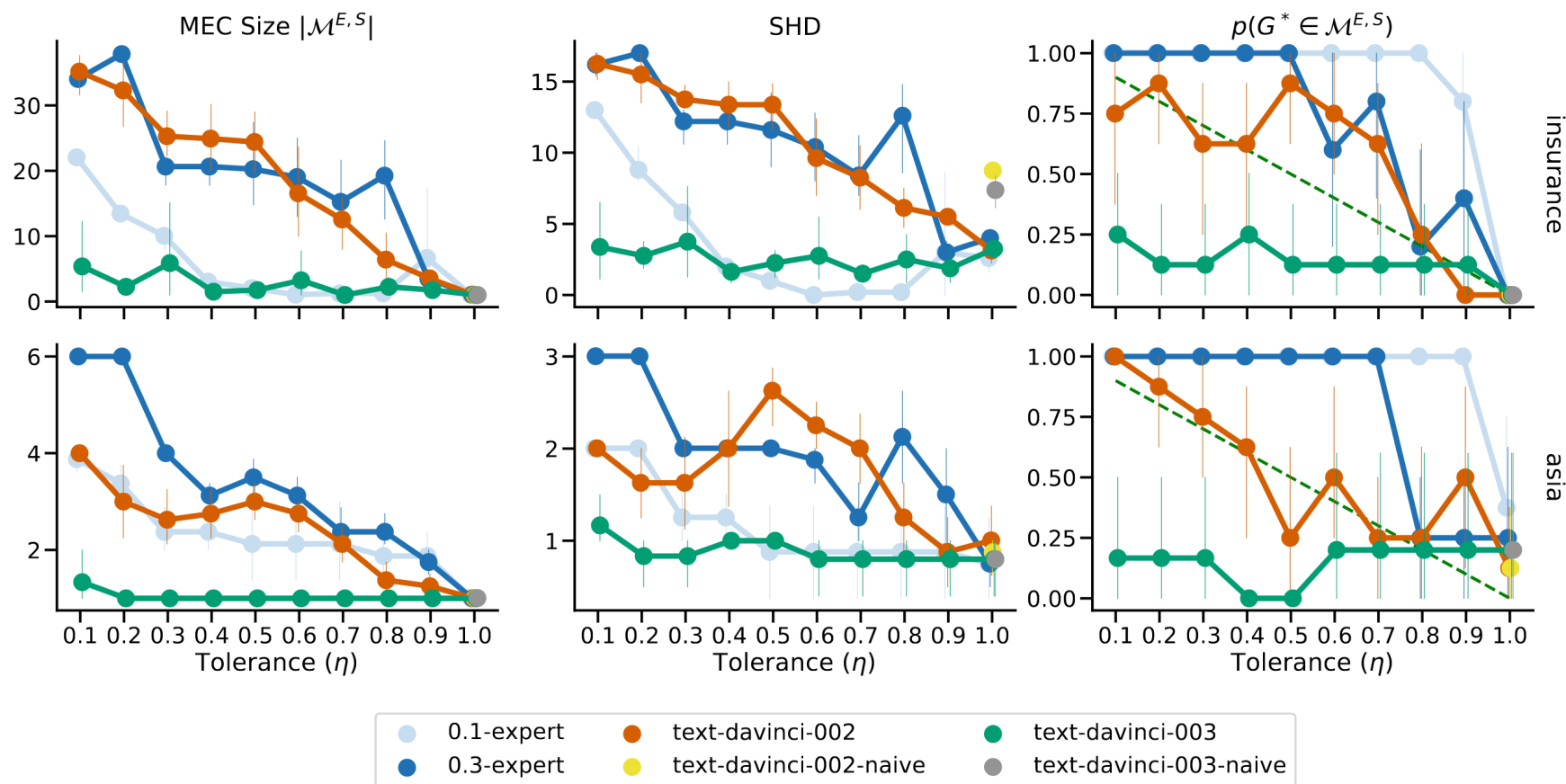$$\text{such that } p\left(G^{\star} \in \mathcal{M}^{E,S}\right) \geq 1 - \eta$$

1. Query expert on all unoriented edges $(E_1, \ldots, E_u)$
2. FOR each potential new orientation $O_i$, we compute the posterior:

$$p(O_i, O_I \mid E_1, \ldots, E_u)$$

   Where $O_I$ is the set of orientations consequential to orienting $O_i$
3. Select ( $o_i, o_I$ ) with the highest posterior
4. IF posterior of updated graph does not satisfy tolerance constraint, STOP
5. ELSE back to 2.

# Results

# Future Work

- Expert model is quite unrealistic

    How to account for systematic errors?


- Computing posterior requires enumerating all graphs in MEC

    How to scale to large number of variables?

# Thanks!

- https://arxiv.org/abs/2307.02390

- https://github.com/StephLong614/Causal-disco-LLM-imperfect-experts