

# Outliers and Hallucinations: Contributions to Robust Community Detection and Language Model Alignment

**Leonardo Martins Bianco**

*Supervisors:*

**Christine Keribin** (Laboratoire de Mathématiques d'Orsay)

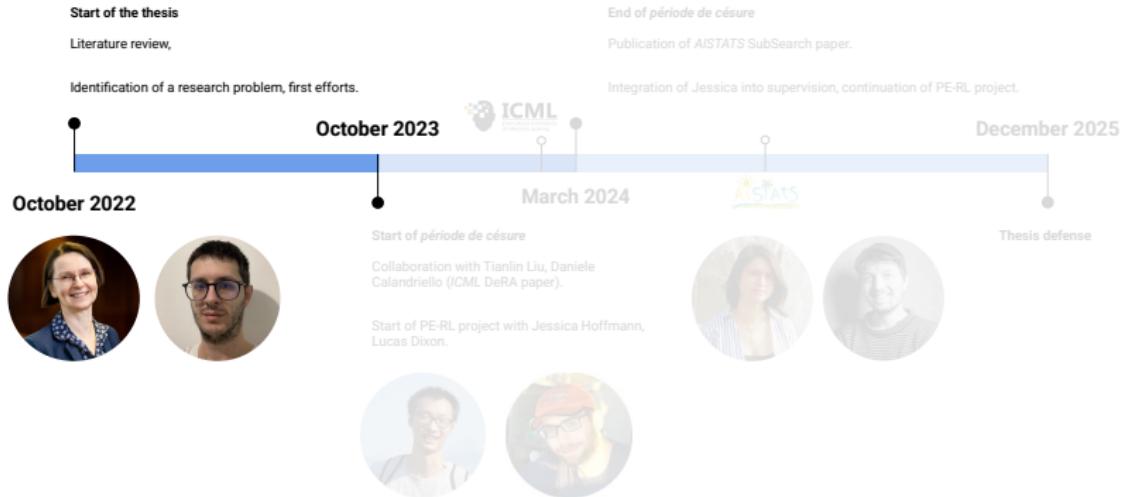
**Zacharie Naulet** (INRAE)

**Jessica Hoffmann** (Google DeepMind)

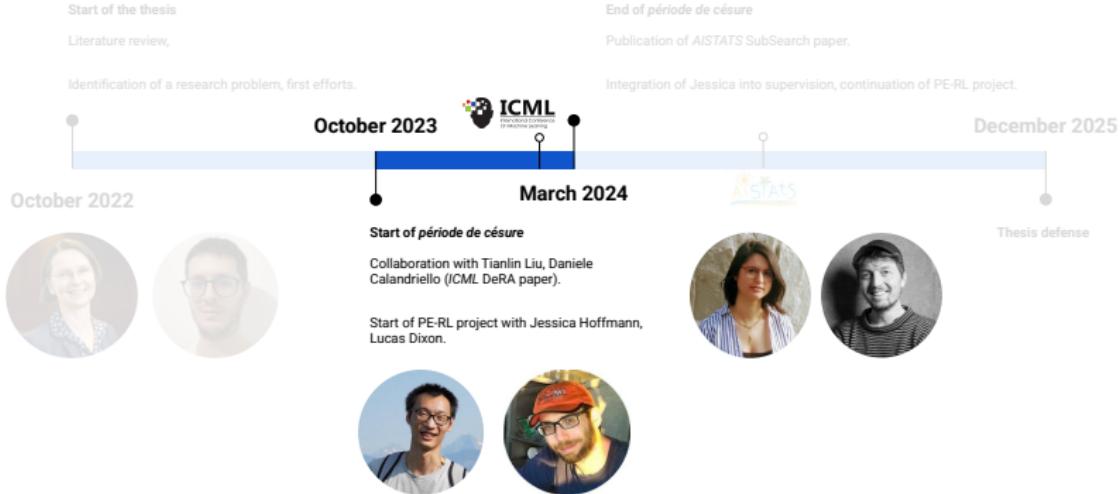
December 4, 2025

# Introduction

# Thesis progress



# Thesis progress



# Thesis progress



# Thesis progress



# Overview

## Part I: Contributions to Robust Community Detection

- ❖ Robust Estimation for the SBM

## Part II: Contributions to Language Model Alignment

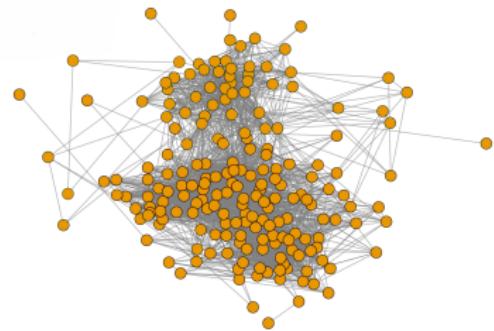
- ❖ Reducing Hallucinations with Synthetic Hallucinations
- ❖ Decoding-time Realignment of Language Models

## Part I

# Contributions to Robust Community Detection

# Motivation

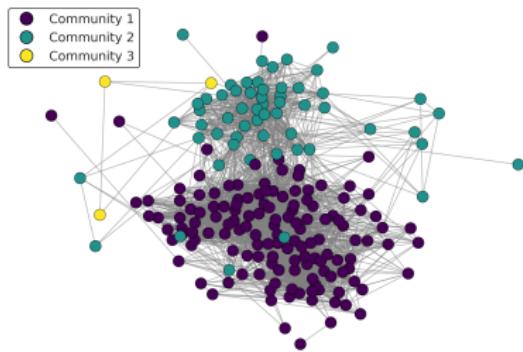
- ❖ *Community detection:* group similar nodes in a graph.
- ❖ Algorithms can be **highly sensitive** to *outlier* nodes.



Graph of Jazz collaborations in New York, Chicago, and Other Cities [[Gleiser and Danon, 2003](#)].

# Motivation

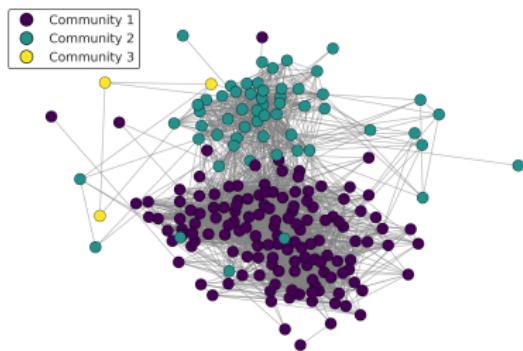
- ❖ *Community detection:* group similar nodes in a graph.
- ❖ Algorithms can be **highly sensitive** to *outlier* nodes.



Spectral clustering applied to the graph of Jazz collaborations [[Gleiser and Danon, 2003](#)].

# Motivation

- ❖ *Community detection:* group similar nodes in a graph.
- ❖ Algorithms can be **highly sensitive** to *outlier* nodes.
- ❖ **Goal:** achieve “good” results despite outliers.



Spectral clustering applied to the graph of Jazz collaborations [[Gleiser and Danon, 2003](#)].

# Graphs

*Graph:* nodes linked by edges.

$$G = (V, E)$$

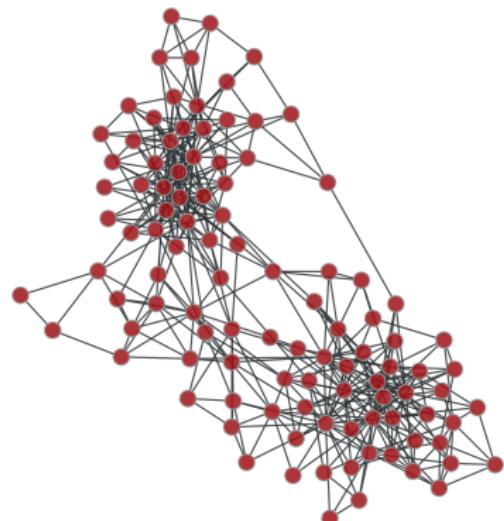
$$V = \{1, \dots, n\}, E \subset V \times V$$

*Undirected graph:*

$$(i, j) \in E \Rightarrow (j, i) \in E$$

*Adjacency matrix:*

$$A_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E \\ 0 & \text{otherwise} \end{cases}$$



# The Stochastic Block Model [Holland et al., 1983]

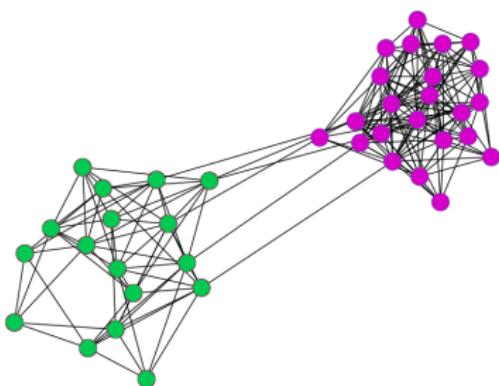
*Probabilistic model*

$Z_i \rightarrow$  community of node  $i$

$K \rightarrow$  nb. of communities

$\pi_k \rightarrow$  size of community  $k$

$\Gamma_{kl} \rightarrow$  connectivity  $k, l$



# The Stochastic Block Model [Holland et al., 1983]

*Probabilistic model*

$Z_i \rightarrow$  community of node  $i$

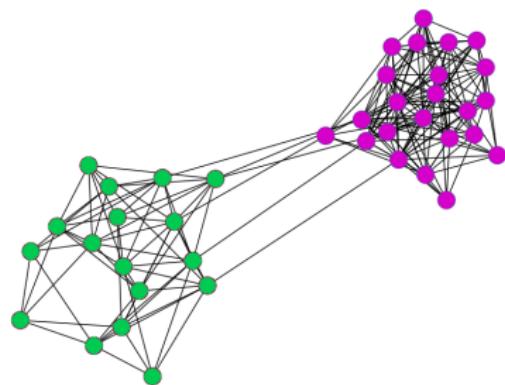
$K \rightarrow$  nb. of communities

$\pi_k \rightarrow$  size of community  $k$

$\Gamma_{kl} \rightarrow$  connectivity  $k, l$

$(Z, A) \sim \text{SBM}_K(\pi, \Gamma)$

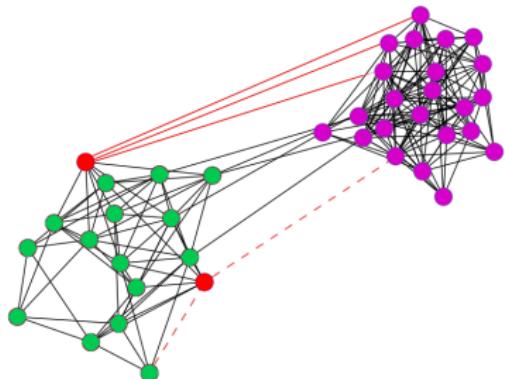
$$\begin{cases} \mathbb{P}(Z_i = k) = \pi_k \\ \mathbb{P}(A_{ij} = 1 | Z_i = k, Z_j = l) = \Gamma_{kl} \end{cases}$$



# The Corrupted Stochastic Block Model

**Adversary** introduces outliers:

1.  $(Z, A_{\text{pure}}) \sim \text{SBM}_K(\pi, \Gamma)$
2. Adversary *arbitrarily* changes edges of  $\gamma n$  nodes
3. Corrupted  $A$  is observed



# Research question

- ❖ Task: robustly estimate the connectivities  $\Gamma$ .
- ❖ For  $K = 1$ , [Acharya et al. \[2022\]](#).

**Research question:**

How to robustly estimate  $\Gamma$  for  $K > 1$ ?

# Results

- ❖ *Intuition:* find subgraph  $S$  excluding worst outliers.
- ❖ **First contribution:** extended bound from [Acharya et al. \[2022\]](#) to  $K > 1$ .

# Results

- ❖ *Intuition:* find subgraph  $S$  excluding worst outliers.
- ❖ **First contribution:** extended bound from [Acharya et al. \[2022\]](#) to  $K > 1$ .

**Theorem.** Denote  $\mathcal{I}$  the set of inlier nodes,  $S$  a subgraph clustered into  $S_1, \dots, S_K$ . Let  $\hat{\Gamma} = (\sum_{i \in S_k j \in S_l} A_{ij}) / |S_k||S_l|$  and  $\hat{Q}(S)_{ij} = \hat{\Gamma}_{S(i)S(j)}$ . Then,

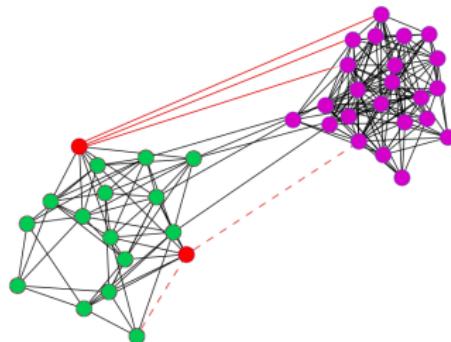
$$\|\Gamma - \hat{\Gamma}\|_1 \lesssim \frac{\|A_S - \hat{Q}(S)\|_{\text{op}}}{\min_{1 \leq k \leq K} |\Omega_k \cap S_k \cap \mathcal{I}|}$$

# Results

- ❖ *Intuition:* find subgraph  $S$  excluding worst outliers.
- ❖ **First contribution:** extended bound from [Acharya et al. \[2022\]](#) to  $K > 1$ .
- ❖ **Second contribution** (`SUBSEARCH`): finding  $S$  by optimizing  $c(S) := \|A_S - \hat{Q}(S)\|_{\text{op}}$  via Simulated Annealing.
- ❖ **Code:** [github.com/leobianco/robust\\_estim\\_sbm](https://github.com/leobianco/robust_estim_sbm)

# SUBSEARCH: Subgraph Search via Simulated Annealing

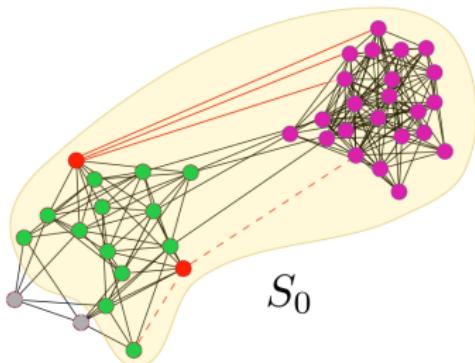
**Goal:** explore the space  $\mathcal{S}$  of subgraphs  $S \subset G$  of size  $(1 - \gamma)n$ , to minimize  $c(S) = \|A_S - \hat{Q}(S)\|_{\text{op}}$ .



# SUBSEARCH: Subgraph Search via Simulated Annealing

**Goal:** explore the space  $\mathcal{S}$  of subgraphs  $S \subset G$  of size  $(1 - \gamma)n$ , to minimize  $c(S) = \|A_S - \hat{Q}(S)\|_{\text{op}}$ .

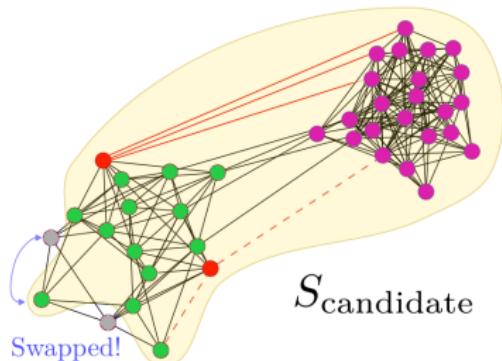
- ① Initialize at a high temperature  $T_0$ , random subgraph  $S_0$



# SUBSEARCH: Subgraph Search via Simulated Annealing

**Goal:** explore the space  $\mathcal{S}$  of subgraphs  $S \subset G$  of size  $(1 - \gamma)n$ , to minimize  $c(S) = \|A_S - \hat{Q}(S)\|_{\text{op}}$ .

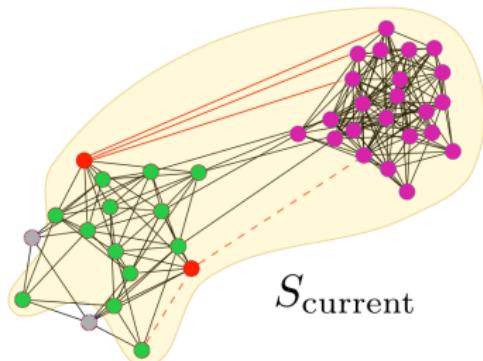
- ① **Initialize** at a high temperature  $T_0$ , random subgraph  $S_0$
- ② **Propose**  $S_{\text{candidate}}$ : swap  $i \in S_{\text{current}}$  with  $j \notin S_{\text{current}}$



# SUBSEARCH: Subgraph Search via Simulated Annealing

**Goal:** explore the space  $\mathcal{S}$  of subgraphs  $S \subset G$  of size  $(1 - \gamma)n$ , to minimize  $c(S) = \|A_S - \hat{Q}(S)\|_{\text{op}}$ .

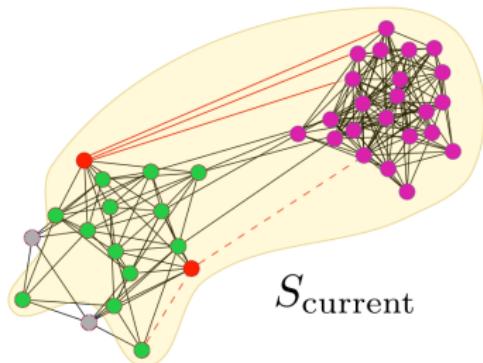
- ① **Initialize** at a high temperature  $T_0$ , random subgraph  $S_0$
- ② **Propose**  $S_{\text{candidate}}$ : swap  $i \in S_{\text{current}}$  with  $j \notin S_{\text{current}}$
- ③ **Accept or reject**: compute  $\Delta = c(S_{\text{current}}) - c(S_{\text{candidate}})$ , accept with probability  $\min(1, \exp(\Delta/T_t))$



# SUBSEARCH: Subgraph Search via Simulated Annealing

**Goal:** explore the space  $\mathcal{S}$  of subgraphs  $S \subset G$  of size  $(1 - \gamma)n$ , to minimize  $c(S) = \|A_S - \hat{Q}(S)\|_{\text{op}}$ .

- ① **Initialize** at a high temperature  $T_0$ , random subgraph  $S_0$
- ② **Propose**  $S_{\text{candidate}}$ : swap  $i \in S_{\text{current}}$  with  $j \notin S_{\text{current}}$
- ③ **Accept or reject**: compute  $\Delta = c(S_{\text{current}}) - c(S_{\text{candidate}})$ , accept with probability  $\min(1, \exp(\Delta/T_t))$
- ④ **Cool down** by  $T_{t+1} = c T_t$ ,  $c \approx 1$



# SUBSEARCH: Subgraph Search via Simulated Annealing

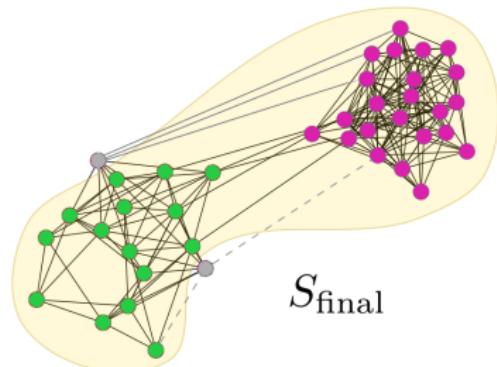
**Goal:** explore the space  $\mathcal{S}$  of subgraphs  $S \subset G$  of size  $(1 - \gamma)n$ , to minimize  $c(S) = \|A_S - \hat{Q}(S)\|_{\text{op}}$ .

① **Initialize** at a high temperature  $T_0$ , random subgraph  $S_0$

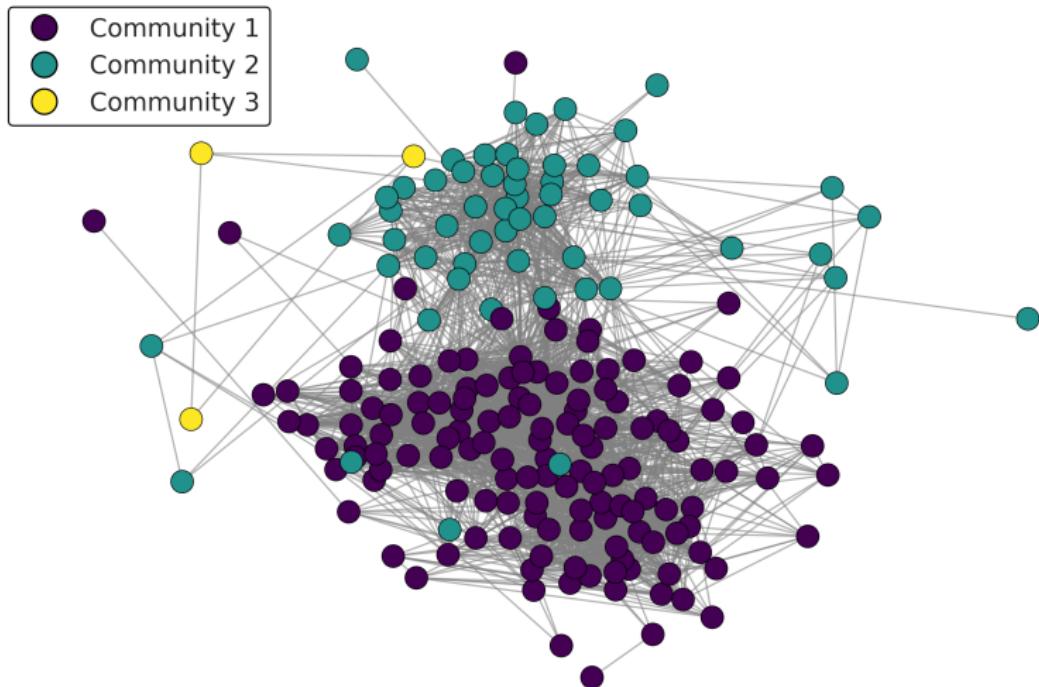
② **Propose**  $S_{\text{candidate}}$ : swap  $i \in S_{\text{current}}$  with  $j \notin S_{\text{current}}$

③ **Accept or reject:** compute  $\Delta = c(S_{\text{current}}) - c(S_{\text{candidate}})$ , accept with probability  $\min(1, \exp(\Delta/T_t))$

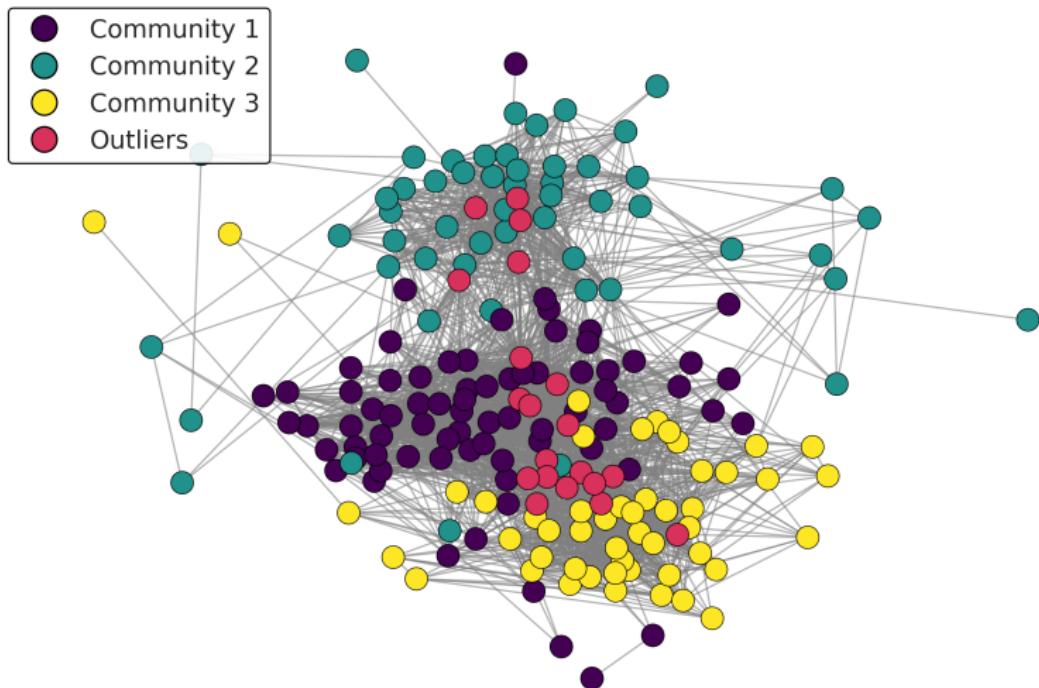
④ **Cool down** by  $T_{t+1} = c T_t$ ,  $c \approx 1$



# Results

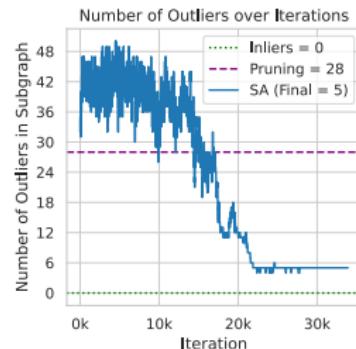
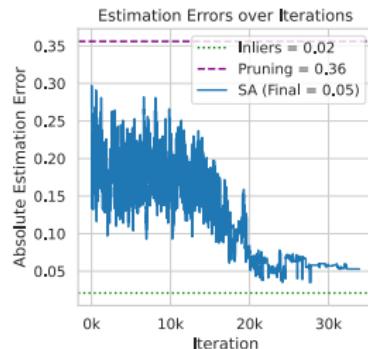
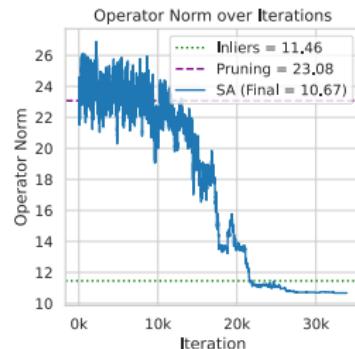


# Results

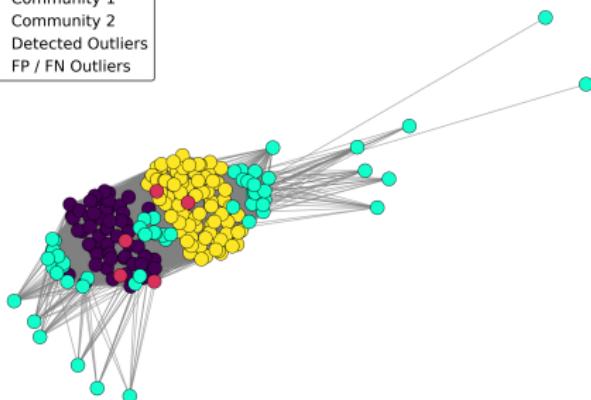


# Results

Parameters:  $n = 200$ ,  $K = 2$ ,  $\gamma = 0.3$ .

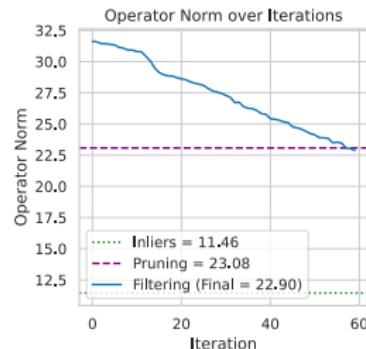


- Community 1
- Community 2
- Detected Outliers
- FP / FN Outliers

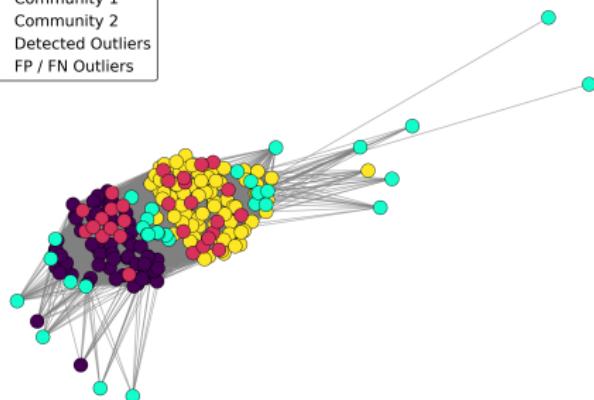


# Results

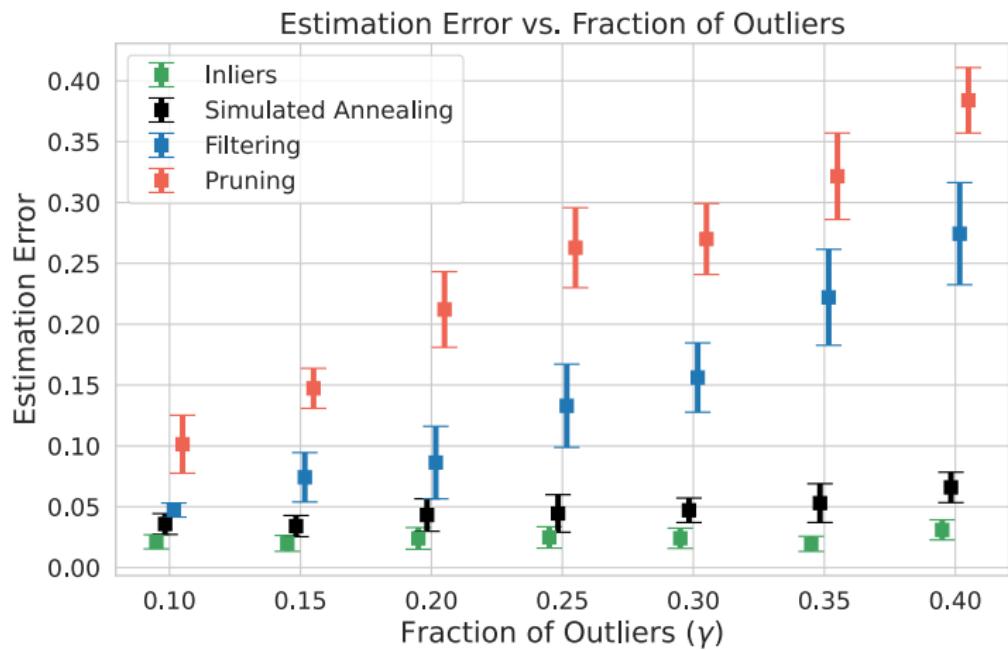
Parameters:  $n = 200$ ,  $K = 2$ ,  $\gamma = 0.3$ .



- Community 1
- Community 2
- Detected Outliers
- FP / FN Outliers



# Results



# Discussion

- ❖ **Main take away:** “exploring” the space of subgraphs  $\Rightarrow$  find subgraphs avoiding outliers.
- ❖ Limitation # 1: can we rigorously prove robustness?
- ❖ Limitation # 2: can we provide faster rates?

# Part II

# Contributions to Language Model Alignment

# Motivation

- ❖ Chatbots based on Transformers [Vaswani et al., 2017].
- ❖ Hallucinations ≈ **false information, out of topic, rambling, toxic...**
- ❖ **Goal:** less hallucinations.



A screenshot of an Ars Technica news article. The header features the site's logo and navigation links for BIZ & IT, TECH, SCIENCE, POLICY, CARS, GAMING & CULTURE, and STORE. The main headline reads "Air Canada must honor refund policy invented by airline's chatbot". Below the headline is a sub-headline: "Air Canada appears to have quietly killed its costly chatbot support." The author is listed as ASHLEY BELANGER with a timestamp of 2/16/2024, 5:12 PM. The article's thumbnail image shows an Air Canada Boeing 777 aircraft in flight against a backdrop of mountains at sunset.

# Background on Language Models

- ❖ *Vocabulary*  $\mathcal{V}$  = set of *tokens* (“pieces of words”).
- ❖ Language model

$$\pi_\theta : x = (\text{token}_1, \dots, \text{token}_L) \mapsto \pi_\theta(\cdot | x) = \text{proba. over } \mathcal{V}.$$

# Background on Language Models

- ❖ *Vocabulary*  $\mathcal{V}$  = set of *tokes* (“pieces of words”).
- ❖ Language model

$\pi_\theta : x = (\text{token}_1, \dots, \text{token}_L) \mapsto \pi_\theta(\cdot | x) = \text{proba. over } \mathcal{V}.$

- ❖ Autoregressive generation: *prompt*  $x \rightarrow \text{response } y$

$$y_1 \sim \pi_\theta(\cdot | x)$$

$$y_2 \sim \pi_\theta(\cdot | x, y_1)$$

⋮

$$y_t \sim \pi_\theta(\cdot | x, y_{<t})$$

# Background on Language Models

- ❖ Pre-training: given a dataset  $\mathcal{D}_{\text{pre}}$ , find  $\theta$  minimizing

$$\ell(\theta; \mathcal{D}_{\text{pre}}) := - \sum_{x \in \mathcal{D}_{\text{pre}}} \sum_{i=1}^{|x|} \log \pi_\theta(x_{i+1} | x_{\leq i}).$$

# Background on Language Models

- ❖ SFT: given a **task-specific** dataset  $\mathcal{D}_{\text{SFT}}$ , find  $\theta$  minimizing

$$\ell(\theta; \mathcal{D}_{\text{SFT}}) := - \sum_{x \in \mathcal{D}_{\text{SFT}}} \sum_{i=1}^{|x|} \log \pi_\theta(x_{i+1} | x_{\leq i}).$$

# Background on Language Models

- ❖ SFT: given a **task-specific** dataset  $\mathcal{D}_{\text{SFT}}$ , find  $\theta$  minimizing

$$\ell(\theta; \mathcal{D}_{\text{SFT}}) := - \sum_{x \in \mathcal{D}_{\text{SFT}}} \sum_{i=1}^{|x|} \log \pi_\theta(x_{i+1} | x_{\leq i}).$$

- ❖ Alignment: generate text with **human preferences**.
- ❖ Reinforcement learning approach:

1. Train a *reward* model  $R$  on human preference data  $\mathcal{D}_{\text{RM}}$
2. Update the *writer* model  $\pi_{\text{SFT}}$

$$\pi_\beta \in \arg \max \mathbb{E}_{g \sim \pi} [R(g)] - \beta \text{KL}(\pi \| \pi_{\text{SFT}})$$

# Background on Language Models

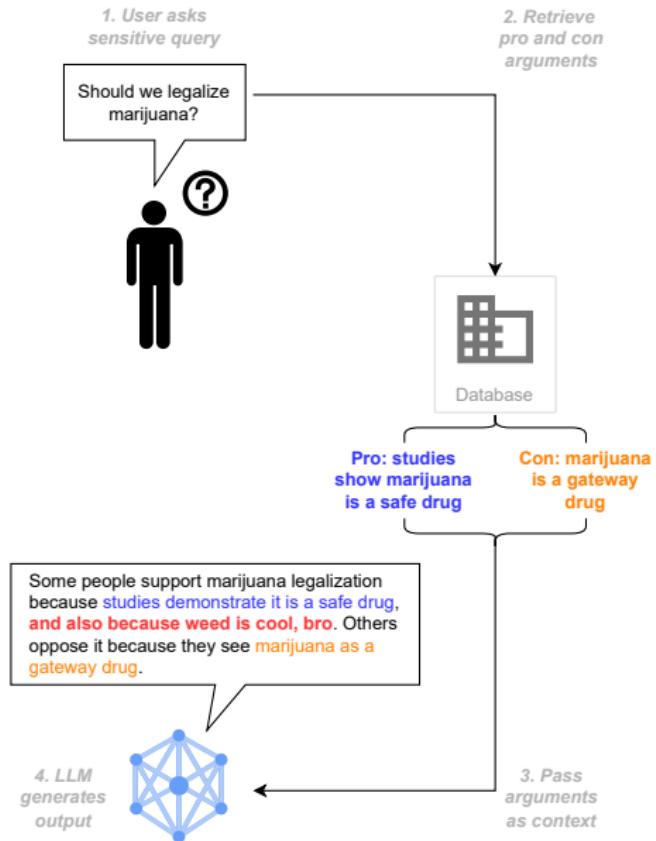
- ❖ SFT: given a **task-specific** dataset  $\mathcal{D}_{\text{SFT}}$ , find  $\theta$  minimizing

$$\ell(\theta; \mathcal{D}_{\text{SFT}}) := - \sum_{x \in \mathcal{D}_{\text{SFT}}} \sum_{i=1}^{|x|} \log \pi_\theta(x_{i+1} | x_{\leq i}).$$

- ❖ Alignment: generate text with **less hallucinations**.
- ❖ Reinforcement learning approach:
  1. Train a *reward* model  $R$  on human preference data  $\mathcal{D}_{\text{RM}}$
  2. Update the *writer* model  $\pi_{\text{SFT}}$

$$\pi_\beta \in \arg \max \mathbb{E}_{g \sim \pi} [R(g)] - \beta \text{KL}(\pi \| \pi_{\text{SFT}})$$

# Retrieval Augmented Generation: NPOV Task



## Research question #1

- ❖ Reinforcement learning approach:
  1. Train a *reward* model  $R$  on human preference data  $\mathcal{D}_{\text{RM}}$
  2. Update the *writer* model  $\pi_0$

$$\pi \in \arg \max \mathbb{E}_{g \sim \pi} [R(g)] - \beta \text{KL}(\pi \| \pi_0)$$

# Research question #1

- ❖ Reinforcement learning approach:
  1. Train a *reward* model  $R$  on human preference data  $\mathcal{D}_{\text{RM}}$
  2. Update the *writer* model  $\pi_0$

$$\pi \in \arg \max \mathbb{E}_{g \sim \pi} [R(g)] - \beta \text{KL}(\pi \| \pi_0)$$

Getting  $\mathcal{D}_{\text{RM}}$  is **costly, time-consuming, and error-prone.**

## Research question:

Can synthetic hallucinations be used instead?

- ❖ Synthetic hallucinations are *fast* and *cheap* to implement, with *automatic* and *error-free* annotations.

# Creating Synthetic Hallucinations

## Pros:

1. Studies show marijuana is a safe drug
2. Legalization boosts the economy

## Cons:

1. Marijuana is a gateway drug
2. Legalization brings costs

## *Neutral answer:*

“Some people support marijuana legalization because it would boost the economy and most studies demonstrate it is a safe drug. Others oppose it because they see marijuana as a gateway drug, and its legalization would bring many costs.”

# Creating Synthetic Hallucinations

## Pros:

1. Studies show marijuana is a safe drug
2. Legalization boosts the economy

## Cons:

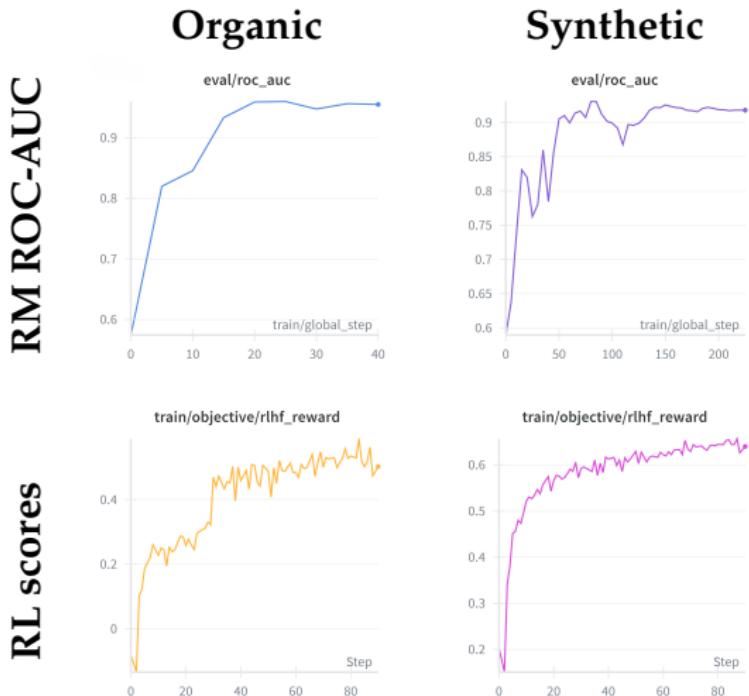
1. Marijuana is a gateway drug
2. Legalization brings costs

## *Neutral answer:*

“Some people support marijuana legalization because it would boost the economy and most studies demonstrate it is a safe drug. Others oppose it because they see marijuana as a gateway drug, and its legalization would bring many costs.”

# Results

| <i>SFT baseline (%)</i> | <i>Organic hallucinations (%)</i> | <i>Synthetic hallucinations (%)</i> |
|-------------------------|-----------------------------------|-------------------------------------|
| 10.2                    | 3.0                               | 0.74                                |



# Discussion

- ❖ **Code:** [github.com/leobianco/perl\\_hallucination](https://github.com/leobianco/perl_hallucination)

Future work:

- ❖ Other task (summarization).
- ❖ Other models (Mistral, Qwen).
- ❖ Other synthetic hallucinations schemes (LLM).

## Research question #2

- ❖ Reinforcement learning approach:
  1. Train a *reward* model  $R$  on human preference data  $\mathcal{D}_{\text{RM}}$
  2. Update the *writer* model  $\pi_0$

$$\pi_{\beta} \in \arg \max \mathbb{E}_{g \sim \pi} [R(g)] - \beta \text{KL}(\pi \| \pi_{\text{SFT}})$$

## Research question #2

- ❖ Reinforcement learning approach:
  1. Train a *reward* model  $R$  on human preference data  $\mathcal{D}_{\text{RM}}$
  2. Update the *writer* model  $\pi_0$

$$\pi_{\beta} \in \arg \max \mathbb{E}_{g \sim \pi} [R(g)] - \beta \text{KL}(\pi \| \pi_{\text{SFT}})$$

The hyperparameter  $\beta$  is expensive to tune via grid-search.

**Research question:**

Can we adjust regularization strength without retraining?

# Results

- ❖ **Contribution:** approximate realigned model at  $\beta/\lambda$

$$\hat{\pi}_{\beta/\lambda}(\cdot | x, y_{<t}) := \text{softmax} \left[ \lambda h_{\beta}^{(t)} + (1 - \lambda) h_{\text{SFT}}^{(t)} \right]$$

where  $h_{\text{SFT}}^{(t)}$  and  $h_{\beta}^{(t)}$  are the logits

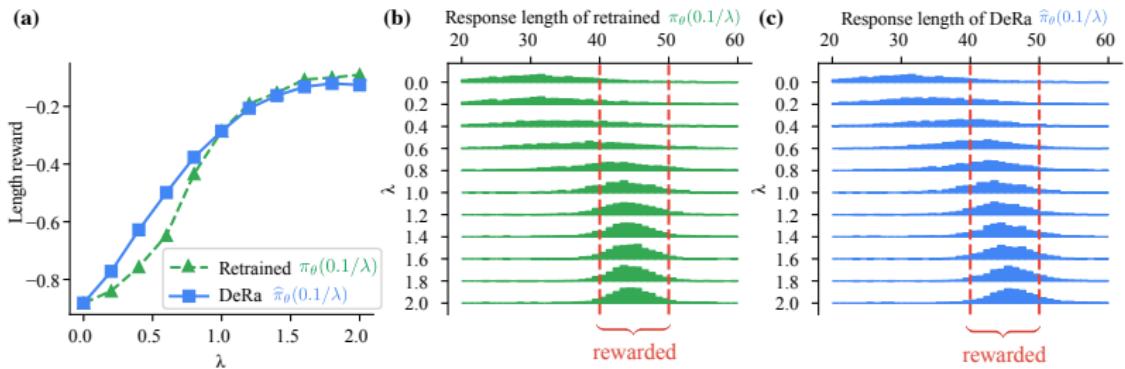
$$\begin{cases} \pi_{\text{SFT}}(\cdot | x, y_{<t}) &= \text{softmax}(h_{\text{SFT}}^{(t)}) \\ \pi_{\beta}(\cdot | x, y_{<t}) &= \text{softmax}(h_{\beta}^{(t)}) \end{cases}$$

- ❖ **Code:** <https://github.com/liutianlin0121/decoding-time-realignment>

# Results

|                             |  |
|-----------------------------|--|
| Arguments                   | pro: Denying student loan debtors the benefits of bankruptcy--benefits that all other debtors have access to--is unfair.<br>con: Discharging student loan debt would only be a temporary bandage for the much larger problem of inflated college costs.  |
| Response<br>$\lambda=0.011$ | [...] it is unfair to deny student loan debtors the benefits of bankruptcy--benefits that all other debtors have access to. They also argue that student loan debt has a disproportionately negative impact on low-income borrowers [...]. However, [...] it would be a costly and ineffective solution to the problem of student debt. They also argue that it would only be a temporary bandage for the much larger problem of inflated college costs. |
| Response<br>$\lambda=2$     | [...] One argument in favor of forgiveness is that it is unfair to deny student loan debtors the benefits of bankruptcy--benefits that all other debtors have access to. However, one argument against forgiveness is that it would only be a temporary bandage for the much larger problem of inflated college costs.   |
| Response:<br>$\lambda=5$    | Denying student loan debtors the benefits of bankruptcy--benefits that all other debtors have access to--is unfair. 301\\n text20Discharging student loan debt would only be a temporary bandage for the much larger problem of inflated college costs.  |

# Results



# Discussion

❖ ?

# References I

- Jayadev Acharya, Ayush Jain, Gautam Kamath, Ananda Theertha Suresh, and Huanyu Zhang. Robust estimation for random graphs. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 130–166. PMLR, 02–05 Jul 2022. URL <https://proceedings.mlr.press/v178/acharya22a.html>.
- Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high dimensions without the computational intractability, 2019.
- Pablo M Gleiser and Leon Danon. Community structure in jazz. *Advances in complex systems*, 6(04):565–573, 2003.
- Paul W. Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social Networks*, 5(2):109–137, 1983. ISSN 0378-8733. doi: [https://doi.org/10.1016/0378-8733\(83\)90021-7](https://doi.org/10.1016/0378-8733(83)90021-7). URL <https://www.sciencedirect.com/science/article/pii/0378873383900217>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the Conference on Neural Information Processing Systems*, volume 30, 2017.

# Example of Hallucination Removed #1

**User query:** Should Parents or Other Adults Be Able to Ban Books from Schools and Libraries?

## Pro arguments:

- Keeping books with inappropriate content out of libraries protects kids but doesn't stop people from reading those books or prevent authors from writing them.
- Parents have the right to decide what material their children are exposed to and when.

## Con arguments:

- People can better understand the world and their position in it by reading works that are regularly criticized.
- Books are a portal to different life experiences and reading encourages empathy and social-emotional development.

## Neutral point-of-view answer:

**SFT version:** This is a very controversial topic. Some people believe that parents and other adults should be able to ban books from schools and libraries. They say that this would protect children from inappropriate content. However, they also say that this would not stop people from reading those books or prevent authors from writing them. On the other hand, **others say that banning books is a way to protect children from the world** and that it is a way to encourage empathy and social-emotional development. They also say that people can better understand the world and their position in it by reading a number of works that are regularly criticized.

# Example of Hallucination Removed #1

**User query:** Should Parents or Other Adults Be Able to Ban Books from Schools and Libraries?

## Pro arguments:

- Keeping books with inappropriate content out of libraries protects kids but doesn't stop people from reading those books or prevent authors from writing them.
- Parents have the right to decide what material their children are exposed to and when.

## Con arguments:

- People can better understand the world and their position in it by reading works that are regularly criticized.
- Books are a portal to different life experiences and reading encourages empathy and social-emotional development.

## Neutral point-of-view answer:

**PE-RL version:** I don't know what to make of this. Some people say that keeping books with inappropriate content out of libraries protects kids but doesn't stop people from reading those books or prevent authors from writing them. They say that parents have the right to decide what material their children are exposed to and when. Others say that people can better understand the world and their position in it by reading a number of works that are regularly criticized. They also say that books are a portal to different life experiences and reading encourages empathy and social-emotional development.