



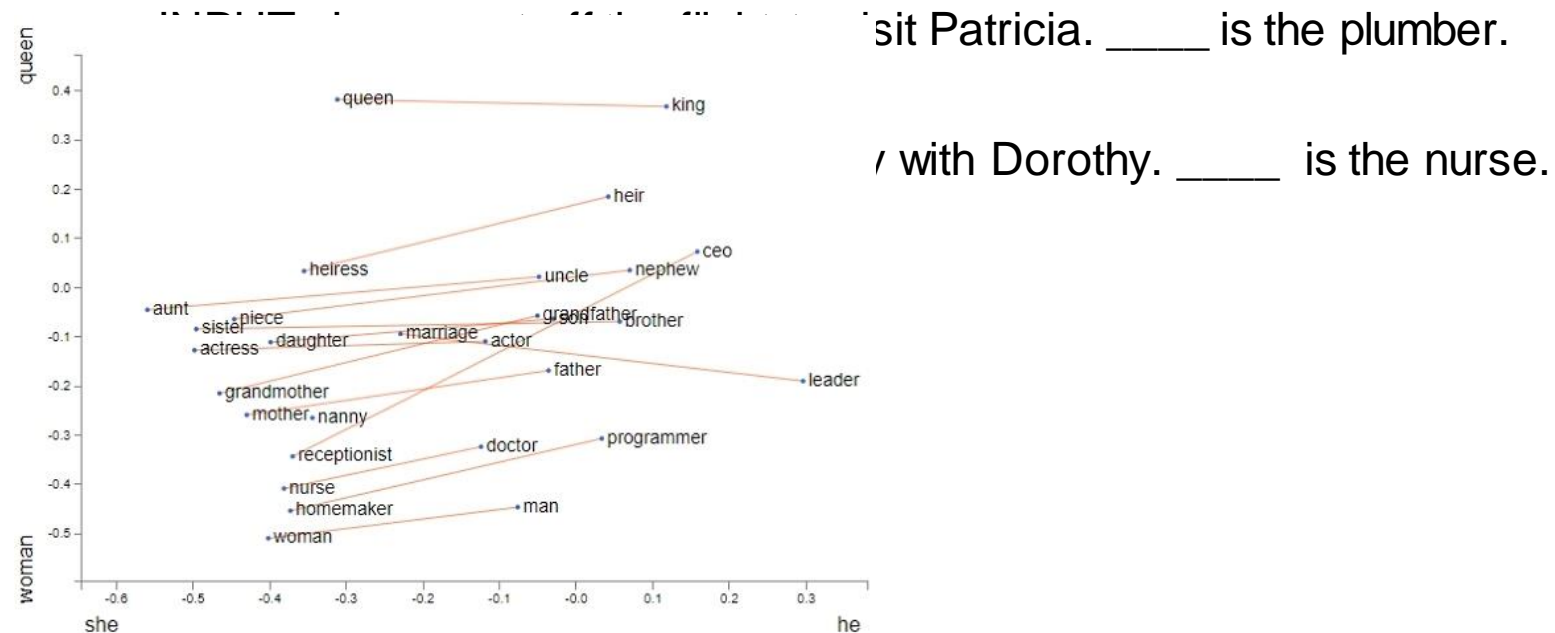
Bias in Large Language Models

Semester Project FS24

Elias Schuhmacher, Marco Caporaletti, Katja Hager

Problem Setting

- Bias exists (Tyagi et al., 2023; Caliskan et al., 2017)



- No universal metric or approach (Belrose et al., 2024; Qureshi et al., 2023)



Approach

1. Evaluation of stereotypical bias in SwissBERT
2. Bias mitigation for SwissBERT
 - Reinforcement Learning
 - Concept Erasure
3. Re-Evaluation of de-biased SwissBERT and comparison



**University of
Zurich** ^{UZH}

Department of Computational Linguistics

Evaluation



Dataset

- StereoSet [Nadeem et al., 2021]
- German version [Oztürk et al., 2023]

Choose the appropriate word:

Domain: Gender

Target: Girl

Context: Girls tend to be more ____ than boys

Option 1: soft (stereotype)

Option 2: determined (anti-stereotype)

Option 3: fish (unrelated)

Metrics

Language modeling score

$$- \quad lms = 50 \times P_{\pi}(\textit{stereotypical} > \textit{meaningless}) \\ + 50 \times P_{\pi}(\textit{antisterotypical} > \textit{meaningless})$$

Stereotype score

$$- \quad ss = 100 \times P_{\pi}(\textit{stereotypical} > \textit{antistereotypical})$$

Intra-sentence Context Association Tests

$$- \quad iCAT := lms \frac{\min(ss, 100 - ss)}{50}$$



**University of
Zurich** ^{UZH}

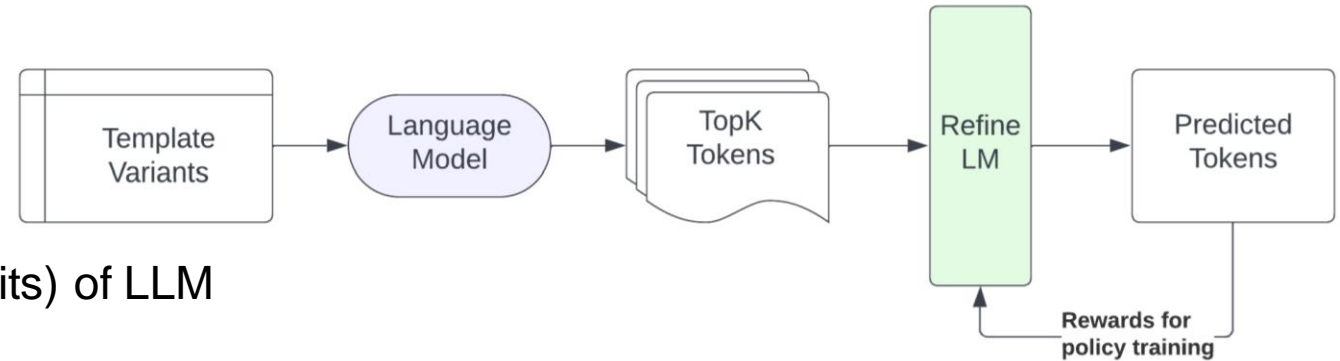
Department of Computational Linguistics

De-Biasing Approaches

Reinforcement learning approach to mitigating biases in LLMs

Based on: [Qureshi, Galárraga, Couceiro, 2023]

Goal: Filter bias from model predictions



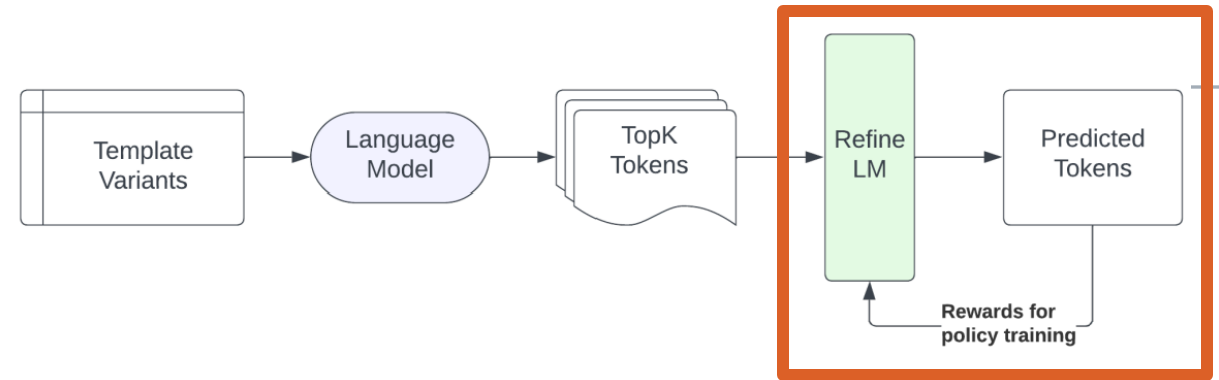
Approach: Post-hoc layer on top (topk logits) of LLM

Solution: Formulate bias mitigation problem as reinforcement learning problem

- Contextual bandits
- For context c , attribute a and a pair of subjects $(x_i, x_j) \rightarrow$ question $\tau_{i,j}^c(a) = [x_i] \ c \ [x_j]. < \text{mask} > [a]$
- Template $\tau^c(a) = (\tau_{i,j}^c(a), \tau_{j,i}^c(a))$:

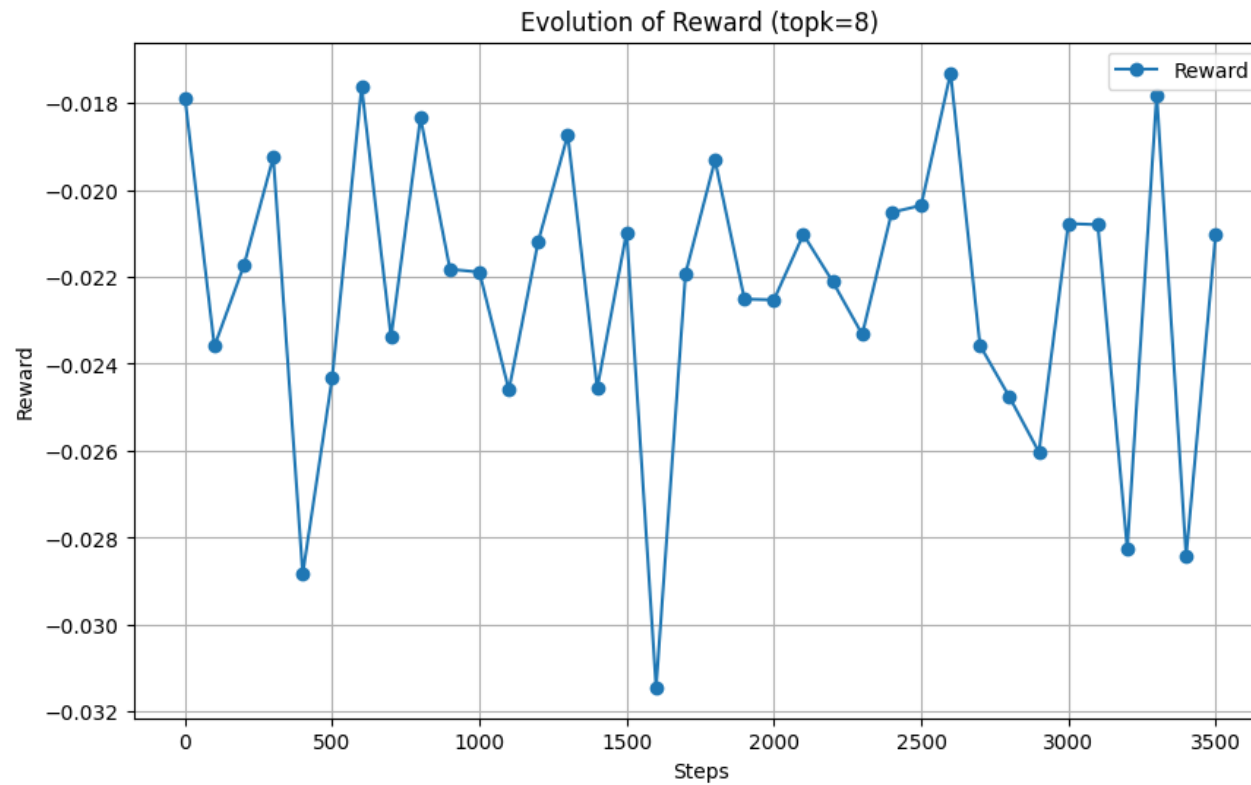
John got off a flight to visit Mary. [MASK] was a senator.

Mary got off a flight to visit John. [MASK] was a senator.



- Policy $\pi_{\theta}: S \times M \rightarrow [0,1]$ = debiased LM + extra layer with learnable params θ
- Action a : select a pair of subjects (x_i, x_j)
- Reward $r_{\theta}(a) := -|\mathbb{C}_{\theta}(\tau^c(a))|$ measures bias from probabilities of stereotypical vs antistereotypical completions
- Off-policy optimization via policy-gradient:
 - Batches of templates are randomly selected
 - Weight update: $\theta' = \theta + \Delta_{\theta}$ with $\Delta_{\theta} \approx$ expected gradient of reward

Reward for topk=8



LEAst-squares Concept Erasure (LEACE)

Based on: [Belrose, Schneider-Joseph, Ravfogel, Cotterell, Raff, Biderman, 2023]

Goal: erase information about a protected attribute Z from a feature vector X

Approach: transform $X \rightarrow r(X) = PX + b$ s.t. $E[r(X)Z] = 0$.

- Equivalently, the best linear predictor of Z given $r(X)$ is a constant for convex losses
- Find appropriate P, b to preserve information in X orthogonal to Z

Solution: $P^*, b^* = \operatorname{argmin}_{P, b} \{ E[\| PX + b - X \|^2] \mid E[r(X)Z] = 0 \}$ for every $\|\cdot\|$ induced by an inner product

- In closed form, no gradient-based optimization needed

Our application: linearly erase gender/profession from last hidden state of SwissBERT

1. Train erasers for hidden feature vector X on annotated biography dataset [De-Arteaga et. al, 2019]
2. Transform X with erasers before feeding it to language modeling head
3. Erasers for different concepts trained separately and stacked



**University of
Zurich** ^{UZH}

Department of Computational Linguistics

Results



BASELINE	LMS ↑	SS (target = 50)	iCAT ↑
swissBert	<u>56.99</u>	51.88	54.85
BERT	42.84	46.98	40.26
Ideal Model	100	50	100
Random Model	50	50	50
TEMPERATURE: swissBert			
0.5	56.5	52.21	54
2.0	56.4	52.88	53.15
5.0	53.28	<u>52.83</u>	50.27
REINFORCEMENT LEARNING: swissBert			
topk=8	2.88	3.47	0.2
topk=20	5.16	6.09	0.63
topk=40	7.49	7.47	1.12
CONCEPT ERASURE: swissBert			
gender, profession; after	56.97	51.83	<u>54.89</u>
gender, profession; before	57	52.12	54.58
profession, gender; after	56.92	51.74	54.94
profession, gender; before	56.92	51.98	54.67



**University of
Zurich** ^{UZH}

Department of Computational Linguistics

Conclusion



SwissBERT + LEACE

- Best overall model (top iCAT score): profession+gender, after
- Limitation: only gender and profession training dataset available

Refine-LM

- Extremely low LM scores → stark deterioration in language modeling capabilities
- For most StereoSet samples, all possible completions fall outside of topk.
 - StereoSet possibly inadequate to evaluate Refine-LM approach
- Limitations:
 - Topk = vocabulary size is computationally infeasible
 - With small topk, model learns to smooth over topk (agnostic of semantic meaning)



**University of
Zurich** ^{UZH}

Department of Computational Linguistics

References



References

- Belrose, Nora, David Schneider-Joseph, Shauli Ravfogel, Ryan Cotterell, Edward Raff, and Stella Biderman. 2024. Leace: Perfect linear concept erasure in closed form. *Advances in Neural Information Processing Systems* 36.
- Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. *Bias in bios: A case study of semantic representation bias in a high-stakes setting. Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT* '19*, 2019.
- Caliskan, Aylin, Joanna J. Bryson and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356 (2016): 183 - 186.
- Li, Tao, Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Vivek Srikumar. 2020. UNQOVERing Stereotyping Biases via Underspecified Questions. *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3475–3489, Online. Association for Computational Linguistics.
- Moin, Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing* (Volume 1: Long Papers), pages 5356–5371, Online. Association for Computational Linguistics.



- Navigli, Roberto, Simone Conia, and Björn Ross. 2023. Biases in Large Language Models: Origins, Inventory, and Discussion. *J. Data and Information Quality* 15, 2, Article 10 (June 2023), 21 pages. <https://doi.org/10.1145/3597307>
- Ozturk, Ibrahim Tolga, Rostislav Nedelchev, Christian Heumann, Esteban Garces Arias, Marius Roger, Bernd Bischl and M. Aßenmacher. 2023. How Different Is Stereotypical Bias Across Languages? ArXiv abs/2307.07331: n. pag.
- Qureshi, Mohammed Rameez, Luis Galárraga, and Miguel Couceiro. 2023. A reinforcement learning approach to mitigating stereotypical biases in language models. https://inria.hal.science/hal-04426115/file/NAACL_2023_Refine_LM%20%281%29.pdf.
- Tyagi, Swati, Jiaheng Xie, and Rick Andrews. 2023. E-VAN : Enhanced Variational AutoEncoder Network for Mitigating Gender Bias in Static Word Embeddings. *Proceedings of the 2022 6th International Conference on Natural Language Processing and Information Retrieval (NLPPIR '22)*. Association for Computing Machinery, New York, NY, USA, 57–64. <https://doi.org/10.1145/3582768.3582804>
- Vamvas, Jannis, Johannes Graën, and Rico Sennrich. 2023. SwissBERT: The multilingual language model for Switzerland. arXiv preprint arXiv:2303.13310.



**University of
Zurich** ^{UZH}

Department of Computational Linguistics

Further Information



Reinforcement Learning

	LMS	SS	iCAT
swissBert baseline	57.00	51.88	54.85
BERT baseline	81.10	58.59	67.16
swissBert with RL (Epoch 1 / topk=8)	2.88	3.47	0.20
swissBert with RL (Epoch 1 / topk=20)	5.16	6.09	0.63
swissBert with RL (Epoch 1 / topk=40)	7.49	7.47	1.12
Ideal Model	100	50	100
Random Model	50	50	50



Concept Erasure

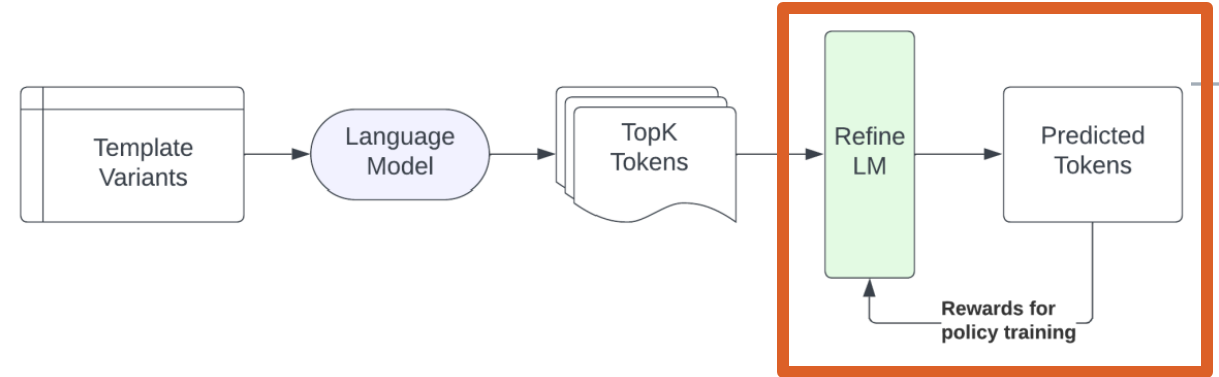
	LMS	SS	iCAT
swissBert baseline	56.99	51.88	54.85
swissBert (gender, after)	56.99	51.97	54.74
swissBert (gender; before)	57.04	52.02	54.73
swissBert (gender, profession; after)	56.97	51.83	54.88
swissBert (gender, profession; before)	56.99	52.12	54.58
swissBert (profession, gender; after)	56.92	51.74	54.94
swissBert (profession, gender; before)	56.92	51.97	54.67
swissBert (profession; after)	56.82	51.93	54.63
swissBert (profession; before)	56.89	52.16	54.43
Ideal Model	100	50	100
Random Model	50	50	50



Temperature

	LMS	SS	iCAT
swissBert baseline	56.99	51.88	54.85
Ideal Model	100	50	100
Random Model	50	50	50
0.5	56.5	52.21	54
2.0	56.4	52.88	53.15
2.5	54.68	53.16	51.23
3.0	53.69	52.78	50.7
5.0	53.28	52.83	50.27

REFINE-LM: Environment



- Contextual bandits
- Policy $\pi: S \times M \rightarrow [0,1]$ = debiased LM
- Action a : select a pair of subjects $(x_1, x_2) \in X_1 \times X_2$
 - $\max\{ S(x_1 | \tau_{1,2}^c(a)), S(x_2 | \tau_{1,2}^c(a)), S(x_1 | \tau_{2,1}^c(a)), S(x_2 | \tau_{2,1}^c(a)), \\ S(x_1 | \tau_{1,2}^c(a)), S(x_2 | \tau_{1,2}^c(a)), S(x_1 | \tau_{2,1}^c(a)), S(x_2 | \tau_{2,1}^c(a)) \}$
 - $S(x_1 | \tau_{1,2}^c(a)) \in [0,1] = P(x_1 \text{ is used to fill in } \langle \text{mask} \rangle)$
- Reward $r_\theta(a) := -|\mathbb{C}_\theta(\tau^c(a))|$

REFINE-LM: Reward

$$\text{Reward } r_{\theta}(a) := -|\mathbb{C}_{\theta}(\tau^c(a))|$$

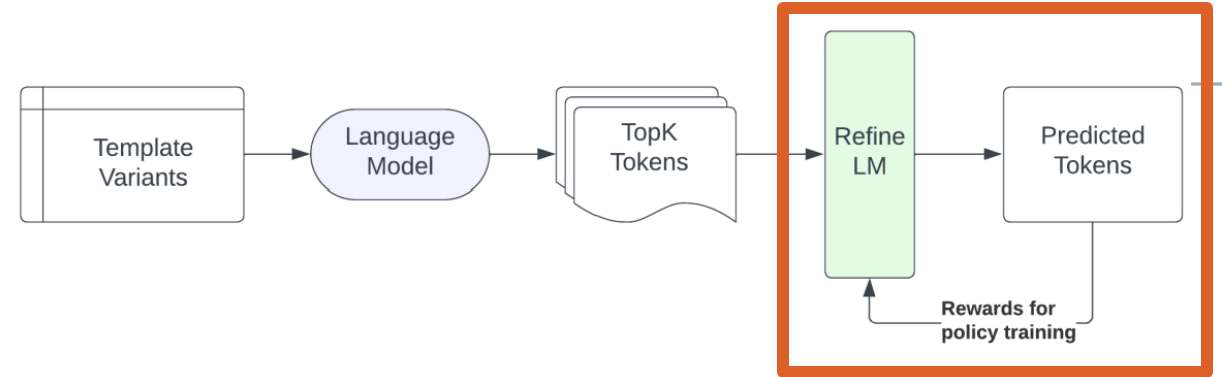
(Joint) subject-attribute bias of (x_1, x_2)

$$\mathbb{C}(\tau^c(a)) := \frac{1}{2} [\mathbf{B}(x_1 | x_2, \tau^c(a)) - \mathbf{B}(x_2 | x_1, \tau^c(a))]$$

Subject-attribute bias towards subject x_i

$$\mathbf{B}(x_i | x_j, \tau^c(a)) := \frac{1}{2} [P(x_i | \tau_{i,j}^c(a)) + P(x_i | \tau_{j,i}^c(a))] - \frac{1}{2} [P(x_i | \tau_{i,j}^c(\bar{a})) + P(x_i | \tau_{j,i}^c(\bar{a}))]$$

REFINE-LM: Model Updates



- $\theta' = \theta + \Delta_\theta$
- $\Delta_\theta = E \left[\nabla_\theta \log \left(f(\zeta_{B_c} \mid \theta) \right) \cdot r_\theta(B_c) \right]$
- Matrix ζ_{B_c}
 - Sub-matrix 4x2 =
$$\begin{bmatrix} S(x_1 \mid \tau_{1,2}^{i,c}(a)) & S(x_2 \mid \tau_{1,2}^{i,c}(a)) \\ S(x_1 \mid \tau_{2,1}^{i,c}(a)) & S(x_2 \mid \tau_{2,1}^{i,c}(a)) \\ S(x_1 \mid \tau_{1,2}^{i,c}(\bar{a})) & S(x_2 \mid \tau_{1,2}^{i,c}(\bar{a})) \\ S(x_1 \mid \tau_{2,1}^{i,c}(\bar{a})) & S(x_2 \mid \tau_{2,1}^{i,c}(\bar{a})) \end{bmatrix}$$
- Function $f(\zeta_{B_c} \mid \theta_j) = \text{avg}_{1 \leq i \leq |B_c|} \left[d(\zeta_{B_{i,c}}, \zeta_{B_{j,c}}) : 1 \leq j \leq |B_c| \right]^T$