# Leveraging Domain Knowledge for Inclusive and Bias-aware Humanitarian Response Entry Classification – Appendices

## A   Hyperparameters

Table 1 reports the results of hyperparameter tunning.

## B   Additional Results

Figures 1-5 report the results of gender and country bias measurement over various backbone LLMs, and architectures, before after applying CDA bias mitigation.

| Hyperparameters | Values |
| --- | --- |
| Number of Epochs | 3 |
| Initial Learning Rate | 1e-4 |
| Dropout Rate | 0.2 |
| Train Batch Size | 8 |
| Validation Batch Size | 16 |
| Optimizer | Adam Weight (with the standard Pytorch (https://pytorch.org/) hyperparameters) |
| Learning Rate Scheduler | Pytorch StepLR (with decay=0.4, step size=1) |
| LLM input text max length | 200 |
| Freezed LLM layers | LLM Embedding and first LLM layer |
| Decision boundary threshold | Finetuned differently for each training setup and tag on the best F1 score validation set after training (from 20 values ranging from the minimum to the maximum probability predicted for each tag). |

Table 1: Hyperparameters used for finetuning Classification models and for Generating final predictions
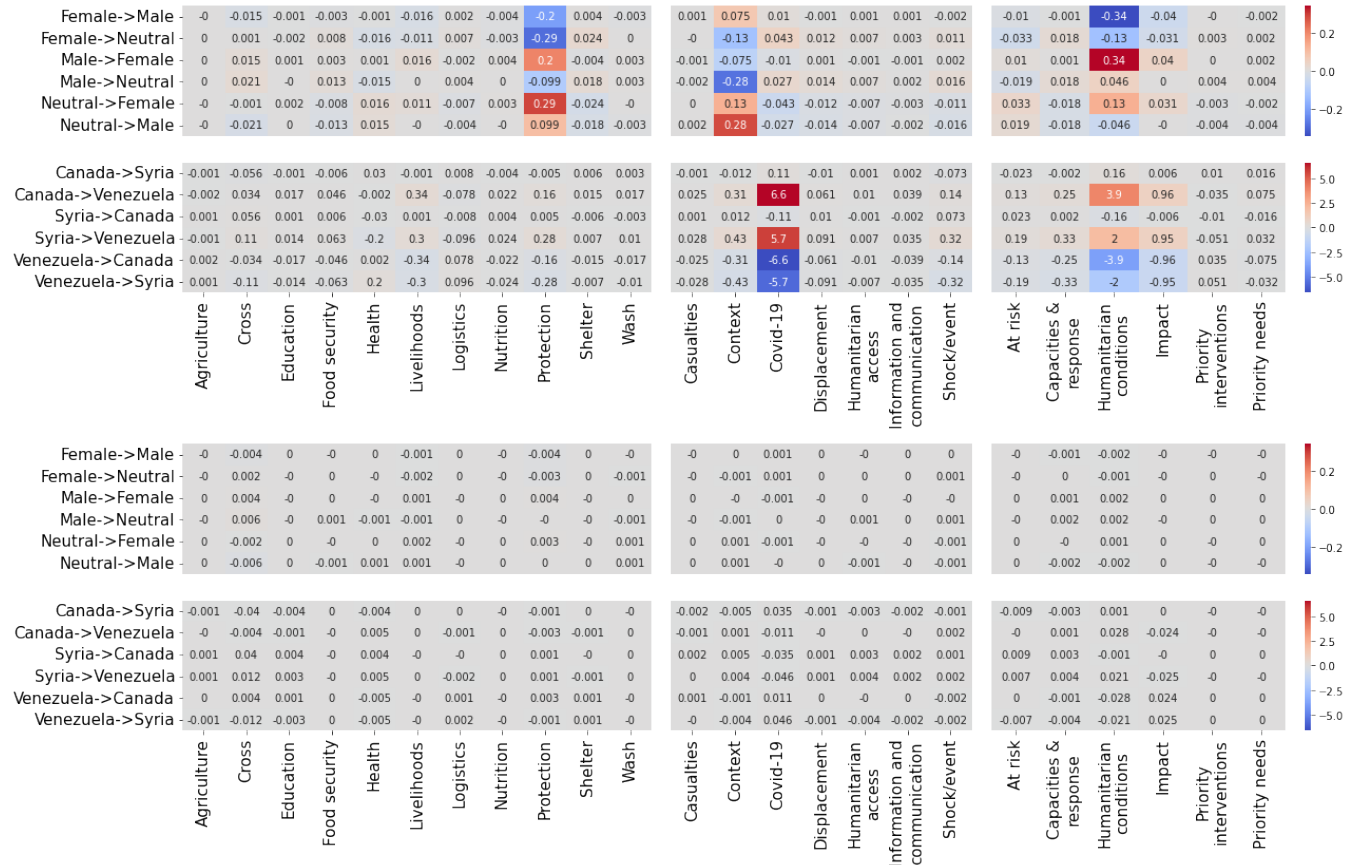


Figure 1: The results of Tag-Shift bias metric for BASE architecture using the HumBERT as the backbone. (Top) Original model without debiasing; (Bottom) Counterfactual debiasing.

**Figure 2 (Top) — Original model without debiasing — Gender**

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Female->Male | 0 | -0.022 | -0.006 | 0.007 | 0.002 | -0.006 | 0.007 | -0.011 | -0.28 | 0.007 | -0.014 |
| Female->Neutral | -0.001 | 0.01 | -0.005 | 0.019 | 0.007 | -0.011 | 0.01 | -0.017 | -0.37 | 0.012 | -0.018 |
| Male->Female | -0 | 0.022 | 0.006 | -0.007 | -0.002 | 0.006 | -0.007 | 0.011 | 0.28 | -0.007 | 0.014 |
| Male->Neutral | -0.001 | 0.02 | -0 | 0.009 | 0.003 | -0.002 | 0.001 | -0.004 | -0.043 | 0.005 | -0.002 |
| Neutral->Female | 0.001 | -0.01 | 0.005 | -0.019 | -0.007 | 0.011 | -0.01 | 0.017 | 0.37 | -0.012 | 0.018 |
| Neutral->Male | 0.001 | -0.02 | 0 | -0.009 | -0.003 | 0.002 | -0.001 | 0.004 | 0.043 | -0.005 | 0.002 |

| | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event |
|---|---|---|---|---|---|---|---|
| Female->Male | 0 | -0.042 | 0.021 | -0.01 | 0.004 | 0.001 | 0.011 |
| Female->Neutral | -0 | -0.16 | 0.076 | 0.005 | 0.008 | 0.004 | 0.043 |
| Male->Female | -0 | 0.042 | -0.021 | 0.01 | -0.004 | -0.001 | -0.011 |
| Male->Neutral | -0.001 | -0.079 | 0.037 | 0.02 | 0.002 | 0.002 | 0.027 |
| Neutral->Female | 0 | 0.16 | -0.076 | -0.005 | -0.008 | -0.004 | -0.043 |
| Neutral->Male | 0.001 | 0.079 | -0.037 | -0.02 | -0.002 | -0.002 | -0.027 |

| | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|
| Female->Male | -0.053 | -0.015 | -0.13 | -0.079 | -0.001 | -0.003 |
| Female->Neutral | -0.062 | -0.005 | -0.094 | -0.031 | 0.002 | -0 |
| Male->Female | 0.053 | 0.015 | 0.13 | 0.079 | 0.001 | 0.003 |
| Male->Neutral | -0.006 | 0.006 | -0.006 | 0.018 | 0.003 | 0.003 |
| Neutral->Female | 0.062 | 0.005 | 0.094 | 0.031 | -0.002 | 0 |
| Neutral->Male | 0.006 | -0.006 | 0.006 | -0.018 | -0.003 | -0.003 |

**Original model without debiasing — Country**

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Canada->Syria | -0.01 | -0.011 | -0.001 | -0.004 | 0.1 | -0.071 | 0.013 | -0.004 | 0.036 | 0.007 | 0.005 |
| Canada->Venezuela | -0.011 | 0.16 | 0.006 | 0.051 | 0.16 | 0.37 | -0.13 | 0.014 | 0.19 | 0.03 | 0.02 |
| Syria->Canada | 0.01 | 0.011 | 0.001 | 0.004 | -0.1 | 0.071 | -0.013 | 0.004 | -0.036 | -0.007 | -0.005 |
| Syria->Venezuela | -0.001 | 0.2 | 0.009 | 0.069 | -0.014 | 0.42 | -0.13 | 0.023 | 0.11 | 0.018 | 0.003 |
| Venezuela->Canada | 0.011 | -0.16 | -0.006 | -0.051 | -0.16 | -0.37 | 0.13 | -0.014 | -0.19 | -0.03 | -0.02 |
| Venezuela->Syria | 0.001 | -0.2 | -0.009 | -0.069 | 0.014 | -0.42 | 0.13 | -0.023 | -0.11 | -0.018 | -0.003 |

| | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event |
|---|---|---|---|---|---|---|---|
| Canada->Syria | 0.007 | -0.074 | 0.077 | -0.007 | -0.002 | 0.008 | -0.18 |
| Canada->Venezuela | 0.055 | 0.56 | 2.7 | 0.15 | 0.024 | -0.002 | 0.49 |
| Syria->Canada | -0.007 | 0.074 | -0.077 | 0.007 | 0.002 | -0.008 | 0.18 |
| Syria->Venezuela | 0.029 | 0.84 | 2.5 | 0.19 | 0.037 | -0.012 | 0.99 |
| Venezuela->Canada | -0.055 | -0.56 | -2.7 | -0.15 | -0.024 | 0.002 | -0.49 |
| Venezuela->Syria | -0.029 | -0.84 | -2.5 | -0.19 | -0.037 | 0.012 | -0.99 |

| | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|
| Canada->Syria | 0.01 | -0.075 | 0.13 | 0.26 | -0.003 | 0.003 |
| Canada->Venezuela | 0.19 | 0.94 | 4.6 | 2.8 | -0.078 | 0.083 |
| Syria->Canada | -0.01 | 0.075 | -0.13 | -0.26 | 0.003 | -0.003 |
| Syria->Venezuela | 0.14 | 1.1 | 2.4 | 1.7 | -0.07 | 0.072 |
| Venezuela->Canada | -0.19 | -0.94 | -4.6 | -2.8 | 0.078 | -0.083 |
| Venezuela->Syria | -0.14 | -1.1 | -2.4 | -1.7 | 0.07 | -0.072 |

**(Bottom) — Counterfactual debiasing — Gender**

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Female->Male | -0 | -0.009 | -0 | 0 | 0.001 | -0.003 | 0 | -0 | -0.017 | 0.001 | -0 |
| Female->Neutral | -0 | -0.002 | 0 | 0 | 0.001 | -0.004 | 0 | -0.001 | -0.013 | 0 | -0.001 |
| Male->Female | -0 | 0.009 | 0 | -0.001 | 0.003 | -0 | 0 | 0 | 0.017 | -0.001 | 0 |
| Male->Neutral | -0 | 0.003 | 0 | 0 | -0 | -0.001 | 0 | 0.001 | 0.001 | -0 | -0 |
| Neutral->Female | 0 | 0.002 | 0 | 0 | -0.001 | 0.004 | -0 | 0.001 | 0.013 | -0 | 0.001 |
| Neutral->Male | 0 | -0.003 | -0 | -0 | 0 | 0.001 | -0 | 0.001 | -0.001 | 0 | 0 |

| | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event |
|---|---|---|---|---|---|---|---|
| Female->Male | -0 | -0.003 | 0.001 | -0 | 0 | -0 | 0.001 |
| Female->Neutral | -0 | -0.005 | 0.009 | -0.003 | 0 | 0 | 0.002 |
| Male->Female | -0 | 0.003 | -0.001 | -0 | 0 | 0 | -0.001 |
| Male->Neutral | -0 | -0.002 | 0.007 | -0.002 | 0 | 0 | 0 |
| Neutral->Female | 0 | 0.005 | -0.009 | 0.003 | -0 | -0 | -0.002 |
| Neutral->Male | 0 | 0.002 | -0.007 | 0.002 | -0 | 0 | -0 |

| | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|
| Female->Male | -0.001 | 0.001 | 0.001 | -0.001 | -0 | -0 |
| Female->Neutral | -0.001 | 0.003 | 0.001 | 0.011 | -0 | 0 |
| Male->Female | 0.001 | -0.001 | -0.001 | 0.001 | 0 | 0 |
| Male->Neutral | 0 | 0.002 | -0.001 | 0.014 | -0 | 0 |
| Neutral->Female | 0.001 | -0.003 | -0.001 | -0.011 | 0 | 0 |
| Neutral->Male | -0 | -0.002 | 0.001 | -0.014 | 0 | -0 |

**Counterfactual debiasing — Country**

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Canada->Syria | -0.002 | -0.014 | -0.001 | 0.001 | 0.055 | -0.005 | 0.002 | -0.001 | 0.001 | -0.002 | -0.001 |
| Canada->Venezuela | -0.001 | -0.007 | 0.004 | 0 | 0.003 | 0.006 | 0.001 | -0 | -0.01 | 0.002 | -0.001 |
| Syria->Canada | 0.002 | 0.014 | 0.001 | -0.001 | -0.055 | 0.005 | -0.002 | 0.001 | -0.001 | 0.002 | 0.001 |
| Syria->Venezuela | 0.001 | 0.003 | 0.005 | -0 | -0.046 | 0.015 | 0.001 | 0 | -0.014 | 0.002 | 0 |
| Venezuela->Canada | 0.001 | 0.007 | -0.004 | -0 | -0.003 | -0.006 | -0.001 | 0 | 0.01 | -0.002 | 0.001 |
| Venezuela->Syria | -0.001 | -0.003 | -0.005 | 0 | 0.046 | -0.015 | -0.001 | 0 | 0.014 | -0.002 | -0 |

| | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event |
|---|---|---|---|---|---|---|---|
| Canada->Syria | 0 | -0.002 | 0.004 | 0.001 | 0 | 0.001 | -0.003 |
| Canada->Venezuela | 0 | 0.015 | -0.015 | 0.005 | 0 | 0.001 | -0.002 |
| Syria->Canada | -0 | 0.002 | -0.004 | -0.001 | 0 | -0.001 | 0.003 |
| Syria->Venezuela | -0 | 0.02 | -0.01 | 0.006 | -0 | 0.001 | 0.001 |
| Venezuela->Canada | -0 | -0.015 | 0.015 | -0.005 | 0 | -0.001 | 0.002 |
| Venezuela->Syria | 0 | -0.02 | 0.01 | -0.006 | 0 | -0.001 | -0.001 |

| | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|
| Canada->Syria | -0.003 | 0.003 | 0.02 | 0.032 | 0.002 | 0.002 |
| Canada->Venezuela | -0.003 | 0.002 | 0.012 | 0.011 | 0 | 0.004 |
| Syria->Canada | 0.003 | -0.003 | -0.02 | -0.032 | -0.002 | -0.002 |
| Syria->Venezuela | -0 | -0.003 | -0.002 | -0.031 | -0.001 | 0.002 |
| Venezuela->Canada | 0.003 | -0.002 | -0.012 | -0.011 | -0 | -0.004 |
| Venezuela->Syria | 0 | 0.003 | 0.002 | 0.031 | 0.001 | -0.002 |

Figure 2: The results of Tag-Shift bias metric for BASE architecture using the XLM-R as the backbone. (Top) Original model without debiasing; (Bottom) Counterfactual debiasing.
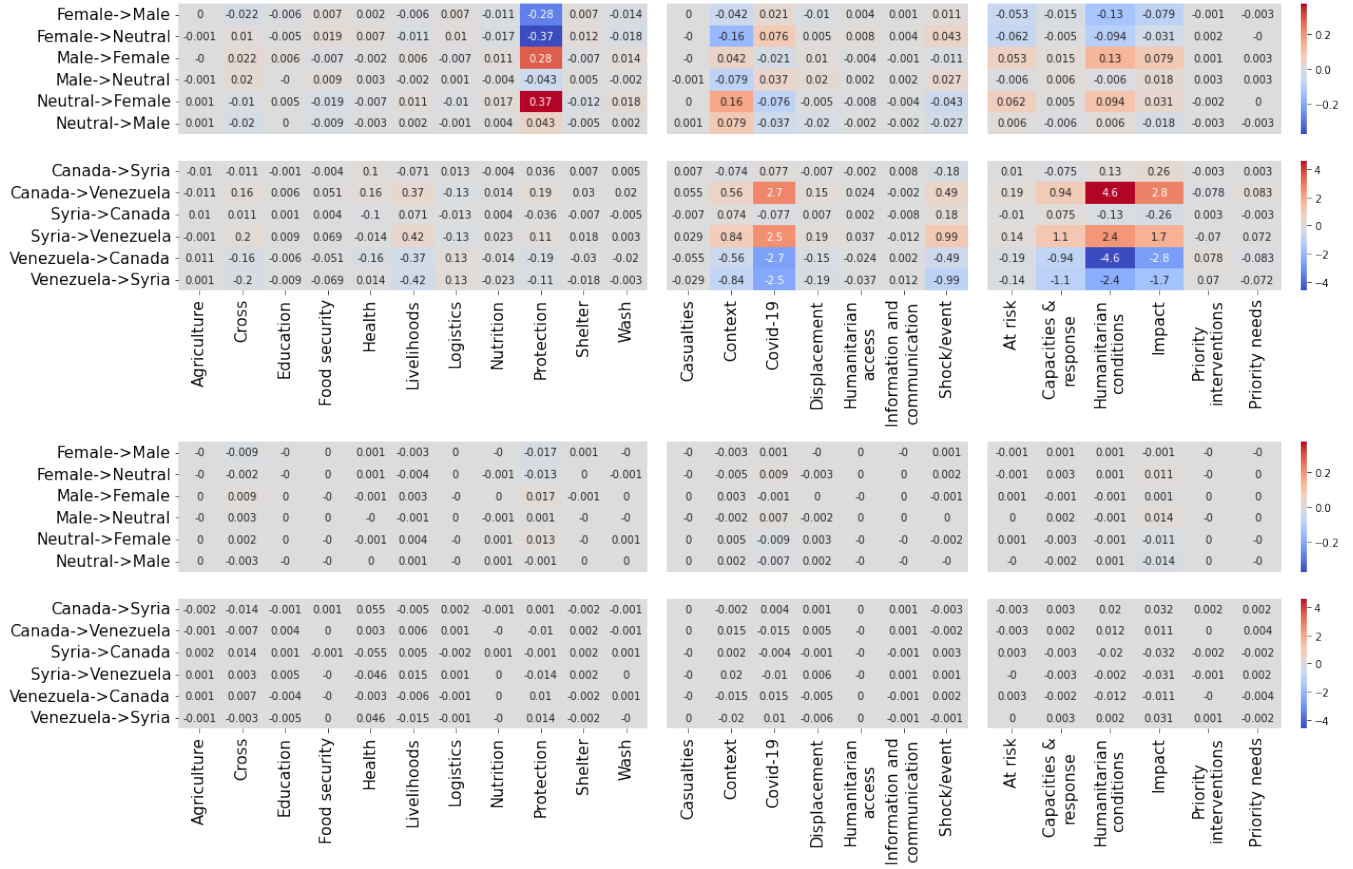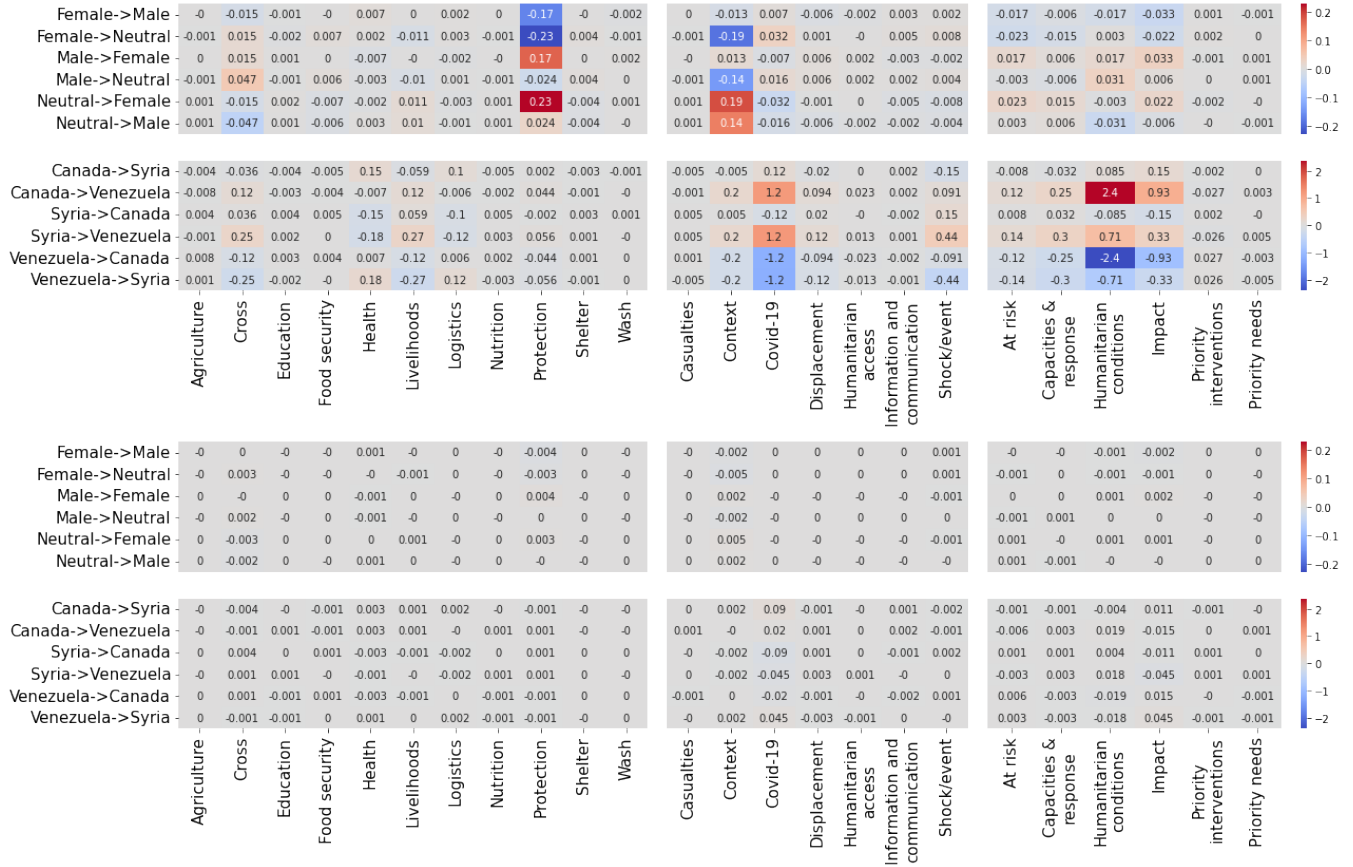
Figure 3: The results of Tag-Shift bias metric for OURS architecture using the XLM-R as the backbone. (Top) Original model without debiasing; (Bottom) Counterfactual debiasing.

**Top — Original model without debiasing (gender):**

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Female->Male | 0 | 0.001 | -0.003 | 0.004 | -0.018 | -0.011 | 0.003 | -0.005 | -0.21 | 0.001 | -0.002 | 0.001 | 0.001 | 0.013 | -0.018 | 0.004 | 0.001 | 0.001 | -0.042 | 0 | -0.097 | -0.038 | 0 | -0.003 |
| Female->Neutral | -0 | -0.05 | -0.009 | 0.019 | -0.048 | -0.012 | 0.004 | -0.017 | -0.3 | 0.009 | -0.002 | -0 | -0.11 | 0.058 | -0.012 | 0.007 | 0.007 | 0.007 | -0.044 | 0.006 | -0.072 | -0.051 | 0.005 | 0 |
| Male->Female | -0 | -0.001 | 0.003 | -0.004 | 0.018 | 0.011 | -0.003 | 0.005 | 0.21 | -0.001 | 0.002 | -0.001 | -0.001 | -0.013 | 0.018 | -0.004 | -0.001 | -0.001 | 0.042 | -0 | 0.097 | 0.038 | -0 | 0.003 |
| Male->Neutral | -0.001 | -0.057 | -0.004 | 0.01 | -0.025 | -0.003 | 0.001 | -0.005 | -0.069 | 0.008 | -0 | -0.001 | -0.11 | 0.031 | 0.002 | 0.003 | 0.004 | 0.003 | -0.006 | 0.009 | -0.007 | -0.011 | 0.003 | 0.002 |
| Neutral->Female | 0 | 0.05 | 0.009 | -0.019 | 0.048 | 0.012 | -0.004 | 0.017 | 0.3 | -0.009 | 0.002 | 0 | 0.11 | -0.058 | 0.012 | -0.007 | -0.007 | -0.007 | 0.044 | -0.006 | 0.072 | 0.051 | -0.005 | -0 |
| Neutral->Male | 0.001 | 0.057 | 0.004 | -0.01 | 0.025 | 0.003 | -0.001 | 0.005 | 0.069 | -0.008 | 0 | 0.001 | 0.11 | -0.031 | -0.002 | -0.003 | -0.004 | -0.003 | 0.006 | -0.009 | 0.007 | 0.011 | -0.003 | -0.002 |

**Top — Original model without debiasing (country):**

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Canada->Syria | -0.009 | -0.029 | -0.003 | -0.004 | 0.19 | -0.096 | 0.028 | -0.016 | -0.003 | -0.005 | 0.01 | -0.001 | -0.043 | 0.4 | -0.063 | -0.016 | 0.017 | -0.1 | 0.008 | -0.015 | 0.028 | -0.27 | -0.003 | -0.011 |
| Canada->Venezuela | -0.017 | 0.28 | 0.036 | 0.04 | 0.037 | 0.3 | -0.051 | 0.023 | 0.25 | 0.17 | 0.002 | 0.017 | 0.23 | 5.4 | 0.23 | 0.017 | 0.026 | 0.98 | 0.31 | 0.12 | 2.4 | 3 | -0.05 | 0.12 |
| Syria->Canada | 0.009 | 0.029 | 0.003 | 0.004 | -0.19 | 0.096 | -0.028 | 0.016 | 0.003 | 0.005 | -0.01 | 0.001 | 0.043 | -0.4 | 0.063 | 0.016 | -0.017 | 0.1 | -0.008 | 0.015 | -0.028 | 0.27 | 0.003 | 0.011 |
| Syria->Venezuela | -0.006 | 0.31 | 0.041 | 0.073 | -0.17 | 0.59 | -0.11 | 0.044 | 0.24 | 0.17 | -0.006 | 0.021 | 0.31 | 3.8 | 0.38 | 0.047 | 0.004 | 1.4 | 0.31 | 0.14 | 1.8 | 3.6 | -0.055 | 0.13 |
| Venezuela->Canada | 0.017 | -0.28 | -0.036 | -0.04 | -0.037 | -0.3 | 0.051 | -0.023 | -0.25 | -0.17 | -0.002 | -0.017 | -0.23 | -5.4 | -0.23 | -0.017 | -0.026 | -0.98 | -0.31 | -0.12 | -2.4 | -3 | 0.05 | -0.12 |
| Venezuela->Syria | 0.006 | -0.31 | -0.041 | -0.073 | 0.17 | -0.59 | 0.11 | -0.044 | -0.24 | -0.17 | 0.006 | -0.021 | -0.31 | -3.8 | -0.38 | -0.047 | -0.004 | -1.4 | -0.31 | -0.14 | -1.8 | -3.6 | 0.055 | -0.13 |

**Bottom — Counterfactual debiasing (gender):**

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Female->Male | -0 | -0 | -0 | 0.001 | 0.001 | -0.002 | 0 | 0 | -0.002 | 0 | -0 | 0 | -0.002 | 0.001 | 0.001 | 0 | 0 | 0.001 | -0 | -0.003 | 0.008 | -0.001 | -0 | 0 |
| Female->Neutral | -0 | 0 | 0 | 0 | 0.001 | -0.002 | 0 | -0 | -0.001 | 0 | -0 | -0 | -0.003 | 0.001 | 0.001 | 0 | 0 | 0.001 | 0 | 0 | 0.003 | -0.002 | 0 | 0 |
| Male->Female | 0 | 0 | 0 | -0.001 | -0.001 | 0.002 | -0 | 0 | 0.002 | -0 | 0 | -0 | 0.002 | -0.001 | -0.001 | -0 | 0 | -0.001 | 0 | 0.003 | -0.008 | 0.001 | 0 | -0 |
| Male->Neutral | -0 | 0.003 | 0 | -0 | -0 | 0 | -0 | 0 | 0.001 | 0 | -0 | -0 | -0.001 | 0.001 | 0 | -0 | 0 | -0 | 0 | 0.003 | -0.004 | -0.001 | 0 | -0 |
| Neutral->Female | 0 | -0 | -0 | -0.001 | 0.002 | -0 | 0 | 0.001 | 0 | 0 | | 0 | 0.003 | -0.001 | -0.001 | -0 | -0 | -0.001 | -0 | -0 | -0.003 | 0.002 | -0 | -0 |
| Neutral->Male | 0 | -0.003 | -0 | 0 | 0 | 0 | -0 | 0 | -0.001 | -0 | -0 | 0 | 0.001 | -0.001 | -0 | 0 | -0 | 0 | -0 | -0.003 | 0.004 | 0.001 | -0 | 0 |

**Bottom — Counterfactual debiasing (country):**

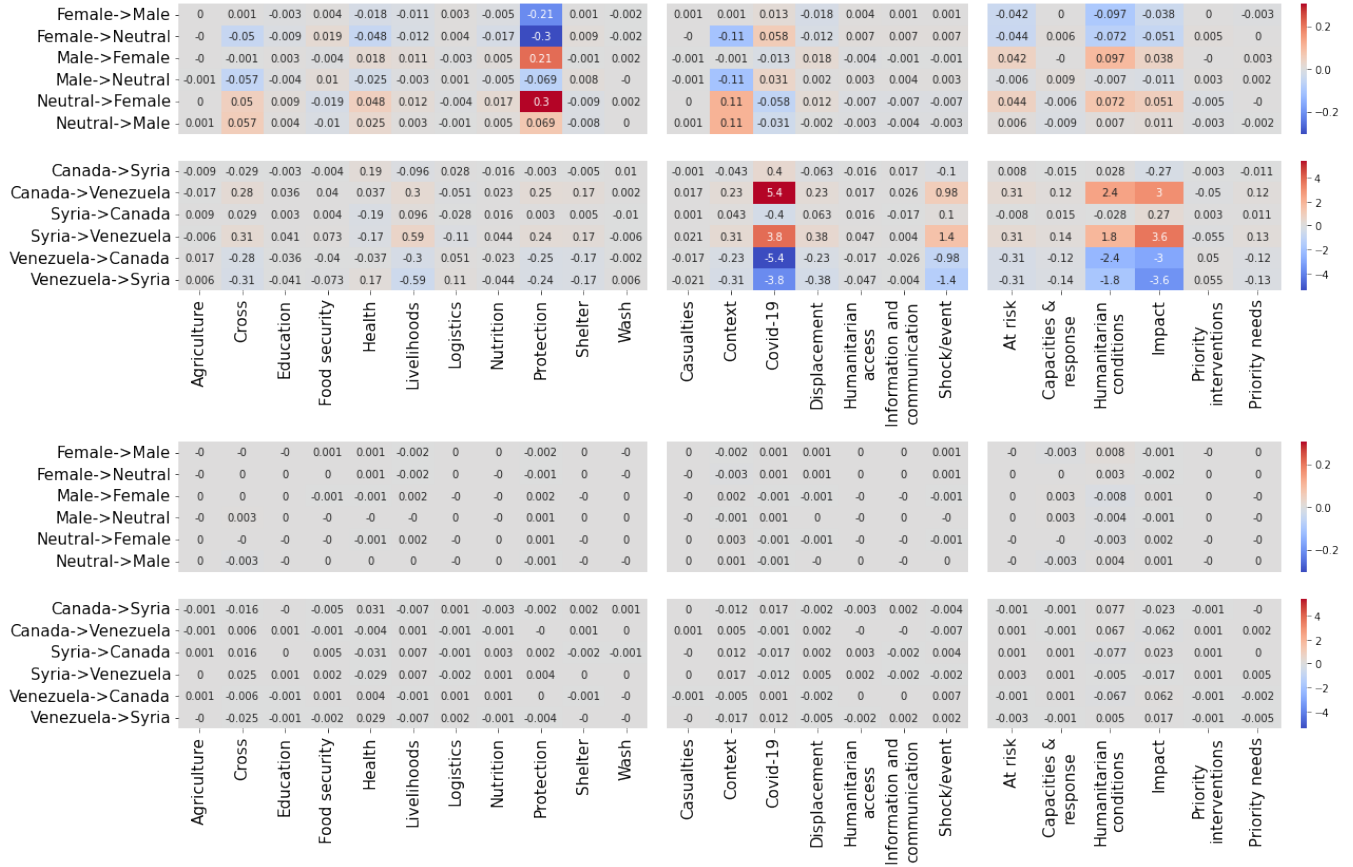| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Canada->Syria | -0.001 | -0.016 | -0 | -0.005 | 0.031 | -0.007 | 0.001 | -0.003 | -0.002 | 0.002 | 0.001 | 0 | -0.012 | 0.017 | -0.002 | -0.003 | 0.002 | -0.004 | -0.001 | -0.001 | 0.077 | -0.023 | -0.001 | -0 |
| Canada->Venezuela | -0.001 | 0.006 | 0.001 | -0.001 | -0.004 | 0.001 | -0.001 | -0.001 | -0 | 0.001 | 0 | 0.001 | 0.005 | -0.001 | 0.002 | -0 | -0 | -0.007 | 0.001 | -0.001 | 0.067 | -0.062 | 0.001 | 0.002 |
| Syria->Canada | 0.001 | 0.016 | 0 | 0.005 | -0.031 | 0.007 | -0.001 | 0.003 | 0.002 | -0.002 | -0.001 | -0 | 0.012 | -0.017 | 0.002 | 0.003 | -0.002 | 0.004 | 0.001 | 0.001 | -0.077 | 0.023 | 0.001 | 0 |
| Syria->Venezuela | 0 | 0.025 | 0.001 | 0.002 | -0.029 | 0.007 | -0.002 | 0.001 | 0.004 | 0 | 0 | 0 | 0.017 | -0.012 | 0.005 | 0.002 | -0.002 | -0.002 | 0.003 | 0.001 | -0.005 | -0.017 | 0.001 | 0.005 |
| Venezuela->Canada | 0.001 | -0.006 | -0.001 | 0.001 | 0.004 | -0.001 | 0.001 | 0.001 | 0 | -0.001 | -0 | -0.001 | -0.005 | 0.001 | -0.002 | 0 | 0 | 0.007 | -0.001 | 0.001 | -0.067 | 0.062 | -0.001 | -0.002 |
| Venezuela->Syria | -0 | -0.025 | -0.001 | -0.002 | 0.029 | -0.007 | 0.002 | -0.001 | -0.004 | 0 | 0 | -0 | -0.017 | 0.012 | -0.005 | -0.002 | 0.002 | 0.002 | -0.003 | -0.001 | 0.005 | 0.017 | -0.001 | -0.005 |

Figure 4: The results of Tag-Shift bias metric for BASE architecture using the m-BERT as the backbone. (Top) Original model without debiasing; (Bottom) Counterfactual debiasing.

**Figure 5: Original model without debiasing (Top)**

Gender Tag-Shift:

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Female->Male | -0 | 0.008 | 0.001 | 0.002 | 0.003 | -0.011 | 0.001 | -0.001 | -0.24 | 0.001 | -0.001 | 0 | 0.011 | 0.001 | -0.011 | 0.002 | 0.001 | 0.003 | -0.026 | -0.008 | -0.069 | -0.045 | 0 | 0 |
| Female->Neutral | -0 | 0.003 | -0.002 | 0.008 | -0 | -0.013 | 0.002 | -0.003 | -0.29 | 0.006 | 0.002 | -0.002 | -0.17 | 0.026 | -0.003 | 0.005 | 0.007 | 0.023 | -0.034 | -0.006 | -0.023 | -0.042 | 0.001 | 0.002 |
| Male->Female | 0 | -0.008 | -0.001 | -0.002 | -0.003 | 0.011 | -0.001 | 0.001 | 0.24 | -0.001 | 0.001 | -0 | -0.011 | -0.001 | 0.011 | -0.002 | -0.001 | -0.003 | 0.026 | 0.008 | 0.069 | 0.045 | -0 | -0 |
| Male->Neutral | 0 | 0 | -0.002 | 0.005 | -0.002 | -0.002 | 0 | -0.001 | -0.027 | 0.004 | 0.003 | -0.002 | -0.2 | 0.02 | 0.004 | 0.002 | 0.003 | 0.017 | -0.003 | 0 | 0.01 | -0.005 | 0 | 0.001 |
| Neutral->Female | 0 | -0.003 | 0.002 | -0.008 | 0 | 0.013 | -0.002 | 0.003 | 0.29 | -0.006 | -0.002 | 0.002 | 0.17 | -0.026 | 0.003 | -0.005 | -0.007 | -0.023 | 0.034 | 0.006 | 0.023 | 0.042 | -0.001 | -0.002 |
| Neutral->Male | 0 | -0 | 0.002 | -0.005 | 0.002 | 0.002 | -0 | 0.001 | 0.027 | -0.004 | -0.003 | 0.002 | 0.2 | -0.02 | -0.004 | -0.002 | -0.003 | -0.017 | 0.003 | -0 | -0.01 | 0.005 | -0 | -0.001 |

Country Tag-Shift:

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Canada->Syria | -0.006 | 0.017 | -0.005 | -0.02 | 0.056 | -0.056 | 0.041 | -0.005 | 0.009 | 0.004 | 0.003 | -0.002 | -0.11 | 0.065 | 0 | 0.002 | 0.003 | -0.017 | -0.001 | -0.018 | 0.15 | -0.27 | 0.003 | 0.003 |
| Canada->Venezuela | -0.009 | 0.062 | 0.015 | -0.002 | -0.013 | 0.11 | -0.042 | 0.012 | 0.077 | 0.026 | 0.004 | 0.006 | 0.22 | 6.2 | 0.066 | 0.049 | 0.013 | 0.25 | 0.076 | 0.15 | 0.69 | 0.93 | -0.011 | 0.065 |
| Syria->Canada | 0.006 | -0.017 | 0.005 | 0.02 | -0.056 | 0.056 | -0.041 | 0.005 | -0.009 | -0.004 | -0.003 | 0.002 | 0.11 | -0.065 | -0 | -0.002 | -0.003 | 0.017 | 0.001 | 0.018 | -0.15 | 0.27 | -0.003 | -0.003 |
| Syria->Venezuela | -0.002 | 0.04 | 0.022 | 0.013 | -0.085 | 0.22 | -0.096 | 0.02 | 0.063 | 0.016 | -0 | 0.007 | 0.51 | 5.8 | 0.07 | 0.041 | 0.011 | 0.31 | 0.08 | 0.19 | 0.57 | 0.92 | -0.015 | 0.051 |
| Venezuela->Canada | 0.009 | -0.062 | -0.015 | 0.002 | 0.013 | -0.11 | 0.042 | -0.012 | -0.077 | -0.026 | -0.004 | -0.006 | -0.22 | -6.2 | -0.066 | -0.049 | -0.013 | -0.25 | -0.076 | -0.15 | -0.69 | -0.93 | 0.011 | -0.065 |
| Venezuela->Syria | 0.002 | -0.04 | -0.022 | -0.013 | 0.085 | -0.22 | 0.096 | -0.02 | -0.063 | -0.016 | 0 | -0.007 | -0.51 | -5.8 | -0.07 | -0.041 | -0.011 | -0.31 | -0.08 | -0.19 | -0.57 | -0.92 | 0.015 | -0.051 |

**Counterfactual debiasing (Bottom)**

Gender Tag-Shift:

| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Female->Male | -0 | 0.001 | 0 | 0 | 0 | -0.001 | 0 | -0 | -0.002 | 0 | -0 | -0 | -0 | 0 | -0 | 0 | 0 | 0 | -0 | 0 | -0 | 0 | 0 | 0 |
| Female->Neutral | -0 | 0.006 | -0 | 0 | -0.001 | -0.001 | 0 | -0.001 | -0.004 | 0 | 0 | -0 | -0.002 | 0 | 0 | 0 | 0 | 0.001 | -0.001 | 0.002 | -0.002 | -0.001 | -0 | -0 |
| Male->Female | -0 | -0.001 | -0 | -0 | 0.001 | -0 | 0 | 0.002 | -0 | 0 | -0 | 0 | 0 | -0 | 0 | 0 | 0 | -0 | 0 | 0 | 0 | 0 | -0 | -0 |
| Male->Neutral | -0 | 0.003 | -0 | 0 | -0.002 | -0.001 | -0 | -0.001 | 0 | 0 | 0 | -0 | -0.001 | 0 | 0 | 0 | 0 | 0.001 | -0 | 0.002 | -0.002 | -0.001 | -0 | -0 |
| Neutral->Female | 0 | -0.006 | 0 | 0 | 0.001 | 0.001 | -0 | 0.001 | 0.004 | -0 | 0 | 0 | 0.002 | -0 | 0 | -0 | 0 | -0.001 | 0.001 | -0.002 | 0.002 | 0.001 | 0 | 0 |
| Neutral->Male | 0 | -0.003 | 0 | -0 | 0.002 | 0.001 | 0 | 0 | 0.001 | -0 | -0 | 0 | 0.001 | -0 | 0 | -0 | 0 | -0.001 | 0 | -0.002 | 0.002 | 0.001 | 0 | 0 |

Country Tag-Shift:

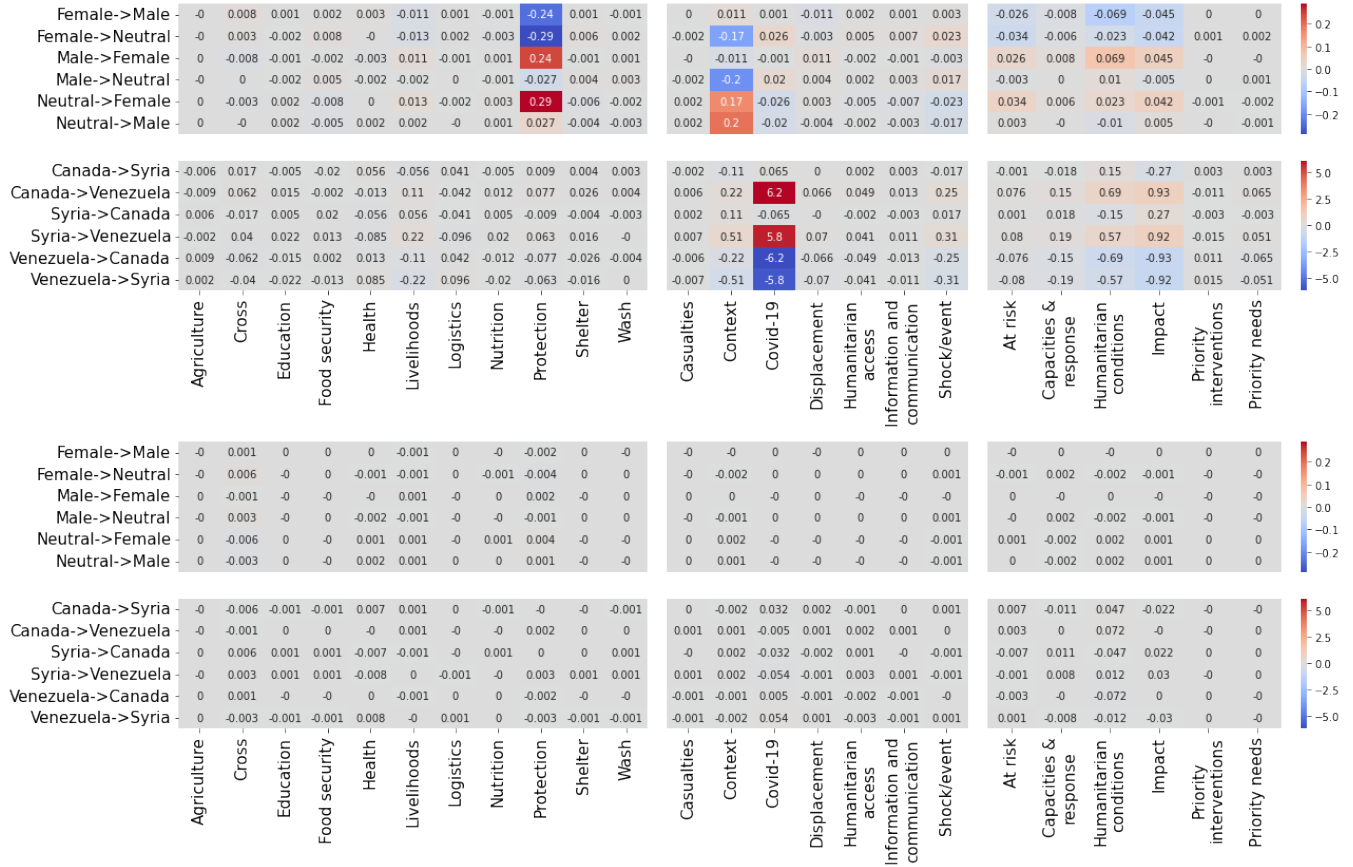| | Agriculture | Cross | Education | Food security | Health | Livelihoods | Logistics | Nutrition | Protection | Shelter | Wash | Casualties | Context | Covid-19 | Displacement | Humanitarian access | Information and communication | Shock/event | At risk | Capacities & response | Humanitarian conditions | Impact | Priority interventions | Priority needs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Canada->Syria | -0 | -0.006 | -0.001 | -0.001 | 0.007 | 0.001 | 0 | -0.001 | -0 | -0 | -0.001 | 0 | -0.002 | 0.032 | 0.002 | -0.001 | 0 | 0.001 | 0.007 | -0.011 | 0.047 | -0.022 | -0 | 0 |
| Canada->Venezuela | -0 | -0.001 | 0 | 0 | -0 | 0.001 | -0 | -0 | 0.002 | 0 | 0 | 0.001 | 0.001 | -0.005 | 0.001 | 0.002 | 0.001 | 0 | 0.003 | 0 | 0.072 | -0 | -0 | 0 |
| Syria->Canada | 0 | 0.006 | 0.001 | 0.001 | -0.007 | -0.001 | -0 | 0.001 | 0 | 0 | 0.001 | 0 | 0.002 | -0.032 | -0.002 | 0.001 | -0 | -0.001 | -0.007 | 0.011 | -0.047 | 0.022 | 0 | 0 |
| Syria->Venezuela | -0 | 0.003 | 0.001 | 0.001 | -0.008 | 0 | -0.001 | 0 | 0.003 | 0.001 | 0.001 | 0.001 | 0.002 | -0.054 | -0.001 | 0.003 | 0.001 | -0.001 | -0.001 | 0.008 | 0.012 | 0.03 | -0 | 0 |
| Venezuela->Canada | 0 | 0.001 | -0 | -0 | 0 | -0.001 | 0 | 0 | -0.002 | -0 | -0 | -0.001 | -0.001 | 0.005 | -0.001 | -0.002 | -0.001 | -0 | -0.003 | -0 | -0.072 | 0 | 0 | -0 |
| Venezuela->Syria | 0 | -0.003 | -0.001 | -0.001 | 0.008 | -0 | 0.001 | 0 | -0.003 | -0.001 | -0.001 | -0.001 | -0.002 | 0.054 | 0.001 | -0.003 | -0.001 | 0.001 | 0.001 | -0.008 | -0.012 | -0.03 | 0 | -0 |

Figure 5: The results of Tag-Shift bias metric for OURS architecture using the m-BERT as the backbone. (Top) Original model without debiasing; (Bottom) Counterfactual debiasing.