

Modeling Commenter Personas in Climate Misinformation Discourse

Anonymous ACL submission

Abstract

Prior work on climate misinformation detection typically classifies individual comments or claims using text alone. This framing overlooks a key property of online climate discourse: misinformation is often produced, reinforced, and sustained through stable speaker-level patterns, where the same individuals repeatedly employ characteristic narratives and rhetorical styles across interactions. We address this gap by introducing personas: data-driven groupings of commenters that capture recurring discourse strategies beyond the content of any single comment. Empirically, modeling persona improves misinformation classification performance beyond text-only baselines in low-capacity and resource-constrained settings, increasing accuracy from 0.79 to 0.81 and macro-F1 from 0.72 to 0.74 in a low-capacity MLP classifier (statistically significant, $p = 0.021$). Beyond detection, we find that personas shape how misinformation is received and propagated. Persona-conditioned language model agents exhibit systematic differences in misinformation acceptance that align with human persona-level susceptibility, and the relationship between misinformation prevalence and content persistence varies by the dominant conversational persona.

1 Introduction

Short-video social media platforms such as TikTok have become major venues for public discussion of climate change, where accurate information, skepticism, and misinformation coexist in the same high-volume comment threads (Pearce et al., 2018; Sultana et al., 2024; Basch et al., 2021; Pera and Aiello, 2024; Baltasar et al., 2024). This setting is particularly challenging for automated analysis: short, algorithmically amplified videos drive rapid diffusion, while discussions are fast-moving, emotionally charged, allowing misinformation to persist and resist correction.

As shown in Figure 1a, most prior work uses text-only classifiers that label each comment or

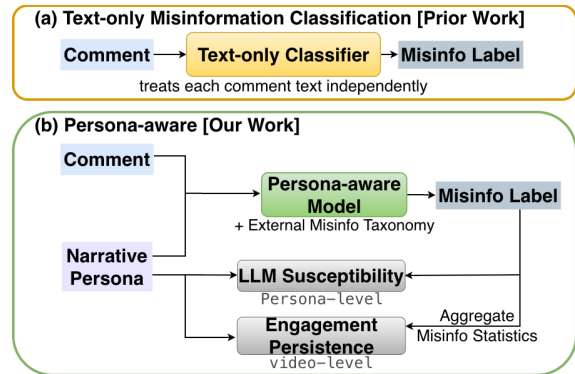


Figure 1: Comparison of text-only vs. persona-aware misinformation modeling. (a) Prior work (Islam et al., 2020; Coan et al., 2021; Upadhyaya et al., 2023; Choi et al., 2023; Rojas et al., 2024) typically predicts a misinformation category from comment text alone, treating each comment independently. (b) We induce narrative personas by clustering comment embeddings and use persona as additional context for misinformation-type prediction; the same persona abstraction also supports persona-conditioned analyses of LLM susceptibility and engagement persistence.

claim in isolation (Islam et al., 2020; Coan et al., 2021; Upadhyaya et al., 2023; Choi et al., 2023; Rojas et al., 2024). However, climate discourse also exhibits speaker-level regularities: recurring narrative frames and rhetorical styles that reappear across threads. This suggests that *how people consistently argue* can provide complementary signal beyond the words in a single comment.

We introduce *personas* to address this challenge: a compact, data-driven abstraction of recurring narrative frames and rhetorical styles in climate discourse. A growing body of work supports the importance of speaker and social context for misinformation and stance modeling. Incorporating user and interaction structure can improve misinformation detection beyond text alone, encouraging representations that aggregate signal across interactions rather than treating each instance independently (Nguyen et al., 2020; Mehta et al., 2022).

063 Stance research similarly argues that beliefs and
064 framing are more faithfully modeled at the author
065 level and are not reliably captured by lexical fea-
066 tures alone (Allaway and McKeown, 2021; Zhang
067 et al., 2024). Communication scholarship further
068 emphasizes that disinformation propagates through
069 identity linked, adversarial narratives that remain
070 stable within groups and are expressed via charac-
071 teristic rhetorical strategies (Díaz Ruiz and Nilsson,
072 2023). Together, these perspectives suggest that
073 *who is speaking* provides a complementary axis for
074 understanding climate misinformation beyond the
075 words in a single comment.

076 Rather than relying on predefined demographics
077 or static author embeddings, we induce personas di-
078 rectly from large-scale, in-the-wild TikTok climate
079 comments by clustering semantic comment embed-
080 dings and consolidating clusters into interpretable
081 persona types. This persona representation is then
082 used in three ways (Figure 1b): (i) to augment fine-
083 grained misinformation-type prediction, with a fo-
084 cus on whether persona provides additional signal
085 in low-capacity or limited-supervision settings, (ii)
086 as a scaffold for a controlled susceptibility probe
087 using persona-conditioned LLM agents¹, and (iii)
088 as an organizing variable for studying how mis-
089 information relates to engagement persistence in
090 comment threads.

091 Our main contributions are:

- 092 • We induce commenter personas from TikTok
093 climate discourse and compress large-scale
094 clusters into interpretable personas.
- 095 • We show that persona context improves fine-
096 grained misinformation-type prediction in
097 low-resource settings. For a low-capacity
098 MLP classifier, adding persona features im-
099 proves accuracy from 0.79 to 0.81 and macro-
100 F1 from 0.72 to 0.74 ($p = 0.021$).
- 101 • We evaluate persona-conditioned LLM agents
102 in a controlled susceptibility setting and show
103 strong correlations with human persona-level
104 misinformation susceptibility.
- 105 • We show that the relationship between mis-
106 information and engagement persistence is
107 persona-conditional: misinformation ampli-
108 fies or suppresses persistence depending on
109 the dominant conversational persona.

110 For the classification component, our goal is not
111 to establish a new state-of-the-art classifier, but to

¹In this work, we use GPT-5 as our LLM agent.

112 isolate and quantify the contribution of persona-
113 level signal under controlled modeling settings.

2 Related Work 114

2.1 Climate Misinfo/Obstruction Narratives 115

116 Prior work has documented that climate misinfor-
117 mation extends beyond outright denial to include
118 more subtle obstruction strategies, such as ques-
119 tioning the feasibility of solutions, emphasizing
120 economic trade-offs, or attacking the credibility
121 of scientific messengers (Lamb et al., 2020; Coan
122 et al., 2021; Holder et al., 2023). These studies
123 motivate the use of fine-grained misinformation
124 taxonomies rather than binary true/false labels. We
125 adopt this perspective by leveraging the QUOTA-
126 CLIMAT taxonomy (QuotaClimat, 2025), which
127 captures diverse climate obstruction narratives. The
128 fine-grained misinformation types in QUOTACLI-
129 MAT and detailed explanations of the categories are
130 provided in Appendix B.1.

131 However, existing work primarily treats misin-
132 formation as isolated textual instances, without
133 modeling stable speaker-level regularities in how
134 these narratives are deployed. In contrast, our
135 work explicitly models recurring rhetorical patterns
136 through data-driven personas and examines how
137 different personas preferentially express specific
138 misinformation types in real-world discourse.

2.2 Persona/Author context in NLP 139

140 In social NLP tasks such as stance and misinforma-
141 tion detection, text alone is often insufficient as it
142 lacks context about the author’s underlying beliefs.
143 Prior work by Benton and Dredze (2018a) shows
144 that stance expressed in a tweet depends on latent
145 author beliefs, and that incorporating auxiliary au-
146 thor embeddings improves stance classification.

147 A growing body of work shows that incorporat-
148 ing author or persona context improves stance and
149 misinformation detection by capturing latent be-
150 liefs and communicative styles (Benton and Dredze,
151 2018a,b). These approaches demonstrate that who
152 is speaking can provide signal beyond text alone.

153 Our work differs in two key ways. First, rather
154 than relying on predefined author attributes or static
155 user embeddings, we induce personas directly from
156 large-scale, in-the-wild climate discourse. Second,
157 we use persona not only as an auxiliary feature
158 for prediction, but as an organizing abstraction for
159 analyzing misinformation strategies, susceptibility,
160 and engagement dynamics.

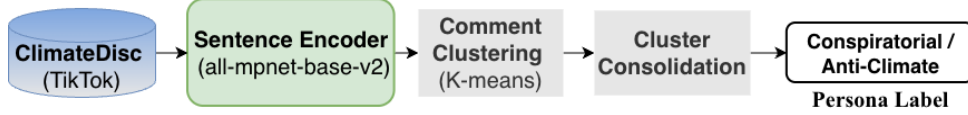


Figure 2: Offline persona induction: encode comments, cluster embeddings, and manually consolidate clusters to assign each comment a persona label p_i (Section 3.1).

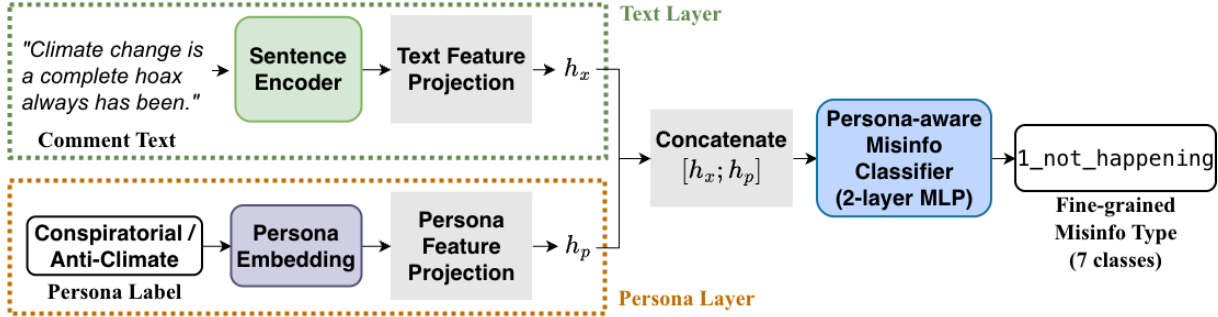


Figure 3: Persona-aware misinformation modeling: obtain fine-grained misinformation-type labels and evaluate prediction models that fuse text and persona representations (Section 3.2). Details of how we obtain comment-level misinformation labels are in Appendix A.

2.3 LLM Personas & Engagement Dynamics

Recent studies show that large language models can be conditioned to simulate behavior associated with specific personas or subpopulations (Park et al., 2023; Krishna et al., 2022; Argyle et al., 2023; Li et al., 2024). We build on this line of work by using persona-conditioned LLM agents as a controlled evaluation tool for climate misinformation susceptibility, rather than as a behavioral simulator per se. We show that persona-level susceptibility patterns in LLM agents align closely with those observed among real human commenters.

Separately, prior work has examined how misinformation diffusion and persistence depend on content properties, emotional framing, and network structure (Vosoughi et al., 2018; Liu et al., 2023; Grinberg et al., 2019; Sun and Xie, 2024). While these studies highlight substantial heterogeneity in engagement, they do not examine whether such dynamics are systematically conditioned on conversational personas. We address this gap by demonstrating persona-dependent persistence effects in climate misinformation discussions on TikTok.

3 Method

Given short-form social media comments discussing climate change, our goal is to characterize and predict fine-grained climate misinformation by explicitly modeling stable speaker-level regularities. For each comment c_i authored by user u_i , we aim to (i) assign a narrative persona p_i capturing re-

curring rhetorical patterns, (ii) obtain a fine-grained misinformation type y_i when applicable, and (iii) analyze how persona conditions misinformation susceptibility and engagement persistence at both the comment and video levels. Dataset details are provided in Section 4.1.

Our pipeline proceeds in two stages. First, we induce narrative personas from large-scale climate comment discourse (Figure 2; Section 3.1). Second, we incorporate persona as an explicit signal for misinformation-type prediction (Figure 3) and for downstream analyses of susceptibility and engagement. Comment-level misinformation labels are obtained using a classifier trained on a curated climate misinformation taxonomy; since the transfer procedure is not central to our main contributions, we defer its details to Appendix A.

3.1 Persona Creation

To capture stable rhetorical and ideological patterns across commenters used in section 3.2, 3.3, 3.4, we induce narrative personas from large-scale climate change comment discourse in CLIMATEDISC (Gao et al., 2025), a large-scale corpus of TikTok climate change videos and associated comments (section 4.1), which serve as a speaker-level abstraction across downstream tasks.

We represent comments as dense semantic embeddings using a pre-trained MPNet-based sentence transformer (Reimers and Gurevych, 2020). We then cluster the embeddings using k -means (Lloyd, 1982), selecting the number of clusters by

| Local / Experiential Observers | Ideological / Political Combatants |
|---|---------------------------------------|
| 1. Local Weather Anecdotes | 0. Anti-Trump / Anti-GOP / Anti-Biden |
| 12. Storms & Flooding | 11. Government Blame / Anarchy |
| 13. Geographical Self-Mentions | 18. Electoral Mobilization |
| Systemic / Structural Critics | Pro-Science / Rationalist |
| 4. Consumption / Plastics Skepticism | 14. Pro-Science / Anti-Misinformation |
| 17. Energy / Policy Debate | |
| 19. Agriculture / Diet | |
| 20. Anti-Corporation / Capitalism | |
| Doomism / Nihilism | Conspiratorial / Anti-Climate |
| 8. "Humans Won't Make It" Doomism | 7. Conspiracy / Muddled Explainers |
| | 16. Denial / Hoax / Manipulation |
| Noise / Low-Signal (Excluded): 2. Emoji / Applause Spam; 3. Minimal Agreement Tokens; 5. Polite Thanks; 6. Short Emotional Bursts; 9. Excited Agreement; 10. Shared Anxiety; 15. @ Mentions. | |

Table 1: Narrative personas obtained by grouping the 21 k -means clusters of comment embeddings into semantically coherent higher-level categories. The Noise / Low-Signal group is excluded from downstream persona-sensitive analyses.

jointly examining Within-Cluster Sum of Squares (WCSS) and Silhouette scores over a candidate range of $k \in [3, 30]$. This analysis identifies $k = 21$ as a practical trade-off between cluster compactness and separation (Section 4.2.1).

To obtain a stable speaker-level abstraction for downstream modeling, we manually validate and consolidate the 21 fine-grained clusters into six higher-level *narrative personas* reflecting shared communicative stances toward climate change, along with a residual "Noise / Low-Signal" group excluded from persona-sensitive analyses. The resulting persona definitions are reported in Table 1.

With personas defined, we next examine whether this abstraction provides predictive signal beyond comment text alone. Specifically, we evaluate whether incorporating persona information improves fine-grained climate misinformation classification.

3.2 Persona-aware Misinfo Classification

We compare a text-only baseline with a persona-aware model to test whether persona provides signal beyond comment text (Figure 3).

Text-only baseline. Model A is a lightweight feed-forward classifier operating over sentence-level semantic representations. Given a comment t in CLIMATEDISC, we first encode it using a pretrained sentence transformer based on MPNet (Reimers and Gurevych, 2020): $\mathbf{x} = \phi(t)$, where $\phi(\cdot)$ maps text to a fixed-dimensional dense vector representation. The encoded representation is then passed through a nonlinear projection, $\mathbf{h}_x = f_x(\mathbf{x})$, where $f_x(\cdot)$ denotes an affine transformation followed by a ReLU nonlinearity and dropout.

Finally, a linear classifier produces logits over fine-grained misinformation types: $\hat{\mathbf{y}} = g(\mathbf{h}_x)$, where $\hat{\mathbf{y}}$ corresponds to a seven-class (excluding 0_not_relevant) misinformation taxonomy from QUOTACLIMAT.

Persona-aware model. Model B extends the text-only baseline by incorporating persona information as an additional latent signal. Each comment is associated with a persona label $c \in \{1, \dots, 6\}$, which is mapped to a learnable persona embedding: $\mathbf{p} = E(c)$, where $E(\cdot)$ denotes a trainable embedding function. The persona embedding is transformed via a nonlinear projection, $\mathbf{h}_p = f_p(\mathbf{p})$, analogous in form to the text projection $f_x(\cdot)$.

The text and persona representations are combined via concatenation,

$$\mathbf{h} = [\mathbf{h}_x; \mathbf{h}_p]$$

and passed through a joint nonlinear transformation, $\mathbf{z} = f_j(\mathbf{h})$, before being mapped to output logits by the same classifier $g(\cdot)$: $\hat{\mathbf{y}} = g(\mathbf{z})$.

This late-fusion formulation isolates the contribution of persona information by introducing \mathbf{p} as an additive latent variable alongside the fixed text representation \mathbf{x} .

We also apply the same formulation using RoBERTa-large representations. We include both lightweight and large pretrained models to examine how the utility of persona varies with text representational capacity, rather than optimize performance.

While this evaluation provides predictive value of persona, it does not address whether the induced personas capture meaningful differences in how misinformation is received or accepted.

3.3 Persona-Level Misinformation Susceptibility with LLM

We use persona-conditioned LLM agents to evaluate whether inferred personas capture systematic differences in misinformation susceptibility. We test whether an LLM prompted with persona descriptions (Appendix C.2) exhibits misinformation acceptance patterns mirroring user comments.

We construct a balanced evaluation set of 350 misinformation claims from the QUOTACLIMAT dataset, sampling 50 claims from each of the seven categories. For each inferred persona, we create a short persona card summarizing characteristic beliefs and rhetorical tendencies (Appendix C.1).

For persona p , we define an LLM susceptibility score as the proportion of misinformation claims not rejected by the persona-conditioned

LLM agent:

$$L_p = \frac{1}{|Q|} \sum_{q \in Q} \mathbb{I}[r_{p,q} \neq \text{reject}],$$

where Q is the set of evaluation claims, $r_{p,q}$ is the LLM’s response for persona p on claim q , and $\mathbb{I}[\cdot]$ denotes an indicator function.

Human–LLM Persona Alignment We assess whether persona-conditioned LLM susceptibility patterns align with those of real human personas. For each persona p , we compute a human misinformation rate H_p , defined as the proportion of comments authored by p that are labeled as misinformation in CLIMATEDISC. We compute the LLM susceptibility score L_p previously defined. We then measure alignment between (H_p) and (L_p) using Pearson and Spearman correlations, where higher correlation indicates closer correspondence between human and LLM persona behavior.

We next examine whether these persona-level differences translate into distinct engagement trajectories.

3.4 Persona-Level Temporal Dynamics

To examine how persona conditions the downstream impact of misinformation, we analyze engagement persistence at the video level. For each video in CLIMATEDISC, we compute (i) the dominant persona among its comments, (ii) the misinformation rate, and (iii) engagement persistence, measured as the time required to reach 60% of total comments (Liu et al., 2023). Our final dataset for this engagement persistence experiment contains 2,925 videos².

Formally, for a video v with N_v total comments, let $t_{v,i}$ denote the timestamp of the i -th comment. We define the cumulative number of comments received by time t as

$$S_v(t) = \sum_{i=1}^{N_v} \mathbb{I}[t_{v,i} \leq t],$$

where $\mathbb{I}[\cdot]$ is the indicator function. Let $t_v^{(0)}$ denote the timestamp of the first observed comment for video v . The time at which the video reaches a fraction α of its total comments is defined as $T_v(\alpha) = \inf\{t : S_v(t) \geq \alpha N_v\}$. We define engagement persistence as $P_v = T_v(0.6) - t_v^{(0)}$, which measures the duration (in hours) required

²We exclude videos whose dominant persona is classified as *Noise / Low-Signal*. We also exclude videos with zero or one comment, since persistence requires at least two comments to be defined.

for video v to accumulate 60% of its total comments. Larger values of P_v indicate more sustained engagement over time.

We test whether the relationship between misinformation and persistence varies by persona by stratifying videos into low, mid, and high misinformation tertiles based on misinformation rate. Within each persona type, we conduct Kruskal–Wallis tests to determine whether persistence differs significantly across misinformation levels.

4 Experiments

4.1 Datasets

QUOTACLIMAT To obtain fine-grained misinformation types, we use QUOTACLIMAT (QuotaClimat, 2025), a curated dataset of 6,091 climate-related statements annotated with one of the misinformation categories specified in Appendix B.1. We use QUOTACLIMAT for two purposes: (i) training a misinformation classifier that is transferred to CLIMATEDISC comments (Appendix A), and (ii) susceptibility probing with persona-conditioned LLM agents using a balanced set of 350 claims (50 per category; Section 3.3).

CLIMATEDISC For persona creation and downstream in-the-wild analyses, we rely on the CLIMATEDISC dataset (Gao et al., 2025), which comprises transcribed English speech from 7,110 TikTok videos and 116,256 corresponding comments tagged with “#climatechange,” collected between January 2024 and November 2024. We use CLIMATEDISC in: (i) persona creation (Section 3.1), (ii) persona-aware misinformation classification on in-the-wild comments (Section 3.2), (iii) human persona-level misinformation rates H_p for LLM alignment (Section 3.3), and (iv) temporal dynamics and engagement persistence (Section 3.4).

For the persona-aware misinformation classification experiments on CLIMATEDISC (Section 3.2), comments are divided into training, validation, and test sets in a 0.8/0.1/0.1 ratio.

4.2 Results

4.2.1 Persona Landscape

Figure 4 reports WCSS and Silhouette scores for $k \in [3, 30]$. WCSS decreases smoothly with k (no decisive elbow), while the Silhouette score attains a local maximum at $k = 21$, indicating comparatively strong separation at this granularity. We therefore set $k = 21$ as a practical trade-off be-

tween separation and interpretability. As a qualitative check, Appendix B.2 lists the ten comments nearest each centroid, which exhibit consistent topical and rhetorical patterns.

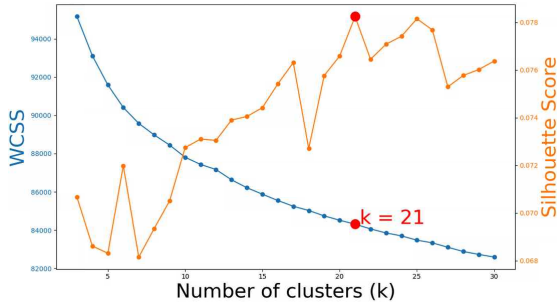


Figure 4: Within-Cluster Sum of Squares (WCSS; lower is better) and Silhouette Score (higher is better) across varying numbers of clusters k . The local Silhouette maximum and WCSS elbow around $k = 21$ motivate our choice of 21 clusters.

For interpretability in downstream analyses, we aggregate the 21 clusters into six higher-level narrative personas plus a residual “Noise / Low-Signal” group. We exclude the Noise group from persona-sensitive analyses (it consists largely of short acknowledgments, emojis, and @-mentions). Table 1 summarizes the mapping.

Figure 5 shows that cluster and persona sizes are relatively balanced: clusters range from 3,798–8,426 comments and personas from 3,992–19,019 comments. This reduces the risk that downstream effects are driven by a single dominant group.

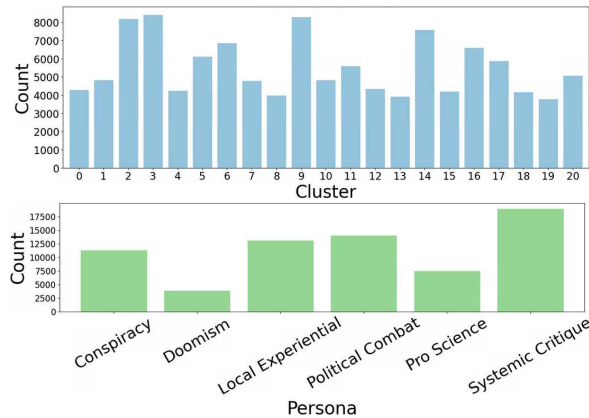


Figure 5: Distribution of cluster/persona sizes.

4.2.2 Text-Based Misinformation Classification with/without Persona

We evaluate whether persona information provides predictive signal beyond text alone (section 3.2).

For the MLP models, incorporating persona improves accuracy from 0.79 to 0.81 and macro-F1

from 0.72 to 0.74. The accuracy gain is statistically significant under an exact paired McNemar test ($p = 0.021$). As shown in Figure 6, the persona-aware model also converges more rapidly in early training, achieving substantially higher validation accuracy and macro-F1 after the first epoch.

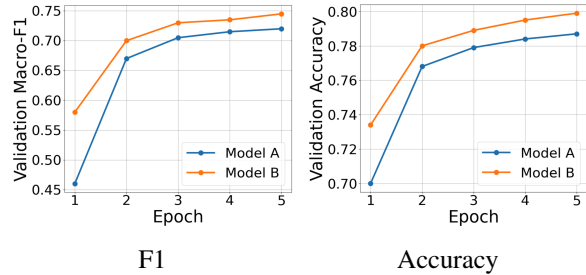


Figure 6: Validation accuracy and F1 across epochs for text-only (A) and persona-aware (B) MLP models.

Using RoBERTa-large, persona yields smaller improvements, increasing accuracy from 0.82 to 0.83 and macro-F1 from 0.76 to 0.78; however, this difference is not statistically significant ($p = 0.11$).

Overall, these results suggest that persona information provides complementary signal in lower-capacity settings, where the text representation alone may be insufficient to capture speaker-level regularities. In contrast, larger pretrained models such as RoBERTa appear able to implicitly encode persona-related cues from text, suggesting a representational ceiling beyond which explicit persona modeling yields diminishing returns.

4.2.3 Persona-Level Susceptibility Probing with LLM

Human vs. LLM Agents Susceptibility We observe a strong positive correlation between human and LLM persona susceptibility patterns. The Pearson correlation between H_p and L_p (specified in section 3.3) is 0.82, indicating substantial linear agreement. The Spearman rank correlation is even higher at 0.94, suggesting that the relative ordering of personas by misinformation susceptibility is highly consistent between humans and LLM agents. This alignment suggests that persona-conditioned LLM agents capture meaningful, persona-specific structure in misinformation acceptance.

Susceptibilities Across Personas Table 2 summarizes overall susceptibility score (specified in section 3.3) per persona. *Conspiratorial / Anti-Climate* shows very high susceptibility (96.3%), while other personas range from 4.9% to 29.4%. Notably, *Ideological / Political Combatants* show

the lowest overall susceptibility (4.9%), suggesting that strong political engagement does not necessarily translate to greater vulnerability to misinformation. Similarly, the Pro-Science / Rationalist persona demonstrates comparatively low susceptibility (10.0%), though this group is not immune to misinformation, indicating that scientific self-identification alone does not guarantee resistance. *Doomism / Nihilism* is comparatively higher (22.0%), highlighting the role of affective orientations such as resignation or fatalism in shaping misinformation vulnerability.

| Persona | Susceptible (%) |
|------------------------------------|-----------------|
| Local / Experiential Observers | 29.4 |
| Ideological / Political Combatants | 4.9 |
| Systemic / Structural Critics | 8.0 |
| Pro-Science / Rationalist | 10.4 |
| Doomism / Nihilism | 22.0 |
| Conspiratorial / Anti-Climate | 96.3 |

Table 2: Overall misinformation susceptibility rate by persona.

4.2.4 Persona-Stratified Engagement Dynamics

We investigate whether misinformation affects content persistence uniformly, or whether this relationship is conditioned on the dominant persona.

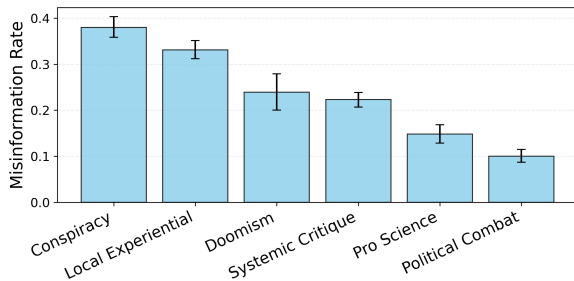


Figure 7: Misinformation rates across persona types (excluding Noise). Bars represent mean proportion of comments containing misinformation for each persona, with error bars showing 95% bootstrap confidence intervals (2,000 iterations).

Misinformation Rates Across Personas Before examining temporal dynamics, we first characterize the baseline misinformation production rates for each persona type. As shown in Figure 7, substantial heterogeneity emerges across persona types. Conspiracy personas exhibit the highest misinformation rate, indicating that approximately 38% of comments from this persona contain misinformation. The Pro Science persona, characterized by explicit appeals to scientific authority, still produces

misinformation in 15% of comments. This finding indicates that scientific self-identification alone does not guarantee resistance to misinformation. The nearly fourfold difference between the highest (Conspiracy) and lowest (Political Combat) personas establishes that communicative style strongly predicts misinformation production.

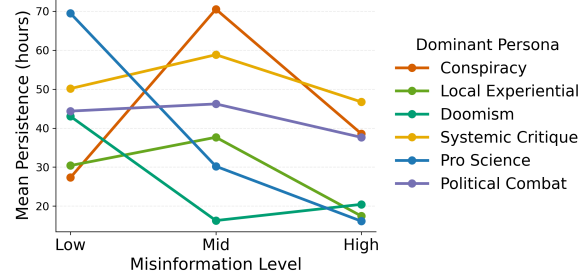


Figure 8: Mean persistence hours across misinformation level (Low, Mid, High) for each persona type.

Persistence Across Personas To examine persona-conditional relationships, we bin videos into three levels (Low, Mid, High) based on misinformation rate and compute mean persistence hours for each persona-by-bin combination. Figure 8 summarizes these persona-conditioned persistence patterns. Pro Science Personas show a strong negative relationship between misinformation and persistence. Videos dominated by Pro Science comments persist substantially longer when misinformation rates are low (69.5 hours) compared to high (16.1 hours). Conspiracy Personas exhibit the opposite pattern. Persistence increases from low-misinformation tertiles (27.3 hours) to high-misinformation tertiles (38.9 hours), even longer persistence at mid-misinformation rate. Systemic Critique Personas exhibit moderate variation across misinformation levels, with persistence remaining within a relatively narrow range compared to other personas. Doomism and Local Experiential Personas show negative relationships similar to Pro Science, with persistence declining at higher misinformation levels. Doomism drops from 43.1 hours (Low) to 20.3 hours (High), a 53% reduction. Local Experiential shows more variability, with an initial increase from Low (30.5 hours) to Mid (37.3 hours) before declining to High (16.8 hours).

To assess whether these observed patterns reflect statistically significant heterogeneity rather than sampling variability, we conduct within-persona Kruskal-Wallis tests examining whether persistence differs across misinformation levels. Table 3

| Persona | H Statistic | p-value |
|--------------------|-------------|------------|
| Systemic Critique | 20.501 | < 0.001*** |
| Local Experiential | 16.336 | < 0.001*** |
| Conspiracy | 9.237 | 0.010** |
| Pro Science | 4.875 | 0.047* |
| Political Combat | 3.850 | 0.146 |
| Doomism | 2.912 | 0.233 |

Table 3: Within-persona Kruskal–Wallis test results examining whether persistence differs across misinformation tertiles (Low, Mid, High). Significance levels are denoted as *** $p < 0.001$, ** $p < 0.01$, and * $p < 0.05$.

presents these results. Four of six personas exhibit statistically significant differences in persistence across misinformation levels at $\alpha = 0.05$. These include Systemic Critique, Local Experiential, Conspiracy, and Pro Science personas, indicating that misinformation prevalence conditions engagement persistence for these communicative styles.

5 Discussion and Conclusion

This work demonstrates that persona constitutes a missing and consequential axis in climate change misinformation research. By modeling who is speaking in addition to what is said, we show that persona provides complementary structure that is not captured by text alone.

The relatively balanced distribution of personas across misinformation types suggests that climate misinformation discourse on TikTok is not dominated by a single narrative or speaker type. This diversity matters for modeling, as it implies that effective detection and intervention strategies must account for multiple recurring rhetorical patterns rather than targeting a single dominant narrative.

Our analysis of the predicted misinformation labels on CLIMATEDISC comments indicates that TikTok climate misinformation remains heavily skewed toward outright denial, but substantial volume also targets the credibility of messengers and promotes anti-solution narratives. These latter strategies focus on delegitimization and motive attacks rather than factual claims, making them more resistant to traditional fact-checking approaches and highlighting the need for alternative mitigation strategies.

Methodologically, our cross-domain matching pipeline from QUOTACLIMAT to CLIMATEDISC provides a practical mechanism for transferring fine-grained climate misinformation taxonomies to large, real-world platforms without requiring extensive human annotation. This scalability is critical for studying evolving misinformation narratives at

platform scale.

Our findings indicate that persona information can aid misinformation classification, but its utility depends on model capacity. In lower-capacity settings, explicitly modeling persona yields consistent and statistically significant gains, suggesting that speaker-level regularities capture information not readily recoverable from individual comments alone. The faster early convergence of persona-aware MLP models further supports the view of persona as a structured inductive bias that helps exploit stable rhetorical patterns rather than sparse lexical cues. In contrast, gains are smaller and not statistically significant for RoBERTa-large, indicating that large pretrained models may already implicitly encode aspects of speaker style, stance, or narrative framing. Together, these results suggest that persona modeling is most beneficial when representational capacity or supervision is limited, offering diminishing returns as models become expressive. Rather than replacing strong text encoders, personas function as a complementary and lightweight inductive bias, particularly valuable in resource-constrained or efficiency-oriented settings.

Additionally, our persona-based LLM simulations demonstrate that such agents can serve as controlled probes of misinformation susceptibility: holding claims constant while varying persona produces systematic differences in acceptance that align with human persona-level vulnerability.

Our temporal dynamics analysis reveals that misinformation-persistence relationships vary substantially by persona type, with four of six personas showing statistically significant effects. These heterogeneous patterns suggest that different persona types are associated with distinct engagement trajectories when misinformation is present, potentially reflecting variation in moderation responses, amplification dynamics, and audience interaction patterns. This has direct implications for platform governance: effective content moderation requires differentiated strategies tailored to specific persona-misinformation combinations rather than uniform approaches across all content types.

Together, our findings suggest that persona-aware modeling offers a principled way to surface speaker-level structure in climate misinformation discourse, unifying prediction, susceptibility analysis, and engagement dynamics within a single abstraction. This perspective complements text-centric approaches and is especially valuable when model capacity or supervision is limited.

607 **Limitations**

608 **Persona Induction Validity and Subjectivity**

609 Our personas are induced via embedding-based k -
610 means clustering and then manually consolidated
611 from 21 clusters into six higher-level types. While
612 this improves interpretability, the resulting taxon-
613 omy is not uniquely determined: different embed-
614 ding models, clustering algorithms, random seeds,
615 or k selection criteria could yield alternative clus-
616 tering results. Further, consolidation relies on re-
617 searcher judgment and may introduce subjectivity
618 in where boundaries are drawn. As such, personas
619 should be interpreted as dataset-specific rhetori-
620 cal groupings rather than definitive or exhaustive
621 categories of users.

622 We infer persona from observed comment text
623 and do not have access to ground-truth user at-
624 tributes (e.g., demographics, ideology, geography,
625 or identity) or off-platform behavior. Consequently,
626 persona labels should not be interpreted as psy-
627 chological traits or real-world subpopulations; they
628 reflect recurring linguistic and rhetorical patterns
629 expressed in public TikTok comments within our
630 collection.

631 **Cross-Domain Misinformation Label Transfer**

632 Our misinformation labels on TikTok comments
633 are produced by a classifier trained on QUOTA-
634 CLIMAT and applied to in-the-wild TikTok dis-
635 course. Although we validate transfer using seman-
636 tic claim–comment matching, this validation relies
637 on a cosine-similarity threshold selected via qual-
638 itative inspection and yields a relatively small set
639 of matched pairs. Moreover, the QUOTACLIMAT
640 taxonomy is claim-centric and may not perfectly
641 align with how misinformation is expressed in con-
642 versational comments (e.g., sarcasm, implicature,
643 partial agreement, or off-topic replies). As a re-
644 sult, predicted misinformation labels likely contain
645 noise, which could affect downstream persona dis-
646 tributions, susceptibility estimates, and persistence
647 analyses.

648 **LLM Agents Susceptibility** Our LLM-based
649 susceptibility results should be interpreted as a
650 controlled probe rather than a behavioral simu-
651 lation of real users. Outputs can be sensitive to
652 prompt wording, decoding settings, and the spe-
653 cific model used, and we evaluate only a single
654 LLM configuration. Furthermore, our metric cap-
655 tures refusal/acceptance patterns under instruction
656 rather than real-world exposure, incentives, or so-

657 cial context. We therefore emphasize alignment in
658 relative persona ordering rather than treating abso-
659 lute LLM susceptibility rates as estimates of human
660 vulnerability.

661 **Temporal Dynamics and Causality** Our per-
662 sistence analysis faces important methodological
663 constraints. First, persistence hours conflate or-
664 ganic sharing velocity, algorithmic promotion, and
665 content moderation decisions, precluding mecha-
666 nistic attribution. Second, our cross-sectional de-
667 sign cannot establish whether misinformation den-
668 sity affects persistence or whether longer-persisting
669 videos accumulate more misinformation over time.
670 Third, video-level aggregation may obscure impor-
671 tant dynamics at finer temporal scales.

672 **Platform and Data Constraints** TikTok’s pro-
673 prietary recommendation algorithm and content
674 moderation policies are opaque, limiting our abil-
675 ity to interpret persistence patterns mechanistically.
676 We cannot observe deleted, shadowbanned, or algo-
677 rithmically suppressed content, potentially biasing
678 our sample toward platform-permitted content.

679 **Ethical Considerations** While our analysis fo-
680 cuses on public comments without identifying in-
681 dividuals, the persona taxonomy could potentially
682 be misused for profiling or targeting. We empha-
683 size that personas describe rhetorical styles, not
684 individual traits, and should not be used to make
685 inferences about individual users’ beliefs or sus-
686 ceptibility.

687 **Acknowledgments**

688 **Personally Identifiable / Offensive Information**

689 We anonymize all potentially personally identifi-
690 able information in our released materials, includ-
691 ing file paths and organization names in the code.
692 In addition, examples of user comments shown in
693 the paper have been sanitized to remove profane or
694 offensive language.

695 **AI Assistants in Writing** We used ChatGPT
696 to assist with language polishing and clarity im-
697 provements throughout the manuscript, to refine
698 the \LaTeX formatting of several figures, and to sup-
699 port code refactoring and readability improvements
700 in the analysis pipeline. The tool was used solely
701 for assistance with writing, formatting, and code
702 refactoring, and did not contribute to the research
703 design, modeling decisions, data analysis, or em-
704 pirical results.

705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759

References

Emily Allaway and Kathleen McKeown. 2021. Zero-shot stance detection: A dataset and model analysis. In *NAACL*.

Lisa Argyle, Ethan Busby, Nancy Fulda, Joshua Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31:1–15.

Clara Baltasar, Sergio D’Antonio Maceiras, Alejandro Martín, and David Camacho. 2024. Analysis of climate change misleading information in tiktok. In *Proceedings of the 1st Workshop on Countering Disinformation with Artificial Intelligence (CODAI), co-located with the 27th European Conference on Artificial Intelligence (ECAI 2024)*, volume 3782 of *CEUR Workshop Proceedings*, pages 54–61, Santiago de Compostela, Spain. CEUR-WS.org. Creative Commons CC BY 4.0.

Corey Basch, Bhavya Yalamanchili, and Joseph Fera. 2021. climate change on tiktok: A content analysis of videos. *Journal of Community Health*, 47:163–167.

Adrian Benton and Mark Dredze. 2018a. Using author embeddings to improve tweet stance classification. In *Proceedings of the 2018 EMNLP Workshop W-NUT: The 4th Workshop on Noisy User-generated Text*, pages 184–194, Brussels, Belgium. Association for Computational Linguistics.

Adrian Benton and Mark Dredze. 2018b. Using author embeddings to improve tweet stance classification. In *Proceedings of the 2018 EMNLP Workshop W-NUT: The 4th Workshop on Noisy User-generated Text*, pages 184–194, Brussels, Belgium. Association for Computational Linguistics.

YeonJung Choi, Lanyu Shang, and Dong Wang. 2023. Climatedist: Climate change misinformation and stance detection dataset.

Travis Coan, Constantine Boussalis, John Cook, and Mirjam Odile Nanko. 2021. Computer-assisted classification of contrarian claims about climate change. *Scientific Reports*, 11.

Carlos Díaz Ruiz and Andreas Nilsson. 2023. Disinformation and echo chambers: Narrative and identity dynamics. *Journal of Marketing*.

Ge Gao, Zhengyang Shan, James Crissman, Ekaterina Novozhilova, YuCheng Huang, Arti Ramanathan, Margrit Betke, and Derry Wijaya. 2025. Insights into climate change narratives: Emotional alignment and engagement analysis on TikTok. In *Proceedings of the Fourth Workshop on NLP for Positive Impact (NLP4PI)*, pages 128–143, Vienna, Austria. Association for Computational Linguistics.

Nir Grinberg, Kenneth Joseph, Lisa Friedland, Briony Swire-Thompson, and David Lazer. 2019. Fake news on twitter during the 2016 u.s. presidential election. *Science*, 363(6425):374–378.

Faye Holder, Sanobar Mirza, Namson-Ngo-Lee, Jake Carbone, and Ruth E. McKie. 2023. Climate obstruction and facebook advertising: how a sample of climate obstruction organizations use social media to disseminate discourses of delay. *Climatic Change*, 176:1–21.

Md Rafiqul Islam, Shaowu Liu, Xianzhi Wang, and Guandong Xu. 2020. Deep learning for misinformation detection on online social networks: a survey and new perspectives. *Social Network Analysis and Mining*, 10.

Ranjay Krishna, Donsuk Lee, Li Fei-Fei, and Michael Bernstein. 2022. Socially situated artificial intelligence enables learning from human interaction. *Proceedings of the National Academy of Sciences*, 119.

William Lamb, Giulio Mattioli, Sebastian Levi, J. Roberts, Stuart Capstick, Felix Creutzig, Jan Minx, Finn Müller-Hansen, Trevor Culhane, and Julia Steinberger. 2020. Discourses of climate delay. *Global Sustainability*, 3.

Junyi Li, Charith Peris, Ninareh Mehrabi, Palash Goyal, Kai-Wei Chang, Aram Galstyan, Richard Zemel, and Rahul Gupta. 2024. The steerability of large language models toward data-driven personas. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7290–7305, Mexico City, Mexico. Association for Computational Linguistics.

Maggie Liu, Jing Wang, and Daniel Preotiuc-Pietro. 2023. Analyzing and predicting persistence of news tweets. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 719–732, Nusa Dua, Bali. Association for Computational Linguistics.

S. Lloyd. 1982. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 28(2):129–137.

S. Mehta and 1 others. 2022. Continually improving social context representations for fake news detection. In *ACL*.

Thanh Tam Nguyen, Ee-Peng Lim, and 1 others. 2020. Fang: Leveraging social context for fake news detection using graph representation. In *CIKM*.

Joon Park, Joseph O’Brien, Carrie Cai, Meredith Morris, Percy Liang, and Michael Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. pages 1–22.

Warren Pearce, Sabine Niederer, Suay Ozkula, and Natalia Sanchez-Querubin. 2018. The social media life of climate change: Platforms, publics, and future imaginaries. *Wiley Interdisciplinary Reviews: Climate Change*, 10.

815 Arianna Pera and Luca Aiello. 2024. [Shifting climates:](#)
816 [Climate change communication from youtube to tik-](#)
817 [tok.](#) pages 376–381.

818 QuotaClimat. 2025. [Frugal ai challenge text train](#)
819 [dataset.](#) Dataset for the Frugal AI Challenge 2025,
820 focused on climate disinformation narratives.

821 Nils Reimers and Iryna Gurevych. 2020. [all-mpnet-](#)
822 [base-v2: Sentence-transformers model.](#) Fine-tuned
823 for sentence embeddings on 1B+ sentence pairs.

824 Cristian Rojas, Frank Algra-Maschio, Mark Andrejevic,
825 Travis Coan, John Cook, and Yuan-Fang Li. 2024.
826 [Hierarchical machine learning models can identify](#)
827 [stimuli of climate change misinformation on social](#)
828 [media.](#) *Communications Earth Environment*, 5.

829 Bebe Chand Sultana, Md Prodhan, Edris Alam, Md So-
830 hel, A. B. M. Mainul Bari, Subodh Pal, Md Islam, and
831 Abu Islam. 2024. [A systematic review of the nexus](#)
832 [between climate change and social media: present](#)
833 [status, trends, and future challenges.](#) *Frontiers in*
834 *Communication*, 9:1301400.

835 Yanqing Sun and Juan Xie. 2024. [Who shares mis-](#)
836 [information on social media? a meta-analysis of](#)
837 [individual traits related to misinformation sharing.](#)
838 *Computers in Human Behavior*, 158:108271.

839 Apoorva Upadhyaya, Marco Fisichella, and Wolfgang
840 Nejd. 2023. [Towards sentiment and temporal aided](#)
841 [stance detection of climate change tweets.](#) *Informa-*
842 *tion Processing Management*, 60:103325.

843 Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018.
844 [The spread of true and false news online.](#) *Science*,
845 359(6380):1146–1151.

846 Y. Zhang and 1 others. 2024. [Modeling user stance with](#)
847 [bipartite graph neural networks.](#) In *ICWSM*.

A Taxonomy Transfer Details

This appendix describes how we obtain fine-grained misinformation-type labels for in-the-wild comments by transferring a curated taxonomy to the CLIMATEDISC corpus.

A.1 Misinformation Classifier Training

We fine-tune a RoBERTa-large classifier on QUOTACLIMAT to predict one of eight climate misinformation categories from text. Training is performed on a single A100 GPU for four epochs with batch size 128 and learning rate 2×10^{-5} . Our model reaches a test accuracy of 0.75 on the eight-class misinformation classification task, indicating that the model is able to reliably distinguish among the fine-grained categories defined in QUOTACLIMAT.

A.2 Cross-domain Validation via Claim–Comment Matching

To validate transfer from curated claims to in-the-wild comments, we embed both QUOTACLIMAT claims and CLIMATEDISC comments using the same sentence transformer and identify high-similarity claim–comment pairs via cosine similarity. Pairs exceeding a fixed threshold of $\tau = 0.85$ are treated as referring to the same underlying claim. We select $\tau = 0.85$ via qualitative inspection: it is the lowest value that consistently yields pairs about the same underlying claim with aligned stance; lower thresholds increasingly admit topic-level matches with stance mismatches. This produces 208 matched comment–claim pairs (Table 4).

On these pairs, the classifier’s prediction on the TikTok comment matches the ground-truth label of the paired QUOTACLIMAT claim 79% of the time (majority baseline 38%), with Cohen’s $\kappa = 0.72$ indicating strong agreement beyond chance (Figure 5). Manual inspection and quantitative agreement metrics further support that matched pairs refer to the same underlying claim.

Based on this validation, we apply the classifier to all comments in CLIMATEDISC to characterize misinformation types in in-the-wild discourse.

To assess label quality beyond the high-similarity regime, where explicit claim alignment is unavailable, we conduct an additional agreement analysis using an independent LLM annotator. We randomly sample 200 comments from CLIMATEDISC that do not exceed the similarity threshold with any QUOTACLIMAT claim. Each comment is annotated independently by ChatGPT using a fixed

| Comment (ClimateDisc) | Matched Claim (Quota-Climat) |
|--|---|
| Climate change is a HOAX | Climate change is a complete hoax always has been. |
| Yes, climate change is real to what extent is arguable, but definitely changing | Climate change is definitely real. It changes every year from winter to autumn to spring to summer. |
| Climate change is very real. The climate has always changed, but in recent years it is changing too fast | Climate change is definitely real. It changes every year from winter to autumn to spring to summer. |

Table 4: Representative matched comment–claim pairs obtained using a cosine-similarity threshold of 0.85.

| Cohen’s Kappa | Quality |
|---------------|--------------------------|
| >0.8 | Almost Perfect Agreement |
| >0.6 | Substantial Agreement |
| >0.4 | Moderate Agreement |
| >0.2 | Fair Agreement |
| $0-0.2$ | Slight Agreement |
| <0 | Almost No Agreement |

Table 5: Interpretation for different ranges of the Cohen’s Kappa values.

prompt (Appendix C.3), and we compare these annotations with our classifier’s predictions.

Across these unmatched comments, agreement reaches 89%, with Cohen’s $\kappa = 0.54$. The substantial raw agreement and moderate κ suggest that the classifier’s predictions remain broadly aligned with independent judgments even when comments cannot be directly grounded to curated claims. Together, these results indicate that the model generalizes beyond claim-level paraphrases to more heterogeneous and implicitly framed misinformation in real-world discourse.

B Datasets

B.1 Climate Misinformation Types

Below are the types of Climate Change Misinformation we used in QUOTACLIMAT.

0_not_relevant: No relevant claim detected or claims that don’t fit other categories.

1_not_happening: Claims denying the occurrence of global warming and its effects - Global warming is not happening. Climate change is NOT leading to melting ice (such as glaciers, sea ice, and permafrost), increased extreme weather, or rising sea levels. Cold weather also shows that climate change is not happening.

923 **2_not_human:** Claims denying human responsi-
924 bility in climate change - Greenhouse gases
925 from humans are not the causing climate
926 change.

927 **3_not_bad:** Claims minimizing or denying nega-
928 tive impacts of climate change - The impacts
929 of climate change will not be bad and might
930 even be beneficial.

931 **4_solutions_harmful_unnecessary:** Claims
932 against climate solutions - Climate solutions
933 are harmful or unnecessary.

934 **5_science_is_unreliable:** Claims questioning cli-
935 mate science validity - Climate science is un-
936 certain, unsound, unreliable, or biased.

937 **6_proponents_biased:** Claims attacking climate
938 scientists and activists - Climate scientists and
939 proponents of climate action are alarmist, bi-
940 ased, wrong, hypocritical, corrupt, and/or po-
941 litically motivated.

942 **7_fossil_fuels_needed:** Claims promoting fossil
943 fuel necessity - We need fossil fuels for eco-
944 nomic growth, prosperity, and to maintain our
945 standard of living.

946 B.2 Representative Comments of Each 947 Cluster

948 Below are the top ten comments closest to each
949 cluster centroid.

950 **Cluster 0** (Partisan Anti-Trump/Anti-GOP
951 Rants):

- 952 1. He's a complete delusional fool, druggy, grapist conspiracy theorist, worm brain
953 lunatic. This POS will never be sworn in. His whole family disowned this crazed
954 weed [clown emoji].
955 2. lol u talkin bout Biden?
956 3. this guy is full of s*** Trump's been helping out more for the American people
957 then Joe Biden and Kamala
958 4. My god Vance just can't answer without lying, there is no way that man should be
959 anywhere near the White House, he's pathetic [skull emoji].
960 5. He's clueless and all his facts. He's the most uneducated person we ever had.
961 6. This man is really ignorant. He doesn't understand what he's even talking about.
962 7. Again, WTF?! He is such a moron!!
963 8. He is part of the problem and can't understand his days are over and his grandchil-
964 dren will suffer for his actions.
965 9. Trump is so stupid
966 10. He is so vile and evil!! Typical conman!

967 **Cluster 1** (Local Weather Anecdotes):

- 968 1. you know you right Cuz. it's unseasonably warm where I live too [laughing emoji].
969 2. This winter in Michigan there was only snow for a few weeks in January and now
970 it's 50 degrees and sunny.
971 3. It was concerning in Michigan yesterday (60s) then it snowed today [laughing
972 emoji].
973 4. I'm in Wisconsin and it was like 90 last week and now it's in the 70s/50s.
974 5. We're having an unseasonably warm winter here in Minnesota.
975 6. we went from rainy to like 65 degrees, what is this summer?
976 7. It's actually scary the climate is warm right now, especially since I am Canadian
977 and not used to this [sweat smile emoji].
978 8. In Utah we hit record lows in November, I'm so over it [crying emoji].
979 9. it's nice right now but Thursday/Friday will be almost 80 in Massachusetts [crying
980 emoji] damn you global warming [victory hand emoji].
981 10. south west area is getting record rain and snow.

982 **Cluster 2** (Emoji/Applause Spam — converted):

1. TRUMP [100 emoji][100 emoji][100 emoji][100 emoji][100 emoji] 983
2. [blue heart emoji][blue heart emoji][blue heart emoji][blue heart emoji] P 984
3. sindio [laughing emoji][laughing emoji][laughing emoji] 985
4. [thumbs up emoji][thumbs up emoji][ok emoji][party emoji] 986
5. [target emoji][target emoji][fire emoji][fire emoji][fire emoji] 987
6. [100 emoji][100 emoji][100 emoji][blue heart emoji][blue heart emoji] 988
7. [screaming emoji] unvorstellbar 989
8. [thumbs up emoji][thumbs up emoji][clap emoji x many] 990
9. [clap emoji][clap emoji][clap emoji][clap emoji] 991
10. [facepalm emoji][facepalm emoji][facepalm emoji] 992

Cluster 3 (Minimal Agreement Tokens):

1. Ya 994
2. ya 995
3. Right 996
4. Right 997
5. RIGHT 998
6. Right 999
7. Right 1000
8. Right 1001
9. Right 1002
10. RIGHT 1003

Cluster 4 (Consumption Plastics / Recycling):

1. clothes or items that can't be recycled get buried or dumped. 1005
2. Stop buying plastic whenever possible. 1006
3. I wish I could buy products not completely covered in plastic. 1007
4. reuse, repurpose, buy less, switch to glass or metal; 2024 no new buying challenge. 1008
5. so much usable stuff goes to the dump. 1009
6. buying second hand clothes to save things from landfill [sweat smile emoji]. 1010
7. buying reusable plastic to replace paper is not a good idea. 1011
8. milk in plastic jugs, bread in plastic bags, lunch meat in plastic [laughing emoji]. 1012
9. all ends up in landfill anyway. 1013
10. I'm sure they don't even recycle any of it. 1014

Cluster 5 (Polite Gratitude Replies):

1. Thank you! 1016
2. Thank you! 1017
3. Thank you! 1018
4. Thank you! 1019
5. Thank you! 1020
6. Thank you! 1021
7. Thank you! 1022
8. Thank you! 1023
9. thank you! 1024
10. Thank you! 1025

Cluster 6 (Short Emotional Reactions):

1. So true [laughing emoji] 1027
2. So cool 1028
3. So scary [crying face emoji] 1029
4. SO TRUE [eyes emoji] 1030
5. So happy [heart emoji] 1031
6. So sad 1032
7. So cute 1033
8. So sad [broken heart emoji] 1034
9. so nice 1035
10. Boosting [blue heart emoji] 1036

Cluster 7 (Explainers + Conspiracies):

1. It's science happens every 10 years and it used to be all water. 1038
2. thermal expansion and added water from land ice. 1039
3. Not a man made phenomenon. 1040
4. land ice and thermal expansion. 1041
5. nothing to do with spraying of skies and geoengineering. 1042
6. water freezes and expands in limestone zones. 1043
7. before this it was acid rain. 1044
8. local event caused by warming Arctic waters. 1045
9. reduced emissions reduced reflective particles, warming the ocean. 1046
10. happens when exiting an ice age. 1047

Cluster 8 (Doomism / Earth Fine, Humans Won't):

1. another great flood, ice age, warming cycle coming. 1050
2. ecological and societal collapse is on the way. 1051
3. the earth will rid itself of us. 1052
4. earth will adapt; we will not survive. 1053
5. Science people, we are doomed. 1054
6. near the environmental point of no return. 1055
7. I'm worried about mankind's future. 1056
8. we will die before it affects earth. 1057
9. perilous times like never before. 1058
10. we must wake up and restore habitability. 1059

Cluster 9 (Excited Agreement Bursts):

1. Omg yes 1061
2. Yessss oh my god 1062

| | | | |
|------|--|---|------|
| 1063 | 3. yesssss | 7. It's weather manipulation. | 1144 |
| 1064 | 4. OMG | 8. climate change is a hoax and scam. | 1145 |
| 1065 | 5. You said it | 9. finally someone tells the truth about climate. | 1146 |
| 1066 | 6. Yesss omg | 10. our idea of climate change is a lie. | 1147 |
| 1067 | 7. OMG | | |
| 1068 | 8. Nailed it | | |
| 1069 | 9. Yes Eva | | |
| 1070 | 10. yes love that | | |
| | Cluster 10 (Shared Anxiety Coping): | Cluster 17 (Energy Policy Debate): | 1148 |
| 1071 | | 1. energy corporations greenwashing rhetoric. | 1149 |
| 1072 | 1. I'm so sad and afraid [blue heart emoji][strong emoji]. | 2. ignoring massive fossil fuel emissions. | 1150 |
| 1073 | 2. we must be aware but also able to live; maybe you're here to bring change. | 3. electric car carbon credits rant. | 1151 |
| 1074 | 3. Say a prayer [pray emoji]. | 4. new refineries would be more efficient. | 1152 |
| 1075 | 4. I can't go on like normal. | 5. switch to renewables and nuclear. | 1153 |
| 1076 | 5. always have hope. | 6. promoting "clean" coal and gas. | 1154 |
| 1077 | 6. I'm coping too. | 7. outsourcing pollution abroad. | 1155 |
| 1078 | 7. same, persist with compassion. | 8. no more fossil fuel projects. | 1156 |
| 1079 | 8. same, I feel terrible. | 9. build nuclear if you care about environment. | 1157 |
| 1080 | 9. I'm scared and sad, I don't want to live in trying times. | 10. greener energy benefits future generations. | 1158 |
| 1081 | 10. Me too, it's too much. | | |
| | Cluster 11 (Government Blame / Anarchy Talk): | Cluster 18 (Electoral Mobilization): | 1159 |
| 1082 | | 1. Republicans are delusional, vote blue. | 1160 |
| 1083 | 1. are you serious? it's our governments. | 2. Vote straight blue ticket. | 1161 |
| 1084 | 2. idiots in charge do nothing. | 3. Vote blue up and down the ballot. | 1162 |
| 1085 | 3. kicking indigenous peoples off lands. | 4. must vote Harris 2024. | 1163 |
| 1086 | 4. governments don't see this as a problem. | 5. Vote Blue. | 1164 |
| 1087 | 5. manipulation of the weak. | 6. Republican voting blue across the board. | 1165 |
| 1088 | 6. HAARP government control. | 7. Vote blue, this is sickening. | 1166 |
| 1089 | 7. this is what anarchy is, people helping each other. | 8. Vote red, blue is crazy. | 1167 |
| 1090 | 8. fight through organizing and political action. | 9. Vote all Democrats. | 1168 |
| 1091 | 9. tell political leaders, not us. | 10. Vote blue across the board. | 1169 |
| 1092 | 10. US government and capitalism accelerated this. | | |
| | Cluster 12 (Storms Flooding): | Cluster 19 (Food / Agriculture Diet): | 1170 |
| 1093 | | 1. need organic and regenerative farming. | 1171 |
| 1094 | 1. back-to-back storms caused flooding. | 2. rainforest deforestation due to animal agriculture. | 1172 |
| 1095 | 2. climatologists warned for decades. | 3. Farming must stop but I need plant-based foods. | 1173 |
| 1096 | 3. storms stronger and more frequent. | 4. raising fuel instead of food. | 1174 |
| 1097 | 4. flash flooding moving north. | 5. plant-based diet is bs. | 1175 |
| 1098 | 5. happened in less than an hour. | 6. want to work in sustainable agriculture. | 1176 |
| 1099 | 6. regular rain causing floods. | 7. kill animals to grow plants. | 1177 |
| 1100 | 7. houses awaiting repairs from hurricanes. | 8. forests cleared for soybean fields. | 1178 |
| 1101 | 8. NYC flooding commentary. | 9. our habits destroy rainforests. | 1179 |
| 1102 | 9. unprecedented Midwest flooding. | 10. plant trees. | 1180 |
| 1103 | 10. every US region has natural disasters. | | |
| | Cluster 13 (Geo-Self Mentions): | Cluster 20 (Anti-Corporation / Anti-Capitalism): | 1181 |
| 1104 | | 1. we cannot count on corporations and the wealthy. | 1182 |
| 1105 | 1. I'm in Tennessee. | 2. fundamental problem with capitalism. | 1183 |
| 1106 | 2. you're in my state now. | 3. corporations keep you poor. | 1184 |
| 1107 | 3. Pennsylvania, same dread. | 4. profit motive stands in the way. | 1185 |
| 1108 | 4. be glad you're not in FL. | 5. corporations get tax breaks. | 1186 |
| 1109 | 5. I live in Visalia. | 6. wealthy keep getting richer. | 1187 |
| 1110 | 6. don't say this, I live there. | 7. someone is buying what corporations sell. | 1188 |
| 1111 | 7. I live here, first time hearing this. | 8. greedflation from corporations. | 1189 |
| 1112 | 8. Southern Illinois. | 9. capitalism is a plague. | 1190 |
| 1113 | 9. here in PA too. | 10. corporations do not care about the people. | 1191 |
| 1114 | 10. I'm in Oregon. | | |
| | Cluster 14 (Pro-Science / Anti-Misinformation): | C LLM Prompts | 1192 |
| 1115 | | C.1 Persona Card Construction | 1193 |
| 1116 | 1. proves you don't know what you're talking about. | To construct persona cards, we prompt ChatGPT | 1194 |
| 1117 | 2. you say peer reviewed studies are wrong. | to synthesize coherent, reusable persona | 1195 |
| 1118 | 3. misinformation spreads easily. | descriptions from clusters of real-world comments. | 1196 |
| 1119 | 4. not okay to think yourself smarter than scientists. | Each cluster is represented by a set of example | 1197 |
| 1120 | 5. will never do their own research. | comments (in Appendix B.2) drawn from our | 1198 |
| 1121 | 6. you trust a convicted felon serial rapist? got it. | dataset, along with a short cluster label | 1199 |
| 1122 | 7. you have no idea but skeptics are dumb? | summarizing its dominant theme. | 1200 |
| 1123 | 8. simple rebuke is always wrong. | The following prompt template is used to generate | 1201 |
| 1124 | 9. these people don't believe in science. | each persona card. The cluster label and | 1202 |
| 1125 | 10. knowledge kills their arguments. | representative comments are injected into the | 1203 |
| | Cluster 15 (@-Mentions): | {cluster_description} field. | 1204 |
| 1126 | | Task: | 1205 |
| 1127 | 1. @sarah lewis | You are given a cluster of social media | 1206 |
| 1128 | 2. @psychicmediumjen | comments that share a common theme. | 1207 |
| 1129 | 3. @Raymel Poche Jr. | Your task is to abstract these comments | 1208 |
| 1130 | 4. @eliza koch | into a single persona description. | 1209 |
| 1131 | 5. @MaKenzie Neas | Instructions: | 1210 |
| 1132 | 6. @H.k.lynne | - Describe this persona's core beliefs | 1211 |
| 1133 | 7. @emily xx | and worldview. | 1212 |
| 1134 | 8. @Rachel Taylor | | |
| 1135 | 9. @aisela | | |
| 1136 | 10. @xottic_daa | | |
| | Cluster 16 (Climate Denial / Manipulation): | | |
| 1137 | | | |
| 1138 | 1. Climate change is a joke. You are brainwashed. | | |
| 1139 | 2. THERE IS NO SUCH THING AS CLIMATE CHANGE!! | | |
| 1140 | 3. f u and climate change it's a f**k** farce; you must be an a**h*** to believe it. | | |
| 1141 | 4. Cut the crap. Nobody believes that mess. | | |
| 1142 | 5. Climate change is a make-believe crisis. | | |
| 1143 | 6. No such thing as climate change. | | |

| | | | |
|------|---|---|------|
| 1213 | - Describe their typical goals or motivations when commenting. | about government inaction. Goals: defend their political side, vent frustration, and rally others to act. | 1278 |
| 1214 | | | 1279 |
| 1215 | - Describe their tone, style, and emotional framing. | Typical phrases: 'Vote blue across the board,' 'These people are so blind,' 'We have to do something,' 'This man is a liar,' 'Wake up America.' | 1280 |
| 1216 | | | 1281 |
| 1217 | - Describe common behaviors, rhetorical patterns, and phrases they use. | | 1282 |
| 1218 | | | 1283 |
| 1219 | - Do not quote individual comments directly. | | 1284 |
| 1220 | | | |
| 1221 | - Write a single coherent persona description suitable for conditioning a language model. | Systemic / Structural Critic | 1285 |
| 1222 | | This persona focuses on systemic issues such as capitalism, corporate greed, consumption habits, fossil fuels, and agriculture. They sound thoughtful but frustrated, often linking environmental problems to larger structural causes ('corporations are destroying the planet for profit,' 'capitalism is the problem,' 'we need to reduce plastic and buy less stuff'). Their tone ranges from earnest to cynical. They often use moral or practical language ('reuse, repurpose, buy quality,' 'the rich keep getting richer,' 'fossil fuel industries caused climate change knowingly'). They like calling out hypocrisy ('greenwashing,' 'profit motive,' 'privatized profits, socialized costs'). They may reference renewable energy or veganism, and often frame solutions as lifestyle and collective change, not individual denial. Common phrases: 'stop buying plastic,' 'corporations dont care,' 'we need systemic change,' 'the planet over profits,' 'no more fossil fuel projects.' | 1286 |
| 1223 | | | 1287 |
| 1224 | Cluster description and example | | 1288 |
| 1225 | comments: | | 1289 |
| 1226 | {cluster_description} | | 1290 |
| 1227 | | | 1291 |
| | C.1.1 Personal Cards | | 1292 |
| | Each persona is represented as a natural-language description capturing characteristic beliefs, goals, tone, and linguistic style. These persona cards are injected verbatim into the LLM prompt to condition responses. | | 1293 |
| 1228 | | | 1294 |
| 1229 | | | 1295 |
| 1230 | | | 1296 |
| 1231 | | | 1297 |
| 1232 | | | 1298 |
| 1233 | Local / Experiential Observer | | 1299 |
| 1234 | This persona speaks from personal experience and local observation, often grounding climate change in everyday life. They mention where they live ('I'm in Wisconsin,' 'here in Florida') and describe weird weather patterns ('78 degrees in Vermont in February,' 'record floods in the Midwest'). Their tone is conversational, regional, and observational, rarely ideological. They might use interjections like 'damn,' 'lol,' or 'it's crazy right now.' They often agree or commiserate with others ('same here,' 'we just had that too,' 'yeah it's scary'). Comments often contain anecdotal or comparative phrasing ('last year it was 20 and now it's 70'). Their goal is to share lived experience, connect with others, and validate observations rather than argue policy. Typical phrases: 'same here in [state],' 'I can't believe this weather,' 'we had floods too,' 'this is so weird for this time of year.' | | 1300 |
| 1235 | | | 1301 |
| 1236 | | | 1302 |
| 1237 | | | 1303 |
| 1238 | | | 1304 |
| 1239 | | | 1305 |
| 1240 | | | 1306 |
| 1241 | | | 1307 |
| 1242 | | | 1308 |
| 1243 | | | 1309 |
| 1244 | | | 1310 |
| 1245 | | | 1311 |
| 1246 | | | 1312 |
| 1247 | | | |
| 1248 | | Pro-Science / Rationalist | 1313 |
| 1249 | | This persona champions evidence, critical thinking, and scientific literacy. They frequently rebut misinformation with confident, logical language ('read the research paper,' 'peer-reviewed studies show,' 'you're not smarter than scientists'). Tone: rational, occasionally sarcastic, and slightly exasperated. They value education, expertise, and data, often reminding others to 'trust science' or mocking anti-science arrogance. Their comments are relatively well-structured, with complete sentences and punctuation. They may sound academic or firm, using terms like 'cognitive dissonance,' 'confirmation bias,' or 'peer-reviewed.' Goal: defend scientific truth and call out ignorance. Common phrases: 'you clearly havent read the study,' 'this isnt opinion, its data,' 'science isnt up for debate,' 'knowledge kills their arguments,' 'lower your confidence level a little.' | 1314 |
| 1250 | | | 1315 |
| 1251 | | | 1316 |
| 1252 | | | 1317 |
| 1253 | | | 1318 |
| 1254 | | | 1319 |
| 1255 | | | 1320 |
| 1256 | | | 1321 |
| 1257 | | | 1322 |
| 1258 | Ideological / Political Combatant | | 1323 |
| 1259 | This persona is fiercely partisan, emotionally charged, and motivated by anger or loyalty toward political figures. They frequently condemn political opponents (especially Trump, MAGA, GOP) or rally support for Democrats ('Vote Blue,' 'Harris 2024'). Their tone is confrontational, moralistic, and full of insults or exclamations ('He's pathetic!!', 'These idiots in charge', 'Wake up!'). They use repetition, capitalization, and emojis (e.g., hearts, 100s) to emphasize points. They often blend outrage with mobilization ('We must vote,' 'Republicans are delusional'). Common behaviors include replying to misinformation with political blame, tagging allies, and expressing outrage | | 1324 |
| 1260 | | | 1325 |
| 1261 | | | 1326 |
| 1262 | | | 1327 |
| 1263 | | | 1328 |
| 1264 | | | 1329 |
| 1265 | | | 1330 |
| 1266 | | | 1331 |
| 1267 | | | 1332 |
| 1268 | | | 1333 |
| 1269 | | | 1334 |
| 1270 | | | 1335 |
| 1271 | | | 1336 |
| 1272 | | | 1337 |
| 1273 | | | 1338 |
| 1274 | | | 1339 |
| 1275 | | Doomist / Nihilist Environmentalist | 1340 |
| 1276 | | This persona expresses despair about humanity's future and sees collapse as inevitable. Tone: emotional, fatalistic, occasionally poetic or | 1341 |
| 1277 | | | 1342 |
| | | | 1343 |
| | | | 1344 |

| | | | |
|------|---|--|------|
| 1345 | philosophical. They believe Earth will | - 1_agree | 1407 |
| 1346 | survive but humans wont ('the earth | - 0_uncertain | 1408 |
| 1347 | will heal, we wont,' 'were doomed,' | - -1_reject | 1409 |
| 1348 | 'mother nature is done with us'). They | | |
| 1349 | use emotive and reflective phrasing | USER: | 1410 |
| 1350 | ('Im scared for the future,' 'its | Claim: {claim} | 1411 |
| 1351 | already too late,' 'the planet will be | Answer: | 1412 |
| 1352 | fine without us'). They may still show | | |
| 1353 | flashes of empathy or communal sorrow | This constrained output format enables direct | 1413 |
| 1354 | ('were all going to die,' 'I just want | quantitative comparison of persona-specific | 1414 |
| 1355 | us all to be okay'). Their goal isnt to | agreement patterns across misinformation | 1415 |
| 1356 | persuade but to lament and emotionally | categories. | 1416 |
| 1357 | process the loss. Typical phrases: | | |
| 1358 | 'were past the point of no return,' | C.3 Misinformation Agreement | 1417 |
| 1359 | 'the earth will rid itself of us,' | The following prompt is used to obtain | 1418 |
| 1360 | 'humans wont survive,' 'its so sad but | independent LLM annotations for assessing | 1419 |
| 1361 | true,' 'were doomed.' | agreement on misinformation labels in comments | 1420 |
| | | without a high-similarity claim match. | 1421 |
| 1362 | Conspiratorial / Anti-Climate This | You are annotating user-generated | 1422 |
| 1363 | persona is skeptical of mainstream | comments for climate change | 1423 |
| 1364 | climate science and frequently frames | misinformation. | 1424 |
| 1365 | climate change as exaggerated, | Each comment must be assigned exactly | 1425 |
| 1366 | fabricated, or caused by hidden | one of the following eight mutually | 1426 |
| 1367 | manipulation. They may claim climate | exclusive categories, based on the | 1427 |
| 1368 | change is a hoax, a scam, or just | primary claim expressed or implied in | 1428 |
| 1369 | natural cycles, and often shift | the comment. | 1429 |
| 1370 | explanations toward weather engineering | Categories: | 1430 |
| 1371 | or conspiracies (for example, 'HAARP,' | 0_not_relevant: No relevant claim | 1431 |
| 1372 | 'geoengineering,' 'chemtrails,' 'cloud | detected or claims that don't fit other | 1432 |
| 1373 | seeding,' or 'they control the | categories | 1433 |
| 1374 | weather'). Their tone is dismissive, | 1_not_happening: Claims denying the | 1434 |
| 1375 | combative, and confident, often using | occurrence of global warming and its | 1435 |
| 1376 | insults or ridicule ('wake up,' | effects (e.g., melting ice, extreme | 1436 |
| 1377 | 'brainwashed,' 'you people are dense'). | weather, sea-level rise) | 1437 |
| 1378 | They may mix partial | 2_not_human: Claims denying human | 1438 |
| 1379 | scientific-sounding terms with | responsibility for climate change | 1439 |
| 1380 | conspiratorial logic (for example, | 3_not_bad: Claims minimizing or denying | 1440 |
| 1381 | 'thermal expansion' alongside 'sky | negative impacts of climate change | 1441 |
| 1382 | spraying'). Their goal is to reject | 4_solutions_harmful_unnecessary: Claims | 1442 |
| 1383 | institutional narratives, challenge | arguing that climate solutions are | 1443 |
| 1384 | authority, and persuade others that the | harmful or unnecessary | 1444 |
| 1385 | mainstream explanation is propaganda or | 5_science_is_unreliable: Claims | 1445 |
| 1386 | control. Typical phrases: 'climate | questioning the validity or reliability | 1446 |
| 1387 | change is a hoax,' 'its just natural | of climate science | 1447 |
| 1388 | cycles,' 'they are manipulating the | 6_proponents_biased: Claims attacking | 1448 |
| 1389 | weather,' 'wake up,' 'do your | climate scientists or climate action | 1449 |
| 1390 | research,' 'follow the money,' 'they | proponents as biased, alarmist, or | 1450 |
| 1391 | want control.' | corrupt | 1451 |
| | | 7_fossil_fuels_needed: Claims asserting | 1452 |
| | | that fossil fuels are necessary for | 1453 |
| | | economic growth or modern living | 1454 |
| | | Instruction: Assign the single category | 1455 |
| | | that best fits the comment. | 1456 |
| | | Output format: One label from | 1457 |
| | | (0_not_relevant, 1_not_happening, | 1458 |
| | | 2_not_human, 3_not_bad, | 1459 |
| | | 4_solutions_harmful_unnecessary, | 1460 |
| | | 5_science_is_unreliable, | 1461 |
| | | 6_proponents_biased, | 1462 |
| | | 7_fossil_fuels_needed) | 1463 |
| | | Comment: {comment text} | 1464 |
| 1392 | C.2 Persona-Conditioned Stance Detection | | |
| 1393 | Prompt | | |
| 1394 | We evaluate whether persona-conditioned LLM | | |
| 1395 | agents systematically agree with different types of | | |
| 1396 | climate misinformation using the following | | |
| 1397 | prompt template. The persona card (Appendix | | |
| 1398 | C.1.1) is injected into the {persona} field and the | | |
| 1399 | misinformation claim into the {claim} field. | | |
| 1400 | Persona: | | |
| 1401 | {persona} | | |
| 1402 | Task: | | |
| 1403 | Would *this persona* agree with, | | |
| 1404 | hesitate about, or reject the following | | |
| 1405 | claim? | | |
| 1406 | Output ONLY one of: | | |