

# The Psychological Effect of Large-Scale Corpus Annotation

Anonymous ACL submission

## Abstract

Large-scale corpus annotation could have a psychological burden on annotators, and the repeated task of coding text could influence annotators' own attitudes and preferences. Existing work is qualitative, correlational, or inconclusive. Here, we use a quasi-experimental design to investigate the causal relationship between prolonged engagement in data labeling and psychological outcomes. We examined how 30 trained annotators who labeled a large dataset were affected across two important psychological dimensions: mental health and Post Traumatic Stress Disorder (PTSD) symptoms. We found no evidence that annotations led to declines in mental health or PTSD symptoms when compared with a non-annotating control group. Bayesian analyses provided additional support for the null effect. Our findings provide recommendations for future annotation efforts, demonstrating that integrating annotator-centric design principles can safeguard against psychological harm in annotation projects.

## 1 Introduction

Annotated datasets are the backbone of supervised machine learning in natural language processing (NLP). Human annotators invest time and effort to label raw text with theoretically defined tags, establishing the “ground truth” for NLP tasks (Kennedy et al., 2022; Geva et al., 2019; Hovy and Lavid, 2010). Specifically, manually annotated corpora are used to train models for applications such as coding linguistic features, sentiment analysis, machine translation, or identifying hate speech on social media platforms (Ljubešić et al., 2023; Pak et al., 2010). Some annotation efforts, often those that require straightforward tasks (e.g., coding word similarity), depend on online crowdsourced workers who receive minimal training or rely on their background knowledge and intuitions to tag text data (Snow et al., 2008). In some cases, however, the process of annotation is complex and sub-

jective (e.g., coding for moral values), making the annotation experience more cognitively and emotionally demanding (Davani et al., 2024; Hoover et al., 2020).

The impact of the annotation process, particularly when it contains emotionally charged or hateful content, on annotators has slowly emerged as a growing focus in the NLP community. Previous research has primarily relied on qualitative interviews and systematic literature reviews to raise awareness about the negative effects of such exposure on annotators (e.g., emotional fatigue, desensitization, and psychological distress, especially in tasks like hate speech detection; (AlEmadi and Zaghouni, 2024; Steiger et al., 2021)). These burdens become pronounced when annotating large-scale psychological corpora, as these tasks demand high levels of cognitive and emotional engagement.

Unlike objective linguistic annotation, coding of text for *psychological* processes often lacks a single “correct” answer. Thus, annotators need to be extensively trained to learn a theoretically-informed typology or an annotation guide to tag text data for its psychological content (Davani et al., 2023). Despite these challenges, such corpora are crucial for diverse applications, including emotion recognition (Poria et al., 2019), moral discourse (Prentiqi et al., 2024), social bias inference (Sap et al., 2020), and building generative language models capable of representing cultural values (Kasirzadeh and Gabriel, 2023). Additionally, psychologically annotated corpora may serve as benchmarks for evaluating the biases of generative language models (Rathje et al., 2024; Kristensen-McLachlan et al., 2025; Abdurahman et al., 2024). However, as large-scale psychological annotation grows more in demand (Wolfe, 2000; Beck et al., 2024), the mental health of annotators emerges as a critical, yet often overlooked, aspect of the data annotation pipeline.

Indeed, there is little experimental research on how the annotation process itself affects annotators'

mental well-being or other psychological attributes (Prabhakaran et al., 2021). To bridge this gap, we investigate the psychological impact of annotations by having Research Assistants (RAs) label the recently compiled Culture and Moral Expressions in Language (CAMEL) Corpus (Anonymous Authors, 2025), a collection of ~51,000 texts, including a wide range of data sources from movie reviews to offensive language datasets (see Table 1 in Appendix). Based on prior work, we focused on two domains: (1) mental health and (2) Post Traumatic Stress Disorder (PTSD) symptoms. A host of other psychological outcomes are also reported in the Appendix. The annotation process was converted into a quasi-experimental design (see Figure 1), with a control group, wherein trained annotators completed self-report measures before and after annotating the corpus. A demographically matched control group also completed the identical measures twice, but did not perform any annotations. In summary, this paper turns the spotlight from the *annotation* to the *annotator*, quantifying the often-overlooked psychological impact of annotations.

## 2 Background

Mental health is characterized by the capacity of individuals to adapt flexibly to stress and adversity, manage emotions effectively, and function well within social contexts (Galderisi et al., 2015). Annotation might impact annotators’ mental well-being by contributing to cognitive fatigue, stress from exposure to sensitive content, and potential emotional strain from prolonged engagement with subjective or ethically complex language (AlEmadi and Zaghouani, 2024; Stoev et al., 2023).

The present experimental set-up is grounded in Self-Determination Theory (STD; (Ryan and Deci, 2017)), a pluralistic theory of motivation and well-being suggesting that individuals have three fundamental psychological needs: autonomy, competence, and relatedness. Satisfying these needs is crucial for enhancing well-being. However, the highly controlled, deskilled, and isolating nature of annotation tasks can hinder the fulfillment of these needs, potentially resulting in a negative impact on well-being. In our study, we designed the annotation process to ensure annotators’ needs were met. Autonomy, the perceived control of one’s actions, could be unmet for annotators due to strict guidelines, such as fixed time limits for each task, rigid hourly quotas, or lack of flexible scheduling

(Sansone and Thoman, 2005). Competence, reflecting an individual’s need to feel efficient (Ryan and Deci, 2017; Kaveti and Akbar, 2020), is undermined by the overly repetitive nature of annotation tasks, offering no performance feedback for skill development. Relatedness, the connection individuals experience when they feel accepted by others (Walton and Cohen, 2007), is impacted as annotators often work in isolation with no community support, making their “invisible labor” concealed behind algorithms (Steiger et al., 2021). The fulfillment of these needs is a robust predictor of mental health (Ng et al., 2012; Slemp et al., 2024).

Annotation practices that make annotators feel undervalued and replaceable may undermine these basic psychological needs and negatively impact mental health. Additionally, annotation work carries the risk of vicarious traumatization. Vicarious Trauma Theory (McCann and Pearlman, 1990) suggests that constant exposure to others’ traumatic experiences can lead to symptoms similar to those of PTSD (known as “the cost of caring” or “compassion fatigue”) (Figley, 2002). Annotators who engage with graphic or hateful material may experience PTSD-like symptoms due to a lack of adequate psychological safeguards. Qualitative interviews (Spence et al., 2023) revealed that moderators frequently encounter trauma-related symptoms, such as intrusive thoughts and sleep disturbances. Crowd annotators typically lack training and receive minimal institutional support, making them particularly vulnerable to vicarious trauma.

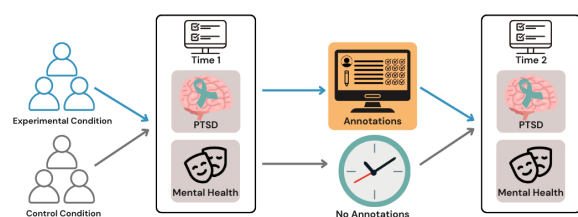


Figure 1: Schematic design of the study

## 3 Experimental Setup

### 3.1 Design

The study employed a quasi-experimental pretest-posttest design with a control group. Individuals in the experimental group annotated a random subset of the CAMEL corpus (see Table 1) over 4 months, whereas the control group did not take part in annotations. A demographically matched control group on gender and race strengthens causal inference by

controlling for potential differences in how mental health is affected across these groups (Wallace et al., 2016; Salk et al., 2017). The pretest-posttest design further allows us to measure change over time. Participants in our experimental group were slightly older, but this difference was minimal and primarily driven by a small subset of older RAs (Appendix C). Both groups were blinded to the study’s aim and were informed that the study examines their general beliefs about the world. All assessments (Appendix A) were conducted twice: (1) at the beginning of the project, before the annotation process began ( $Time_1$ ), and (2) upon completion of the project ( $Time_2$ ). Annotators were RAs in our institution and received research or course credits. Our Institutional Review Board approved the experiment (IRB #5291; Appendix D).

### 3.2 Training of Research Assistants

To prepare RAs for the annotations, we incorporated principles from STD into our training process. We developed a comprehensive annotation manual that included detailed explanations for each label. To enhance competence, the training included thorough manual reviews and small-group discussions, each with five annotators, all supervised by experienced graduate students, to clarify any ambiguities. This structured approach promoted a shared understanding of the guidelines. Following the training, RAs consistently completed between 800 and 1,000 annotations per week. To enhance relatedness, weekly meetings were held with supervisors to address emerging challenges. To support autonomy, we emphasized the importance of independent judgment. Annotators were not permitted to confer on labeling any specific text (Appendix E).

### 3.3 Participants

Two cohorts, Cohort 1 (February – May 2024) and Cohort 2 (August – December 2024), were included to ensure that observed effects were not artifacts of a particular time in the academic year. We observed attrition rates of ~10% in both groups at  $Time_2$  assessments, likely due to students’ limited availability to retake the post-test psychological measures. In the aggregate sample, the experimental group consisted of 30 annotators ( $M_{age} = 20.7$ ,  $SD_{age} = 1.7$ ,  $Male = 6$ ,  $Female = 22$ ,  $Nonbinary = 2$ ;  $White = 14$ ). Annotators met weekly with their supervisor to address any discomfort they may have experienced during the annotation work and had access to a communication channel outside of meet-

ings to reach out if they felt distressed. The control group consisted of 35 undergraduate students ( $M_{age} = 19.9$ ),  $SD_{age} = 1.4$ ,  $Male = 10$ ,  $Female = 23$ ,  $Nonbinary = 2$ ;  $White = 16$ ). Data from Cohorts 1 and 2 were combined to increase statistical power.

### 3.4 Measures & Procedure

Mental health and PTSD symptoms were assessed using psychometrically validated scales (Lukat et al., 2016; Weathers, 2013). Participants in an experimental lab setting (see Figure 1), and completed a 30-minute survey with these scales along with additional measures reported in Appendix.

### 3.5 Statistical Analyses

We conducted a linear mixed model (group as a between-subject factor and time as a within-subject factor) to examine the psychological impact of annotating a large-scale corpus. Equation 1 demonstrates the effect of annotations on outcome variables (mental health or PTSD symptoms) for annotator  $i$  at time  $t$  where  $u_i$  represents the annotator-specific random intercept (accounting for individual variation) and  $\epsilon_{it}$  represents residual. Our key parameter of interest,  $\beta_3$ , estimates the difference in pre–post change in the outcome between the experimental and control groups. It tells us how much the experimental group’s outcome changed over time compared with the control group’s change.

$$\text{Outcome}_{it} = \beta_0 + \beta_1 \text{Group}_i + \beta_2 \text{Time}_t + \beta_3 (\text{Group}_i \times \text{Time}_t) + u_i + \epsilon_{it} \quad (1)$$

We also conducted a Bayesian hierarchical multi-level model using the  $\mathcal{N}(0, 10^2)$  distribution as our priors on fixed effects (Bürkner, 2017; Browne and Draper, 2006). This Bayesian approach allows for estimations of the full posterior distributions of the fixed effects and annotator–level random intercepts, derivation of 95% Credible Intervals (CrI) for each parameter estimate, and comparison of alternative model specifications (e.g., with and without the  $\text{Time} \times \text{Group}$  interaction) using Bayes factors. Notably, Bayes Factors ( $BF$ s) compare the likelihood of the data under the null vs. the alternative, allowing us to quantify evidence for both null and alternative hypotheses, not just a binary “reject” vs. “do not reject” decision (Wagenmakers et al., 2018).

## 4 Results

Linear mixed-effects models revealed no significant treatment effect for mental health, ( $\beta_3 =$

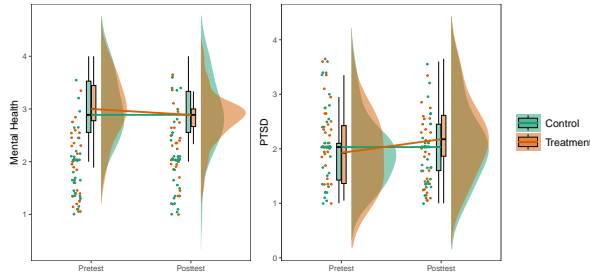


Figure 2: Effect of Group on Changes in Mental Health and PTSD Over Time.

0.11,  $SE = 0.09$ ,  $p = 0.223$ ), or for PTSD, ( $\beta_3 = -0.12$ ,  $SE = 0.16$ ,  $p = .455$ ) (see Figure 2 & Table 2). Bayesian model posterior estimates for the Time $\times$ Group interaction suggested a small positive change in mental health for the treatment group ( $\beta_3 = 0.11$ , 95%CrI [-0.06,0.29]). Bridge-sampling yielded “very strong evidence” for the null hypothesis ( $BF_{01} = 64.0$ ) (Wagenmakers et al., 2018), indicating that treatment and control changed similarly over time. Moreover, posterior estimates for the time $\times$ group interaction depicted a small positive change in PTSD for the treatment group ( $\beta_3 = -0.14$ , 95%CrI [-0.45,0.20]). Bridge-sampling retained “very strong evidence” for the null hypothesis ( $BF_{01} = 45.9$ ) (Wagenmakers et al., 2018), highlighting that treatment and control groups changed similarly (see Table 3). Also, no reliable change in any of the additional psychological outcomes among annotators was seen, as shown in Tables 4 - 7.

## 5 Discussion

Do human annotators pay a psychological price for NLP progress? Overall, our experimental findings offer reassuring evidence that large-scale corpus annotation does not inherently pose a mental health toll on annotators, at least as long as safeguards are embedded in the workflow to meet annotators’ psychological needs. Over two sets of four-month periods of intensive data annotation, we observed no meaningful declines in annotators’ self-reported mental health and PTSD symptoms. Although our focus was on mental health and PTSD symptoms, we also found null effects for other psychological variables, including social biases, political attitudes, cognitive flexibility, and moral values.

Notably, null effects merit cautious interpretation as two features of our protocol likely mitigated potential stressors. First, the corpus was intentionally balanced: emotionally neutral movie-review

snippets were interleaved with aversive content (e.g., hate-speech excerpts). Randomizing text order prevented prolonged exposure to any single content type, plausibly buffering annotators against cumulative negative emotional reactions. Second, our training and annotator-centric procedures were grounded in theories of motivation and well-being, embedding safeguards that proactively address the psychological needs of annotators. Our protocol thus offers practical recommendations demonstrating how mental health costs of large-scale annotation can be avoided by theory-driven safeguards. By embedding these psychological principles into the annotation pipeline, researchers can not only produce high-quality labeled data, but also uphold their ethical responsibility to the annotators.

This approach improves the transparency of datasets and also allows for further analysis of individual biases and mental health indicators (Prabhakaran et al., 2021). Moreover, we recommend that researchers incorporate steps (e.g., an iterative feedback loop in which annotators can report the psychological impact of their work) to minimize negative effects (Kennedy et al., 2022). Large Language Models (LLMs) are another method of remediation, increasingly taking on annotator roles alongside annotators (Holter and Ell, 2023) or independently (Rathje et al., 2024; Tan et al., 2024). However, researchers should ensure the validation of annotations performed by LLMs, as they may not capture the subjectiveness of certain psychological constructs (Abdurahman et al., 2024). As the field continues to build classifiers and fine-tune existing models using annotated data, it is essential to keep in mind that the data (Atari et al., 2023b) and how annotators interact with it (Davani et al., 2023) shape the quality of machine learning models.

## 6 Limitations

Our study focused on the psychological impact of large-scale annotations on annotators; however, we did not assess the extent to which annotators engaged with the annotation tasks or how accurate their annotations were. Although we implemented rigorous protocols, we did not directly measure the level of attention given to the annotation process. Future research could incorporate additional measures (e.g., sustained attention checks, task persistence, and response time analysis) (Langner and Eickhoff, 2013) to capture the extent to which annotators efficiently annotated the corpora.

Our annotation corpus also spanned a broad spectrum of texts ranging from often-neutral movie reviews to highly offensive hate-speech excerpts, which contrasts with the decontextualized, token-based annotations that many NLP projects employ (e.g., part-of-speech tagging, syntactic parsing). Importantly, we observed no reliable decline in annotators' mental health when processing such emotionally taxing and cognitively demanding content. Given these results, it is plausible that annotating less emotionally taxing, more linguistically focused corpora would similarly have a minimal psychological impact. Future research should directly compare these different annotation tasks to confirm that lighter workflows remain manageable using similar protocols. Moreover, unlike corpora focused solely on offensive-language detection (Albanyan et al., 2023; Albanyan and Blanco, 2022; Toraman et al., 2022), our design combined neutral texts with offensive and hateful language excerpts; while offensive-language corpora might elicit greater distress, our findings suggest that buffering emotionally heavy content with neutral snippets could potentially mitigate the psychological toll.

Additionally, our study employed a quasi-experimental design with a total sample of 65 participants. While this relatively small sample size limits statistical power and generalizability, we were still able to detect small to medium effect sizes, underscoring the relevance of our findings. Additionally, our sample consisted exclusively of English-speaking undergraduate students from a large public university in the U.S., which may not fully represent the more diverse populations typically engaged in online annotation platforms (Henrich et al., 2010; Park et al., 2025; Nguyen et al., 2020; Lee et al., 2024; Hovy and Prabhumoye,

2021). Moreover, there is a potential selection bias in the experimental group due to non-random assignment (Joyce et al., 2017). For instance, individuals who volunteered as research assistants may have been more hard-working, ambitious, or interested in culture and computational social science, compared with the control group, which we recruited from the subject pool in our department. Future research should address these limitations by implementing true randomization to enhance the external validity of these findings.

We also did not have sufficient statistical power to examine treatment effects within specific subgroups (e.g., we only had two nonbinary individuals in our annotators). For instance, it is possible that individuals from minority demographic or cultural backgrounds may be especially vulnerable to the mental health toll of large-scale annotation tasks. Future studies with larger annotator samples would allow for more robust subgroup analyses and a better understanding of differential effects across populations.

Lastly, our study focused on a short-term evaluation of the psychological effects of the annotation task, without assessing long-term outcomes. This limitation prevents us from determining whether these (null) effects persist beyond task completion. Future research should expand on these findings to investigate the durability of these psychological outcomes over time to better understand their lasting impact.

## References

- Suhaib Abdurahman, Mohammad Atari, Farzan Karimi-Malekabadi, Mona J Xue, Jackson Trager, Peter S Park, Prezi Golazizian, Ali Omrani, and Morteza Dehghani. 2024. *Perils and opportunities in using large language models in psychological research*. *PNAS nexus*, 3(7):pgae245.
- CJ Adams, Jeffrey Sorensen, Julia Elliott, Lucas Dixon, Mark McDonald, nithum, and Will Cukierski. 2017. Toxic comment classification challenge. <https://kaggle.com/competitions/jigsaw-toxic-comment-classification-challenge>. Kaggle.
- Senem Aktas. 2023. *Nostalgic sentiment analysis of youtube comments dataset*.
- Abdullah Albanyan and Eduardo Blanco. 2022. *Pinpointing fine-grained relationships between hateful tweets and replies*. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):10418–10426.

449	Abdullah Albanyan, Ahmed Hassan, and Eduardo Blanco. 2023. <a href="#">Not all counterhate tweets elicit the same replies: A fine-grained analysis</a> . In <i>Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)</i> , pages 71–88, Toronto, Canada. Association for Computational Linguistics.	501
450		502
451		503
452		504
453		
454		505
455		506
456	Maryam M. AlEmadi and Wajdi Zaghouni. 2024. <a href="#">Emotional toll and coping strategies: Navigating the effects of annotating hate speech data</a> . In <i>Proceedings of the Workshop on Legal and Ethical Issues in Human Language Technologies @ LREC-COLING 2024</i> , pages 66–72, Torino, Italia. ELRA and ICCL.	507
457		508
458		509
459		510
460		
461		511
462	Bob Altemeyer. 1996. <i>The Authoritarian Specter</i> . Harvard University Press.	512
463		513
464		514
465		515
466	Anonymous Authors. 2025. The cultural and moral expressions in language (camel) corpus. Manuscript under review.	
467	Mohammad Atari, Jesse Graham, and Morteza Dehghani. 2020. <a href="#">Foundations of morality in iran</a> . <i>Evolution and Human Behavior</i> , 41(5):367–384.	516
468		517
469		518
470	Mohammad Atari, Jonathan Haidt, Jesse Graham, Sena Koleva, Sean T Stevens, and Morteza Dehghani. 2023a. <a href="#">Morality beyond the weird: How the nomological network of morality varies across cultures</a> . <i>Journal of Personality and Social Psychology</i> .	519
471		
472		520
473		521
474		522
475	Mohammad Atari, Mona J Xue, Peter S Park, Damián E Blasi, and Joseph Henrich. 2023b. <a href="#">Which humans?</a>	523
476		524
477	Mohammad Atari and Aliah Zewail. 2025. <a href="#">Morality across time and space</a> . <i>European Review of Social Psychology</i> , pages 1–36.	525
478		526
479		527
480	Jordan R. Axt. 2018. <a href="#">The best way to measure explicit racial attitudes is to ask about them</a> . <i>Social Psychological and Personality Science</i> , 9(8):896–906.	528
481		
482		529
483	Mahzarin R. Banaji and Anthony G. Greenwald. 2013. <i>Blindspot: Hidden Biases of Good People</i> . Delacorte Press, New York, NY.	530
484		531
485		532
486	Jacob Beck, Stephanie Eckman, Bolei Ma, Rob Chew, and Frauke Kreuter. 2024. <a href="#">Order effects in annotation tasks: Further evidence of annotation sensitivity</a> . In <i>Proceedings of the 1st Workshop on Uncertainty-Aware NLP (UncertainNLP 2024)</i> , pages 81–86, St Julians, Malta. Association for Computational Linguistics.	533
487		534
488		535
489		536
490		
491		537
492		538
493	Matt Boraske. 2022. <a href="#">Aita subreddit dataset</a> .	539
494		540
495	William J. Browne and David Draper. 2006. <a href="#">A comparison of bayesian and likelihood-based methods for fitting multilevel models</a> . <i>Bayesian Analysis</i> , 1(3):473–514.	541
496		542
497		543
498		544
499	Paul-Christian Bürkner. 2017. <a href="#">brms: An r package for bayesian multilevel models using stan</a> . <i>Journal of statistical software</i> , 80:1–28.	545
500		546
		547
		548
		549
		550
		551
		552
		553
		554
		555

556	Jesse Graham, Jonathan Haidt, Sena Koleva, Matt Motyl, Ravi Iyer, Sean P Wojcik, and Peter H Ditto. 2013. <a href="#">Moral foundations theory: The pragmatic validity of moral pluralism</a> . In <i>Advances in experimental social psychology</i> , volume 47, pages 55–130. Elsevier.	609
557		610
558		611
559		612
560		613
561		
562	Anthony G. Greenwald, Debbie E. McGhee, and Jordan L. K. Schwartz. 1998. <a href="#">Measuring individual differences in implicit cognition: The implicit association test</a> . <i>Journal of Personality and Social Psychology</i> , 74(6):1464.	614
563		615
564		616
565		617
566		
567	Social Grep. 2022. <a href="#">The 2022 trucker strike dataset</a> .	618
568		619
569	James J Gross and Oliver P John. 2003. <a href="#">Individual differences in two emotion regulation processes: implications for affect, relationships, and well-being</a> . <i>Journal of personality and social psychology</i> , 85(2):348.	620
570		621
571		622
572		623
573	James J. Gross and Oliver P. John. 2012. <a href="#">Emotion regulation questionnaire</a> . <i>Journal of Personality and Social Psychology</i> .	624
574		625
575		626
576	Christian Haerpfer, Ronald Inglehart, Alberto Moreno, Christian Welzel, Katerina Kizilova, Juan Diez-Medrano, and Birgitta Puranen. 2020. World values survey wave 7 (2017-2020) cross-national data-set. Retrieved from <a href="https://www.worldvaluessurvey.org/">https://www.worldvaluessurvey.org/</a> .	627
577		628
578		629
579		630
580		631
581	Jonathan Haidt and Craig Joseph. 2004. <a href="#">Intuitive ethics: How innately prepared intuitions generate culturally variable virtues</a> . <i>Daedalus</i> , 133(4):55–66.	632
582		633
583		634
584	Joseph Henrich, Steven J Heine, and Ara Norenzayan. 2010. <a href="#">The weirdest people in the world?</a> <i>Behavioral and brain sciences</i> , 33(2-3):61–83.	635
585		636
586		637
587	Ole Magnus Holter and Basil Ell. 2023. <a href="#">Human-machine collaborative annotation: A case study with GPT-3</a> . In <i>Proceedings of the 4th Conference on Language, Data and Knowledge</i> , pages 193–206, Vienna, Austria. NOVA CLUNL, Portugal.	638
588		639
589		640
590		641
591		642
592	Joe Hoover, Gwenyth Portillo-Wightman, Leigh Yeh, Shreya Havaladar, Aida Mostafazadeh Davani, Ying Lin, Brendan Kennedy, Mohammad Atari, Zahra Kamel, Madelyn Mendlen, Gabriela Moreno, Christina Park, Tingyee E. Chang, Jenna Chin, Christian Leong, Jun Yen Leung, Arineh Mirinjian, and Morteza Dehghani. 2020. <a href="#">Moral foundations twitter corpus: A collection of 35k tweets annotated for moral sentiment</a> . <i>Social Psychological and Personality Science</i> , 11(8):1057–1071.	643
593		644
594		645
595		646
596		647
597		648
598		649
599		650
600		651
601		652
602	Dirk Hovy and Shrimai Prabhumoye. 2021. <a href="#">Five sources of bias in natural language processing</a> . <i>Language and Linguistics Compass</i> , 15(8):e12432.	653
603		654
604		655
605	Eduard Hovy and Julia Lavid. 2010. <a href="#">Towards a ‘science’ of corpus annotation: a new methodological challenge for corpus linguistics</a> . <i>International journal of translation</i> , 22(1):13–36.	656
606		657
607		658
608		659
	Ted Joyce, Dahlia K. Remler, David A. Jaeger, Onur Altindag, Stephen D. O’Connell, and Sean Crockett. 2017. <a href="#">On measuring and reducing selection bias with a quasi-doubly randomized preference trial</a> . <i>Journal of Policy Analysis and Management</i> , 36(2):438–459.	660
		661
		662
		663
	Atoosa Kasirzadeh and Iason Gabriel. 2023. <a href="#">In conversation with artificial intelligence: aligning language models with human values</a> . <i>Philosophy &amp; Technology</i> , 36(2):27.	664
		665
	Pushyami Kaveti and Md Navid Akbar. 2020. <a href="#">Role of intrinsic motivation in user interface design to enhance worker performance in amazon mturk</a> . In <i>Proceedings of the 13th ACM International Conference on Pervasive Technologies Related to Assistive Environments</i> , pages 1–7.	666
		667
	B. Kennedy, M. Atari, A. M. Davani, L. Yeh, A. Omrani, Y. Kim, and M. Dehghani. 2022. <a href="#">Introducing the gab hate corpus: defining and applying hate-based rhetoric to social media posts at scale</a> . <i>Language Resources and Evaluation</i> , pages 1–30.	668
		669
	M. J. Kirton. 1981. <a href="#">A reanalysis of two scales of tolerance of ambiguity</a> . <i>Journal of Personality Assessment</i> , 45(4):407–414.	670
		671
	Ross Deans Kristensen-McLachlan, Miceal Canavan, Marton Kárdos, Mia Jacobsen, and Lene Aarøe. 2025. <a href="#">Are chatbots reliable text annotators? sometimes</a> . <i>PNAS nexus</i> , 4(4):pgaf069.	672
		673
	Elizabeth J. Krumrei-Mancuso and Stephen V. Rouse. 2016. <a href="#">The development and validation of the comprehensive intellectual humility scale</a> . <i>Journal of Personality Assessment</i> , 98:209–221.	674
		675
	Nour Kteily, Arnold K Ho, and Jim Sidanius. 2012. <a href="#">Hierarchy in the mind: The predictive power of social dominance orientation across social contexts and domains</a> . <i>Journal of Experimental Social Psychology</i> , 48(2):543–549.	676
		677
	Robert Langner and Simon B Eickhoff. 2013. <a href="#">Sustaining attention to simple tasks: a meta-analytic review of the neural mechanisms of vigilant attention</a> . <i>Psychological bulletin</i> , 139(4):870.	678
		679
	Mark R Leary, Kate J Diebels, Erin K Davisson, Katrina P Jongman-Sereno, Jennifer C Isherwood, Kaitlin T Raimi, Samantha A Deffler, and Rick H Hoyle. 2017. <a href="#">Cognitive and interpersonal features of intellectual humility</a> . <i>Personality and Social Psychology Bulletin</i> , 43(6):793–813.	680
		681
	Nayeon Lee, Chani Jung, Junho Myung, Jiho Jin, Jose Camacho-Collados, Juho Kim, and Alice Oh. 2024. <a href="#">Exploring cross-cultural differences in english hate speech annotations: From dataset construction to analysis</a> . <i>Preprint</i> , arXiv:2308.16705.	682
		683
	N. Ljubešić, I. Mozetič, and P. K. Novak. 2023. <a href="#">Quantifying the impact of context on the quality of manual hate speech annotation</a> . <i>Natural Language Engineering</i> , 29(6):1481–1494.	684
		685



771	Theodore M Singelis, Harry C Triandis, Dharm PS Bhawuk, and Michele J Gelfand. 1995. <a href="#">Horizontal and vertical dimensions of individualism and collectivism: A theoretical and measurement refinement.</a> <i>Cross-cultural research</i> , 29(3):240–275.	826
772		827
773		828
774		829
775		830
776	Gavin R Slemp, James G Field, Richard M Ryan, Vivien W Forner, Anja Van den Broeck, and Kelsey J Lewis. 2024. <a href="#">Interpersonal supports for basic psychological needs and their relations with motivation, well-being, and performance: A meta-analysis.</a> <i>Journal of Personality and Social Psychology</i> .	831
777		832
778		833
779		834
780		835
781		836
782	Rion Snow, Brendan O’connor, Dan Jurafsky, and Andrew Y Ng. 2008. <a href="#">Cheap and fast—but is it good? evaluating non-expert annotations for natural language tasks.</a> In <i>Proceedings of the 2008 conference on empirical methods in natural language processing</i> , pages 254–263.	837
783		838
784		839
785		840
786		841
787		842
788	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. <a href="#">Recursive deep models for semantic compositionality over a sentiment treebank.</a> In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	843
789		844
790		845
791		846
792		847
793		848
794		849
795		850
796	Ruth Spence, Antonia Bifulco, Paula Bradbury, Elena Martellozzo, and Jeffrey DeMarco. 2023. <a href="#">The psychological impacts of content moderation on content moderators: A qualitative study.</a> <i>Cyberpsychology: Journal of Psychosocial Research on Cyberspace</i> , 17(4).	851
797		852
798		853
799		854
800		855
801		856
802	Miriah Steiger, Timir J Bharucha, Sukrit Venkatagiri, Martin J Riedl, and Matthew Lease. 2021. <a href="#">The psychological well-being of content moderators: the emotional labor of commercial moderation and avenues for improving support.</a> In <i>Proceedings of the 2021 CHI conference on human factors in computing systems</i> , pages 1–14.	857
803		858
804		859
805		860
806		861
807		862
808		863
809	Teodor Stoev, Kristina Yordanova, and Emma L. Tonkin. 2023. <a href="#">Experiencing annotation: Emotion, motivation and bias in annotation tasks.</a> In <i>2023 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom Workshops)</i> , pages 534–539.	864
810		865
811		866
812		867
813		868
814		869
815	Zhen Tan, Dawei Li, Song Wang, Alimohammad Beigi, Bohan Jiang, Amrita Bhattacharjee, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. 2024. <a href="#">Large language models for data annotation and synthesis: A survey.</a> <i>Preprint</i> , arXiv:2402.13446.	870
816		871
817		872
818		873
819		874
820	Cagri Toraman, Furkan Şahinuç, and Eyup Yilmaz. 2022. <a href="#">Large-scale hate speech detection with cross-domain transfer.</a> In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 2215–2225, Marseille, France. European Language Resources Association.	875
821		876
822		877
823		
824		
825		
	Jackson Trager, Alireza S. Ziabari, Aida Mostafazadeh Davani, Preni Golazizian, Farzan Karimi-Malekabadi, Ali Omrani, Zhihe Li, Brendan Kennedy, Nils Karl Reimer, Melissa Reyes, Kelsey Cheng, Mellow Wei, Christina Merrifield, Arta Khosravi, Evans Alvarez, and Morteza Dehghani. 2022. <a href="#">The moral foundations reddit corpus.</a> <i>Preprint</i> , arXiv:2208.05545.	
	Harry C. Triandis and Michele J. Gelfand. 1998. <a href="#">Converging measurement of horizontal and vertical individualism and collectivism.</a> <i>Journal of Personality and Social Psychology</i> , 74(1):118–128.	
	Harry C Triandis and Michele J Gelfand. 2012. <a href="#">A theory of individualism and collectivism.</a> <i>Handbook of theories of social psychology</i> , 2.	
	Ayşe K Uskul, Susan E Cross, and Ceren Günsoy. 2023. <a href="#">The role of honour in interpersonal, intrapersonal and intergroup processes.</a> <i>Social and Personality Psychology Compass</i> , 17(1):e12719.	
	Eric-Jan Wagenmakers, Jonathon Love, Maarten Marsman, Tahira Jamil, Alexander Ly, Josine Verhagen, Ravi Selker, Quentin F Gronau, Damian Dropmann, Bruno Boutin, et al. 2018. <a href="#">Bayesian inference for psychology. part ii: Example applications with jasp.</a> <i>Psychonomic bulletin &amp; review</i> , 25:58–76.	
	Stephanie Wallace, James Nazroo, and Laia Bécáres. 2016. <a href="#">Cumulative effect of racial discrimination on the mental health of ethnic minorities in the united kingdom.</a> <i>American journal of public health</i> , 106(7):1294–1300.	
	Gregory M Walton and Geoffrey L Cohen. 2007. <a href="#">A question of belonging: race, social fit, and achievement.</a> <i>Journal of personality and social psychology</i> , 92(1):82.	
	Frank W. Weathers. 2013. <a href="#">The ptsd checklist for dsm-5 (pcl-5).</a>	
	Joanna L. Wolfe. 2000. <a href="#">Effects of annotations on student readers and writers.</a> In <i>Proceedings of the Fifth ACM Conference on Digital Libraries</i> , pages 19–26.	
	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2016. <a href="#">Character-level convolutional networks for text classification.</a> <i>Preprint</i> , arXiv:1509.01626.	
	<b>Appendix</b>	
	<b>A Primary Measures</b>	
	<b>A.1 Mental Health</b>	
	<b>Mental Health.</b> Mental health is characterized by individuals’ general emotional, psychological, and social well-being (Lukat et al., 2016). We measured mental health using nine items on a 5-point scale anchored by “Not True” (1) to “True” (5). An example item is, “I feel that I am actually well equipped to deal with life and its difficulties.”	

878 Responses to these items were averaged to create a  
879 composite score (Cronbach’s  $\alpha = 0.90$ ).

880 **Post Traumatic Stress Disorder.** We used  
881 the 20-item PTSD Checklist for DSM-5 (PCL-5)  
882 (Weathers, 2013) to assess symptoms of anxiety  
883 resulting from stressful or distressing events, com-  
884 monly associated with PTSD. An example of the  
885 item is “Having strong negative feelings such as  
886 fear, horror, anger, guilt, or shame.” All items were  
887 rated on a 5-point scale ranging from “Not at All”  
888 (1) to “A great deal” (5). Responses to these items  
889 were averaged to create a composite score (Cron-  
890 bach’s  $\alpha = 0.94$ ).

891 **B Additional Measures**

892 In addition to our primary focus on mental  
893 health and PTSD-related symptoms, we included a  
894 broader set of psychological measures to explore  
895 whether large-scale annotation tasks might subtly  
896 influence other aspects of cognition, identity, or  
897 social perception. We also measured these vari-  
898 ables because we intended to collect a variety of  
899 individual-difference measures from annotators,  
900 inspired by prior recommendations (Prabhakaran  
901 et al., 2021). These measures were organized into  
902 four superordinate domains, each targeting a the-  
903 oretically meaningful dimension of psychological  
904 functioning.

905 First, we assessed psychological flexibility, a col-  
906 lection of constructs with significant importance  
907 in clinical psychology, reflecting an individual’s  
908 ability to adapt to difficult thoughts and feelings  
909 while maintaining goal-directed behavior. Second,  
910 we included a battery of measures tapping into so-  
911 cial biases. These included both explicit self-report  
912 items and implicit tests designed to capture prefer-  
913 ences toward particular racial and gender groups.  
914 Third, to explore how participants relate to cultural  
915 norms and self-construal, we measured cultural ori-  
916 entation. Finally, we assessed moral values. These  
917 values can shape how annotators interpret morally  
918 charged language, and may themselves be influ-  
919 enced by prolonged exposure to polarized or norm-  
920 violating content. Together, these four domains  
921 allowed us to probe whether the annotation experi-  
922 ence has downstream psychological effects beyond  
923 immediate well-being and PTSD symptoms.

924 **B.1 Psychological Flexibility**

925 **Tolerance of Ambiguity.** This construct refers to  
926 the degree to which individuals prefer ambiguity

927 over familiarity (Kirton, 1981). We measured toler-  
928 ance of ambiguity using 13 items on a 5-point scale,  
929 ranging from “Strongly Disagree” (1) to “Strongly  
930 Agree” (5). An example item is, “I generally prefer  
931 novelty over familiarity.” Responses were averaged  
932 to create a composite score (Cronbach’s  $\alpha = 0.77$ ).

933 **Emotion Regulation.** This construct refers to  
934 the ability of individuals to express and modulate  
935 both and negative emotions (Gross and John, 2003).  
936 We used a 10-item scale to measure emotion reg-  
937 ulation (Gross and John, 2012). An example of  
938 the item is “I control my emotions by changing the  
939 way I think about the situation I’m in.” All items  
940 were rated on a 7-point scale anchored by “Strongly  
941 Disagree” (1) to “Strongly Agree” (7). Responses  
942 to these items were averaged to create a total score  
943 (Cronbach’s  $\alpha = 0.79$ ).

944 **Perspective Taking.** This construct refers to the  
945 tendency of individuals to understand the thoughts,  
946 feelings, and motivations of others in interpersonal  
947 situations (Davis, 1983). We used a 14-item scale  
948 to measure perspective taking. An example of the  
949 item is “Before criticizing somebody, I try to imag-  
950 ine how I would feel if I were in their place.” All  
951 items were rated on a 5-point scale anchored by  
952 “Does not describe me well” (1) to “Describes me  
953 well” (5). Responses to these items were averaged  
954 to create a composite score (Cronbach’s  $\alpha = 0.85$ ).

955 **B.2 Social Biases**

956 **Implicit gender and racial biases.** Implicit bi-  
957 ases are automatic, unconscious associations that  
958 influence human attitudes and behavior (Banaji and  
959 Greenwald, 2013). We used a Science–Gender Im-  
960 plicit Association Test to assess the automatic asso-  
961 ciation between science (e.g., “hypothesis,” “data”)  
962 versus liberal-arts concepts (e.g., “essay,” “gram-  
963 mar”) and male (“he,” “male,” “boy”) versus fe-  
964 male pronouns (“she,” “woman,” “girl”) (Green-  
965 wald et al., 1998). After three practice rounds,  
966 participants completed four alternating test blocks  
967 of stereotype-congruent and incongruent pairings.  
968 The order of these blocks was counterbalanced fol-  
969 lowing the procedure (Nosek et al., 2002). Positive  
970 *D*-scores indicated a stronger implicit stereotype,  
971 resulting from faster responses when male terms  
972 were paired with science and female terms with  
973 liberal arts. Conversely, negative scores reflected  
974 faster responses for female-science/male-arts pair-  
975 ings. We also administered a Race IAT to evaluate  
976 automatic associations between racial categories  
977 (“European American” vs. “African American”)

and evaluative attributes (“pleasant” vs. “unpleasant”) (Dasgupta et al., 2000). Participants categorized faces and words across seven blocks: three for practice and four for testing. Congruent (e.g., European American + pleasant) and incongruent (e.g., African American + pleasant) trial types were counterbalanced across participants. An individual’s *D*-score was calculated as the difference in average response latencies between congruent and incongruent blocks, standardized by the pooled variance. Higher positive scores reflected a stronger implicit preference for European American over African American stimuli (Dasgupta, 2009).

**Explicit Racial Bias.** We used a 2-item scale (Axt, 2018) to measure individuals’ beliefs about Black and White Americans. An example item is: “The How warm or cold do you feel toward Black Americans?” All items were rated on a 10-point scale, ranging from “Extremely Cold” (1) to “Extremely Warm” (5). Responses were averaged to create a composite score (Cronbach’s  $\alpha = 0.84$ ).

**Right-Wing Authoritarianism.** This construct refers to an ideological orientation that prefers stricter authority, traditional norms, and punitive measures (Altemeyer, 1996; Osborne et al., 2023). We used a 15-item scale (Altemeyer, 1996) to measure Right-Wing Authoritarianism. An example item is: “The “old-fashioned ways” and “old-fashioned” values’ still show the best way to live.” All items were rated on a 5-point scale, ranging from “Strongly Disagree” (1) to “Strongly Agree” (5). Responses were averaged to create a composite score (Cronbach’s  $\alpha = 0.84$ ).

**Social Dominance Orientation.** This construct refers to the extent to which individuals prefer group hierarchies and inequality (Kteily et al., 2012). We used a 17-item scale (Pratto et al., 1994) to measure social dominance orientation. An example item is: “Some groups of people are simply not the equals of others.” All items were rated on a 7-point scale, ranging from “Strongly Disagree” (1) to “Strongly Agree” (7). Responses were averaged to create a composite score (Cronbach’s  $\alpha = 0.86$ ).

### B.3 Cultural Orientations

**Individualism & Collectivism.** The cultural orientations of individualism and collectivism are central to understanding how people define themselves, relate to others, and make sense of the social world (Markus and Kitayama, 2014; Oyserman, 2017, 2011; Triandis and Gelfand, 2012). These two constructs can be categorized into four key com-

ponents: horizontal individualism (the extent to which individuals perceive themselves as independent from their groups; e.g., “I want to do my own thing”; (Cronbach’s  $\alpha = 0.82$ )), vertical individualism (the extent to which individuals seek power and status through competition; e.g., “I want to be the best”; (Cronbach’s  $\alpha = 0.58$ )), horizontal collectivism (the extent to which individuals see themselves as part of a group; e.g., “The well-being of my co-workers is important to me”; (Cronbach’s  $\alpha = 0.62$ )), and vertical collectivism (the extent to which individuals are willing to sacrifice personal goals for the benefit of their group; e.g., “I would sacrifice an enjoyable activity if my family disapproved”; (Cronbach’s  $\alpha = 0.80$ )) (Triandis and Gelfand, 1998; Singelis et al., 1995). We used a 15-item scale (Triandis and Gelfand, 1998) to measure these components. All items were rated on a 10-point scale, ranging from “Never” (1) to “Always” (10). Responses were averaged to create composite scores.

**Honor.** This construct assess the extent to which individuals experience distress over damaging their own reputation or that of their family, as well as their perceived inability to protect women’s honor and maintain moral integrity (Uskul et al., 2023). We used a 30-item scale (Ceylan Batur and Sakalli, 2019) to measure honor. An example of the item is “If you were unable to protect your partner when someone harasses him/her.” All items were rated on a 5-point scale anchored by “Not at All bad” (1) to “Extremely bad” (5). Responses to these items were averaged to create a composite measure (Cronbach’s  $\alpha = 0.83$ ).

#### Moral Values

**Moral Foundations.** Moral Foundations Theory posits that there are several dimensions to human morality, each of which serves an evolutionary purpose (Haidt and Joseph, 2004). These moral foundations are intuitive, culturally shaped, and deeply rooted in evolved psychological systems that enabled cooperation (Graham et al., 2013). We used the 36-item Moral Foundation Questionnaire-2 (Atari et al., 2023a) to assess moral concerns: care (e.g., “It pains me when I see someone ignoring the needs of another human being.”; (Cronbach’s  $\alpha = 0.92$ )), equality (e.g., “I believe that everyone should be given the same quantity of resources in life.”; (Cronbach’s  $\alpha = 0.75$ )), loyalty (e.g., “I think children should be taught to be loyal to their country.”; (Cronbach’s  $\alpha = 0.69$ )), authority (e.g., “I believe that one of the most important

values to teach children is to have respect for authority.”; (Cronbach’s  $\alpha = 0.79$ )), purity (e.g., “I believe chastity is an important virtue.”; (Cronbach’s  $\alpha = 0.59$ )), and proportionality (e.g., “It makes me happy when people are recognized for their merits.”; (Cronbach’s  $\alpha = 0.75$ )). All items were rated on a 5-point scale, ranging from “Does not describe me at all” (1) to “Describes me extremely well” (5).

**Qeirat.** This construct is characterized by moral concerns for one’s family, especially female kin, extended family members, and by extension, their country. This construct, closely related to but different from honor, can be considered an additional moral foundation (Atari and Zewail, 2025). We used an 8-item scale (Atari et al., 2020) to quantify participants’ endorsement of Qeirat values. An example of the item is “A man should be responsive to his family’s needs.” All items were rated on a 7-point scale anchored by “Strongly Disagree” (1) to “Strongly Agree” (7). Responses to these items were averaged to create a composite score (Cronbach’s  $\alpha = 0.80$ ).

**Disgust.** Disgust is a moral emotion and has been found to be highly related to moral judgments (Schnall et al., 2008). We used a 13-item scale (Olatunji et al., 2007) to assess the extent to which individuals feel an aversion to various unpleasant or aversive contexts. An example item is: “It would not upset me at all to watch a person with a glass eye take the eye out of the socket.” All items were rated on a 5-point scale, ranging from “Strongly Disagree” (1) to “Strongly Agree” (5). Responses were averaged to create a composite score (Cronbach’s  $\alpha = 0.86$ ).

**Trust.** This construct refers to the extent to which individuals expect others to act in an honest and reliable manner. We used a 6-item scale (Haerper et al., 2020) to assess the extent to which individuals trust other people. An example item is: “Most people are basically good and kind.” All items were rated on a 5-point scale, ranging from “Strongly Disagree” (1) to “Strongly Agree” (5). Responses were averaged to create a composite measure (Cronbach’s  $\alpha = 0.78$ ).

**Intellectual Humility.** Intellectual humility is characterized by the recognition that one’s knowledge and beliefs might be limited or incorrect, and being open to revise one’s views in light of new evidence (Leary et al., 2017). We used a 22-item scale (Krumrei-Mancuso and Rouse, 2016) to quantify intellectual humility. An example of the item

is “When I am really confident in a belief, there is very little chance that belief is wrong.” All items were rated on a 5-point scale anchored by “Strongly Disagree” (1) to “Strongly Agree” (5). Responses to these items were averaged to create a composite score (Cronbach’s  $\alpha = 0.86$ ).

## C Demographic Match Analysis

To assess whether the experimental and control groups were comparable in terms of demographic characteristics, chi-square tests were conducted for gender and race, and an independent samples t-test was performed for age. A chi-square test of independence indicated that gender distribution did not significantly differ between the experimental and control groups,  $\chi^2(1, N = 65) = 0.169, p = .681$ . Similarly, race distribution did not significantly differ between groups,  $\chi^2(1, N = 65) = 0.482, p = .488$ . However, an independent-samples t-test revealed a significant difference in age between the experimental group ( $M = 20.80, SD = 1.73$ ) and the control group ( $M = 19.91, SD = 1.36$ ),  $t(63) = 2.31, p = .024$ , suggesting that participants in the experimental group were slightly older than those in the control group. This difference was primarily driven by a subset of research assistants (10% of the sample) in the experimental group who were between 24 and 26 years old.

## D Institutional Review Board (IRB)

Our IRB includes the anonymity of the participants’ data and ensures minimal risk to them during and after the experiment.

## E Data & Code

Data, code and annotation materials are publicly available at <https://anonymous.4open.science/r/Annotations-DCA6>.

Table 1: The CAMEL Corpus contains data extracted from various datasets. Below is a summary of the sources.

Corpus	Description	<i>N</i>	Source
Moral Foundations Reddit Corpus	Reddit posts related to moral discourse	8,000	Trager et al. (2022)
AITA Subreddit	Submissions from the ‘Am I The A-hole;’ subreddit along with the top 10 comments for each submission	2,000	Boraske (2022)
The 2022 Trucker Strike	Reddit posts capturing the discourse of the 2022 trucker strike	1,000	Grep (2022)
Hate Speech and Offensive Language Dataset	Dataset containing X (Twitter) posts considered racist, sexist, homophobic, or generally offensive.	2,000	Davidson et al. (2017)
Jigsaw Toxic Comment Classification Challenge	A subset of Wikipedia comments focusing on negative online behaviors such as toxic comment	2,000	Adams et al. (2017)
AG News Dataset	News articles spanning four topics: ‘World,’ ‘Sports,’ ‘Business,’ ‘Science and Technology’	4,000	Zhang et al. (2016)
English Historical Quotes Dataset	A collection of random historical quotes	5,000	Roucher (2023)
Stanford Sentiment Treebank	Movie reviews	8,000	Socher et al. (2013)
Chatbot Instruction Prompts Dataset	Prompts and responses for chatbot instructions	2,000	Palla (2024)
Nostalgic Sentiment Analysis of YouTube Comments Dataset	Nostalgic sentiment expressed in Youtube comments	1,000	Aktas (2023)
Moral Stories Dataset	A series of short stories showcasing moral and ethical dilemmas	4,000	Emelin et al. (2021)
Candia et al. Dataset	online posts from X (Twitter), 8Chan, and Reddit	9,000	Candia et al. (2022)
Latin-English Parallel Corpus	English translations of historical Latin texts	1,000	Rosenthal (2023)
Trump Tweets Dataset	A collection of tweets from former President Donald Trump	1,000	Schlatt (2023)
GoEmotions Dataset	Texts spanning 27 emotion categories such as admiration, disappointment, and relief	1,000	Demszky et al. (2020)

Table 2: Effects of Annotation Tasks on Annotator’s Psychological Outcomes

	<i>B</i>	<i>SE</i>	<i>p</i>
Mental Health	0.11	0.09	0.223
PTSD	-0.12	0.16	0.455

\**p*<0.05; \*\**p*<0.01

*Note.* Frequentist analysis of the interaction effects of time×group on the listed variables.

Table 3: Bayesian Hierarchical Model Results: Estimates, 95% Credible Intervals, Bayes Factors, and Interpretation

Outcome	Estimate	95% Credible Interval	BF <sub>10</sub>	BF <sub>01</sub>	Interpretation
Mental Health	0.11	[−0.06, 0.29]	0.02	50.00	Very strong evidence for null
PTSD	−0.12	[−0.45, 0.20]	0.02	50.00	Very strong evidence for null

*Note.* BF<sub>10</sub> = Bayes factor favoring the alternative hypothesis; BF<sub>01</sub> = 1/BF<sub>10</sub> favoring the null hypothesis. 1–3 “anecdotal,” 3–10 “moderate,” 10–30 “strong,” 30–100 “very strong,” >100 “extreme.”

Table 4: Frequentist Treatment Effects in Additional Psychological Measures

Variable	$\beta_3$	SE	<i>p</i>
<b>Psychological Flexibility</b>			
Tolerance of Ambiguity	−0.16	0.12	0.188
Emotion Regulation	0.01	0.22	0.961
Intellectual Humility	0.00	0.10	0.917
Perspective Taking	−0.00	0.10	0.990
<b>Social Biases</b>			
Right-Wing Authoritarianism	−0.19	0.07	0.008**
Social Dominance Orientation	−0.15	0.14	0.375
Explicit Black Preference	0.03	0.36	0.913
Explicit White Preference	−0.14	0.41	0.500
Implicit Gender Bias	0.01	0.10	0.917
Implicit Racial Bias	−0.00	0.10	0.990
<b>Culture Orientations</b>			
Horizontal Individualism	−0.16	0.43	0.765
Vertical Individualism	−0.17	0.38	0.975
Horizontal Collectivism	−0.12	0.27	0.760
Vertical Collectivism	−0.88	0.38	0.038*
Male Honor	0.17	0.11	0.094
Female Honor	0.06	0.10	0.642
<b>Moral Values</b>			
Care	0.03	0.17	0.846
Equality	0.17	0.14	0.303
Loyalty	0.17	0.13	0.305
Proportionality	−0.07	0.16	0.500
Authority	0.04	0.14	0.925
Purity	−0.09	0.12	0.901
Qeirat	−0.09	0.22	0.569

*Note.* Each entry shows the time × group interaction effect. Significance levels: \**p* < 0.05, \*\**p* < 0.01.

Table 5: Bayesian Treatment Effects in Additional Psychological Measures

Variable	Estimate	95% Credible Interval	BF <sub>10</sub>	BF <sub>01</sub>	Interpretation
<b>Psychological Flexibility</b>					
Tolerance of Ambiguity	-0.16	[-0.39, 0.08]	0.03	33.33	Very strong evidence for null
Emotion Regulation	0.02	[-0.53, 0.58]	0.02	50.00	Very strong evidence for null
<b>Social Biases</b>					
Right-Wing Authoritarianism	-0.20	[-0.34, -0.06]	0.25	4.00	Moderate evidence for null
Social Dominance Orientation	-0.14	[-0.42, 0.15]	0.02	50.00	Very strong evidence for null
Explicit Black Preference	0.00	[-0.74, 0.75]	0.03	33.33	Very strong evidence for null
Explicit White Preference	-0.27	[-1.09, 0.56]	0.04	25.00	Strong evidence for null
Implicit Gender Bias	0.01	[-0.17, 0.19]	0.02	50.00	Very strong evidence for null
Implicit Racial Bias	0.00	[-0.18, 0.18]	0.03	33.33	Very strong evidence for null
<b>Culture Orientations</b>					
Horizontal Individualism	-0.16	[-0.44, -1.02]	0.05	20.00	Strong evidence for null
Vertical Individualism	-0.03	[-0.78, 0.71]	0.04	25.00	Strong evidence for null
Horizontal Collectivism	-0.10	[-0.87, 0.67]	0.04	25.00	Strong evidence for null
Vertical Collectivism	-0.78	[-1.54, -0.02]	0.29	3.45	Moderate evidence for null
Male Honor	0.18	[-0.04, 0.39]	0.04	25.00	Strong evidence for null
Female Honor	0.04	[-0.16, 0.25]	0.01	100.00	Very strong evidence for null
<b>Moral Values</b>					
Care	0.03	[-0.32, 0.39]	0.02	50.00	Very strong evidence for null
Equality	0.21	[-0.09, 0.50]	0.04	25.00	Strong evidence for null
Loyalty	0.20	[-0.09, 0.49]	0.04	25.00	Strong evidence for null
Proportionality	-0.07	[-0.41, 0.27]	0.02	50.00	Very strong evidence for null
Authority	0.05	[-0.26, 0.35]	0.02	50.00	Very strong evidence for null
Purity	-0.05	[-0.29, 0.20]	0.01	100.00	Very strong evidence for null
Qeirat	-0.12	[-0.57, 0.33]	0.02	50.00	Very strong evidence for null

Note. BF<sub>10</sub>= Bayes factor favoring the alternative; BF<sub>01</sub> = 1/BF<sub>10</sub> favoring the null. 1–3 “anecdotal,” 3–10 “moderate,” 10–30 “strong,” 30–100 “very strong,” >100 “extreme.”

Table 6: Descriptive Statistics for Mental Health, Psychological Flexibility, and Social Biases

	Mental Health and Psychological Flexibility				Social Biases				
	<i>M</i> <sub>Time1</sub>	<i>SD</i> <sub>Time1</sub>	<i>M</i> <sub>Time2</sub>	<i>SD</i> <sub>Time2</sub>	<i>M</i> <sub>Time1</sub>	<i>SD</i> <sub>Time1</sub>	<i>M</i> <sub>Time2</sub>	<i>SD</i> <sub>Time2</sub>	
Mental Health	3.04	.56	2.88	.49	Right- Wing Authoritarianism	2.22	.52	2.15	.57
PTSD	1.91	.64	2.17	.71	Social Dominance Orientation	2.11	.73	2.24	.75
Intellectual Humility	3.76	0.41	3.73	.41	Explicit White Preference	7.94	2.18	7.79	1.81
Perspective Taking	3.27	.32	3.19	0.41	Explicit Black Preference	8.56	2.15	8.34	1.72
Emotion Regulation	4.85	0.86	4.98	0.84	Implicit Gender Bias	0.32	0.36	0.25	0.53
Perspective Taking	3.27	1.11	3.19	1.26	Implicit Racial Bias	0.07	0.44	0.06	0.43

Note. Well-being was measured on a 1-4 scale, while PTSD symptoms, Intellectual Humility, Perspective Taking, and Right-Wing Authoritarianism (RWA) were measured on a 1-5 scale. Emotion Reappraisal, Suppression, and Social Dominance Orientation (SDO) were measured on a 1-7 scale. Explicit racial preferences were measured on a 1-10 scale, and implicit biases were measured as association scores ranging from -1 to +1. For all variables, lower values indicate low endorsement and higher values indicate high endorsement.

Table 7: Descriptive Statistics for Moral Values and Cultural Orientations

	Moral Values				Cultural Orientation				
	$M_{Time_1}$	$SD_{Time_1}$	$M_{Time_2}$	$SD_{Time_2}$	$M_{Time_1}$	$SD_{Time_1}$	$M_{Time_2}$	$SD_{Time_2}$	
Care	3.98	.86	3.91	0.72	Horizontal Individualism	8.13	1.55	8.12	1.49
Equality	2.65	.83	2.67	.84	Vertical Individualism	8.56	1.04	8.36	1.10
Proportionality	3.67	.64	3.60	.70	Horizontal Collectivism	7.10	1.79	7.06	1.61
Loyalty	2.83	.64	2.86	.66	Vertical Collectivism	5.98	1.59	5.97	1.28
Authority	2.65	0.67	2.64	.84	Male Honor	4.35	0.27	4.14	.37
Purity	2.29	0.66	2.23	.66	Female Honor	3.41	.53	3.36	.49
Trust	3.36	.67	3.33	.70	Moral Integrity	4.35	.31	4.18	.42
Disgust	3.05	.51	3.10	.52	Qeirat	4.75	.94	4.56	1.02

*Note.* Moral Values, Trust, Disgust, Honor, and Qeirat were measured on a 1–5 scale, and Individualism and Collectivism were measured on a 1–10 scale. For all variables, lower values indicate low endorsement and higher values indicate high endorsement.