

Building Knowledge-Guided Lexica to Model Cultural Variation

Anonymous ACL submission

Abstract

Cultural variation exists between regions (e.g., the United States vs. China), but also *within* regions (e.g., California vs. Texas, Los Angeles vs. San Francisco). Measuring this regional cultural variation can illuminate how and why people think and behave differently. Historically, it has been difficult to computationally model cultural variation due to a lack of training data and scalability constraints. In this work, we introduce a new research problem for the NLP community: *How do we measure variation in cultural constructs across regions using language?* We then provide a scalable solution: building knowledge-guided lexica to model cultural variation, encouraging future work at the intersection of NLP and cultural understanding.

1 Introduction

People think and behave differently around the world. This is partly due to *cultural variation*, or the differences among individuals that exist due to some form of social learning (Cohen, 2001). Having a computational method that utilizes language to measure cultural variation could help us better understand humans (Tsai et al., 2006; Oishi et al., 2009), build more culturally-aware NLP systems (Hovy and Yang, 2021), and advance interdisciplinary research in anthropology, cultural psychology, etc. However, due to a lack of data and scalability constraints, few such methods exist.

In this paper, we present *measuring regional variation in culture* as a problem of interest for the NLP community and build a knowledge-guided lexical model as a scalable solution. Specifically, we focus on measuring *individualism and collectivism*¹ across the United States (US) using geolocated Tweets.

¹Cultural psychologists have quantified axes on which culture differs, also called *cultural dimensions*. A key cultural dimension that influences behaviors like voting, donating, etc. is *individualism vs. collectivism* (Hofstede, 2011). Collectivism stresses the importance of the community, while individualism focuses on each person’s rights and concerns.

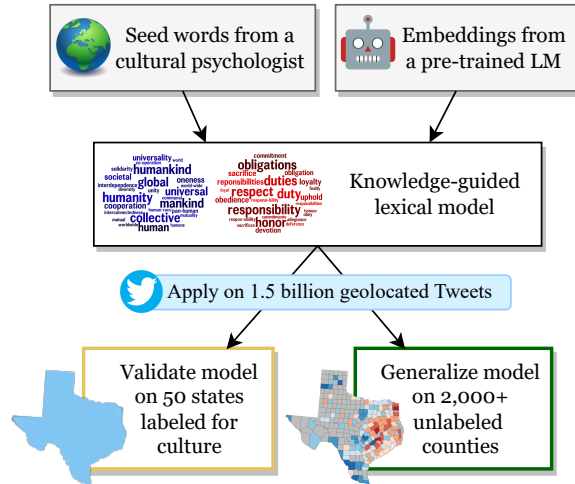


Figure 1: We build knowledge-guided lexica to model cultural variation using two types of domain knowledge: seed words based on cultural psychology theory and embeddings from a pre-trained language model.

Historically, measuring cultural dimensions across regions has been mostly done through questionnaires, such as the World Values Survey (WVS) (Haerpfer et al., 2020). However, questionnaires are time-consuming and heavily restricted in scope; the most recent WVS wave required 4 years and averaged 52 participants per US state. Recent work probes language models (LMs) for cultural values (Arora et al., 2023), but these LMs do not reflect all cultures equally (Havaladar et al., 2023).

The overhead of traditional survey-based approaches and inconsistent cultural awareness of existing LMs motivates scalable, computational methods that use *existing language data* to measure cultural variation instead. For the example problem addressed in this paper, we seek to measure individualism and collectivism across US counties using the following resources:

- Domain expertise from cultural psychologists.
- An open-source corpus (see Appendix A) of 1.5 billion geolocated Tweets from 6 million

058 US users (Giorgi et al., 2018).
059 • Individualism and collectivism scores for fifty
060 US states (Vandello and Cohen, 1999).

061 Pre-LM era solutions to measure culture use single
062 words (Giorgi et al., 2020) or manually curated
063 lexica (Graham et al., 2009), thus relying on a small
064 number of highly specific words. A more modern
065 NLP solution would take the form of either (1)
066 training a model, or (2) prompting an LM. How-
067 ever, to classify 1.5 billion Tweets, (1) requires a
068 sizable amount of labeled training data, and (2) is
069 not computationally scalable. For instance, run-
070 ning this corpus through GPT-4 would cost roughly
071 \$900,000 (see Appendix B).

072 Additionally, building a Tweet-level deep learn-
073 ing model to predict culture is impractical. Most of
074 an individual’s language does not indicate their cul-
075 tural beliefs; therefore, it is prohibitively expensive
076 to label enough Tweets to train an adequate model.

077 Our method builds upon a line of work in
078 NLP called lexicon induction (Araque et al., 2020;
079 Buechel et al., 2020; Geng et al., 2022), which
080 analyzes massive corpora in NLP without solely
081 relying on deep learning. Past work mainly builds
082 lexica for sentiment, emotion, etc. We uniquely fo-
083 cus on the domain where *little training data exists*
084 and *not every utterance can be relevantly labeled*.

085 **Leveraging domain knowledge.** Our proposed
086 method to model cultural variation utilizes both
087 domain expertise from cultural psychology (via a
088 set of expert-curated seed words) and knowledge
089 implicit in LMs (via word embeddings) to build
090 scalable lexical models.

091 We validate our method against past collectivism
092 research at the US state-level (Pelham et al., 2022;
093 Vandello and Cohen, 1999) and extend the anal-
094 ysis of individualism and collectivism across US
095 counties, allowing for a more fine-grained spatial
096 analysis (i.e., understanding how large areas like
097 states are culturally heterogeneous).

098 We also show how county-level analyses of cul-
099 ture can obtain new insights into existing popu-
100 lations, via a taxonomy of *communities* (socio-
101 demographic clusters of counties) from the Amer-
102 ican Communities Project (Chinni and Gimpel,
103 2011). This taxonomy has been previously utilized
104 to understand differences in health behaviors and
105 outcomes (Aggarwal et al., 2023; Guntuku et al.,
106 2021), and we use it to better understand how com-
107 munities vary in individualism and collectivism.

2 Building Knowledge-Guided Lexica

108 Lexica, or sets of curated words, are a highly scal-
109 able method for analyzing large datasets. However,
110 building lexica linked to cultural theory from the
111 ground up is also a time-intensive process.

112 To mitigate this, we propose a method that com-
113 bines two types of domain knowledge to efficiently
114 create lexica that can measure cultural variation.
115 We first ask an expert psychologist to generate two
116 small sets of seed words that capture individualism
117 and collectivism respectively based on their knowl-
118 edge of cultural psychology. We next leverage
119 knowledge implicitly present in language models
120 (e.g., word associations, word similarity, etc.) to
121 expand these small sets of seed words into high-
122 validity lexical models.

123 Using this method, we can measure regional vari-
124 ation for any cultural construct using language from
125 those regions. Our approach has two components:
126 Expansion and Purification. Figure 2 details this
127 approach for our example problem – measuring
128 individualism and collectivism across US counties.

129 **Step 1: Expansion** Given a set of seed words
130 from a cultural psychologist (see Appendix for seed
131 words), we utilize word embeddings² to expand the
132 set of seed words in two ways: we locate all words
133 that are similar to each individual seed word (*syn-
134 onym expansion*), as well as locate the words that
135 are similar to the overall construct described by the
136 complete set of seed words (*concept expansion*).

137 For synonym expansion, we find the nearest
138 neighbors for each individual seed word in em-
139 bedding space and add these neighboring words
140 to our lexica. For concept expansion, we average
141 the embeddings of each seed word set (e.g. indi-
142 vidualism) to find the *centroid embeddings*. We
143 then find the nearest neighbors of each centroid
144 embedding. By using embedding space to expand
145 our lexica, we additionally calculate a weight for
146 each expanded word, i.e., the cosine similarity be-
147 tween the expanded word and the corresponding
148 seed word or centroid embedding. The weight for
149 each seed word is 1.

150 This method is highly tunable – any embeddings
151 can be used, and the number of nearest neighbors
152 returned during expansion can be adjusted based
153 on desired length of the final lexicon.

²We use FastText (Bojanowski et al., 2017) due to its fixed vocabulary size, efficient nearest neighbors functionality, and ability to find synonyms in context-free scenarios, but our methods are more general and agnostic to embedding type.

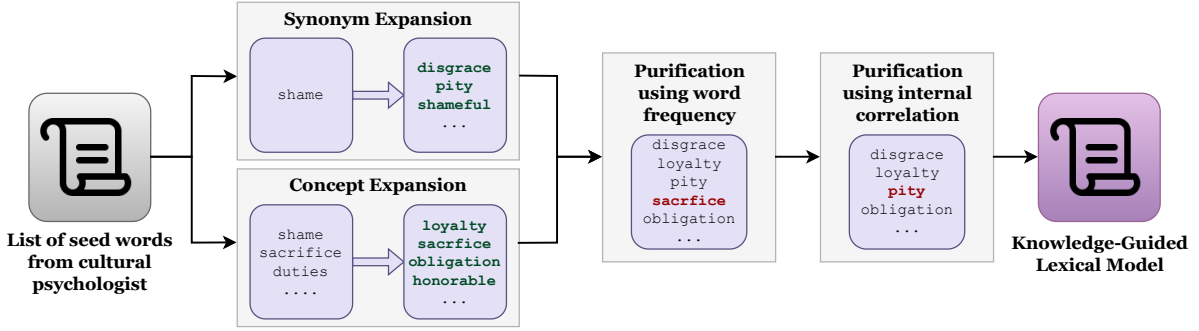


Figure 2: Our knowledge-guided lexica creation method. The first stage, *Expansion*, consists of synonym expansion and concept expansion, done in parallel. The second stage, *Purification*, includes frequency-based and correlation-based pruning, done sequentially.

Step 2: Purification Upon aggregating the words returned from both expansion types, we want to ensure that the resulting lexica are both pertinent and internally correlated.

To ensure pertinence, we filter out rare words, or any words below a given usage frequency (Bojanowski et al., 2017). Next, we ensure internal correlation. We apply our lexica to our US Twitter Corpus and compute the weighted frequencies for each word at two granularities: county-level and state-level. This produces scores that reflect the individualism and collectivism tendencies of every US county and state. We then check how each individual word’s frequency correlates with the corresponding overall individualism/collectivism score. Any word that doesn’t show a significant positive correlation (product-moment correlation coefficient $r < 0.15$) is removed from the lexica. This step ensures that every word contributes correctly to measuring the relevant cultural dimension.

Figure 5 visualizes our knowledge-guided individualism and collectivism lexica.

3 Validation

Upon expanding and purifying the lexica, we validate our results using the collectivism scores from Vandello and Cohen (1999) for each of the 50 US states. We see a significant positive correlation of our collectivism scores with Vandello & Cohen’s collectivism scores (Table 1).

We also use relevant collectivism indicators from the Global Collectivism Index (Pelham et al., 2022) – religiosity, living arrangements (i.e. grandparent living in the household), and in-group bias. Using corresponding questions from the 2017 U.S. census and the 2017 wave of the World Values Survey (Haerper et al., 2020), we get data for all of these

	Individualism Lexicon Score	Collectivism Lexicon Score
Vandello & Cohen’s Collectivism Scores	-0.374	0.388
Living Arrangements	-0.291	0.200
Religiosity	-0.658	0.400
Ingroup Bias	-0.513	0.464

Table 1: Pairwise product-moment correlations between our individualism and collectivism lexica applied to our Twitter Corpus and validation variables. We use Vandello & Cohen’s Collectivism Scores and GCI indicators, at the US state-level, for validation. All correlations are significant ($p < 0.05$).

indicators at the US state level. We also see a significant positive correlation with all of these indicators (Table 1). Further details on validation are given in Appendix C.

4 Results

We apply our validated knowledge-guided lexica to county-level geolocated Tweets, to gain a more fine-grained understanding of how individualism and collectivism vary regionally.³

Figure 4 illustrates this variation, plotting the difference between individualism and collectivism score. The deep south shows high levels of collectivism (dark red) and low levels of individualism (light blue). Conversely, the West coast and the Northeast show low levels of collectivism (light red) and high levels of individualism (dark blue). Counties with under 100 users are poorly represented (Giorgi et al., 2018) and are colored gray.

³We release our lexica, county-level and state-level scores, and relevant code at https://anonymous.4open.science/r/knowledge_driven_lexica-E8EE/

Individualism and Collectivism across ACP Communities

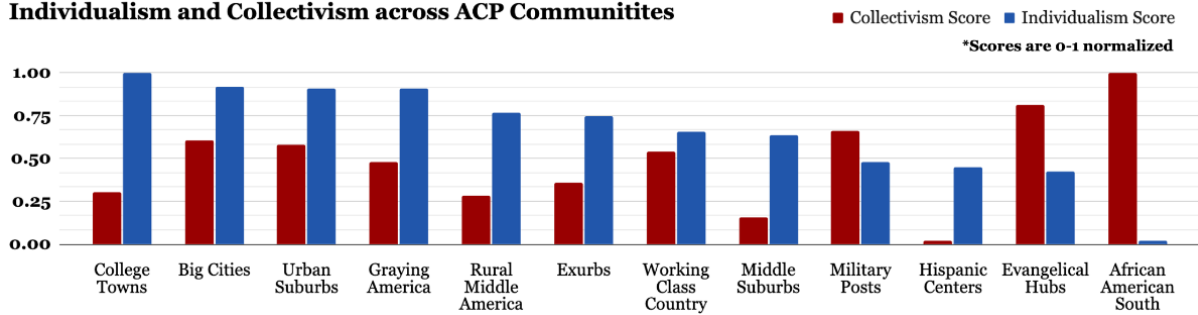


Figure 3: A comparison of collectivism (red) and individualism (blue) scores across communities defined by the American Communities Project, ordered from most individualistic (left) to least individualistic (right). We only analyze communities with over 40 included counties. Scores are 0-1 normalized.

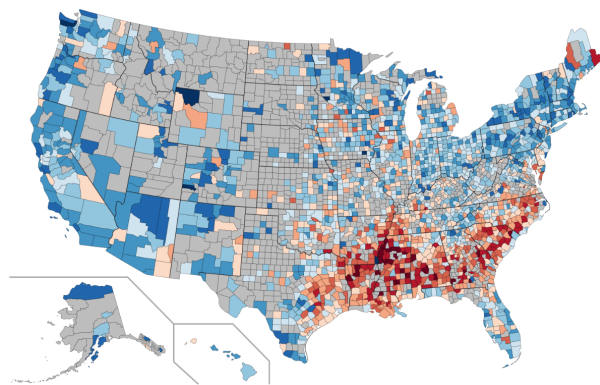


Figure 4: Collectivism (red) and individualism (blue) across US counties. Dark red = higher collectivism and dark blue = higher individualism. Gray counties have insufficient Tweets to estimate a score.

and Big Cities are highly individualist. These areas are also more affluent and have higher rates of education (ACP, 2023). This fits with prior research finding that people who are wealthy or educated tend to be more individualistic (Binder, 2019). In contrast, the data shows that Evangelical Hubs and the African American South are highly collectivist. These communities are tight-knit and religious areas (ACP, 2023), which have been linked to collectivism (Pelham et al., 2022). Military Posts are also more collectivist, which fits with the tight ties in military service and “duty to one’s troop.” This insight is helpful because we know of no cultural psychology research comparing military communities with civilian communities.

5 Conclusion & Future Work

We present a method to efficiently measure cultural variation by leveraging domain knowledge from cultural psychology and language models to create knowledge-guided lexica. These lexica, applied to social media language, can estimate cultural differences at fine-grained geographic levels, such as states, counties, and communities.

Future work could build on this method to get deeper insights into communities and cultures. For example, our method could be used to identify more types of Tweets that mark cultural differences; we encourage researchers to build more sophisticated models on these identified Tweets. Additionally, our method is easily extendable to other cultural dimensions, such as tightness/looseness, future orientation, etc. This method could also measure cultural variation globally, which requires analyzing different languages. Since our method is language agnostic, it can easily extend to non-English settings by leveraging multilingual embeddings.

Community-level insights. Cultural similarity is not always based on geographical proximity; two cities hundreds of miles apart may be more similar than a city and a rural farm a few miles away (Guntuku et al., 2021). To show how county-level analyses of culture can help us better understand communities, we additionally use 15 community types (e.g., College Towns, Urban Suburbs) identified by the American Communities Project (ACP). The ACP identified these communities based on socio-demographic attributes, not spatial clusters of counties. Previous studies have used these community types to identify cultural variation in excessive alcohol consumption (Giorgi et al., 2020) and self-reported physical and mental health (Aggarwal et al., 2023; Mangalik et al., 2023).

Figure 3 shows county-level individualism and collectivism scores grouped into their corresponding ACP community (see Table 3 for counts.) These results provide novel insights into how culture varies regionally. For example, College Towns

6 Limitations

While we label each county for individualism and collectivism, we note that regions do not have a single culture. Within all regions, there is heterogeneity of cultural values and beliefs. Since we use an open-source Twitter corpus, we also have poor coverage of counties with little to no Twitter data. Additionally, not all aspects of culture are revealed in language – we are limited to analyzing only what people say online.

In our analyses, we do not control for race, income, or other demographic variables. We know cultural values are correlated to some demographic variables. For example, collectivism and individualism vary with income. Future work can improve upon these estimates by accounting for individual demographics. Additionally, it is unclear if this method of measuring cultural variation will work for all cultural dimensions. For example, power distance (Hofstede, 2011) involves the relationship dynamics of two people, which might make it difficult to capture with lexica.

7 Ethical Considerations

The goal of studying cultural variation is to better understand cultures, not individuals. Nonetheless, the characterization of culture has the danger of stereotyping individuals. Individuals within each culture vary greatly. Studying culture can help us understand differences in psychology, but we should not assume that a cultural average will definitely apply to a particular individual from that culture.

All data used in this study are publicly available. While geolocated Twitter data is used, only aggregated spatial-level data is reported. That is, no person-level identifiable information is used or released for this study.

References

ACP. 2023. Home - american communities project - americancommunities.org. <https://www.americancommunities.org/>. [Accessed 14-08-2023].

Arnav Aggarwal, Sunny Rai, Salvatore Giorgi, Shreya Havaladar, Garrick Sherman, Juhi Mittal, and Sharath Chandra Guntuku. 2023. A cross-modal study of pain across communities in the united states. In *Companion Proceedings of the ACM Web Conference 2023*, pages 1050–1058.

Oscar Araque, Lorenzo Gatti, and Kyriaki Kalimeri. 2020. Moralstrength: Exploiting a moral lexicon and embedding similarity for moral foundations prediction. *Knowledge-based systems*, 191:105184.

Arnav Arora, Lucie-aimée Kaffee, and Isabelle Augenstein. 2023. Probing pre-trained language models for cross-cultural differences in values. In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 114–130, Dubrovnik, Croatia. Association for Computational Linguistics.

Carola Conces Binder. 2019. Redistribution and the individualism–collectivism dimension of culture. *Social Indicators Research*, 142(3):1175–1192.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the association for computational linguistics*, 5:135–146.

Sven Buechel, Susanna Rücker, and Udo Hahn. 2020. Learning and evaluating emotion lexicons for 91 languages. *arXiv preprint arXiv:2005.05672*.

Dante Chinni and James Gimpel. 2011. *Our patchwork nation: The surprising truth about the "real" America*. Penguin.

Dov Cohen. 2001. Cultural variation: considerations and implications. *Psychological bulletin*, 127(4):451.

Yilin Geng, Zetian Wu, Roshan Santhosh, Tejas Srivastava, Lyle Ungar, and João Sedoc. 2022. Inducing generalizable and interpretable lexica. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4430–4448.

Salvatore Giorgi, Daniel Preoțiu-Pietro, Anneke Buffone, Daniel Rieman, Lyle Ungar, and H. Andrew Schwartz. 2018. The remarkable benefit of user-level aggregation for lexical-based population-level predictions. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1167–1172, Brussels, Belgium. Association for Computational Linguistics.

Salvatore Giorgi, David B Yaden, Johannes C Eichstaedt, Robert D Ashford, Anneke EK Buffone, H Andrew Schwartz, Lyle H Ungar, and Brenda Curtis. 2020. Cultural differences in tweeting about drinking across the us. *International journal of environmental research and public health*, 17(4):1125.

Jesse Graham, Jonathan Haidt, and Brian A Nosek. 2009. Liberals and conservatives rely on different sets of moral foundations. *Journal of personality and social psychology*, 96(5):1029.

Sharath Chandra Guntuku, Alison M. Buttenheim, Garrick Sherman, and Raina M. Merchant. 2021. Twitter discourse reveals geographical and temporal variation in concerns about covid-19 vaccines in the united states. *Vaccine*, 39(30):4034–4038.

367	Christian Haerper, Ronald Inglehart, Alejandro
368	Moreno, Christian Welzel, Kseniya Kizilova, Jaime
369	Diez-Medrano, Marta Lagos, Pippa Norris, Eduard
370	Ponarin, Bi Puranen, et al. 2020. World values sur-
371	vey: Round seven–country-pooled datafile. <i>Madrid,</i>
372	<i>Spain & Vienna, Austria: JD Systems Institute &</i>
373	<i>WVSA Secretariat</i> , 7:2021.
374	Shreya Havaldar, Bhumika Singhal, Sunny Rai,
375	Langchen Liu, Sharath Chandra Guntuku, and Lyle
376	Ungar. 2023. Multilingual language models are not
377	multicultural: A case study in emotion . In <i>Proceed-</i>
378	<i>ings of the 13th Workshop on Computational Ap-</i>
379	<i>proaches to Subjectivity, Sentiment, & Social Media</i>
380	<i>Analysis</i> , pages 202–214, Toronto, Canada. Associa-
381	tion for Computational Linguistics.
382	Geert Hofstede. 2011. Dimensionalizing cultures: The
383	hofstede model in context. <i>Online readings in psy-</i>
384	<i>chology and culture</i> , 2(1):8.
385	Dirk Hovy and Diyi Yang. 2021. The importance of
386	modeling social factors of language: Theory and
387	practice . In <i>Proceedings of the 2021 Conference</i>
388	<i>of the North American Chapter of the Association</i>
389	<i>for Computational Linguistics: Human Language</i>
390	<i>Technologies</i> , pages 588–602, Online. Association
391	for Computational Linguistics.
392	Siddharth Mangalik, Johannes C Eichstaedt, Salva-
393	tore Giorgi, Jihu Mun, Farhan Ahmed, Gilvir Gill,
394	Adithya V Ganesan, Shashanka Subrahmanya, Nikita
395	Soni, Sean AP Clouston, et al. 2023. Robust
396	language-based mental health assessments in time
397	and space through social media. <i>arXiv preprint</i>
398	<i>arXiv:2302.12952</i> .
399	Shigehiro Oishi, Ed Diener, Richard E Lucas, and
400	Eunkook M Suh. 2009. Cross-cultural variations
401	in predictors of life satisfaction: Perspectives from
402	needs and values. <i>Culture and well-being: The col-</i>
403	<i>lected works of Ed Diener</i> , pages 109–127.
404	Brett Pelham, Curtis Hardin, Damian Murray, Mitsuru
405	Shimizu, and Joseph Vandello. 2022. A truly global,
406	non-weird examination of collectivism: The global
407	collectivism index (gci). <i>Current Research in Eco-</i>
408	<i>logical and Social Psychology</i> , 3:100030.
409	Jeanne L Tsai, Brian Knutson, and Helene H Fung.
410	2006. Cultural variation in affect valuation. <i>Journal</i>
411	<i>of personality and social psychology</i> , 90(2):288.
412	Joseph A Vandello and Dov Cohen. 1999. Patterns
413	of individualism and collectivism across the united
414	states. <i>Journal of personality and social psychology</i> ,
415	77(2):279.

A Open-Source Twitter Corpus 416

We use the County Tweet Lexical Bank, an open source data set of features extracted from a corpus of 1.5 billion tweets from approximately 6 million US county mapped users (Giorgi et al., 2018). While the full details of the dataset can be found in the original paper, we give a high-level summary to aid the reader. The dataset is built from a larger corpus which is a 10% sample of Twitter from 2009-2015 (over 30 billion tweets). These tweets are then mapped to US counties via latitude and longitude coordinates associated with the tweets or self-reported location information in the Twitter user’s profile (a free text field). A Twitter user is included in this data set if they have posted at least 30 or more English tweets, and a county is included if at least 100 such users are mapped to that respective county. This process resulted in 1.5 billion tweets mapped to over 2000 US counties.

B Scalability Calculations 435

We outline the proposed costs of using various LM-based techniques to label our corpus of 1.5 billion Tweets:

Proposed cost of GPT-4 As of August 2023, the OpenAI API rate for GPT-4 is \$0.06 cents per 1,000 tokens. Assuming 10 tokens per Tweet, we get:

$$1.5e9 \text{ Tweets} \times \frac{10 \text{ Tokens}}{\text{Tweet}} \times \frac{\$0.06}{1,000 \text{ Tokens}} \quad (1) \quad 442$$

This yields a total cost of \$900,000. 443

Proposed cost of GPT-3.5 As of August 2023, the OpenAI API rate for GPT-3.5 is \$0.002 cents per 1,000 tokens. Assuming 10 tokens per Tweet, we get:

$$1.5e9 \text{ Tweets} \times \frac{10 \text{ Tokens}}{\text{Tweet}} \times \frac{\$0.002}{1,000 \text{ Tokens}} \quad (2) \quad 448$$

This yields a total cost of \$30,000. 449

C Validation: Additional Details 450

All six variables in the Global Collectivism Index – total fertility rate, living arrangements (% households with people over 60 and children under 14), stability of marriage (divorce rate to marriage rate ratio), religiosity, collective transportation, and in-group bias (approximated by compatriotism due to lack of state-level data) – are replicable at the state-level using US census data and WVS data. 451
452
453
454
455
456
457
458

Collectivism Seed Words	duties, responsibilities, role, fit in, community, sacrifice, shame, required, rules, honor, support, rely, loyal, respect, obedience
Individualism Seed Words	humans, humanity, worldwide, universal, mankind, everyone, collective, global, equity, imagination, cooperate, cooperation, shared, joint, identity, guilt, diversity

Table 2: Seed words hypothesized to identify individualism and collectivism on social media, provided by a domain expert in cultural psychology.

ACP Community	Num Counties
Exurbs	207
Graying America	164
African American South	252
Evangelical Hubs	269
Working Class Country	159
Military Posts	70
Urban Suburbs	103
College Towns	151
Big Cities	46
Hispanic Centers	87
Rural Middle America	403
Middle Suburbs	77

Table 3: Number of included counties for each ACP community included in the analysis in Figure 3.

Note that when aggregating US census data from county-level to state-level, we treat each county as being weighted equally, due to disproportionate amounts of data coming from big cities.

In order to determine which of these six replicated variables also measure collectivism within the United States, we sample subsets of the six variables and use Cronbach’s alpha to measure internal consistency. We limit the subsets to size three or larger, following Pelham and colleagues’ (Pelham et al., 2022) validation of three collectivism indicators per nation. The set of living arrangements, religiosity, and compatriotism yielded the highest Cronbach’s alpha (0.702), so we chose these three variables as a validation metric.

Table 1 shows the correlations between each of the three validation variables, the collectivism lexicon score, and the individualism lexicon score for US states. Collectivism word use positively correlates with all validation outcomes, and individualism word use correlates negatively. We further validate against median income at the state-level. Prior research has found that income is negatively correlated with collectivism (Pelham et al., 2022) Similarly, income was negatively correlated with our collectivism lexicon scores (-0.273) and positively with our individualism lexicon scores (0.424). We also observe a strong negative correlation (-0.470)



Figure 5: Word clouds visualizing our individualism lexica (blue, top) and collectivism lexica (red, bottom). Larger words have a higher weight, while smaller words have a lower weight.

between our individualism and collectivism scores at the US state level.

We also validate against Vandello and Cohen’s collectivism scores (Vandello and Cohen, 1999). We see a positive correlation (0.388) with our collectivism lexicon scores and a negative correlation (-0.374) with our individualism lexicon scores. This suggests that our lexica measurements are indeed tapping into real cultural differences.