

Ground-truthing perspectives on highly subjective text: basic human values perceived in song lyrics

Anonymous ACL submission

Abstract

We present an interdisciplinary approach to creating a dataset on a highly subjective text annotation task. The task thus requires explicit insight into broader human annotator perspectives and perceptions, and conscious curation of what will be annotated. In this, with strong inspiration from best practices in the social sciences, we add to emerging and increasing calls for greater accountability with regard to data and its quality. For our task, we choose the annotation of perceived human values in song lyrics. Drawing from a representative US population sample, we present our strategy to select song lyrics to be annotated, estimate the amount of annotators needed, and assess data quality. Based on this, we obtain a dataset of 360 richly annotated song lyrics. We substantiate the benefit of having more annotators, and show how annotations show promising consistency with earlier insights on personal value proximity from a validated cross-cultural instrument study. Finally, we give a first illustration of how our data can be employed in connection to applied machine learning approaches.

1 Introduction

With growing interest in AI and the rising popularity of Large Language Models, AI advances appear to push for larger datasets to train models, which ideally need few human annotations. At the same time, language is a cultural phenomenon, in which human interpretation plays a key role in transmission and understanding.

In broader situations in which applied machine learning techniques may automate and scale up actions that formerly relied on human perception and judgement, the question of what makes for good data and ‘ground truth’ to depart from has been less articulated and appreciated than the promise of generalizability and scalability by the applied machine learning techniques (Birhane et al., 2022; Sambasivan et al., 2021). However, calls for data-centric AI

have recently been emerging¹, and recognition that human annotator disagreement can be a meaningful signal, rather than noise suggesting unreliable annotation is increasing (Aroyo and Welty, 2015). Furthermore, awareness is rising on the need for more explicit data documentation, mostly from the perspective of higher accountability on responsible data handling and reporting, and out of concern for potential societally harmful consequences of irresponsible practice (Gebru et al., 2021; Mitchell et al., 2019; Geiger et al., 2020). Efforts to more strongly institutionalize responsible practice also are visible in the *ACL communities (Rogers et al., 2021), where the completion of a Responsible NLP Checklist presently has become a mandatory element of manuscript submissions.

Computational researchers historically may not have been trained to be aware of data quality considerations. As such, standardized checklists, forms and best practice ‘rules of thumb’ help lowering the threshold to report and discuss these. At the same time, a simple rule of thumb may not stimulate critical reflection on current common practice. For example, as for the question, “How many annotators would be needed for NLP corpus ground truth?”, a well-cited book on natural language annotation for machine learning (Pustejovsky and Stubbs, 2013) suggests to “have your corpus annotated by at least two people (more is preferable, but not always practical)” before being ready to move on to gold standard data. This is a remarkably low number, without clear substantiation of whether this indeed would be sufficient.

Beyond the computational domain, other domains and disciplines have (typically for a much longer time than the computational domain) been building expertise on how to properly curate for data, and capture aspects of the data that may not trivially be measurable. For example, both

¹see <https://datacentricai.org/>

081 in archives and museums, long-standing traditions
082 of purposeful and well-documented curation exist
083 (Jo and Gebru, 2020; Huang and Liem, 2022).
084 Furthermore, in the quantitative social sciences,
085 well-established best practices exist for situations
086 in which *constructs*, i.e., phenomena that cannot
087 directly physically be measured, are to be quan-
088 tified in valid and reliable ways, based on hu-
089 man responses. For this, the discipline of psycho-
090 metrics (Furr and Bacharach, 2014) is taught as
091 an entry-level course to students in psychology,
092 whereas long-standing expertise from survey sci-
093 ence (Groves et al., 2009) can further assess in
094 robust sampling, and designing for robust human
095 responses. While this expertise has been referred to
096 in several works targeting computationally oriented
097 research communities (Welty et al., 2019; Jacobs
098 and Wallach, 2021; Kern et al., 2023), to the best
099 of our knowledge, institutionalized methodological
100 uptake of the expertise remains rare.

101 The current paper resulted from work in an in-
102 terdisciplinary team, including members with dis-
103 ciplinary backgrounds in the computational and
104 social sciences. Departing from an interest in de-
105 veloping well-substantiated ground-truthing pro-
106 cedures for (highly) subjective annotation tasks,
107 we describe the creation of a dataset with annota-
108 tions of perceived human values in song lyrics. We
109 approach this task in a way that social scientists
110 would: we increase the odds of obtaining valid and
111 reliable measurements by being purposeful about
112 strata in data sampling, gaining annotations from a
113 representative human population sample, explicitly
114 investigating the impact of higher numbers of an-
115 notations per item, and relying on evidence-backed
116 theory, in the sense that phenomena to be studied
117 were shown in earlier literature to have scientific
118 and general validity, and our results can be related
119 to earlier established results.

120 2 Background

121 2.1 Perspectivist Ground Truthing

122 Automated systems often rely on manually anno-
123 tated reference data for training and evaluation.
124 Multiple labels from multiple annotators are gath-
125 ered for reasons associated with the annotators
126 themselves, e.g. a lack of trust in crowdsourcing
127 or annotations from non-experts, or because there
128 is an expectation that people will vary in their re-
129 sponses to the phenomenon of interest (Cabitza
130 et al., 2023; Basile et al., 2020). These annotations

131 are then aggregated to produce a single label that
132 is used to train and / or evaluate systems, as it is
133 often incumbent on automated systems to produce
134 a singular response.

135 Thus, most problems are treated as ‘classifica-
136 tion’ problems. Variance occurring in the reported
137 annotations is removed, usually by taking the la-
138 bel chosen most often by the annotators. Although
139 the sources that give rise to the variance observed
140 in the data may legitimately vary (as is possible
141 even in ‘objective’ problems where annotators are
142 medical experts (Kompa et al., 2021)), variance is
143 often treated as an error even when a case can be
144 made that there are indeed multiple ways to inter-
145 pret the phenomenon of interest (Aroyo and Welty,
146 2015), e.g., when different groups of annotators
147 reliably label media differently (Prabhakaran et al.,
148 2023; Homan et al., 2022), or when the task itself
149 is ambiguous (Artstein and Poesio, 2008).

150 A growing movement in the field of ground-
151 truthing has taken to viewing this variation in some
152 instances as being a necessary part of the phe-
153 nomenon of interest². More specifically, it is ar-
154 gued that the annotation projects occur on a contin-
155 uum: on one end are objective phenomena whose
156 interpretation is not expected to vary based on the
157 perspective of the annotator, and on the other are
158 phenomena where it is indeed expected to vary
159 based on the lived experience, feelings etc. of the
160 annotator (Cabitza et al., 2023). In some instances,
161 the expectation is that there will be multiple valid
162 labels for an item, that will systematically vary
163 based on the social group of the person who is
164 labelling it e.g. (Prabhakaran et al., 2023).

165 Although determining the degree of subjectivity
166 of a task is a challenge, and research is ongoing in
167 terms of appropriate methods and metrics to extract,
168 the Perspectivist approach advocates creating and
169 reporting disaggregated data to allow for a continu-
170 ous update as to knowledge on the dataset (Liem
171 and Demetriou, 2023).

172 2.2 Human values

173 Basic human values can be used to describe peo-
174 ple or groups: social science theory suggests that
175 each person uses a hierarchical list of values as
176 life-guiding principles (Rokeach, 1973). Most
177 widely used in social and cultural psychology is
178 the Schwartz theory, whose formal definition is
179 that values "(1) are concepts or beliefs, (2) pertain

²Sometimes referred to as the *Perspectivist manifesto*.

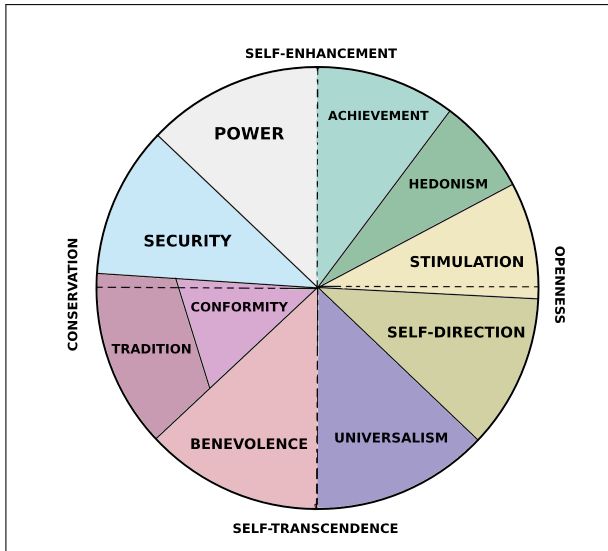


Figure 1: Visualization of the Schwartz 10-value inventory from (Schwartz, 1992) used in this paper, such that more abstract values of Conservation, vs. Openness to Change, and Self-transcendence vs. Self-enhancement form 4 higher-order abstract values. Illustration adapted from (Maio, 2010).

to desirable end states or behaviors, (3) transcend specific situations, (4) guide selection or evaluation of behavior and events, and (5) are ordered by relative importance" (Schwartz, 1992). Broadly speaking, values are abstract desirable goals that guide and motivate actions towards them across contexts (Sagiv and Schwartz, 2022).

The modern study of human values spans over 500 samples in nearly 100 countries over the past 30 years, and has shown a relatively stable structure (Sagiv and Schwartz, 2022), as illustrated in Figure 1. This has been observed across cultures in terms of the specific values present, and which values are prioritized together. Obtained scores across cultures also correlate with a broad range of impactful phenomena: e.g., cultures that value conservation and conforming to authority tend towards religiosity and away from openness and self-direction, altruistic behavior correlates with self-transcendent values like benevolence and universalism where competitiveness and unethical behavior correlate with self-enhancement goals like achievement and power, and right-wing political ideology correlates with tradition, conformity and security, where universalism values better predict left-wing ideology (Sagiv and Schwartz, 2022).

As such, the structure can be used to understand what individuals use to guide their actions, but also what entire populations prioritize when repre-

sentative samples are aggregated. In addition, the relative stability of the structure allows for a convenient method to estimate the reliability and validity of measurements in novel contexts: these should, in principle, show similar structure.

2.3 Human Values in Text

Schwartz suggests that we communicate our values in order to gain cooperation and coordinate our efforts (Schwartz, 1992). As such, this communication will likely manifest in the form of words in speech and text (Boyd and Pennebaker, 2017). A vast amount of text and speech is produced and consumed: every minute in 2022 an estimated 1 million hours of content were streamed, and over 350,000 tweets were shared³. Thus, studying values perceived in text both is relevant from a social sciences perspective (in order to understand behaviors and priorities of diverse groups of people), and from a computational perspective (given how much text is available).

Although some work estimating the values of the authors of text has been conducted, work on how values in text are perceived is lacking: novel attempts have been made to measure the values of individuals who have written personal essays and social media posts e.g. (Maheshwari et al., 2017; Ponizovskiy et al., 2020), and in arguments abstracted from various forms of public facing text (Kiesel et al., 2022). However, we have not observed work on how to measure values perceived in text, nor work that treats the estimated values in text as a hierarchical list, in line with theory (Rokeach, 1973). Further, how language is perceived may vary substantially depending on what group is perceiving it: e.g. perceptions have been shown to vary widely by group in terms of what language is harmful (Solaiman et al., 2023; Prabhakaran et al., 2023), and how emotions are described even when there is a common structure (Jackson et al., 2019).

2.4 Music Lyrics

Music listening is an extremely popular activity. Over 616 million people subscribe to streaming services worldwide⁴, and out of the music industry's reported 31.2 billion USD⁵ revenue, more

³<https://web-assets.domo.com/miyagi/images/product/product-feature-22-data-never-sleeps-10.png>

⁴<https://www.musicbusinessworldwide.com/files/2022/12/f23d5bc086957241e6177f054507e67b.png>

⁵<https://midiareserach.com/blog/recorded-music-market-2022-reality-bites>

than 17 billion comes from music streaming⁶. Out of over 1400 number-1 singles in the UK charts, only 30 were instrumental⁷. Lyrics were shown to be a salient component of music (Demetriou et al., 2018), and thus are likely to be a widely consumed form of text of importance to a broad audience. As reported in Appendix A.1, the responses of our annotation participants, who were drawn from a representative sample of the US population, quantitatively confirm the prevalence and importance of lyrics to them as music listeners, while lyrics at the same time indeed are a good data source for soliciting subjective judgements.

Although the relationship of values within the context of contemporary music remains relatively understudied, an interview study showed participant confirmation that personal values indeed play a role in people’s music preferences (Manolios et al., 2019). In addition, (Gardikiotis and Baltzis, 2012) showed correlations between higher order value scores derived from the Schwartz Value Survey (Schwartz, 1992) and a dimension-reduced set of musical preference scores. Although their model explained only a small portion of the variance, (Preniqi et al., 2022) showed correlations between participants’ self-reported moral values, and moral values estimated from the lyrical content of artists whose Facebook pages participants had liked.

3 Primary Lyric Data

Our aim is to collect a sample of lyric data where the lyrics are as accurate as possible, and our sample is as representative as possible. We sampled from the population of songs in the Million Playlist Dataset (MPD)⁸ as it is large and recent compared to other similar datasets. The lyrics themselves were obtained through the API of Musixmatch⁹, a lyrics and music language platform. Musixmatch lyrics are crowdsourced by users who add, correct, sync, and translate them. Musixmatch then engages in several steps to verify quality of content, including spam detection, formatting, spelling and translation checking, as well as manual verification by over 2000 community curators, and a local team of Musixmatch editors.

⁶https://cms.globalmusicreport.ifpi.org/uploads/Global_Music_Report_State_of_The_Industry_5650fff4fa.pdf

⁷https://en.wikipedia.org/wiki/List_of_instrumental_number_ones_on_the_UK_Singles_Chart

⁸<https://research.atspotify.com/2020/09/the-million-playlist-dataset-remastered/>

⁹<https://www.musixmatch.com/>

3.1 Fuzzy Stratified Song Sampling

An initial challenge is determining how to represent the population of songs when it is known to be very large¹⁰. An ideal scenario would be one in which we randomly sampled a known number of songs from a set of clearly defined strata (e.g., relevant subsets within the overall sample). However, for music, we do not know how many songs we would need to sample in order to reach saturation, what the relevant strata to randomly sample within should be, and how to measure relevant parameters from each stratum.

Some measurable strata that affect the use of language in the song lyrics are clear (e.g., the year of release, which may reflect different events or time-specific colloquial slang). Others are less clear: e.g., there is no single metric of popularity, although it can be estimated from various sources such as hit charts. Some may be subjective, such as genre, for which there may be some overlap of human labelling, but no clear taxonomy exists in the eyes of musicological domain experts (Liem et al., 2012). Based upon these considerations, we aim for a stratified random sampling procedure, based on strata that we acknowledge to be justifiable given our purpose, yet in some cases conceptually ‘fuzzy’: (1) release date; (2) popularity, as estimated via artist playlist frequency from the MPD (Chen et al., 2018); (3) genre, estimated from topic modeling on Million Song Dataset artist tags (Schindler et al., 2012); (4) topic, through a bag-of-words representation of the lyrics data.

To draw an initial subpopulation of songs, we first uniformly sub-sampled 60k out of 300k artists from the Million Playlist Dataset (MPD) (Chen et al., 2018). We then queried the Musixmatch API to determine if the lyrics for each of the songs of the 60k sample of artists was available.

3.2 Bias Correction

We expect that our dataset requires a bias correction. Given the skewness of data concentration with regard to several of our strata, songs that are recent and widely popular will most likely be drawn. To correct for this and get a more representative sample of an overall song catalogue, we oversample from less populated bins. For this, we use the maximum-a-posteriori (MAP) estimate of the categorical distribution of each stratum. The over-

¹⁰e.g., Spotify reports over 100 million songs in its catalogue <https://newsroom.spotify.com/company-info/>

sampling is controlled by concentration parameter a of the symmetric Dirichlet distribution. We heuristically set this parameter such that songs in under-populated bins still will make up up 5-10 % of our overall pool¹¹. Through this method, we subsampled 2200 songs with lyrics.

3.3 Inclusion Criteria

As the annotation of highly subjective perceived values in lyrics has not been studied yet, it is unclear whether any valid and reliable annotations can be obtained from it. As such, together with the ambition to investigate many annotations from a representative population sample, it may be unwise to immediately annotate thousands of songs, but rather focus on rich insights on smaller well-curated data. For this, the following screening procedure was followed. Three members of the research team manually screened several hundreds of songs randomly sampled from our 2200 songs. They verified the match of songs to lyrics, the available metadata, and rejected songs that had words that were not English, contained very few words, were only onomatopoeic, or were only repetitions. As a consequence, we finally kept 380 songs.

4 Survey Measures

Our annotation procedure seeks to obtain human perceptions on perceived values in song lyrics. To obtain such perceptions, we design a survey, in similar fashion to how in psychometrics, a survey will be designed as a measurement instrument for psychological constructs.

With many lyrics being in the English language, we choose to obtain our annotations from representative samples of the US population in terms of self-reported sex, ethnicity and age. Such samples can be obtained through the Prolific¹² platform. We follow Prolific's guidelines on fair compensation to set our compensation rates. Survey design and data handling were pre-discussed with our institutional data management and research ethics advisors, we obtained formal data management plan and human research ethics approval, and participants give informed consent before proceeding with the survey.

4.1 Lyrics affinity and subjectivity perception

To gain further measurable evidence on the degree to which song lyrics are important yet subjective

¹¹Full code of our sampling procedure is at https://anonymous.4open.science/r/lyrics-value-estimators-CE33/1_stimulus_sampling/stratified_sampling.py

¹²<https://prolific.co>

to a representative population sample, our survey starts with 16 general questions about song lyric preferences. Furthermore, after participants performed their annotations, we also ask them to rate how subjective they considered the task to be. As gaining a general understanding of these phenomena is not the main purpose of our current study though, but rather further substantiation supporting our current work, we will report on findings from these questions in Appendix A.1.

4.2 Short Schwartz Values Survey

Our primary annotations involve impressions of the values expressed in song lyrics. To this end, we adapted the Short Schwartz Values Survey (SSVS) (Lindeman and Verkasalo, 2005) to determine the wording of the questions, as it is the shortest instrument that has shown adequate reliability. The original wording of the questionnaire displays the name of the value being rated, followed by a number of words to describe it e.g. "POWER (social power, authority, wealth)"¹³. Original instructions can be found in Appendix A.2.

We made three adaptations to this questionnaire. First, we adjust the question text to ask not for ratings of life-guiding principles for the individual responding to the survey, but rather for the respondent's impressions of the 'speaker' of the lyrics. This 'speaker' is the someone or something whose perspective is reflected through the lyrics, and this may not be the author or artist expressing the lyrics. For example, the speaker in the song 'I gave you power' by the artist Nas is a gun, and the speaker in 'Rosetta Stoned' by the rock band Tool is a person hallucinating from psychedelics. In other words, the creator may use a persona in the writing of song lyrics for artistic purposes, which may not directly represent their values. As such, an annotator's impression of the creator may differ from their impression of the person represented by the lyrics. As we are interested in the values perceived in the text, we explicitly ask participants to respond with the perspective of the speaker in mind, and not the author. Further illustration of our explicit instructions is given in Appendix A.2.

Secondly, the original SSVS uses a 9 point Likert-Scale where 0="Opposed to My Principles", 1="Not Important", and 8="Of Supreme Importance". In order to gather continuous measure-

¹³Actual wording of items was retrieved from <https://blogs.helsinki.fi/everyday-thinking/files/2015/11/The-Short-Schwartzs-Value-Survey.docx>.

ments, we aimed for a continuous scale where 0 essentially indicates that the value is either not discussed or otherwise not estimatable from the lyric text. Next to this, we balance the scale to have maximum opposition to a given value be at -100, where maximum importance will be at 100. As such, in contrast to the original SSVS, our scale is symmetric with a rating of 0 indicating neutrality.

Thirdly, (Cabitza et al., 2020) suggested that a rater’s confidence is an indication of intra-rater reliability. Thus, we also asked participants “How confident are you in your ratings of these lyrics?”, to which they responded on a scale of 0 (Not at all Confident) to 100 (Completely Confident).

4.3 Annotation Interface

The survey was implemented on an instance of the Qualtrics¹⁴ platform. The annotations were collected using the response format shown as illustrated in Figure 10, following explicit instructions as discussed in Appendix A.2. More specifically, a set of lyrics are displayed, with a clickable interface below them. The interface contains brief descriptions of each of the 10 Schwartz values, followed by a vertical bar on which participants can indicate a continuous response, as described in Section 4.2. An option to select “Not Applicable” was also available for each value. We considered that “0” and “Not Applicable” responses both indicate that the importance of that given value to the speaker based on the lyrics could not be determined by the participant (i.e., they were either not discussed in the song lyrics, or were otherwise unclear). As we expect that not all songs will discuss all values, and most songs may discuss very few values, we initialize the rating bar at “0”.

With this, we now have the setup to gather annotation data. In the remainder of this paper, we discuss how this was done to research three questions: (1) How many annotator ratings are needed for stable annotations to emerge? (2) Do our obtained value perception annotations relate to existing validated knowledge on stable structures among values? and (3) Can our refined annotations be used in computational NLP setups?

5 Number of Ratings

Our procedure to determine the number of ratings to gather was inspired by (DeBruine and Jones, 2018). Specifically, we first recruited a represen-

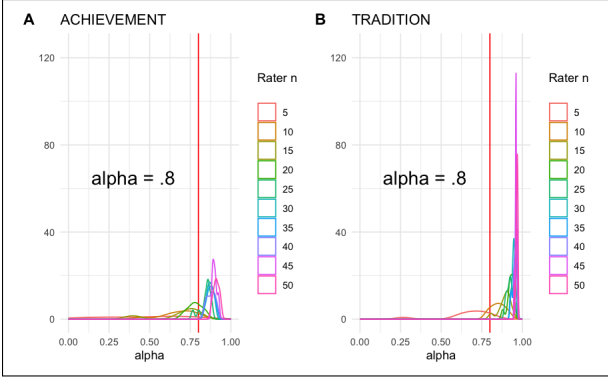


Figure 2: Distribution of Cronbach’s α from a representative US Sample (N=505) rating 20 songs, for the values Achievement and Tradition. Vertical line represents the α threshold for comparison.

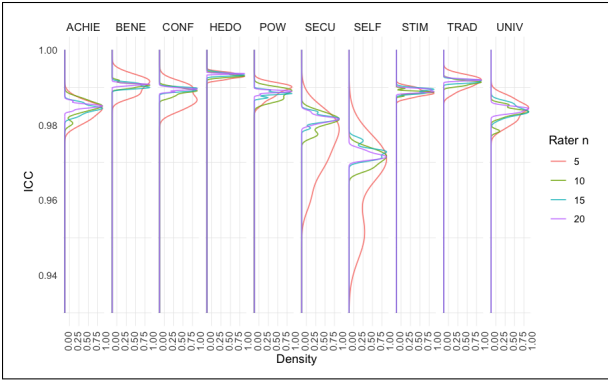


Figure 3: Rotated scaled density plots of ICC for subsamples from annotations on the 360 songs.

tative pilot sample (N=505), in which respondents used our interface to annotate perceived values for a fixed set of 20 songs. From these annotations, we computed canonical mean ratings per value, per song. For each of the values, we then estimated Cronbach’s α for a range of subsample sizes (5 to 50 participants, in increments of 5), repeating this procedure 10 times per increment. Following this, we visually examined density plots of the distribution of Cronbach’s α . In the social sciences, an $\alpha \geq 0.7$ is commonly considered an acceptable level of reliability. Taking a conservative estimate, we choose to obtain 25 ratings per song lyric in our main study; for that amount of ratings, Cronbach’s α in our pilot data would comfortably exceed 0.8.

From this, we perform our main study data collection. We recruit a new representative US population sample (N=600), where each participant goes through our survey questions, and receives 18 randomly selected song lyrics to annotate for perceived values. As a result, we obtained 22-30

¹⁴<https://www.qualtrics.com/>

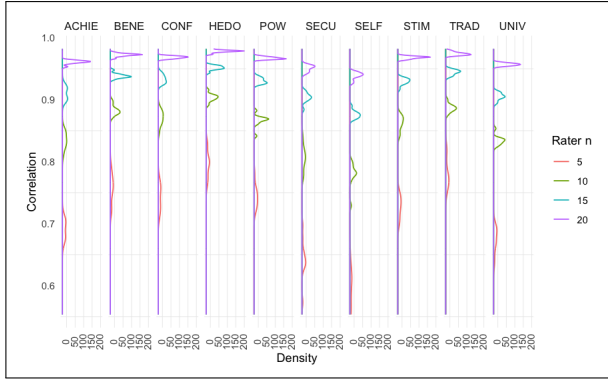


Figure 4: Rotated scaled density plots of Pearson correlations between canonical mean and subsample means, from a mean 27 ratings per each of the 360 songs.

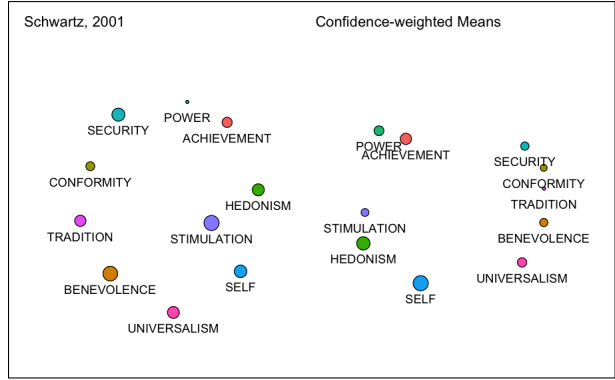


Figure 5: MDS plots derived from the correlation plot reported in (Schwartz et al., 2001), and our participant responses as confidence-weighted means¹⁵.

508 annotations per song, with an average of 27.

509 From these, checking for the reliability of our
 510 annotations from this sample, we repeatedly sub-
 511 sample 5, 10, 15 and 20 ratings for each value
 512 within each song, and calculate and visualize intra-
 513 class correlations (ICC, Figure 3), as well as Person
 514 correlations between subsample means and canonical
 515 means (Figure 4). From this, we see that higher
 516 numbers of raters lead to higher ICCs and Pearson
 517 correlation, thus indicating more stable outcomes
 518 that approach the canonical mean that would be
 519 obtained from a large sample of annotators. Seeing
 520 Pearson correlation to the canonical mean already
 521 exceeds 0.9 for all values from 15 subsampled rat-
 522 ings, our target of 25 annotations per songs indeed
 523 can be considered too conservative, and in future
 524 work, 15 ratings on average will likely suffice.

525 For our further analysis, we will have to aggre-
 526 gate the subjective labels. Being unaware of a single
 527 ideal method to achieve this, for the purposes
 528 of this work, we report results using an aggregation
 529 method inspired by (Cabitzza et al., 2020). Specifi-
 530 cally, we estimate confidence-weights by dividing
 531 participant’s self-reported confidence of a given
 532 rating by the highest possible response (100), and
 533 then compute aggregated means weighted by these.

534 6 Structural comparison

535 As a first attempt to assess the relative validity of
 536 our procedure, we depart from the earlier observa-
 537 tions that cross-cultural stable structure was found
 538 on what values are likely to cluster together. We
 539 compare distances as derived from the upper trian-
 540 gle of a correlation matrix reported in (Schwartz
 541 et al., 2001) to those derived from proximity in
 542 ratings obtained in our study. For both, we gen-

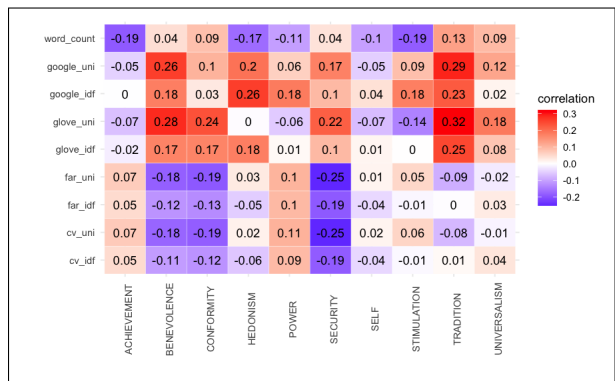


Figure 6: Pearson correlations between NLP systems / word counts, and participant ratings of songs, by value.

543 erate a multi-dimensional scaling plot (MDS) for
 544 visual comparison, which has previously been used
 545 as method to assess confirmation of earlier theory
 546 (Ponizovskiy et al., 2020). From these plots
 547 (Figure 5, in as little as our 360 annotated lyrics,
 548 we surprisingly indeed see similar clusters and relative
 549 positioning relations emerging as those obtained
 550 from a formal cross-cultural study.

551 7 Comparisons with NLP models

552 Finally, as first step towards ways in which our data
 553 may be connected to computational NLP methods,
 554 we perform a preliminary comparison on how compu-
 555 tational NLP-based value assessment of lyrics
 556 data compares to the way in which our annotators
 557 annotated perceived human values.

558 We depart from a dictionary of words associated
 559 with the 10 Schwartz values (Ponizovskiy et al.,
 560 2020). With this dictionary as reference, we
 561 computationally perform assessments to the degree
 562 to which each value would be reflected in the lyrics
 563 text according to traditional word counting (Poni-

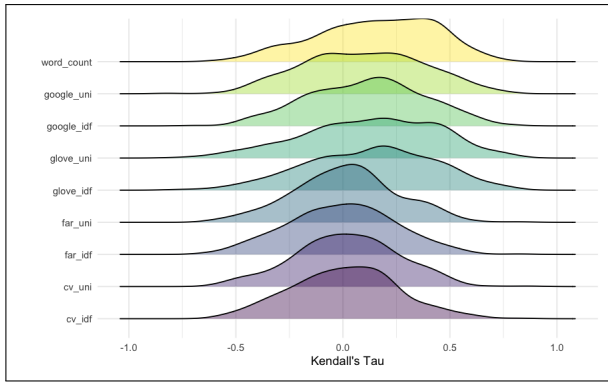


Figure 7: Rank correlations between NLP systems / word counts and confidence-weighted participant means transformed to rankings

zovskiy et al., 2020), as well as by assessing cosine similarity between dictionary words and lyrics texts using four classes of pre-trained word embeddings: word2vec-google-new, a generic English word embedding trained on Google News dataset (Mikolov et al., 2013); glove-common-crawl, another generic English word embedding trained on Common Crawl dataset (Pennington et al., 2014); faruqui-mxm-[1~10], trained on the collected initial lyrics candidate pool, employing the Glove model (Pennington et al., 2014) (using ten models populated from ten cross-validation folds, whose parameters are tuned based on English word similarity judgement data (Faruqui and Dyer, 2014).); and cv-mxm-[1~10], ten variants of lyrics based word-embeddings from cross-validation folds selected by Glove loss values on the validation set.

We weigh terms in the lyrics texts in two different ways: uniformly and weighted by Inverse Document Frequency (IDF). Then, we compare value assessments from these computational methods to the ones obtained from our annotators, in two ways. First, we consider Pearson correlations of computational value similarity assessments to our raw participant song value annotations (Figure 6). Second, we take the earlier-theorized perspective that value assessments should be seen as ranked lists, and we consider rank correlations between the machine and human value assessments based on Kendall's τ (Figure 7).

In earlier work (Richard et al., 2003), Pearson correlations of 0.1-0.2 were considered as moderate evidence of the validity of a proposed dictionary in relation to a psychometrically valid instrument. As such, considering our data as good reference data,

only the more generic Glove and Google news embeddings seem to reach those levels of correlation. From a rank correlation perspective, the word count methods and these two embedding models hint at slightly positive rank correlations. This may be promising in terms of the degree of specialization needed to assess values; at the same time, neither of the methods presented here have thoroughly been optimized, and as such, these results should not be seen as strong benchmarking evidence. Future work will be needed to more deeply connect computational NLP techniques with our data.

8 Limitations and future work

In this paper, we described our procedure for ground-truthing perspectives on highly subjective text. By paying attention to grounding in social sciences theory and purposeful sampling strategies, the discussion of how to get to 'good data' has been much more extensive than commonly is done in computational domains. With this, we hope to have illustrated how beyond (welcome) completion of checklists and data sheets, being purposeful about data can pro-actively shape annotation design.

As for limitations to our work, while we are committed to open science practices, we cannot share the primary lyric data due to copyright prohibitions. However, we do release metadata of the songs of interest, together with our participant annotations, and the code used for the analyses and plots in our paper¹⁶. We acknowledge our current sample of 360 lyrics is still small and may need expansion, and that, while we had a representative population sample, not every member of the sample rated every song. We thus did gather diverse opinions, but cannot claim they fully represent the target population. We also did not assess whether variations on the annotation instrument might result in substantial differences in the annotations we received (Kern et al., 2023), nor did we repeat our procedure (Inel et al., 2023). In addition, we can further connect our work to related research on examining how participants from different groups will annotate corpora (Homan et al., 2022; Prabhakaran et al., 2023). Finally, while we only provide a preliminary comparison to computational NLP methods, it will be worthwhile to use our data in the context of more sophisticated state-of-the-art NLP systems.

¹⁶https://anonymous.4open.science/r/values_in_lyrics-8F3F/

References

Lora Aroyo and Chris Welty. 2015. Truth is a lie: Crowd truth and the seven myths of human annotation. *AI Magazine*, 36(1):15–24.

Ron Artstein and Massimo Poesio. 2008. Inter-coder agreement for computational linguistics. *Computational linguistics*, 34(4):555–596.

Valerio Basile et al. 2020. It’s the end of the gold standard as we know it. on the impact of pre-aggregation on the evaluation of highly subjective tasks. In *CEUR WORKSHOP PROCEEDINGS*, volume 2776, pages 31–40. CEUR-WS.

Abeba Birhane, Pratyusha Kalluri, Dallas Card, William Agnew, Ravit Dotan, and Michelle Bao. 2022. The values encoded in machine learning research. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 173–184.

Ryan L Boyd and James W Pennebaker. 2017. Language-based personality: A new approach to personality in a digital world. *Current opinion in behavioral sciences*, 18:63–68.

Federico Cabitza, Andrea Campagner, and Valerio Basile. 2023. Toward a perspectivist turn in ground truthing for predictive computing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 6860–6868.

Federico Cabitza, Andrea Campagner, and Luca Maria Sconfienza. 2020. As if sand were stone. new concepts and metrics to probe the ground on which to build trustable ai. *BMC Medical Informatics and Decision Making*, 20(1):1–21.

Ching-Wei Chen, Paul Lamere, Markus Schedl, and Hamed Zamani. 2018. Recsys challenge 2018: Automatic music playlist continuation. In *Proceedings of the 12th ACM Conference on Recommender Systems*, pages 527–528.

Lisa M DeBruine and Benedict C Jones. 2018. [Determining the number of raters for reliable mean ratings](#).

Andrew Demetriou, Andreas Jansson, Aparna Kumar, and Rachel M Bittner. 2018. Vocals in music matter: the relevance of vocals in the minds of listeners. In *ISMIR*, pages 514–520.

Manaal Faruqui and Chris Dyer. 2014. [Community evaluation and exchange of word vectors at wordvectors.org](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, System Demonstrations*, pages 19–24. The Association for Computer Linguistics.

R. Michael. Furr and Verne R. Bacharach. 2014. *Psychometrics : an introduction*, second edition edition. SAGE Publications.

Antonis Gardikiotis and Alexandros Baltzis. 2012. ‘rock music for myself and justice to the world!’: Musical identity, values, and music preferences. *Psychology of Music*, 40(2):143–163. 699
700
701
702

Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. [Datasheets for datasets](#). *Communications of the ACM*, 64(12):86–92. 703
704
705
706
707

R Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang. 2020. Garbage in, garbage out? do machine learning application papers in social computing report where human-labeled training data comes from? In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 325–336. 708
709
710
711
712
713
714

Robert M Groves, Floyd J Fowler Jr, Mick P Couper, James M Lepkowski, Eleanor Singer, and Roger Tourangeau. 2009. *Survey methodology*, volume 561. John Wiley & Sons. 715
716
717
718

Christopher Homan, Tharindu Cyril Weerasooriya, Lora Aroyo, and Chris Welty. 2022. Annotator response distributions as a sampling frame. In *Proceedings of the 1st Workshop on Perspectivist Approaches to NLP@ LREC2022*, pages 56–65. 719
720
721
722
723

Han-Yin Huang and Cynthia C. S. Liem. 2022. Social inclusion in curated contexts: Insights from museum practices. In *FACCT ’22: 2022 ACM Conference on Fairness, Accountability, and Transparency*. 724
725
726
727

Oana Inel, Tim Draws, and Lora Aroyo. 2023. Collect, measure, repeat: Reliability factors for responsible ai data collection. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 11, pages 51–64. 728
729
730
731
732

Joshua Conrad Jackson, Joseph Watts, Teague R Henry, Johann-Mattis List, Robert Forkel, Peter J Mucha, Simon J Greenhill, Russell D Gray, and Kristen A Lindquist. 2019. Emotion semantics show both cultural variation and universal structure. *Science*, 366(6472):1517–1522. 733
734
735
736
737
738

Abigail Z. Jacobs and Hanna Wallach. 2021. Measurement and fairness. In *FACCT ’21: 2021 ACM Conference on Fairness, Accountability, and Transparency*. 739
740
741

Eun Seo Jo and Timnit Gebru. 2020. Lessons from archives: Strategies for collecting sociocultural data in machine learning. In *Conference on Fairness, Accountability and Transparency (FAT ’20) January 27-30 2020, Barcelona, Spain*. ACM, New York, NY, USA. 742
743
744
745
746
747

Christoph Kern, Stephanie Eckman, Jacob Beck, Rob Chew, Bolei Ma, and Frauke Kreuter. 2023. Annotation sensitivity: Training data collection methods affect model performance. *arXiv preprint arXiv:2311.14212*. 748
749
750
751
752

753	Johannes Kiesel, Milad Alshomary, Nicolas Handke, Xiaoni Cai, Henning Wachsmuth, and Benno Stein. 2022. Identifying the human values behind arguments. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4459–4471. Association for Computational Linguistics.	811
754		812
755		813
756		814
757		815
758		816
759		817
760	Benjamin Kompa, Jasper Snoek, and Andrew L Beam. 2021. Second opinion needed: communicating uncertainty in medical machine learning. <i>NPJ Digital Medicine</i> , 4(1):4.	818
761		819
762		820
763		821
764	Cynthia C. S. Liem, Andreas Rauber, Thomas Lidy, Richard Lewis, Christopher Raphael, Joshua D. Reiss, Tim Crawford, and Alan Hanjalic. 2012. Music Information Technology and Professional Stakeholder Audiences: Mind the Adoption Gap . In <i>Dagstuhl Follow-Ups</i> , volume 3. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik.	822
765		823
766		824
767		825
768		826
769		827
770		828
771	Cynthia CS Liem and Andrew M Demetriou. 2023. Treat societally impactful scientific insights as open-source software artifacts. In <i>2023 IEEE/ACM 45th International Conference on Software Engineering: Software Engineering in Society (ICSE-SEIS)</i> , pages 150–156. IEEE.	829
772		830
773		831
774		832
775		833
776		834
777	Marjaana Lindeman and Markku Verkasalo. 2005. Measuring values with the short schwartz’s value survey. <i>Journal of personality assessment</i> , 85(2):170–178.	835
778		836
779		837
780	Tushar Maheshwari, Aishwarya N Reganti, Samiksha Gupta, Anupam Jamatia, Upendra Kumar, Björn Gambäck, and Amitava Das. 2017. A societal sentiment analysis: Predicting the values and ethics of individuals by analysing social media content. In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers</i> , pages 731–741. Association for Computational Linguistics.	838
781		839
782		840
783		841
784		842
785		843
786		844
787		845
788		846
789	Gregory R Maio. 2010. Mental representations of social values. In <i>Advances in experimental social psychology</i> , volume 42, pages 1–43. Elsevier.	847
790		848
791		849
792	Sandy Manolios, Alan Hanjalic, and Cynthia CS Liem. 2019. The influence of personal values on music taste: towards value-based music recommendations. In <i>Proceedings of the 13th ACM Conference on Recommender Systems</i> , pages 501–505.	850
793		851
794		852
795		853
796		854
797	Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality . In <i>Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States</i> , pages 3111–3119.	855
798		856
799		857
800		858
801		859
802		860
803		861
804		862
805	Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In <i>Proceedings of the Conference on Fairness, Accountability, and Transparency (FAT* ’19)</i> . ACM.	863
806		864
807		865
808		866
809		867
810		868
	Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL</i> , pages 1532–1543. ACL.	869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

864 Shalom H Schwartz. 1992. Universals in the content
 865 and structure of values: Theoretical advances and
 866 empirical tests in 20 countries. In *Advances in exper-*
 867 *imental social psychology*, volume 25, pages 1–65.
 868 Elsevier.

869 Shalom H Schwartz, Gila Melech, Arielle Lehmann,
 870 Steven Burgess, Mari Harris, and Vicki Owens. 2001.
 871 Extending the cross-cultural validity of the theory
 872 of basic human values with a different method of
 873 measurement. *Journal of cross-cultural psychology*,
 874 32(5):519–542.

875 Irene Solaiman, Zeerak Talat, William Agnew, Lama
 876 Ahmad, Dylan Baker, Su Lin Blodgett, Hal
 877 Daumé III, Jesse Dodge, Ellie Evans, Sara Hooker,
 878 et al. 2023. Evaluating the Social Impact of Generative
 879 AI Systems in Systems and Society. *arXiv*
 880 *preprint arXiv:2306.05949*.

881 Chris Welty, Praveen K. Paritosh, and Lora Aroyo. 2019.
 882 Metrology for AI: from benchmarks to instruments.
 883 *arXiv preprint arXiv:1911.01875*.

884 A Appendix

885 A.1 Lyrics affinity and subjectivity perception

886 Our data collection protocols allowed us to gather
 887 self reports on the importance of lyrics. Our initial
 888 pool of questions was inspired by the Preference
 889 Intensity scale in (Schäfer and Sedlmeier, 2009),
 890 and consisted of Likert-type questions. We turn
 891 these into 16 question statements on the partici-
 892 pants’ relation to music lyrics by also adding our
 893 own suggested questions. Participants respond to
 894 the questions using a 5-point Likert scale, which in-
 895 cluded the points “Strongly Disagree”, “Somewhat
 896 Disagree”, “Neither Agree nor Disagree”, “Some-
 897 what Agree”, and “Strongly Agree”. Percentages
 898 in the table below indicate the proportion of respon-
 899 dents that indicated either “Somewhat Agree” and
 900 “Strongly Agree”.

901 We currently report on responses given to these
 902 questions from our two data collection rounds: our
 903 pilot study (N=505), whose primary aim was the
 904 estimation of the number of ratings needed per
 905 song lyric, and our actual annotation collection
 906 study (N=600) on which the main outcomes in our
 907 paper are reported.

908 In Figure 8, we visualize self-reported percent-
 909 ages of respondents’ music libraries containing
 910 lyrics, for both our respondent samples. Here, we
 911 see that respondents’ music overwhelmingly con-
 912 tains lyrics, with a median of 90%. Furthermore,
 913 in Table 1, for both samples of respondents, we
 914 indicate the percentages of users that indicated to
 915 somewhat or strongly agree with given statements.

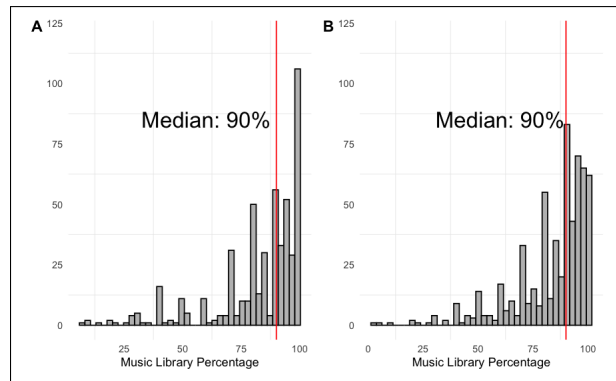


Figure 8: Distribution of self-reported percentage of music library containing lyrics from two representative US samples, N=505 and N=600 respectively.

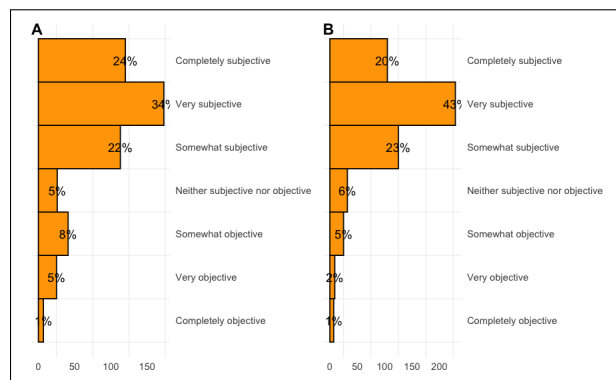


Figure 9: Distribution of self-reported subjectivity of lyric annotation task, N=505 and N=532 respectively.

916 From this, we again observe a strong preference
 917 for songs with lyrics (>70% on many of the state-
 918 ments).

919 Finally, at the end of our survey, we also ask par-
 920 ticipants to self-report a rating of the subjectivity of
 921 the lyrics annotation task we gave to them. Distri-
 922 butions are visualized in Figure 9. From these, we
 923 see confirmed the task indeed is perceived as highly
 924 subjective in the eyes of our sample population.

925 As gaining a general understanding of music
 926 lyrics affinity is not our current main goal, we chose
 927 not to iteratively validate and refine our questions
 928 as a formal psychometric instrument at this stage
 929 (for this, more explicit iterative analysis would be
 930 needed on the instrument being capable of distin-
 931 guishing between different types of users by mak-
 932 ing use of the full scale). However, we did start
 933 analyzing to what extent the current questions may
 934 be used as an instrument, or at least as a way to fur-
 935 ther characterize sub-populations of human respon-
 936 dents. Here, given the large preference towards
 937 music that contains lyrics, asking for lyrics vs. non-

Question	Pilot	Main
I prefer music that contains lyrics, as opposed to music that does not	72%	72%
I always pay attention to the lyrics of a song, if the song has them	70%	72%
If a song has lyrics that I don't like for any reason, I don't listen to it	49%	43%
If I am not sure about the lyrics of a song, I search them on the internet	76%	77%
I memorize the lyrics to the songs I listen to	70%	75%

Table 1: Question wording, and proportion of respondents rounded to the nearest whole number, that indicated either 'somewhat agree' or 'strongly agree' in two surveys, N=505, and N=600 respectively.

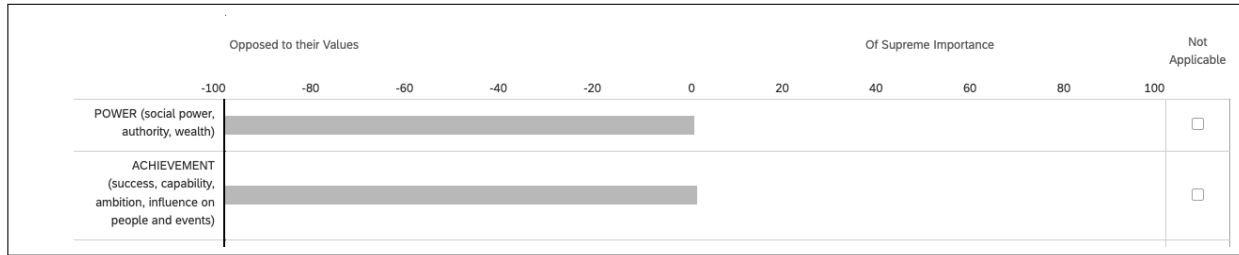


Figure 10: Visualization of the annotation interface on Qualtrics for two of ten annotated values

lyrics music preference will not allow for us to be able to distinguish between respondents. At the same time, responses to the degree to which a respondent pro-actively engages with lyrics (e.g. by actively searching for them, writing about them, or writing lyrics themselves) may yield interpretable factors on which respondents can be distinguished. However, we leave a deeper analysis of this for future work.

A.2 Adjusted Short Schwartz Value Survey

The original Schort Schwartz Value survey appears in (Lindeman and Verkasalo, 2005). The original question wording¹⁷ was:

"Please, rate the importance of the following values as a life-guiding principle for you. Use the 8-point scale in which 0 indicates that the value is opposed to your principles, 1 indicates that the values is not important for you, 4 indicates that the values is important, and 8 indicates that the value is of supreme importance for you."

- POWER (social power, authority, wealth)
- ACHIEVEMENT (success, capability, ambition, influence on people and events)
- HEDONISM (gratification of desires, enjoyment in life, self-indulgence)

- STIMULATION (daring, a varied and challenging life, an exciting life)
- SELF-DIRECTION (creativity, freedom, curiosity, independence, choosing one's own goals)
- UNIVERSALISM (broad-mindedness, beauty of nature and arts, social justice, a world at peace, equality, wisdom, unity with nature, environmental protection)
- BENEVOLENCE (helpfulness, honesty, forgiveness, loyalty, responsibility)
- TRADITION (respect for tradition, humbleness, accepting one's portion in life, devotion, modesty)
- CONFORMITY (obedience, honoring parents and elders, self-discipline, politeness)
- SECURITY (national security, family security, social order, cleanliness, reciprocation of favors)

In our survey, participants were initially shown a set of instructions designed to explain how to use the instrument, and explain our working definitions of 'artist' as separate from the 'speaker' of the lyrics, see(Figure 11). We then presented our adjusted question wording:

"Between the quotation marks below are some song lyrics. Please take a moment to read them

¹⁷retrieved from <https://blogs.helsinki.fi/everyday-thinking/files/2015/11/The-Short-Schwartzs-Value-Survey.docx>.

990 and think about the SPEAKER the lyrics. Please
991 remember that this SPEAKER might be a the AU-
992 THOR themselves, or someone or something else:",
993 after which lyrics were displayed, along with the
994 annotation instrument.

Thanks!

You will now be shown parts of song lyrics from 18 songs, and asked to complete some questions about how you perceive them.

IMPORTANT: Lyrics can be written from different perspectives, some of which are not the same as the writer of the lyrics. In other words, the **AUTHOR** of the lyrics may choose a **SPEAKER** for their lyrics that is not themselves.

The **SPEAKER** of the lyrics could be a fictional character, a real person from history or the present, or even an imaginary object. And of course it could be the **AUTHOR** themselves. Please answer the questions while thinking about the **SPEAKER**.

WARNING: These lyrics are drawn from popular music, some of which use offensive language or describe offensive situations.



Click on the bar (green arrow) to indicate how important you think the value is to the **SPEAKER**.

Clicking further to the right indicates that you think it is a supremely important guiding principle in their life.

Further to the left indicates that you think it is the opposite of that guiding principle of their life.

Click on 'Not Applicable' (red arrow) or 0 if you don't think the **SPEAKER** has not expressed anything about that value. They indicate the same thing.

Please click on the arrow below when you're ready to proceed.

Figure 11: Visualization of the instructions page of the annotation interface on Qualtrics.