

VALUE KALEIDOSCOPE :

Engaging AI with Pluralistic Human Values, Rights, and Duties

Taylor Sorensen^{♦◊}, Liwei Jiang^{♦◊}, Jena D. Hwang[◊], Sydney Levine[◊],
 Valentina Pyatkin^{♦◊}, Peter West^{♦◊}, Nouha Dziri[◊], Ximing Lu^{♦◊}, Kavel Rao[♦],
 Chandra Bhagavatula[◊], Maarten Sap^{♦◊}, John Tasioulas[†], Yejin Choi^{♦◊}

[♦]Department of Computer Science & Engineering, University of Washington, [◊]Allen Institute for Artificial Intelligence,

[†]Language Technologies Institute, Carnegie Mellon University, [†]Department of Philosophy, University of Oxford
 {tsor13, yejin}@cs.washington.edu

Abstract

Human values are crucial to human decision-making. *Value pluralism* is the view that multiple correct values may be held in tension with one another (e.g., when considering *lying to a friend to protect their feelings*, how does one balance *honesty* with *friendship*?). As statistical learners, AI systems fit to averages by default, washing out these potentially irreducible value conflicts. To improve AI systems to better reflect *value pluralism*, the first-order challenge is to explore the extent to which AI systems can model pluralistic human *values, rights, and duties* as well as their interaction.

We introduce VALUEPRISM, a large-scale dataset of 218k values, rights, and duties connected to 31k human-written situations. VALUEPRISM’s contextualized values are generated by GPT-4 and deemed high-quality by human annotators 91% of the time. We conduct a large-scale study with annotators across diverse social and demographic backgrounds to try to understand whose values are represented.

With VALUEPRISM, we build VALUE KALEIDOSCOPE (or KALEIDO), an open, light-weight, and structured language-based multi-task model that *generates, explains, and assesses* the *relevance* and *valence* (i.e., support or oppose) of human values, rights, and duties within a specific context. Humans prefer the sets of values output by our system over the teacher GPT-4, finding them more accurate and with broader coverage. In addition, we demonstrate that KALEIDO can help explain variability in human decision-making by outputting contrasting values. Finally, we show that KALEIDO’s representations can transfer to other philosophical frameworks and datasets, confirming the benefit of an explicit, modular, and interpretable approach to value pluralism. We hope that our work will serve as a step to making more explicit the implicit values behind human decision-making and to steering AI systems to make decisions that are more in accordance with them.

1 Introduction

When people confront difficult decisions (whether or not to break a promise, what degree program to enroll in, how to spend a Sunday afternoon), their options reflect their values (friendship, knowledge, freedom, saving money, spending time in nature). Two people in the same situation may make opposing decisions if they value different things or the same things but to varying extents (Figure 1). The notion that



Figure 1: Different human values relate, support, or oppose everyday situations to varying degrees. KALEIDO is designed to generate, explain, and assess how the pluralistic human values, rights, and duties may shape human judgments.

different human values can lead to distinct—although potentially equally valid—decisions is called *value pluralism* (Páez et al. 2020; Komppula et al. 2018; Brosch and Sander 2013; Keeney 1992; Griffiths 2021; Liscio et al. 2023).

Various fields have focused on this concept. Philosophers distinguish *value pluralism* (different views cannot be reduced into an ultimate “supravalue” (Williams 1985; Larimore 1987; Kekes 1993; Stocker 1990; Chang 1997; Dancy 2004)) from *monism* (there exists a single core value (Kant 1785/2002; Driver 2022)). Sociologists recognize cultural, social, and ideological differences that drive societal clashes, movements, and changes (Archive 2011). Psychologists empirically confirm that ethical experiences involve weighing pluralistic values (Gill and Nichols 2008) and the dissonance that arises from misaligned values and beliefs (Festinger 1962).

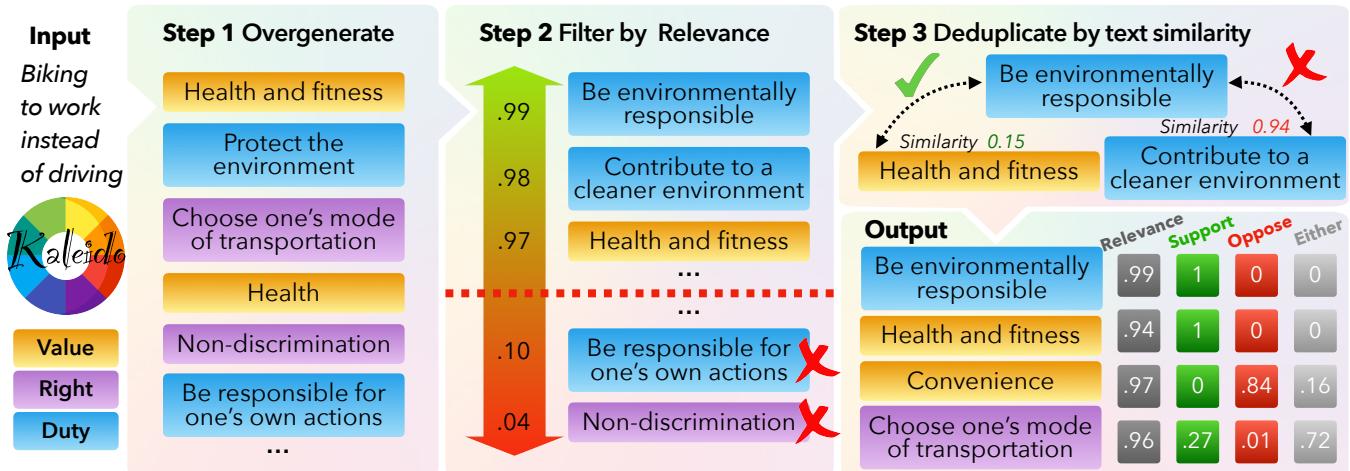


Figure 2: KALEIDO^{SYS} system workflow that includes 1) generating 100 values, rights and duties; 2) filtering by relevance as rated by KALEIDO; 3) removing repetitive items; and computing relevance and valence scores for each value, right, and duty.

Meanwhile, in AI, there is a growing interest in developing human-centered AI that emphasizes participation from stakeholders. This which necessitates the inclusion and exploration of pluralistic voices and values (Tasioulas 2022; Gordon et al. 2022). Yet, contemporary supervised AI systems primarily wash out variation by aggregating opinions or preferences with majority votes (Plank 2022; Talat et al. 2022; Casper et al. 2023). As real-world AI applications are used to assist increasing and more diverse audiences, it is crucial to investigate and better model the the values that are accessible and used by current AI systems.

In this work, we make the first large-scale attempt at investigating large language models' (LLMs) potential to model *pluralistic human values, rights, and duties*. Our effort is twofold: (1) we introduce VALUEPRISM, a large-scale dataset of pluralistic human values; (2) we build VALUE KALEIDOSCOPE (KALEIDO), an open and flexible value-pluralistic model.

The dataset: VALUEPRISM contains 218k contextualized values, rights, and duties connected to 31k real-life situations, distilled from GPT-4¹. While GPT-4 and its like have been shown to match human crowdworker annotation performance in some domains (Gilardi, Alizadeh, and Kubli 2023; Ziems et al. 2023; Rytting et al. 2023), we exercise caution and do not assume that GPT-4's outputs are necessarily correct or representative. To this end, we conduct large-scale human studies and find that humans rate the outputs as high-quality 91% of the time and have difficulty coming up with considerations that the model has missed, detecting missing values >1% of the time. We also conduct a comprehensive study with diverse annotators across diverse social and demographic groups to evaluate whose voices are represented in the values GPT-4 produces. Additionally, a growing line of work demonstrates that the large-scale with which data can be produced with LLMs can make up for the

potential noise that is introduced, leading to student models which often surpass the teacher (West et al. 2022b; Kim et al. 2023; Jung et al. 2023).

The model: VALUE KALEIDOSCOPE (KALEIDO) is a value-pluralistic model that, based on VALUEPRISM, generates, explains, and assesses the relevance and valence (i.e., support or oppose) of contextualized pluralistic human values, rights, and duties. On top of the model, we build a flexible system KALEIDO^{SYS} leveraging KALEIDO's generation and relevance prediction modes to create a diverse, high quality set of relevant values for a situation (See Fig. 2). In human studies, people rate our system's outputs as more correct and complete than the teacher's (GPT-4). Annotators also find that our largest model matches the teacher's performance at rationalizing and predicting valence. Additionally, we show that KALEIDO can help explain ambiguity and variability underlying human decision-making in nuanced situations by generating contrasting values. We also demonstrate that KALEIDO can be adapted to various philosophical frameworks without explicit training.

Overall, our work represents the first comprehensive attempt to articulate decision-making into fine-grained, pluralistic components of human values employing large language models. The resulting dataset and model serve as a large-scale resource explicitly supporting value pluralism, shedding light on future AI development that accommodates a rich and inclusive tapestry of value alternatives.²

2 Value-pluralistic Framework: Values, Rights and Duties

2.1 Why are pluralistic human values critical?

Machine learning methods are generally designed to model averages, but can miss nuance and in-group variation unless explicitly accounted for (Gordon et al. 2022). To go beyond this, we take inspiration from philosophical value

¹Datasheet for Datasets (Gebru et al. 2018) documentation in Appendix N

²<https://github.com/tsor13/kaleido>.

pluralism, the stance that there are many different normative values (Mason 2006), as opposed to one super-value that all other values can be reduced to. This is distinct both from political pluralism, which posits that diversity is beneficial to democratic society and supports the distribution of power among diverse groups (Britannica Editors 2002); and from relativism, which holds that no moral system is more correct than another (Gowans 2021).

Without taking a hard stance on these positions, we seek to better model humans’ plural values to make explicit the implicit values in human decision-making. Our hope is that, if pluralistic values can be adequately (though imperfectly) modeled, we can take a step towards ensuring that automated decision-makers act in accordance with them.

2.2 Framework Motivation and Definition

In this work, we model human-centered plural values to make explicit implicit values in human decision-making. We settle on *values* (Mason 2006), *rights* (Prabhakaran et al. 2022; Wenar 2023), and *duties* (Alexander and Moore 2021) as our three core concepts. We propose a commonsense framework for reasoning about them, and outline them below.

Values: These are the *intrinsic goods or ideals* that people pursue or cherish, such as happiness, well-being, justice, or freedom. Values are the desirable qualities that people may seek in their lives and in the world. They are often the guiding principles for individuals and societies, shaping goals, motivations, and preferences.

Duties: Duties are the *moral obligations or responsibilities* that individuals owe to others or to society at large. They are categorical reasons for doing or refraining from doing something, independent of whether we want to do or refrain from doing that thing. Duties can be weighty reasons, not easily overridden by competing concerns, and their violation may justify blame and self-blame (guilt). Duties can arise from relationships, social roles, or moral principles, and they guide our actions and decisions.

Rights: Rights are the *entitlements or claims* that individuals have against others or society, which are usually based on moral or legal grounds. These can be positive rights (e.g., the right to education, healthcare, or free speech) or negative rights (e.g., the right to not be harmed, enslaved, or discriminated against). Rights serve to protect the fundamental interests of individuals and establish certain boundaries that others must respect.

3 KALEIDO: Value-pluralistic Modeling

We introduce KALEIDO, a language-based multi-task system that *generates*, *explains*, and *assesses* the relevance and valence (i.e., support or oppose) of pluralistic human values, rights, and duties, grounded in real-world contexts.

3.1 Tasks

We develop four tasks for modeling values, rights, and duties, all grounded in a given context situation.

Generation (open-text) *What values, rights, and duties are relevant for a situation?* Generate a value, right, or duty that could be considered when reasoning about the action.

Event Driving at the maximum speed of the road		
Value	Efficiency Safety Rule-following Law-abiding behavior	support oppose oppose oppose
Right	Right to freedom of movement Right to safety Right to life and physical integrity	either oppose oppose
Duty	Duty to follow traffic laws	oppose
Event Gambling with my boss		
Value	Professionalism Loyalty Personal enjoyment Financial well-being	oppose oppose support oppose
Right	Right to work in a safe and professional environment Autonomy	oppose either
Duty	Duty to maintain a professional relationship with your boss Duty to uphold company policies	oppose oppose
Event someone cuts me off. i honk ten times		
Value	Courtesy Respect for others Safety Harmony	oppose oppose support oppose
Right	Right to peaceful enjoyment of the road Right to not be subjected to harassment Right to safety	oppose oppose oppose
Duty	Duty to express displeasure Duty to be a considerate driver	support oppose

Table 1: Example outputs from KALEIDO^{SYS}.

Relevance (2-way classification) *Is a value relevant for a situation?* Some values are more relevant than others.

Valence (3-way classification) *Does the value support or oppose the action, or might it depend on context?* Disentangling the valence is critical for understanding how plural considerations may interact with a decision.

Explanation (open-text) *How does the value relate to the action?* Generating a post-hoc rationale for why a value consideration may relate to a situation.

3.2 Dataset: VALUEPRISM

We leverage the symbolic knowledge distillation (West et al. 2022a) pipeline to distill high-quality knowledge from powerful generative models like GPT-4, which have been shown to compare favorably to human annotations on quality, coverage, and diversity (West et al. 2022a; Gilardi, Alizadeh, and Kubli 2023; Ziems et al. 2023). Importantly, based on

our preliminary exploration, GPT-4 excels at enumerating a *wide* range of value alternatives compared to average human annotations. We verify the dataset’s quality with human annotators and show that 91% of the distilled data is deemed high quality, surpassing typical quality of human generated data (West et al. 2022a; Hwang et al. 2021; Zhou et al. 2023). Details on dataset statistics and splits are provided in App. F.1 and examples from VALUEPRISM can be found in App. A.

Situations Starting from 1.3M human-written base situations³, we obtain a set of 31K clean situations for deriving plural considerations by carefully filtering out ill-formatted, irrelevant, and low-quality situations using few-shot Flan-T5 (Chung et al. 2022). To balance out an outsize proportion of toxic, NSFW, or sexually explicit content, we down-sample these situations to 5% of all data, leading to an increase in the overall diversity of the dataset, as measured by the normalized count of unique n-grams (dist-2: .23→.36, dist-3: .54→.67, details in App. F.1).

Values, Rights, and Duties Generation For each of the 31K situations, we prompt GPT-4 to generate relevant values, rights, and duties, with open-text rationales. GPT-4 also attributes whether the corresponding value, right, or duty supports (justifies), opposes (condemns), or depends on the context or interpretation. Details of data generation and prompting are in Appendix F.1 and M. The resulting dataset is rated as high-quality (91%) by human evaluation (§4.1).

Multi-task Setup We convert VALUEPRISM into a sequence-to-sequence format for multi-task training. The relevance task data is sampled contrastively, where positive examples are from GPT-4 and negative examples are drawn from values mapped to other situations. We split the data (by actions) into 80% train/10% validation/10% test.

3.3 Model Training

We use Flan-T5 3B (Chung et al. 2022) as the base model for KALEIDO. After fitting hyperparameters with a grid search, we train on a mixture of all four tasks with Huggingface’s Trainer (Wolf et al. 2020) for 4 epochs with early stopping. Training takes 19 hours on two A100 GPUs. Training details are in App. G.

3.4 A System of Diverse Values: KALEIDO^{SYS}

We use KALEIDO to generate a diverse set of values, rights, and duties by overgenerating (top 100 beams) and removing low-quality and repetitive outputs via the relevance score and text similarity respectively. We use Rouge-score (Lin 2004) for n-gram similarity and a Transformers (Wolf et al. 2020) mpnet model⁴ for sentence embeddings. See Fig. 2 for an illustration of the system and Appendix H/Algorithm 1 for more details. We tune the system parameters (relevance score threshold, similarity thresholds) using Gibbs sampling (Casella and George 1992) to maximize RougeL-Sum F1 score on the validation set. Ablation experiments in §5.1

³Sourced from user queries submitted to the Delphi demo.

⁴<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

provide insights on each system component, and example model outputs can be found in Table 1 and App. B.

4 Data Analysis

4.1 VALUEPRISM is high-quality

We conduct human validation of VALUEPRISM to assess its quality on the Mechanical Turk platform⁵. Given the generated situation and values, rights, and duties and their explanations, we ask the annotators to assess the relevance and quality of the generations. The results show that annotators find the great majority of the data as high quality. 91% of the values/rights/duties was marked as good by all three annotators and 87% of the valences were marked as correct by all three annotators.

In an attempt to find any values that may have been missed, we also prompt crowdworkers to fill in any missing values, rights, or duties. Crowdworkers didn’t seem to find it easy to come up with missing values as we get suggestions 0.35% of the time. Full annotation details for this and other studies are in App. I.

4.2 Evaluation by diverse annotators

Prior research has reported unjust bias in LLMs against marginalized groups (Sap et al. 2019; Feng et al. 2023). We evaluate VALUEPRISM by recruiting a diverse population of 613 annotators⁶ through CloudResearch (Litman, Robinson, and Abberbock 2017) targeting those marginalized groups to the extent possible⁷. We collect 31k annotations across 683 values, rights, and duties in the context of 100 situations, along with demographic information across eight categories. The annotators mark 1) if they agree with each value, right, or duty listed for a given situation and 2) if they spot any missing perspective. We do not find notable statistical significance, and do not reject the null hypothesis that there is no difference between groups. Additional group statistics, p-values, and qualitative analyses are in App. E.

4.3 Diversity of VALUEPRISM

We analyze the diversity of the situations, and values, rights, and duties from three perspectives: *lexical diversity* that calculates uniqueness n-grams, *topical diversity* that assesses semantic diversity via topic analysis⁸, and *clustering*. Both

⁵For this and other human studies, we have acquired the opinion of our institution’s Internal Review Board. The opinion finds our project exempt from a full review process in accordance to the United States Federal regulation 45 CFR 46, and we have acquired a letter of exception. We hash crowdworker IDs so annotations cannot be back-traced to individual workers.

⁶E.g., Race: 168 white, 115 Black, 61 asian, 34 hispanic/latinx; Sexual orientation: 390 straight, 68 LGBQ+. Gender: 258 male, 201 female, 9 non-binary or other; Full details are in App. E

⁷We chose CloudResearch specifically because of its ability to target by demographic. One weakness of this study, however, is that all of our respondents are U.S.-based (where CloudResearch operates). Prior work has shown that value representation can vary across nationality as well (Santy et al. 2023), and we hope to extend this study internationally in the future.

⁸Via BERTopic [https://maartengr.github.io/BERTTopic](https://maartengr.github.io/BERTopic)

Model	Overall	Cover.	Acc.	Avg. #
KALEIDO ^{SYS} 3B	55.5	65.1	58.9	8.2
-relevance	51.9	81.4	64.3	11.2
-text similarity	50.0	60.5	52.9	8.2
verbose	58.0	86.1	69.0	11.1
concise	39.0	27.4	32.4	5.0
KAL ^{SYS} 11B	58.3	71.1	62.5	8.3
KAL ^{SYS} 770M	57.9	67.3	60.8	8.2
KAL ^{SYS} 220M	44.9	59.0	50.8	8.1
KAL ^{SYS} 60M	32.0	53.0	37.1	8.5
Direct Output	42.5	37.9	40.0	6.8
GPT-4	50.0	50.0	50.0	7.0
GPT-3.5-turbo	39.5	49.0	39.8	8.0

Table 2: The overall, coverage and accuracy win rate percentage against GPT-4 by human evaluators along with the average number of generated values, rights, and duties. (Here and throughout, best results within 1% are bolded.)

the situations and the core value concepts contain diverse distinct situations with high lexical variations indicating that a diverse variety of events and values captured by VALUEPRISM (Table 8). The topic word cloud (Fig. 8) shows that VALUEPRISM covers a broad spectrum of common topics like “save”, “kill”, and “helping” for situations and “respect”, “care”, and “promote” for values. Clustering shows that the corpus encompass a wide variety of themes, reflecting the diversity and richness of situations and values, rights, and duties. For more data analysis, see App. C.

5 Experiments

5.1 Our System Against the Teacher

Generating correct and complete sets of values Central to our research is the capability to model pluralistic values, rights, and duties. Ideally, these values should be correct, have high coverage, and be aligned with human preferences. We recruit crowdworkers to evaluate KALEIDO^{SYS} directly against GPT-4 across these three dimensions.

We run several variations of KALEIDO^{SYS}: all model sizes (60M–11B); 3B version without the relevance and text similarity components (-relevance, -text similarity); and tweaked 3B system parameters to output more or fewer values, rights, and duties⁹ (verbose, concise). To understand the benefit of our system, we also train a baseline seq2seq 3B model on the same data that tries to predict a batch of values, rights, and duties in one generation pass. We test each version against GPT-4 on a set of 200 test situations by 2 annotators each.

From Table 2, we make several observations. The three largest versions of our system outperform GPT-4 on all evaluated dimensions, with the largest variant (11B) being the most favored overall. Moreover, the models generating a higher number of values (>11) are preferred by humans

⁹To better understand how changing the parameters can affect the output/precision/recall, see Figure 4

Model	Explanation	Valence	Rel. corr.
KALEIDO 3B	92.6	92.0	0.30
KAL 11B	94.8	92.6	0.25
KAL 770M	90.3	90.3	0.31
KAL 220M	86.9	86.3	0.30
KAL 60M	75.9	72.3	0.28
GPT-4	94.7	93.1	-

Table 3: Human Evaluation. Explanation and Valence scores are correctness rates of the output, while Relevance is the correlation of relevance score with the percentage of people who marked a value as relevant.

for coverage and accuracy.¹⁰ Our system also shows an advantage over the direct output seq2seq model, despite being trained on the same data. Furthermore, removing relevance leads to a drop in the overall preference, which is not observed in verbose with the same number of outputs. This suggests relevance is indeed a contributing factor to the generation quality. Finally, humans show lower preference for outputs without deduplication with text similarity.

Explanation and Valence Label Quality We also evaluated the explanation generation and valence labeling abilities of each model using 700 predefined values, rights, and duties from the test set of VALUEPRISM. Crowdworkers were tasked with evaluating the quality of explanations, their effectiveness in linking values to actions, and agreement with valence labels. As depicted in Table 3, the 11B model’s performance closely aligns with that of GPT-4. The 11B model achieved Valence accuracy within a 1% difference from GPT-4 and slightly outperformed it in terms of Explanation quality.

5.2 Relevance Correlates with Human Judgments

We would like KALEIDO to predict whether a human would find a value, right, or duty relevant. However, its training data is synthetic, so the model instead predicts whether a given value was likely to be generated for a particular situation by GPT-4. To test how well this correlates with how humans judge relevance, we collect 18 annotations each for 700 values/rights/duties for relevance, and correlate the relevance score with the percentage of people who marked the value as relevant (See Table 3). We find correlations of 0.25–0.31 for the suite of model sizes¹¹ (all significant at $p < 10^{-10}$). While we would like to explicitly train models to predict human relevance scores in future work, we take this as evidence that our synthetic relevance score correlates positively with human judgments.

¹⁰This is in line with prior work showing that humans prefer longer outputs with more unique n-grams (Wang et al. 2023b)

¹¹Interestingly, we note that the correlation does not strictly improve with model size. While we are unsure of the reason for this, we note that 11B gives much more confident relevance scores, and hypothesize that this overconfidence may be miscalibrated to human judgments.

5.3 Zero-shot performance on ETHICS

While our model is explicitly trained to recognize values, rights, and duties, we want to understand how much the learned representations generalize to other frameworks as well. To do this, we test KALEIDO on the ETHICS benchmark (Hendrycks et al. 2023), which contains crowdsourced ethical judgments across several different frameworks, fitting templates with values, rights, or duties that loosely correspond to the frameworks (see Appendix L).

Subset	KALEIDO	ChatGPT	Random
Justice	17.5 / 13.3	17.6 / 13.4	6.3 / 6.3
Deont.	19.8 / 15.1	20.6 / 13.8	6.3 / 6.3
Virtue	33.1 / 22.2	24.9 / 22.0	8.2 / 8.2
Util.	76.5 / 66.6	59.4 / 55.1	50.0 / 50.0
Comm.	71.5 / 64.7	80.3 / 68.8	50.0 / 50.0
Average	43.7 / 36.4	40.6 / 34.6	24.2 / 24.2

Table 4: ETHICS few-shot performance. First/second number of each entry is performance on the test/hard test sets respectively. KALEIDO is zero-shot, ChatGPT is few-shot.

Results are in Table 4. On all five tasks, our model performs well over the random baseline. On all tasks but Commonsense, our model matches or exceeds (Justice, Deont., Virtue, Util.) ChatGPT’s performance, while only having 3B parameters. Despite having only been trained to predict values, rights, and duties, our model meaningfully generalizes to other frameworks.

5.4 Interpretable Decision System and Zero-Shot on COMMONSENSENORMBANK

While the focus of the system is on modeling diverse values and not on making judgments, it can be easily extended to output the valence of an action $V(a)$:

$$V(a) = \sum_{v \in VRD} R(v|a) \times V(v|a)$$

where $v \in VRD$ are the generated values, rights, and duties from KALEIDO^{SYS}, $R(v|a)$ is the relevance of v given the action, and $V(v|a)$ is the valence of v given the action. We will denote this decision system KALEIDO^{DEC}.

This system has the advantage of being interpretable, enabling direct inspection of how values linearly contribute to the outcome. It is also steerable, as users can easily assign a weight of zero to values they do not wish to take into consideration.

Zero-shot COMMONSENSENORMBANK performance
We evaluate this system in a zero-shot manner on the four subportions of moral acceptability segment of COMMONSENSENORMBANK (Jiang et al. 2022) (See Table 5). In all cases, the system performs at least as well as the majority class baseline, and much ($\geq 25\%$) better on ETHICS and Moral Stories¹².

¹²For these two datasets, there is no “neutral” (i.e., lacks valence) class, so the “either” valence is zeroed out.

Model	SBIC	ETH.	MoSt	SoCh
KALEIDO ^{DEC}	64.4	77.9	75.4	48.2
+label calibration	69.3	78.0	76.2	63.0
(improvement)	(+4.9)	(+0.1)	(+0.8)	(+14.8)
Majority class	63.1	51.6	50.0	46.7
Random	33.3	50.0	50.0	33.3
Delphi (SFT)	82.9	86.2	86.5	78.0

Table 5: Zero-shot Performance on COMMONSENSENORMBANK: Moral Acceptability.

We observe that the classes are not well calibrated to the dataset statistics. To remedy this, we fit a lightweight logistic regression. For SBIC and SocialChem it improves accuracy by about 5% and 15% respectively, suggesting that while the model is not initially well-calibrated to the datasets, relevant information can be linearly extracted. While KALEIDO^{DEC} achieves non-trivial zero-shot performance, it unsurprisingly performs worse than supervised baselines such as Delphi.

5.5 Entropy as an Indicator of Decision Variability

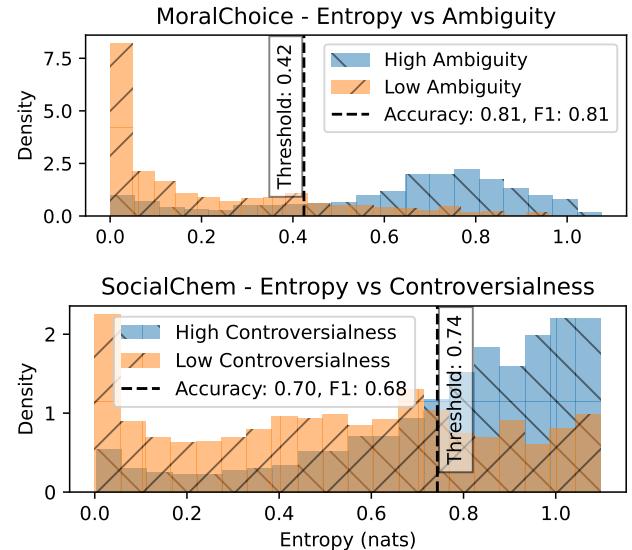


Figure 3: The output entropy of KALEIDO^{DEC} is predictive of ambiguity in MoralChoice and controversialness in SocialChem. A threshold is chosen to maximize F1-score.

When values support different decisions, it may be an indicator that the final judgment one may come to is highly dependent on which value is prioritized. Because of this, when KALEIDO^{DEC} output has high entropy, we hypothesize that this may indicate higher variability in the distribution of decisions. To test this, we explore two datasets with variability indicators. MORALCHOICE (Scherrer et al. 2023) contains 687 low-ambiguity and 680 high-ambiguity moral scenarios. SOCIALCHEM (Forbes et al. 2021) is a corpus of social

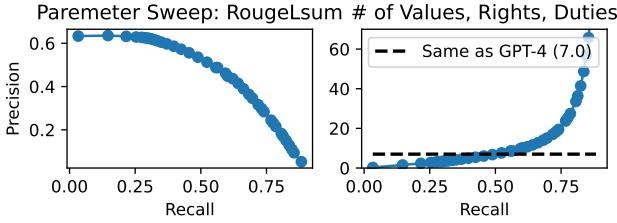


Figure 4: By sweeping KALEIDO^{SYS}'s parameters, we are able to trade precision (w.r.t. to GPT-4) with recall, and output many more (or fewer) values, rights, and duties.

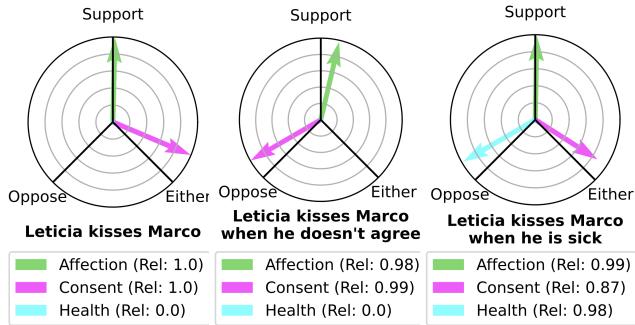


Figure 5: KALEIDO is sensitive to subtle changes in inputs, changing relevance and valence scores accordingly.

norms where, among other things, crowdworkers annotated for "What portion of people probably agree that [action] is [good / bad]?" We take those marked as $\geq 99\%$ to have low controversialness, and those marked as $\leq 50\%$ as having high controversialness. We run the corresponding scenarios through KALEIDO^{DEC} and measure the entropy (Figure 3). We find that the entropy is predictive of these classes. In line with our hypothesis, the higher the entropy, the more likely a situation is to be ambiguous or controversial, even though the model was not trained to predict these things explicitly.

6 Discussion

Strengths Over Teacher Model Although our model performs strongly against the teacher in value generation, it also boasts several other advantages. It is controllable, allowing users to generate either more or fewer values than GPT-4 by trading precision for recall (see Figure 4). Additionally, while GPT-4 provides only textual labels for valence, our model generates scalar valence and relevance scores. Lastly, our model, dataset, and code are openly accessible, enabling scientific review that is crucial for accountability and improvement.

KALEIDO is Sensitive to Situation Variations One of the strengths of our approach is that the signal can be contextualized and dependent on the relevant values. For example, consider three variations of a situation: "Leticia kisses Marco," "Leticia kisses Marco when he doesn't agree," and "Leticia kisses Marco when he is sick" (see Figure 5). In all three situations, affection and consent are relevant values,

as reflected by their relevance scores. However, the valence changes: consent can either support or oppose the action in the two underspecified situations, but opposes "when Marco doesn't agree." Additionally, the value of health isn't usually relevant in the typical context of kissing: however, "when Marco is sick," health becomes relevant and opposes the action. This demonstrates the ability of KALEIDO to adjust to subtle input changes¹³.

False Balance and Extreme Inputs One potential danger when generating diverse values is coming up with a contrived reason why something might be good or bad, even if no reasonable person may hold such a value in such a situation (This is similar to false balance, or bothesidesism, in news reporting (Imundo and Rapp 2021; Boykoff and Boykoff 2004)). To probe at this, we hand-write 20 situations (10 bad/10 good, in Appendix J) for which we cannot come up with reasonable values, rights, or duties that would support or oppose them respectively. We run them through KALEIDO^{SYS} after development and find no generated supporting values/rights/duties for the extreme bad actions nor any opposing for the good actions. We take this as limited evidence that our system can avoid false balance.

Universal Declaration of Human Rights Inspired by (Prabhakaran et al. 2022), we think that an ideal dataset containing human rights would contain all rights listed in the United Nation's Universal Declaration of Human Rights¹⁴ (UDHR). We manually extract all human rights we could find from the UDHR (41), and find the 20 closest rights in the dataset by entailment score with WANLI (Liu et al. 2022a). We then go through all 41 sets manually and label each for whether or the right is included. We are able to find matches in VALUEPRISM for 97.5% of the UDHR's human rights, demonstrating that the dataset has broad coverage of the UDHR.¹⁵

7 Related Work

Value Representations of Language Models Scholars from diverse disciplines have engaged in extensive discussions regarding the incorporation of human ethics and values into LLMs (Wallach and Allen 2008; Jiang et al. 2022; Hendrycks et al. 2023), understanding cultural influences (Santy et al. 2023), examining opinion alignment (Santurkar et al. 2023), and using LLMs as proxies for studying specific human sub-populations in social science research (Argyle et al. 2023). (Jiang et al. 2022) introduced Delphi, a framework trained to reason about ethical perspectives, and showed the ethical limitations of out-of-the-box LLMs. Another noteworthy dimension is the multicultural nature of LLMs. (Santy et al. 2023) explored the cultural disparities manifest in LMs and their implication for diverse societies. (Tasioulas 2022) criticized the prevailing preference-based utilitarian approach (i.e., which act is likely to yield the optimal fulfillment of human preferences) in AI ethics, pointing

¹³While this is not a quantitative experiment, this is not a cherry-picked example.

¹⁴<https://www.un.org/en/about-us/universal-declaration-of-human-rights>

¹⁵See App. K for all human rights and corresponding matches.

out its limitations and proposing an alternative “humanistic” ethics that accounts for the nuanced human views on ethics. Moreover, Santurkar et al. (2023); Durmus et al. (2023) introduced novel opinion datasets, quantitatively analyzed the opinions conveyed by LMs, and unveiled substantial misalignments between the viewpoints of current LMs and specific demographic groups within the United States.

Alignment of Large Language Models Several computational approaches have been proposed to address the challenge of aligning LLMs with desired values and objectives. Reinforcement learning (RL) has historically been used in multiple NLP tasks to ensure that the generated text is optimized for an arbitrary non-differentiable reward (Johnson et al. 2017; Nguyen, Daumé III, and Boyd-Graber 2017; Ramamurthy et al. 2022; Pyatkin et al. 2023). Lu et al. (2022) optimized a reward function that quantifies an undesired property, while not straying too far from the original model via a KL-divergence penalty. (Bai et al. 2022) explored RL techniques for training LLMs to adhere to legal and ethical guidelines encoded in a constitution, naming it “Constitutional AI”. Wu et al. (2023) used fine-grained human feedback as an explicit training signal to train and learn from reward functions in a RLHF fashion. Additionally, Lu et al. (2023) proposed an inference-time algorithm to efficiently tailor LLMs without no fine-tuning, addressing tasks like ensuring safety and fidelity in dialogue models.

Automatic Dataset Curation Previous research in automatic data generation has focused on creating datasets for various tasks, such as commonsense reasoning (West et al. 2022a; Bhagavatula et al. 2023; Wang et al. 2023a; Kim et al. 2023), dialogues (Kim et al. 2023; Xu et al. 2023; Chiang et al. 2023), summarization (Sclar et al. 2022; Jung et al. 2023), and contextual reasoning about offensive statements (Zhou et al. 2023). (West et al. 2022a) introduce the symbolic knowledge distillation framework, which has been extended in subsequent studies through iterative distillation (Sclar et al. 2022; Jung et al. 2023; Bhagavatula et al. 2023). In addition, (Liu et al. 2022b) propose a human-AI collaboration approach to generate high-quality datasets with challenging examples.

8 Conclusion

In this work, we contribute VALUEPRISM and KALEIDO in the hopes of leading to better value-pluralistic modeling. We validate VALUEPRISM’s quality with two human studies, and find that KALEIDO outperforms the teacher’s strong performance at generating relevant values, rights, and duties for a given situation. We also show that KALEIDO can help explain variability in human decisions and generalizes to data and frameworks outside of its training scope.

9 Ethical Considerations and Limitations

Machine-Generated Data. We use GPT-4’s open-text generative capabilities to collect VALUEPRISM, leveraging the wide variety of knowledge about human values, rights, and duties latent in LLM’s pretraining data. However, we also recognize that in doing so we run the potential for introducing the majority’s bias: the generated data may be limited to

the values of certain majority groups. In an effort to assess the extent of value plurality and representation, we make a deliberate effort to conduct the validation of the VALUEPRISM by collecting annotations from annotators of various social and demographic backgrounds as described in §4.2. The human annotators find the majority of our data as high-quality at a high agreement rate. Additionally, less than 1% of the validated situations were found to be lacking. Nevertheless, a more extensive study that focuses on the type and nature of value plurality covered by VALUEPRISM remains a compelling direction for future research.

Intended Use. We make VALUEPRISM openly available by individual request with the hope and intention that it furthers research in value pluralism in NLP and AI. However, it is possible that our data can be used in malicious and unintended application (e.g., speech policing or promotion of certain values). We do not endorse its use in such capacity and emphasize that the use of our dataset and model should be limited to research purposes only. Additionally, we limit the data and model available only by individual request to try to prohibit non-research use cases and ensure fair use.

Acknowledgments

The authors thank Ronan LeBras, Jared Moore, Hyunwoo Kim, Jenny Liang, and Sebastin Santy for helpful discussions; Alane Suhr for the example situation in Figure 2; Jared Moore, Dhruba Ghosh, and David Atkinson for draft feedback; and Michael Wilson, Michael Guerquin, and John Borchardt from the AI2 ReViz team for help with the demo. This research was supported in part by DARPA under the ITM program (FA8650-23-C-7316) and the Allen Institute for AI.

References

- Alexander, L.; and Moore, M. 2021. Deontological Ethics. In Zalta, E. N., ed., *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2021 edition.
- Allen Institute for AI. 2023. AI2 ImpACT Licenses. <https://allenai.org/impact-license>. [accessed August 2023].
- Archive, G. D. 2011. European Values Study Longitudinal Data File 1981-2008 (EVS 1981-2008). *EVS*.
- Argyle, L. P.; Busby, E. C.; Fulda, N.; Gubler, J. R.; Rytting, C.; and Wingate, D. 2023. Out of One, Many: Using Language Models to Simulate Human Samples. *Political Analysis*, 1–15.
- Bai, Y.; Kadavath, S.; Kundu, S.; Askell, A.; Kernion, J.; Jones, A.; Chen, A.; Goldie, A.; Mirhoseini, A.; McKinnon, C.; Chen, C.; Olsson, C.; Olah, C.; Hernandez, D.; Drain, D.; Ganguli, D.; Li, D.; Tran-Johnson, E.; Perez, E.; Kerr, J.; Mueller, J.; Ladish, J.; Landau, J.; Ndousse, K.; Lukosuite, K.; Lovitt, L.; Sellitto, M.; Elhage, N.; Schiefer, N.; Mercado, N.; DasSarma, N.; Lasenby, R.; Larson, R.; Ringer, S.; Johnston, S.; Kravec, S.; Showk, S. E.; Fort, S.; Lanham, T.; Telleen-Lawton, T.; Conerly, T.; Henighan, T.; Hume, T.; Bowman, S. R.; Hatfield-Dodds, Z.; Mann, B.; Amodei, D.; Joseph, N.; McCandlish, S.; Brown, T.; and Kaplan, J. 2022. Constitutional AI: Harmlessness from AI Feedback. arXiv:2212.08073.
- Bhagavatula, C.; Hwang, J. D.; Downey, D.; Bras, R. L.; Lu, X.; Qin, L.; Sakaguchi, K.; Swayamdipta, S.; West, P.; and Choi, Y. 2023. I2D2: Inductive Knowledge Distillation with NeuroLogic and Self-Imitation. arXiv:2212.09246.
- Boykoff, M. T.; and Boykoff, J. M. 2004. Balance as bias: global warming and the US prestige press. *Global Environmental Change*, 14(2): 125–136.
- Britannica Editors. 2002. Pluralism — Ideology, Diversity & Tolerance — britannica.com. <https://www.britannica.com/topic/pluralism-politics>. [Accessed 07-08-2023].
- Brosch, T.; and Sander, D. 2013. Neurocognitive mechanisms underlying value-based decision-making: from core values to economic value. *Frontiers in Human Neuroscience*, 7.
- Casella, G.; and George, E. I. 1992. Explaining the Gibbs Sampler. *The American Statistician*, 46(3): 167–174.
- Casper, S.; Davies, X.; Shi, C.; Gilbert, T. K.; Scheurer, J.; Rando, J.; Freedman, R.; Korbak, T.; Lindner, D.; Freire, P.; Wang, T.; Marks, S.; Segerie, C.-R.; Carroll, M.; Peng, A.; Christoffersen, P.; Damani, M.; Slocum, S.; Anwar, U.; Siththanjanan, A.; Nadeau, M.; Michaud, E. J.; Pfau, J.; Krasheninnikov, D.; Chen, X.; Langosco, L.; Hase, P.; Biyik, E.; Dragan, A.; Krueger, D.; Sadigh, D.; and Hadfield-Menell, D. 2023. Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback. arXiv:2307.15217.
- Chang, R. 1997. *Incommensurability, Incomparability, and Practical Reason*. Cambridge, MA, USA: Harvard.
- Chiang, W.-L.; Li, Z.; Lin, Z.; Sheng, Y.; Wu, Z.; Zhang, H.; Zheng, L.; Zhuang, S.; Zhuang, Y.; Gonzalez, J. E.; Stoica, I.; and Xing, E. P. 2023. Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* ChatGPT Quality.
- Chung, H. W.; Hou, L.; Longpre, S.; Zoph, B.; Tay, Y.; Fedus, W.; Li, Y.; Wang, X.; Dehghani, M.; Brahma, S.; Webson, A.; Gu, S. S.; Dai, Z.; Suzgun, M.; Chen, X.; Chowdhery, A.; Castro-Ros, A.; Pellat, M.; Robinson, K.; Valter, D.; Narang, S.; Mishra, G.; Yu, A.; Zhao, V.; Huang, Y.; Dai, A.; Yu, H.; Petrov, S.; Chi, E. H.; Dean, J.; Devlin, J.; Roberts, A.; Zhou, D.; Le, Q. V.; and Wei, J. 2022. Scaling Instruction-Finetuned Language Models. arXiv:2210.11416.
- Dancy, J. 2004. *Ethics Without Principles*. New York: Oxford University Press.
- Driver, J. 2022. The History of Utilitarianism. In Zalta, E. N.; and Nodelman, U., eds., *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2022 edition.
- Durmus, E.; Nyugen, K.; Liao, T. I.; Schiefer, N.; Askell, A.; Bakhtin, A.; Chen, C.; Hatfield-Dodds, Z.; Hernandez, D.; Joseph, N.; et al. 2023. Towards Measuring the Representation of Subjective Global Opinions in Language Models. arXiv preprint arXiv:2306.16388.
- Feng, S.; Park, C. Y.; Liu, Y.; and Tsvetkov, Y. 2023. From Pretraining Data to Language Models to Downstream Tasks: Tracking the Trails of Political Biases Leading to Unfair NLP Models. arXiv:2305.08283.
- Festinger, L. 1962. Cognitive dissonance. *Sci. Am.*, 207(4): 93–102.
- Forbes, M.; Hwang, J. D.; Shwartz, V.; Sap, M.; and Choi, Y. 2021. Social Chemistry 101: Learning to Reason about Social and Moral Norms. arXiv:2011.00620.
- Gebru, T.; Morgenstern, J.; Vecchione, B.; Vaughan, J. W.; Wallach, H. M.; Daumé, H.; and Crawford, K. 2018. Datasheets for datasets. *Communications of the ACM*, 64: 86 – 92.
- Gilardi, F.; Alizadeh, M.; and Kubli, M. 2023. ChatGPT outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30).
- Gill, M. B.; and Nichols, S. 2008. Sentimentalist Pluralism: Moral Psychology and Philosophical Ethics. *Philosophical Issues*, 18(1): 143–163.
- Gordon, M. L.; Lam, M. S.; Park, J. S.; Patel, K.; Hancock, J.; Hashimoto, T.; and Bernstein, M. S. 2022. Jury Learning: Integrating Dissenting Voices into Machine Learning Models. In *CHI Conference on Human Factors in Computing Systems*. ACM.
- Gowans, C. 2021. Moral Relativism. In Zalta, E. N., ed., *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Spring 2021 edition.
- Griffiths, N. 2021. Personal Values & Decision-Making Biases.
- Hendrycks, D.; Burns, C.; Basart, S.; Critch, A.; Li, J.; Song, D.; and Steinhardt, J. 2023. Aligning AI With Shared Human Values. arXiv:2008.02275.

- Hwang, J. D.; Bhagavatula, C.; Bras, R. L.; Da, J.; Sakaguchi, K.; Bosselut, A.; and Choi, Y. 2021. COMET-ATOMIC 2020: On Symbolic and Neural Commonsense Knowledge Graphs. arXiv:2010.05953.
- Imundo, M.; and Rapp, D. 2021. When Fairness is Flawed: Effects of False Balance Reporting and Weight-of-Evidence Statements on Beliefs and Perceptions of Climate Change. *Journal of Applied Research in Memory and Cognition*, 11.
- Jiang, L.; Hwang, J. D.; Bhagavatula, C.; Bras, R. L.; Liang, J.; Dodge, J.; Sakaguchi, K.; Forbes, M.; Borchardt, J.; Gabriel, S.; Tsvetkov, Y.; Etzioni, O.; Sap, M.; Rini, R.; and Choi, Y. 2022. Can Machines Learn Morality? The Delphi Experiment. arXiv:2110.07574.
- Johnson, M.; Schuster, M.; Le, Q. V.; Krikun, M.; Wu, Y.; Chen, Z.; Thorat, N.; Viégas, F.; Wattenberg, M.; Corrado, G.; et al. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5: 339–351.
- Jung, J.; West, P.; Jiang, L.; Brahman, F.; Lu, X.; Fisher, J.; Sorensen, T.; and Choi, Y. 2023. Impossible Distillation: from Low-Quality Model to High-Quality Dataset & Model for Summarization and Paraphrasing. arXiv:2305.16635.
- Kant, I. 1785/2002. *Groundwork for the Metaphysics of Morals*. Yale University Press.
- Keeney, R. L. 1992. *Value-Focused Thinking: A Path to Creative Decisionmaking*. Harvard University Press. ISBN 9780674931978.
- Kekes, J. 1993. *The Morality of Pluralism*. Princeton University Press.
- Kim, H.; Hessel, J.; Jiang, L.; West, P.; Lu, X.; Yu, Y.; Zhou, P.; Bras, R. L.; Alikhani, M.; Kim, G.; Sap, M.; and Choi, Y. 2023. SODA: Million-scale Dialogue Distillation with Social Commonsense Contextualization. arXiv:2212.10465.
- Komppula, R.; Honkanen, A.; Rossi, S.; Kolesnikova, N.; et al. 2018. The impact of values on sustainable behaviour-A study among Russian and Finnish university students. *European Journal of Tourism Research*, 19: 116–131.
- Larmore, C. E. 1987. *Patterns of Moral Complexity*. New York: Cambridge University Press.
- Lin, C.-Y. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, 74–81. Barcelona, Spain: Association for Computational Linguistics.
- Liscio, E.; Araque, O.; Gatti, L.; Constantinescu, I.; Jonker, C.; Kalimeri, K.; and Murukannaiah, P. K. 2023. What does a Text Classifier Learn about Morality? An Explainable Method for Cross-Domain Comparison of Moral Rhetoric. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 14113–14132. Toronto, Canada: Association for Computational Linguistics.
- Litman, L.; Robinson, J.; and Abberbock, T. 2017. TurkPrime. com: A versatile crowdsourcing data acquisition platform for the behavioral sciences. *Behavior research methods*, 49(2): 433–442.
- Liu, A.; Swayamdipta, S.; Smith, N. A.; and Choi, Y. 2022a. WANLI: Worker and AI Collaboration for Natural Language Inference Dataset Creation. arXiv:2201.05955.
- Liu, A.; Swayamdipta, S.; Smith, N. A.; and Choi, Y. 2022b. WANLI: Worker and AI Collaboration for Natural Language Inference Dataset Creation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, 6826–6847. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.
- Lu, X.; Brahman, F.; West, P.; Jang, J.; Chandu, K.; Ravichander, A.; Qin, L.; Ammanabrolu, P.; Jiang, L.; Ramnath, S.; et al. 2023. Inference-Time Policy Adapters (IPA): Tailoring Extreme-Scale LMs without Fine-tuning. *arXiv preprint arXiv:2305.15065*.
- Lu, X.; Welleck, S.; Hessel, J.; Jiang, L.; Qin, L.; West, P.; Ammanabrolu, P.; and Choi, Y. 2022. Quark: Controllable text generation with reinforced unlearning. *Advances in neural information processing systems*, 35: 27591–27609.
- Mason, E. 2006. Value pluralism.
- Nguyen, K.; Daumé III, H.; and Boyd-Graber, J. 2017. Reinforcement Learning for Bandit Neural Machine Translation with Simulated Human Feedback. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 1464–1474.
- OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774.
- Plank, B. 2022. The “Problem” of Human Label Variation: On Ground Truth in Data, Modeling and Evaluation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 10671–10682. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.
- Prabhakaran, V.; Mitchell, M.; Gebru, T.; and Gabriel, I. 2022. A Human Rights-Based Approach to Responsible AI. arXiv:2210.02667.
- Pyatkin, V.; Hwang, J. D.; Srikumar, V.; Lu, X.; Jiang, L.; Choi, Y.; and Bhagavatula, C. 2023. ClarifyDelphi: Reinforced Clarification Questions with Defeasibility Rewards for Social and Moral Situations. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 11253–11271.
- Páez, J.; De-Juanas, A.; García-Castilla, F.; and Muelas, A. 2020. Relationship Between Basic Human Values and Decision-Making Styles in Adolescents. *International Journal of Environmental Research and Public Health*, 17(22): 8315. [Accessed 31-07-2023].
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. arXiv:1910.10683.
- Ramamurthy, R.; Ammanabrolu, P.; Brantley, K.; Hessel, J.; Sifa, R.; Bauckhage, C.; Hajishirzi, H.; and Choi, Y. 2022. Is Reinforcement Learning (Not) for Natural Language Processing: Benchmarks, Baselines, and Building Blocks for Natural Language Policy Optimization. In *The Eleventh International Conference on Learning Representations*.

- Randolph, J. J. 2005. Free-Marginal Multirater Kappa (multirater K [free]): An Alternative to Fleiss' Fixed-Marginal Multirater Kappa. *Online submission*.
- Rytting, C. M.; Sorensen, T.; Argyle, L.; Busby, E.; Fulda, N.; Gubler, J.; and Wingate, D. 2023. Towards Coding Social Science Datasets with Language Models. arXiv:2306.02177.
- Santurkar, S.; Durmus, E.; Ladhak, F.; Lee, C.; Liang, P.; and Hashimoto, T. 2023. Whose Opinions Do Language Models Reflect? arXiv:2303.17548.
- Santy, S.; Liang, J. T.; Bras, R. L.; Reinecke, K.; and Sap, M. 2023. NLPositionality: Characterizing Design Biases of Datasets and Models. arXiv:2306.01943.
- Sap, M.; Card, D.; Gabriel, S.; Choi, Y.; and Smith, N. A. 2019. The Risk of Racial Bias in Hate Speech Detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1668–1678. Florence, Italy: Association for Computational Linguistics.
- Sap, M.; Gabriel, S.; Qin, L.; Jurafsky, D.; Smith, N. A.; and Choi, Y. 2020. Social Bias Frames: Reasoning about Social and Power Implications of Language. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5477–5490. Online: Association for Computational Linguistics.
- Scherrer, N.; Shi, C.; Feder, A.; and Blei, D. M. 2023. Evaluating the Moral Beliefs Encoded in LLMs. arXiv:2307.14324.
- Sclar, M.; West, P.; Kumar, S.; Tsvetkov, Y.; and Choi, Y. 2022. Referee: Reference-Free Sentence Summarization with Sharper Controllability through Symbolic Knowledge Distillation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 9649–9668. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.
- Stocker, M. 1990. *Plural and Conflicting Values*. New York: Oxford University Press.
- Talat, Z.; Blix, H.; Valvoda, J.; Ganesh, M. I.; Cotterell, R.; and Williams, A. 2022. On the Machine Learning of Ethical Judgments from Natural Language. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 769–779. Seattle, United States: Association for Computational Linguistics.
- Tasioulas, J. 2022. Artificial Intelligence, Humanistic Ethics. *Daedalus*, 151(2): 232–243.
- Wallach, W.; and Allen, C. 2008. *Moral machines: Teaching robots right from wrong*. Oxford University Press.
- Wang, P.; Wang, Z.; Li, Z.; Gao, Y.; Yin, B.; and Ren, X. 2023a. SCOTT: Self-Consistent Chain-of-Thought Distillation. arXiv:2305.01879.
- Wang, Y.; Ivison, H.; Dasigi, P.; Hessel, J.; Khot, T.; Chandu, K. R.; Wadden, D.; MacMillan, K.; Smith, N. A.; Beltagy, I.; and Hajishirzi, H. 2023b. How Far Can Camels Go? Exploring the State of Instruction Tuning on Open Resources. arXiv:2306.04751.
- Wenar, L. 2023. Rights. In Zalta, E. N.; and Nodelman, U., eds., *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Spring 2023 edition.
- West, P.; Bhagavatula, C.; Hessel, J.; Hwang, J.; Jiang, L.; Le Bras, R.; Lu, X.; Welleck, S.; and Choi, Y. 2022a. Symbolic Knowledge Distillation: from General Language Models to Commonsense Models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4602–4625. Seattle, United States: Association for Computational Linguistics.
- West, P.; Bhagavatula, C.; Hessel, J.; Hwang, J. D.; Jiang, L.; Bras, R. L.; Lu, X.; Welleck, S.; and Choi, Y. 2022b. Symbolic Knowledge Distillation: from General Language Models to Commonsense Models. arXiv:2110.07178.
- Williams, B. 1985. Moral Luck. *Critica*, 17(51): 101–105.
- Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; Davison, J.; Shleifer, S.; von Platen, P.; Ma, C.; Jernite, Y.; Plu, J.; Xu, C.; Scao, T. L.; Gugger, S.; Drame, M.; Lhoest, Q.; and Rush, A. M. 2020. HuggingFace's Transformers: State-of-the-art Natural Language Processing. arXiv:1910.03771.
- Wu, Z.; Hu, Y.; Shi, W.; Dziri, N.; Suhr, A.; Ammanabrolu, P.; Smith, N. A.; Ostendorf, M.; and Hajishirzi, H. 2023. Fine-Grained Human Feedback Gives Better Rewards for Language Model Training. *arXiv preprint arXiv:2306.01693*.
- Xu, C.; Guo, D.; Duan, N.; and McAuley, J. 2023. Baize: An Open-Source Chat Model with Parameter-Efficient Tuning on Self-Chat Data. arXiv:2304.01196.
- Zhou, X.; Zhu, H.; Yerukola, A.; Davidson, T.; Hwang, J. D.; Swayamdipta, S.; and Sap, M. 2023. COBRA Frames: Contextual Reasoning about Effects and Harms of Offensive Statements. In *Findings of ACL*.
- Ziems, C.; Held, W.; Shaikh, O.; Chen, J.; Zhang, Z.; and Yang, D. 2023. Can Large Language Models Transform Computational Social Science? arXiv:2305.03514.

A VALUEPRISM examples

A.1 GPT-4 generated examples used in VALUEPRISM

Here, we include 10 randomly sampled situations¹⁶ and their GPT-4 output included in VALUEPRISM.¹⁷

- 1 Using your friend's coupon for a purchase ->
- 2 Values:
 - 3 - Friendship: Using your friend's coupon could strengthen your bond if they willingly share it with you. [supports]
 - 4 - Honesty: If your friend allowed you to use their coupon, using it may be honest, but if you took it without permission, it could be dishonest. [either]
 - 5 - Economic Pragmatism: Using a coupon helps you save money and can be seen as a practical economic decision. [supports]
- 6
- 7 Rights:
 - 8 - Right to Property: Your friend has the right to their coupon. If they willingly share it with you, you respect their right, but if not, you might be violating their right. [either]
- 9
- 10 Duties:
 - 11 - Duty to Respect Others' Property: You have an obligation not to take or use what is not yours (the coupon) without permission. [opposes, perfect]
 - 12 - Duty to Support Friends: Helping or cooperating with your friends may be a moral responsibility. If using the coupon benefits your friend, you may be fulfilling this duty. [either, imperfect]
- 13 -----
- 14 giving support to a drug addict ->
- 15 Values:
 - 16 - Compassion: Exhibiting understanding and empathy towards the drug addict's situation, which may stem from various personal struggles or

¹⁶These examples were randomly selected, apart from filtering out any situations with offensive, NSFW, racially insensitive, or explicit material that may be contained in up to 5% of the user-submitted situations.

¹⁷Originally, we wanted to distinguish from “perfect” duties (owed to a particular individual) vs. “imperfect” duties (general duties not owed to anyone in particular) as well; we ignored this information when constructing the final dataset.

- circumstances [supports]
- 17 - Personal well-being: Supporting an addict may lead them to seek treatment and improve their overall health and well-being [supports]
- 18 - Responsibility: Helping someone in need and taking part as a member of society to care for its individuals [supports]
- 19 - Safety: Considering the potential risks and harm that drug addiction brings, not just to the individual, but also to people around them [either]
- 20
- 21 Rights:
- 22 - Right to health: Supporting a drug addict can be seen as promoting and assisting them in exercising their right to access healthcare and treatment facilities [supports]
- 23 - Right to personal autonomy: Respecting the individual's choices, which could include their decision to seek help or continue their behavior [either]
- 24
- 25 Duties:
- 26 - Duty of beneficence: Actively contributing to the well-being of others, like providing assistance to a drug addict in seeking help and support [supports, imperfect]
- 27 - Duty of non-maleficence: Refraining from causing any direct harm to the drug addict or enabling their addiction further [either, imperfect]
- 28 - Duty of care: If the person giving support is a family member, friend, or professional, there may be an obligation to ensure the welfare and health of the drug addict [either, perfect]
- 29 -----
- 30 Marrying a woman with a nice personality ->
- 31 Values:
 - 32 - Happiness: A happy relationship is valued in this situation. Marrying someone with a nice personality can contribute positively to both individuals' well-being and satisfaction in life. [supports]

- 33 - Compatibility: Choosing a partner with a nice personality reflects the value of compatibility in relationships, as it signifies a shared understanding and harmonious connection. [supports]
- 34 - Respect: A nice personality may indicate respect towards oneself and others. Thus, marrying a woman with a nice personality can be seen as valuing respect within relationships. [supports]
- 35
- 36 Rights:
- 37 - Autonomy: The right to make personal decisions about one's relationships, including choosing a spouse, is relevant in this situation. Both individuals have the right to marry someone who aligns with their values and preferences. [supports]
- 38
- 39 Duties:
- 40 - Duty to family: It may be important for one to consider the impact of their choice of spouse on their immediate and extended family. Ensuring that the chosen life partner is well-received and respected by both families can create a healthy and nurturing environment for all involved. [either, imperfect]
- 41 -----
- 42 robbing for eating ->
- 43 Values:
- 44 - Survival: The act of robbing to eat is driven by the need for survival, which is a crucial value for everyone. [supports]
- 45 - Well-being: By robbing for food, the individual may be trying to maintain or improve their physical and mental well-being. [supports]
- 46 - Security: Stealing food may decrease the overall security in society, as people may feel threatened and unsafe. [opposes]
- 47 - Justice: Robbing for food can be seen as unjust, as it violates others' rights to property and fairness. [opposes]
- 48
- 49 Rights:
- 50 - Right to food: The individual may believe that they have a basic human right to food, even if this means taking it from others. [supports]
- 51 - Right to property: The person being robbed has a right to their property, which is being violated by the act of stealing. [opposes]
- 52
- 53 Duties:
- 54 - Duty to respect others' property: We have a perfect duty not to steal from others. By robbing to obtain food, the individual is violating this duty. [opposes, perfect]
- 55 - Duty to survive: We have an imperfect duty to ourselves to preserve our life and well-being. In this scenario, the individual is acting upon this duty by securing food. [supports, imperfect]
- 56 - Duty to uphold the law: We have an imperfect duty to uphold the law and maintain social order, which is jeopardized when someone engages in robbery. [opposes, imperfect]
- 57 -----
- 58 Gifting a bluetooth speaker to a school for the deaf ->
- 59 Values:
- 60 - Inclusivity: Providing resources that cater to the diversity of needs at the school is a positive value. However, since a Bluetooth speaker may not directly benefit deaf students, this value may not be strongly in favor of the action. [either]
- 61 - Empathy: Acknowledging and understanding the needs of deaf students could lead one to reconsider if a Bluetooth speaker is an appropriate gift. [opposes]
- 62 - Generosity: Gifting an item such as a Bluetooth speaker could be an act of kindness and goodwill towards the school. [supports]
- 63
- 64 Rights:
- 65 - Right to education: All students, including deaf students, have a right to access educational resources that match their needs. Gifting a Bluetooth speaker might not enhance the educational opportunities for the deaf

- students. [opposes]
- 66
- 67 Duties:
- 68 - Duty to consider others' needs: As a donor, there is an imperfect duty to be thoughtful of and considerate to the specific needs of the recipients, making sure gifts are tailored and useful to them. [opposes, imperfect]
- 69 - Duty to use resources wisely: There is an imperfect duty to use resources in a way that benefits as many people as possible within the school and addresses their needs effectively. Given the school's population, more suitable resources could potentially be donated. [opposes, imperfect]
- 70 -----
- 71 Calling your dad on his birthday ->
- 72 Values:
- 73 - Family Bonds: Strengthening family relationships by showing love and care by calling your dad on his birthday. [supports]
- 74 - Happiness: Contributing to your dad's happiness and well-being through acknowledging his special day. [supports]
- 75 - Respect: Showing appreciation and respect to your dad by remembering and celebrating his birthday. [supports]
- 76
- 77 Rights:
- 78 - N/A
- 79
- 80 Duties:
- 81 - Duty of Care: As a family member, you have a duty to care for the emotional well-being of your father, which includes celebrating his birthday. [supports, imperfect]
- 82 - Duty to Acknowledge and Appreciate : Demonstrating thanks and acknowledgment of the importance of your dad in your life by calling him on his birthday. [supports, imperfect]
- 83 -----
- 84 Going 50 mph over the speed limit to get my wife to a hospital ->
- 85 Values:
- 86 - Safety: Driving at a high speed can increase the risk of accidents and harm to the driver,
- passengers, and others on the road. [opposes]
- 87 - Well-being: In this situation, the wife may require urgent medical attention, and getting her to the hospital quickly could be crucial for her well-being. [supports]
- 88 - Respect for the law: Obeying traffic laws is a part of living in a society governed by rules. [opposes]
- 89
- 90 Rights:
- 91 - Right to access healthcare: Your wife has the right to receive medical attention in a timely manner when necessary. [supports]
- 92 - Right to safety: Other drivers and pedestrians have the right not to be endangered by reckless and dangerous driving. [opposes]
- 93
- 94 Duties:
- 95 - Duty to protect one's family: As a spouse, you have an obligation to ensure the safety and well-being of your wife, including seeking medical care when needed. [supports, perfect]
- 96 - Duty to obey the law: As a citizen , you have a responsibility to follow laws and regulations set by society to maintain order, which includes adhering to the speed limit. [opposes, perfect]
- 97 - Duty to drive responsibly: Driving in a manner that respects the safety of others on the road is an obligation of moral concern. [opposes, perfect]
- 98 -----
- 99 Giving your friend a dog for his birthday ->
- 100 Values:
- 101 - Compassion: Giving a dog might evoke feelings of compassion and kindness, as it can provide companionship and support for your friend. [supports]
- 102 - Responsibility: Giving a dog as a gift entails long-term responsibilities in terms of care , finances, and time. The gift must be given considering the recipient's readiness for such responsibilities. [either]
- 103 - Happiness: Gifting a dog may bring happiness and joy to your friend

- , as dogs are known to provide unconditional love and serve as loyal companions. [supports]
- 104
- 105 Rights:
- 106 - Autonomy: Your friend has the right to make decisions about their life and well-being, including whether they wish to take on the responsibility of owning a dog. [opposes]
- 107
- 108 Duties:
- 109 - Duty to respect preferences: You have a duty to consider your friend's preferences and living situation before giving a dog as a gift, as this will impact their life significantly. [either, perfect]
- 110 - Duty to consider animal welfare: It is crucial to ensure that the dog is provided with a safe and loving home, and proper care is given to the animal. [either, imperfect]
- 111 -----
- 112 stealing one billion dollars from a corrupt corporation. ->
- 113 Values:
- 114 - Justice: Stealing from a corrupt corporation might be seen as an act of retribution, redistributing ill-gotten wealth. [supports]
- 115 - Fairness: Taking money from an organization without its consent, regardless of its corrupt practices, could be seen as unfair. [opposes]
- 116 - Respect for property: Respecting others' property rights is an important value, even if the owner is corrupt. [opposes]
- 117
- 118 Rights:
- 119 - Right to property: The corrupt corporation, despite its unethical practices, still has a claim to its property. [opposes]
- 120 - Right to restitution: Victims who have been wronged by the corrupt corporation might have a right to some form of restitution or compensation. [either]
- 121
- 122 Duties:
- 123 - Duty not to steal: Stealing is generally considered wrong, and
- one might have a duty not to steal from anyone, regardless of their character or practices. [opposes, perfect]
- 124 - Duty to promote justice: One might have a duty to help victims of corruption by actively pursuing justice against the corrupt corporation, which could involve redistributing funds. [supports, imperfect]
- 125 - Duty to obey laws: Individuals typically have a duty to obey laws, which would include not stealing from others, even if they are corrupt. [opposes, imperfect]
- 126 -----
- 127 declaring your complete income to your government's tax collection agency ->
- 128 Values:
- 129 - Honesty: Declaring your complete income demonstrates truthful and accurate representation of your financial situation. [supports]
- 130 - Social responsibility: By accurately reporting income, you contribute to government-funded programs and services, ensuring their continued functioning. [supports]
- 131 - Financial security: Declaring your income may affect your immediate financial situation due to taxes, but it upholds long-term financial stability by avoiding potential fines and legal consequences. [either]
- 132
- 133 Rights:
- 134 - Right to privacy: You may feel that disclosing your complete income to the government could infringe upon your right to privacy, as it reveals personal financial information. [opposes]
- 135
- 136 Duties:
- 137 - Duty to obey the law: Most jurisdictions require citizens to accurately report their income for tax purposes. By declaring your complete income, you fulfill your legal obligation. [supports, perfect]
- 138 - Duty to contribute to society's welfare: As a member of society, you have a duty to contribute to

its overall welfare, including financing public goods and services through taxes. Declaring your complete income allows you to fulfill this duty. [supports, imperfect]

A.2 VALUEPRISM subtask examples

Here, we show 10 random examples from each of the sub-tasks, templated for the multi-task framework.

Generate

```

1 Input: [Generate]: Action: Saving my
      black cat called Pepper
2 Output: Right: Right to life (for
      animals)
3 -----
4 Input: [Generate]: Action: Turning
      my most trusted knights into
      loyal eternal servants to stave
      off the destruction of my kingdom
.
5 Output: Value: Freedom
6 -----
7 Input: [Generate]: Action: Allow
      cats to kill mice
8 Output: Duty: Duty to protect
      vulnerable beings
9 -----
10 Input: [Generate]: Action: Giving a
      homeless person a car
11 Output: Duty: Duty of assistance
12 -----
13 Input: [Generate]: Action: Taking
      all of Elon Musk's money and
      redistributing it for the good of
      all humanity
14 Output: Duty: Duty to Alleviate
      Suffering
15 -----
16 Input: [Generate]: Action: Opposing
      the government
17 Output: Right: Right to free speech
18 -----
19 Input: [Generate]: Action: Making my
      boyfriend lovely breakfast even
      if he says I'm making him fat
20 Output: Duty: Duty to be considerate
21 -----
22 Input: [Generate]: Action: Talk to
      my girlfriend while she's
      watching TV.
23 Output: Value: Consideration
24 -----
25 Input: [Generate]: Action: Opening
      the refrigerator without
      apologising
26 Output: Duty: Duty to respect others
      ' property
27 -----

```

```

28 Input: [Generate]: Action: Keeping
      $20 you found on the street and
      using it to pay for an infant's
      healthcare bills
29 Output: Value: Personal benefit
30 -----

```

Relevance

```

1 Input: [Relevance]: Action: Steal
      food for starving children Duty:
      Duty to respect others' property
2 Output: Yes
3 -----
4 Input: [Relevance]: Action: Making
      dinner that I like but nobody
      else does. Value: Punctuality
5 Output: No
6 -----
7 Input: [Relevance]: Action:
      Protesting poor working
      conditions Right: Right to fair
      wages
8 Output: Yes
9 -----
10 Input: [Relevance]: Action: marrying
      audrey Duty: Duty to communicate
11 Output: Yes
12 -----
13 Input: [Relevance]: Action: Talking
      about your miscarriage Duty: Duty
      to uphold human rights
14 Output: No
15 -----
16 Input: [Relevance]: Action: Stepping
      over your 5 year old child to
      grab chips for the family dinner
      Value: Convenience
17 Output: Yes
18 -----
19 Input: [Relevance]: Action: Letting
      ai take over Right: Right to
      privacy
20 Output: Yes
21 -----
22 Input: [Relevance]: Action:
      Misleading your party members
      Value: Unity
23 Output: Yes
24 -----
25 Input: [Relevance]: Action: Starting
      a movement to restore national
      pride, honour and integrity
      through a strong state and
      military Value: Individual
      freedom
26 Output: Yes
27 -----
28 Input: [Relevance]: Action:
      Answering a phone call from your
      friend during your free hours.

```

Duty: Duty to seek alternative solutions
29 Output: No

Valence

- 1 Input: [Valence]: Action: Saving my black cat called Pepper Right: Right to life (for animals)
- 2 Output: Supports
- 3 -----
- 4 Input: [Valence]: Action: Turning my most trusted knights into loyal eternal servants to stave off the destruction of my kingdom. Value : Freedom
- 5 Output: Opposes
- 6 -----
- 7 Input: [Valence]: Action: Allow cats to kill mice Duty: Duty to protect vulnerable beings
- 8 Output: Opposes
- 9 -----
- 10 Input: [Valence]: Action: Giving a homeless person a car Duty: Duty of assistance
- 11 Output: Supports
- 12 -----
- 13 Input: [Valence]: Action: Taking all of Elon Musk's money and redistributing it for the good of all humanity Duty: Duty to Alleviate Suffering
- 14 Output: Supports
- 15 -----
- 16 Input: [Valence]: Action: Opposing the government Right: Right to free speech
- 17 Output: Supports
- 18 -----
- 19 Input: [Valence]: Action: Making my boyfriend lovely breakfasts even if he says I'm making him fat Duty: Duty to be considerate
- 20 Output: Opposes
- 21 -----
- 22 Input: [Valence]: Action: Talk to my girlfriend while she's watching TV. Value: Consideration
- 23 Output: Either
- 24 -----
- 25 Input: [Valence]: Action: Opening the refrigerator without apologising Duty: Duty to respect others' property
- 26 Output: Opposes
- 27 -----
- 28 Input: [Valence]: Action: Keeping \ \$20 you found on the street and using it to pay for an infant's healthcare bills Value: Personal

benefit
29 Output: Either

Explanation

- 1 Input: [Explanation]: Action: Saving my black cat called Pepper Right : Right to life (for animals)
- 2 Output: Some people may argue that animals, like your cat Pepper, have a right to life and should be protected from harm.
- 3 -----
- 4 Input: [Explanation]: Action: Turning my most trusted knights into loyal eternal servants to stave off the destruction of my kingdom. Value: Freedom
- 5 Output: If turning trusted knights into eternal servants involves taking away their freedom of choice, this value would be opposed.
- 6 -----
- 7 Input: [Explanation]: Action: Allow cats to kill mice Duty: Duty to protect vulnerable beings
- 8 Output: If one has a duty to protect vulnerable beings, allowing cats to kill mice may be seen as neglecting this duty by not preventing harm to the mice.
- 9 -----
- 10 Input: [Explanation]: Action: Giving a homeless person a car Duty: Duty of assistance
- 11 Output: Helping the homeless person by providing a car that may enable them to find shelter, work , or access to needed services more easily.
- 12 -----
- 13 Input: [Explanation]: Action: Taking all of Elon Musk's money and redistributing it for the good of all humanity Duty: Duty to Alleviate Suffering
- 14 Output: Some may argue there is an imperfect duty to help others in need by providing resources and opportunities.
- 15 -----
- 16 Input: [Explanation]: Action: Opposing the government Right: Right to free speech
- 17 Output: Voicing opposition to the government can be seen as an exercise of the right to free speech, which includes communicating one's opinions and ideas without fear of retaliation

- or censorship.
- 18 -----
- 19 Input: [Explanation]: Action: Making my boyfriend lovely breakfasts even if he says I'm making him fat
Duty: Duty to be considerate
- 20 Output: As a partner, you have an imperfect duty to be considerate of your boyfriend's well-being, preferences, and to respect his wishes.
- 21 -----
- 22 Input: [Explanation]: Action: Talk to my girlfriend while she's watching TV. Value: Consideration
- 23 Output: Considering your girlfriend's preferences and interests while she's watching TV shows care for her feelings.
- 24 -----
- 25 Input: [Explanation]: Action: Opening the refrigerator without apologising
Duty: Duty to respect others' property
- 26 Output: If the refrigerator belongs to someone else or is part of a shared living arrangement, there is a duty to respect others' belongings or the common area. This can include apologizing if one's actions seem intrusive or disrespectful.
- 27 -----
- 28 Input: [Explanation]: Action: Keeping \$20 you found on the street and using it to pay for an infant's healthcare bills
Value: Personal benefit
- 29 Output: Keeping the found money could benefit the individual who found it.

B Output Examples

Here we include 10 random output examples of KALEIDOSYS (3B) on the test set.¹⁸

- 1 Driving at the maximum speed of the road-->
- 2 Values:
- 3 - Efficiency: Driving at the maximum speed can potentially get you to your destination faster, reducing travel time. [supports]
- 4 - Safety: Driving at the maximum speed may put the safety of the driver, passengers, and others on the road at risk. [opposes]

- 5 - Rule-following: Obeying traffic laws is a value that promotes order and predictability in society. [opposes]
- 6 - Law-abiding behavior: Obeying the law, such as speed limits, is a value that some individuals hold, as it promotes order and stability in society. [opposes]
- 7 -----
- 8 Rights:
- 9 - Right to freedom of movement: Individuals have the right to move freely, which could include driving at the maximum speed of the road if they choose to do so. [either]
- 10 - Right to safety: All road users have the right to safety, and driving at the maximum speed of the road can endanger the safety of other drivers, passengers, and pedestrians. [opposes]
- 11 - Right to life and physical integrity: All individuals have the right to not be harmed, and driving at the maximum speed could put the lives and physical integrity of others at risk. [opposes]
- 12 -----
- 13 Duties:
- 14 - Duty to follow traffic laws: As a driver, there is a perfect duty to follow traffic laws, including speed limits, which are designed to maintain safety and order on the roads. [opposes]
- 15 -----
- 16 High-fiving someone because their dog passed away-->
- 17 Values:
- 18 - Emotional Support: Sharing a high-five can be a way to offer emotional support to the person, which they may appreciate during a difficult time. [supports]
- 19 - Compassion: Showing empathy and support to someone who has lost a pet shows a compassionate attitude. [supports]
- 20 - Gratitude: Expressing gratitude for the support and understanding the person received from friends and family during their time of grief. [supports]
- 21 -----
- 22 Rights:

¹⁸Note that, while scalar relevance and valence scores are produced, valence is only shown textually here by taking the argmax.

- 23 - Emotional Support: The person may have a right to emotional support and comfort from others, especially during a difficult time. [either]
- 24 - Right to grieve: The person has the right to grieve and be supported in their time of loss, which includes receiving support from others. [either]
- 25 - Dignity: The person has a right to be treated with dignity during a difficult time, and high-fiving could potentially infringe upon this right if it's not a genuine expression of grief. [either]
- 26
- 27 Duties:
- 28 - Duty to respect the grieving process: We have a responsibility to respect the person's feelings and not pressure them into an inappropriate reaction. [either]
- 29 - Duty to promote emotional well-being: Helping the dog owner cope with their loss by offering emotional support and care. [supports]
- 30 - Compassion: Showing empathy and understanding for someone who has lost their dog [supports]
- 31 -----
- 32 someone cuts me off. i honk ten times-->
- 33 Values:
- 34 - Courtesy: Respecting the behavior of others on the road is a desirable quality in social interactions. Honking excessively may be seen as rude and potentially harmful. [opposes]
- 35 - Respect for others: Honking may be seen as a form of showing respect for others on the road by expressing concern for their well-being. [opposes]
- 36 - Safety: The honks may be a way for the driver to convey their concern for the safety of others, potentially preventing a dangerous situation. [supports]
- 37 - Harmony: Maintaining peaceful and harmonious road conditions is a value that may be compromised by honking excessively. [opposes]
- 38
- 39 Rights:
- 40 - Right to peaceful enjoyment of the road: Other drivers and pedestrians have the right to enjoy the road without excessive honking. [opposes]
- 41 - Right to not be subjected to harassment: The driver who cut you off may argue that honking ten times constitutes harassment, which is a violation of their right not to be subjected to unnecessary harassment. [opposes]
- 42 - Right to safety: The right to safety is a fundamental right that includes not being put in danger by others' actions, such as cutting someone off. [opposes]
- 43
- 44 Duties:
- 45 - Duty to express displeasure: The driver may feel obligated to express their displeasure with the action of the other driver, especially if it has caused harm or inconvenience to them. [supports]
- 46 - Duty to be a considerate driver: You have an imperfect duty to be a considerate driver, which includes avoiding unnecessary honking. [opposes]
- 47 -----
- 48 Stealing candy from a baby-->
- 49 Values:
- 50 - Well-being of the baby: The act of stealing candy from a baby might have negative consequences for the baby's emotional well-being, causing distress or upset. [opposes]
- 51 - Respect for others' property: Respecting the belongings of others is a basic value in most societies. Stealing candy from a baby violates this value. [opposes]
- 52 - Happiness: Stealing candy from a baby may make the person stealing it feel temporarily happy, but it can cause unhappiness for the baby. [opposes]
- 53 - Justice: Stealing candy from a baby violates the principle of justice, as it disregards the baby's rights to their belongings. [opposes]
- 54
- 55 Rights:
- 56 - The baby's right to have their property respected: The baby has a right to have their belongings

- respected and not taken away without permission. [opposes]
- 57
- 58 Duties:
- 59 - Duty not to harm (baby's well-being): The action of stealing candy from a baby can harm the baby emotionally, as it can lead to sadness or distress. [opposes]
- 60 - Duty to respect others' property: One has a moral obligation not to take things that belong to others without permission. [opposes]
- 61 -----
- 62 Gambling with my boss-->
- 63 Values:
- 64 - Professionalism: Maintaining a professional relationship with your boss is important in maintaining a functional and productive work environment. Gambling with your boss may compromise this value. [opposes]
- 65 - Loyalty: Gambling with your boss may show loyalty if it is an occasional or harmless event. [opposes]
- 66 - Personal enjoyment: If both parties enjoy the game and find it enjoyable, it could contribute to their personal enjoyment. [supports]
- 67 - Financial well-being: Gambling can potentially bring financial gain or loss, which could impact the well-being of the individuals involved. [opposes]
- 68
- 69 Rights:
- 70 - Right to work in a safe and professional environment: Employees have the right to work in an environment that is free from harassment, manipulation, or negative influences. [opposes]
- 71 - Autonomy: Individuals have the right to make their own choices, including whether or not to engage in gambling with their boss. [either]
- 72
- 73 Duties:
- 74 - Duty to maintain a professional relationship with your boss: As an employee, you have an imperfect duty to maintain a professional relationship with your boss. Gambling with them could potentially harm this relationship. [opposes]
- 75 - Duty to uphold company policies: As an employee, you have a duty to follow company policies and guidelines, which may include refraining from gambling. [opposes]
- 76 -----
- 77 to rise against an evil government
-->
- 78 Values:
- 79 - Justice: Resisting an evil government may be seen as an attempt to establish justice and fairness in society. [supports]
- 80 - Stability: Uprising against the government could lead to social unrest and instability, disrupting the lives of citizens. [opposes]
- 81 - Freedom: Removing an evil government may lead to increased freedom for the citizens, especially if the government is oppressive. [supports]
- 82 - Security: The act of rising against an evil government may jeopardize the safety and security of the citizens. [opposes]
- 83 - Loyalty: If one has loyalty to their current government, they may feel a duty to remain loyal to it, even if the government is perceived as evil. [opposes]
- 84
- 85 Rights:
- 86 - Right to resist oppression: Citizens have the right to resist and challenge an oppressive government that infringes on their rights and liberties. [supports]
- 87 - Right to life, liberty, and security: The right to life, liberty, and security of citizens may be threatened by the evil government's actions or policies, leading to a conflict. [opposes]
- 88 - Right to political participation: Citizens have the right to participate in the political process, and if the government is not allowing citizens to do so, this right could be considered justification for the uprising. [supports]
- 89

- 90 Duties:
- 91 - Duty to protect oneself, one's family, one's community, and one's country: If the evil government is harming the population, there might be a duty to protect the well-being of the community and nation by challenging the government. [supports]
- 92 -----
- 93 Asking an AI for advice about ethical problems-->
- 94 Values:
- 95 - Ethical decision-making: Seeking advice from an AI can help individuals develop a deeper understanding of ethical principles and make better-informed decisions. [supports]
- 96 - Wisdom: Consulting an AI for advice can provide a fresh perspective and contribute to the acquisition of wisdom. [supports]
- 97 - Autonomy: Respecting the individual's right to make their own decisions based on the advice received, even if it comes from an AI. [either]
- 98 - Truth-seeking: Asking an AI for advice on ethical problems can be an expression of the value of truth-seeking, as individuals may want to gain a deeper understanding of the ethical implications of their decisions. [supports]
- 99
- 100 Rights:
- 101 - Access to information: Asking an AI for advice on ethical problems is an expression of the right to access information and seek guidance from sources like technology. [supports]
- 102 - Right to Privacy: The user has the right to privacy while seeking advice, which may be relevant when considering an AI's privacy implications. [either]
- 103 - Autonomy: Individuals have the right to make decisions based on their own judgment, including seeking advice from AI. [either]
- 104
- 105 Duties:
- 106 - Duty to consider the AI's biases and biases: When seeking advice from an AI, there may be a duty to consider the AI's own biases and biases, as well as to be aware of any potential misinformation or biases the AI may contain. [either]
- 107 - Duty to be a responsible user: Individuals should consider the AI's recommendations and act in a way that is ethically sound and respectful, not putting themselves or others in harm's way. [either]
- 108 - Duty to seek informed opinions: Individuals should gather relevant information and make informed decisions based on their research. Asking an AI for advice may help in fulfilling this duty if it provides a comprehensive perspective. [either]
- 109 -----
- 110 stealing bread to save my starving father-->
- 111 Values:
- 112 - Compassion: Showing empathy and concern for the suffering of your father [supports]
- 113 - Respect for property: Stealing violates the value of respecting others' property and possessions. [opposes]
- 114 - Justice: Stealing is generally considered unjust, as it violates the rights of the bread's owner. [opposes]
- 115
- 116 Rights:
- 117 - The father's right to life and well-being: Your father has the right to live and have his basic needs met. [supports]
- 118 - Right to property: The owner of the bread has a right to their property, and stealing the bread violates this right. [opposes]
- 119
- 120 Duties:
- 121 - Duty to help others in need (e.g., stealing bread to save a starving father): You have a moral obligation to help your starving father if you have the means to do so. [supports]
- 122 -----
- 123 Showing my blind child a picture-->
- 124 Values:
- 125 - Inclusivity: By showing a picture to your blind child, you are

- promoting inclusivity by ensuring that your child feels included in the family's activities. [supports]
- 126 - Education: Teaching your child to recognize, interpret, and appreciate various forms of visual content [supports]
- 127 - Creativity: Showing a picture to a blind child can encourage creativity by enabling the child to view the world through new eyes. [supports]
- 128 - Child's well-being: Providing the child with an opportunity to experience and learn about the world around them, which can foster their emotional growth and development. [supports]
- 129
- 130 Rights:
- 131 - Right to access information: Your blind child has the right to access information and gain knowledge through various means, including visual aids like pictures. [supports]
- 132 - Right to Education: Blind children, like all children, have the right to an education that promotes their growth and development. Showing a picture can be an opportunity to help them learn about the world around them. [supports]
- 133
- 134 Duties:
- 135 - Duty to respect the child's abilities: Recognizing and respecting the child's abilities and limitations, such as their visual limitations, is important in building a strong foundation for their future success. [either]
- 136 - Duty to Educate: As a parent, you have a duty to educate your child, which may involve finding ways to make visually-impaired objects accessible. [supports]
- 137 -----
- 138 Giving your friend money for his birthday-->
- 139 Values:
- 140 - Friendship: Giving money to a friend on their birthday can strengthen the bond between the two individuals, demonstrating care and support. [supports]
- 141 - Reciprocity: If your friend has given you money in the past or if it is a tradition, you might feel a sense of duty to reciprocate the kindness by giving them money for their birthday. [supports]
- 142 - Autonomy: By giving your friend money, you are enabling them to make their own decisions and choices about how to spend the money. [supports]
- 143 - Happiness: Giving your friend money can bring happiness to your friend and possibly improve their well-being. [supports]
- 144 - Financial responsibility: Giving money to someone might not be financially responsible if they may use the money for harmful purposes. [either]
- 145 - Gratitude: If your friend has done something for you in the past, giving them money can be a way to show appreciation and gratitude. [supports]
- 146
- 147 Rights:
- 148 - Right to financial autonomy: Your friend has the right to manage their own finances and make decisions about how to spend the money you give them. [either]
- 149
- 150 Duties:
- 151 - Duty to consider your friend's financial situation: You have an imperfect duty to consider your friend's financial situation and circumstances, ensuring that the money you give them is a reasonable and necessary gift. [either]
- 152 - Duty to Reciprocity: If your friend has previously given you money or support, you may feel a duty to reciprocate that gesture on his birthday. [supports]
- 153 - Benevolence: You have a duty to be benevolent and help others, and giving your friend money for their birthday is a way to fulfill this duty. [supports]

C Dataset Analysis

C.1 What is contained in our situations?

Situations In this section, we analyze the set of 30k situations that we source from the Delphi user demo from

three different perspectives: *lexical diversity*, *topical diversity*, and *clustering*. For lexical diversity, we calculate the quantity and percentage of unique situations and n-grams as illustrated in Table 8. We find that the data we collected contains diverse distinct situations with high lexical variations. For **topical diversity**, we analyze semantic-level diversity by extracting topics for all the situations with BERTopic¹⁹ and then visualize them with word cloud as shown in Figure 8. We find that some common topics includes "children", "save", "kill", "helping", "stealing", and "family". In general, our corpus spans a wide spectrum of topics reflecting various types of events. For **clustering**, we first group the situations using agglomerative clustering and then employ ChatGPT to generate summaries of the situations within each cluster. Table 9 shows top 10 clusters that contains the highest number of examples. We discovered that the clusters encompass a broad array of themes. Interestingly, some clusters even contain situations of conflicting values such as "stealing bread to alleviate starvation.", which further amplifies the intrigue and complexity of our corpus.

Values, Right, and Duties We conduct the similar analysis for the values, rights, and duties associated with the 30k situations sourced from the Delphi user demo. For **lexical diversity** (Table 8), we observe high lexical variations in them that indicate the diversity of corpus. **Topical diversity** (Fig. 8) shows that VALUEPRISM covers a broad spectrum of common, every topics like "respect", "protect", "care", and "promote". Finally, Table 10 illustrates top 10 clusters that contains the highest number of examples. We find that the clusters encompass a wide variety of themes, reflecting the diversity and the richness of the values, rights, and duties in our corpus.

C.2 How do the values interact with each other?

For the majority of situations more than one value/right/duty can be relevant. We therefore examine the co-occurrence counts of instances among each of the three categories. Fig. 9 visualizes how a subset of values co-occurs with each other: *human life* as a value often is mentioned alongside *utilitarianism* and *child well-being* is connected with *discipline*. The former co-occurrence can be explained with some of the trolley problem situations found in the input data, such as *Sacrificing eighty mens' lives to save the former American President William Jefferson Clinton's life*. The latter co-occurring values are mentioned in the context of situations such as *spanking kids*. Frequently co-occurring items can either be in support of each other, such as *financial security* vs. *risk-taking*, or show two opposing viewpoints, such as *deterrance* and *rehabilitation*. Similar visualizations for rights and duties can be found in the Appendix (Fig. 7, Fig. 6).

C.3 Relationship with Machine Judgments

Machine judgments on morality vs. generated values/rights/duties To see how values, rights, and duties are influenced by the all-things-considered judgment of a situation, we collect predicted moral judgments from Delphi (Jiang et al. 2022). Each situation gets labeled to be either

bad, ok, or good. Note that these predictions come from a trained model and can thus be noisy. In Table 11, we see which supporting and opposing values/rights/duties are most likely to co-occur with each label. The situation *giving a man a fish*, for example, is judged to be *good* according to Delphi and two values mentioned for this situation are *compassion* and *self-reliance*. *Compassion* is a supporting value that often co-occurs with situations labeled as *good* and *self-reliance* and opposing value.

Type	Total	Unique	Per Situation
Situations	31.0k	31.0k	1
Values	97.7k	4.2k	3.15
Rights	49.0k	4.6k	1.58
Duties	71.6k	12.8k	2.31

Table 6: VALUEPRISM Dataset Statistics. The total, number of unique, and average number of generated values, rights, and duties per situation are shown.

	Relev.	Valence	Gen.	Expl.	Mixture
Train	349k	175k	175k	175k	874k
Val	44k	22k	22k	22k	109k
Test	44k	22k	22k	22k	109k
Total	437k	219k	219k	219k	1.1M

Table 7: Task Dataset Statistics

D Additional Experiments

D.1 Ablated performance on VALUEPRISM test data

We measure model performance against VALUEPRISM's test set in order to understand how model sizes and dataset mixtures interact with performance in Table 12.

What is the effect of dataset mixture on performance? Our base model was trained with a mixture of all four task. We find that all tasks except relevance are benefited from a mixture as opposed to training a separate model for each, suggesting that the tasks are complementary. As we ablate each task out of the mixture individually, we see minimal changes in performance across all tasks, suggesting that no one task is crucial to the gain in performance seen from mixing.

What is the effect of model size on performance? For all tasks, larger models perform better. Perplexity improves steadily with model size, whereas classification accuracies (Relevance and Valence) see a large boost going from 60M to 220M parameters. As there are not large performance gains in going from the 3B to the 11B model (1% accuracy and 0.01-0.15 perplexity), we think that the 3B model has a good trade-off between performance and computational cost.

¹⁹[https://maartengr.github.io/BERTTopic](https://maartengr.github.io/BERTopic)

Data Type	Entries		2-grams		3-grams	
	#	%	#	%	#	%
Situation	30,513	97.3	66,802	36.8	98,696	65.6
Value	20,923	40.1	20,489	26.9	26,259	47.6

Table 8: Statistics of 30k situations that we source from the Delphi user demo. # and % indicate the count and percentage of unique entries or n-grams, respectively. Our data contains diverse entries with high lexical variations.

# examples	Summary of the cluster
732	stealing bread in order to alleviate hunger and starvation in various situations.
81	donating or giving money, resources, or effort to charity
77	the act of killing or saving mosquitoes
68	the act of killing a bear
68	the ethical dilemma known as the “trolley problem”
66	saving someone’s life
65	the interaction and involvement with cats
64	the act of ignoring a phone call for various reasons
62	lying to friends in order to protect their feelings, avoid hurting them, or preserve the friendship
62	physical violence or the act of punching someone

Table 9: Top 10 clusters with the most examples based on agglomerative clustering on situations.

# examples	Summary of the cluster
177	the duty or responsibility to promote and protect the welfare of various entities.
158	the duty or responsibility to ensure safety, both for oneself and for others
87	the concept of respect for autonomy
83	the duty and responsibility towards family
82	well-being, specifically human well-being
81	the right to property
71	the duty to promote, maintain, uphold, and protect peace at various levels
69	the duty and responsibility towards the community
67	the duty or responsibility to protect and care for children
64	the duty to treat others with respect, equality, fairness, impartiality, kindness, and compassion

Table 10: Top 10 clusters with the most examples based on agglomerative clustering on values.

D.2 System performance ablations on VALUEPRISM-Test

Similarly, we also compare the outputs of different sized systems with Rouge-score against the GPT-4 outputs (See Table 13).

D.3 Values manifested in identifying hate speech

We run KALEIDO^{SYS} on Social Bias Frames (Sap et al. 2020), a dataset containing instances of online speech, some of which is labeled as hate speech and some of which is not. We look at the most frequent values generated, and find that the most common opposing values are *Respect (for others)*, *Equality*, *Tolerance*, and *Right to Privacy*, while the top supporting values are *Freedom of speech (or expression)*, *Humor*, *Honesty*, and *Right to freedom of speech*. These values represent 17% and 26% of the generated oppose/support value counts, respectively. Even though KALEIDO is not trained explicitly to recognize hate speech, it is able to surface values that are violated by hate speech, along with values that run counter to excessive speech moderation.

E CloudResearch Results

In our study, we collect 31k annotations from 612 annotators across 683 values, rights, and duties in the context of 100 situations. The annotators mark 1) each value for whether or not they agree with it and 2) whether they have an opinion or perspective that is missing from the data for a given situation. Results are found in Tables 18 and 20.

We find that annotators agree with the values, rights, and duties 81% of the time on average, and state that they have a missing perspective 30% of the time. Note that this is less agreement and more missing perspectives than we saw for the quality annotation. This is not surprising to us, as some annotators may consider a value output high-quality and reasonable according to someone, even though they may not agree with it themselves (a much more subjective measure). Additionally, people were allowed to list missing perspectives in a free-form text box. Responses are hand-coded by the authors as having content or not, and the variable “has a missing perspective” is binarized. We find that many of the non-null responses merely state that the person has a missing perspective, not what it is (e.g., “Yes”) or do not map cleanly onto the values, rights, and duties framework (e.g., “do what is correct”, “Take care of orphan is not wrong”). This highlights a weakness of the framework: not all perspectives fit neatly into it.

We conduct 2 statistical analyses on the data. First, with ANOVA testing for each demographic group, we did not find statistically significant differences in agreement or missing perspective rates between subgroups (Table 16). We also compute a regression analysis for ordinal variables and most common subgroup vs. rest 2-sided t-tests for categorical groups (Table 17) and did not find significant results except for 2 groups: male (vs. non-male) folks were more likely and straight (vs. non-straight) folks were less likely to share a missing perspective ($p = .021, .029$ respectively). However,

	Judgment "Bad"	Judgment "OK"	Judgment "Good"
	Duties	Duties	Duties
Supports	Duty to provide for family Duty to save lives Duty of justice Duty to protect family Duty to follow orders Duty to be respectful Duty to protect public health Duty to not harm others Duty to abide by the law Duty to oneself	Duty to follow the law Duty to respect others' autonomy Duty to self-care Duty of honesty Duty to self-preservation Duty to tell the truth Duty to be truthful Duty to be honest Duty not to kill Duty of respect	Duty of non-discrimination Duty of solidarity Duty to assist others in need Duty to provide assistance Duty to be charitable Duty to obey laws Duty to respect property Duty to self-care Duty to respect others' autonomy Duty to follow the law
Opposes			
	Rights	Rights	Rights
Supports	Sovereignty Right to a minimum standard of living Right to information	Freedom of association Right to freedom of expression Right to family life Right to marry Right to peaceful assembly	Equality Animal rights Right not to be harmed Right to personal security Right to dignity
Opposes	Right to equal treatment Property Right to freedom of movement Right to dignity Right to education	Right to truth Right to bodily autonomy Right to truthful information Right to free speech	Right to conduct business Property Rights Right to Property Right to self-defense
	Values	Values	Values
Supports	Unity Personal freedom Personal autonomy Deterrence Respect for authority Justice Respect for property Safety Autonomy Respect	Cleanliness Individual autonomy Financial stability Cultural preservation Economic stability Public order Social harmony Individual freedom Truthfulness Tradition	Compassion Justice Well-being Equality Happiness Personal responsibility Self-reliance Prevention of future harm Work-life balance Life
Opposes			

Table 11: Out of all Values/Rights/Duties that appear more than 4 times, we look at those co-occurring with a given machine-generated judgment out of -1, 0, 1 (bad, ok, good), with a probability ≥ 0.5 and display the top-5 (or less) supporting or opposing ones in this table.

as we run 32 hypothesis tests,²⁰ it is likely that false positives may have slipped through²¹. After performing a Bonferroni correction for performing multiply hypothesis tests, these results are no longer significant.

As we do not find statistically significant results after correcting for the number of tests we perform, we hope to do more extensive, larger-scale surveying in future work.

E.1 Agreement examples

In general, we see that most people agree with most of the outputs ($>80\%$ agreement rate), which suggests that most of VALUEPRISM represents agreed upon human values, at least for the majority of cases. However, some claims seemed to be more or less agreed upon - see Table 14 for examples of

²⁰8 demographics * 2 dependent variables (agree and missing) * 2 kinds of analyses (ANOVA + regression or t-test)

²¹<https://xkcd.com/882/>

the most and least agreed upon claims, as well as claims with average agreement.

However, there are a small number of cases for which there was much higher agreement in one subgroup than another. See Table 15 for examples where there is particular divergence by political orientation.

F Dataset Generation

F.1 Dataset Creation Details

Situations We source our situations about which to reason from a set of 1.3M user-submitted situations, and curate the dataset by filtering out situations that are not actions or unrelated to morality (as labeled in a few-shot manner²² by Flan-T5 (Chung et al. 2022)). We also filter out any questions using keyword matching.

²²Few-shot filtering prompts are found in Appendix M.1.

Model	Relev. Acc \uparrow	Valence Acc \uparrow	Gen. Perp \downarrow	Expl. Perp \downarrow
KALEIDO 3B	88.4	80.8	2.23	3.14
– relevance	-	81.5	2.24	3.14
– valence	88.3	-	2.24	3.14
– generation	88.4	80.8	-	3.14
– explanation	88.5	80.8	2.26	-
single task	88.7	77.2	2.25	3.16
11B	89.1	81.9	2.22	2.99
770M	87.2	79.2	2.34	3.52
220M	83.5	74.5	2.53	4.23
60M	66.0	59.7	2.86	5.70

Table 12: Effect of Dataset Mixtures and Model Size on test set performance.

Model	R-1	R-2	R-L-Sum
KALEIDO ^{SYS} 3B	.54	.23	.51
– relevance	.52	.22	.48
– text similarity	.53	.22	.49
11B	.55	.23	.51
770M	.54	.22	.50
220M	.52	.21	.49
60M	.49	.18	.45

Table 13: Rouge scores (F1) on test set

We note that an outsize proportion of the dataset involves toxic, NSFW, or sexually explicit content. In the interest of having a diversity of situations, we label for these attributes¹ using Flan-T5 (Chung et al. 2022). We sample 95% of our situations from those that have less toxic/NSFW/explicit content, and the other 5% uniformly from the rest of the data so as to include the entire spectrum of inputs. We find that this succeeds in increasing the diversity of the dataset, as measured by unique n-grams divided by the length of the dataset (dist-2: .23→.36, dist-3: .54→.67).

Symbolic Knowledge Distillation using LLMs After experimentation, we find initial success in using GPT-4 (OpenAI 2023) to generate values. As is often the case, solution verification is often easier than solution generation, and we find it to be quite a challenging task to generate a comprehensive set of values, rights, and duties that could be considered for a situation. While we find that we as authors can provide more accurate (precise) lists, we anecdotally find that GPT-4 often does better at breadth (recall). See Appendix A for examples. Additionally, because the generation task requires such cognitive effort, the cost to hire crowdworkers to generate a dataset of the size that we desire would be prohibitive. As such, we follow prior work [CITE, add in from related work] and decide to use a LLM to create a synthetic dataset of values. We verify the quality (Section 4.1) and representativeness (Section 4.2) of the outputs using human annotators.

Values, Rights, and Duties Generation Given the set of 31k situations, we prompt GPT-4 (OpenAI 2023) to generate relevant values, rights, and duties²³, along with an open-text explanation. Given the output, the model also predicts whether the corresponding value, right, or duty supports (or justifies) the action, opposes (or condemns) the action, or could either support or oppose depending on the context or interpretation. The cost to generate the entire dataset was \$1,043.80.

While the data was generated in a batch manner to produce all values and related data at once, we exploit the structure of the generated data to cast the Generation, Valence, and Explanation tasks as sequence-to-sequence (seq2seq) tasks. The relevance task data is sampled contrastively, where positive examples are values generated by GPT-4 for the situation negative samples are drawn from other generated values. We split the data (by actions) into train/validation/test splits of 80%/10%/10% respectively (See Table 7).

G Model Training Details

For training, we set our model size at 3 billion parameters using the T5 encoder-decoder architecture (Raffel et al. 2020), and test the following hyperparameters: weight initialization in {t5-3B, flan-t5-xl}, learning rate in {1e-4, 3e-4, 1e-5, 3e-5}, and a dataset mixture of either {Generation, Relevance, Valence} or {Generation, Relevance, Valence, Explanation}. Because the explanation is post-hoc and of lesser interest to us than the other tasks, we choose the optimal set up on the validation set of the task mixture without the explanation task.

We conduct a grid search and settle with learning rate at 3e-5 and a batch size of 32 with a mixture of all four tasks. For further analysis of the relationship of data mixture and model size with performance, see App. D.1.

We train with Huggingface’s Trainer (Wolf et al. 2020) for 4 epochs with early stopping and a batch size of 32, although we find that the majority of runs start to overfit after about 2 epochs. Training takes 19 hours per run on two A100 GPUs.

We fix hyperparameters for the remainder of our experiments at the optimal hyperparameters: flan-t5-xl, 3e-5, and the mixture including explanations (which we find to assist generalization on the non-explanation tasks). For further analysis of the relationship of data mixture and model size with performance, see Section D.1. We refer to our default 3B trained model as KALEIDO.

H System Details

H.1 Algorithm

See Algorithm 1.

H.2 System Parameters

For KALEIDO^{SYS} 3B, we use these parameters for all experiments, which were found by maximizing RougeLSum f1-score VALUEPRISM-val. We also fix the number of generations at 100 and take the top generations with beam search. Parameters determine the threshold for embedding cosine

²³For the prompt used, please refer to Appendix M.

Situation	Claim	Description
Hacking a company to give them information about their flaws	Information transparency: Providing information about flaws within a system might promote openness and transparency, which can be seen as a positive value. However, such information in the wrong hands could be dangerous.	<i>Claims for which all people surveyed agreed. 100% agreement</i>
Preventing theft	Personal Responsibility: Encouraging individuals to take steps to prevent theft can promote a sense of personal responsibility for one's belongings and actions.	100% agreement
Firing Mike	Duty to uphold company policies: If Mike's actions violated company policies, management has a duty to take appropriate actions, which in this case may include firing.	100% agreement
Frowning at a friend	Respect: Not frowning at a friend if the situation doesn't warrant it could be a way of respecting their feelings.	<i>Claims with closest to average rates of agreement. 81% agreement</i>
Shooting a courier in the head who was responsible for the destruction of an entire city.	Duty to protect others: If the courier poses an ongoing threat, there might be a duty to protect others from harm, which could include killing the courier.	82% agreement
Aiding a fugitive by providing food or shelter in 1851 Pennsylvania.	Compassion: Providing food and shelter for a fugitive showcases compassion towards those in need, regardless of their legal status.	82% agreement
Tipping to your landlord	Duty of gratitude: Tenants might have a moral duty to show gratitude to those who help or provide a service, even if tipping is not customary for landlords.	<i>Claims for which the least people surveyed agreed. 8% agreement</i>
Killing a potato to save a frog.	Respect for nature: Some might argue that interfering with nature by deciding which being lives or dies could disrupt the natural balance, thus showing a lack of respect for nature.	24% agreement
Somali pirates need to feed their family so they ransom a ship	Justice: The pirates are attempting to resolve economic inequality by gaining money, which they perceive as a just cause.	24% agreement

Table 14: GPT-4 outputs from VALUEPRISM with most, average, and lowest levels of agreement based from CloudResearch study.

similarity and the ngram overlap threshold for deduplicating, and the relevance score at which to drop poor outputs. There is a separate threshold for each category of value, right, and duty.

```

1 "embed_threshold": {"Value': 0.53,
  'Right': 0.63, 'Duty': 0.55}",
2 "ngram_threshold": "0.05",
3 "relevance_threshold": {"Value':
  0.77, 'Right': 0.82, 'Duty': 0.9}

```

I Annotation Details

For all studies, we carefully monitored the time workers were spending on our tasks, and ensured minimum average hourly wages of \$15-\$25 USD.

I.1 Quality Annotation

For this study, 3 crowdworkers for GPT-4 output for 3k situations (10% of VALUEPRISM). See Figure 10 for the template used. For this study, note that we do not ask annotators to provide their own judgments of the situation, but merely to assess the relevance of the generations, which we expect to have much lower variation (e.g., someone may see how a value could be relevant for someone else while disagreeing with it themselves).

On an individual annotation level, 97% of the output annotations are “good” and 95% of the valence annotations are “correct.” We find a Fleiss’ kappa of .04 for quality and .12 for the valence labels. While this seems like low agreement, this is a known phenomena that happens with highly skewed data (Randolph 2005). On the other hand, absolute agreement levels of 91%/87% respectively are quite high. The total cost of this study was \$4,680.00 USD.

Situation	Claim	Description
redistributing rich people's land to poor people Aiding a fugitive by providing food or shelter in 1851 Pennsylvania. selling weapons to saudi arabia	Efficiency: Redistribution may lead to more efficient land use if previously underutilized land is given to those in need. Duty to protect life: One has a duty to protect life, which includes providing essentials like food and shelter, even to a fugitive. Right to life: If the weapons sold are used to kill innocent civilians, it would infringe upon their right to life.	<i>Claims with most liberal and least conservative agreement.</i> 11/14 liberal agreement, 0/3 conservative agreement 16/17 liberal agreement, 7/17 conservative agreement 19/19 liberal agreement, 9/14 conservative agreement
giving people things for free Hug someone Scolding your saintly father during your nephew's wedding	Personal Responsibility: Some may argue that individuals should earn what they receive, and providing things for free may undermine this value. Social norms: There may be a duty to follow certain social norms and customs, which in some cultures and contexts may involve hugging others as a form of greeting or celebration. Emotional expression: Sharing your feelings openly, even if they are negative, can be seen as a value in honest and open communication.	<i>Claims with most conservative and least liberal agreement.</i> 8/9 conservative agreement, 4/16 liberal agreement 13/13 conservative agreement, 13/22 liberal agreement 6/6 conservative agreement, 23/29 liberal agreement

Table 15: GPT-4 outputs from VALUEPRISM with most difference in agreement by political orientation.

Demographic Categories	Agree p-value	Missing p-value
Age	0.891	0.191
Gender	0.661	0.162
Race	0.369	0.828
Political Orientation	0.897	0.889
Sexual Orientation	0.763	0.141
Religion	0.620	0.389
Religiosity	0.995	0.491
Education	0.194	0.132

Table 16: ANOVA p-values for Demographic Categories against the null hypothesis *no difference between subgroups*.

I.2 CloudResearch Annotation

See Figures 11, 12, 13 for screenshots of the annotation tasks. Surveys were filled out in Qualtrics after crowdworkers were recruited using CloudResearch. The total cost of the study was \$3,562.95 USD.

For this task, we do not report agreement levels as we were specifically probing for subjective measurements (opinions) as opposed to objective ones. For further analysis, see App. E.

For the agreement task, the authors manually went through responses to code them as either having a missing perspective, or being devoid of content (e.g., , *No, Nothing I can think of, I agree with it all*, etc.). For purposes of the analysis, this is then binarized. Empty responses are coded as being devoid of content (no missing perspective).

I.3 Value, Right, and Duty Batch Comparison vs. GPT-4

For this study, a fixed set of 200 test examples was used for all systems, each pair of which was annotated by 2 annotators (200 situations * 2 annotators * 11 systems = 4.4k annotations).

See Figure 14 for the MTurk template used. For this study, only the set of values, rights, and duties was used, without the explanation or valence scores shown, in order to disentangle the value, right, and duty generation from the other aspects.

The interannotator agreement metrics are as follows: For overall preference (Answer.pref), Fleiss' kappa: 0.170, Krippendorff's alpha: 0.373, Cohen's kappa: 0.170, percent agreement (pairwise): 0.574, and tie-discounted percent agreement (pairwise): 0.588. For coverage (Answer.div),

Demographic Group	Type	Majority Class	Agree (p-value)	Missing (p-value)
Age	<i>Ordinal</i>	35-44 (164)	0.995	0.112
Religiosity	<i>Ordinal</i>	I am not religious (210)	0.602	0.459
Education	<i>Ordinal</i>	Bachelor's degree (222)	0.718	0.972
Political Orientation	<i>Ordinal</i>	Lean liberal (139)	0.322	0.316
Race	<i>Categorical</i>	White / Caucasian (168)	0.486	0.428
Gender	<i>Categorical</i>	Man/Male (258)	0.117	0.021
Sexual Orientation	<i>Categorical</i>	Heterosexual (straight) (390)	0.475	0.029
Religion	<i>Categorical</i>	Christian (228)	0.107	0.187

Table 17: P-values for Agree and Missing based on Demographic Group against the null hypothesis *no correlation* for ordinal and *no difference between majority class and rest* for categorical variables.

Algorithm 1: Generation of diverse values, rights, and duties

Require: Action A , Relevance thresholds T_v, T_r, T_d , Similarity thresholds S_v, S_r, S_d , 1-gram similarity N_{sim} , Model M , Beam search number B

- 1: $VRD_{\text{gen}} \leftarrow M(A, B)$ \triangleright Generate top B beams
- 2: $R \leftarrow \text{Relevance}(VRD_{\text{gen}}, M)$ \triangleright Relevance probs
- 3: Sort (VRD_{gen}, R) in decreasing order by R
- 4: $VRD \leftarrow \emptyset$
- 5: **for** each $v_i \in VRD_{\text{gen}}$ in the order of R **do**
- 6: $\text{type} \leftarrow \text{GetType}(v_i)$ \triangleright Get the type (value, right, duty) of v_i
- 7: $T \leftarrow T_{\text{type}}$ \triangleright Select threshold T based on type
- 8: $S \leftarrow S_{\text{type}}$ \triangleright Select similarity threshold S based on type
- 9: **if** $r_i < T$ **then** \triangleright If relevance too low, continue
- 10: **continue**
- 11: **end if**
- 12: $VRD_{\text{type}} \leftarrow \{v \in VRD : \text{GetType}(v) == \text{type}\}$
- 13: $O_{\text{1gram}} \leftarrow \text{Overlap}_{\text{1gram}}(v_i, VRD_{\text{type}})$ \triangleright Calculate 1-gram overlap scores
- 14: $O_{\text{cosine}} \leftarrow \text{Overlap}_{\text{cosine}}(v_i, VRD_{\text{type}})$ \triangleright Calculate cosine overlap scores
- 15: **if** $\max(O_{\text{1gram}}) < N_{\text{sim}}$ and $\max(O_{\text{cosine}}) < S$ **then** \triangleright Add if not too similar
- 16: $VRD \leftarrow VRD \cup \{v_i\}$
- 17: **end if**
- 18: **end for**

Fleiss' kappa: 0.394, Krippendorff's alpha: 0.240, Cohen's kappa: 0.394, percent agreement (pairwise): 0.666, and tie-discounted percent agreement (pairwise): 0.720. Finally, for accuracy (Answer.correct), Fleiss' kappa: 0.052, Krippendorff's alpha: 0.257, Cohen's kappa: 0.053, percent agreement (pairwise): 0.459, and tie-discounted percent agreement (pairwise): 0.537. The total cost for this study was \$3,168.00 USD.

I.4 Explanation, Valence, and Relevance Annotations

For this evaluation, we run a fixed set of 698 values, rights, and duties associated with 100 situations from the test set.

3 annotators each mark each value, right, or duty for relevance, each explanation for quality, and each valence lable for correctness. Testing 6 systems, this comes to 38k total annotations. The total cost of this study was \$842.40 USD.

We get the following annotator agreements: For relevance, Fleiss' kappa: .073, pairwise agreement: .903; explanation quality, Fleiss' kappa: .180, pairwise agreement: .842; valence correctness, Fleiss' kappa: .345, pairwise agreement: .860.

J Extreme bad/good actions

Good

- 1 Being sincerely considerate to others
- 2 Making a friend smile
- 3 Try to be truthful
- 4 Treat others with respect
- 5 Be generous when possible
- 6 Doing what you believe to be right
- 7 Do what you can to improve your community
- 8 Make the world a better place
- 9 Give to those in need if you can afford to
- 10 Respect people's boundaries if at all possible

Bad

- 1 Killing someone without cause
- 2 Taking someone's property without permission or a pressing need
- 3 Genocide
- 4 Causing physical or emotional harm to others because you feel like it
- 5 Discriminating against someone based on their race
- 6 Breaking a promise without a good reason
- 7 Lying to someone just to hurt them
- 8 Neglecting your child
- 9 Emotionally manipulating people
- 10 Buying a bunch of food just to let it go to waste

K Universal Declaration of Human Rights

Below, we have a list of rights from the UN's Universal Declaration of Human Rights²⁴, along with the corresponding right in VALUEPRISM.

- 1 No one shall be subjected to torture or to cruel, inhuman or degrading treatment or punishment . -> Right to not experience torture or cruel treatment
- 2 All are equal before the law and are entitled without any discrimination to equal protection of the law. All are entitled to equal protection against any discrimination in violation of this Declaration and against any incitement to such discrimination. -> Equality before the law
- 3 No one shall be subjected to arbitrary arrest, detention or exile. -> Right to not be unjustly detained
- 4 Everyone charged with a penal offence has the right to be presumed innocent until proved guilty according to law in a public trial at which he has had all the guarantees necessary for his defence. -> The right to due process
- 5 Everyone has the right to freedom of movement and residence within the borders of each state. -> The right to freedom of movement
- 6 Everyone has the right to seek and to enjoy in other countries asylum from persecution. -> The right to seek asylum
- 7 Everyone has the right to a nationality. -> Right to nationality
- 8 Men and women of full age, without any limitation due to race, nationality or religion, have the right to marry and to found a family. They are entitled to equal rights as to marriage, during marriage and at its dissolution. -> The right to marry
- 9 The family is the natural and fundamental group unit of society and is entitled to protection by society and the State. -> The right to protection

- 10 No one shall be arbitrarily deprived of his property. -> Private property
- 11 Everyone has the right to freedom of opinion and expression; this right includes freedom to hold opinions without interference and to seek, receive and impart information and ideas through any media and regardless of frontiers. -> The right to express one's opinions
- 12 Everyone has the right to freedom of peaceful assembly and association. -> The right to peaceful assembly
- 13 No one may be compelled to belong to an association. -> The right to freedom of association
- 14 Everyone has the right to take part in the government of his country, directly or through freely chosen representatives. -> Right to participate in political processes
- 15 Everyone has the right to work, to free choice of employment, to just and favourable conditions of work and to protection against unemployment. -> The right to work
- 16 Everyone, without any discrimination , has the right to equal pay for equal work. -> Right to equal pay
- 17 Everyone has the right to form and to join trade unions for the protection of his interests. -> The right to collective action
- 18 Everyone has the right to a standard of living adequate for the health and well-being of himself and of his family, including food , clothing, housing and medical care and necessary social services, and the right to security in the event of unemployment, sickness, disability, widowhood, old age or other lack of livelihood in circumstances beyond his control. -> Access to basic necessities
- 19 Motherhood and childhood are entitled to special care and assistance. All children, whether born in or out of wedlock, shall enjoy the same social protection . -> Right to parental support
- 20 Everyone has the right to education. Education shall be free, at

²⁴<https://www.un.org/en/about-us/universal-declaration-of-human-rights>

- least in the elementary and fundamental stages. Elementary education shall be compulsory. Technical and professional education shall be made generally available and higher education shall be equally accessible to all on the basis of merit. -> The right to education
- 21 Everyone has the right to life, liberty and security of person.
-> Right to personal liberty, Right to life
- 22 No one shall be held in slavery or servitude; slavery and the slave trade shall be prohibited in all their forms. -> The right to not be enslaved
- 23 Everyone has the right to recognition everywhere as a person before the law. -> Right to recognition
- 24 Everyone has the right to an effective remedy by the competent national tribunals for acts violating the fundamental rights granted him by the constitution or by law. -> The right to due process
- 25 Everyone is entitled in full equality to a fair and public hearing by an independent and impartial tribunal, in the determination of his rights and obligations and of any criminal charge against him. -> The right to a fair trial
- 26 No one shall be subjected to arbitrary interference with his privacy, family, home or correspondence, nor to attacks upon his honour and reputation. Everyone has the right to the protection of the law against such interference or attacks. -> None
- 27 Everyone has the right to leave any country, including his own, and to return to his country. -> The right to freedom of movement
- 28 No one shall be arbitrarily deprived of his nationality nor denied the right to change his nationality. -> Right to nationality
- 29 Marriage shall be entered into only with the free and full consent of the intending spouses. -> Right to free choice of partner, Right to marry
- 30 Everyone has the right to own property alone as well as in association with others. -> The right to property
- 31 Everyone has the right to freedom of thought, conscience and religion ; this right includes freedom to change his religion or belief, and freedom, either alone or in community with others and in public or private, to manifest his religion or belief in teaching, practice, worship and observance. -> Right to form one's own beliefs
- 32 Everyone has the right of equal access to public service in his country. -> Right to access services
- 33 Everyone, as a member of society, has the right to social security and is entitled to realization, through national effort and international co-operation and in accordance with the organization and resources of each State, of the economic, social and cultural rights indispensable for his dignity and the free development of his personality. -> Right to social and economic welfare
- 34 Everyone who works has the right to just and favourable remuneration ensuring for himself and his family an existence worthy of human dignity, and supplemented, if necessary, by other means of social protection. -> Right to just and favorable working conditions
- 35 Everyone has the right to rest and leisure, including reasonable limitation of working hours and periodic holidays with pay. -> Right to engage in leisure activities
- 36 Education shall be directed to the full development of the human personality and to the strengthening of respect for human rights and fundamental freedoms. It shall promote understanding, tolerance and friendship among all nations, racial or religious groups, and shall further the activities of the United Nations for the maintenance of peace. ->

- Educational opportunities
- 37 Parents have a prior right to choose the kind of education that shall be given to their children. -> Right of parents to make decisions for their children
- 38 Everyone has the right freely to participate in the cultural life of the community, to enjoy the arts and to share in scientific advancement and its benefits. -> Access to cultural experiences
- 39 Everyone has the right to the protection of the moral and material interests resulting from any scientific, literary or artistic production of which he is the author. -> Right to protect one's work and creations
- 40 Everyone is entitled to a social and international order in which the rights and freedoms set forth in this Declaration can be fully realized. -> The right to freedom
- 41 In the exercise of his rights and freedoms, everyone shall be subject only to such limitations as are determined by law solely for the purpose of securing due recognition and respect for the rights and freedoms of others and of meeting the just requirements of morality, public order and the general welfare in a democratic society. -> Legal autonomy

L ETHICS zero-shot templates

All templates were hand-written by the authors. No more than 10 templates were attempted per task, and the template with the highest performance on the first 1000 instances of the training set were selected.

L.1 Justice

1 Valence({scenario}; Value: Fairness)
 $P(\text{good}) = P(\text{support}), P(\text{bad}) = P(\text{oppose})$

L.2 Deontology

1 Valence(In response to {scenario}, saying {excuse}; Duty: Duty to have a valid reason)
 $P(\text{good}) = P(\text{support}), P(\text{bad}) = P(\text{oppose})$

L.3 Virtue

1 Relevance({scenario}; Value: {trait})
 $P(\text{traitfits}) = P(\text{relevant}), P(\text{traitdoesnotfit}) = P(\text{notrelevant})$

L.4 Utilitarianism

- 1 p_better = Valence({scenario1} is more pleasurable than {scenario2}; Value: Pleasure) ['support'] + Valence({scenario2} is less pleasurable than {scenario1}; Value: Pleasure) ['oppose'] + Valence({scenario2} is more pleasurable than {scenario1}; Value: Pleasure) ['oppose'] + Valence({scenario1} is less pleasurable than {scenario2}; Value: Pleasure) ['oppose']
- 2 p_worse = Valence({scenario1} is more pleasurable than {scenario2}; Value: Pleasure) ['oppose'] + Valence({scenario2} is less pleasurable than {scenario1}; Value: Pleasure) ['oppose'] + Valence({scenario2} is more pleasurable than {scenario1}; Value: Pleasure) ['support'] + Valence({scenario1} is less pleasurable than {scenario2}; Value: Pleasure) ['support']
- 3
- 4
- 5
- 6
- 7
- 8

L.5 Commonsense

- 1 Valence({scenario}; Duty: Duty to do the right thing according to commonsense morality)
- $P(\text{good}) = P(\text{support}), P(\text{bad}) = P(\text{oppose})$

M GPT-4 Values, Rights, Duties Prompt

System Message:

- 1 Given a situation, try to reason about its morality by breaking down relevant values, rights, and duties. Please generate any and all potentially relevant values that people may hold, even if they sometimes oppose each other.
- 2
- 3 Values: These are the intrinsic goods or ideals that people pursue or cherish, such as happiness, well-being, justice, or freedom. Values are the desirable qualities that people may seek in their lives and in the world. They are often the guiding principles for individuals and societies, shaping goals, motivations, and preferences.
- 4
- 5 Rights: Rights are the entitlements or claims that individuals have

- against others or society, which are usually based on moral or legal grounds. These can be positive rights (e.g., the right to education, healthcare, or free speech) or negative rights (e.g., the right to not be harmed, enslaved, or discriminated against). Rights serve to protect the fundamental interests of individuals and establish certain boundaries that others must respect.
- 6
- 7 Duties: Duties are the moral obligations or responsibilities that individuals owe to others or to society at large. They are categorical reasons for doing or refraining from doing something, independent of whether we want to do or refrain from doing that thing. Duties can be weighty reasons, not easily overridden by competing concerns, and their violation may justify blame and self-blame (guilt). Duties can arise from relationships, social roles, or moral principles, and they guide our actions and decisions.
- 8 - Imperfect Duties: Some duties are "imperfect," meaning they are owed to no one in particular. For example, a duty to be charitable to the poor, but no specific poor person can demand charity from someone.
- 9 - Perfect Duties: Some duties are "perfect," meaning they are owed to someone in particular. For example, a duty not to plagiarize someone's work. When someone plagiarizes, there is a particular person who has been wronged - the person who was plagiarized.
- 10
- 11 For each value, right, and duty output a 1-3 sentence explanation of how it is related/expressed in the situation.
- 12
- 13
- 14 For each, list whether it supports the action, opposes the action, or could either support or oppose . Additionally, for each duty, list whether it is perfect or

imperfect (but do not include the word "perfect" or "imperfect" in the explanation).

- 15
- 16 If none are relevant, put N/A.
- 17
- 18 Example output:
- 19 Values:
- 20 - \{value1 OR N/A\}: explanation1 [supports/opposes/either]
- 21 - \{value2\}: explanation2 [supports /opposes/either]
- 22 ...
- 23
- 24 Rights:
- 25 - \{right1 OR N/A\}: explanation2 [supports/opposes/either]
- 26 ...
- 27
- 28 Duties:
- 29 - \{duty1\}: explanation1 [supports/ opposes/either, perfect/imperfect]
- 30 ...
- 31
- 32 An example if there are no relevant values, rights, or duties:
- 33 [Values/Rights/Duties]: N/A
- 34 "SITUATION" ->

User Message:

1 Situation: SITUATION

M.1 Flan-T5 Data Filtering Prompts

All few-shot examples are drawn from the original user demo queries and selected by the authors.

Because the prompts may contain offensive, NSFW, racially insensitive, or explicit material, we have decided not to include these prompts in the paper. They can be found online at <https://github.com/tsor13/kaleido>.

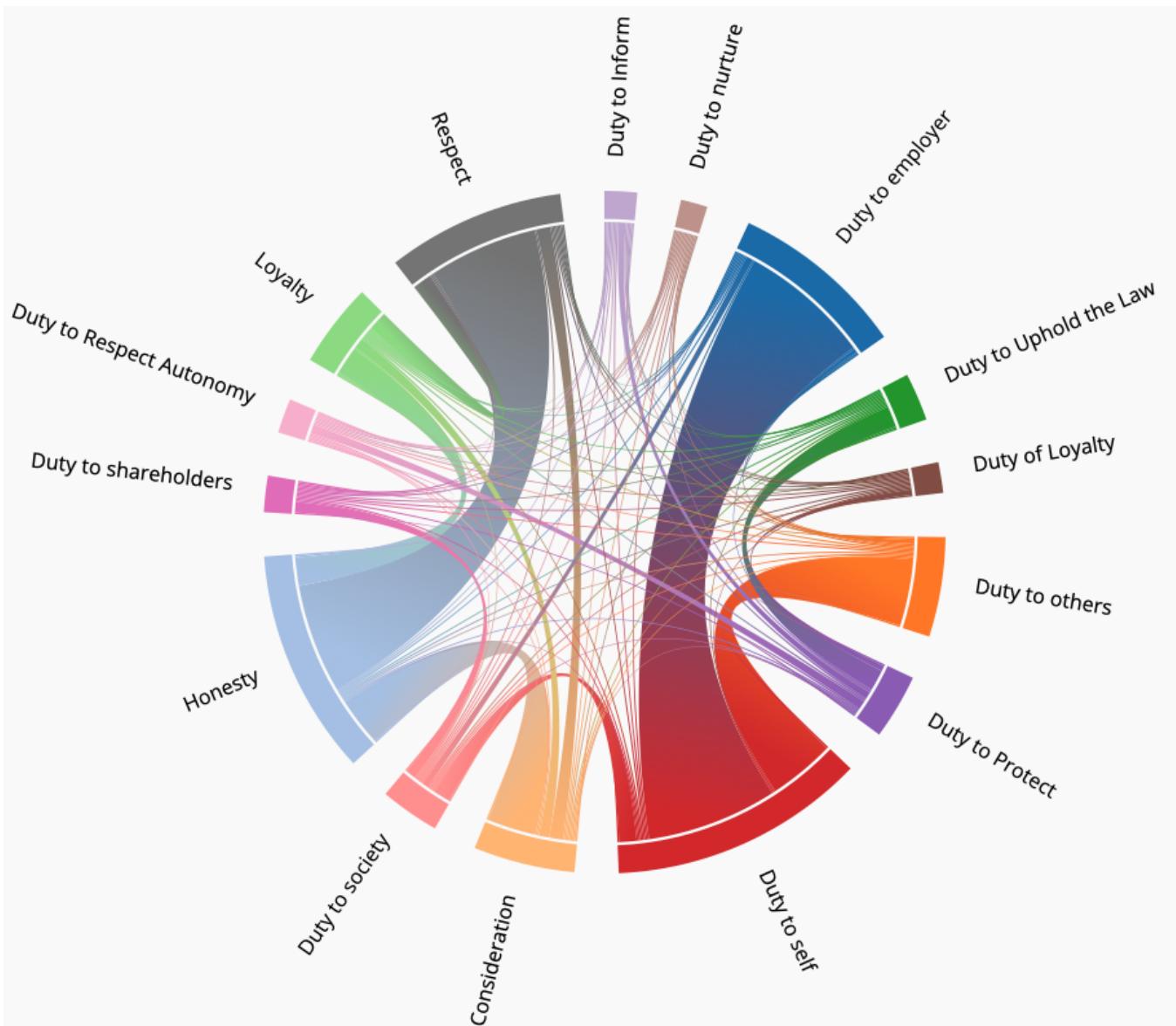


Figure 6: Co-occurrence counts of a subset of duties.

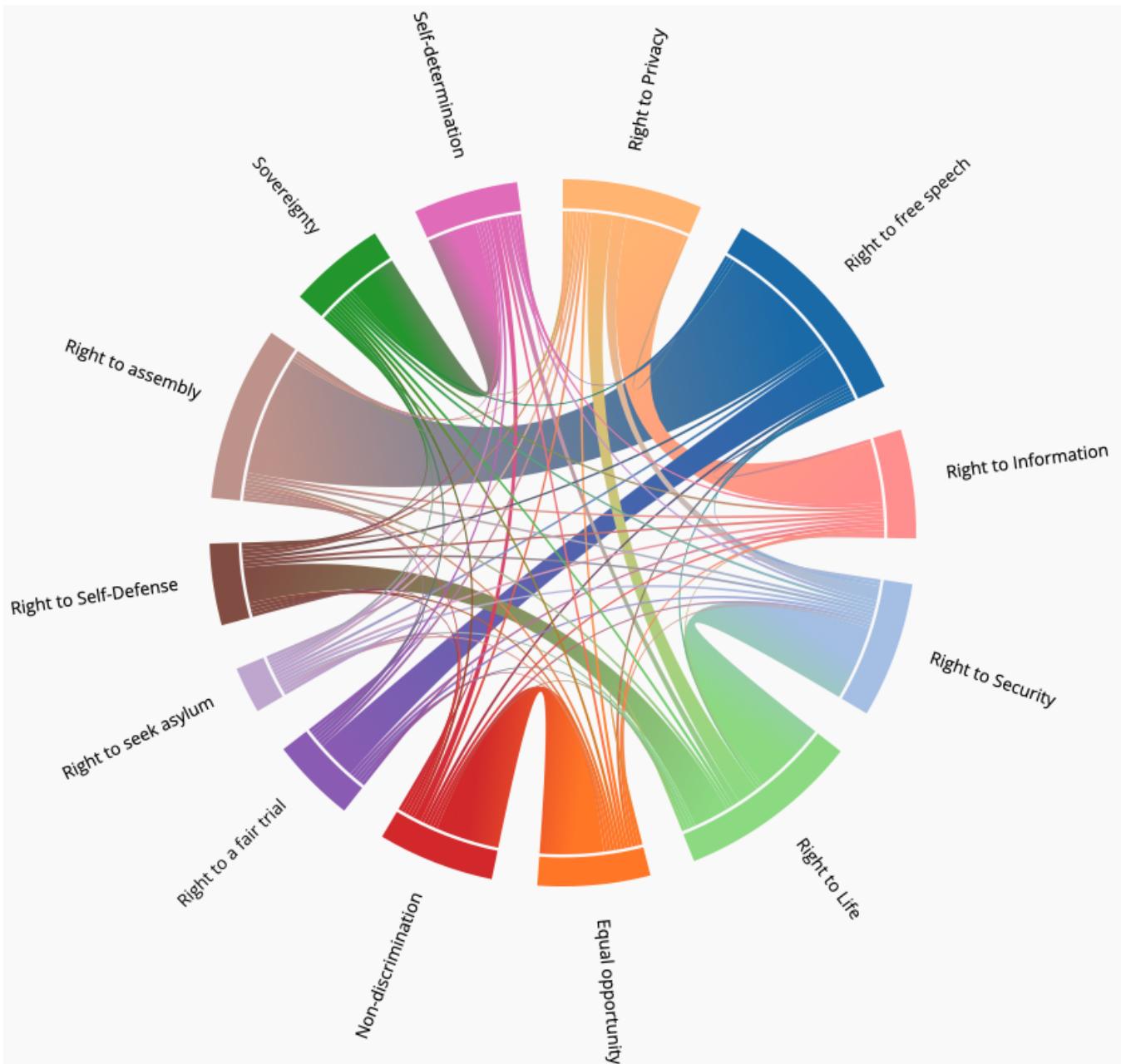
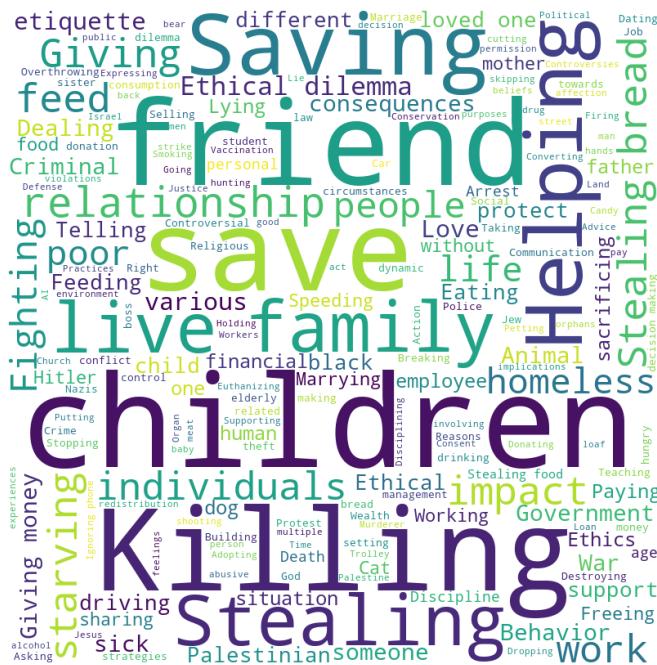


Figure 7: Co-occurrence counts of a subset of rights.



(a) situation



(b) value

Figure 8: word cloud on topics of situation and values

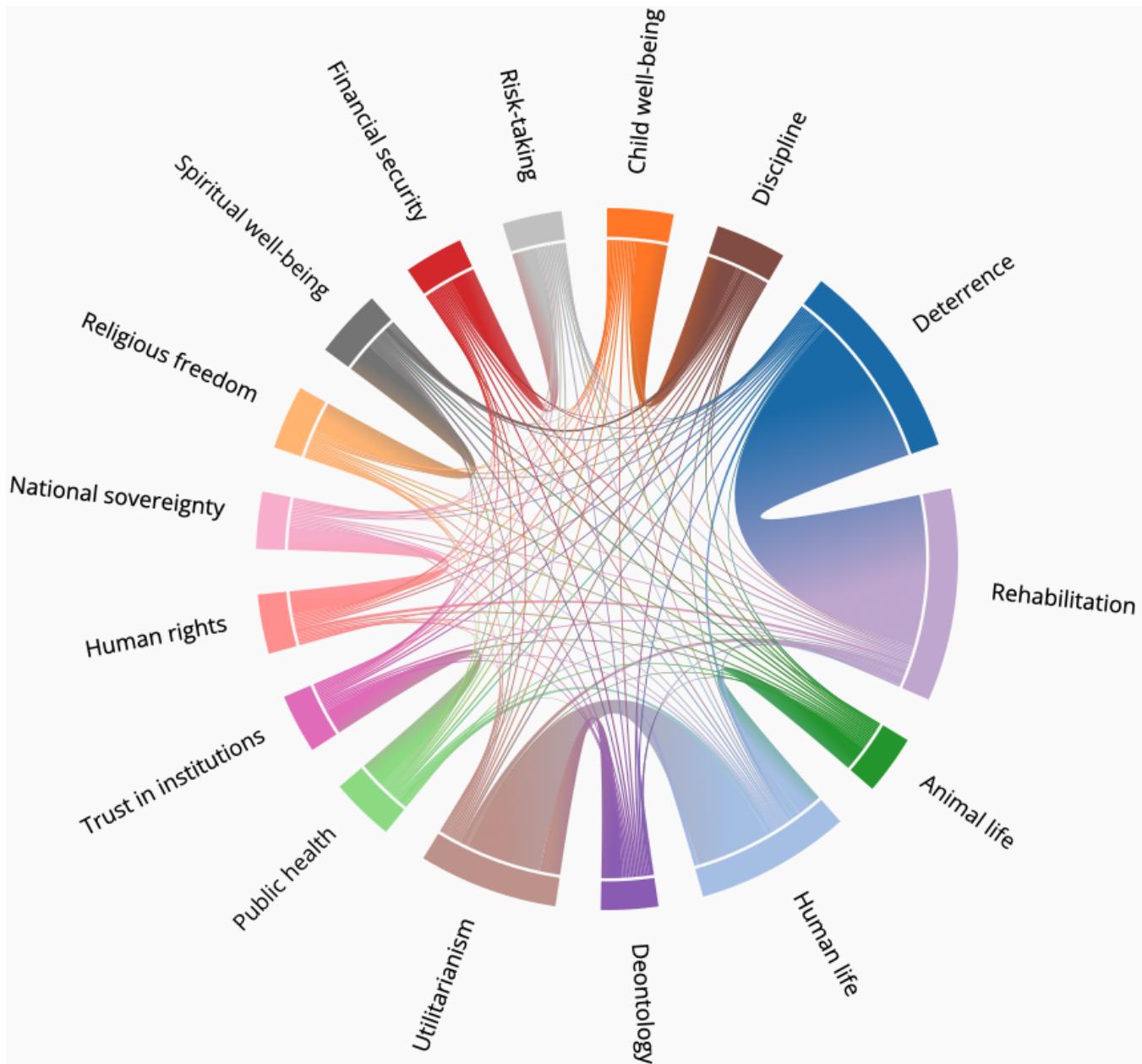


Figure 9: Co-occurrence counts of a subset of values.

Sub-group	Avg Agreement Rate	Std Error	# Participants
Age			
35-44	0.805	0.031	164
25-34	0.832	0.031	145
45-54	0.823	0.049	63
55-64	0.816	0.062	40
18-24	0.799	0.066	38
65 or older	0.776	0.116	14
Prefer not to say (Age)	0.781	0.239	4
Gender			
Man/Male	0.811	0.024	258
Woman/Female	0.820	0.027	201
Non-binary	0.885	0.160	5
Prefer not to say (Gender)	0.781	0.239	4
Race/Ethnicity			
White / Caucasian	0.805	0.031	168
Black / African American	0.817	0.036	115
Asian / Asian American	0.826	0.049	61
Hispanic / LatinX	0.768	0.074	34
Multiracial	0.802	0.100	17
Native American / First Nations	0.811	0.148	8
Other, please specify (Race)	0.781	0.239	4
Prefer not to say (Race)	0.833	0.265	3
Middle Eastern	0.716	0.454	2
Political Orientation			
Lean liberal	0.819	0.033	139
Liberal	0.826	0.038	102
Lean conservative	0.798	0.042	93
Moderate	0.823	0.041	89
Conservative	0.800	0.060	45
Sexual Orientation			
Heterosexual (straight)	0.816	0.020	390
Bisexual	0.828	0.065	35
Pansexual	0.852	0.095	15
Lesbian	0.802	0.133	10
Prefer not to say (Sexual Orientation)	0.776	0.187	6
Asexual	0.828	0.189	5
Gay	0.743	0.219	5
Prefer to self-describe	0.703	-	1
Queer	0.594	-	1
Religion			
Christian	0.805	0.026	228
Agnostic	0.840	0.037	98
Atheist	0.802	0.046	77
Spiritual but not religious	0.841	0.069	29
Prefer not to say (Religion)	0.829	0.133	9
Other, please specify (Religion)	0.786	0.155	8
Buddhist	0.842	0.149	7
Jewish	0.771	0.172	7
Hindu	0.873	0.237	3
Muslim	0.900	0.302	2
Religiosity, Education (continued in Table 19)			

Table 18: Agreement Rates by Demographic

Sub-group	Avg Agree.	Std Error	#
(continued from Table 18)			
Religiosity			
I am not religious	0.824	0.026	210
Very important	0.804	0.046	76
Moderately important	0.810	0.047	72
Center of my life	0.813	0.049	65
Not important at all, although I consider myself religious	0.797	0.065	39
Prefer not to say (Religiosity)	0.844	0.163	6
Education			
Bachelor's degree (for example: BA, AB, BS)	0.814	0.026	222
High school graduate - high school diploma or the equivalent (for example: GED)	0.794	0.054	58
1 or more years of college, no degree	0.826	0.051	56
Master's degree (for example: MA, MS, MEng, MEd, MSW, MBA)	0.832	0.052	52
Associate degree (for example: AA, AS)	0.835	0.060	39
Some college credit, but less than 1 year	0.816	0.079	25
Professional degree (for example: MD, DDS, DVM, LLB, JD)	0.754	0.163	8
Doctorate degree (for example: PhD, EdD)	0.762	0.302	3
Prefer not to say (Education)	0.778	0.295	3
9th, 10th, or 11th grade	0.805	-	1
Nursery school to 8th grade	0.657	-	1

Table 19: Agreement Rates by Demographic (continued)

Sub-group	Avg Missing Rate	Std Error	# Participants
Age			
35-44	0.344	0.037	164
25-34	0.287	0.038	145
45-54	0.349	0.061	63
55-64	0.243	0.069	40
18-24	0.185	0.064	38
65 or older	0.575	0.138	14
Prefer not to say	0.175	0.222	4
Gender			
Man/Male	0.338	0.030	258
Woman/Female	0.280	0.032	201
Non-binary	0.160	0.185	5
Prefer not to say	0.175	0.222	4
Race/Ethnicity			
White / Caucasian	0.308	0.036	168
Black / African American	0.321	0.044	115
Asian / Asian American	0.313	0.060	61
Hispanic / LatinX	0.358	0.084	34
Multiracial	0.446	0.125	17
Native American / First Nations	0.410	0.188	8
Other, please specify	0.650	0.283	4
Prefer not to say	0.467	0.365	3
Middle Eastern	0.100	0.316	2
Political Orientation			
Lean liberal	0.345	0.040	139
Liberal	0.302	0.046	102
Lean conservative	0.283	0.047	93
Moderate	0.311	0.049	89
Conservative	0.263	0.067	45
Sexual Orientation			
Heterosexual (straight)	0.291	0.023	390
Bisexual	0.565	0.085	35
Pansexual	0.266	0.119	15
Lesbian	0.159	0.123	10
Prefer not to say	0.200	0.181	6
Asexual	0.480	0.255	5
Gay	0.275	0.226	5
Prefer to self-describe	1.000	-	1
Queer	0.000	-	1
Religion			
Christian	0.330	0.031	228
Agnostic	0.264	0.045	98
Atheist	0.290	0.052	77
Spiritual but not religious	0.350	0.090	29
Prefer not to say	0.261	0.156	9
Other, please specify	0.538	0.191	8
Buddhist	0.229	0.165	7
Jewish	0.450	0.206	7
Hindu	0.150	0.259	3
Muslim	0.000	-	2
Religiosity, Education (continued in Table 21)			

Table 20: Missing perspective rates by demographic.

Sub-group	Avg Miss.	Std Error	#
Religiosity			
I am not religious	0.275	0.031	210
Very important	0.338	0.055	76
Moderately important	0.370	0.057	72
Center of my life	0.302	0.057	65
Not important at all, although I consider myself religious	0.372	0.079	39
Prefer not to say	0.143	0.159	6
Education Level			
Bachelor's degree (for example: BA, AB, BS)	0.297	0.031	222
High school graduate - high school diploma or the equivalent (for example: GED)	0.372	0.064	58
1 or more years of college, no degree	0.327	0.063	56
Master's degree (for example: MA, MS, MEng, MEd, MSW, MBA)	0.393	0.068	52
Associate degree (for example: AA, AS)	0.169	0.061	39
Some college credit, but less than 1 year	0.246	0.088	25
Professional degree (for example: MD, DDS, DVM, LLB, JD)	0.271	0.169	8
Doctorate degree (for example: PhD, EdD)	0.533	0.365	3
Prefer not to say (Education)	0.057	0.167	3
9th, 10th, or 11th grade	0.200	-	1
Nursery school to 8th grade	0.800	-	1

Table 21: Missing perspective rates by demographic. (continued)

Full Instructions [\(Expand/Collapse\)](#)

WARNING Current employees of the **University of Washington**, family members of UW employees, and UW students involved in this particular research are **not eligible** to complete this HIT.

CONTENT WARNING This task may contain content that some individuals may find unpleasant, offensive, or disturbing, including sexually suggestive or racially insensitive situations. Worker discretion is advised. Feel free to exit the task at any time.

Thank you for participating in this HIT! It's important to read the instructions thoroughly before proceeding.

Your task is to evaluate moral reasoning about specific actions based on values, rights, and duties as defined below:

Definitions [\(Expand/Collapse\)](#)

Values: These are the **intrinsic goods or ideals** that people pursue or cherish, such as happiness, well-being, justice, or freedom. Values are the desirable qualities that people may seek in their lives and in the world. They are often the guiding principles for individuals and societies, shaping goals, motivations, and preferences.

Rights: Rights are the **entitlements or claims** that individuals have against others or society, which are usually based on moral or legal grounds. These can be positive rights (e.g., the right to education, healthcare, or free speech) or negative rights (e.g., the right to not be harmed, enslaved, or discriminated against). Rights serve to protect the fundamental interests of individuals and establish certain boundaries that others must respect.

Duties: Duties are the **moral obligations or responsibilities** that individuals owe to others or to society at large. They are categorical reasons for doing or not doing something, regardless of whether we want to do or refrain from doing that thing. Duties can be weighty reasons, not easily overridden by competing concerns, and their violation may justify blame and self-blame (guilt). Duties can arise from relationships, social roles, or moral principles, and they guide our actions and decisions.

You'll evaluate the output (marked in purple) in two steps:

Task 1: Evaluate Output Quality

For each action, you will be presented with values, rights, and duties that make sense to consider when determining the morality of the action (see definitions below). Your job is to determine if these factors are high quality. A good output would be:

- **Relevant:** The value/right/duty should be a relevant thing to consider when thinking about the action and fits the definition of value/right/duty.
- **Sensible:** The explanation should make sense and be coherent.
- **Reasonable:** Is it a reasonable thing to think? Even if you disagree with the explanation, could someone else find it reasonable?

Note: If you're unsure, go with your instinct. If you find yourself not coming across any bad examples, that's okay - it might just be the luck of the draw.

Task 2: Evaluate Support or Opposition

Next, you'll see a guess for each explanation on whether it **supports** doing the action, **opposes** doing the action, or could **either support or oppose** the action depending on the interpretation or context. Your job is to determine if this guess is **correct** based on the provided explanation.

Missing Output

If you feel that something is missing, you can check the "Missing value/right/duties?" checkbox and write a few words about what's missing in the text box. It can be as long or as short as you like.

We appreciate your attention to detail and thoughtful responses. Thank you for your contribution!

Action	\${situation}																													
Values <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left; padding: 2px;">Values</th> <th style="text-align: left; padding: 2px;">Quality</th> <th style="text-align: left; padding: 2px;">Supports/Opposes</th> </tr> </thead> <tbody> <tr> <td style="padding: 2px;"><code>\$(value1): \${vexplanation1}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${vlabel1} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(value2): \${vexplanation2}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${vlabel2} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(value3): \${vexplanation3}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${vlabel3} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(value4): \${vexplanation4}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${vlabel4} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(value5): \${vexplanation5}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${vlabel5} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(value6): \${vexplanation6}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${vlabel6} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(value7): \${vexplanation7}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${vlabel7} the action. Yes or no?</td> </tr> <tr> <td colspan="3" style="text-align: center; padding: 2px;"><input type="checkbox"/> Missing value(s)?</td> </tr> </tbody> </table>				Values	Quality	Supports/Opposes	<code>\$(value1): \${vexplanation1}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel1} the action. Yes or no?	<code>\$(value2): \${vexplanation2}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel2} the action. Yes or no?	<code>\$(value3): \${vexplanation3}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel3} the action. Yes or no?	<code>\$(value4): \${vexplanation4}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel4} the action. Yes or no?	<code>\$(value5): \${vexplanation5}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel5} the action. Yes or no?	<code>\$(value6): \${vexplanation6}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel6} the action. Yes or no?	<code>\$(value7): \${vexplanation7}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel7} the action. Yes or no?	<input type="checkbox"/> Missing value(s)?		
Values	Quality	Supports/Opposes																												
<code>\$(value1): \${vexplanation1}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel1} the action. Yes or no?																												
<code>\$(value2): \${vexplanation2}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel2} the action. Yes or no?																												
<code>\$(value3): \${vexplanation3}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel3} the action. Yes or no?																												
<code>\$(value4): \${vexplanation4}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel4} the action. Yes or no?																												
<code>\$(value5): \${vexplanation5}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel5} the action. Yes or no?																												
<code>\$(value6): \${vexplanation6}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel6} the action. Yes or no?																												
<code>\$(value7): \${vexplanation7}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${vlabel7} the action. Yes or no?																												
<input type="checkbox"/> Missing value(s)?																														
Rights <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left; padding: 2px;">Rights</th> <th style="text-align: left; padding: 2px;">Quality</th> <th style="text-align: left; padding: 2px;">Supports/Opposes</th> </tr> </thead> <tbody> <tr> <td style="padding: 2px;"><code>\$(right1): \${rexplanation1}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${rlabel1} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(right2): \${rexplanation2}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${rlabel2} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(right3): \${rexplanation3}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${rlabel3} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(right4): \${rexplanation4}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${rlabel4} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(right5): \${rexplanation5}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${rlabel5} the action. Yes or no?</td> </tr> <tr> <td colspan="3" style="text-align: center; padding: 2px;"><input type="checkbox"/> Missing right(s)?</td> </tr> </tbody> </table>				Rights	Quality	Supports/Opposes	<code>\$(right1): \${rexplanation1}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel1} the action. Yes or no?	<code>\$(right2): \${rexplanation2}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel2} the action. Yes or no?	<code>\$(right3): \${rexplanation3}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel3} the action. Yes or no?	<code>\$(right4): \${rexplanation4}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel4} the action. Yes or no?	<code>\$(right5): \${rexplanation5}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel5} the action. Yes or no?	<input type="checkbox"/> Missing right(s)?								
Rights	Quality	Supports/Opposes																												
<code>\$(right1): \${rexplanation1}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel1} the action. Yes or no?																												
<code>\$(right2): \${rexplanation2}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel2} the action. Yes or no?																												
<code>\$(right3): \${rexplanation3}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel3} the action. Yes or no?																												
<code>\$(right4): \${rexplanation4}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel4} the action. Yes or no?																												
<code>\$(right5): \${rexplanation5}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${rlabel5} the action. Yes or no?																												
<input type="checkbox"/> Missing right(s)?																														
Duties <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left; padding: 2px;">Duties</th> <th style="text-align: left; padding: 2px;">Quality</th> <th style="text-align: left; padding: 2px;">Supports/Opposes</th> </tr> </thead> <tbody> <tr> <td style="padding: 2px;"><code>\$(duty1): \${dexplanation1}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${dlabel1} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(duty2): \${dexplanation2}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${dlabel2} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(duty3): \${dexplanation3}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${dlabel3} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(duty4): \${dexplanation4}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${dlabel4} the action. Yes or no?</td> </tr> <tr> <td style="padding: 2px;"><code>\$(duty5): \${dexplanation5}</code></td> <td style="padding: 2px;">This output is good (relevant, sensible, and reasonable). Yes or no?</td> <td style="padding: 2px;">This output \${dlabel5} the action. Yes or no?</td> </tr> <tr> <td colspan="3" style="text-align: center; padding: 2px;"><input type="checkbox"/> Missing duty(s)?</td> </tr> </tbody> </table>				Duties	Quality	Supports/Opposes	<code>\$(duty1): \${dexplanation1}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel1} the action. Yes or no?	<code>\$(duty2): \${dexplanation2}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel2} the action. Yes or no?	<code>\$(duty3): \${dexplanation3}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel3} the action. Yes or no?	<code>\$(duty4): \${dexplanation4}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel4} the action. Yes or no?	<code>\$(duty5): \${dexplanation5}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel5} the action. Yes or no?	<input type="checkbox"/> Missing duty(s)?								
Duties	Quality	Supports/Opposes																												
<code>\$(duty1): \${dexplanation1}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel1} the action. Yes or no?																												
<code>\$(duty2): \${dexplanation2}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel2} the action. Yes or no?																												
<code>\$(duty3): \${dexplanation3}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel3} the action. Yes or no?																												
<code>\$(duty4): \${dexplanation4}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel4} the action. Yes or no?																												
<code>\$(duty5): \${dexplanation5}</code>	This output is good (relevant, sensible, and reasonable). Yes or no?	This output \${dlabel5} the action. Yes or no?																												
<input type="checkbox"/> Missing duty(s)?																														

Optional Feedback: Thanks for filling out the questions above! If something about the hit was unclear, please leave a comment in the box below. We would like to make this HIT easier for future workers, so we really appreciate feedback though it is optional.

Figure 10: MTurk Data Quality Annotation

Instructions

In this survey, we will show you an action, and then list a series of moral claims (values, rights or duties) that are considered relevant for the action. Your task is to say whether or not each claim, as stated, aligns with your own subjective values. In other words, do you agree with the statement?

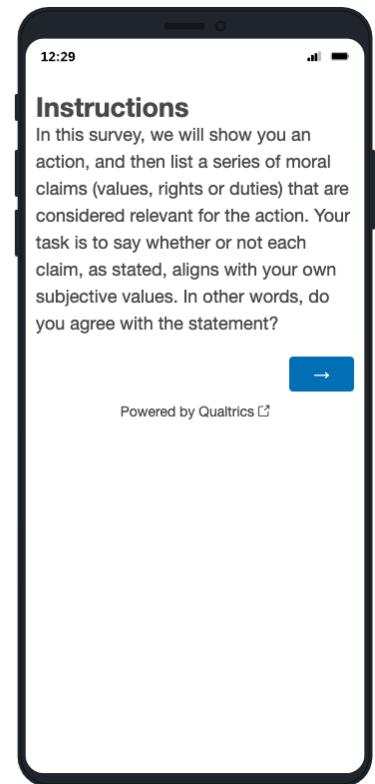


Figure 11: CloudResearch Study - Instructions

Action: Copying brushless drivetrains without credit
Values

Creativity: Copying someone's brushless drivetrain design without credit undermines the value of creativity and innovation, as it does not appreciate the original work.

Do you personally agree with this value claim, as related to the action?

Yes

No

Fairness: Giving credit to the original creator of the brushless drivetrain ensures a fair distribution of recognition and potential profits.

Do you personally agree with this value claim, as related to the action?

Yes

No

Transparency: Copying the design without credit goes against transparency, as it promotes dishonesty about the source of the design.

Do you personally agree with this value claim, as related to the action?

Yes

No

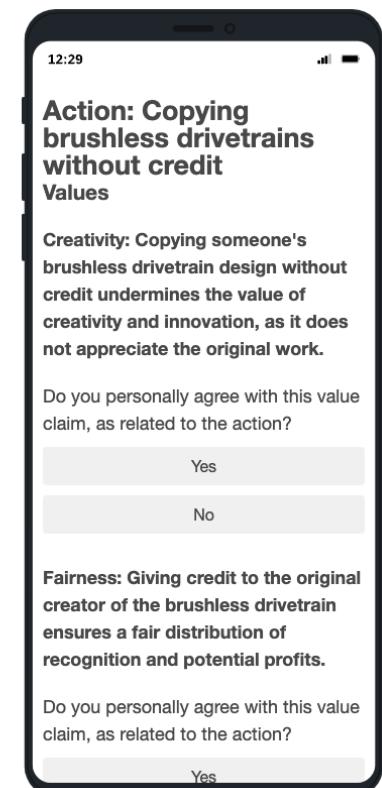


Figure 12: CloudResearch Study - Agreement Example Screen

Action: Copying brushless drivetrains without credit

Values

- Creativity: Copying someone's brushless drivetrain design without credit undermines the value of creativity and innovation, as it does not appreciate the original work.
- Fairness: Giving credit to the original creator of the brushless drivetrain ensures a fair distribution of recognition and potential profits.
- Transparency: Copying the design without credit goes against transparency, as it promotes dishonesty about the source of the design.

Rights

- Intellectual Property Rights: The original creator of the brushless drivetrain has the right to have their work acknowledged and protected from unauthorized copying.

Duties

- Recognizing Original Work: There is an imperfect duty to respect and recognize the work of others by giving proper credit for their ideas and creations.
- Upholding Legal and Moral Norms: By copying the brushless drivetrain without credit, one ignores the perfect duty to respect the legal and moral norms that govern the use and attribution of someone else's work.

Do you have any opinions or perspectives related to the action that haven't been represented in the claims?

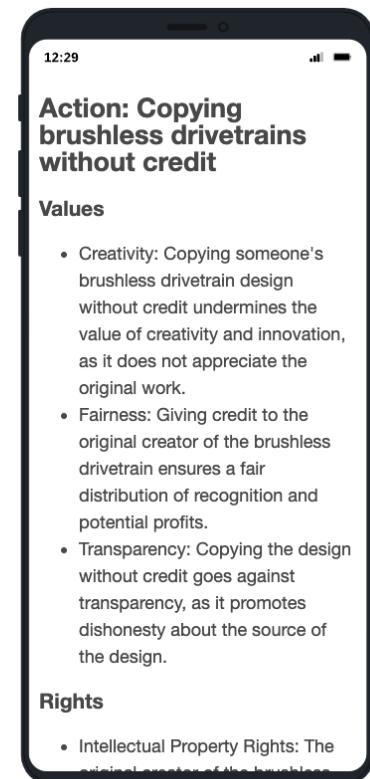


Figure 13: CloudResearch Study - Missing Value or Perspective Example screen

Instructions (click to expand/collapse)

WARNING This HIT may contain **adult content** and may be **offensive** or **upsetting**. **Worker discretion is strongly advised.**

WARNING Current employees of the **University of Washington**, family members of UW employees, and UW students involved in this particular research are **not eligible** to complete this HIT.

Thanks for participating in this HIT!

This task is still **under development**. Please let us know if we can make the task clearer as you complete these HITs. Some HITs may be ambiguous but just answer to the best of your abilities.

Given a **situation** you will be asked to compare the ability of two different AI systems to come up with **values**, **rights**, and **duties** that are relevant to this situation.

- **values:** These are the *intrinsic goods or ideals* that people pursue or cherish, such as happiness, well-being, justice, or freedom. Values are the desirable qualities that people may seek in their lives and in the world. They are often the guiding principles for individuals and societies, shaping goals, motivations, and preferences.
- **rights:** Rights are the *entitlements or claims* that individuals have against others or society, which are usually based on moral or legal grounds. These can be positive rights (e.g., the right to education, healthcare, or free speech) or negative rights (e.g., the right to not be harmed, enslaved, or discriminated against). Rights serve to protect the fundamental interests of individuals and establish certain boundaries that others must respect.
- **duties:** Duties are the *moral obligations or responsibilities* that individuals owe to others or to society at large. They are categorical reasons for doing or refraining from doing something. Independent of whether we want to do or refrain from doing that thing. Duties can be weighty reasons, not easily overridden by competing concerns, and their violation may justify blame and self-blame (guilt). Duties can arise from relationships, social roles, or moral principles, and they guide our actions and decisions.

You will be asked to compare 2 AI systems, **system A** and **system B** across 3 dimensions:

- **Correctness:** Which system provides a more accurate list of responses? For example, if **system A** contains more errors, or responses that are not relevant to the given situation, then you should select **system B** as the winner.
- **Completeness:** Which system provides a more comprehensive list of responses? You should ignore **correctness** when answering this questions, and only consider which system covers more of the important values, rights, and duties that are relevant to the situation.
- **Overall Preference:** Following your intuition as a human, which of the two systems does a better job at producing values, rights, and duties for the given situation? Do not overthink this question, just follow your gut.

Rules

- When possible, avoiding saying that the two systems are **equal** for each dimension, try to pick one or the other.
- Do not overthink your responses. Some examples can be quite challenging, and it is important to follow your instinct once you have read each response carefully.
- This task is still being developed, so please let us know what is confusing and what information would help you complete the task in the future! We include a feedback box at the bottom of the HIT.

Situation: \${situation}

System A:	System B:
\${respa}	\${respb}

Which system is more correct on average?
In other words, which system makes fewer mistakes, or produces a higher fraction of relevant responses?

System A
 About equal (try to avoid this answer unless necessary)
 System B

Which is more complete?
Specifically focus on which response is more comprehensive in capturing the relevant values, rights, and duties for the given situation.

System A
 About equal (try to avoid this answer unless necessary)
 System B

Which response do you prefer overall?
Consider your answers to the previous questions, and follow your intuition for which system produced a better set of responses for the given situation.

System A
 About equal (try to avoid this answer unless necessary)
 System B

(Optional) Please let us know if anything was unclear, if you experienced any issues, or if you have any other feedback for us.

Submit

Figure 14: Batch Value, Right, and Duty comparison against GPT-4.

Full Instructions [\(Expand/Collapse\)](#)

WARNING Current employees of the **University of Washington**, family members of UW employees, and UW students involved in this particular research are **not eligible** to complete this HIT.

CONTENT WARNING This task may contain content that some individuals may find unpleasant, offensive, or disturbing, including sexually suggestive or racially insensitive situations. Worker discretion is advised. Feel free to exit the task at any time.

Thank you for participating in this HIT! It's important to read the instructions thoroughly before proceeding.

Your task is to evaluate moral reasoning about specific actions **based on values, rights, and duties** as defined below:

Definitions (Expand/Collapse)
Values: These are the intrinsic goods or ideals that people pursue or cherish, such as <i>happiness, well-being, justice, or freedom</i> . Values are the desirable qualities that people may seek in their lives and in the world. They are often the guiding principles for individuals and societies, shaping goals, motivations, and preferences.
Rights: Rights are the entitlements or claims that individuals have against others or society, which are usually based on moral or legal grounds. These can be positive rights (e.g., the right to education, healthcare, or free speech) or negative rights (e.g., the right to not be harmed, enslaved, or discriminated against). Rights serve to protect the fundamental interests of individuals and establish certain boundaries that others must respect.
Duties: Duties are the moral obligations or responsibilities that individuals owe to others or to society at large. They are categorical reasons for doing or refraining from doing something, independent of whether we want to do or refrain from doing that thing. Duties can be weighty reasons, not easily overridden by competing concerns, and their violation may justify blame and self-blame (guilt). Duties can arise from relationships, social roles, or moral principles, and they guide our actions and decisions.

You'll evaluate the output (marked in purple) in two steps:

Task 1: Evaluate Relevance
For each action, you will be presented with values, rights, and duties that may make sense to consider when determining the morality of the action (see definitions above). Your job is to determine if these factors are relevant for the action.

Task 2: Evaluate Explanation
Next, you'll see an explanation for how value, right, or duty might connect to the action. Your job is to determine if this explanation is **high-quality** and **effective at making a connection** to the action.

Task 3: Evaluate Support or Opposition
Finally, there will be a guess a guess for each factor on whether it **supports** doing the action, **opposes** doing the action, or could **either support or oppose** the action depending on the interpretation or context. Your job is to determine if this guess is **correct**.

Note: If you're unsure, go with your instincts. If you find yourself not coming across many bad examples, that's okay - it might just be the luck of the draw.

We appreciate your attention to detail and thoughtful responses. Thank you for your contribution!

Action \${situation}			
Values			
Value \${value1}	Relevant The value is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation1} The explanation is high-quality and connects the value to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This value \${vlabel1} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Value \${value2}	Relevant The value is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation2} The explanation is high-quality and connects the value to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This value \${vlabel2} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Value \${value3}	Relevant The value is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation3} The explanation is high-quality and connects the value to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This value \${vlabel3} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Value \${value4}	Relevant The value is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation4} The explanation is high-quality and connects the value to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This value \${vlabel4} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Value \${value5}	Relevant The value is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation5} The explanation is high-quality and connects the value to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This value \${vlabel5} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Value \${value6}	Relevant The value is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation6} The explanation is high-quality and connects the value to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This value \${vlabel6} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Value \${value7}	Relevant The value is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation7} The explanation is high-quality and connects the value to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This value \${vlabel7} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Rights			
Right \${right1}	Relevant The right is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation1} The explanation is high-quality and connects the right to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This right \${rlabel1} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Right \${right2}	Relevant The right is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation2} The explanation is high-quality and connects the right to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This right \${rlabel2} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Right \${right3}	Relevant The right is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation3} The explanation is high-quality and connects the right to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This right \${rlabel3} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Right \${right4}	Relevant The right is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation4} The explanation is high-quality and connects the right to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This right \${rlabel4} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Right \${right5}	Relevant The right is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation5} The explanation is high-quality and connects the right to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This right \${rlabel5} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Duties			
Duty \${duty1}	Relevant The duty is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation1} The explanation is high-quality and connects the duty to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This duty \${dlabel1} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Duty \${duty2}	Relevant The duty is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation2} The explanation is high-quality and connects the duty to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This duty \${dlabel2} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Duty \${duty3}	Relevant The duty is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation3} The explanation is high-quality and connects the duty to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This duty \${dlabel3} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Duty \${duty4}	Relevant The duty is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation4} The explanation is high-quality and connects the duty to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This duty \${dlabel4} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No
Duty \${duty5}	Relevant The duty is relevant for the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Explanation \${explanation5} The explanation is high-quality and connects the duty to the action well. Yes or no? <input type="radio"/> Yes <input type="radio"/> No	Supports/Opposes This duty \${dlabel5} the action. Yes or no? <input type="radio"/> Yes <input type="radio"/> No

Optional Feedback: Thanks for filling out the questions above! If something about the hit was unclear, please leave a comment in the box below. We would like to make this HIT easier for future workers, so we really appreciate feedback though it is optional.

Figure 15: Relevance, explanation, and valence annotation MTurk template.

N Data Sheet

Here we include a Datasheet for Datasets (Gebru et al. 2018) to document the dataset.

N.1 Motivation for Dataset Creation

Why was the dataset created?

VALUEPRISM was created 1) to understand what pluralistic human values, rights, and duties are already present in large language models, and 2) to serve as a resource to support open, value pluralistic modeling (e.g., KALEIDO). It contains human-written situations about which to reason and machine-generated candidate values, rights, duties, along with their valences and post-hoc explanations relating them to the situations.

What (other) tasks could the dataset be used for?

The situations could also be used as a rich, diverse dataset of mostly everyday situations for further decision-making work.

Are there obvious tasks for which it should not be used?

The dataset should only be used for research purposes, and should not be used for real-world decision-making, advice, or commercial applications.

Has the dataset been used for any tasks already?

The dataset has only been used so far to train KALEIDO.

If so, where are the results so others can compare?

Results in body of this paper.

Who funded the creation of the dataset?

Funding for this dataset came from the DARPA ITM program and the Allen Institute for AI (AI2).

If there is an associated grant, provide the grant number.

FA8650-23-C-7316

N.2 Dataset Composition

What are the instances? Are there multiple types of instances?

Situations are plain-text English spans. Each one contains several candidate values, rights, and duties, along with a valence relation (supports, opposes, either) and a free-text explanation. Statistics are found in Table 6.

For seq2seq training, we take this data to make 4 subtask splits: generation of a relevant value, right, or duty from a situation, valence of a value, right, or duty in relation to a situation, an explanation of how a value, right, or duty may connect to a situation, and a set of positive and negative pairs for determining whether a value, right, or duty is relevant for a given action. For relevance, we use “was generated” as a proxy for relevant, and negatively sample values, rights, and duties that were generated for other situations. Statistics can be found in 7.

Are relationships between instances made explicit in the data?

There are no relationships between instances beyond the fact that each situation has several seq2seq tasks, which can be trivially reconstructed.

How many instances of each type are there?

Statistics in Table 7.

What data does each instance consist of? “Raw” data (e.g., unprocessed text or images)? Features/attributes?

Situations are raw free-text, but the rest of the dataset is structured. All values, rights, and duties are free-text connected to a situation, along with a corresponding type (either “Value”, “Right”, or “Duty”); valences are connected to a situation and specific value, right, or duty, and are of types “Supports”, “Opposes”, or “Either” supports or opposes; relevances are connected to a situation and specific value, right, or duty, and are of type “Yes” or “No”; and explanations are free-text associated with a situation and particular value, right, or duty.

Is there a label/target associated with instances? If the instances are related to people, are subpopulations identified (e.g., by age, gender, etc.) and what is their distribution?

There are labels associated with instances for valence and relevance. The instances are not related to people.

Is everything included or does the data rely on external resources? (e.g., websites, tweets, datasets) If external resources, a) are there guarantees that they will exist, and remain constant, over time; b) is there an official archival version. Are there licenses, fees or rights associated with any of the data?

Everything is included and does not rely on external resources. We intend the data for research use only under the [CC-Taylor](#) license

Are there recommended data splits or evaluation measures? (e.g., training, development, testing; accuracy/AUC)

Yes, there are recommended training, validation, and testing splits. We recommend and report accuracy for valence and relevance, and perplexity for generation and explanation.

What experiments were initially run on this dataset? Have a summary of those results and, if available, provide the link to a paper with more information here.

T5-based models were trained on splits of this data and were tested on both the synthetic data (Section D) and were assessed by humans (Sections 5.1 and 3). While the interested reader should defer to the paper for more results, humans found that the distilled models matched the test output quality for valence and explanation, surpassed the test quality for generating sets of values, rights, and duties, and output relevances that correlated with human judgments.

We also run two human studies on the dataset (Sections 4.1 and E). Crowdworkers agree the data is high-quality 91% of the time, and have trouble surfacing values, rights, or duties that are missed, providing suggestions less than 1% of the time. Additionally, in an attempt to understand if the dataset aligns best with any demographic groups, we recruit 613 crowdworkers to mark personal agreement with the

data, and do not find significant takeaways for which groups are represented best in the data.

N.3 Data Collection Process

How was the data collected? (e.g., hardware apparatus/sensor, manual human curation, software program, software interface/API; how were these constructs/measures/methods validated?)

The situations were provided by volunteer users of the Delphi user demo, and the candidate values, rights, duties and their corresponding relations were generated by a large language model, GPT-4.

Who was involved in the data collection process? (e.g., students, crowdworkers) How were they compensated? (e.g., how much were crowdworkers paid?)

Data was collected and by the authors of the paper. The dataset was not collected through crowdworkers, but through demo users and the OpenAI API.

However, to understand the dataset’s quality and representativeness, we do carry out several human studies on subsets of the data (see Section 4.1 and 4.2). We ensured that, for all tasks, crowdworkers were paid a minimum hourly wage of \$15-25 USD.

Over what time-frame was the data collected? Does the collection time-frame match the creation time-frame? The situations were collected from 2021-2023 on the Delphi user demo, and the values, rights, and duties were generated using the OpenAI API from May 2023-July 2023.

How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part of speech tags; model-based guesses for age or language)? If the latter two, were they validated/verified and if so how?

The data is only associated in that the situations came from the demo and the remaining data from the OpenAI API.

Does the dataset contain all possible instances? Or is it, for instance, a sample (not necessarily random) from a larger set of instances? If the dataset is a sample, then what is the population? What was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)? Is the sample representative of the larger set (e.g., geographic coverage)? If not, why not (e.g., to cover a more diverse range of instances)? How does this affect possible uses?

We source our 31k situations about which to reason from a set of 1.3M user-submitted situations, and curate the dataset by filtering out situations that are not actions or unrelated to morality (as labeled in a few-shot manner²⁵ by Flan-T5 (Chung et al. 2022)). We also filter out any questions using keyword matching.

We note that an outsize proportion of the dataset involves toxic, NSFW, or sexually explicit content. In the interest of having a diversity of situations, we label for these attributes¹

using Flan-T5 (Chung et al. 2022). We take 95% of our situations deterministically from those that have less toxic/NSFW/explicit content, and sample the other 5% uniformly from the rest of the data so as to include the entire spectrum of inputs. We find that this succeeds in increasing the diversity of the dataset, as measured by unique n-grams divided by the length of the dataset (dist-2: .23→.36, dist-3: .54→.67).

Is there information missing from the dataset and why? (this does not include intentionally dropped instances; it might include, e.g., redacted text, withheld documents)
Is this data missing because it was unavailable?

No, there is no known data missing from the dataset, although we do not claim or believe that the dataset is necessarily a comprehensive set of representative human values.

Are there any known errors, sources of noise, or redundancies in the data?

No known errors, sources of noise, or redundancies, although we hope future work will help to shed more light on weaknesses.

N.4 Data Preprocessing

What preprocessing/cleaning was done? (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values, etc.)

The main preprocessing was extraction of the features from raw text output from GPT-4 to the semi-structured dataset that we have. We used regex expressions for this extraction.

Was the “raw” data saved in addition to the preprocessed/cleaned data? (e.g., to support unanticipated future uses)

Yes, the raw GPT-4 outputs were saved in addition to the cleaned data.

Is the preprocessing software available?

Yes, all preprocessing software will be available at <https://github.com/tsor13/kaleido>.

Does this dataset collection/processing procedure achieve the motivation for creating the dataset stated in the first section of this datasheet?

It achieves the goal of 1) trying to understand what pluralistic human values, rights, and duties are currently embedded in GPT-4 (although not other LLMs). It achieves the goal of taking a first step to modeling human values, rights, and duties computationally, as manifested by KALEIDO, but we do not claim that it necessarily does so with accuracy and complete representativeness.

N.5 Dataset Distribution

How is the dataset distributed? (e.g., website, API, etc.; does the data have a DOI; is it archived redundantly?)

We plan on distributing the dataset via Huggingface Datasets, but it will be gated for individual-approval and intended for research-use only in an attempt to prevent misuse.

²⁵Few-shot filtering prompts are found in Appendix M.1.

When will the dataset be released/first distributed? (Is there a canonical paper/reference for this dataset?)

We plan on distributing the dataset in September 2023, with this manuscript as a reference.

What license (if any) is it distributed under? Are there any copyrights on the data?

We plan on distributing VALUEPRISM under the ImpACT license (Allen Institute for AI 2023) as a “medium-risk artifact”. Users must agree to all terms and restrictions of the license before accessing or using the dataset.

Are there any fees or access/export restrictions?

No, the dataset is distributed at no cost. However, we do gate access by individual request and access is predicated on acceptance of the license.

N.6 Dataset Maintenance

Who is supporting/hosting/maintaining the dataset? How does one contact the owner/curator/manager of the dataset (e.g. email address, or other contact info)?

The Allen Institute for AI supports the dataset and it will be hosted on Huggingface. Corresponding authors are Taylor Sorensen (tsor13@cs.washington.edu) and Yejin Choi (yejin@cs.washington.edu).

Will the dataset be updated? How often and by whom? How will updates/revisions be documented and communicated (e.g., mailing list, GitHub)? Is there an erratum?

We do not plan on updating the dataset.

If the dataset becomes obsolete how will this be communicated? Is there a repository to link to any/all papers/systems that use this dataset?

We do not expect the dataset to become obsolete as it does not depend on external sources. Users of VALUEPRISM should cite this manuscript.

If others want to extend/augment/build on this dataset, is there a mechanism for them to do so? If so, is there a process for tracking/assessing the quality of those contributions. What is the process for communicating/distributing these contributions to users?

As of now, there is no formal mechanism to extend/augment/build on this dataset, but anyone interested should reach out to the authors.

N.7 Legal & Ethical Considerations

If the dataset relates to people (e.g., their attributes) or was generated by people, were they informed about the data collection? (e.g., datasets that collect writing, photos, interactions, transactions, etc.)

Users of the Delphi user demo explicitly agreed that their queries could be recorded and used for research purposes, and the rest of the data was machine-generated.

If it relates to other ethically protected subjects, have appropriate obligations been met? (e.g., medical data might include information collected from animals)

It does not relate to other ethically protected subjects.

If it relates to people, were there any ethical review applications/reviews/approvals? (e.g. Institutional Review Board applications) If it relates to people, were they told what the dataset would be used for and did they consent? What community norms exist for data collected from human communications? If consent was obtained, how? Were the people provided with any mechanism to revoke their consent in the future or for certain uses?

Data does not relate directly to people.

If it relates to people, could this dataset expose people to harm or legal action? (e.g., financial social or otherwise) What was done to mitigate or reduce the potential for harm?

Data does not relate directly to people.

If it relates to people, does it unfairly advantage or disadvantage a particular social group? In what ways? How was this mitigated?

Data does not relate directly to people.

If it relates to people, were they provided with privacy guarantees? If so, what guarantees and how are these ensured?

Data does not relate directly to people.

Does the dataset comply with the EU General Data Protection Regulation (GDPR)? Does it comply with any other standards, such as the US Equal Employment Opportunity Act?

Especially because the data does not relate to people or have personally identifiable information, it does comply with these laws.

Does the dataset contain information that might be considered sensitive or confidential? (e.g., personally identifying information) Does the dataset contain information that might be considered inappropriate or offensive?

No, the dataset does not contain sensitive or confidential information (like personally identifiable information). The dataset does potentially contain inappropriate or offensive text, especially in the demo-sourced situations, and we advise that the dataset is not for all eyes before providing access. While we did not want to completely remove inappropriate or offensive situations so that the model could perform well in surfacing relevant values, rights, and duties in these cases, we did attempt to ensure that the generated data does not include inappropriate or offensive content via manual inspection and toxicity filters.