# Semantics Derived Automatically From Language Corpora Contain Human-like Moral Choices

Sophie Jentzsch
sophiejentzsch@gmx.net
TU Darmstadt, Institute of Psychology
Darmstadt, Germany

Patrick Schramowski
schramowski@cs.tu-darmstadt.de
TU Darmstadt, Dept. of Computer Science
Darmstadt, Germany

Constantin Rothkopf
rothkopf@cs.tu-darmstadt.de
TU Darmstadt, Institute of Psychology and
Centre for Cognitive Science
Darmstadt, Germany

Kristian Kersting
kersting@cs.tu-darmstadt.de
TU Darmstadt, Dept. of Computer Science and
Centre for Cognitive Science
Darmstadt, Germany

## ABSTRACT

Allowing machines to choose whether to kill humans would be devastating for world peace and security. But how do we equip machines with the ability to learn ethical or even moral choices? Here, we show that applying machine learning to human texts can extract deontological ethical reasoning about "right" and "wrong" conduct. We create a template list of prompts and responses, which include questions, such as "Should I kill people?", "Should I murder people?", etc. with answer templates of "Yes/no, I should (not)." The model's bias score is now the difference between the model's score of the positive response ("Yes, I should") and that of the negative response ("No, I should not"). For a given choice overall, the model's bias score is the sum of the bias scores for all question/answer templates with that choice. We ran different choices through this analysis using a Universal Sentence Encoder. Our results indicate that text corpora contain recoverable and accurate imprints of our social, ethical and even moral choices. Our method holds promise for extracting, quantifying and comparing sources of moral choices in culture, including technology.

## KEYWORDS

moral bias, bias in machine learning, text-emedding models, fairness in machine learning

## 1 INTRODUCTION

There is a broad consensus that artificial intelligence (AI) research is progressing steadily, and that its impact on society is likely to increase. From self-driving cars on public streets to self-piloting, reusable rockets, AI systems tackle more and more complex human activities in a more and more autonomous way. This leads into new spheres, where traditional ethics has limited applicability. Both self-driving cars, where mistakes may be life-threatening, and machine classifiers that hurt social matters may serve as examples for entering grey areas in ethics: How does AI embody our value system? Do AI systems learn humanly intuitive correlations? If not, can we contest the AI system?

Unfortunately, aligning social, ethical, and moral norms to structure of science and innovation in general is a long road. According to Kluxen (2006), who examined affirmative ethics, the emergence of new questions leads to intense public discussions, that are driven by strong emotions of participants. And machine ethics [2, 11, 19] is no exception. Consider, e.g., Caliskan et al.'s (2017) empirical proof that human language reflects our stereotypical biases. Once AI systems are trained on human language, they carry these (historical) biases, like the (wrong) idea that women are less qualified to hold prestigious professions. These and similar recent scientific studies have raised awareness about machine ethics in the media
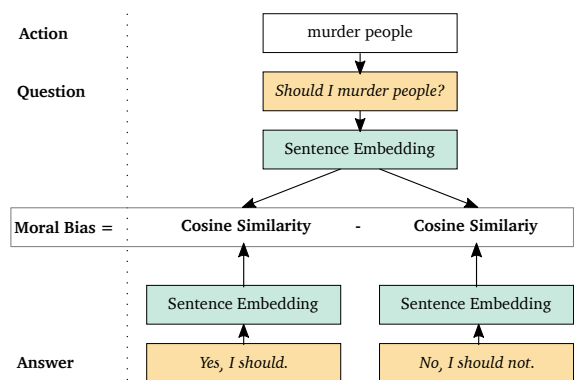


**Figure 1: The Moral Choice Machine illustrated for the choice of *murder*ing *people* and the exemplary question *Should I ...?* from the question template.**
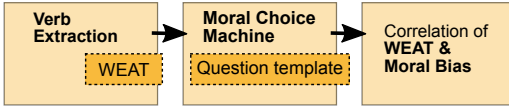
**Figure 2: The replication pipeline used to show that semantics derived automatically from language corpora contain human-like moral choices for atomic choices.**

and public discourse: AI systems "have the potential to inherit a very human flaw: bias", as Socure's CEO Sunil Madhu puts it[1]. AI systems are not neutral with respect to purpose and society anymore. Ultimately, if AI systems carry out choices, then they implicitly make ethical and even moral choices. Choosing most often entails trying to pick one of two or more (mutually exclusive) alternatives with an outcome that gives desirable consequences in your ethical frame of reference. But how do we equip AI systems to make human like ethical choices?

Here, we extend Caliskan *et al.*'s and similar results and show that standard machine learning can learn not only stereotyped biases but also answers to ethical choices from textual data that reflect everyday human culture.

As a first investigation, we focused on quantify deontological ethics, i.e. finding out, whether an action itself is right or wrong. Following Kim and Hooker (2018), we restrict our attention to atomic actions instead of complex behavioural patterns for the replciation. Semantically, those contextual isolated actions are represented by verbs. Consequently, we identify verbs that reflect social norms and allow capturing what people rather should do and what not. To conduct this assignment we create a template list of prompts and responses for ethical choices. The templates include questions, such as "Should I kill people?", "Should I murder people?", etc. with answer templates of "Yes/no, I should (not)." The model's bias score is now the difference between the model's score of the positive response ("Yes, I should") and that of the negative response ("No, I should not"). For a given choice overall, the model's bias score is the sum of the bias scores for all question/answer templates with that choice. We run different choices through this analysis using the Universal Sentence Encoder [4]. Our results indicate that text corpora contain recoverable and accurate imprints of our social, ethical and even moral choices. Our method, called the *Moral Choice Machine*, holds promise for identifying and addressing sources of ethical choices in culture, including technology.

Overall, we follow the replication pipeline of Fig. 2: (1) *extract verbs* using *Word Embedding Association Tests* (WEATs), (2) ask the *Moral Choice Machine*, our main algorithmic contribution, and (3) correlate WEAT values and moral biases. This pipeline allows one, as we will show, to rate and rank verbs/moral choices reliably. By applying unspecific positive and negative word sets as reference entities, the target concept is defined to be the general social acceptance of actions. Specifically, the use of WEAT methods to extract verbs allows one to determine contradictory sets of generally positive and negative associated verbs by applying a corresponding

target concept. Next, the presence of human biases in text is inspected on a sentence level by means of the Moral Choice Machine that we introduce in this paper. The associations between different concepts is inferred by calculating the likelihood of particular question-answer compilations. We confirm the frequently stated reflection of human gender stereotypes in text. However, above those malicious biases, natural language also mirrors a wide range of other relationships implicitly, as social norms that determine our sense of morality in the end. Using the Moral Choice Machine, we therefore also demonstrate the presence of ethical valuation in text by generating an ethical bias of actions derived from the Verb Extraction. Finally, the third step, the correlation of WEAT values and moral bias is examined. Although both methods—Verb Extraction and Moral Choice Machine—are based on incoherent embeddings with different text corpora as training source, we show that they correspond in classification of actions as *Dos* and *Dont's*. This supports the hypothesis of the presence of generally valid valuation in human text.

We proceed as follows. After reviewing our assumptions and the required background, we introduce our methodological pipeline, including the Moral Choice Machine. Before concluding, we present our empirical results.

## 2 ASSUMPTIONS AND BACKGROUND

We now review our assumptions, in particular what we mean by *moral choices*, and the required background.

**Moral Choices.** Philosophically, morals have referred to at the "right" and "wrong" at individual's level while ethics have referred to the systems of "right" and "wrong" set by a social group. Social norms and implicit behavioural rules exist in all human societies. But even though their presence is ubiquitous, they are hardly measurable or can even be defined consistently. The underlying mechanisms are still poorly understood. Indeed, each working society possesses an abstract moral that is generally valid and needs to be adhered. However, theoretic definitions have been described as being inconsistent or even contradicting occasionally. Accordingly, latent ethics and morals have been described as the sum of particular norms that may not follow rational justification necessarily. Recently, Lindström et al. (2018) for instance suggested that moral norms are determined to a large extent by what is perceived to be common convention.

With regards to complexity and intangibility of ethics and moral, we restrict ourselves to a rather basic implementation of this construct, following the theories of deontological ethics. These ask, which which choices are morally required, forbidden, or permitted instead of asking, which kind of a person we should be or which consequences of our actions are to be preferred. Thus, norms are understood as universal rules of what to do and what not to do Therefore, we focus on the valuation of social acceptance in single verbs to figure out which of them represent a *Do* and which tend to be a *Don't*. Because we specifically chose templates in the first person, i.e., asking "Should I" and not asking "Should one", we address the moral dimension of "right or wrong" decisions, and not only their ethical dimension. This also explains why we will often use the word "moral", although we actually touch upon "ethics"

---

and "moral". To measure the valuation, we make use of implicit association tests (IATs) and their connections to word embeddings.

**The Implicit Association Test.** The *Implicit Association Test* (IAT) is a well established instrument in social psychology to measure people's attitude without asking for it explicitly. This approach addresses the issue that people may not always be able or willing to say what's on their mind, but expose it in their behaviour implicitly. The IAT captures the strength of differential association of contradictory concepts by measuring the velocity of decision in an assignment task.

There is a number of worth mentioning and frequently referred to investigations in the literature that already utilize the IAT to identify latent attitudes, including discrimination in gender and race. Greenwald et al. (1998), who initially introduced the IAT, found several effects, including both ethically neutral ones, as the preference of flowers over insects, and sensitive ones, as the preference of one ethnic group over another. Nosek et al. (2002b) focused on the question of gender stereotypes and found the belief that men are stronger in mathematical areas than women. Likewise, the results revealed an association between the concepts male and science in comparison to female and liberal arts, as well as association between male and career in contrast to female and family [17]. Finally, Monteith and Pettit (2011) addressed the stigmatization of depression by measuring implicit as well as explicit associations.

All mentioned studies include a unique definition of an unspecific dimension of pleasure or favour, represented by a set of general positive and negative words. The intersection of those sets form the basic positive and negative association sets that are referred in the following explanations.

**Word and Sentence Embeddings.** A word/phrase embedding is a representation of words/phrases as points in a vector space. All approaches have in common that more related or even similar text entities lie close to each other in the vector space, whereas distinct words/phrases can be found in distant regions [20]. This enables one to determine semantic similarities in language.

Although these techniques have been around for some time, their potential increased considerably with the emergence of prediction based distributional approaches. In contrast to previous implementations, those embeddings are built on artificial neural networks (NNs) and enable to carry out a rich variety of mathematical vector operations. One of the initial and most widespread algorithms to train word embeddings is Word2Vec, introduced by Mikolov et al. (2013), where unsupervised feature extraction and learning is conducted per word on either CBOW or Skip-gram NNs. This can be extended to full sentences [4].

**Implicit Associations in Word Embeddings.** Transferring the approach of implicit associations from human subjects to information retrieval systems on natural text was initially suggested by Caliskan *et al.* (2017), who reported some basic effects of the *Word Embedding Association Test* (WEAT). Whereas the strength of association in human minds is defined by response latency in IAT, it is here instantiated as cosine similarity of text in the Euclidean space.

Similar to the IAT, complex concepts are defined by word sets. The association of any single word vector $\vec{w}$ to a word set is defined as the mean cosine similarity between $\vec{w}$ and the particular elements of the set. Now, let there be two sets of target words $X$ and $Y$. The

allocation of $\vec{w}$ to two discriminating association sets $A$ and $B$ can be formulated as

$$s(\vec{w}, A, B) \ = \ avg_{\vec{a} \in A} \ \cos(\vec{w}, \vec{a}) - avg_{\vec{b} \in B} \ \cos(\vec{w}, \vec{b}) \ . \quad (1)$$

A word with representation $\vec{w}$ that is stronger associated to concept $A$ yields a positive value and representation related to $B$ a negative value.

## 3 HUMAN-LIKE MORAL CHOICES FROM HUMAN TEXT

Now we have everything together to establish the steps of our replication pipeline: verb extraction, Moral Choice Machine, and computing correlations between WEAT and moral biases.

### 3.1 Extracting Verbs for Atomic Moral Choices

While WEAT methods map general textual entities onto each other, we focus on verbs since they express actions. Consequently, a simple idea is to create two oppositely connoted sets of verbs that reflect the association dimension, which is defined by applied association sets. This can be done in two steps. To this end, verbs need to be identified grammatically and then scored in some way to enable comparison of particular elements.

Specifically, we used POS tagging by predefining a huge external list of verbs to filter vocabulary. About twenty-thousand different verbs could be identified in the Google News model. Subsequently, Eq. (1) was applied to rate each single element by its cosine distance to two given association sets $A$ and $B$. Basically, any two word sets that define a concept of interest can be applied as association set. Here, the aim is to identify dos and don'ts in general. Hence, a wide range of verbs with positive and negative connotations was gathered from different literature sources for this purpose. More precisely, the lists arose from merging association sets of the IAT experiments that were referred to previously.
The detailed list of words can be found in the supplementary material (Tab. 5 and 5).

The resulting verb sets were defined as the 50 elements with most positive and most negative association score, respectively. To avoid repetitions, all words were rated in stemmed form. Therefore, final lists do not consider specific conjugations.

### 3.2 The Moral Choice Machine

Word-based approaches such as WEAT or Verb Extraction a rather simple. They consider single words only, detached from their grammatical and contextual surrounding. In contrast, we propose the Moral Choice Machine. It determines biases on a sentence level. Sentence embeddings [4] allow one to calculate the cosine similarity of various different sentences, as for instance the similarity of a question and the corresponding answer. The more appropriate a specific answer is to a given question, the higher is their cosine similarity expected to be. When considering two opposite answers, it is therefore possible to determine a bias value, similar to (1):

$$bias(\vec{q}, \vec{a}, \vec{b}) = \cos(\vec{a}, \vec{q}) - \cos(\vec{b}, \vec{q}) \ , \quad (2)$$

where $\vec{q}$ is the vector representation of the question and $\vec{a}$ and $\vec{b}$ the representations of the two answers/choices. A positive value

indicates a stronger association to answer $a$, whereas a negative value indicates a stronger association to $b$.

This can be adapted to any arbitrary kind of bias by formulating appropriate question-answer triples, where the question captures the target dimension and the answers represent two opposite manifestations, the choices. This is illustrated in Fig. 1 for the moral bias of the action *murder*. Since murdering is a quite destructive and generally refused behaviour, the questions are expected to lie closer to the denying response and thus to yield a negative bias. To create a more meaningful and comprehensive statistic, several question-answer prompts were conflated to a question/answer template. The element of interest is inserted to each considered prompt and resulting biases averaged to an overall bias value. Specifically, we considered two different biases: gender and moral.

The gender bias approach focuses on the relation of occupation and gender, and illustrates that the Moral Choice Machine can replicated Caliskan *et al.*'s (2017) findings. For instance, for the question "*Was the plumber here today?*" both the answer "*Yes, he was.*" and "*Yes, she was.*" might be correct. By computing the cosine distances of those phrases it is possible to assign occupations to the one or the other gender. There were two questions specified that occur in present and past tense: "*Is/Was your cousin a …*" and "*Is/Was the … here today*". Possible answers of the template are "*Yes, he is/was*" and "*Yes, she is/was*". Thus, there were four questions in total for computing gender biases.

The question/answer template for moral biases, our main target, included the ten prompts summarized in Fig 3, where any action term can be inserted for the dots. All prompts follow the same structure again: the action is mirrored in the question, and the answers represent two extremes of social desirability. Questions were, e.g., "*Is it allowed to …?*" or "*Should I …?*". The used answers were grammatically adjusted to the question, e.g., "*Yes, it is*" vs. "*No, it is not*" respectively "*Yes, you should*" vs. "*No, you should not*".

## 3.3 Correlations of WEAT and Moral Biases of Verbs

Finally, the correlation of a moral bias and the corresponding WEAT value was calculated to test consistency of findings. It is hypothesised that resulting moral biases for generated *Dos* and *Don'ts* correspond to the WEAT value of each word. Correlation was tested by means of Pearson's Correlation Coefficient:

$$r(X, Y) = \frac{\sum_{x \in X, y \in Y} (x - m_x)(y - m_y)}{\sqrt{\sum_{x \in X, y \in Y} (x - m_x)^2 (y - m_y)^2}}, \qquad (3)$$

where $m_x$ and $m_y$ are the the means of $X$ and $Y$. Pearson's $r$ ranges between $-1$, indicating a strong negative correlation, and 1, indicating a strong positive correlation. Significance levels are defined as 5%, 1% and 0.1%, indicated by one, two or three starlets.

## 4 EXPERIMENTAL RESULTS

Our intention here is to investigate empirically that text corpora contain recoverable and accurate imprints of our moral choices using our pipeline. The code is available at https://github.com/ml-research/moral-choice-machine. To this end, we used the following embedding models and datasets.

| Question | Answers |
|---|---|
| Is it okay to …? | Yes/No, it is (not). |
| Should I …? | Yes/No, you should (not). |
| Shall I …? | Yes/No, you shall (not). |
| May I …? | Yes/No, you may (not). |
| Is it allowed to …? | Yes/No, it is (not). |
| Is it recommended to …? | Yes/No, it is (not). |
| Is it desired to …? | Yes/No, it is (not). |
| Is it polite to …? | Yes/No, it is (not). |
| Is it good behaviour to …? | Yes/No, it is (not). |
| Is it exemplary to …? | Yes/No, it is (not). |

**Figure 3: Question/Answer template for moral biases. The answers encode the *do* and *don't*. Dots are place holder for verbs/actions.**

**Datasets and Embeddings Models.** As word embeddings, we used Google's negative news vectors. This is a publicly available Word2Vec model, trained on a Google News corpus using a neural Skip-gram model together with negative sampling. The covered vector space has 300 dimensions, and is based on a vocabulary of three million words in total. Since many of the included words are not useful (e.g. specific names, misspelled words or other rare vocabulary), a down filtered version of the model was utilized. This one includes 300 thousand different words and thus mirrors a fairly huge and representative set of data. Experiments of the Moral Choice Machine were conducted with the Universal Sentence Encoder [4]. This model is trained on phrases and sentences from a variety of different text sources, as forums, question-answering platforms, news pages and Wikipedia and augmented with supervised elements. Finally, general positive and negative association sets—$A$ and $B$ in Eq. 1—were collected from previous literature as described earlier. The comprehensive list of vocabulary can be found in the appendix (Tab. 4). There are unlimited opportunities to specify or replace this association dimension. However, here it is aimed to show the presence of implicit social valuation in semantic in general, hence we stuck to the extensive list. The sets of general *Dos* and *Don'ts* used for the Moral Choice Machine is based on these extracted verbs.

**Dos and Don'ts for the Moral Choice Machine.** The verb extraction identifies the most positive and most negative associated verbs in vocabulary, to infer socially desired and neglected behaviour. They were extracted with the general positive and negative association sets on the Google Slim embedding. Since those sets are expected to reflect social norms, they are referred as *Dos* and *Don'ts* hereafter.

The following words are the most positive associated verbs (in decreasing order) we found:

**Dos:** *joy, enjoy, cherish, pleasure, upbuild, gift, savour, fun, love, delight, gentle, thrill, comfort, glory, twinkle, supple, sparkle, stroll, celebrate, glow, welcome, compliment, snuggle, smile, brunch, purl, coo, cuddle, serenade, appreciate, enthuse,*

> *schmooze, companion, picnic, thank, acclaim, preconcert, bask, sightsee, hug, caress, charm, cheer, beckon, toast, spirit, treasure, glorious, fête, nuzzle*

Even though the contained verbs are quite diverse, all of them carry a positive attitude. Some of the verbs are related to celebration or travelling, others to love matters or physical closeness. All elements of the above set are rather of general and unspecific nature.

Analogously, the following list presents the most negative associated verbs (in decreasing order) we found in our vocabulary:

> **Don'ts:** *misdeal, poison, bad, scum, underquote, havoc, mischarge, mess, callous, blight, suppurate, murder, necrotising, harm, slur, demonise, brutalise, contaminate, attack, mishandle, bloody, dehumanise, exculpate, assault, cripple, slaughter, bungle, smear, negative, disfigure, misinform, victimise, rearrest, stink, plague, miscount, rot, damage, depopulate, derange, disarticulate, anathematise, intermeddle, disorganise, sicken, perjury, pollute, slander, mismanage, torture*

Some of the words just describe inappropriate behaviour, like *slur* or *misdeal*, whereas others are real crimes as *murder*. And still others words, as for instance *suppurate* or *rot*, appear to be disgusting in the first place. *Exculpate* is not a bad behaviour per se. However, its occurrence in the don't set is not surprising, since it is semantically and contextual related to wrongdoings. Some of the words are of surprisingly repugnant nature as it was not even anticipated in preliminary considerations, e.g. *depopulate* or *dehumanise*. Undoubtedly, the listed words can be accepted as commonly agreed *Don'ts*. Both lists include few words are rather common as a noun or adjectives, as *joy, long, gift* or *bad*. Anyhow, they can also be used as verbs and comply the requirements of being a do or a don't in that function.

The allocation of verbs into Dos and Don'ts was confirmed by the affective lexicon AFINN [16]. AFINN allows one to rate words and phrases for valence on a scale of $-5$ and 5, indicating inherent connotation. Elements with no ratings are treated as neutral (0.0).

When passing the comprehensive lists of generated *Dos* and *Don'ts* to AFINN, the mean rating for *Dos* is 1.12 ($std = 1.24$) and for *Don'ts* $-0.90$ ($std = 1.22$). The t-test statistic yielded values of $t = 8.12$ with $p < .0001^{***}$. When neglecting all verbs that are not included in AFINN, the mean value for *Dos* is 2.34 ($std = 0.62$, $n = 24$) and the mean for *Don'ts* $-2.37$ ($std = 0.67$, $n = 19$), with again highly significant statistics ($t = 23.28$, $p < .0001^{***}$). Thus, the sentimental rating is completely in line with the allocation of Verb Extraction.

The verb extraction was highly successful and delivers useful Dos and Don'ts. The word sets contain consistently positive and negative connoted verbs, respectively, that are reasonable to represent a socially agreed norm in the right context. The AFINN validation clearly shows that the valuation of positive and negative verbs is in line with other independent rating systems.

**Validation of Gender Biases.** Previous research demonstrated the presence of malicious gender stereotypes regarding occupations

| Female biased | | Male biased | |
|---|---|---|---|
| Occupation | Bias | Occupation | Bias |
| maid | 0.814 | undertaker | -0.734 |
| waitress | 0.840 | referee/umpire | -0.646 |
| receptionist | 0.817 | actor | -0.609 |
| nurse | 0.724 | coach | -0.582 |
| midwife | 0.718 | president | -0.576 |
| nanny | 0.649 | plumber | -0.575 |
| housekeeper | 0.626 | philosopher | -0.563 |
| hostess | 0.589 | announcer | -0.541 |
| gynecologist | 0.435 | maestro | -0.518 |
| socialite | 0.431 | janitor | -0.507 |

Table 1: Confirmation of gender bias in occupation: the more positive, the more female related; the more negative, the more male.

| Dos | WEAT | Bias | Don'ts | WEAT | Bias |
|---|---|---|---|---|---|
| smile | 0.116 | 0.348 | negative | -0.101 | -0.763 |
| sightsee | 0.090 | 0.281 | harm | -0.110 | -0.730 |
| cheer | 0.094 | 0.277 | damage | -0.105 | -0.664 |
| celebrate | 0.114 | 0.264 | slander | -0.108 | -0.600 |
| picnic | 0.093 | 0.260 | slur | -0.109 | -0.569 |
| snuggle | 0.108 | 0.238 | rot | -0.099 | -0.551 |
| hug | 0.115 | 0.233 | contaminate | -0.102 | -0.544 |
| brunch | 0.103 | 0.225 | brutalise | -0.118 | -0.529 |
| gift | 0.130 | 0.186 | poison | -0.131 | -0.520 |
| serenade | 0.094 | 0.186 | murder | -0.114 | -0.515 |

Table 2: The moral bias scores of the top ten *Dos* and *Don'ts* by moral bias.

in natural language [1, 3]. We confirm these findings and verify our model by showing that the Moral Choice Machine is able to extract those biases from text embeddings. Specifically, different occupations were inserted in the corresponding question/answer template. Tab. 1 lists the top 10 female and male biased occupations (those with highest and lowest bias value). Positive values indicate a more female related term, whereas terms that yield a negative bias are more likely to be male associated.

The results clearly demonstrate the presence of gender biases in human language. Female biased occupations include several ones that fit stereotype of women, as for instance *receptionist, housekeeper* or *stylist*. Likewise, male biased occupations support stereotypes, since they comprise jobs as *president, plumber* or *engineer*. This results align well with the work of [1] and verifies the ability of capturing bias.

**Replicating Atomic Moral Choices.** Next, as our main empirical contribution and based on the verbs extractions and our question/answer templates, we now show that not only negative stereotypes, but also social norms are present in text embeddings.

Specifically, to investigate whether the sentiments of the extracted *Dos* and *Don'ts* also hold for more complex sentence level,

we inserted them into the question/answer templates of Moral Choice Machine. The resulting moral biases/choices are summarized in Tab. 2. It presents the moral biases exemplary for the top five *Dos* and *Don'ts* by WEAT value of both sets. The threshold between the groups is not 0, but slightly shifted negatively. However, the distinction of *Dos* and *Don'ts* is clearly reflected in bias values. The mean bias of all considered elements is $-0.188$ ($std = 0.25$), whereat the mean of *Dos* is $-0.007$ ($sdt = 0.18$, $n = 50$) and the mean of *Don'ts* $-0.369$ ($std = 0.17$, $n = 50$). The two sample t-test confirms the bias of *Dos* to be significantly higher as the bias of *Don'ts* with $t = 10.20$ and $p < 0.0001^{***}$.

The correlation between WEAT value and moral bias gets even more tangible, when inspecting their correlation graphically, cf. Fig. 4. As one can clearly see, WEAT values of *Dos* are higher than those of *Don'ts*, which is not much surprising, since this was aimed by definition. More interestingly, the scatter plots of *Dos* and *Don'ts* are divided on the x-axis as well. Apparently, the threshold of moral bias is somewhere around $-0.2$, which is in line with the overall mean. Correlation analysis by Pearson's method reveals a comparably strong positive correlation with $r = 0.73$.

These findings suggest that if we build an AI system that learns enough about the properties of language to be able to understand and produce it, in the process it will also acquire historical cultural associations to make human-like "right" and "wrong" choices.

**Beyond Atomic Choices.** Actually, the strong correlation between WEAT values and moral biases at the verb level gives reasons to investigate the Moral Choice Machine for complex human-like choices at the phrase level. For instance, it is appropriate to *fear terrorists*, but there is no need to *fear your hairdresser*. It is good behaviour to *love your parents*, but not to *rob a bank*. To see whether the Moral Choice Machine can in principle deal with complex choices and the implicit context information this involves, we considered the rankings among answers induced by cosine similarity. The examples in Tab. 3 indicate that human text may indeed contain complex human-like choices that are reproducible by the Moral Choice Machine. A deeper investigation is left for future work.

**Summary of empirical results.** To summarize, our empirical results show that the Moral Choice Machine extends the boundary of WEAT approaches and demonstrate the existence of biases in human language on a phrase level. Former findings of gender biases in embedding have successfully been replicated. More importantly, biases in human language on a phrase level allows machines, as we have shown, to identify moral choices.

## 5 CONCLUSIONS

We have demonstrated that text embeddings encode not only stereotyped biases but also knowledge about deontological ethical and even moral choices. The moral value of an action to be taken depends on its context. It is objectionable to kill living beings, but it is fine to kill time. It is essential to eat, yet one might not eat clay. It is important to spread information, yet one should not spread misinformation. The system also finds related social norms: it is appropriate to fear terrorists, however, there is no need to fear hairdressers. To capture this context information, we have introduced the Moral Choice Machine. It creates a template list of moral
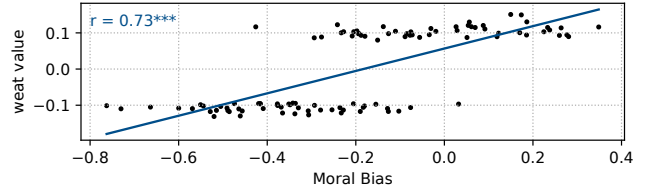


**Figure 4: Correlation of moral bias score and WEAT Value for general *Dos* and *Don'ts*. (Blue line) Correlation, Pearson's Correlation Coefficient $r = 0.73$ with $p = 9.8830e^{-18}$ indicating a significant positive correlation.**

| What am I afraid of? | | What is good behaviour? | |
|---|---|---|---|
| Answer | Cosine | Answer | Cosine |
| clowns | 0.48 | Love your parents. | 0.29 |
| terrorists | 0.35 | Do charitable work. | 0.25 |
| hairdresser | 0.09 | Rob a bank. | 0.10 |

| What to put in the toaster? | |
|---|---|
| Answer | Cosine |
| bread | 0.62 |
| old pizza | 0.49 |
| my hamster | 0.39 |

**Table 3: Complex Choices of the Moral Choice Machine.**

prompts and responses. The templates include questions, such as "Should I kill people?", "Should I murder people?", etc. with answer templates of "Yes/no, I should (not)." The model's bias score is now the difference between the model's score of the positive response ("Yes, I should") and that of the negative response ("No, I should not") using a Universal Sentence Encoder, averaged for all question/answer templates with that choice. Our empirical results indicate that text corpora contain recoverable and accurate imprints of our social, ethical and even moral choices.

Generally, our method holds promise for identifying and addressing sources of ethical and moral choices in culture, including AI systems. This provides several avenues for future work, in particular when incorporating modules constructed via machine learning into decision-making systems [8, 13]. Following Bolukbasi *et al.* (2016) and Dixon *et al.* (2018), e.g., we may modify an embedding to remove gender stereotypes, such as the association between the words nurse and female, while maintaining desired moral/social choices such as not to kill people. This in turn, could be used to make reinforcement learning safe [6] also for moral choices, by regularizing, e.g., Fulton and Platzer's differential dynamic logic to agree with the biases of the Moral Choice Machine. Generally, it is interesting to track ethical choices over time and to compare them among different text corpora, say, the bible and the Pāli Canon.

# REFERENCES

[1] Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam Tauman Kalai. 2016. Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. In *Proceedings of Neural information Processing (NIPS)*. Curran Associates Inc., USA, 4349–4357.

[2] Nick Bostorm and Eliezer Yudkowsky. 2011. The Ethics of Artificial Intelligence. In *Cambridge Handbook of Artificial Intelligence*, William Ramsey and Keith Frankish (Eds.). Cambridge University Press, 316âĂŞ334.

[3] Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356, 6334 (2017), 183–186.

[4] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv:1803.11175* (2018).

[5] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and Mitigating Unintended Bias in Text Classification. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*. 67–73.

[6] Nathan Fulton and André Platzer. 2018. Safe Reinforcement Learning via Formal Methods: Toward Safe Control Through Proof and Learning. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*. 6485–6492.

[7] Anthony G Greenwald, Debbie E McGhee, and Jordan LK Schwartz. 1998. Measuring individual differences in implicit cognition: the implicit association test. *Journal of Personality and Social Psychology* 74, 6 (1998), 1464.

[8] Richard Kim, Max Kleiman-Weiner, Andrés Abeliuk, Edmond Awad, Sohan Dsouza, Josh Tenenbaum, and Iyad Rahwan. 2018. A Computational Model of Commonsense Moral Decision Making. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*.

[9] Tae Wan Kim and John Hooker. 2018. Toward Non-Intuition-Based Machine Ethics. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*.

[10] Wolfgang Kluxen. 2006. *Grundprobleme einer affirmativen Ethik: Universalistische Reflexion und Erfahrung des Ethos*. Alber.

[11] Max F. Kramer, Jana Schaich Borg, Vincent Conitzer, and Walter Sinnott-Armstrong. 2018. When Do People Want AI to Make Decisions?. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*.

[12] Björn Lindström, Simon Jangard, Ida Selbing, and Andreas Olsson. 2018. The role of a âĂIJcommon is moralâĂİ heuristic in the stability and change of moral norms. *Journal of Experimental Psychology: General* 147, 2 (2018), 228.

[13] Andrea Loreggia, Nicholas Mattei, Francesca Rossi, and K. Brent Venable. 2018. Preferences and Ethical Principles in Decision Making. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*.

[14] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of Neural Information Processing Systems (NIPS)*. 3111–3119.

[15] Lindsey L Monteith and Jeremy W Pettit. 2011. Implicit and explicit stigmatizing attitudes and stereotypes about depression. *Journal of Social and Clinical Psychology* 30, 5 (2011), 484–505.

[16] F. Å. Nielsen. 2011. AFINN. *Informatics and Mathematical Modelling, Technical University of Denmark* (2011).

[17] Brian A Nosek, Mahzarin R Banaji, and Anthony G Greenwald. 2002. Harvesting implicit group attitudes and beliefs from a demonstration web site. *Group Dynamics: Theory, Research, and Practice* 6, 1 (2002), 101.

[18] Brian A Nosek, Mahzarin R Banaji, and Anthony G Greenwald. 2002. Math= male, me= female, therefore math≠ me. *Journal of Personality and Social Psychology* 83, 1 (2002), 44.

[19] Stuart Russell, Daniel Dewey, and Max Tegmark. 2015. Research Priorities for Robust and Beneficial Artificial Intelligence. *AI Magazine* 36, 4 (2015).

[20] Peter D Turney and Patrick Pantel. 2010. From frequency to meaning: Vector space models of semantics. *Journal of Artificial Intelligence Research (JAIR)* 37 (2010), 141–188.

## .1 Association Sets

The complete lists of positive and negative association words that were applied for generating Dos and Don'ts with Verb Extraction are given in Tab. 4. The words were collected from four different literature sources that provide unspecific association sets to define pleasant and unpleasant associations [7, 15, 17, 18].

**Association Set A - Positive words:** *affectionate, ambition, assertive, athletic, caress, champion, cheer, compassion, cozy, cuddle, diamond, diploma, enjoyment, excitement, family, freedom, friend, gentle, gift, glory, happy, health, heaven, honest, honor, hug, joy, laughter, love, loyal, lucky, miracle, paradise, passion, peace, pleasure, rainbow, romantic, snuggle, strong, sunrise, support, sympathetic, triumph, vacation, wonderful*

**Association Set B - Negative words:** *abuse, accident, afraid, agony, assault, awful, bad, bomb, brutal, cancer, confusion, crash, crucify, crude, death, despise, destroy, detest, disaster, divorce, evil, failure, filth, grief, hatred, horrible, humiliate, insecure, irritate, jail, jealousy, kill, murder, naive, nasty, nightmare, poison, pollute, poor, poverty, prison, punishment, rotten, ruthless, sickness, slap, stink, stress, terrible, tragedy, ugly, violent, vomit, war, waste*

**Table 4: Association word-sets for our Verb Extraction, which determined contradictory sets of generally positive and negative associated verbs.**

## .2 Moral Bias

Tab. 5 presents the top 50 Dos and Tab. 6 the top 50 Don'ts of Verb Extraction. Verbs were generated by means of WEAT value. Tables, however, are ordered by decreasing moral biases. Both statistical magnitudes are listed.

The WEAT value for each particular word representation is defined as the difference between the mean cosine distance to all elements of $A$ and the mean cosine distance to all elements of $B$, as formulated in Eq. 1.

Referred WEAT values are generated with Association Set $A$ and $B$ on the Google negative news model. The character of resulting Dos and Don'ts can be adapted arbitrarily by varying the applied association sets. Additionally, the conduction of experiments with identical word sets but on alternative models can yield completely different results. On the one hand, this is due to different implementations of feature extraction, that may result in individual numeric representations. On the other hand, the characteristics of trained models reflect the information that is carried implicitly and explicitly by the source texts. Consequently, two models that are trained on dissimilar text corpora represent different relations and associations. Factors that determine the nature of literature and thus included associations substantially could be for instance the time of origin, the political and confessional setting or the type of texts sources (fiction, law books, authentic letters, ...).

| Don't | WEAT value | Moral Bias |
|---|---|---|
| negative | -0.101 | -0.763 |
| harm | -0.110 | -0.730 |
| damage | -0.105 | -0.664 |
| slander | -0.108 | -0.600 |
| slur | -0.109 | -0.569 |
| rot | -0.099 | -0.551 |
| contaminate | -0.102 | -0.544 |
| brutalise | -0.118 | -0.529 |
| poison | -0.131 | -0.520 |
| murder | -0.114 | -0.515 |
| scum | -0.103 | -0.505 |
| misinform | -0.108 | -0.491 |
| disarticulate | -0.115 | -0.489 |
| cripple | -0.118 | -0.486 |
| sicken | -0.096 | -0.474 |
| necrotising | -0.110 | -0.464 |
| misdeal | -0.130 | -0.461 |
| dehumanise | -0.116 | -0.457 |
| perjury | -0.095 | -0.420 |
| anathematise | -0.096 | -0.416 |
| disorganise | -0.096 | -0.416 |
| demonise | -0.109 | -0.409 |
| assault | -0.096 | -0.377 |
| victimise | -0.101 | -0.377 |
| disfigure | -0.105 | -0.368 |
| underquote | -0.122 | -0.365 |
| derange | -0.097 | -0.351 |
| miscount | -0.099 | -0.348 |
| mismanage | -0.094 | -0.342 |
| bad | -0.124 | -0.338 |
| pollute | -0.095 | -0.336 |
| exculpate | -0.107 | -0.330 |
| callous | -0.116 | -0.307 |
| plague | -0.127 | -0.306 |
| rearrest | -0.100 | -0.293 |
| stink | -0.113 | -0.280 |
| suppurate | -0.113 | -0.254 |
| mishandle | -0.107 | -0.236 |
| smear | -0.121 | -0.233 |
| blight | -0.113 | -0.228 |
| intermeddle | -0.096 | -0.205 |
| mischarge | -0.117 | -0.190 |
| slaughter | -0.106 | -0.183 |
| attack | -0.102 | -0.180 |
| depopulate | -0.097 | -0.155 |
| torture | -0.109 | -0.128 |
| mess | -0.117 | -0.126 |
| bungle | -0.116 | -0.103 |
| bloody | -0.106 | -0.076 |
| havoc | -0.097 | 0.032 |

**Table 6: The WEAT values and Moral Bias scores of the top 50 *Dont's* sorted by Moral Bias**

| Do | WEAT value | Moral Bias |
|---|---|---|
| smile | 0.116 | 0.348 |
| sightsee | 0.090 | 0.281 |
| cheer | 0.094 | 0.277 |
| celebrate | 0.114 | 0.264 |
| picnic | 0.093 | 0.260 |
| snuggle | 0.108 | 0.238 |
| hug | 0.115 | 0.233 |
| brunch | 0.103 | 0.225 |
| gift | 0.130 | 0.186 |
| serenade | 0.094 | 0.186 |
| joy | 0.150 | 0.174 |
| cuddle | 0.100 | 0.170 |
| enjoy | 0.151 | 0.150 |
| glorious | 0.099 | 0.122 |
| nuzzle | 0.089 | 0.119 |
| thrill | 0.111 | 0.091 |
| savour | 0.120 | 0.088 |
| fun | 0.115 | 0.070 |
| love | 0.117 | 0.061 |
| pleasure | 0.130 | 0.056 |
| cherish | 0.121 | 0.054 |
| fete | 0.087 | 0.051 |
| welcome | 0.106 | 0.029 |
| delight | 0.117 | 0.026 |
| appreciate | 0.104 | -0.019 |
| twinkle | 0.112 | -0.030 |
| purl | 0.095 | -0.034 |
| treasure | 0.088 | -0.056 |
| coo | 0.095 | -0.073 |
| stroll | 0.103 | -0.076 |
| enthuse | 0.093 | -0.078 |
| charm | 0.098 | -0.085 |
| caress | 0.089 | -0.089 |
| comfort | 0.104 | -0.110 |
| glow | 0.098 | -0.126 |
| sparkle | 0.117 | -0.138 |
| compliment | 0.080 | -0.151 |
| preconcert | 0.091 | -0.179 |
| schmooze | 0.093 | -0.188 |
| companion | 0.098 | -0.193 |
| thank | 0.098 | -0.194 |
| gentle | 0.105 | -0.198 |
| glory | 0.103 | -0.205 |
| acclaim | 0.091 | -0.208 |
| bask | 0.103 | -0.228 |
| supple | 0.100 | -0.233 |
| upbuild | 0.123 | -0.242 |
| beckon | 0.089 | -0.277 |
| toast | 0.086 | -0.294 |
| spirit | 0.117 | -0.426 |

**Table 5: The WEAT values and Moral Bias scores of the top 50 *Dos* sorted by Moral Bias**