# A framework for the computational linguistic analysis of dehumanization

**Julia Mendelsohn** [1,*]**, Yulia Tsvetkov** [2] **and Dan Jurafsky** [3]

[1]*School of Information, University of Michigan, Ann Arbor , MI , USA*
[2]*Language Technologies Institute, Carnegie Mellon University, Pittsburgh, PA, USA*
[3]*Department of Linguistics, Stanford University, Stanford, CA, USA*

Correspondence*:
Julia Mendelsohn
juliame@umich.edu

## ABSTRACT

Dehumanization is a pernicious psychological process that often leads to extreme intergroup bias, hate speech, and violence aimed at targeted social groups. Despite these serious consequences and the wealth of available data, dehumanization has not yet been computationally studied on a large scale. Drawing upon social psychology research, we create a computational linguistic framework for analyzing dehumanizing language by identifying linguistic correlates of salient components of dehumanization. We then apply this framework to analyze discussions of LGBTQ people in the *New York Times* from 1986 to 2015. Overall, we find increasingly humanizing descriptions of LGBTQ people over time. However, we find that the label *homosexual* has emerged to be much more strongly associated with dehumanizing attitudes than other labels, such as *gay*. Our proposed techniques highlight processes of linguistic variation and change in discourses surrounding marginalized groups. Furthermore, the ability to analyze dehumanizing language at a large scale has implications for automatically detecting and understanding media bias as well as abusive language online.

Keywords: computational sociolinguistics, dehumanization, lexical variation, language change, media, New York Times, LGBTQ

Trigger Warning: this paper contains offensive material that some may find upsetting, especially in Table 4 and Table 7.

## 1 INTRODUCTION

Despite the American public's increasing acceptance of LGBTQ people and recent legal successes, LGBTQ individuals remain targets of hate and violence (Dinakar et al., 2012; Silva et al., 2016; Gallup, 2019). At the core of this issue is dehumanization, "the act of perceiving or treating people as less than human" (Haslam and Stratemeyer, 2016), a process that heavily contributes to extreme intergroup bias (Haslam, 2006). Language is central to studying this phenomenon; like other forms of bias (Wiebe et al., 2004; Greene and Resnik, 2009; Recasens et al., 2013; Voigt et al., 2017; Breitfeller et al., 2019), dehumanizing attitudes are expressed through subtle linguistic manipulations, even in carefully-edited texts. It is crucial to understand the use of such linguistic signals in mainstream media, as the media's representation of marginalized social groups has far-reaching implications for social acceptance, policy, and safety.

While small-scale studies of dehumanization and media representation of marginalized communities provide valuable insights (e.g. Esses et al. (2013)), there exist no known large-scale analyses, likely due to difficulties in quantifying such a subjective and multidimensional psychological process. However, the ability to do large-scale analysis is crucial for understanding how dehumanizing attitudes have evolved over long periods of time. Furthermore, by being able to account for a greater amount of media discourse, large-scale techniques can provide a more complete view of the media environment to which the public is exposed.

Linguistics and computer science offer valuable methods and insights on which large-scale techniques might be developed for the study of dehumanization. By leveraging more information about the contexts in which marginalized groups are discussed, computational linguistic methods enable large-scale study of a complex psychological phenomenon and can even reveal linguistic variations and changes not easily identifiable through qualitative analysis alone.

We develop a computational linguistic framework for analyzing dehumanizing language, with a focus on lexical signals of dehumanization. Social psychologists have identified numerous components of dehumanization, such as *negative evaluations of a target group*, *denial of agency*, *moral disgust*, and *likening members of a target group to non-human entities such as vermin*. Drawing upon this rich body of literature, we first identify linguistic analogs for these components and propose several computational techniques to measure these linguistic correlates. We then apply this general framework to explore changing representations of LGBTQ groups in the *New York Times* over thirty years and both qualitatively and quantitatively evaluate our techniques within this case study. We additionally use this lens of dehumanization to investigate differences in social meaning between the denotationally-similar labels *gay* and *homosexual*. We focus on a single case study in order to conduct an in-depth analysis of our methodology, but our framework can generalize to study representations of other social groups, which we briefly explore in our discussion.

This paper aims to bridge the gaps between computational modeling, sociolinguistics, and dehumanization research with implications for several disciplines. In addition to enabling large-scale studies of dehumanizing language and media representation of marginalized social groups, these techniques can be built into systems that seek to capture both conscious and unconscious biases in text. Furthermore, this work has implications for improving machines' abilities to automatically detect hate speech and abusive language online, which are typically underpinned by dehumanizing language. Finally, our case study demonstrates that such computational analyses of discussions about marginalized groups can provide unique insights into language variation and change within sensitive sociopolitical contexts, and help us understand how people (and institutions) use language to express their ideologies and attitudes towards certain social groups.

## 2 BACKGROUND

### 2.1 Dehumanization

Our lexical semantic analysis involves quantifying linguistic correlates of component psychological processes that contribute to dehumanization. Our approaches are informed by social psychology research on dehumanization, which is briefly summarized here. Prior work has identified numerous related processes that comprise dehumanization (Haslam, 2006). One such component is *likening members of the target group to non-human entities*, such as machines or animals (Haslam, 2006; Goff et al., 2008; Kteily et al., 2015). By perceiving members of a target group to be non-human, they are "outside the boundary in which moral values, rules, and considerations of fairness apply" (Opotow, 1990), which thus leads to violence and other forms of abuse. Metaphors and imagery relating target groups to vermin are particularly insidious and played a prominent role in the genocide of Jews in Nazi Germany and Tutsis in Rwanda (Harris and Fiske, 2015). More recently, the vermin metaphor has been invoked by the media to discuss terrorists and political leaders of majority-Muslim countries after September 11 (Steuter and Wills, 2010). According to Tipler and Ruscher (2014), the vermin metaphor is particularly powerful because it conceptualizes the target group as "engaged in threatening behavior, but devoid of thought or emotional desire".

*Disgust* underlies the dehumanizing nature of these metaphors and is itself another important element of dehumanization. Disgust contributes to members of target groups being perceived as less-than-human and of negative social value (Sherman and Haidt, 2011). It is often evoked (both in real life and experimental settings) through likening a target group to animals. Buckels and Trapnell (2013) find that priming participants to feel disgust facilitates "moral exclusion of out-groups". Experiments by Sherman and Haidt (2011) and Hodson and Costello (2007) similarly find that disgust is a predictor of dehumanizing perceptions of a target group. Both moral disgust towards a particular social group and the invocation of non-human metaphors are facilitated by *essentialist* beliefs about groups, which Haslam (2006) presents as a necessary component of dehumanization. In order to distinguish between human and non-human, dehumanization requires an exaggerated perception of intergroup differences. Essentialist thinking thus contributes to dehumanization by leading to the perception of social groups as categorically distinct, which in turn emphasizes intergroup differences (Haslam, 2006)

According to Haslam (2006), prior work describes *"extremely negative evaluations of others"* as a major component of dehumanization. This is especially pronounced in Bar-Tal's account of delegitimization, which involves using negative characteristics to categorize groups that are "excluded from the realm of

acceptable norms and values" (1990). While Bar-Tal defines delegitimization as a distinct process, he considers dehumanization to be one means of delegitimization. Opotow (1990) also discusses broader processes of moral exclusion, one of which is dehumanization. A closely related process is *psychological distancing*, in which one perceives others to be objects or nonexistent (Opotow, 1990). Nussbaum (1999) identifies elements that contribute to the objectification (and thus dehumanization) of women, one of which is *denial of subjectivity*, or the habitual neglect of one's experiences, emotions, and feelings.

Another component of dehumanization is the *denial of agency* to members of the target group (Haslam, 2006). According to Tipler and Ruscher, there are three types of agency: "the ability to (1) experience emotion and feel pain (affective mental states), (2) act and produce an effect on their environment (behavioral potential), and (3) think and hold beliefs (cognitive mental states) (2014). Dehumanization typically involves the denial of one or more of these types of agency (Tipler and Ruscher, 2014).

In Section 3, we introduce computational linguistic methods to quantify several of these components.

## 2.2 Related Computational Work

While this is the first known computational work that focuses on dehumanization, we draw upon a growing body of literature at the intersection of natural language processing and social science. We are particularly inspired by the area of automatically detecting subjective language, largely pioneered by Janyce Wiebe and colleagues who developed novel lexical resources and algorithms for this task (Wiebe et al., 2004). These resources have been used as linguistically-informed features in machine learning classification of biased language (Recasens et al., 2013). Other work has expanded this lexicon-based approach to account for the role of syntactic form in identifying the writer's perspective towards different entities (Greene and Resnik, 2009; Rashkin et al., 2016).

These methods have been used and expanded to analyze pernicious, but often implicit social biases (Caliskan et al., 2017). For example, Voigt et al. analyze racial bias in police transcripts by training classifiers with linguistic features informed by politeness theory (2017), and Garg et al. investigate historical racial biases through changing word embeddings (2018). Other studies focus on how people's positions in different syntactic contexts affect power and agency, and relate these concepts to gender bias in movies (Sap et al., 2017) and news articles about the #MeToo movement (Field et al., 2019). There is also a growing focus on identifying subtle manifestations of social biases, such as condescension (Wang and Potts, 2019), microaggressions (Breitfeller et al., 2019), and "othering" language (Burnap and Williams, 2016; Alorainy et al., 2019). In addition, our focus on dehumanization is closely related to the detection and analysis of hate speech and abusive language (Schmidt and Wiegand, 2017; ElSherief et al., 2018).

Gender and racial bias have also been identified within widely-deployed NLP systems, for tasks including toxicity detection (Sap et al., 2019), sentiment analysis (Kiritchenko and Mohammad, 2018), coreference resolution (Rudinger et al., 2018), language identification (Blodgett and O'Connor, 2017), and in many other areas (Sun et al., 2019). Given the biases captured, reproduced, and perpetuated in NLP systems, there is a growing interest in mitigating subjective biases (Sun et al., 2019), with approaches including modifying embedding spaces (Bolukbasi et al., 2016; Manzini et al., 2019), augmenting datasets (Zhao et al., 2018), and adapting natural language generation methods to "neutralize" text (Pryzant et al., 2019).

A related line of research has developed computational approaches to investigate language use and variation in media discourse about sociopolitical issues. For example, some work has drawn upon political communication theory to automatically detect an issue's framing (Entman, 1993; Boydstun et al., 2013; Card et al., 2015) through both supervised classification (Boydstun et al., 2014; Baumer et al., 2015) and unsupervised methods, such as topic modeling and lexicon induction (Tsur et al., 2015; Field et al., 2018; Demszky et al., 2019). Scholars have also developed computational methods to identify lexical cues of partisan political speech, political slant in mass media, and polarization in social media (Monroe et al., 2008; Gentzkow and Shapiro, 2010; Demszky et al., 2019).

## 2.3 Attitudes Towards LGBTQ Communities in the United States

Some background about LGBTQ communities is necessary for our case study of LGBTQ dehumanization in the *New York Times*. Bias against LGBTQ people is longstanding in the United States. Overall, however, the American public has become more accepting of LGBTQ people and supportive of their rights. In 1977, equal percentages of respondents (43%) agreed and disagreed with the statement that gay or lesbian

relations between consenting adults should be legal (Gallup, 2019). Approval of gay and lesbian relations then decreased in the 1980s; in 1986, only 32% of respondents believed they should be legal. According to Gallup, attitudes have become increasingly positive since the 1990s, and in 2019, 73% responded that gay or lesbian relations should be legal. The Pew Research center began surveying Americans about their beliefs about same-sex marriage in 2001 and found similar trends (Pew Research Center, 2017). Between 2001 and 2019, support for same-sex marriage jumped from 35% to 61%.

In addition to the public's overall attitudes, it is important to consider the specific words used to refer to LGBTQ people. Because different group labels potentially convey different social meanings, and thus have different relationships with dehumanization, our case study compares two LGBTQ labels: *gay* and *homosexual*. The Gallup survey asked for opinions on legality of "homosexual relations" until 2008, but then changed the wording to "gay and lesbian relations". This was likely because many gay and lesbian people find the word *homosexual* to be outdated and derogatory. According to the LGBTQ media monitoring organization GLAAD, *homosexual*'s offensiveness originates in the word's dehumanizing clinical history, which had falsely suggested that "people attracted to the same sex are somehow diseased or psychologically/emotionally disordered" [1]. Beyond its outdated clinical associations, some argue that the word *homosexual* is more closely associated with sex and all of its negative connotations simply by virtue of containing the word *sex*, while terms such as *gay* and *lesbian* avoid such connotations (Peters, 2014). Most newspapers, including the *New York Times*, almost exclusively used the word *homosexual* in articles about gay and lesbian people until the late 1980s (Soller, 2018). The *New York Times* began using the word *gay* in non-quoted text in 1987. Many major newspapers began restricting the use of the word *homosexual* in 2006 (Peters, 2014). As of 2013, the *New York Times* has confined the use of *homosexual* to specific references to sexual activity or clinical orientation, in addition to direct quotes and paraphrases [2].

Beyond differences in how LGBTQ people perceive the terms *gay* or *lesbian* relative to *homosexual*, the specific choice of label can affect attitudes towards LGBTQ people. In 2012, Smith et al. (2017) asked survey respondents about either "gay and lesbian rights" or "homosexual rights". Respondents who read the word "homosexual" showed less support for LGBTQ rights. This effect was primarily driven by high authoritarians, people who show high sensitivity to intergroup distinctions. The authors posit that *homosexual* makes social group distinctions more blatant than *gay* or *lesbian*. This leads to greater psychological distancing, thus enabling participants to remove LGBTQ people from their realm of moral consideration (Smith et al., 2017). Based on prior research and evolving media guidelines, we expect our computational analysis to show that *homosexual* occurs in more dehumanizing contexts than the label *gay*.

## 3    OPERATIONALIZING DEHUMANIZATION

In Section 2.1, we discussed multiple elements of dehumanization that have been identified in social psychology literature. Here we introduce and quantify lexical correlates to operationalize four of these components: *negative evaluations of a target group*, *denial of agency*, *moral disgust*, and *use of vermin metaphors*.

### 3.1    Negative Evaluation of a Target Group

One prominent aspect of dehumanization is extremely negative evaluations of members of a target group (Haslam, 2006). Attribution of negative characteristics to members of a target group in order to exclude that group from "the realm of acceptable norms and values" is specifically the key component of *delegitimization*, a process of moral exclusion closely related to dehumanization. We hypothesize that this negative evaluation of a target group can be realized by words and phrases whose connotations have extremely low valence, where valence refers to the dimension of meaning corresponding to positive/negative (or pleasure/displeasure) (Osgood et al., 1957; Mohammad, 2018). Thus, we propose several valence lexicon-based approaches to measure this component: paragraph-level valence analysis, Connotation Frames of perspective, and word embedding neighbor valence. Each technique has different advantages and drawbacks regarding precision and interpretability.

---

[1]  https://www.glaad.org/reference/lgbtq

[2]  https://www.glaad.org/reference/style

### 3.1.1 Paragraph-level Valence Analysis

One dimension of affective meaning is *valence*, which corresponds to an individual's evaluation of an event or concept, ranging from negative/unpleasant to positive/pleasant (Osgood et al., 1957; Russell, 1980). A straightforward lexical approach to measure *negative evaluations of a target group* involves calculating the average valence of words occurring in discussions of the target group. We obtain valence scores for 20,000 words from the NRC VAD lexicon, which contains real-valued scores ranging from zero to one for valence, arousal and dominance. A score of zero represents the lowest valence (most negative emotion) and a score of one is the highest possible valence (most positive emotion) (Mohammad, 2018). Words with the highest valence include *love* and *happy*, while words with the lowest valence include *nightmare* and *shit*.

We use paragraphs as the unit of analysis because a paragraph represents a single coherent idea or theme (Hinds, 1977). This is particularly true for journalistic writing (Shuman, 1894), and studies on rhetoric in journalism often treat paragraphs as the unit of analysis (e.g. Barnhurst and Mutz, 1997; Katajamaki and Koskela, 2006). Furthermore, by looking at a small sample of our data, we found that paragraphs were optimal because full articles often discuss unrelated topics while single sentences do not provide enough context to understand how the newspaper represents the target group. We calculate paragraph-level scores by taking the average valence score over all words in the paragraph that appear (or whose lemmas appear) in the NRC VAD lexicon.

### 3.1.2 Connotation Frames of Perspective

While paragraph-level valence analysis is straightforward, it is sometimes too coarse because we aim to understand the sentiment *directed towards* the target group, not just nearby in the text. For example, suppose the target group is named "B". A sentence such as "A violently attacked B" would likely have extremely negative valence, but the writer may not feel negatively towards the victim, "B".

We address this by using Rashkin et al.'s Connotation Frames Lexicon, which contains rich annotations for 900 English verbs (2016). Among other things, for each verb, the Connotation Frames Lexicon provides scores (ranging from -0.87 to 0.8) for the writer's perspective towards the verb's subject and object. In the example above for the verb *attack*, the lexicon lists the writer's perspective towards the subject "A", the attacker, as -0.6 (strongly negative) and the object "B" as 0.23 (weakly positive).

We extract all subject-verb-object tuples containing at least one target group label using the Spacy dependency parser [3]. For each subject and object, we capture the noun and the modifying adjectives, as group labels (such as *gay*) can often take either nominal or adjectival forms. For each tuple, we use the Connotation Frames lexicon to determine the writer's perspective towards the noun phrase containing the group label. We then average perspective scores over all tuples.

### 3.1.3 Word Embedding Neighbor Valence

While a Connotation Frames approach can be more precise than word-counting valence analysis, it limits us to analyzing SVO triples, which excludes a large portion of the available data about the target groups. This reveals a conundrum: broader context can provide valuable insights into the implicit evaluations of a social group, but we also want to directly probe attitudes towards the group itself.

We address this tension by training vector space models to represent the data, in which each unique word in a large corpus is represented by a vector (embedding) in high-dimensional space. The geometry of the resulting vector space captures many semantic relations between words. Furthermore, prior work has shown that vector space models trained on corpora from different time periods can capture semantic change (Kulkarni et al., 2015; Hamilton et al., 2016). For example, diachronic word embeddings reveal that the word *gay* meant "cheerful" or "dapper" in the early 20[th] century, but shifted to its current meaning of sexual orientation by the 1970s. Because word embeddings are created from real-world data, they contain real-world biases. For example, Bolukbasi et al. (2016) demonstrated that gender stereotypes are deeply ingrained in these systems. Though problematic for the widespread use of these models in computational systems, these revealed biases indicate that word embeddings can actually be used to identify stereotypes about social groups and understand how they change over time (Garg et al., 2018).

---

[3] spacy.io

This technique can similarly be applied to understand how a social group is negatively evaluated within a large text corpus. If the vector corresponding to a social group label is located in the semantic embedding space near words with clearly negative evaluations, that group is likely negatively evaluated (and possibly dehumanized) in the text.

We first preprocess the data by lowercasing, removing numbers, and removing punctuation. We then use the word2vec skip-gram model to create word embeddings (Mikolov et al., 2013). We use Gensim's default parameters with two exceptions; we train our models for ten iterations in order to ensure that the models converge to the optimal weights and we set the window size to 10 words, as word vectors trained with larger window sizes tend to capture more semantic relationships between words (Levy and Goldberg, 2014)[4]. For our diachronic analysis, we first train word2vec on the entire corpus, and then use the resulting vectors to initialize word2vec models for each year of data in order to encourage coherence and stability across years. After training word2vec, we zero-center and normalize all embeddings to alleviate the hubness problem (Dinu et al., 2014).

We then identify vectors for group labels by taking the centroid of all morphological forms of the label, weighted by frequency. For example, the vector representation for the label *gay* is actually the weighted centroid of the words *gay* and *gays*. This enables us to simultaneously account for adjectival, singular nominal, and plural nominal forms for each social group label with a single vector. Finally, we estimate the valence for each group label by identifying its 500 nearest neighbors via cosine similarity, and calculating the average valence of all neighbors that appear in the NRC VAD Valence Lexicon [5].

We also induce a valence score directly from a group label's vector representation by adapting the regression-based sentiment prediction from Field and Tsvetkov (2019) for word embeddings. This approach yielded similar results as analyzing nearest neighbor valence but was difficult to interpret. More details for and results from this technique can be found in the Supplementary Material.

## 3.2 Denial of Agency

*Denial of agency* refers to the lack of attributing a target group member with the ability to control their own actions or decisions (Tipler and Ruscher, 2014). Automatically detecting the extent to which a writer attributes cognitive abilities to a target group member is an extraordinarily challenging computational task. Fortunately, the same lexicons used to operationalize *negative evaluations* provide resources for measuring lexical signals of *denial of agency*.

### 3.2.1 Connotation Frames

As in Section 3.1, we use Connotation Frames to quantify the amount of agency attributed to a target group. We use Sap et al.'s extension of Connotation Frames for agency (2017). Following Sap et al.'s interpretation, entities with high agency exert a high degree of control over their own decisions and are active decision-makers, while entities with low agency are more passive (2017). This contrast is particularly apparent in example sentences such as *X searched for Y* and *X waited for Y*, where the verb *searched* gives X high agency and *waited* gives X low agency (Sap et al., 2017). Additionally, Sap et al.'s released lexicon for agency indicates that subjects of verbs such as *attack* and *praise* have high agency, while subjects of *doubts* and *needs* have low agency (2017).

This lexicon considers the agency attributed to subjects of nearly 2000 transitive and intransitive verbs. To use this lexicon to quantify *denial of agency*, we extract all sentences' head verbs and their subjects, where the subject noun phrase contains a target group label. Unlike Rashkin et al.'s real-valued Connotation Frames lexicon for perspective, the agency lexicon only provides binary labels, so we calculate the fraction of subject-verb pairs where the subject has high agency.

### 3.2.2 Word Embedding Neighbor Dominance

The NRC VAD Dominance Lexicon provides another resource for quantifying dehumanization Mohammad (2018). The NRC VAD lexicon's dominance dimension contains real-valued scores between zero and one for 20,000 English words. However, the dominance lexicon primarily captures power, which

---

[4] https://radimrehurek.com/gensim/models/word2vec.html

[5] We conducted additional analyses by considering 25, 50, 100, 250, and 1000 nearest neighbors, which yielded similar results and can be found in the Supplementary Material.

is distinct from but closely related to agency. While power refers to one's control over others, agency refers to one's control over oneself. While this lexicon is a proxy, it qualitatively appears to capture signals of *denial of agency*; the highest dominance words are *powerful*, *leadership*, *success*, and *govern*, while the lowest dominance words are *weak*, *frail*, *empty*, and *penniless*. We thus take the same approach as in Section 3.1.3, but instead calculate the average dominance of the 500 nearest neighbors to each group label representation [6].

As in Section 3.1.3, we also induced a dominance score directly from a group label's vector representation by adapting the regression-based sentiment prediction from Field and Tsvetkov (2019) for word embeddings. More details and results for this technique can be found in the Supplementary Material.

## 3.3 Moral Disgust

To operationalize *moral disgust* with lexical techniques, we draw inspiration from Moral Foundations theory, which postulates that there are five dimensions of moral intuitions: care, fairness/proportionality, loyalty/ingroup, authority/respect, and sanctity/purity (Haidt and Graham, 2007). The negative end of the sanctity/purity dimension corresponds to moral disgust. While we do not directly incorporate Moral Foundations Theory in our framework for dehumanization, we utilize lexicons created by Graham et al. (2009) corresponding to each moral foundation. The dictionary for moral disgust includes over thirty words and stems, including *disgust\**, *sin*, *pervert*, and *obscen\** (the asterisks indicate that the dictionary includes all words containing the preceding prefix)[7].

We opt for a vector approach instead of counting raw frequencies of moral disgust-related words because such words are rare in our news corpus. Furthermore, vectors capture associations with the group label itself, while word counts would not directly capture such associations. Using the word embeddings from Section 3.1.3, we construct a vector to represent the *concept* of moral disgust by averaging the vectors for all words in the "Moral Disgust" dictionary, weighted by frequency. This method of creating a vector from the Moral Foundations dictionary resembles that used by Garten et al. (2016). We identify implicit associations between a social group and moral disgust by calculating cosine similarity between the group label's vector and the Moral Disgust concept vector, where a higher similarity suggests closer associations between the social group and moral disgust.

## 3.4 Vermin as a Dehumanizing Metaphor

Metaphors comparing humans to vermin have been especially prominent in dehumanizing groups throughout history (Haslam, 2006; Steuter and Wills, 2010). Even if a marginalized social group is not directly equated to vermin in the press, this metaphor may be invoked in more subtle ways, such as through the use of verbs that are also associated with vermin (like *scurry* as opposed to the more neutral *hurry*) (Marshall and Shapiro, 2018). While there is some natural language processing work on the complex task of metaphor detection (e.g. Tsvetkov et al., 2014), these systems cannot easily quantify such indirect associations.

We thus quantify the metaphorical relationship between a social group and vermin by calculating similarities between these concepts in a distributional semantic vector space. As with *moral disgust*, we create a *Vermin* concept vector by averaging the following vermin words' vectors, weighted by frequency: *vermin*, *rodent(s)*, *rat(s) mice*, *cockroaches*, *termite(s)*, *bedbug(s)*, *fleas* [8]. We do not include the singular *mouse* or *flea* because non-vermin senses of those words were more frequent, and word2vec does not account for polysemy. We calculate cosine similarity between each group label and the *Vermin* concept vector, where a high cosine similarity suggests that the group is closely associated with vermin.

Table 1 provides an overview of the four elements of dehumanization that we study and the lexical techniques used to quantify them.

---

[6] We conducted additional analyses by considering 25, 50, 100, 250, and 1000 nearest neighbors, which yielded similar results and can be found in the Supplementary Material

[7] https://www.moralfoundations.org/othermaterials

[8] Largely inspired by https://en.wikipedia.org/wiki/Vermin

| Dehumanization Element | Operationalization |
|---|---|
| Negative evaluation of target group | Paragraph-level sentiment analysis<br>Connotation frames of perspective<br>Word embedding neighbor valence |
| Denial of agency | Connotation frames of agency<br>Word embedding neighbor dominance |
| Moral disgust | Vector similarity to disgust |
| Vermin metaphor | Vector similarity to vermin |

Table 1: Overview of linguistic correlates and our operationalizations for four elements of dehumanization.

## 4  DATA

The data for our case study spans over thirty years of articles from the *New York Times*, from January 1986 to December 2015, and was originally collected by Fast and Horvitz (2016). The articles come from all sections of the newspaper, such as "World", "New York & Region", "Opinion", "Style", and "Sports". Our distributional semantic methods rely on all of the available data in order to obtain the most fine-grained understanding of the relationships between words possible. For the other techniques, we extract paragraphs containing any of the following words from a predetermined list of **LGTBQ terms**: *gay(s), lesbian(s), bisexual(s), homosexual(s), transgender(s), transsexual(s), transexual(s), transvestite(s), transgendered, asexual, agender, aromantic, lgb, lgbt, lgbtq, lgbtqia, glbt, lgbtqqia, genderqueer, genderfluid, intersex, pansexual.*

Each acronym label is matched insensitive to case and punctuation. Some currently prominent LGBTQ terms, such as *queer* and *trans* are not included in this study, as other senses of these words were more frequent in earlier years. We filter out paragraphs from sections that typically do not pertain to news, such as "Arts", "Theater", and "Movies". While these sections could provide valuable information, we focus on representation of LGBTQ groups in more news-related contexts.

A challenging question when analyzing mass media for subjective attitudes is deciding whose perspective we want to capture: an individual reporter, the institution, or society at large? In this case study, we aim to identify the institution's dehumanizing attitudes towards LGBTQ people. We represent the *New York Times* institution as a combination of the journalists' words in news articles, direct quotes, paraphrases from interviews, and published opinion articles. Therefore, despite our news focus, we include data from "Opinion" sections; while opinion articles are stylistically different from traditional journalistic reporting due to more overt biases and arguments, these articles are important in constructing the institution's perspective. In addition, we consider all text in each relevant paragraph, including quotes and paraphrases, because they are important to a newspaper's framing of an issue, as particular quotes representing specific stances are intentionally included or excluded from any given article (Niculae et al., 2015).

We refer to the remaining subset of the *New York Times* data after filtering as the *LGBTQ corpus*. The *LGBTQ* corpus consists of 93,977 paragraphs and 7.36 million tokens. A large increase in reporting on LGBTQ-related issues has led to a skewed distribution in the amount of data over years, with 1986 containing the least data (1144 paragraphs and 73,549 tokens) and 2012 containing the most (5924 paragraphs and 465,254 tokens).

For all experiments, we also include results for the terms *American* and *Americans*. We include *American(s)* to contrast changes in LGBTQ labels' representation with another social group label. This ensures that the changes we find in dehumanizing language towards LGBTQ groups do not apply uniformly to all social groups, and are thus not merely an artifact of the publication's overall language change. While a natural "control" variable would be labels such as *straight* or *heterosexual*, these terms only occurred within discussions of LGBTQ communities because they name socially unmarked categories. We also considered comparing LGBTQ labels to *person/people*, but because word embedding-based experiments are sensitive to syntactic forms, we opt for a label that behaves more syntactically similar to *gay* and *homosexual*, particularly with both nominal and adjectival uses. Nevertheless, *American(s)* is by no means a neutral control variable. Because of its in-group status for the *New York Times* (a U.S. institution), we
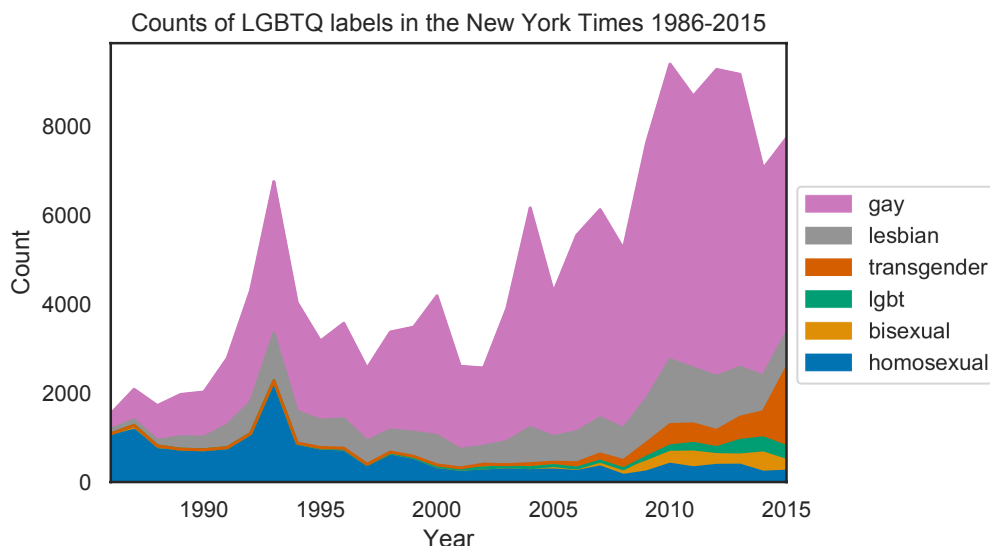
Figure 1: Counts for the six most frequent LGBTQ labels in each year of the *New York Times* data.

expect our measurements to show that *American(s)* appears in more humanizing contexts than LGBTQ labels; however, we do not expect to find substantial changes in the use of *American(s)* over time.

Figure 1 shows the counts of group labels for each year in the *New York Times* from 1986 to 2015. For visualization purposes, only words with a total count greater than 1000 are shown. The relative frequency of *homosexual* decreased substantially over time, while *gay*, *lesbian*, and *bisexual* are more frequent in later years. The terms *lgbt* and *transgender* also emerged after 2000. Counts for all LGBTQ labels can be found in the Supplementary Material.

## 5 RESULTS

### 5.1 Word Embeddings

Using all of the *New York Times* data, we create word2vec models for each year using the methods described in Section 3.1.3. Because our computational techniques rely upon these word2vec models, it is useful to gain a sense of how LGBTQ terms are semantically represented within these models. We thus inspect the ten nearest neighbors, or most similar words, to LGBTQ terms in different years. Note that the neighboring words in Tables 2 and 3 are shown purely for qualitative investigation; our measures for quantifying each dehumanization component incorporate far more information from the word2vec models beyond the top ten neighbors.

Table 2 shows the ten nearest neighbors (by cosine similarity) to our vector representation of all LGBTQ terms, which is the weighted average of the embeddings of all LGBTQ terms considered. For visual convenience, we filter out words occurring fewer than ten times, proper names, as well as other LGBTQ labels and forms of the word *heterosexual*, which are common neighbors for all terms across all years.

Table 2 shows that in 1986, LGBTQ groups were most highly associated with words that often convey a sense of sexual deviancy, including *promiscuity*, *promiscuous*, *polygamy*, *bestiality*, and *pornography*. These associations suggest that LGBTQ people were dehumanized to some extent at this time, and their identities were not fully recognized or valued. This shifted by 2000, where we no longer see associations between LGBTQ groups and ideas that evoke moral disgust. Instead, the 2000 vector space shows that LGBTQ people have become more associated with civil rights issues (suggested by *interracial*, *homophobia*, and *discrimination*). The words *ordination* and *ordaining* likely appear due to major controversies that arose at this time about whether LGBTQ people should be permitted to be ordained. We also see some indications of self-identification with the term *openly*. Finally, we see a slight shift towards associations with identity in 2015, with nearby words including *nontransgender*, *closeted*, *equality*, and *sexuality*.

| 1986 | 2000 | 2015 |
|---|---|---|
| sex | interracial | sex |
| premarital | openly | nontransgender |
| sexual | unwed | unmarried |
| abortion | homophobia | interracial |
| promiscuity | premarital | closeted |
| polygamy | ordination | equality |
| promiscuous | nonwhites | couples |
| vigilantism | ordaining | abortion |
| bestiality | discrimination | sexuality |
| pornography | abortion | antiabortion |

Table 2: Nearest words to weighted average of all LGBTQ terms' vectors in 1986, 2000, and 2015

| 1986 | | 2000 | | 2015 | |
|---|---|---|---|---|---|
| Gay | Homosexual | Gay | Homosexual | Gay | Homosexual |
| homophobia | premarital | interracial | premarital | interracial | premarital |
| women | abortion | openly | openly | sex | sexual |
| feminist | sexual | homophobia | deviant | couples | bestiality |
| vigilante | sex | unwed | interracial | mormons | pedophilia |
| vigilantism | promiscuity | ordination | promiscuity | marriage | adultery |
| suffrage | polygamy | premarital | immoral | closeted | infanticide |
| sexism | anal | abortion | sexual | equality | abhorrent |
| a.c.l.u. | intercourse | antigay | criminalizing | abortion | sex |
| amen | consenting | discrimination | polygamy | unmarried | feticide |
| queer | consensual | marriagelike | consensual | openly | fornication |

Table 3: Nearest words to vector representations of *gay* and *homosexual* in 1986, 2000, and 2015

Curiously, the word *abortion* is a nearby term for all three years. Perhaps this is because opinions towards abortion and LGBTQ rights seem to be divided along similar partisan lines.

Table 3 shows the ten nearest neighboring words to our representations of *gay* and *homosexual* after filtering out proper names, words appearing less than ten times that year, other LGBTQ terms, and forms of *heterosexual*. Table 3 reveals variation in social meaning between *gay* and *homosexual* despite denotational similarity, and these differences intensify over time. In 1986, *gay* is associated with terms of discrimination, civil rights and activism, such as *homophobia*, *feminist*, *suffrage*, *sexism*, and *a.c.l.u*. On the other hand, *homosexual* is primarily associated with words related to sexual activity (e.g. *promiscuity*, *anal*, *intercourse*, *consenting*).

In 1986, this pattern may be due to discussions about sexual transmission of AIDS, but the pejoration of *homosexual* continues over time. While *gay* becomes associated with issues related to marriage equality and identity in 2015, *homosexual* becomes extremely associated with moral disgust and illicit activity, with nearest neighbors including *bestiality*, *pedophilia*, *adultery*, *infanticide*, and *abhorrent*.

This qualitative analysis of word embedding neighbors reveals significant variation and change in the social meanings associated with LGBTQ group labels, with clear relationships to dehumanizing language. We will now present our quantitative results for measuring each component of dehumanization.

## 5.2 Negative Evaluation Towards Target Group

### 5.2.1 Quantitative Results

#### 5.2.1.1 *Paragraph-level Valence Analysis*

Figure 2A shows the average valence for paragraphs containing LGBTQ labels (and *American(s)* for comparison), where a paragraph's valence is simply the average valence over its words (or lemmas) that appear in the NRC VAD Valence Lexicon. The NRC VAD lexicons actually contain several LGBTQ terms, which all have lower than the average valence score of 0.5: *transsexual* (0.264), *homosexual* (0.333), *lesbian*
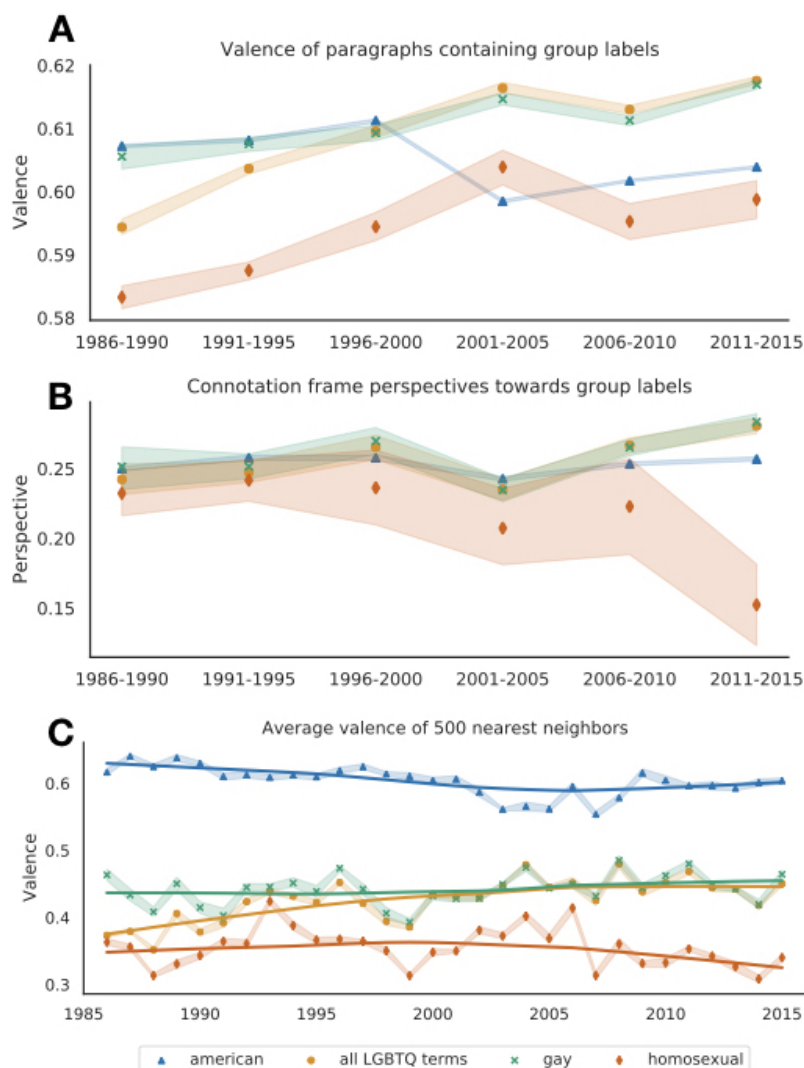
Figure 2: **(A)** Average paragraph-level valence for paragraphs containing *gay*, *homosexual*, any LGBTQ term, and *American*, grouped into 5-year intervals. Paragraph-level scores are calculated as the average valence over all words that appear in the NRC VAD Valence Lexicon, which range from 0 (most negative) to 1 (most positive) (Mohammad, 2018). Paragraphs containing LGBTQ labels become more positive over time. Paragraphs containing *homosexual* are significantly more negative than those containing other LGBTQ labels. **(B)** Average connotation frame perspective scores over five-year intervals. Scores are calculated for each subject-verb-object triple containing these group labels as the writer's perspective based on the head verb's entry in the Connotation Frames lexicon (Rashkin et al., 2016). **(C)** Average valence of 500 nearest words to vector representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American*, averaged over 10 word2vec models trained on *New York Times* data from each year. The solid lines are Lowess curves for visualization purposes. Words' valence scores are from the NRC VAD Valence Lexicon. For all plots, the shaded bands represent 95% confidence intervals.

(0.385), *gay* (0.388), and *bisexual* (0.438). These values contrast starkly with more positively-valenced entries in the lexicon, such as *heterosexual* (0.561), *person* (0.646), *human* (0.767), *man* (0.688), and *woman* (0.865). These disparities likely reveal biases among the human annotators whose judgments were used to construct the NRC VAD lexicon (Mohammad, 2018). While the lexicon may itself be an interesting artifact of dehumanizing attitudes towards LGBTQ people, we remove these terms before calculating paragraph-level valence scores in order to isolate linguistic signals in the *New York Times* data from annotation biases. Without this preprocessing step, the temporal trends and relative differences

between *all LGBTQ terms*, *gay*, and *homosexual* remain roughly the same, but all LGBTQ labels occur in significantly more negative paragraphs than *American*.

Figure 2A shows the average paragraph valence. For visualization purposes, we present the results over five-year intervals due to data sparsity in later years for *homosexual* (there were just 208 paragraphs containing *homosexual* in 2014, relative to 3669 containing *gay* in the same year). Analysis of overlapping confidence intervals and Wilcoxon signed-rank tests over the means for each of the thirty years indicates that *gay* and *all LGBTQ terms* occur in significantly more positive paragraphs than *homosexual* ($p < 0.0001$). A linear regression analysis over all years reveals that *all LGBTQ terms*, *gay*, and *homosexual* all significantly increase in paragraph-level valence over time ($p < 0.0001$). However, when considering just the last 15 years, *gay* still significantly increases in paragraph-level valence, while *homosexual* may be trending downward, although this trend does not reach significance in our data ($p = 0.078$).

The paragraph-level valence analysis shown in Figure 2A suggests that LGBTQ groups have become increasingly positively evaluated over time, and thus likely less dehumanized in the *New York Times*. However, the slight downward trend in valence for paragraphs containing *homosexual* between 2001 and 2015 suggests that evaluations of people described as *homosexual* have not improved in the same way as those described by other labels.

Finally, this measurement does not support our initial hypothesis that LGBTQ groups have been more negatively evaluated than *American(s)*, but still reveals that the observed trends for LGBTQ labels are not merely artifacts of changing reporting styles, since paragraphs containing *American(s)* show a very different pattern. Overall, this result demonstrates substantial language change in the *New York Times*'s discussion of LGBTQ people as well as variation in the contexts where different group labels appear, particularly *homosexual*.

### 5.2.1.2 Connotation Frames of Perspective

Figure 2B shows the writer's average perspective (valence) towards noun phrases containing either any LGBTQ labels, *gay(s)*, *homosexual(s)*, or the comparison group *American(s)* using the Connotation Frames lexicon (Rashkin et al., 2016). The wide variation, particularly for *homosexual*, is likely due to sparsity, as limiting the connotation frames analysis to verbs' immediate subject and direct object noun phrase dependents (consisting of only determiners, adjectives, and nouns) greatly reduced the amount of data for each year; there were only 39 triples for *homosexual* in 2015. We thus show results aggregated over five-year intervals.

As with paragraph-level valence, the writer's perspective towards the label *homosexual* is significantly more negative than towards *gay* ($p < 0.001$). Linear regression indicates that perspectives towards noun phrases named by any LGBTQ term, *gay*, and *American* have all significantly increased over time ($p < 0.01$). However, the trends are still quite different, as the slopes for *gay* and *all LGBTQ terms* are an order of magnitude greater than *American* ($m = (1.1 \pm 0.39) \times 10^{-4}$ for *American*, $m = (1.4 \pm 0.18) \times 10^{-3}$ for *all LGBTQ terms*, and $m = (1.1 \pm 0.22) \times 10^{-3}$ for *gay*). Furthermore, the writer's perspective towards noun phrases containing *homosexual* have significantly decreased over time ($p < 0.0001$).

Overall, Connotation Frames' perspective scores reveal a similar pattern as the paragraph-level valence analysis, where LGBTQ groups overall appear to be more positively evaluated in the *New York Times* over time. Unlike *gay* and the aggregated *all LGBTQ terms*, the label *homosexual* undergoes pejoration, as *homosexual* becomes increasingly used when (implicitly) expressing negative attitudes towards LGBTQ people.

### 5.2.1.3 Word Embedding Neighbor Valence

Figure 2C shows the average valence scores of the 500 nearest neighbors to the vector representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American* for each year. In contrast to our other techniques to quantify *negative evaluations of a target group*, this measurement notably shows that the valence of *American*'s neighboring words is significantly greater than any of the LGBTQ group representations' neighbors every year (Wilcoxon's signed-rank test, $p < 0.0001$), indicating that *American* is used in more positive contexts than LGBTQ terms. Furthermore, all LGBTQ vectors' neighbors have an average valence below the neutral 0.5. The average valence for neighboring words of *gay* and the aggregated *all*

*LGBTQ terms* representation significantly increase over time ($p < 0.0001$), suggesting some increasing humanization in the language used in discussions of LGBTQ people.

Figure 2C also reveals dramatic connotational differences between *gay* and *homosexual*. As shown by non-overlapping confidence intervals and a Wilcoxon signed-rank test, the average valence for *homosexual*'s neighbors is significantly lower than *gay*'s neighbors over all years ($p < 0.0001$). Furthermore, while *gay*'s average neighbor valence increases over time ($p < 0.0001$), *homosexual*'s neighboring words become slightly but significantly more negative over time ($p < 0.001$). Analyzing the valence of the nearest neighbors indicates that *homosexual* has long been used in more negative (and potentially dehumanizing) contexts than *gay*, and that these words' meanings have further diverged as the label *homosexual* has been used in increasingly negative contexts over time.

### 5.2.2 Qualitative Analysis

#### 5.2.2.1 Paragraph-level Valence Analysis

How well does paragraph-level valence analysis capture *negative evaluations of a target group*? To facilitate a qualitative evaluation of this technique, we identify several hundred paragraphs with the highest and lowest average valence. Most paragraphs with high valence scores appear to express positive evaluations of LGBTQ individuals, and those with low scores express negative evaluations.

Table 4 contain examples with extremely high and low valence. We identify several major themes from these results. Most paragraphs with high valence scores emphasize equal rights, while some focus on the activities of advocacy organizations. On the other end, paragraphs with extremely low valence often focus on violence against LGBTQ people, disease (especially AIDS), and LGBTQ issues internationally. Other themes that emerge in low-valence paragraphs include reports on (and direct quotes from) public figures who dehumanized LGBTQ people and portrayals of LGBTQ people as reckless, irresponsible, and angry.

While this technique accurately captures the valence for many paragraphs, we also identify several shortcomings. Some extreme outliers are extremely short paragraphs, including subtitles within articles which are included as paragraphs in the data. Table 5 shows several examples that were mischaracterized by our paragraph-level valence analysis technique. In addition, there are several paragraphs with highly positive average valence that actually express negative evaluations of LGBTQ people. The valence of the third paragraph in Table 5 is skewed by the positive words *supported* and *marriage* even though the paragraph is actually discussing low support for gay marriage. While the fourth paragraph argues that gay couples would be subpar parents relative to straight couples, it uses positive terms such as *love* and *ideal*. Furthermore, kinship terms tend to be assigned highly positive values in the NRC VAD Valence Lexicon, including *child* and *family*. Similarly, even though the final example describes discrimination based on sexual orientation, the paragraph's average valence is impacted by positive kinship terms such as *father* (0.812) and *mother* (0.931) [9].

Overall, our qualitative analysis shows that highly positive valence often accompanies expressions of positive evaluation towards LGBTQ groups, and low valence often accompanies expressions of negative evaluation. However, paragraph-level valence scores are also impacted by specific words cued by various topics; paragraphs about same-sex marriage tend to be more positive because words like *marriage*, *marry*, and *couple* have high valence scores while paragraphs reporting on hate crimes tend to be more negative because they contain low-valence words related to crime, violence, and injury. Furthermore, this method cannot disentangle perspectives within the text; although there are linguistic signals of dehumanization expressed in reports on anti-LGBTQ violence and homophobic speech, these dehumanizing attitudes are not necessarily from the viewpoint of the journalist or the institution. Nevertheless, there could be an overall dehumanizing effect if the media's discussions of a marginalized social group emphasizes such events that harm people. Repeated associations between LGBTQ labels and such negative contexts could potentially contribute to negative evaluations of LGBTQ groups.

#### 5.2.2.2 Connotation Frames of Perspective

To qualitatively analyze how well the connotation frames' lexicon capture *negative evaluation of a target group*, we identify SVO tuples where the verb indicates that the writer has extremely positive or negative

---

[9] We also conducted paragraph-level sentiment analysis using binary positive vs. negative emotion lexicons such as LIWC (Pennebaker et al., 2001), but found similar quantitative results and no qualitative improvement over the VAD lexicon

| Valence | Score | Text | Year | Interpretation |
|---|---|---|---|---|
| High | 0.853 | All Americans, **gay** and non-**gay**, deserve respect and support for their families and basic freedoms. | 2004 | Equality |
| High | 0.804 | The experience of the joy and peace that comes with that — it was a clear indication to me that **homosexual** love was in itself a good love and could be a holy love,' Father McNeill said in the film. | 2015 | Equality |
| High | 0.801 | The Straight for Equality in Sports Award is given by PFLAG National, a non-profit organization for families, friends and allies of **gay**, **lesbian**, **bisexual** and **transgender** people. | 2013 | Advocacy |
| High | 0.780 | What do you consider the most interesting and important **LGBT** organizations working today in the city, with youth or more generally? How about more nationally? | 2010 | Advocacy |
| Low | 0.266 | "We kill the women. We kill the babies, we kill the blind. We kill the cripples. We kill them all. We kill the faggot. We kill the **lesbian**... When you get through killing them all, go to the goddamn graveyard and dig up the grave and kill them a-goddamn-gain because they didn't die hard enough." | 1993 | Direct Quote |
| Low | 0.364 | A 21-year-old college student pleaded guilty yesterday to fatally stabbing a **gay** man in Queens in what prosecutors termed a vicious burst of anti-**homosexual** violence. | 1991 | Violence |
| Low | 0.403 | One of his most difficult clients was a **transsexual** prostitute and drug addict who was infected with the AIDS virus and presumably spreading it to her customers and fellow addicts. | 1987 | AIDS |
| Low | 0.373 | Enabling promiscuity, indeed! Burroughs Wellcome is as responsible for the reckless abuse of amyl nitrate by **homosexuals** as the manufacturers of narcotic analgesics are for the horrors of opiate addiction. | 1996 | Recklessness |
| Low | 0.402 | The activists from Africa shrugged with resignation and sank back down on the benches. The **gay** Americans absolutely exploded at the poor woman from the airline. | 2011 | Recklessness |
| Low | 0.397 | Homosexuality is forbidden in Iran. Last year a United Nations report on human rights in Iran expressed concern that **gays** "face harassment, persecution, cruel punishment and even the death penalty. | 2012 | International |

Table 4: Example paragraphs with extremely high and low valence scores, along with an interpretation of the patterns we find. Words with extremely high valence scores (greater than 0.85) appear in blue, and somewhat high-valence words (scores between 0.7 and 0.85) appear in light blue. Words with extremely low valence scores (less than 0.15) appear in red, and somewhat low-valence words (scores between 0.15 and 0.3) appear in pink.

| Valence | Score | Text | Year | Explanation |
|---|---|---|---|---|
| High | 0.929 | Blessing of **Homosexuals** | 1990 | Subtitle |
| Low | 0.031 | Hate for Liberals and **Gays** | 2008 | Subtitle |
| High | 0.777 | Of the seven in attendance, only the Rev. Al Sharpton and Representative Dennis J. Kucinich supported **gay** marriage unambiguously. | 2003 | Marriage |
| High | 0.765 | And I believe children can receive love from **gay** couples, but the ideal is – and studies have shown that the ideal is where a child is raised in a married family with a man and a woman. | 2005 | Marriage Family |
| High | 0.776 | Ms. Bright, now a college sophomore, grew up in her mother's home but regularly visited her **gay** father, Lee, in Cartersville, Ga. She remembers when a friend was not allowed to visit her father's home because he was gay. | 1993 | Family |

Table 5: Examples mischaracterized by paragraph-level valence analysis. Words with extremely high valence scores (greater than 0.85) appear in blue, and somewhat high-valence words (scores between 0.7 and 0.85) appear in light blue. Words with extremely low valence scores (less than 0.15) appear in red, and somewhat low-valence words (scores between 0.15 and 0.3) appear in pink.

| Perspective | Score | Text | SVO | Year |
|---|---|---|---|---|
| Negative | -0.83 | The most forceful comment came from Cardinal Anthony J. Bevilacqua of Philadelphia, who said his archdiocese screened out gay candidates. 'We feel a person who is homosexual-oriented is not a suitable candidate for the priesthood, even if **he** had never **committed any homosexual act**,' the cardinal said. | S: he<br>V: committed<br>O: any homosexual act | 2002 |
| Positive | +0.80 | 'Gays are accepted here and respected here,' said Mayor Tony Tarracino. '**The gays saved a lot** of the oldest parts of town, and they brought in art and culture. They deserve a lot of credit for what Key West is today.' | S: the gays<br>V: saved<br>O: a lot | 1990 |
| Positive | +0.80 | In his speech, **he praised gay rights advocates** for their hard work and also thanked many elected officials, including his predecessor, Gov. David A. Paterson, and the four Republican state senators who provided the critical votes to pass the marriage bill and whom Mr. Cuomo named one by one to some of the loudest applause of the evening. | S: he<br>V: praised<br>O: gay rights advocates | 2011 |

| Assigned Perspective | Score | Text | SVO | Year |
|---|---|---|---|---|
| Negative | -0.87 | Previously, Judge Vaughn Walker, who ruled the ban against same-sex unions unconstitutional in federal court, had said that ProtectMarriage could not appeal his decision to the Ninth Circuit, because they were never able to prove that **gay marriage harmed them** in any way. | S: gay marriage<br>V: harmed<br>O: them | 2011 |
| Positive | +0.73 | Following are excerpts from opinions by the Supreme Court today in its decision that **the Constitution** does not **protect private homosexual relations** between consenting adults (...) Justice Stevens wrote a separate dissenting opinion, joined by Justices Brennan and Marshall. | S: the Constitution<br>V: protect<br>O: private homosexual relations | 1986 |
| Positive | +0.70 | Do you know there is a Congressional candidate from Missouri who is saying that allowing **gays** into the military could **strengthen Al Qaeda**? I'm thinking, how exactly would that work? 'They dance better than me, and they know how to accessorize. I'm very, very angry. It's time for jihad.' | S: gays<br>V: strengthen<br>O: Al Qaeda | 2010 |

Table 6: Examples of paragraphs where the writer expresses highly positive and negative perspective towards LGBTQ groups, according to the Connotation Frames lexicon. Below the double line are examples of paragraphs where the writer's perspective is mischaracterized by the Connotation Frames lexicon.

perspective towards either the subject or object. The first paragraph in Table 6 contains an SVO tuple where the writer has the most negative perspective towards the noun phrases containing a group label. Inside a direct quote, this paragraph uses the phrase *any homosexual act* as the object to the verb *committed*, which has the effect of framing homosexuality as a crime. By deeming gay candidates unworthy of the priesthood, the speaker clearly negatively evaluates LGBTQ people. On the opposite end, many paragraphs labeled as containing extremely positive perspectives towards LGBTQ groups do appear to have positive evaluations of these groups. The second and third paragraphs of Table 6 illustrate this, where *the gays* are viewed positively for having *saved* a town, and *gay rights advocates* are *praised* for their work.

However, we found several instances where paragraphs are mislabeled, shown in the bottom half of Table 6. In the fourth paragraph of Table 6, our technique identifies *gay marriage* as the subject of the negative-perspective verb *harmed*, but does not account for the preceding text, which actually contradicts the premise that *gay marriage* causes harm, and thus does not overtly negatively evaluate of LGBTQ groups (although this particular example reveals the difficulty of operationalizing this component because ProtectMarriage groups strongly oppose same-sex marriage and may have negative evaluations of LGBTQ people). The second example similarly shows that this method does not adequately account for various forms of negation, as the positive-perspective verb *protect* is actually negated. The last example in Table 6 presents a complex case that is even challenging for qualitative analysis. Our method identifies *gays* as the subject of the verb *strengthen*, even though the subject should be the gerund *allowing gays (into the military)*, and the lexicon's entry for the writer's perspective towards the subject of of *strengthen* is a highly positive 0.7. However, the object of this verb is the terrorist organization *Al Qaeda*; our background knowledge suggests that the capacity to *strengthen* Al Qaeda would reflect negative perspectives. However, this additional context provided by the rest of the paragraph indicates that the writer is being sarcastic and considers the proposition that gays have any impact on strengthening Al Qaeda to be ridiculous. Finally, the writer emphasizes their own stance in opposition to the Missouri congressional candidate by calling upon common stereotypes of gay people being good at dancing and accessorizing.

Measuring the connotation frames' lexicon perspective scores over verbs' subjects and direct objects cannot leverage as much context as measuring valence over paragraphs using the NRC VAD lexicon labeled

for 20,000 words. However, this technique can make more fine-grained distinctions regarding the writer's (and institution's) attitudes directed towards LGBTQ people and is not as dramatically impacted by the emotional valence of the topic discussed. Neither technique can disentangle the journalist's perspective from those expressed by others and simply reported by the journalist. While removing direct quotations may partially address this issue, we deliberately do not remove text from direct quotes or paraphrases. The journalists and newspaper make intentional decisions about what text to include and exclude from quotations, which could still meaningfully represent their perspectives and values (Niculae et al., 2015).

### 5.2.2.3 Word Embedding Neighbor Valence

Compared to the previous methods, one limitation of using word embeddings to quantify *negative evaluations of a target group* is that embeddings are not easily interpretable by analyzing a small sample of data. Instead, we assess this technique by identifying LGBTQ terms' nearest neighbors in several outlier years. To facilitate this qualitative analysis, we identify a set of *unique nearest neighbors* for each LGBTQ label in each outlier year, where a word is a unique nearest neighbor for a given LGBTQ term and year if it is not in that term's top 500 nearest neighbors in any other year.

| Valence | Year | Example |
|---------|------|---------|
| Low | 1999 | Matthew Shepard, a **gay** college student in Wyoming, had been pistol-whipped and left to die after being tied to a fence on Oct. 7, 1998. Aaron McKinney, who was charged with first-degree murder and other crimes in connection with Mr. Shepard's killing, went on trial Monday, denying that the act was a hate crime, but rather connected to drug use and outrage at a sexual advance he said Mr. Shepard made. |
| Low | 2014 | Uganda's vehement anti-**gay** movement began in 2009 after a group of American preachers went to Uganda for an anti-**gay** conference and then worked with Ugandan legislators to draft a bill that called for putting **gay** people to death. While the bill was being debated, attacks against **gay** Ugandans began to increase. In early 2011, David Kato, a slight, bespectacled man and one of the country's most outspoken **gay** rights activists, was beaten to death with a hammer. |
| Low | 2014 | 'Hey, @McDonalds: You're sending #CheersToSochi while goons wearing Olympic uniforms assault **LGBT** people,' read one comment last week, from the author and activist Dan Savage. |
| High | 1993 | The regulations, which are to take effect Feb. 5, codify the Administration's policy that was worked out as a compromise with the Joints Chiefs of Staff, who had defended the 50-year-old ban, arguing that allowing **homosexuals** to serve openly would hurt the morale of troops, and thus hurt military readiness. |

Table 7: Example paragraphs from years where LGBTQ terms' nearest neighbors had exceptionally high and low valence.

Table 7 contains several example paragraphs that illustrate overarching themes for the outlier years 1993, 1999, and 2014. In 1999, *gay*, *homosexual* and the aggregated representation of *all LGBTQ terms* were all more more closely associated with low-valence words than in almost any other year. We connect this finding to a period of intense reporting in the months following the October 1998 murder of a gay Wyoming college student, Matthew Shepard, which drew national attention to anti-LGBTQ violence. Because LGBTQ labels frequently co-occurred with text about this incident, terms related to Matthew Shepard's case had closer representations to LGBTQ terms in this year. For example, *gay* and *all LGBTQ terms*'s 500 nearest neighbors include *wyoming* in 1999 and *shepard* from the years 1998-2000. Unique nearest neighbors for *gay* in 1999 include other terms that could be connected to this incident, including *homicidal*, *imprisoned*, and *hatred*. Not only was Shepard's murder rooted in the dehumanization of LGBTQ people, but the media's emphasis on the gruesome details of Shepard's death further dehumanized him (Ott and Aoki, 2002). Ott and Aoki argue that the media's framing of this case actually further stigmatized LGBTQ people.

Our word embedding neighbor valence measure reveals that the most negative year for *gay* and *all LGBTQ terms* since 1999 was 2014, the second most-recent year of data. We identify several major themes in 2014 that co-occurred with LGBTQ group labels and possibly led to this distributional semantic pattern, primarily reporting on anti-LGBTQ laws and attitudes in Uganda and Russia (particularly in light of the 2014 Winter Olympics in Sochi). The terms *athletes* and *winterolympics* appeared in *gay*'s nearest neighbors in 2014. In addition, the terms *Uganda*, *Ugandan*, and *Mugisha* (a Ugandan LGBT advocate) are among *gay*'s unique nearest 500 neighbors in 2014.
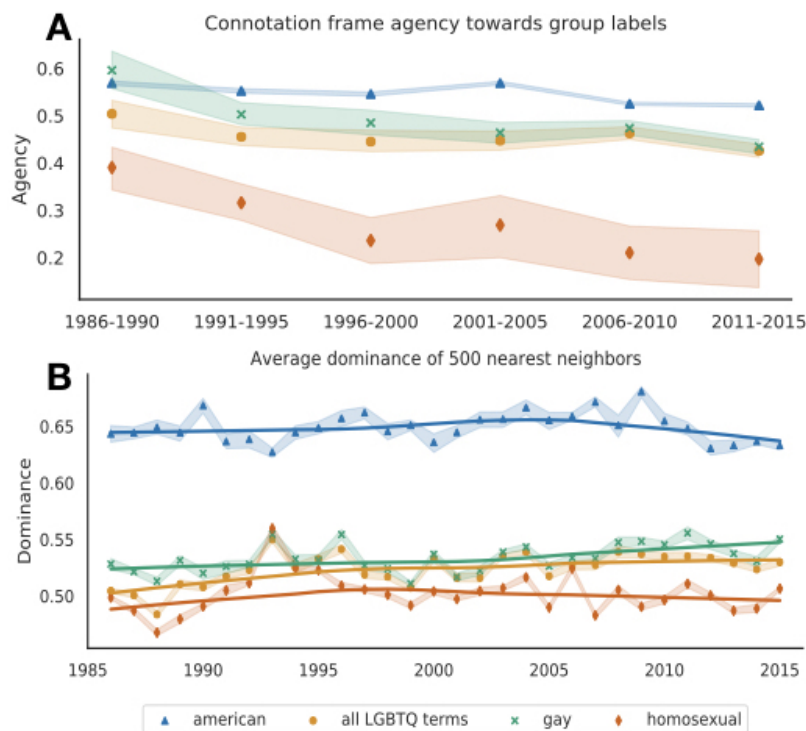
Figure 3: **(A)** Agency of *gay*, *homosexual*, *all LGBTQ terms*, and *American* using the Connotation Frames lexicon for agency for all subject-verb-object tuples containing each group label (Sap et al., 2017), calculated over five-year intervals. An SVO tuple received a score of 1 if the label appears in a positive agency position relative to its head verb and 0 if it does not. **(B)** Average dominance of 500 nearest words to our representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American*, averaged over 10 word2vec models trained on *New York Times* data for each year. Dominance scores for each word come from the word's entry in the NRC VAD Dominance Lexicon (Mohammad, 2018), which range from 0 (least dominance) to 1 (most dominance). For both plots, the shaded bands represent 95% confidence intervals and the solid lines in (B) are Lowess curves for visualization purposes.

Unlike in 1999 and 2014, LGBTQ terms in 1993 are associated with higher-valence words, especially *homosexual*. *Homosexual*'s unique nearest neighbors in 1993 include the high-valence words *pledge*, *civilian*, *readiness*, and *inclusion*. These words are likely connected with numerous stories in 1993 covering the controversy over whether LGBTQ people should be allowed to serve in the military.

## 5.3  Denial of Agency

### 5.3.1  Quantitative Results

#### 5.3.1.1  Connotation Frames of Agency

Figure 3A shows the agency of each group label based on its head verb's entry in the Connotation Frames lexicon for agency (Sap et al., 2017). As in Figure 2B, there is large variance due to data sparsity when using the Connotation Frames lexicon, particularly for *homosexual*, which is considerably less frequent than *gay* or other LGBTQ terms in later years. In order to maximize precision when extracting subject-verb pairs, we extract only nouns and their immediate adjectival modifiers, which limits the amount of data. We thus show average agency over five-year intervals.

Wilcoxon signed-rank tests on the means for each group labels over all years indicate that *gay* occurs in contexts with significantly higher agency than *homosexual* ($p < 0.0001$). All four group labels significantly decrease in agency over time according to linear regressions over all 30 years ($p < 0.001$), but the slope for *homosexual* is much greater ($m = (-7.9 \pm 1.3) \times 10^{-3}$ for *homosexual*, compared to $m = (-3.9 \pm .55) \times 10^{-3}$ for *gay*, and $m = (-1.5 \pm .46) \times 10^{-3}$ for *all LGBTQ terms*). Furthermore in the most

| Agency | Text | SVO | Year |
|---|---|---|---|
| High | Within the close-knit world of professional childbearers, many of whom share their joys and disillusionments online and in support groups, **gay couples** have **developed a reputation** as especially grateful clients... | S: gay couples<br>V: developed<br>O: a reputation | 2005 |
| High | Tonight, **the gay rights group** Stonewall Democrats will **endorse a candidate** for A.G. It's a relatively big prize in the four-man Democratic primary, given that liberal city voters will have relatively serious sway... | S: the gay rights group<br>V: endorse<br>O: a candidate | 2006 |
| Low | Nigeria's **gay men** and lesbians regularly **face harassment** and arrest, gay activists here say. The criminal code bans acts 'against the order of nature,' and imposes sentences of up to 14 years for those convicted... | S: gay men<br>V: face<br>O: harassment | 2005 |
| Low | Much of the debate among military and civilian officials is now focusing on some version of an approach called 'don't ask, don't tell.' (...) But under the 'don't tell' element, there would be restrictions on the extent to which **homosexuals** could **acknowledge their homosexuality**. | S: homosexuals<br>V: acknowledge<br>O: their homosexuality | 1993 |

Table 8: Examples where the writer attributes high and low agency towards LGBTQ groups, according to the Connotation Frames lexicon for agency.

recent 15 years, *gay* and *all LGBTQ terms* show no significant change ($p = 0.097$ for *gay* and $p = 0.14$ for *all LGBTQ terms*), but *homosexual* still decreases significantly in agency ($p < 0.05$).

Figure 3A suggests that LGBTQ groups experience greater denial of agency in the *New York Times* than the institution's in-group identifier *American*. Furthermore, people described as *homosexual* experience even more denial of agency than people who are described as *gay*. Unlike the improving attitudes indicated by our analysis of *negative evaluations of a target group*, it appears that *denial of agency* increased over time for all LGBTQ groups. However, the relatively rapid decrease in agency for *homosexual* is consistent with other results suggesting *homosexual*'s pejoration.

### 5.3.1.2 Word Embedding Neighbor Dominance

Figure 3B shows the average dominance of each group label's 500 nearest neighbors. *American* is significantly associated with greater dominance than *gay*, *homosexual*, and *all LGBTQ terms* (Wilcoxon signed-rank test; $p < 0.0001$), and *gay* has significantly higher dominance than *homosexual* ($p < 0.0001$). While the dominance associated with *gay* and *all LGBTQ terms* significantly increased over time ($p < 0.0001$), the dominance associated with *homosexual* did not significantly change ($p = 0.65$). Furthermore, the average nearest neighbor dominance for *homosexual* decreased in the most recent 15 years ($p < 0.01$).

Even though dominance may more directly encode *power* rather than *agency*, the NRC VAD Dominance Lexicon is useful for operationalizing *denial of agency* because of the close relationship between these concepts. As with Connotation Frames of agency, these results suggest that LGBTQ groups experience greater *denial of agency* than the *New York Times*'s in-group *American*. Both techniques show differences between the labels *gay* and *homosexual*, where *homosexual* is consistently associated with lower agency than *gay* and further decreases over time. However, these two measurements suggest different temporal dynamics for the *denial of agency* of LGBTQ people; Connotation Frames' agency slightly decreases for *all LGBTQ terms* over time, but increases with word embedding neighbor dominance.

### 5.3.2 Qualitative Analysis

#### 5.3.2.1 Connotation Frames of Agency

We qualitatively investigate the labels assigned by this technique for a sample of paragraphs. In general, the binary labels for positive and negative agency seem reasonably accurate, as shown by the first four example in Table 8. Verbs that attribute high agency to the subject include *develop* and *endorse*, suggesting that the LGBTQ-aligned subjects are in control and actively making their own decisions. On the other end, LGBTQ people have low agency when they are the subjects of passive verbs such as *face* and *acknowledge*.

The Connotation Frames lexicon for agency seems to be especially accurate for low agency; we could not find counterexamples in our sample where LGBTQ people were portrayed with high agency but labeled with low agency. However, we found several mischaracterizations where LGBTQ people were labeled as having high agency but are not portrayed as agentive or in control of their own actions. Our Connotation Frames technique considers the example below to attribute high agency to LGBTQ people
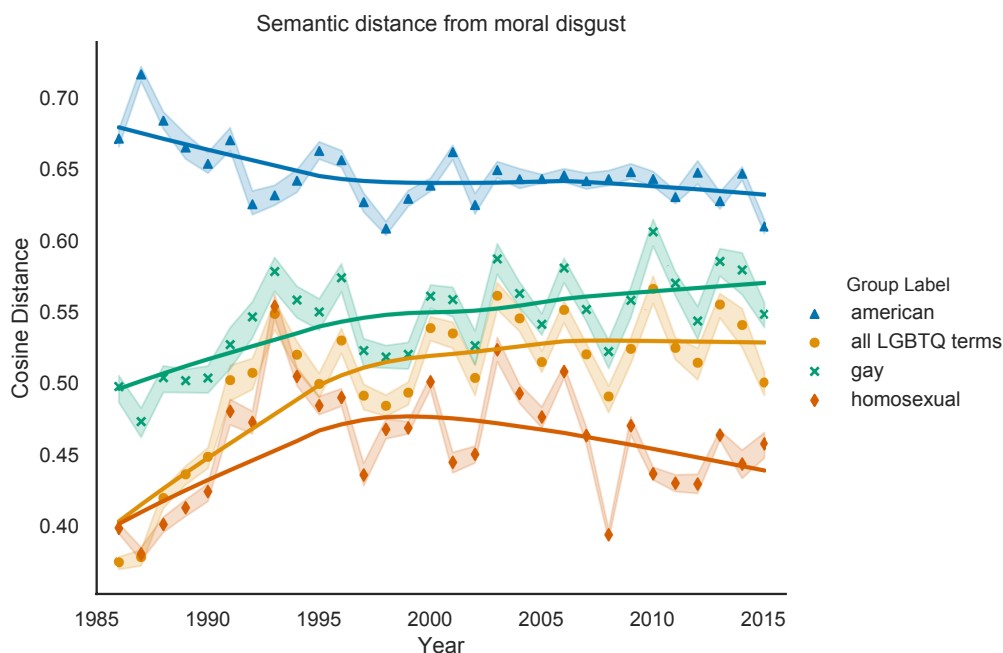
Figure 4: Cosine distance between our representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American* and the vector representation of the *Moral Disgust* concept, averaged over 10 word2vec models trained on *New York Times* data for each year. Increases in cosine distance indicate decreases in *Moral Disgust*; possible values range from 0 (most closely associated with Moral Disgust) to 1 (least associated with Moral Disgust). Shaded bands represent 95% confidence intervals and the solid lines are Lowess curves for visualization purposes.

because *homosexual* appears in the subject of the high-agency verb *violate*; however, *homosexual* actually modifies *relationships*, not people themselves. Furthermore, this debate within religion appears to be devoid of input from LGBTQ people and does not portray them as particularly agentive.

- At the same time, it underscored a stark division in Judaism over the place of homosexuals in society. Orthodox rabbinical groups believe that **homosexual relationships violate Jewish law**... (1996)

### 5.3.2.2 Word Embedding Neighbor Dominance

Using the VAD Dominance Lexicon to calculate average dominance of each social group label corresponds well to our notion of *denial of agency*. Because *gay*'s nearest neighbors have a much higher average dominance than *homosexual*'s for most years, we compare words that are nearby neighbors for *gay* and not *homosexual* for multiple years' word2vec spaces. Words frequently among the 500 words nearest to *gay* and not *homosexual* include high-agency words such as *activist*, *liberation*, *advocate*, and *advocacy*, which have dominance scores of 0.877, 0.857, 0.818, and 0.731, respectively. Words frequently among *homosexual*'s 500 nearest neighbors and not *gay*'s include low-agency words such as *submissive* (0.173), *degrading* (0.232), *enslavement* (0.302), and *repressed* (0.311).

We additionally investigate the word2vec models corresponding to several outlier years. *Homosexual*'s neighbors have the highest average dominance in 1993, which is likely due to military-related language in debates surrounding the "Don't Ask, Don't Tell" legislation. High-dominance words unique to *homosexual*'s nearest neighbors in 1993 include *forces* (0.886), *military* (0.875), *enforce* (0.836) and *troops* (0.804). *Gay*'s neighbors' in 1999 have the lowest average dominance than any other year, which is likely connected to Matthew Shepard's death and the subsequent outrage; unique neighbors for *gay* in 1999 include *imprisoned* (.302) and *repressed* (0.311).

## 5.4 Moral Disgust

### 5.4.1 Quantitative Results

Figure 4 shows the changing relationships between *all LGBTQ terms*, *gay*, *homosexual* and the dehumanizing concept of *Moral Disgust*. Because the cosine distance between *American* and *Moral Disgust* is significantly greater over all years than any LGBTQ representation (Wilcoxon signed-rank test; $p < 0.0001$), *American* is the least associated with *Moral Disgust*. Furthermore, the cosine distance between *gay* and *Moral Disgust* is significantly greater than the distance between *homosexual* and *Moral Disgust* for every year ($p < 0.0001$), indicating that *homosexual* is more closely associated with *Moral Disgust* than *gay* is. Linear regression analyses show that *all LGBTQ terms* and *gay* significantly increase in cosine distance from the *Moral Disgust* vector ($p < 0.0001$), indicated weakening associations between LGBTQ people and moral disgust over time. On the other hand, the distance between *homosexual* and *Moral Disgust* does not change significantly over time ($p = 0.54$), and even decreases after 2000 ($p < 0.05$).

Overall, these measurements of associations between LGBTQ people and *Moral Disgust* are consistent with our other operationalizations of dehumanization. All LGBTQ labels are more closely associated with *Moral Disgust* than the newspaper's in-group term *American*, but these associations weaken over time, suggesting increased humanization. Notably, the term *homosexual* has always been more associated with *Moral Disgust* than the denotationally-similar term *gay*, and *homosexual* actually becomes more closely associated with this dehumanizing concept in recent years.

### 5.4.2 Qualitative Analysis

Our analysis of *homosexual*'s changing semantic neighbors from Table 3 has shown that this term has become more associated with immoral concepts, suggesting that moral disgust is a mechanism by which LGBTQ people are dehumanized. Although rarely directly invoked, the connection between LGBTQ people and disgust is supported by the data, such as in the examples shown below, where words belonging to the moral disgust lexicon are in bold. Figure 4 indicates that late 1980s and early 1990s, LGBTQ labels rapidly became more semantically distant from *Moral Disgust*. This likely reflects decreasing attention to the AIDS epidemic, as many disease-related words are included in the moral disgust lexicon.

- Senator Jesse Helms, the North Carolina Republican who has vigorously fought homosexual rights, wants to reduce the amount of Federal money spent on AIDS sufferers, because, he says, it is their "deliberate, **disgusting**, revolting conduct" that is responsible for their **disease**. (1995)

- A lawyer named G. Sharp, address unknown, called the cover picture "utterly **repulsive**." Donald Ingoglia of Sacramento was equally outraged. "Showing two smiling gays on the cover illustrates how **sick** our society has become," he wrote. "You have my nonlawyer friends falling off their chairs." (1992)

- ...Mr. Robison could be harsh – he yelled in the pulpit and referred to gay men and lesbians as **perverts** – but Mr. Huckabee was a genial ambassador ... (2008)

- ...When bishops started telling parishioners that their gay and lesbian siblings were **sinners**, and that family planning was a grievous wrong, people stopped listening to them – for good reason. (2013)

## 5.5 Vermin as a Dehumanizing Metaphor

### 5.5.1 Quantitative Results

Figure 5 shows the relationships between LGBTQ labels (and *American*) and the dehumanizing vermin metaphor, quantified as the cosine distance between the labels' word2vec vectors and a *Vermin* concept representation, which is the centroid of multiple vermin-related words. As with *Moral Disgust*, the in-group term *American* is further away from *Vermin* over all years than any LGBTQ term (Wilcoxon signed-rank test; $p < 0.0001$). The cosine distance between *gay* and *Vermin* is also greater than between *homosexual* and *Vermin* ($p < 0.0001$), indicating that *homosexual* is more closely associated with the dehumanizing vermin metaphor than *gay* is. Furthermore, while *all LGBTQ terms* and *gay* become more semantically distant from *Vermin* over time, ($p < 0.0001$), the association between *Vermin* and *homosexual* does not significantly change over time ($p = 0.13$).

This measure of the implicit *vermin metaphor* reveals similar patterns as the other dehumanization measures. Overall, LGBTQ groups are more associated with vermin than the comparison group *American*,
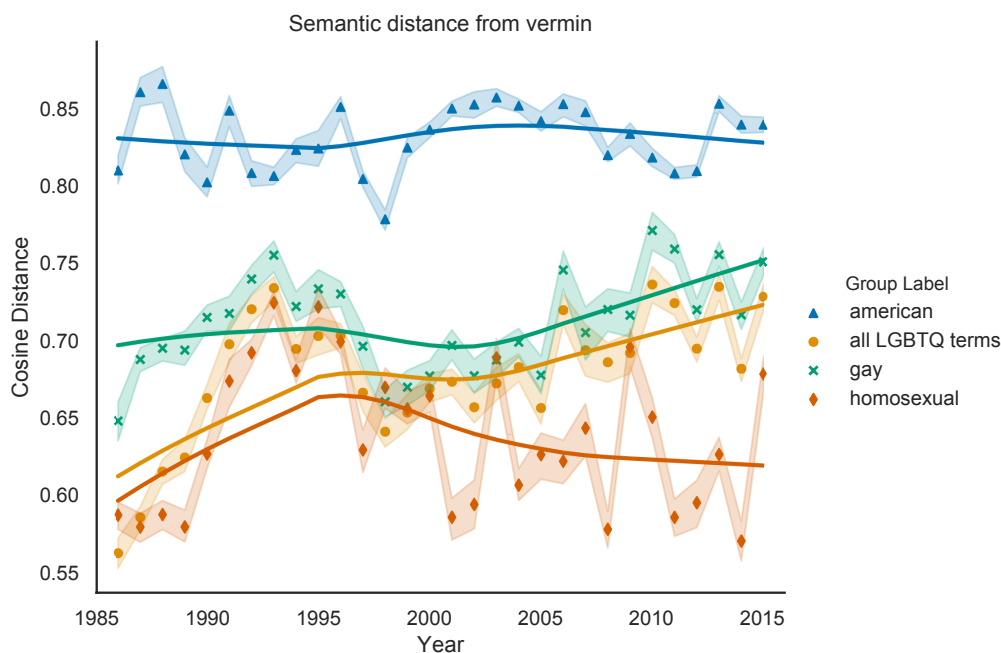
Semantic distance from vermin

Figure 5: Cosine distance between our representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American* and the vector representation of the *Vermin* concept, averaged over 10 word2vec models trained on *New York Times* data for each year. Possible values for cosine distance range from 0 (most closely associated with *Vermin*) to 1 (least associated with *Vermin*). Shaded bands represent 95% confidence intervals, and the solid lines are Lowess curves for visualization purposes.

but this association weakens over time, suggesting increased humanization. In addition, *homosexual* has become a more dehumanizing term, with stronger associations with vermin than other LGBTQ labels.

### 5.5.2 Qualitative Analysis

Metaphors comparing humans to vermin have been especially prominent in dehumanizing groups throughout history (Haslam, 2006; Steuter and Wills, 2010). Although no New York Times writers directly compare LGBTQ people to vermin, this metaphor may be invoked in more subtle ways. There are only three paragraphs in the *LGBTQ corpus* that explicitly mention vermin in order to criticize the LGBTQ people-as-vermin metaphor. Nevertheless, these paragraphs point to the existence of this metaphor.

- Since gay women can't be stigmatized en masse with AIDS, the council had to use real ingenuity to prove that they, too, are vermin at "much greater risk from one another" than from gay-bashers ... (1998)

- The equating of gay men to vermin is appalling," Addessa said from Philadelphia. "We need to encourage the Eagles and Owens to make a public apology and for the Eagles to publicly discipline Owens. These comments that equate gay men to some inferior life form do real harm, creating a cultural environment which justifies violence against gay and lesbian people. (2004)

- In three hours at training camp Tuesday, he hustled vigorously through practice, eagerly signed autographs for visiting military personnel and tried to explain incendiary remarks that appeared in a magazine regarding the sexual orientation of a former teammate in San Francisco, words that seemed to compare gays to rodents. (2004)

## 6 HUMAN EVALUATION OF VECTOR-BASED MEASURES

Our vector-based methods can directly capture associations between LGBTQ people and dehumanizing concepts. However, findings from these methods are difficult to interpret, as discussed in earlier qualitative analysis sections. Furthermore, while the NRC VAD Lexicon and the Connotation Frames Lexicons have been evaluated in prior work (Rashkin et al., 2016; Sap et al., 2017; Mohammad, 2018), our vector-based methods have not. Thus, we recruit humans from Amazon Mechanical Turk (MTurk) to quantitatively evaluate our four vector-based measures: word embedding neighbor valence (for *negative evaluation of a*

*target group*), word embedding neighbor dominance (for *denial of agency*), semantic distance from the concept of *moral disgust*, and semantic distance from the concept of *vermin*.

Although these four measures rely on vector representations of LGBTQ labels and not individual paragraphs, we use paragraphs as the unit of analysis for our evaluation in order for the task to be feasible for human annotators. In Section 6.1, we describe how we use our vector-based methods to obtain the most and least dehumanizing paragraphs for each dehumanization component. We discuss the MTurk task design in Section 6.2 and results in Section 6.3.

## 6.1 Identifying the most (de)humanizing paragraphs

### 6.1.1 Word embedding neighbor valence and dominance

Our word embedding neighbor valence and dominance methods are proxies for measuring the *negative evaluation of the target group* and *denial of agency* dimensions of dehumanization, respectively. They directly estimate the valence and dominance scores for LGBTQ terms based on NRC VAD entries for each term's semantic neighbors.

To obtain full paragraphs corresponding to the most and least dehumanizing extremes of *negative evaluation of a target group*, we first train word2vec on the entire *New York Times* dataset using the same hyperparameters as in Section 3.1.3. Let $N$ be the nearest 500 words to the representation of *all LGBTQ terms* in this vector space, and let $V$ and $D$ be the full NRC Valence and Dominance Lexicons. We define subset lexicons, $V_s = N \cap V$ and $D_s = N \cap D$; $V_s$ and $D_s$ are the subsets of the NRC Valence and Dominance Lexicons containing only words that neighbor *all LGBTQ terms*. We calculate *neighbor valence* scores for each paragraph $P$ as $\frac{1}{|P|}\Sigma_{w \in P}V_s[w]$, where $|P|$ is the total number of tokens in $P$ and $V_s[w]$ is the valence score of $w$. Similarly, we calculate *neighbor dominance* scores as $\frac{1}{|P|}\Sigma_{w \in P}D_s[w]$.

For human evaluation, we consider paragraphs with the highest and lowest scores for *neighbor valence* and *neighbor dominance*. We remove paragraphs containing fewer than 15 or more than 75 words. Because our case study focuses on the words *gay(s)* and *homosexual(s)*, we further restrict our sample to paragraphs containing these terms.

### 6.1.2 Moral disgust and vermin metaphor

We measure implicit associations of LGBTQ groups with *moral disgust* and *vermin* by calculating the cosine distance between LGBTQ terms' vectors and vector representations of *moral disgust* and *vermin*. Thus, we identify paragraphs corresponding to the most and least dehumanizing extremes by comparing the cosine distance between paragraph embeddings and the *Moral Disgust* and *Vermin* concept vectors. We create each paragraph's embedding by calculating the tfidf-weighted average of all words' vectors and removing the first principal component, which improves the quality of sentence and document embeddings (Arora et al., 2017).

We select the paragraphs that are the closest (most semantically similar) and furthest from the *Moral Disgust* and *Vermin* vectors based on cosine distance. As in Section 6.1.1, we limit our sample to paragraphs containing between 15 and 75 words and either the term *gay(s)* or *homosexual(s)*.

## 6.2 MTurk task design

As discussed in our qualitative analyses, journalistic text captures numerous perspectives, not only from journalists themselves, but also from people quoted and people or groups described within the text. While our current computational methods do not disambiguate these perspectives, human evaluation can provide insights into whose perspectives primarily drive our findings about dehumanization. Thus, we manually divide each measure's most and least dehumanizing paragraphs into three categories based on whose views are most prominent: the author, a person quoted or paraphrased, or a person/group mentioned or described within the text. For each measure, our final sample for human evaluation consists of the 20 most humanizing and 20 most dehumanizing paragraphs within each of the three "viewpoint" categories, yielding 120 paragraphs for each vector-based measure.

MTurk workers read a paragraph and answered a question about the attitudes of the author, person quoted, or people mentioned/described in the text. Table 9 shows four examples, the dehumanization component that they correspond to, whether they are ranked high (most dehumanizing) or low (least dehumanizing),

| Paragraph | Component | Extreme | Viewpoint | Question |
|---|---|---|---|---|
| Some people think that equality can be achieved by offering gays civil unions in lieu of marriage. Civil unions are not a substitute for marriage. Separate rights are never equal rights. | negative evaluation | low | author | How does the author feel about gay people? |
| "I also learned it was possible to be black and gay," Mr. Freeman said. "The first black gay I met, I didn't believe it. I thought you could only be a member of one oppressed minority." | denial of agency | high | person quoted | To what extent does Mr. Freeman think that gay people are able to control their own actions and decisions? |
| In a speech exceptional for its deep emotion and sharp message, Ms. Fisher implicitly rebuked those in her party who have regarded the sickness as a self-inflicted plague earned by immoral behavior – homosexual sex or intravenous drug abuse. | moral disgust | high | person mentioned | To what extent does Ms. Fisher's party consider gay people to be disgusting or repulsive? |
| The Supreme Court on Tuesday was deeply divided over one of the great civil rights issues of the age, same-sex marriage. But Justice Anthony M. Kennedy, whose vote is probably crucial, gave gay rights advocates reasons for optimism based on the tone and substance of his questions. | vermin | low | person mentioned | Vermin are animals that carry disease or cause other problems for humans. Examples include rats and cockroaches. To what extent does [the author] consider gay people to be vermin-like? |

Table 9: Examples of four paragraphs annotated by MTurk workers, one for each dehumanization component. *Extreme* refers to whether the paragraph is ranked as the most dehumanizing (*high*) or least dehumanizing (*low*) for each measure. *Viewpoint* refers to whose perspective workers are asked to reason about. The question that MTurk workers answer is modified based on both the dehumanization component and the viewpoint.

the most prominent viewpoint, and the exact question that workers answered. The question depends on which dehumanization component's measure is being evaluated. In addition, we include the actual name of people quoted or mentioned in order to simplify the task. Each question is answered with a 5-point Likert scale with endpoints specified in the task. For the *negative evaluation* and *denial of agency* questions, 1 is the most dehumanizing option and 5 is the most humanizing option, but the opposite is the case for *vermin* and *moral disgust*. As a postprocessing step, we reverse the scale for the latter so higher values always correspond to more humanizing views.

Three MTurk workers completed each task. Workers were located in the United States, already completed at least 1000 MTurk tasks, and have an approval rate of at least 98%. Each task took approximately 20-25 seconds and workers were compensated $0.05. To avoid confusion with multiple question formulations, we published the tasks for each dehumanization component separately.

## 6.3 Human evaluation results

The results from the MTurk study, shown in Figure 6, largely support our use of vector-based measures. Paragraphs with the highest *neighbor valence* were judged to hold more positive evaluations of gay people ($p < 0.0001$). Paragraphs whose embeddings are nearest to the *Moral Disgust* concept vector are judged to express stronger views of gay people as "disgusting" or "repulsive" compared to the furthest paragraphs ($p < 0.0001$). Similarly, paragraphs nearest to *Vermin* concept consider gay people to be more vermin-like than the paragraphs furthest away ($p < 0.0001$).

The only component that does not follow these expected results is *denial of agency*, where paragraphs with highest and lowest *neighbor dominance* are not judged to be significantly different ($p = 0.19$). This may reflect that using a lexicon for dominance is not a perfect proxy for the more nuanced concept of agency. Another possible explanation is the inherent complexity in measuring *denial of agency*. While the other components are already challenging by requiring an annotator to reason about another person's attitudes towards the target group, assessing *denial of agency* is even more complicated, as it requires an annotator to reason about another person's perceptions of the cognitive capabilities of members of the target group.

The bottom row of Figure 6 separates the results based on whose viewpoint MTurk workers are asked to reason about: the paragraph's author, the people quoted, or the people mentioned in the text. This reveals a strikingly consistent pattern; the difference between the two extremes is largest when workers are asked
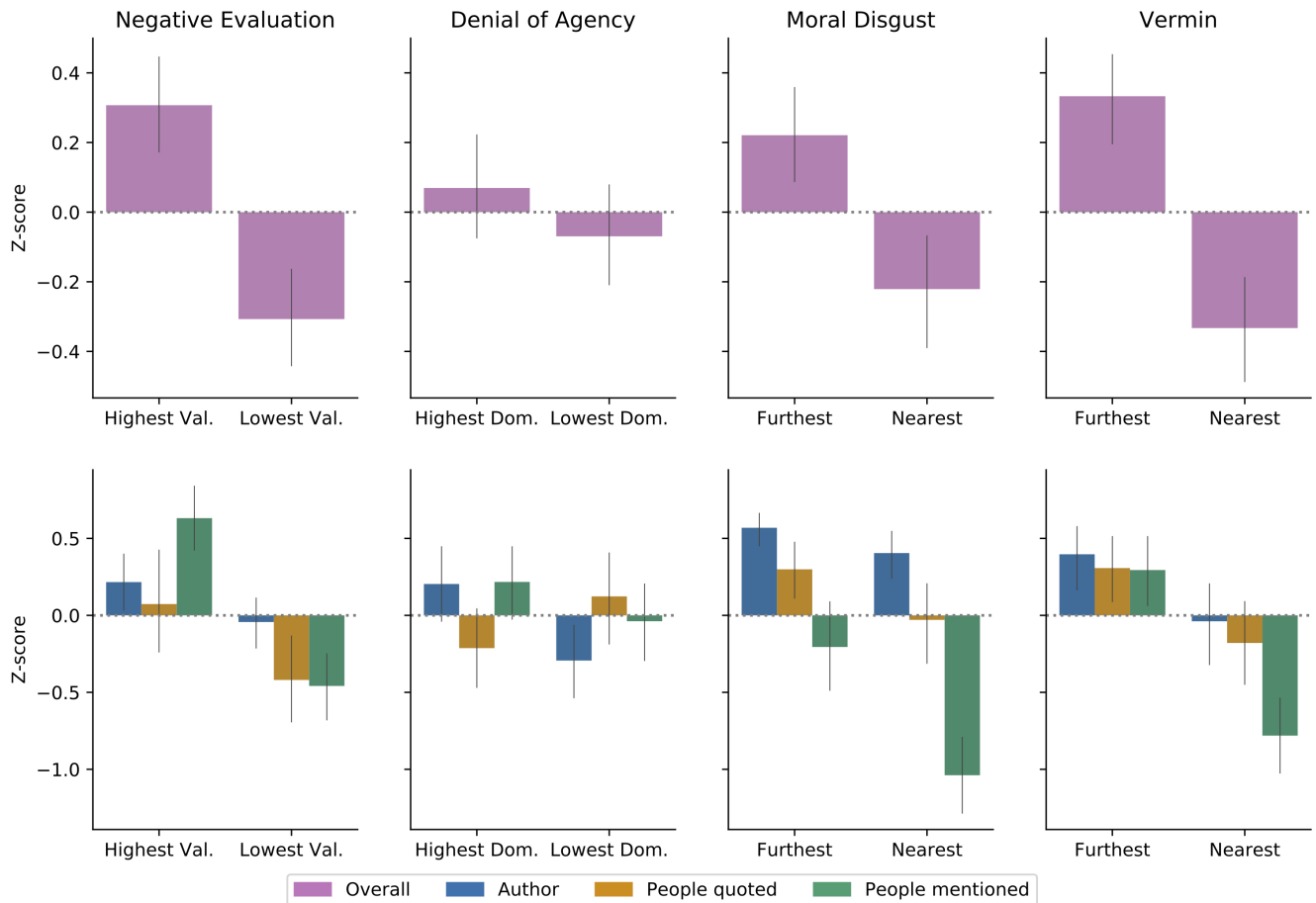
Figure 6: Results from human evaluation of our vector-based methods for quantifying *negative evaluation of the target group*, *denial of agency*, *moral disgust*, and *vermin metaphor*. Higher values are more humanizing (more positive evaluation, greater agency, less association with moral disgust or vermin) and lower values are more dehumanizing. The top row shows overall ratings after z-score normalization for each component and the bottom row separates ratings by the viewpoint workers are asked to judge.

about the *people mentioned*, smallest when asked about *the author*, and in-between when asked about *people quoted*. This suggests that dehumanizing representations of LGBTQ people in the *New York Times* may be most driven by descriptions about other people's attitudes, and to a lesser extent, direct quotes and paraphrases.

## 7 DISCUSSION

Our framework for the computational linguistic analysis of dehumanization involves identifying major dimensions of dehumanization from social psychology literature, proposing linguistic correlates for each dimension, and developing robust and interpretable computational methods to quantify these correlates. We apply this framework to study the dehumanization of LGBTQ people in the *New York Times* from 1986 to 2015. We measure four dimensions of dehumanization: *negative evaluations of a target group*, *denial of agency*, *moral disgust*, and (implicit) invocations of *vermin metaphors*. Overall, we find increasingly humanizing descriptions of LGBTQ people over time. LGBTQ people have become more associated with positive emotional language, suggesting that *negative evaluations of the target group* have diminished. LGBTQ people have become more associated with higher-dominance words, suggesting decreasing *denial of agency*, although this finding was not replicated with the verb-centric "Connotation Frames" measurement. Furthermore, labels for LGBTQ people have moved further away from the concepts of *moral disgust* and *vermin* within distributional semantic representations, suggesting that harmful associations between LGBTQ people and these dehumanizing concepts have weakened over time.

Despite these trends, the labels *gay* and *homosexual* exhibit strikingly different patterns. *Homosexual* is associated with more negative language than *gay*, suggesting more negative evaluations of people described as *homosexual* than *gay*. *Homosexual* is also associated with greater *denial of agency*, and has stronger connections to *moral disgust* and *vermin* than *gay*. Unlike for other LGBTQ labels, discussions of *homosexual* people have not become more humanizing over time, and several techniques even suggest that *homosexual* has become used in more dehumanizing contexts since 2000. Through its repeated use in these contexts, the use of the word *homosexual* appears to have emerged as an index of more dehumanizing attitudes towards LGBTQ people than other labels. Despite the denotational similarity between *homosexual* and *gay*, our computational techniques tracks the stark divergence in social meanings.

We restrict our analysis to the lexical level for ease of interpretability, and leveraged a diverse array of existing resources, including the NRC VAD lexicon (Mohammad, 2018), Connotation Frames lexicons (Rashkin et al., 2016; Sap et al., 2017), and the Moral Foundations Dictionary (Graham et al., 2009). For *negative evaluations of a target group* and *denial of agency*, we propose multiple different techniques that vary in accuracy and interpretability. Word-counting methods are often inaccurate due to their simplicity but their results are easily interpretable, while embedding-based methods suffer the opposite problem. Carefully considering the tradeoff between model quality and interpretability is especially important in work that seeks to characterize complex and sensitive social phenomena such as dehumanization.

## 7.1 Limitations and Future Work

As the first attempt to computationally analyze dehumanization, this work has many limitations. While we demonstrate how the proposed techniques capture linguistic signals of dehumanization, our qualitative and quantitative evaluation suggest that the findings may be driven more by events and attitudes of people described in the text rather than the journalists' own views. An exciting area of future work could involve developing more sophisticated methods to disambiguate the writer's attitudes, attitudes of people mentioned or quoted, and events, while recognizing that each of these could contribute to the overall representation of marginalized groups in the media. In addition, the present work uses word2vec since all known affective lexicons are type-level, but contextualized embedding-based methods have great potential for more nuanced analyses of dehumanizing language by leveraging token-level representations (Devlin et al., 2018; Peters et al., 2018).

Our framework could be expanded to include more insights from dehumanization theory. Beyond the four components discussed in this article, social psychology research has identified other cognitive processes that contribute to dehumanization, including *psychological distancing*, *essentialism* (the perception that the target group has some essence that makes them categorically and fundamentally different), and *denial of subjectivity* (neglect of a target group member's personal feelings and experiences) (Rothbart and Taylor, 1992; Nussbaum, 1999; Graf et al., 2013; Haslam and Stratemeyer, 2016). Scholars also differentiate between two forms of dehumanization, *animalistic* (likening humans to animals) and *mechanistic* (likening humans to inanimate objects or machines), which may differ substantially in their linguistic expressions (Haslam, 2006).

For simplicity and ease of interpretation, we quantify lexical cues of dehumanization. However, our understanding of dehumanizing language would be enriched by considering linguistic features beyond the lexicon. For example, Acton (2014) has shown that definite plurals in English (e.g. *the gays*) have a unique social and pragmatic effect compared to bare plurals (e.g. *gays*) by packaging individual entities into one monolith and setting this group apart from the speaker. Indexing a speaker's non-membership in the group being discussed creates social distance between the speaker and group (Acton, 2014), which makes it likely that using definite plurals to label marginalized social groups plays an important role in dehumanization. Similarly, examining non-lexical signals could help us capture elements of dehumanization not easily identifiable with lexical resources alone. For example, a group label's word class (e.g. *gay* as a noun or adjective) may have implications for *essentialism*, as adjectives simply name attributes to some entity, while nouns explicitly state the entity's category membership and encapsulates other stereotypes associated with that category (Wierzbicka, 1986; Hall and Moore, 1997; Graf et al., 2013; Palmer et al., 2017). We furthermore believe that incorporating discourse-level analysis, such as examining the role of direct quotes in an article and who is being quoted, could provide informative insights that could address some limitations discussed earlier.

We support our proposed framework with a case study of LGBTQ representation in the *New York Times*. This case study is limited as an analysis of the dehumanization of LGBTQ people in the media. We only investigate one data source, which does not capture the entirety of media discourse about LGBTQ people. Furthermore, we only study newspaper articles written in (Standard) American English. Future work could focus on cross-linguistic comparisons of dehumanizing language and assess how well our measures generalize to other languages. Finally, the case study focuses on the labels *gay* and *homosexual* due to data availability. As a consequence, we have less understanding about the differences and changes in representations of LGBTQ people who do not identify with these labels.

The primary aim of this paper is to develop a computational framework for analyzing dehumanizing language towards targeted groups. While our in-depth case study focuses on one particular social group, this framework can be generalized to study dehumanization across a wide variety of social groups, and this could be a fruitful area of future work. For example, Asians have faced increased prejudice and dehumanization since the beginning of the COVID-19 pandemic (Van Bavel et al., 2020; Vidgen et al., 2020; Ziems et al., 2020). Our framework could be applied to understand who dehumanizes these populations in both news and social media, and how the degree of dehumanization changes over time or varies by region. This framework could provide a nuanced view into the shifting nature of dehumanization towards Asians. For example, the "Asians are good at math" stereotype may have led dehumanization via *denial of agency* or *denial of subjectivity* (Shah, 2019). However, stereotypes of Asians as COVID-19 carriers may have made *moral disgust* and *associations with vermin* more salient mechanisms of dehumanization. In our case study, we use computational measures of dehumanizing language to show how the terms *gay* and *homosexual* have diverged in meaning. This method of demonstrating how denotationally similar items differ in connotation can also generalize to other issues and social groups. For example, we may expect labeling COVID-19 as the *Wuhan Virus* or *Chinese Virus* may be associated with greater dehumanization of Asians than the names *COVID-19* or *SARS-CoV-2* (Van Bavel et al., 2020; Xu and Liu, 2020).

## 7.2 Ethical Implications

We hope to draw attention to issues of media representation of marginalized groups and to eventually help make the online world a safer and kinder place. An important part of this mission is to acknowledge the ethical implications and potential issues of our own work.

The methods that we use to quantify dehumanization are themselves biased and potentially harmful. For example, we show in Section 5.2.1.1 that the lexicon used to measure valence contains its own anti-LGBTQ biases by considering LGBTQ group labels to be primarily negative/unpleasant. We also train word2vec models on *New York Times* data, which presents biases. Though models trained on biased data are typically concerning due to harmful downstream effects (Bolukbasi et al., 2016), we leverage this data as a resource for uncovering human biases and understanding *how* biases emerge in the media.

Another concern of this work in our computational methods to represent human beings. Representing people as sequences of numbers (especially in our vector-based experiments) is inherently dehumanizing. While we hope that this work will humanize and empower marginalized groups, we acknowledge that it can also have effect of perpetuating their dehumanization.

Other ethical implications of this project appear within our case study. We do not include LGBTQ labels such as *queer* or *trans*, which often had different meanings and were found in unrelated contexts in earlier years. Furthermore, our analysis uses an aggregated representation for LGBTQ people, which unintentionally minimizes the vast diversity of social identities encompassed within this umbrella. We highlight striking temporal changes and differences between *gay* and *homosexual*, which were chosen because these labels were well-represented in all years. However, emphasizing these labels at the expense of others may contribute to the erasure of people who are marginalized even within LGBTQ communities.

## 8 CONCLUSION

This work is the first known computational linguistic study of dehumanization, and provides contributions to multiple fields. The proposed framework and techniques to quantify salient components of dehumanization can shed light on linguistic variation and change in discourses surrounding marginalized groups. Furthermore, these tools for large-scale analysis have potential to complement smaller-scale psychological

studies of dehumanization. Finally, this work has implications for automatically detecting and understanding media bias and abusive language online.

## 9 DATA AVAILABILITY STATEMENT

We obtained the *New York Times* data through direct communication with one of the creators of the dataset, Ethan Fast (Fast and Horvitz, 2016).

## 10 AUTHOR CONTRIBUTIONS

JM, YT, and DJ collaborated on the conception and design of the study. JM prepared the data, conducted the case study analysis, and conducted statistical analysis. A first draft of the paper was written by JM. JM, YT, and DJ all read and revised the manuscript, and all authors approved the submitted version.

## 11 CONFLICT OF INTEREST

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## 12 ACKNOWLEDGEMENTS

We would like to thank Rob Voigt, Chris Potts, Penny Eckert, David Jurgens, Ceren Budak, Anna-Marie Sprenger, and Ismael Mendoza for their insightful contributions and feedback. We would also like to thank our audience at NWAV-46, who provided helpful feedback on an earlier version of this project. Finally, we thank our anonymous reviewers for their thoughtful and valuable comments.

## REFERENCES

Acton, E. K. (2014). *Pragmatics and the social meaning of determiners*. Ph.D. thesis, Stanford University

Alorainy, W., Burnap, P., Liu, H., and Williams, M. L. (2019). "the enemy among us" detecting cyber hate speech with threats-based othering language embeddings. *ACM Transactions on the Web (TWEB)* 13, 1–26

Arora, S., Liang, Y., and Ma, T. (2017). A simple but tough-to-beat baseline for sentence embeddings

Bar-Tal, D. (1990). Causes and consequences of delegitimization: Models of conflict and ethnocentrism. *Journal of Social issues* 46, 65–81

Barnhurst, K. G. and Mutz, D. (1997). American journalism and the decline in event-centered reporting. *Journal of Communication* 47, 27–53

Baumer, E., Elovic, E., Qin, Y., Polletta, F., and Gay, G. (2015). Testing and comparing computational approaches for identifying the language of framing in political news. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 1472–1482

Blodgett, S. L. and O'Connor, B. (2017). Racial disparity in natural language processing: A case study of social media african-american english. *arXiv preprint arXiv:1707.00061*

Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in Neural Information Processing Systems*. 4349–4357

Boydstun, A. E., Card, D., Gross, J., Resnick, P., and Smith, N. A. (2014). Tracking the development of media frames within and across policy issues

Boydstun, A. E., Gross, J. H., Resnik, P., and Smith, N. A. (2013). Identifying media frames and frame dynamics within and across policy issues. In *New Directions in Analyzing Text as Data Workshop, London*

Breitfeller, L., Ahn, E., Jurgens, D., and Tsvetkov, Y. (2019). Finding microaggressions in the wild: A case for locating elusive phenomena in social media posts. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (Hong Kong, China: Association for Computational Linguistics), 1664–1674. doi:10.18653/v1/D19-1176

Buckels, E. E. and Trapnell, P. D. (2013). Disgust facilitates outgroup dehumanization. *Group Processes & Intergroup Relations* 16, 771–780

Burnap, P. and Williams, M. L. (2016). Us and them: identifying cyber hate on twitter across multiple protected characteristics. *EPJ Data science* 5, 11

Caliskan, A., Bryson, J. J., and Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science* 356, 183–186

Card, D., Boydstun, A., Gross, J. H., Resnik, P., and Smith, N. A. (2015). The media frames corpus: Annotations of frames across issues. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. 438–444

Demszky, D., Garg, N., Voigt, R., Zou, J., Shapiro, J., Gentzkow, M., et al. (2019). Analyzing polarization in social media: Method and application to tweets on 21 mass shootings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 2970–3005

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*

Dinakar, K., Jones, B., Havasi, C., Lieberman, H., and Picard, R. (2012). Common sense reasoning for detection, prevention, and mitigation of cyberbullying. *ACM Transactions on Interactive Intelligent Systems (TiiS)* 2, 18:1–18:30

Dinu, G., Lazaridou, A., and Baroni, M. (2014). Improving zero-shot learning by mitigating the hubness problem. *arXiv preprint arXiv:1412.6568*

ElSherief, M., Kulkarni, V., Nguyen, D., Wang, W. Y., and Belding, E. (2018). Hate lingo: A target-based linguistic analysis of hate speech in social media. In *Twelfth International AAAI Conference on Web and Social Media*

Entman, R. M. (1993). Framing: Toward clarification of a fractured paradigm. *Journal of communication* 43, 51–58

Esses, V. M., Medianu, S., and Lawson, A. S. (2013). Uncertainty, threat, and the role of the media in promoting the dehumanization of immigrants and refugees. *Journal of Social Issues* 69, 518–536

Fast, E. and Horvitz, E. (2016). Long-term trends in the public perception of artificial intelligence. *CoRR* abs/1609.04904

Field, A., Bhat, G., and Tsvetkov, Y. (2019). Contextual affective analysis: A case study of people portrayals in online# metoo stories. In *Proceedings of the International AAAI Conference on Web and Social Media*. vol. 13, 158–169

Field, A., Kliger, D., Wintner, S., Pan, J., Jurafsky, D., and Tsvetkov, Y. (2018). Framing and agenda-setting in russian news: a computational analysis of intricate political strategies. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. 3570–3580

Field, A. and Tsvetkov, Y. (2019). Entity-centric contextual affective analysis. In *57th Annual Meeting of the Association for Computational Linguistics (ACL 2019)*

Gallup (2019). Gay and lesbian rights. http://news.gallup.com/poll/1651/gay-lesbian-rights.aspx

Garg, N., Schiebinger, L., Jurafsky, D., and Zou, J. (2018). Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences* 115, E3635–E3644

Garten, J., Boghrati, R., Hoover, J., Johnson, K. M., and Dehghani, M. (2016). Morality between the lines: Detecting moral sentiment in text. In *Proceedings of IJCAI 2016 workshop on Computational Modeling of Attitudes, New York, NY. Retrieved from http://mortezadehghani. net/wp-content/uploads/morality-lines-detecting. pdf*

Gentzkow, M. and Shapiro, J. M. (2010). What drives media slant? evidence from us daily newspapers. *Econometrica* 78, 35–71

Goff, P. A., Eberhardt, J. L., Williams, M. J., and Jackson, M. C. (2008). Not yet human: implicit knowledge, historical dehumanization, and contemporary consequences. *Journal of personality and social psychology* 94, 292

Graf, S., Bilewicz, M., Finell, E., and Geschke, D. (2013). Nouns cut slices: Effects of linguistic forms on intergroup bias. *Journal of Language and Social Psychology* 32, 62–83

Graham, J., Haidt, J., and Nosek, B. A. (2009). Liberals and conservatives rely on different sets of moral foundations. *Journal of personality and social psychology* 96, 1029

Greene, S. and Resnik, P. (2009). More than words: Syntactic packaging and implicit sentiment. In *Proceedings of human language technologies: The 2009 annual conference of the north american*

*chapter of the association for computational linguistics* (Association for Computational Linguistics), 503–511

Haidt, J. and Graham, J. (2007). When morality opposes justice: Conservatives have moral intuitions that liberals may not recognize. *Social Justice Research* 20, 98–116

Hall, D. G. and Moore, C. E. (1997). Red bluebirds and black greenflies: Preschoolers' understanding of the semantics of adjectives and count nouns. *Journal of Experimental Child Psychology* 67, 236–267

Hamilton, W. L., Leskovec, J., and Jurafsky, D. (2016). Diachronic word embeddings reveal statistical laws of semantic change. *CoRR* abs/1605.09096

Harris, L. T. and Fiske, S. T. (2015). Dehumanized perception. *Zeitschrift für Psychologie*

Haslam, N. (2006). Dehumanization: An integrative review. *Personality and social psychology review* 10, 252–264

Haslam, N. and Stratemeyer, M. (2016). Recent research on dehumanization. *Current Opinion in Psychology* 11, 25–29

Hinds, J. (1977). Paragraph structure and pronominalization. *Paper in Linguistics* 10, 77–99

Hodson, G. and Costello, K. (2007). Interpersonal disgust, ideological orientations, and dehumanization as predictors of intergroup attitudes. *Psychological Science* 18, 691–698

Katajamaki, H. and Koskela, M. (2006). The rhetorical structure of editorials in english, swedish and finnish business newspapers. In *Teoksessa Proceedings of the 5th International Aelfe Conference*. 215–19

Kiritchenko, S. and Mohammad, S. (2018). Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*. 43–53

Kteily, N., Bruneau, E., Waytz, A., and Cotterill, S. (2015). The ascent of man: Theoretical and empirical evidence for blatant dehumanization. *Journal of personality and social psychology* 109, 901

Kulkarni, V., Al-Rfou, R., Perozzi, B., and Skiena, S. (2015). Statistically significant detection of linguistic change. In *Proceedings of the 24th International Conference on World Wide Web* (International World Wide Web Conferences Steering Committee), 625–635

Levy, O. and Goldberg, Y. (2014). Dependency-based word embeddings. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 302–308

Manzini, T., Lim, Y. C., Tsvetkov, Y., and Black, A. W. (2019). Black is to criminal as caucasian is to police: Detecting and removing multiclass bias in word embeddings. *arXiv preprint arXiv:1904.04047*

Marshall, S. R. and Shapiro, J. R. (2018). When "scurry" vs. "hurry" makes the difference: Vermin metaphors, disgust, and anti-immigrant attitudes. *Journal of Social Issues* 74, 774–789

Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119

Mohammad, S. M. (2018). Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words. In *Proceedings of The Annual Conference of the Association for Computational Linguistics (ACL)* (Melbourne, Australia)

Monroe, B. L., Colaresi, M. P., and Quinn, K. M. (2008). Fightin' words: Lexical feature selection and evaluation for identifying the content of political conflict. *Political Analysis* 16, 372–403

Niculae, V., Suen, C., Zhang, J., Danescu-Niculescu-Mizil, C., and Leskovec, J. (2015). Quotus: The structure of political media coverage as revealed by quoting patterns. In *Proceedings of the 24th International Conference on World Wide Web* (International World Wide Web Conferences Steering Committee), 798–808

Nussbaum, M. C. (1999). *Sex and social justice* (Oxford University Press)

Opotow, S. (1990). Moral exclusion and injustice: An introduction. *Journal of social issues* 46, 1–20

Osgood, C. E., Suci, G. J., and Tannenbaum, P. H. (1957). *The measurement of meaning*. 47 (University of Illinois press)

Ott, B. L. and Aoki, E. (2002). The politics of negotiating public tragedy: Media framing of the matthew shepard murder. *Rhetoric & Public Affairs* 5, 483–505

Palmer, A., Robinson, M., and Phillips, K. K. (2017). Illegal is not a noun: Linguistic form for detection of pejorative nominalizations. In *Proceedings of the First Workshop on Abusive Language Online*. 91–100

Pennebaker, J. W., Francis, M. E., and Booth, R. J. (2001). Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates* 71, 2001

Peters, J. W. (2014). The decline and fall of the 'h' word. The New York Times

Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., et al. (2018). Deep contextualized word representations. In *Proc. of NAACL*

Pew Research Center (2017). Changing attitudes on gay marriage. http://www.pewforum.org/fact-sheet/changing-attitudes-on-gay-marriage/

Pryzant, R., Martinez, R. D., Dass, N., Kurohashi, S., Jurafsky, D., and Yang, D. (2019). Automatically neutralizing subjective bias in text. *arXiv preprint arXiv:1911.09709*

Rashkin, H., Singh, S., and Choi, Y. (2016). Connotation frames: A data-driven investigation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 311–321

Recasens, M., Danescu-Niculescu-Mizil, C., and Jurafsky, D. (2013). Linguistic models for analyzing and detecting biased language. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 1650–1659

Rothbart, M. and Taylor, M. (1992). Category labels and social reality: Do we view social categories as natural kinds?

Rudinger, R., Naradowsky, J., Leonard, B., and Van Durme, B. (2018). Gender bias in coreference resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. 8–14

Russell, J. A. (1980). A circumplex model of affect. *Journal of personality and social psychology* 39, 1161

Sap, M., Gabriel, S., Qin, L., Jurafsky, D., Smith, N. A., and Choi, Y. (2019). Social bias frames: Reasoning about social and power implications of language. *arXiv preprint arXiv:1911.03891*

Sap, M., Prasettio, M. C., Holtzman, A., Rashkin, H., and Choi, Y. (2017). Connotation frames of power and agency in modern films. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2329–2334

Schmidt, A. and Wiegand, M. (2017). A survey on hate speech detection using natural language processing. In *Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media*. 1–10

Shah, N. (2019). "asians are good at math" is not a compliment: Stem success as a threat to personhood. *Harvard Educational Review* 89, 661–686

Sherman, G. D. and Haidt, J. (2011). Cuteness and disgust: the humanizing and dehumanizing effects of emotion. *Emotion Review* 3, 245–251

Shuman, E. L. (1894). *Steps into journalism: Helps and hints for young writers* (Correspondence School of Journalism)

Silva, L., Mondal, M., Correa, D., Benevenuto, F., and Weber, I. (2016). Analyzing the targets of hate in online social media. In *Tenth International AAAI Conference on Web and Social Media* (AAAI), 687–690

Smith, B. A., Murib, Z., Motta, M., Callaghan, T. H., and Theys, M. (2017). "gay" or "homosexual"? the implications of social category labels for the structure of mass attitudes. *American Politics Research* , 1532673X17706560

Soller, K. (2018). Six times journalists on the paper's history of covering aids and gay issues. The New York Times

Steuter, E. and Wills, D. (2010). 'the vermin have struck again': Dehumanizing the enemy in post 9/11 media representations. *Media, War & Conflict* 3, 152–167

Sun, T., Gaut, A., Tang, S., Huang, Y., ElSherief, M., Zhao, J., et al. (2019). Mitigating gender bias in natural language processing: Literature review. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 1630–1640

Tipler, C. and Ruscher, J. B. (2014). Agency's role in dehumanization: Non-human metaphors of out-groups. *Social and Personality Psychology Compass* 8, 214–228

Tsur, O., Calacci, D., and Lazer, D. (2015). A frame of mind: Using statistical models for detection of framing and agenda setting campaigns. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 1629–1638

Tsvetkov, Y., Boytsov, L., Gershman, A., Nyberg, E., and Dyer, C. (2014). Metaphor detection with cross-lingual model transfer. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 248–258

Van Bavel, J. J., Baicker, K., Boggio, P. S., Capraro, V., Cichocka, A., Cikara, M., et al. (2020). Using social and behavioural science to support covid-19 pandemic response. *Nature Human Behaviour* , 1–12

Vidgen, B., Botelho, A., Broniatowski, D., Guest, E., Hall, M., Margetts, H., et al. (2020). Detecting east asian prejudice on social media. *arXiv preprint arXiv:2005.03909*

Voigt, R., Camp, N. P., Prabhakaran, V., Hamilton, W. L., Hetey, R. C., Griffiths, C. M., et al. (2017). Language from police body camera footage shows racial disparities in officer respect. *Proceedings of the National Academy of Sciences* 114, 6521–6526

Wang, Z. and Potts, C. (2019). Talkdown: A corpus for condescension detection in context. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 3702–3710

Wiebe, J., Wilson, T., Bruce, R., Bell, M., and Martin, M. (2004). Learning subjective language. *Computational linguistics* 30, 277–308

Wierzbicka, A. (1986). What's in a noun?(or: How do nouns differ in meaning from adjectives?). *Studies in Language. International Journal sponsored by the Foundation "Foundations of Language"* 10, 353–389

Xu, C. and Liu, M. Y. (2020). Social cost with no political gain: The" chinese virus" effect

Zhao, J., Wang, T., Yatskar, M., Ordonez, V., and Chang, K.-W. (2018). Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. vol. 2

Ziems, C., He, B., Soni, S., and Kumar, S. (2020). Racism is a virus: Anti-asian hate and counterhate in social media during the covid-19 crisis. *arXiv preprint arXiv:2005.12423*

## SUPPLEMENTARY MATERIAL

| Label | Count | Label | Count | Label | Count | Label | Count |
|---|---|---|---|---|---|---|---|
| gay(s) | 96977 | lgbt | 1783 | lgbtq | 129 | agender | 10 |
| lesbian(s) | 20233 | transvestite(s) | 625 | glbt | 68 | aromantic | 5 |
| homosexual(s) | 16638 | tran(s)sexual(s) | 627 | genderqueer | 51 | lgbtqia | 4 |
| transgender(s/ed) | 6066 | asexual | 255 | lgb | 29 | genderfluid | 0 |
| bisexual(s) | 3464 | intersex | 210 | pansexual | 22 | lgbtqqia | 0 |

Table 10: Overall counts for all LGBTQ terms in the New York Times from 1986-2015. Each label includes its morphological and orthographic variants.

## Nearest neighbor valence with different thresholds

Figure 7: Average valence of 25, 50, 100, 250, 500 and 1000 nearest words to vector representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American*, averaged over 10 word2vec models trained on *New York Times* data from each year. The solid lines are Lowess curves for visualization purposes. Words' valence scores are from the NRC VAD Valence Lexicon. For all plots, the shaded bands represent 95% confidence intervals.

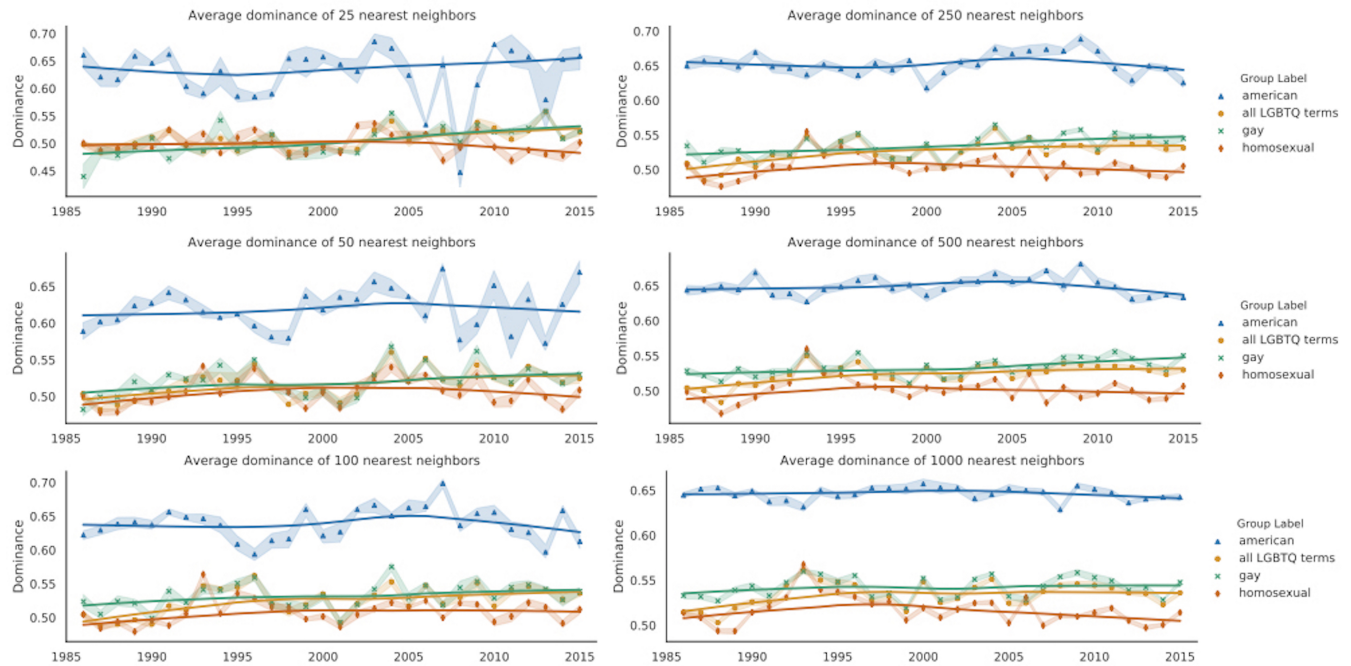## Nearest neighbor dominance with different thresholds



Figure 8: Average dominance of 25, 50, 100, 250, 500 and 1000 nearest words to vector representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American*, averaged over 10 word2vec models trained on *New York Times* data from each year. The solid lines are Lowess curves for visualization purposes. Words' dominance scores are from the NRC VAD Dominance Lexicon. For all plots, the shaded bands represent 95% confidence intervals.

## Valence prediction results

In addition to quantifying the *negative evaluation of a target group* by calculating the average valence score of a group label's vector representation's nearest neighbors according to the NRC VAD lexicon, we also directly induced a valence score from the vector representation itself. We use the zero-centered, normalized, word embeddings created for each year as features in a regression model to directly predict valence (Field et al., 2019). Specifically, we train ridge regression models for each year, where each year's Word2Vec representation for words from the NRC VAD lexicon are features and the lexicon's valence scores are labels. 85% of words from the VAD lexicon were kept as the training set, and the remaining 15% was used as a test set to evaluate performance. We then use this set of regression models to predict a group label's valence from its vector representation.

Figure 9 shows the directly-induced valence score for each set of group labels from the ridge regression fit to the NRC VAD valence lexicon. Because we trained a different Word2Vec model for each year, we trained a different ridge regression model for each year. Over all thirty years, the Pearson correlation between predicted valence and actual valence on the test set ranged from 0.617 to 0.675, and $R^2$ values ranged from 0.423 to 0.451. The predicted scores show similar trends to the average neighbor valence. *Homosexual* has the most negative valence for every year, followed by *gay* and the aggregate over all LGBTQ terms, followed by *American* with the most positive valence. *American* is significantly more positive than all LGBTQ labels over all years (Wilcoxon signed-rank test; $p < 0.0001$), and *gay* is significantly more positive than *homosexual* for every year (Wilcoxon signed-rank test; $p < 0.0001$). Despite the stability in differences between these terms across experiments, the regression analysis suggests different temporal dynamics. Figure 9 shows that the predicted valence for *gay*, *homosexual*, and all LGBTQ terms all increase over time, but *homosexual*'s predicted valence decreases from 2001 to 2015 ($p < 0.01$). This result is consistent with the other findings in this article in illustrating the pejoration of *homosexual* in recent years.
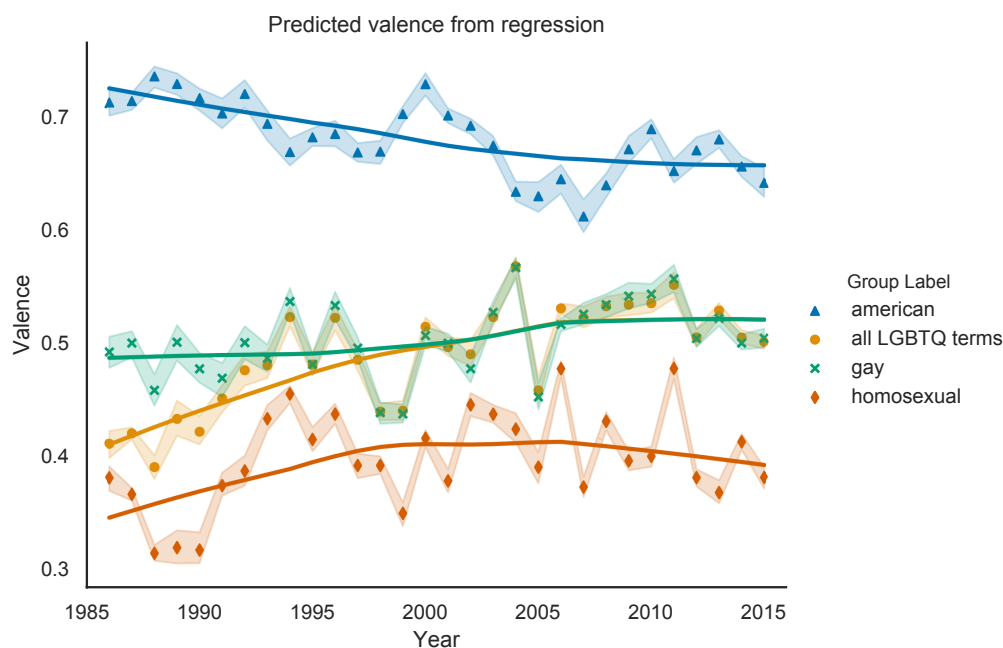
Figure 9: Predicted valence of our representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American* directly induced by fitting ridge regression models to the NRC VAD Valence Lexicon with the lexicon's words' Word2Vec vectors as features for each year. Results are averaged over 10 Word2Vec models trained over each year's data. Shaded bands represent 95% confidence intervals and the solid lines are Lowess curves as visual aids. Higher scores represent more positive valence.

## Agency prediction results

Because we use the NRC VAD dominance lexicon to quantify *denial of agency* in the same way we use the valence lexicon to quantify *negative evaluations of a target group*, we again directly induce scores directly from target group label representation. We train another set of ridge regression models with word embedding features for each year, but now we fit the model to the NRC VAD dominance lexicon's scores rather than the valence scores.

Figure 10 shows the predicted dominance for each group label, which is calculated by fitting ridge regression models to the NRC VAD Dominance Lexicon using the lexicon's words' Word2Vec representations as features for each year. Pearson correlations between predicted and actual dominance scores for all regression models ranged from 0.561 to 0.614 on the test set, and $R^2$ values range from 0.338 to 0.361. Consistent with the average neighbor dominance approach , *American* has significantly greater dominance than any of the other LGBTQ terms (Wilcoxon signed-rank test; $p < 0.0001$). However a Wilcoxon signed-rank test over each year's means shows that there is no significant difference between the terms *gay* and *homosexual* ($p = 0.21$). The predicted dominance of *all LGBTQ terms* and *gay* significantly decrease ($p < 0.0001$), but not in the last 15 years ($p = 0.85$ for *all LGBTQ terms* and $p = 0.51$ for *gay*). *Homosexual* does not significantly change in predicted dominance in either the full 30 years ($p = 0.96$) or in the last 15 years ($p = 0.89$).

Why do *gay* and *homosexual* show such different patterns in directly-induced predicted dominance from the regression than average dominance based on their neighbors' entries in the NRC VAD lexicon? While the average dominance over the nearest neighbors showed significant differences, they were small in magnitude (often corresponding to differences of less than 0.025 points on a scale from 0 to 1). Perhaps because the word2vec features could only predict just over a third of the variance in dominance scores, they were not able to capture subtle semantic distinctions that could characterize differences in dominance scores
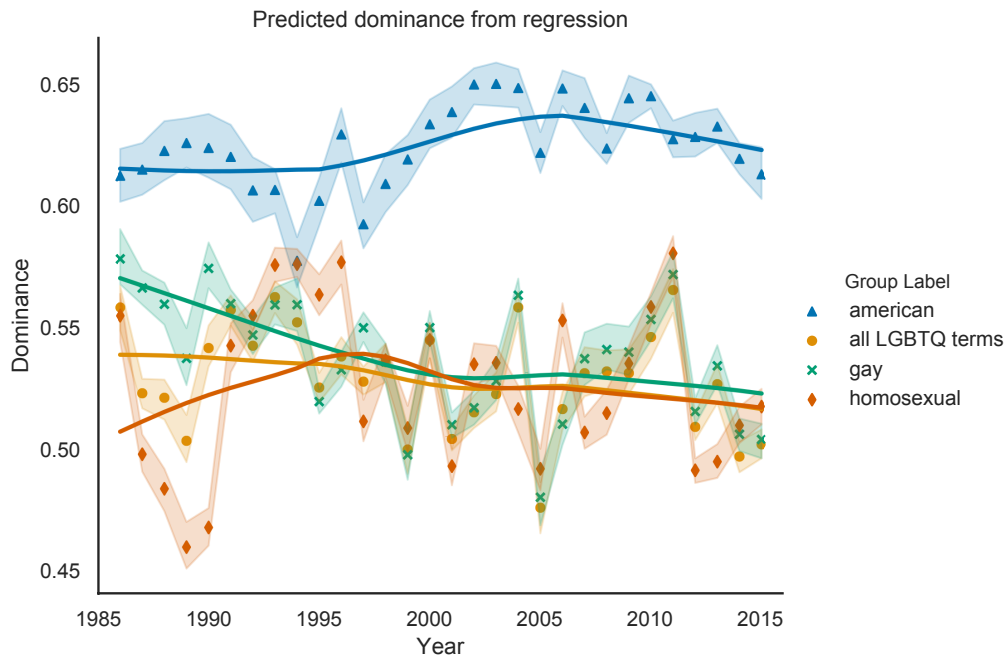
Figure 10: Predicted dominance of our representations of *gay*, *homosexual*, *all LGBTQ terms*, and *American* directly induced by fitting ridge regression models to the NRC VAD Dominance Lexicon with the lexicon's words' Word2Vec representations as features for each year. Results are averaged over 10 Word2Vec models trained over each year, shaded bands represent 95% confidence intervals and the solid lines are Lowess curves for visualization purposes. Higher scores represent greater predicted dominance.