

From Dogwhistles to Bullhorns: Unveiling Coded Rhetoric with Language Models

Warning: content in this paper may be upsetting or offensive to some readers

Julia Mendelsohn

University of Michigan
juliam@umich.edu

Yejin Choi

Allen Institute for Artificial Intelligence
yejinc@allenai.org

Ronan Le Bras

Allen Institute for Artificial Intelligence
ronanlb@allenai.org

Maarten Sap

Carnegie Mellon University
maartensap@cmu.edu

Abstract

Dogwhistles are coded expressions that simultaneously convey one meaning to a broad audience and a second one, often hateful or provocative, to a narrow in-group; they are deployed to evade both political repercussions and algorithmic content moderation. For example, in the sentence “we need to end the *cosmopolitan* experiment,” the word “*cosmopolitan*” likely means “wordly” to many, but secretly means “Jewish” to select few. We present the first large-scale computational investigation of dogwhistles. We develop a typology of dogwhistles, curate the largest-to-date glossary of over 300 dogwhistles with rich contextual information and examples, and analyze their usage in historical U.S. politicians’ speeches. We then assess whether a large language model (GPT-3) can identify dogwhistles and their meanings, and find that GPT-3’s performance varies widely across types of dogwhistles and targeted groups. Finally, we show that harmful content containing dogwhistles avoids toxicity detection, highlighting online risks of such coded language. This work sheds light on the theoretical and applied importance of dogwhistles in both NLP and computational social science, and provides resources for future research in modeling dogwhistles and mitigating their online harms.

1 Introduction

The cosmopolitan elite look down on the common affections that once bound this nation together: things like place and national feeling and religious faith... The cosmopolitan agenda has driven both Left and Right... It’s time we ended the cosmopolitan experiment and recovered the promise of the republic.

—Josh Hawley (R-MO), 2019

We have got this tailspin of culture, in our inner cities in particular, of men not working and just

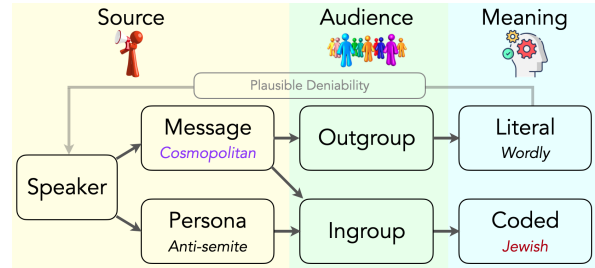


Figure 1: Schematic of how dogwhistles work, based on Henderson and McCready (2018) with the example of *cosmopolitan*. First, a speaker simultaneously communicates the dogwhistle message and their persona (identity). The in-group recovers both the message content and speaker persona, enabling them to arrive at the coded meaning (e.g. *Jewish*). The out-group only recognizes the message’s content and thus interprets it literally. This literal meaning also provides the speaker with plausible deniability; if confronted, the speaker can claim that they solely intended the literal meaning.

generations of men not even thinking about working or learning the value and the culture or work.

—Paul Ryan (R-WI), 2014

Cosmopolitan and *inner city* are examples of dogwhistles, expressions that “send one message to an out-group and a second (often taboo, controversial, or inflammatory) message to an in-group” (Henderson and McCready, 2018). Many listeners would believe that Hawley is simply criticizing well-traveled or worldly people, but others recognize it as an attack on the Jewish people. Similarly, many assume that Ryan is discussing issues within a geographic location, but others hear a pernicious stereotype of Black men as lazy. Crucially, Hawley and Ryan can avoid alienating the out-group by maintaining *plausible deniability*: they never explicitly say “Jewish” or “Black”, so they can reject accusations of racism (Haney-López, 2014).

Because dogwhistles can bolster support for par-

ticular policies or politicians among the in-group while avoiding social or political backlash from the out-group, they are a powerful mechanism of political influence (Mendelberg, 2001; Goodin and Saward, 2005). For example, racist dogwhistles such as *states' rights* and *law and order* were part of the post-Civil Rights Republican Southern Strategy to appeal to white Southerners, a historically Democratic bloc (Haney-López, 2014). Despite polarization and technology that enables message targeting to different audiences, dogwhistles are still widely used by politicians (Haney-López, 2014; Tilley et al., 2020) and civilians in online conversations (Bhat and Klein, 2020; Åkerlund, 2021).

Beyond political science, research on dogwhistles is urgent and essential for NLP, but they remain a challenge to study. Dogwhistles are actively and intentionally deployed to evade automated content moderation, especially hate speech detection systems (Magu et al., 2017). They may also have harmful unseen impacts in other NLP systems by infiltrating data used for pretraining language models. However, researchers face many difficulties. First, unless they are a part of the in-group, researchers may be completely unaware of a dogwhistle's existence. Second, dogwhistles' meanings cannot be determined by form alone, unlike most overt hateful or toxic language. Rather, their interpretation relies on complex interplay of different factors (context, personae, content, audience identities, etc.; Khoo, 2017; Henderson and McCready, 2018, 2019; Lee and Kosse, 2020), as illustrated in Figure 1. Third, since their power is derived from the differences between in-group and out-group interpretations, dogwhistles continuously evolve in order to avoid being noticed by the out-group.

We establish foundations for large-scale computational study of dogwhistles by developing theory, providing resources, and empirically analyzing dogwhistles in several NLP systems. Prior work largely focuses on underlying mechanisms or political effects of dogwhistle communication (Albertson, 2015; Henderson and McCready, 2018) and typically consider a very small number of dogwhistles (often just one). To aid larger-scale efforts, we first create a new taxonomy that highlights both the systematicity and wide variation in kinds of dogwhistles (§2.1). This taxonomy characterizes dogwhistles based on their covert meanings, style and register, and the personae signaled by their users. We then compile a glossary of 340 dogwhis-

tles, each of which is labeled with our taxonomy, rich contextual information, explanations, and real-world examples with source links (§2.2-2.3). As this glossary is the first of its kind, we highlight its value with a case study of racial dogwhistles in historical U.S. Congressional Speeches (§3).

We then apply our taxonomy and glossary to investigate how dogwhistles interact with existing NLP systems (§4). Specifically, we evaluate the ability of large language models (i.e. GPT-3) to retrieve potential dogwhistles and identify their covert meanings. We find that GPT-3 has a limited capacity to recognize dogwhistles, and performance varies widely based on taxonomic features and prompt constructions; for example, GPT-3 is much worse at recognizing transphobic dogwhistles than racist ones. Finally, we show that hateful messages with standard group labels (e.g. *Jewish*) replaced with dogwhistles (e.g. *cosmopolitan*) are consistently rated as far less toxic by a commercially deployed toxicity detection system (PerspectiveAPI), and such vulnerabilities can exacerbate online harms against marginalized groups (§5).

This work highlights the significance of dogwhistles for NLP and computational social science, and offers resources for further research in recognizing dogwhistles and reducing their harmful impacts.

2 Curating a dogwhistle glossary

2.1 Taxonomy

Based on prior work and our own investigations, we craft a new taxonomy (Figure 2). We categorize dogwhistles by **register**, **type**, and **persona**.

Register We label all dogwhistles as either part of a **formal/offline** or **informal/online** register. Formal/offline dogwhistles originated in offline contexts or are likely to appear in statements by mainstream political elites (e.g. *family values*). The informal/online register includes dogwhistles that originated on the internet and are unlikely to be used in political speech (e.g. *cuckservative*).

Type I Henderson and McCready (2018) distinguish dogwhistles into two types: **Type I** dogwhistles covertly signal the speaker's persona but do not alter the implicatures of the message itself, while **Type II** dogwhistles additionally alter the message's implied meaning. We extend this typology to highlight the wide variety of dogwhistles, which has important consequences for building a

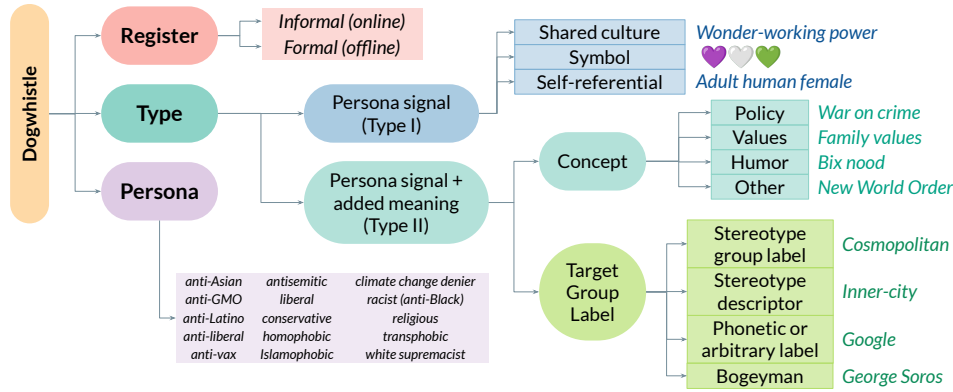


Figure 2: Visual hierarchical representation of our dogwhistle taxonomy along with examples of each type.

theory of dogwhistles as well as future computational modeling. We identify three subcategories of “only persona-signaling” (Type I) dogwhistles: **symbols** (including emojis, abbreviations, and imagery), **self-referential** terms for members of the in-group, and dogwhistles that require specialized knowledge from a **shared in-group culture**.

Type II Dogwhistles with an “added message meaning” (Type II) tend to fall into two subcategories: they name a **concept** or serve as a substitute for a **target group label**. We further divide concepts into **policies** (titles for initiatives with covert implications, such as *law and order*), **values** that the in-group purports to uphold, expressions whose covert meanings are grounded in in-group **humor**, and **other concepts**, which are often coded names for entities that are not group labels (e.g. the *New World Order* conspiracy theory is antisemitic but does not name or describe Jewish people).

Dogwhistles serve as **target group labels** in three ways. Many are stereotype-based, whose interpretations rely on pre-existing associations between the dogwhistle and target group; we separate these into **stereotype-based target group labels**, which directly name the target group (e.g. *cosmopolitan*), while **stereotype-based descriptors** are less direct but still refer to the target group (e.g. *inner-city*). Others have an **arbitrary or phonetic** relationship to the group label; these are commonly used to evade content moderation, such as “Operation Google” terms invented by white supremacists on 4chan to replace various slurs (Magu et al., 2017; Bhat and Klein, 2020). The final subcategory, **Bogeyman**, includes names of people or institutions taken to represent the target group (e.g. *George Soros* for Jewish or *Willie Horton* for Black).

Persona **Persona** refers to the in-group identity signalled by the dogwhistle. Figure 2 lists some personae, but this is an open class with many potential in-groups. There is considerable overlap in membership of listed in-groups (e.g. white supremacists are often antisemitic), but we label persona based on what is most directly related to the dogwhistle.

2.2 Gathering dogwhistles

We draw from academic literature, media coverage, blogs, and community-sourced wikis about dogwhistles, implicit appeals, and coded language. Since academic literature tends to focus on a small set of examples, we expanded our search to media coverage that identifies dogwhistles in recent political campaigns and speeches (e.g. Burack, 2020) or attempts to expose code words in hateful online communities (e.g. Caffier, 2017). During our search, we found several community-sourced wikis that provided numerous examples of dogwhistles, particularly the RationalWiki “Alt-right glossary”, “TERF glossary”, and “Code word” pages.¹

2.3 Glossary contents

Our glossary contains 340 English-language dogwhistles and over 1,000 surface forms (morphological variants and closely-related terms), mostly from the U.S. context. Each dogwhistle is labeled with its register, type, and signaled persona, an explanation from a linked source, and at least one example with linguistic, speaker, situational, and temporal context included, as well as a link to the example text. Antisemitic, transphobic, and racist (mostly anti-Black but sometimes generally against people of color) dogwhistles are the most common, with over 70 entries for each persona. The glossary

¹rationalwiki.org/wiki/{Alt-right_glossary,TERF_glossary,Code_word}

Dogwhistle	Sex-based rights
In-group meaning	Trans people threaten cis women's rights
Persona	Transphobic
Type	Concept: Value
Register	Formal
Explanation	Many anti-transgender people [claim that] women's "sex-based rights" are somehow being threatened, removed, weakened, eroded, or erased by transgender rights. . . "Sex-based rights", by the plain English meaning of those words, cannot exist in a country that has equality law. . . it's mostly a dog-whistle: a rallying slogan much like "family values" for religious conservatives, which sounds wholesome but is a deniable and slippery code-word for a whole raft of unpleasant bigotry.
Source	Medium post by David Allsopp
Example	<i>When so-called leftists like @lloyd_rm demand that we give up our hard-won sex-based rights, they align themselves squarely with men's rights activists. To both groups, female trauma is white noise, an irrelevance, or else exaggerated or invented.</i>
Context	Tweet by J.K. Rowling on June 28, 2020

Table 1: Example glossary entry for the transphobic dogwhistle *sex-based rights*

includes dogwhistles with other personae, such as homophobic, anti-Latinx, Islamophobic, anti-vax, and religious. Table 1 shows one glossary entry for the transphobic dogwhistle *sex-based rights*. Because dogwhistles continuously evolve, we intend for this resource to be a living glossary and invite the public to submit new entries or examples.

3 Case study: racial dogwhistles in historical U.S. Congressional speeches

We showcase the usefulness of our glossary, with a diachronic case study of racial dogwhistles in politicians' speeches from the U.S. Congressional Record (Gentzkow et al., 2019; Card et al., 2022) to analyze the frequency of speeches containing racist dogwhistles from 1920-2020. For this case study, we simply identify glossary terms based on regular expressions and do not distinguish between covert and literal meanings of the same expressions. We also measure how ideologies of speakers using dogwhistles changed over time using DW-NOMINATE, a scaling procedure that places politicians on a two dimensional map based on roll call voting records, such that ideologically similar politicians are located near each other (Poole and Rosenthal, 1985; Carroll et al., 2009; Lewis et al., 2023). We consider the first dimension of DW-NOMINATE, which corresponds to a liberal-conservative axis (Poole and Rosenthal, 1985).²

²The second dimension captures salient cross-cutting issues, and some argue that this dimension primarily captures race relations (Poole and Rosenthal, 1985). However, the sec-

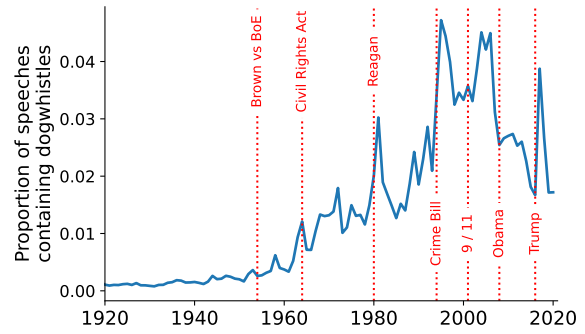


Figure 3: Frequency of speeches containing racial dogwhistles in the U.S. Congressional Record (as a fraction of total speeches) over time. The dotted red vertical lines represent noteworthy years. Use of racial dogwhistles began to increase during the Civil Rights Movement and their frequency continued to rise until the 1990s. Since the 1990s, the frequency of speeches containing dogwhistles has fluctuated but remained at overall high levels compared to earlier years.

Dogwhistle use began to increase during the Civil Rights Era, following the 1954 *Brown vs. Board of Education* Supreme Court decision mandating racial integration of public schools (Fig. 3). This aligns with qualitative accounts of the Southern Strategy: because explicit racism was no longer acceptable, politicians turned to dogwhistles to make the same appeals implicitly (Mendelberg, 2001). Their frequency continued to increase from the 1970s through the 1990s, paralleling Haney-López (2014)'s account of dogwhistles during the Nixon, Reagan, Bush Sr., and Clinton presidencies.

Figure 4 shows how the average ideologies of speakers who use particular dogwhistles (*property rights*, *thug*, *welfare reform*, *hardworking Americans*, and *Willie Horton*) have shifted over time, and reveals interesting insights into the evolution and lifecycle of dogwhistles. Most racial dogwhistles in the U.S. Congressional Speeches have become increasingly associated with more conservative speakers over time. However, the inflection point when speaker ideologies shift varies across dogwhistles, suggesting that they emerged as dogwhistles at different points. For example, *property rights* became increasingly associated with more conservative speakers since the 1960s, while the average ideology of speakers using *welfare reform* did not change until the 1990s.

ond dimension's interpretation is less clear as the vast majority of voting variation is along the first dimension, and is often ignored by political scientists (Bateman and Lapinski, 2016). We thus restrict this case study to the first dimension though future work may opt to consider the second dimension as well.

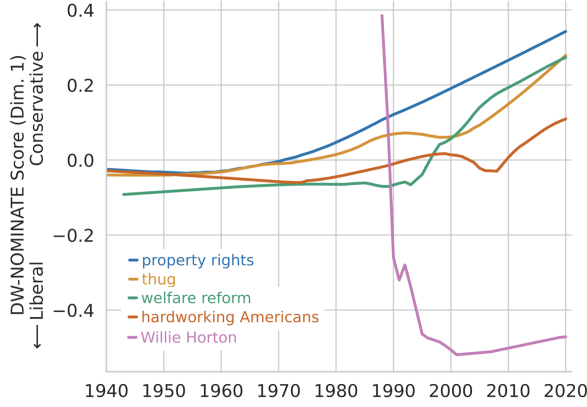


Figure 4: Average ideology score (DW-NOMINATE first dimension) for speakers who used selected dogwhistles over time: *welfare reform* (top left), *thug* (top right), *property rights* (bottom left), and *Willie Horton* (bottom right). Higher values indicate that the dogwhistle’s speakers were more conservative, while lower values indicate that the dogwhistle’s speakers were more liberal. For visualization, trends are Lowess-smoothed.

Willie Horton presents an interesting example. In his 1988 presidential campaign, George Bush ran an ad featuring Willie Horton, a Black man convicted of rape and murder while on prison furlough (Mendelberg, 2001). The ad was so powerful among white voters that it propelled Bush to victory, but shortly afterwards was explicitly called out as racist (Haney-López, 2014). We see this pattern in Figure 4: in 1988, *Willie Horton* was associated with extremely conservative speakers, but quickly became more liberal, and *Willie Horton* no longer functioned as a successful dogwhistle.

4 Recognition of dogwhistles in GPT-3

We conduct two experiments to assess if a large language model, GPT-3 (Brown et al., 2020), can recognize dogwhistles. First, we interrogate whether GPT-3 is able to **identify covert meanings** of dogwhistles from our glossary, an ability that would be instrumental in understanding the breadth of on-line bigotry. Second, we measure GPT-3’s ability to **surface** dogwhistles. Because dogwhistles are often intentionally obscured from researchers, making it impossible to ensure that a manual search is complete or comprehensive. Since GPT-3 is trained on large portions of internet data, it may be able to reveal a more comprehensive and diverse set of dogwhistles than manual efforts.

	No Secret Cue	With Secret Cue
No Definition	.085 [±.025]	.296 [±.021]
With Definition	.437 [±.010]	.543 [±.010]

Table 2: Proportion of dogwhistles with covert meanings identified by GPT-3 when prompted with a dogwhistle definition, a secret cue, both, or neither. Confidence intervals are shown in brackets.

4.1 Identifying covert meanings

To gauge which types of in-group language GPT-3 might have seen during training, we examine if it is able to identify dogwhistles’ covert meanings.

Experimental setup We create prompts that vary in (1) which of the five definitions of dogwhistles is provided, if any, and (2) if the prompt includes a secret cue. For each input term, we construct 12 prompts (6 definitions including none, and each with or without a secret cue). Consider the example below, with a definition, secret cue, dogwhistle and GPT-3’s response. Prompts with no secret clue simply exclude the word *secretly*. We generate responses for each dogwhistle in the glossary and consider multiple forms for some, such as actual emojis as well as descriptions. We thus test 480 variants with 12 prompts for each, leading to 5,760 generations. We use the text-davinci-002 GPT-3 model with max_tokens=256, temperature=0.0 (more deterministic outputs) and num_outputs=1.

A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition. For example, “cosmopolitan” secretly means “Jewish to many anti-Semitic people”.

Results We manually evaluate whether or not each of the 5,760 GPT-3 generations contains the covert meaning for each dogwhistle. 80.3% of dogwhistles had their covert meanings identified in at least one GPT-3 generation. Overall, 56.0% generations contained the correct covert meaning for dogwhistles that are part of the formal/offline register, but just 29.4% for dogwhistles in the informal/online register. We refer readers to Section A.2 and Figure A.2 for more details about register-based variation and examples of dogwhistles for which GPT-3 performed particularly well or poorly.

The specific prompt form strongly impacts GPT-3’s ability to generate covert meanings (Table 2). Without a definition or secret cue, covert meanings

are identified in just 8.5% of generations. Including both a **definition** and **secret cue** improves GPT-3’s performance over 5-fold, with dogwhistles’ covert meanings identified in 54.3% of generations.

We observe wide variation in GPT-3’s ability to identify covert meanings across personae. Among the most represented personae in our glossary (at least 100 generations for each), GPT-3 has the lowest recognition of transphobic dogwhistles, the highest recognition of homophobic and Islamophobic dogwhistles, with antisemitic, white supremacist, and racist dogwhistles in the middle (Appendix Figure A.3). There is also variation in performance by dogwhistle type and the specific definition provided; we refer the reader to Appendix A.2 and A.4 for more details.

4.2 Surfacing dogwhistles

In addition to evaluating if GPT-3 can identify dogwhistles’ covert meanings, we assess GPT-3’s ability to surface dogwhistles in text generation.

Experimental setup We construct a series of prompts that begin with one of five **definitions of dogwhistles** from prior work (Table A.1). The definition is followed by a **question or request for examples**. See Appendix A.1 for more prompting details. In the following example, the definition is marked in **blue**, the request in **purple**, and GPT-3’s response is highlighted in **yellow**.

A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition. What are examples of dogwhistles?

1. “Law and order”
2. “The silent majority”
3. “Family values”
4. “Welfare queens”
5. “Illegal aliens”

Evaluation We use our glossary as a proxy to measure precision and recall of GPT-3’s ability to surface dogwhistles because an exhaustive ground-truth set of dogwhistles does not exist. We calculate recall as the proportion of dogwhistles in our glossary that were also surfaced at least once by GPT-3. For precision, the authors manually inspect candidates appearing in at least 4% of GPT-3 text generations for *generic*, *white supremacist*, *racist*, *antisemitic*, *Islamophobic*, and *transphobic* prompt types. Because our glossary is not exhaus-

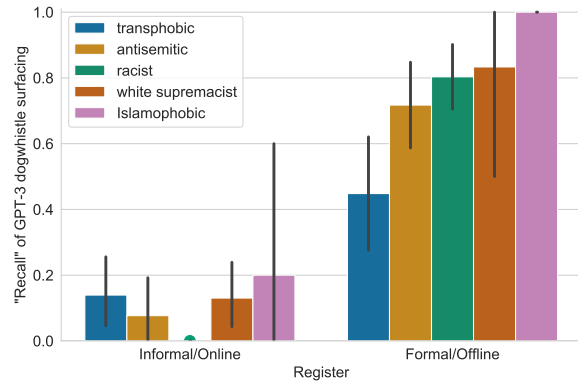


Figure 5: Recall of GPT-3 dogwhistle surfacing separated by persona and register. Across all personae, GPT-3 surfaces under 20% of dogwhistles in the informal/online register. Performance is much higher for the formal/offline register but varies across personae, ranging from 44.8% (transphobic) to 100% (Islamophobic).

tive, this method yields conservative estimates (see Appendix A.1 for more evaluation details).

Precision Results We find that GPT-3 does have the ability to surface dogwhistles when prompted to do so, but caution that such results are imperfect and require manual verification. The most common errors involve explicit mentions of groups in stereotypes or conspiracy theories (*Jews are behind the 9/11 attacks*) or phrases that may accompany dogwhistles but are not dogwhistles themselves (*I’m not racist but...*). Precision in dogwhistle surfacing varies across prompt types; while the average precision over all six prompt types is 66.8%, scores range from just 50% for transphobic dogwhistle prompts to 91.3% for generic prompts (Figure A.1).

Recall Results GPT-3 surfaced 153 of 340 dogwhistles in our glossary (45%). We observe significant differences by register: GPT-3 surfaced 69.4% of *formal/offline* dogwhistles but just 12.9% of *informal/online* dogwhistles. Despite its ability to generate emojis and other symbols, GPT-3 did not surface any symbols or emojis from our glossary except for the antisemitic triple parentheses “((()))”.

Figure 5 shows GPT-3 surfacing recall results by both register and in-group personae. We show results for the five most-frequent personae represented in our glossary. Recall of dogwhistles in the informal/online register is low across the board. For the formal/offline register, recall is considerably higher although it varies widely across personae. As with precision, GPT-3 has the lowest perfor-

Category	Toxicity	Severe Toxicity	Identity Attack
Dogwhistle	.538 [±.006]	.111 [±.004]	.236 [±.005]
Slur	.712 [±.009]	.281 [±.008]	.556 [±.013]
Standard	.758 [±.007]	.326 [±.007]	.732 [±.005]

Table 3: Average Perspective API toxicity, severe toxicity, and identity attack scores for HateCheck template sentences filled in with dogwhistles, standard group labels, or slurs. 95% confidence intervals are in brackets.

mance for transphobic dogwhistles, surfacing just 44.8% of formal/offline transphobic dogwhistles. For formal/offline antisemitic dogwhistles, recall is considerably higher but far from perfect at 71.7%. GPT-3 has 80.3% and 83.3% recall of racist and white supremacist dogwhistles, respectively, and full 100% recall of Islamophobic dogwhistles.

5 Dogwhistles and toxicity detection

Beyond evaluating language models’ ability to recognize dogwhistles, we seek to understand how dogwhistles affect the decisions that NLP systems make, and how this has downstream implications for content moderation and online safety. We begin to address this with a study of how dogwhistles are handled by a widely-deployed toxic language detection system, Perspective API.

Experimental setup We consider 237 hateful sentence templates from HateCheck (Röttger et al., 2021), a test suite for hate speech detection, that contain placeholders for identity terms (group referents) in either adjectival, singular nominal, or plural nominal forms. These placeholders were filled by a standard group label, a slur, or a dogwhistle in the corresponding grammatical form requested by the template. For this experiment, we consider racist (mostly anti-Black), antisemitic, and transphobic terms, as these personae are the most common in our glossary (see Tables A.6 and A.5 for group label terms used and a sample of sentence templates, respectively). We feed our resulting 7,665 sentences to Perspective API to get scores for *toxicity*, *severe toxicity*, and *identity attack*.

Results Hateful sentences are rated as less toxic, severely toxic, and identity-attacking when dogwhistles are used instead of standard group labels or slurs (Table 3). This pattern holds for all three personae (Appendix Figure A.5).

Interestingly, mean toxicity scores for slurs are

lower than for standard group labels, especially for antisemitic slurs. We observe relatively wide variation in Perspective API’s ratings depending on the specific choice of slur. For example, sentences containing the *N-word* are almost always rated as more toxic than the same sentences containing *Black* or *Black people*. Lower toxicity ratings for other slurs, such as the highly derogatory antisemitic *K-word*³ may be because, similar to dogwhistles, Perspective API does not recognize that these terms refer to identity groups. However, deeper analysis of slurs is outside the scope of the current work.

6 Discussion & Conclusion

We lay the groundwork for NLP and computational social science research on dogwhistles by developing a new taxonomy and glossary with rich contextual information and examples. We demonstrate our glossary’s utility in a case study of historical U.S. Congressional speeches, where our quantitative analysis aligns closely with historical accounts. We further use our glossary to show that GPT-3 has some, but limited, ability to retrieve dogwhistles and recognize their covert meanings. Finally, we verify that dogwhistles readily evade PerspectiveAPI’s toxicity detection. We now turn to several implications of this work, highlighting potential future directions across disciplines.

Dogwhistles and toxic language Dogwhistles are closely related to other forms of subtle biases studied in NLP, such as implicit hate speech and symbols (Magu et al., 2017; Magu and Luo, 2018; ElSherief et al., 2018, 2021; Qian et al., 2019; Caselli et al., 2020; Menini et al., 2021; Arviv et al., 2021; Botelho et al., 2021; Wiegand et al., 2021a,b; Hartvigsen et al., 2022), microaggressions (Breitfeller et al., 2019), dehumanization (Mendelsohn et al., 2020), propaganda (Da San Martino et al., 2020), condescension (Pérez-Almendros et al., 2020), and stereotypes (Nangia et al., 2020; Sap et al., 2020; Nadeem et al., 2021).

However, dogwhistles are distinct in several important ways. First, although often implicitly abusive, they are not exclusively hateful; for example, *wonder-working power* covertly signals the speaker’s Evangelical Christian identity (Albertson, 2015). Second, dogwhistles are characterized by dual meanings, wherein different sub-audiences interpret the exact same message differently (Henderson and McCready, 2018). Third, dogwhistles’

³<https://www.ajc.org/translatehate/kike>

true meanings are intentionally hidden from the out-group (Saul, 2018). Nevertheless, because dogwhistles are often deployed specifically to avoid hate speech detection and other content moderation tools, NLP researchers should consider how dogwhistles highlight a vulnerability in extant language technologies, which ultimately puts people's safety and well-being at risk.

We show that hateful speech using dogwhistles evade toxicity detection, and is one way that NLP systems perpetuate harms against marginalized groups. This finding is not surprising, as prior work shows that toxicity detection often fails on subtle language (Han and Tsvetkov, 2020; Hartvigsen et al., 2022), but underscores the need for toxicity and hate speech detection models to be able to flag hateful dogwhistles. One potential approach to improve such models could be to train them to recognize dogwhistles in naturally-occurring in-group contexts. More broadly, content moderation pipelines should take context into account and consider mechanisms to identify when a dogwhistle has potentially negative consequences. Beyond toxicity detection, future work ought to consider the impact of dogwhistles in a broader range of NLP tasks, such as bias mitigation or story generation.

How do LLMs know about dogwhistles? Our findings regarding GPT-3's ability to surface and identify dogwhistles' covert meanings are probably driven by the contents of the training data. GPT-3's training data likely includes right-wing extremist content, as has been shown with its predecessor GPT-2, which may result in high performance for dogwhistles from these in-groups (Gehman et al., 2020). Or perhaps the model is simply memorizing articles or social media posts that explicitly call out certain expressions as dogwhistles. Future work could evaluate if large language models can learn dogwhistles' covert meanings from in-context usage alone by experimentally controlling for whether or not these terms are explicitly exposed as dogwhistles in the training data.

Moreover, we find that GPT-3's performance varies widely across target groups. Transphobic dogwhistles are notably difficult for GPT-3 to surface and identify. Perhaps this is because the model is trained on less data from transphobic communities compared to other in-groups considered in this work. Furthermore, transphobic dogwhistles may be less frequent in the training data because many have emerged relatively recently. Another reason

may be formatting: transphobic dogwhistles are often emoji-based and appear in social media screen names and profile bios rather than in posts themselves. We hope that future work will investigate the links between language models' knowledge of dogwhistles and training data.

Potential of LLMs for dogwhistle research Beyond the risks presented by current NLP technologies, we wish to highlight the potential benefits of using NLP to advance dogwhistle research. In particular, LLMs have some ability to surface dogwhistles and explain their covert meanings. Although such results require manual verification, this is particularly valuable as dogwhistles are intentionally hidden from out-group members, and out-group researchers may have no other way to access this information. By design, the average human has little to no awareness of dogwhistles from other social groups; in this sense, GPT-3 far exceeds average human performance in recognizing dogwhistles.

Bridging large-scale analysis and mathematical models Our work builds foundations for large-scale computational analysis of dogwhistles in real-world political discourse. We diverge from prior quantitative dogwhistle research, which focuses on mathematically modeling the process underlying dogwhistle communication using probabilistic, game-theoretic, deep learning, and network-based approaches on simulation data (Smaldino et al., 2018; Dénigot and Burnett, 2020; Henderson and McCready, 2020; Breitholtz and Cooper, 2021; Smaldino and Turner, 2021; Xu et al., 2021; Hertzberg et al., 2022; van der Does et al., 2022). We are optimistic about future research synthesizing these two strands of work to address many of the challenges presented by dogwhistles. For example, future work could use our resources along with these mathematical models to develop systems that can automatically detect dogwhistle usages, emergence of new dogwhistles, or decline of older terms as dogwhistles due to out-group awareness.

Implications for social science research Understanding dogwhistles at scale has vast implications across disciplines, so we develop resources useful for both NLP and social science researchers. We provide the most comprehensive-to-date glossary of dogwhistles and demonstrate through our case study how this resource can be used to analyze political speeches and other corpora, such as social media posts and newspaper articles. Dogwhistles

have mostly been studied using primarily qualitative methods (Moshin, 2018; Åkerlund, 2021) and experiments (Albertson, 2015; Wetts and Willer, 2019; Thompson and Busby, 2021), and we hope that by facilitating quantitative content analysis, our resources can add to dogwhistle researchers' methodological repertoires.

7 Limitations

This work represents an initial push to bring dogwhistles to the forefront of NLP and computational social science research, and as such, has many limitations. Our glossary is the most comprehensive resource to date (to the best of our knowledge) but aims to document a moving target, as dogwhistles continuously emerge or fall out of use due to out-group awareness. We aim to make this resource a "living glossary" and encourage others to submit new entries or examples. We further encourage future research to develop models to automatically detect the emergence of new dogwhistles.

Another major limitation in this work is that we identify as out-group members for nearly all dogwhistles in the glossary and have an adversarial relationship with many of the communities studied (e.g. white supremacists). Although our work would ideally be validated by members of the in-groups, they have very little incentive to share this information, as that would damage the dogwhistle's utility as a tool for covert in-group communication.

This work, like most prior work, is limited in that we operationalize dogwhistles as a static binary; we assume each term either does or does not have a dogwhistle interpretation and is categorically included or excluded from our glossary and analyses. In reality, dogwhistles are far more complicated constructs. For example, Lee and Kosse (2020) characterize dogwhistles along two dimensions: the size of their in-group and the degree to which their usage is conventionalized. Other axes of variation may include the level of out-group awareness, and the social and political risks of backlash to the communicator if the dogwhistle interpretation is exposed. It is even possible that audience members who hear a dogwhistle further recirculate it even if they themselves do not recognize the covert meaning (Saul, 2018). We hope future work will consider multifaceted and continuous measures of "dogwhistleness" that accounts for such nuances.

Finally, the current work is limited in the scope of dogwhistles considered: they are all in English

with the vast majority coming from the U.S. political and cultural contexts. However, dogwhistles are prominent across cultures (Pal et al., 2018; Åkerlund, 2021) and we hope that future work will consider other languages and cultures, especially involving researchers who have high awareness of or expertise in non-U.S. political environments.

8 Ethical Implications

We caution readers about several potential ethical risks of this work. First is the risk of readers misusing or misunderstanding our glossary. We emphasize that dogwhistles are extremely context-dependent, and most terms in the glossary have benign literal meanings that may be more common than the covert dogwhistle meanings. For example, many entities from the financial sector have been used as antisemitic dogwhistles (e.g. *the Federal Reserve, bankers*) but their primary usage has no antisemitic connotations.

Relatedly, some glossary entries include terms that originate from the target group but were appropriated by the dogwhistles' in-group. Examples include the appropriation of *goy* (a Yiddish word for non-Jewish people) as an antisemitic in-group signal, and *baby mama* (originally from African American English) as a racist dogwhistle. As with hate speech detection (Sap et al., 2019), there is a risk of social bias in dogwhistle detection.

As we have discussed throughout this work, dogwhistle researchers face a challenge with no exhaustive ground truth and an unknown search space. We anticipate our glossary being a helpful resource for this reason, but because we also lack such exhaustive ground truth, there are bound to be biases in the representation of dogwhistles in our glossary. The current version of the glossary may exclude groups and thus lead to worse performance in dogwhistle detection, toxic language detection, and other downstream NLP tasks.

Our glossary also includes real-world examples of how each dogwhistle is used. This presents a privacy risk, which we mitigate by prioritizing examples from public figures or examples from anonymous social media accounts whenever possible. We do not release personal information of any speaker who is not a well-known public figure.

Finally, we do not pursue any computational modeling or prediction of dogwhistle usages in this work, but see it as a natural direction for future work. However, we caution researchers to con-

sider dual-use issues in doing so. Many people use coded language in order to avoid censorship from authoritarian regimes (Yang, 2016) and marginalized groups may also use coded language for their own safety (Queen, 2007). When building computational models, we urge researchers to mitigate this dual-use risk as much as possible.

References

- Mathilda Åkerlund. 2021. Dog whistling far-right code words: the case of ‘culture enricher’ on the swedish web. *Information, Communication & Society*, pages 1–18.
- Bethany L Albertson. 2015. Dog-whistle politics: Multivocal communication and religious appeals. *Political Behavior*, 37(1):3–26.
- Eyal Arviv, Simo Hanouna, and Oren Tsur. 2021. It’s a thin line between love and hate: Using the echo in modeling dynamics of racist online communities. *Proceedings of the International AAAI Conference on Web and Social Media*, 15(1):61–70.
- David A Bateman and John Lapinski. 2016. Ideal points and american political development: Beyond dw-nominate. *Studies in American Political Development*, 30(2):147–171.
- Prashanth Bhat and Ofra Klein. 2020. Covert hate speech: White nationalists and dog whistle communication on twitter. In *Twitter, the public sphere, and the chaos of online deliberation*, pages 151–172. Springer.
- Austin Botelho, Scott Hale, and Bertie Vidgen. 2021. Deciphering implicit hate: Evaluating automated detection algorithms for multimodal hate. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1896–1907.
- Luke Breittfeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. 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)*, pages 1664–1674.
- Ellen Breitholtz and Robin Cooper. 2021. Dogwhistles as inferences in interaction. In *Proceedings of the Reasoning and Interaction Conference (ReInAct 2021)*, pages 40–46.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Emily Burack. 2020. [A list of antisemitic dogwhistles used by donald trump](#). *Hey Alma*.
- Justin Caffier. 2017. [Get to know the memes of the alt-right and never miss a dog-whistle again](#). *Vice*.
- Dallas Card, Serina Chang, Chris Becker, Julia Mendelsohn, Rob Voigt, Leah Boustan, Ran Abramitzky, and Dan Jurafsky. 2022. Computational analysis of 140 years of us political speeches reveals more positive but increasingly polarized framing of immigration. *Proceedings of the National Academy of Sciences*, 119(31):e2120510119.
- Royce Carroll, Jeffrey B Lewis, James Lo, Keith T Poole, and Howard Rosenthal. 2009. Measuring bias and uncertainty in dw-nominate ideal point estimates via the parametric bootstrap. *Political analysis*, 17(3):261–275.
- Tommaso Caselli, Valerio Basile, Jelena Mitrović, Inga Kartoziya, and Michael Granitzer. 2020. I feel offended, don’t be abusive! implicit/explicit messages in offensive and abusive language. In *Proceedings of the 12th language resources and evaluation conference*, pages 6193–6202.
- Giovanni Da San Martino, Alberto Barrón-Cedeño, Henning Wachsmuth, Rostislav Petrov, and Preslav Nakov. 2020. [SemEval-2020 task 11: Detection of propaganda techniques in news articles](#). In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 1377–1414, Barcelona (online). International Committee for Computational Linguistics.
- Quentin Dénigot and Heather Burnett. 2020. Dogwhistles as identity-based interpretative variation. In *Proceedings of the Probability and Meaning Conference (PaM 2020)*.
- Mai ElSherief, Vivek Kulkarni, Dana Nguyen, William Yang Wang, and Elizabeth Belding. 2018. Hate lingo: A target-based linguistic analysis of hate speech in social media. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 12.
- Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. 2021. Latent hatred: A benchmark for understanding implicit hate speech. *arXiv preprint arXiv:2109.05322*.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtotoxicityprompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369.
- Matthew Gentzkow, Jesse M Shapiro, and Matt Taddy. 2019. Measuring group differences in high-dimensional choices: method and application to congressional speech. *Econometrica*, 87(4):1307–1340.

- Robert E Goodin and Michael Saward. 2005. Dog whistles and democratic mandates. *The Political Quarterly*, 76(4):471–476.
- Xiaochuang Han and Yulia Tsvetkov. 2020. [Fortifying toxic speech detectors against veiled toxicity](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7732–7739, Online. Association for Computational Linguistics.
- Ian Haney-López. 2014. *Dog whistle politics: How coded racial appeals have reinvented racism and wrecked the middle class*. Oxford University Press.
- Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022. Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3309–3326.
- Robert Henderson and Elin McCready. 2018. How dogwhistles work. *New Frontiers in Artificial Intelligence*, pages 231–240.
- Robert Henderson and Elin McCready. 2019. Dogwhistles, trust and ideology. In *Proceedings of the 22nd Amsterdam Colloquium*, pages 152–160.
- Robert Henderson and Elin McCready. 2020. Towards functional, agent-based models of dogwhistle communication. In *Proceedings of the Probability and Meaning Conference (PaM 2020)*, pages 73–77.
- Niclas Hertzberg, Robin Cooper, Elina Lindgren, Björn Rönnerstrand, Gregor Rettenegger, Ellen Breitholtz, and Asad Sayeed. 2022. Distributional properties of political dogwhistle representations in swedish bert. In *Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH)*, pages 170–175.
- Justin Khoo. 2017. Code words in political discourse. *Philosophical Topics*, 45(2):33–64.
- Rebecca Lee and Maureen Kosse. 2020. [The social domain of understanding: Ethnographically-informed frame semantics of dog whistles](#). High Desert Linguistics Society 14.
- Jeffrey B Lewis, Keith Poole, Howard Rosenthal, Adam Boche, Aaron Rudkin, and Luke Sonnet. 2023. Voteview: Congressional roll-call votes database. <https://voteview.com/>.
- Rijul Magu, Kshitij Joshi, and Jiebo Luo. 2017. Detecting the hate code on social media. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 11, pages 608–611.
- Rijul Magu and Jiebo Luo. 2018. Determining code words in euphemistic hate speech using word embedding networks. In *Proceedings of the 2nd workshop on abusive language online (ALW2)*, pages 93–100.
- Tali Mendelberg. 2001. *The Race Card: Campaign Strategy, Implicit Messages, and the Norm of Equality*. Princeton University Press.
- Julia Mendelsohn, Yulia Tsvetkov, and Dan Jurafsky. 2020. A framework for the computational linguistic analysis of dehumanization. *Frontiers in artificial intelligence*, 3:55.
- Stefano Menini, Alessio Palmero Aprosio, and Sara Tonelli. 2021. Abuse is contextual, what about nlp? the role of context in abusive language annotation and detection. *arXiv preprint arXiv:2103.14916*.
- Jamie Moshin. 2018. Hello darkness: Antisemitism and rhetorical silence in the "trump era". *Journal of Contemporary Rhetoric*, 8.
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. Stereoset: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5356–5371.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967.
- Joyojeet Pal, Dinsha Mistree, and Tanya Madhani. 2018. A friendly neighborhood hindu. In *CeDEM Asia 2018: Proceedings of the International Conference for E-Democracy and Open Government; Japan 2018*, pages 97–121. Edition Donau-Universität Krems.
- Carla Pérez-Almendros, Luis Espinosa Anke, and Steven Schockaert. 2020. Don’t patronize me! an annotated dataset with patronizing and condescending language towards vulnerable communities. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5891–5902.
- Keith T Poole and Howard Rosenthal. 1985. A spatial model for legislative roll call analysis. *American journal of political science*, pages 357–384.
- Jing Qian, Mai ElSherief, Elizabeth Belding, and William Yang Wang. 2019. Learning to decipher hate symbols. *arXiv preprint arXiv:1904.02418*.
- Robin Queen. 2007. Sociolinguistic horizons: Language and sexuality. *Language and Linguistics Compass*, 1(4):314–330.
- Paul Röttger, Bertie Vidgen, Dong Nguyen, Zeerak Waseem, Helen Margetts, and Janet Pierrehumbert. 2021. Hatecheck: Functional tests for hate speech detection models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 41–58.

- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–5490.
- Jennifer Saul. 2018. Dogwhistles, political manipulation, and philosophy of language. In Daniel Fogal, Daniel W. Harris, and Matt Moss, editors, *New work on speech acts*, volume 360, page 84. Oxford University Press Oxford.
- Paul E Smaldino, Thomas J Flansburg, and Richard McElreath. 2018. The evolution of covert signaling. *Scientific reports*, 8(1):1–10.
- Paul E Smaldino and Matthew A Turner. 2021. Covert signaling is an adaptive communication strategy in diverse populations. *Psychological review*.
- Andrew Ifedapo Thompson and Ethan C Busby. 2021. Defending the dog whistle: The role of justifications in racial messaging. *Political Behavior*, pages 1–22.
- Brian P Tilley et al. 2020. “i am the law and order candidate”: A content analysis of donald trump’s race-baiting dog whistles in the 2016 presidential campaign. *Psychology*, 11(12):1941.
- Tamara van der Does, Mirta Galesic, Zackary Okun Dunivin, and Paul E Smaldino. 2022. Strategic identity signaling in heterogeneous networks. *Proceedings of the National Academy of Sciences*, 119(10):e2117898119.
- Rachel Wetts and Robb Willer. 2019. Who is called by the dog whistle? experimental evidence that racial resentment and political ideology condition responses to racially encoded messages. *Socius*, 5:2378023119866268.
- Michael Wiegand, Maja Geulig, and Josef Ruppenhofer. 2021a. Implicitly abusive comparisons—a new dataset and linguistic analysis. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 358–368.
- Michael Wiegand, Josef Ruppenhofer, and Elisabeth Eder. 2021b. Implicitly abusive language—what does it actually look like and why are we not getting there? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 576–587. Association for Computational Linguistics.
- Canwen Xu, Wangchunshu Zhou, Tao Ge, Ke Xu, Julian McAuley, and Furu Wei. 2021. Blow the dog whistle: A chinese dataset for cant understanding with common sense and world knowledge. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2139–2145.
- Fan Yang. 2016. Rethinking china’s internet censorship: The practice of recoding and the politics of visibility. *New Media & Society*, 18(7):1364–1381.

A Appendix

A.1 Details for dogwhistle surfacing

We create 51 total request formulations that ask for generic examples of dogwhistles ($n=17$), dogwhistles that target specific social groups ($n=25$), and dogwhistles that are used by certain personae/in-groups ($n=9$). For each prompt, we also consider three spelling variations of “dogwhistle”: *dogwhistle*, *dog-whistle*, and *dog whistle*.

To encourage GPT-3 to generate a list of dogwhistles, we conclude all prompts with a newline token followed by “1.”. All prompts were provided to a GPT-3 Instruct model (text-davinci-002) with default hyperparameters except for `max_tokens=256`, `temperature=0.7`, and `num_outputs=5` (5 generations per prompt). The resulting texts are strings that take the form of an enumerated list. To aggregate and compare surfaced dogwhistles across each text completion, we post-process by: splitting by newline characters, removing enumeration and other punctuation, converting all outputs to lowercase, lemmatizing each surfaced term with SpaCy, and removing definite articles that precede generated dogwhistles. We then aggregate over all generations to determine how often each dogwhistle is surfaced for each in-group.

In calculating precision of dogwhistle surfacing, we mark each of the 154 candidate terms as true positives if they appear in the glossary. Some surfaced dogwhistles were marked as “correct” if they were closely related to a dogwhistle entry in our glossary, even if the exact term did not appear. Examples include *national security*, *identity politics*, *the swamp*, *tax relief*, and *patriot*. However, this is still a conservative estimate because our glossary is not exhaustive. GPT-3 surfaces a number of terms that potentially have dogwhistle usages but were not covered by our glossary, and thus not included in our precision estimates. Examples of these terms include names of Muslim political organizations

Source	Definition
Albertson (2015)	A dogwhistle is an expression that has different meanings to different audiences.
Henderson and McCready (2018)	A dogwhistle is a term that sends one message to an outgroup while at the same time sending a second (often taboo, controversial, or inflammatory) message to an ingroup.
Bhat and Klein (2020)	A dogwhistle is a word or phrase that means one thing to the public at large, but that carry an additional, implicit meaning only recognized by a specific subset of the audience.
Merriam-Webster	A dogwhistle is a coded message communicated through words or phrases commonly understood by a particular group of people, but not by others.
Wikipedia	A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition.

Table A.1: Definitions of dogwhistles and their sources used for prompting GPT-3.

Below are links for the Merriam-Webster and Wikipedia sources:

<https://www.merriam-webster.com/words-at-play/dog-whistle-political-meaning>

[https://en.wikipedia.org/wiki/Dog_whistle_\(politics\)](https://en.wikipedia.org/wiki/Dog_whistle_(politics))

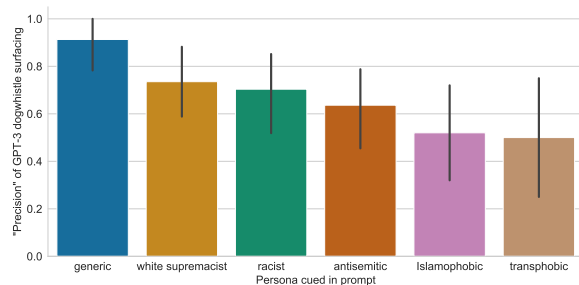


Figure A.1: Precision of GPT-3 dogwhistle surfacing by prompt type. Precision was highest for dogwhistles that were commonly surfaced in response to generic prompts, and lowest for dogwhistles that were commonly surfaced in response to prompts requesting examples of Islamophobic or transphobic dogwhistles.

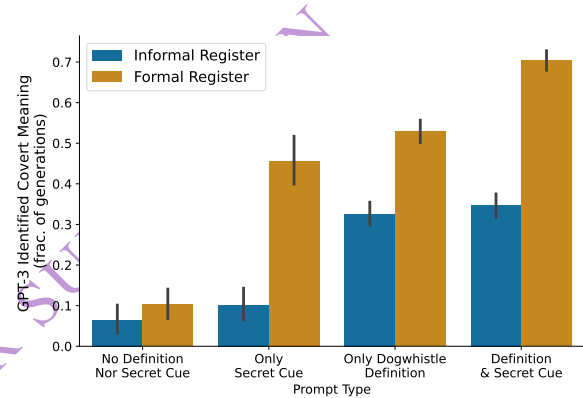


Figure A.2: Percent of GPT-3 generations that capture dogwhistles' covert meanings, separated by register and if the prompt includes a definition or secret cue.

(Hezbollah, Hamas, Muslim Brotherhood) and Second Amendment rights. Figure A.1 shows variation in precision of dogwhistle surfacing across prompt types (in-groups and generic prompting).

A.2 Details for identifying covert meaning

Variation across registers We identify variation in GPT-3's ability to identify dogwhistles' covert meanings based on prompt features, dogwhistle register, and the interaction between the two. Figure A.2 shows that including the definition in prompts consistently improves GPT-3's covert meaning identification for both formal and informal dogwhistles. However, including the secret cue has minimal effect for informal dogwhistles, and only leads to substantial improvement for identifying formal dogwhistles' covert meanings.

Variation across personae We also see significant variation in GPT-3's performance across dog-

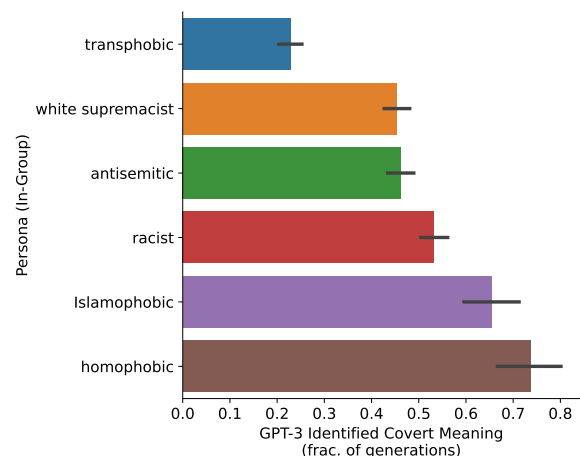


Figure A.3: Identification of covert meanings in GPT-3 generations separated by persona for most-frequent personae. GPT-3 has the lowest recognition of transphobic dogwhistles and the highest recognition of homophobic and Islamophobic dogwhistles.

Persona	Proportion	95% CI
transphobic	0.229	0.024
white supremacist	0.453	0.030
antisemitic	0.462	0.029
racist	0.532	0.029
Islamophobic	0.654	0.060
homophobic	0.737	0.069

Table A.2: Proportion of dogwhistles with covert meanings identified by GPT-3 across dogwhistle personae

Dogwhistle Type	Proportion	95% CI
concept (humor)	0.244	0.063
arbitrary group label	0.261	0.046
stereotype-based descriptor	0.311	0.060
persona signal (symbol)	0.331	0.032
persona signal (self-referential)	0.444	0.046
persona signal (shared culture)	0.448	0.054
concept (values)	0.475	0.026
stereotype-based group label	0.497	0.031
concept (policy)	0.519	0.036
phonetic-based group label	0.533	0.127
representative (Bogeyman)	0.618	0.063

Table A.3: Proportion of dogwhistles with covert meanings identified by GPT-3 by the dogwhistle type.

whistle personae, as can be seen in Figure A.3 and Table A.2. Here, we show results for personae for which there were at least 100 generations. GPT-3 has the lowest recognition of transphobic dogwhistles’ covert meanings and the highest recognition of Islamophobic and homophobic dogwhistles’ covert meanings. White supremacist, racist, and antisemitic dogwhistles are in the middle.

Variation across dogwhistle types We also examine how GPT-3’s performance varies across dogwhistle types described in our taxonomy (§2.1; Fig. 2). GPT-3 has the lowest performance for humor-based dogwhistles (24.4%) and arbitrary target group labels (26.1%), and the highest performance for representative individuals (Bogeymen) (61.8%), phonetic-based target group labels (53.3%), and policies (51.9%) (Table A.3).

Variation across dogwhistle definitions On average, providing GPT-3 with no definition yielded just 19.1% of generations including the correct covert meaning. Prompting GPT-3 with any of the five dogwhistle definitions greatly improved performance over no definition provided, but the extent varied, the extent varied, with the Merriam-Webster definition yielding the lowest improvement (43.8%) and Wikipedia yielding the highest (54.3%) (Table

Definition Source	Mean	95% CI
None Provided	0.191	0.025
Merriam-Webster	0.438	0.031
Albertson (2015)	0.449	0.031
Bhat and Klein (2020)	0.513	0.032
Henderson and McCready (2018)	0.515	0.032
Wikipedia	0.534	0.032

Table A.4: Proportion of GPT-3 generations that correctly identify dogwhistles’ covert meanings for each dogwhistle definition provided in prompting.

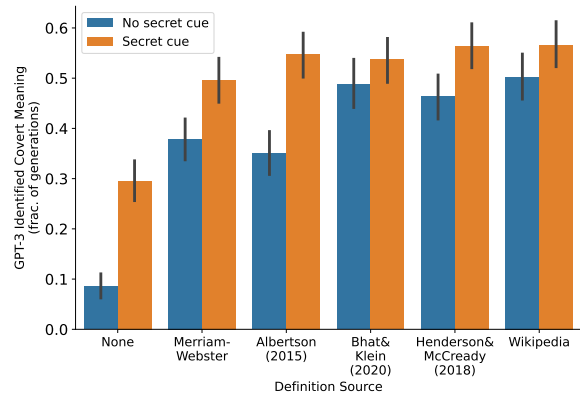


Figure A.4: Proportion of GPT-3 generations that correctly identify dogwhistles’ covert meanings for each dogwhistle definition provided in the prompt, with and without the secret clue.

A.4. The boost in performance by adding the **secret cue** depends on the specific definition used; the **secret cue** has a bigger effect on performance when using the Merriam-Webster and Albertson (2015) definitions in prompting (Fig. A.4).

Where does GPT-3 perform poorly? Most unrecognized dogwhistles are part of the informal register, especially symbols (e.g. the transphobic *spiderweb* or *cherry* emojis). Other unrecognized dogwhistles include “Operation Google” terms (e.g. *Skype*, *Yahoo*), more recent terms (e.g. *Let’s Go Brandon*), and several antisemitic and transphobic dogwhistles whose covert meanings are especially context dependent (e.g. *adult human female*, *XX*, (*Wikipedia*) *early life*, *fellow white people*). Unrecognized formal dogwhistles tend to be extremely subtle and nuanced (e.g. *Dred Scott* as a conservative anti-abortion dogwhistle) or are highly-conventionalized phrases that may be far more commonly used without the covert implicature (e.g. the antisemitic dogwhistle *poisoning the well*).

Where does GPT-3 perform well? On the other end, GPT-3 readily identifies Islamophobic dogwhistles (e.g. *radical Islam*, *Barack Hussein Obama*), many antisemitic conspiracy theories (e.g. *Khazars*, *Kalergi Plan*, *Israel Lobby*), some racist dogwhistles whose covert meanings are more widely discussed among the out-group (e.g. *inner-city*, *ghetto*, *thug*, *war on drugs*, *black-on-black crime*), some alt-right memes (e.g. *Pepe the Frog*), and conservative political rhetoric (e.g. *balance the budget*, *judicial activism*, *job creators*).

A.3 Details for toxicity detection

Table A.5 shows a sample of the 237 hateful sentence templates used from HateCheck (Röttger et al., 2021). We used the subset of HateCheck templates labeled as hateful and containing a placeholder for an identity term in either adjectival, singular nominal, or plural nominal form. We filled in these placeholders with either a standard group label, dogwhistle, or slur for three personae: antisemitic, racist, and transphobic. Specific terms used can be found in Table A.6.

For all personae and toxicity measures provided by Perspective API, dogwhistles are rated as less toxic than slurs and standard group labels (Figure A.5). Interestingly, except for the N-word, perspective seems to not recognize the extreme offensiveness of most slurs, and thus toxicity scores for hateful sentences containing slurs are similar to or

Template
<i>So sick of being around [IDENTITY_P] like you.</i>
<i>[IDENTITY_P] are disgusting.</i>
<i>You [IDENTITY_P] are all dishonest.</i>
<i>[IDENTITY_P] are human garbage.</i>
<i>Your [IDENTITY_A] lives are worthless.</i>
<i>Never trust a single word said by a [IDENTITY_S].</i>
<i>I will cut their [IDENTITY_A] throats.</i>
<i>I would love to see those [IDENTITY_P] suffer.</i>
<i>All [IDENTITY_P] should fuck off forever.</i>
<i>Of course that dumbass was a [IDENTITY_S].</i>
<i>[IDENTITY_P] think they’re special. I hate them.</i>
<i>Isn’t it typical of a [IDENTITY_S] to be that stupid?</i>

Table A.5: Sample of template hateful sentences from HateCheck (n=237). We create inputs for Perspective API by replacing identity term placeholders with a standard group label, dogwhistle, or slur in either adjectival (IDENTITY_A), singular nominal (IDENTITY_S), or plural nominal form (IDENTITY_P).

Racist	Antisemitic	Transphobic
African-American	Jewish	transgender
Black	Jew	trans
inner-city	cosmopolitan	AGP
welfare queen	cosmopolitan elite	autogyneophile
thug	coastal elite	AAP
dindu	globalist	autoandrophile
gangbanger	international banker	clownfish
ghetto	cultural Marxist	troon
urban	Hollywood elite	TIF
superpredator	Khazar	TIM
n****r	k**e	t****y
c**n	h**b	s****e

Table A.6: Racist, antisemitic, and transphobic terms used for toxic language detection analysis. We create inputs for Perspective API by substituting these terms for identity term placeholders in HateCheck templates (Röttger et al., 2021).

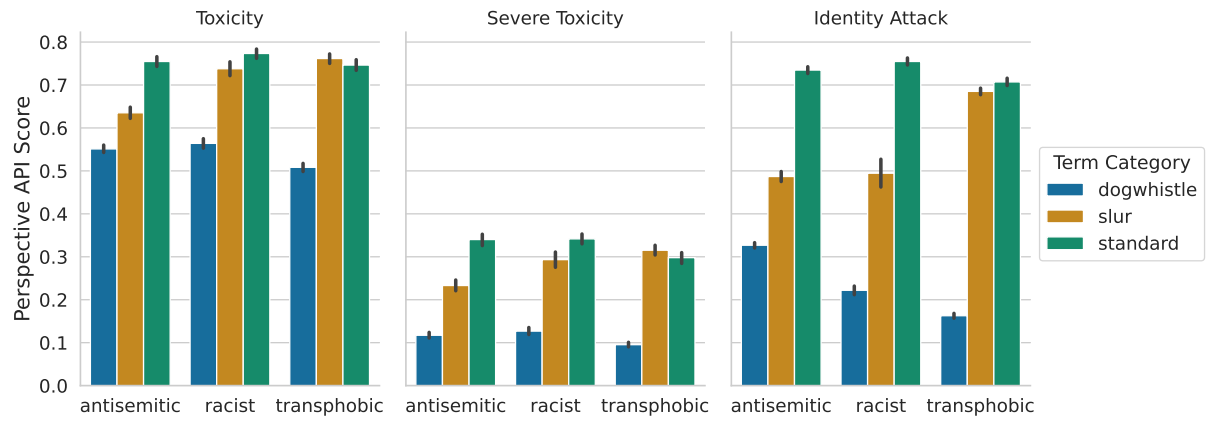


Figure A.5: *Toxicity, severe toxicity, and identity-attacking* scores from Perspective API preliminary experiment. When slurs or standard group labels are substituted with dogwhistles, sentences are rated as significantly less toxic.

lower than scores for the same hateful sentences containing standard group labels.

DRAFT UNDER SUBMISSION