# Re-evaluating the Tomes for the Times

**Ryan Brate,[†] Marieke van Erp,[†] Antal van den Bosch[⊕]**

[†]DHLab, [⊕]Utrecht University

KNAW Humanities Cluster, DHLab, Amsterdam, the Netherlands

Utrecht University, Institute for Language Sciences, Utrecht, the Netherlands

{ryan.brate, marieke.van.erp}@dh.huc.knaw.nl

a.p.j.vandenbosch@uu.nl

## Abstract

Literature is to some degree a snapshot of the time it was written in and the societal attitudes of the time. Not all depictions are pleasant or in-line with modern-day sensibilities; this becomes problematic when the prevalent depictions over a large body of work are negatively biased, leading to their normalisation. Many much-loved and much-read classics are set in periods of heightened social inequality: slavery, pre-womens' rights movements, colonialism, etc. In this paper, we exploit known text co-occurrence metrics with respect to token-level level contexts to identify prevailing themes associated with known problematic descriptors. We see that prevalent, negative depictions are perpetuated by classic literature. We propose that such a methodology could form the basis of a system for *making explicit* such problematic associations, for interested parties: such as, sensitivity coordinators of publishing houses, library curators, or organisations concerned with social justice.

**Keywords:** bias in literature, charged terms, corpus linguistics

**Disclaimer:** This paper contains derogatory words and phrases. They are provided solely as illustrations of the research results and do not reflect the opinions of the authors or their organisations.

## 1. Introduction

> But Mrs. Tome Gallien's Adventure? A woman like Mrs. Tome Gallien wouldn't stop at anything! It might be a pair of llamas from Peru! Or a greasy witchy-gypsy to tell his fortune! Or a homeless little jet-black pickaninny with a banjo and–consumption!

The quotation above is an excerpt from "The Stingy Receiver", by Eleanor Hallowell Abbott. It is impressive in that it manages to propagate negative connotations for two people groups in only a few short lines. It also exemplifies, that in at least in some regions of the collective body of literature, there very much exists biased and potentially harmful depictions of peoples. Rare problematic depictions are one thing, but this is particularly problematic when such depictions are highly prevalent in a corpus, encouraging the reader towards negatively biased concept associations. Writing *witchy-gypsy* in a novel is not helpful, but it is particularly problematic if the association of 'witchy' and 'gypsy' are relatively high on aggregate.

This issue is already recognised: there is a growing effort acknowledging and challenging such highly prevalent and prejudiced depictions, promoting greater awareness in learners of biased themes.[1] Publications such as "Anti-Bias Curriculum: Tools for Empowering Young Children" ([Denman-Sparks, 1993](#)), which highlight the portrayal of 'happy-go-lucky blacks', the 'fat eye-rolling mammy'; the 'inscrutable, slant-eyed' Oriental; the 'naked and savage' Native American; womanhood as 'domesticated motherhood', the 'demure young woman', the 'doll-loving' girl or the 'wicked' step-mother. Such endeavours can be helped by data-driven context analysis of concept depictions, to better understand prevailing contexts associated with particular concepts, and to understand how prevalent they are. In this paper, we ask: *In applying transparent descriptor:context–feature metrics, as the basis for a trope ranking system; to what extent observe known tropes in highly ranked regions?* The remainder of this paper is organised as follows. In Section 2, we discuss related work. In Section 3, we present our datasets. In Section 4, we detail the steps in our analysis methodology, followed by evaluation and results in Section 5. We present our conclusions in Section 6.

## 2. Related Work

Instances of casually used language with racist overtones in literature, particularly in non–current work, are a recognised phenomenon ([Betensky, 2019](#)) and have recently received mainstream media attention.[2] Harmful biases may also be subtler,

---

[1] https://www.teachingforchange.org/

[2] https://www.theguardian.com/books/2023/mar/26/agatha-christie-novels-reworked-to-remove-potentially-offensive-language Last accessed: 20 October 2023

for example, correlations between goodness and beauty (Rees, 1988; Yacovone, 2020). Much of the work in the field is expert–analysis of subject matter, highlighting specific instances, general themes, typically in regards to small number of literature examples.

Automatic detection of charged language has largely focused on hate speech detection on contemporary texts and in particular social media. For example several workshops on toxic language were organised[3,4] as well as two SemEval 2022 shared tasks namely Task 4: Patronizing and Condescending Language Detection (Perez-Almendros et al., 2022) and Task 5: MAMI - Multimedia Automatic Misogyny Identification (Fersini et al., 2022). Various datasets have been made available, either manually annotated from comments such as Jigsaw (cjadams et al., 2019) covering a range of offensive comments such as pertaining to disability, gender, race or ethnicity, religion, or sexual orientation, or machine-generated such as ToxiGen (Hartvigsen et al., 2022) focusing on minority groups.

What makes the issue of detecting toxic, charged or offensive language difficult is that the language can take many different forms: at times clearly offensive terms are present, at other times emphasis (for example highlighting a stereotype) make the comment offensive (Perez Almendros and Schockaert, 2022). Therefore, it is important to focus on detecting both explicit and implicit charged language (Lin, 2022). Furthermore, the research community is also still investigating what type of evaluation metrics are best suited (Chen et al., 2023; Jourdan et al., 2023)for which situation.

To the best of our knowledge, we are the first to automatically assess chargedness of literary texts. Previous work has demonstrated the potential for co-occurrence statistics of extracted contextual features, to identify the prevailing associations of derogatory terms in text (Brate et al., 2023). In this paper, we apply this approach of noun–context co-occurrence statistics as the basic of a negative bias recommender system in the English language literature domain.

## 3.  Literature Corpora

The entire English and American literature data set is harvested from Project Gutenberg,[5] according to their listed Library of Congress Classification system (LoCC)[6] labels. A total of 9,185 English literature (LoCC of 'PR') and 10,873 American Literature (LoCC of 'PS') books are identified, to date. Four of the English literature set, had no corresponding text file and were discarded. A further two English literature books could not be processed and were discarded. A complete list of the utilised 9,179 English literature and 10,873 Project Gutenberg book ids available on our Github repository. For each book, only the text between Project Gutenberg standard start and end tags is considered.

## Data processing

Separately for the English and American literature corpora, co-occurrence frequencies of proper nouns and nouns, and the contexts features with which they appear are extracted. Where context features consist of co-occurrent verbs for which the noun is the agent; and verbs for which the noun is the patient and adjectives applied to the nouns. This is done via pattern matching against spaCy[7] dependency parsings of the corpora at the sentence part level, as separated by ; or : or , punctuation marks. This results in 29,763,170 noun–adjective pairs, 16,157,119 agent–verb pairs and 8,996,080 patient–verb pairs for the English literature corpus. 22,592,953 noun–adjective pairs, 14,266,666 agent–verb pairs and 7,725,577 patient–verb pairs for American.

These noun–context frequencies are then converted to relative affinity scores via the Log Likelihood Ratio (LLR) test statistic (Dunning, 1993). These noun-context affinity scores are used as the basis for understanding what are the most powerful connotations from the data set. That is, for every context instances in a corpus, there is a corresponding matrix, $M$, of affinity scores of dimension $(|nouns|, |contextfeatures|)$.

As per Table 3, the co–occurrence of some noun and some context feature, can be considered in terms a contingency table of binomial outcomes. As binomial outcomes, they can be considered as the outcomes of binomial generative processes, for which a joint likelihood, $L$, can be formulated: $L(feature = 1|noun = 1, p_1) \times L(feature = 1|noun = 0, p_2)$. The Log Likelihood Ratio, is a factored ratio of this likelihood with respect to two generative assumptions, as reflected in parameters $p_1$ and $p_2$: the null assumption, that the occurrence of some feature is irrespective of the occurrence or absence of the noun in questions ($p_1=p_2$); and the alternative assumption, that $p_1 \neq p_2$. In both cases, the values of parameters are taken according to their maximum likelihood assumptions with respect to the observed values given by the contingency table. The net effect of this score, is a numeric indication of the relative affinity of noun and feature.

---

| | some feature | some feature' |
|---|---|---|
| some noun | C(n,f) | C(n,f') |
| some noun' | C(n',f) | C(n',f') |

Table 1: Illustrative contingency table of the co–occurrence of some noun and some feature, as the basis noun–feature LLR calculation.

$$LLR(\text{noun, feature}) = -2.\lambda(\text{noun, feature})$$
$$\lambda(\text{noun, feature}) = \frac{L_{null}}{L_{alt}} \quad (1)$$

## 4. Methodology

We first validate the applicability of the noun–context affinity scores, for the context types, as the basis for identifying prevalent and negatively biased characterisations. We then apply this methodology to the literature corpora to identify the problematic connotations for words known to be somewhat derogatory. All supporting code is available in the corresponding GitHub repository.[8]

We define a negative trope as a depiction which is both negatively biased in its connotations, and also highly prevalent in the corpus or sub-corpus being considered.

We select the following words, recognised in broad terms as having derogatory use-cases: *cripple, dwarf, gypsy, native, negress, negro, nigger, oriental, servant, slave, tribe, and tribesman*. Additionally, we explore asymmetric gender biases via the descriptors: *bachelor, spinster*. We supplement these with problematic pairs identified in the validation exercise.

We would expect noun–feature pairs, which are indicative of some trope, to be highly mutually predictive of one another. We define some noun and some feature being highly mutually predictive where both: i) the noun–feature LLR score is highly ranked with respect to the noun; and ii) the noun–feature LLR score is highly ranked with respect to the feature. In terms of the LLR matrix representing the corpus of dimensions $(|nouns|, |context features|)$; as illustrated by Figure 1, the mutual affinity is a product of some metric proportional to the within–row and within–column rank. These respective metrics are given by the factors in Equation 2, where $R_r$ is the within–row rank with respect to $N_c$ features and and $R_c$ is the within–column rank with respect to $N_r$ rows, and the overall mutual affinity score is given by their geometric mean. The noun–feature mutual affinity
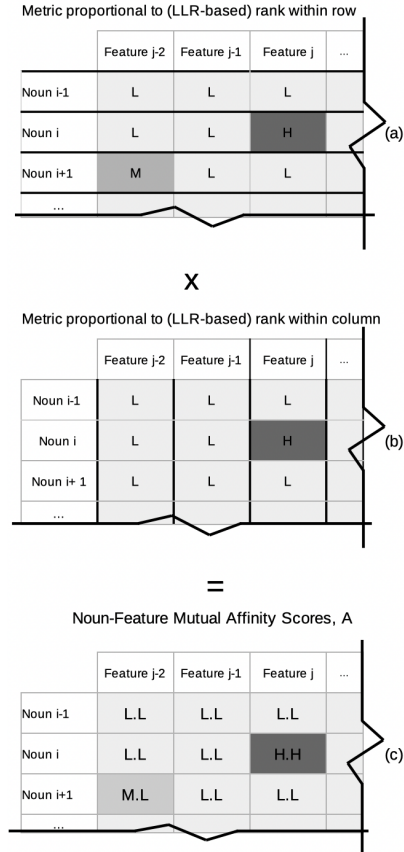
Figure 1: Pictorial representation the noun-feature mutual affinity score (c), a product of: (a) a metric proportional to the LLR-based rank of the feature given the noun; and (b) a metric proportional to the LLR-based rank of the noun given the feature.

scores are between 0 and 1.

$$\begin{matrix}\text{noun–feature}\\\text{mutual affinity score}\end{matrix} = \sqrt{\left(1 - \frac{(R_r - 1)}{N_c}\right)\left(1 - \frac{(R_c - 1)}{N_r}\right)} \quad (2)$$

Geometric mean is used as the aggregate ranking function, owing to its tendency to promote *more similar* fractions. The very quality we want given the emphasis on high *mutual* association. I.e., $\sqrt{0.9 \times 0.9} > \sqrt{0.91 \times 0.89}$.

Consider the part, $\left(1 - \frac{(R_r - 1)}{N_c}\right)$: This approximates the percentile position of some noun, with respect to all other nouns coincident with a feature. The form of the equation is a correction over $\left(1 - \frac{R_r}{N_c}\right)$, in that it correctly gives gives a high score indicative of a high affinity to the scenario, Rf=1, Nf=1; and a low score indicative of a low affinity to the scenario, Rf=1000, Nf=1000.

## 5. Evaluation and Results

We validate the methodology according to the seed words and expected outcomes from Table 2. For known derogatory descriptors of section 4. We then

apply the methodology to profile their prevailing negative biases in the corpora, as the basis for ranking a small subset of the worst-offending books and authors.

## 5.1. Validation of the methodology

We first ask whether the targeted context features and the mutual affinity scores yield meaningful connotations. Table 2 lists known problematic characterisations of certain words (Denman-Sparks, 1993). For each of these word and connotation sets we test whether the methodology as described in Section 4 is able to extract these characterisations.

| Word | Known problematic connotations |
| --- | --- |
| Orientals | slit-eyed, inscrutible |
| Chicano | sombrero-wearing peon, fiesta-loving, macho bandito |
| savage, primitive, lazy, conniving, superstitious treacherous, wily, crafty inscrutable, docile, backwards | loaded language (typically racist) |

Table 2: Known–problematic connotations from "Anti-Bias Curriculum: Tools for Empowering Young Children" (Denman-Sparks, 1993)

Tables 3 and 4 list the strength of association of selected noun–adjective pairs from Table 2 in the American literature corpus. Table 5 is similarly for the English literature corpus. High mutual affinity scores clearly correspond with instances of race–related loaded language.

| Noun | Corresponding adjectives | Adjective rank by noun–feature mutual affinity score |
| --- | --- | --- |
| Oriental | skew–eyed | 2 / 65 |
|  | not slant-eyed | 4 / 65 |
| Chinese | slant–eyed | 7 / 192 |
|  | slit-eyed | 30 / 192 |
| savage | naked | 2/902 |
|  | african | 9/902 |

Table 3: Nouns and corresponding adjectives indicative of Table 2 known problematic connotations, and their mutual affinity score rankings according to co–occurrence frequencies in the data.

## 5.2. Problematic depictions in literature

For the known-problematic descriptors listed in 4, Tables 6 and 7 give the resulting features with highly ranked noun–feature mutual affinity scores; where we deem the feature in question, compounded by its high affinity, to be negative in its connotations. We observe a variety of negative associations and tropes: an over-emphasise on the physical size of black peoples referenced by derogatory terms; their association with servile roles and de–humanising treatment; the hideous dwarf; the helpless cripple, the savage, warlike tribe. We also observe indications of the out–dated misogynistic trope of the eligible–bachelor versus the unwanted spinster.

| Adjectives | Corresponding nouns | noun rank by noun–feature mutual affinity score |
| --- | --- | --- |
| crafty | Jew | 57 / 842 |
| docile | slave | 6 / 335 |
|  | Javanese | 21 / 335 |
| savage | tribe | 3 / 2376 |
|  | Indians | 10 / 2376 |
|  | race | 22 / 2376 |
| superstitious | negro | 15 / 644 |
|  | native | 18 / 644 |
|  | Indians | 20 / 644 |
|  | black | 30 / 644 |
|  | Islamites | 37 / 644 |
| primitive | race | 16 / 1610 |
|  | culture | 22 / 1619 |
| wily | savage | 1. 691 |
|  | Oriental | 3 / 691 |
|  | Greek | 7 / 691 |
|  | Indian | 10 / 691 |
|  | Italian | 22 / 691 |
|  | redskin | 36 / 691 |
|  | Jew | 37 / 691 |
| treacherous | savage | 13 / 1258 |
|  | Indians | 36 / 1258 |
|  | native | 37 / 1258 |

Table 4: The loaded adjectives of Table 2, together nouns demonstrative of the known problematic connotations and their corresponding noun–adjective mutual affinity score ranks. Rankings are with respect to the total number of nouns co–occurrent with each adjective in the American literature corpus.

| Adjectives | Corresponding nouns | noun rank by noun–feature mutual affinity score |
| --- | --- | --- |
| crafty | chief | 16 / 652 |
|  | savage | 18 / 652 |
|  | Somalis | 39 / 652 |
|  | redskin | 45 / 652 |
| docile | slave | 11 / 335 |
|  | Javanese | 21 / 335 |
| superstitious | Spaniards | 14 / 911 |
|  | native | 30 / 911 |
|  | race | 98 / 911 |
| lazy | nigger | 8 / 1661 |
| primitive | savage | 8 / 1450 |
|  | race | 11 / 1450 |
|  | Christian | 17 / 1450 |
| wily | Jesuit | 1 / 847 |
|  | Italian | 3 / 847 |
|  | Greek | 5 / 847 |
|  | Teuton | 8 / 847 |
|  | Frenchman | 10 / 847 |
|  | Jew | 15 / 847 |
|  | savage | 19 / 847 |
|  | Russian | 37 / 847 |
|  | Chinaman | 38 / 847 |
| treacherous | savage | 27 / 1540 |
|  | malay | 55 / 1540 |

Table 5: The loaded adjectives of Table 2, together nouns demonstrative of the known problematic connotations and their corresponding noun–adjective mutual affinity score ranks. Rankings are with respect to the total number of nouns co–occurrent with each adjective in the English literature corpus.

## 6. Conclusion

In this paper, we proposed a corpus analysis method that produces ranked lists of contextual relations most commonly found with charged terms. The method uses a parser for extracting basic re-

| Noun | Associated adjectives |
|---|---|
| barbarian | negro (7 / 271), dark–skinned (17) |
| cripple | helpless (1 / 221), poor (2), hopeless (3), wretched (4), miserable (6), deformed (11), hateful (13) |
| dwarf | hideous (1 / 278), ugly (2) , misshapen (3), deformed (4), apish (6), hunchbacked (14), crippled (21) |
| native | african (1 / 748), dark–skinned (4), savage (8), uncivilised (10), intelligent (17), wretched (20), ignorant (22), cannibal (33) |
| negress | surly (2 / 87), thick–lipped (6), hideous (8), shrivelled (11), deformed (13), tall (14), gigantic (18) |
| negro | naked (1 / 438), half-naked (3) , gigantic (4), big (7), huge (13), free (14), fugitive (11), faithful (12), diminutive (18), giant (21), |
| servant | native (15 / 1995), negro (20) |
| slave | negro (4 / 1200), nubian (5), black (10), abyssinian (18) |
| tribe | savage (1 / 929), hostile (2), warlike (8), wild (13), barbarous (15) |
| tribesman | wild (3 / 91), savage (5 / 91) |
| spinster | gaunt (5 / 405), sour (7), dour–faced (16), unwanted (27), flat–bosomed (29), |
| bachelor | eligible (2 / 507) |
| | **Associated verbs for which the noun is the agent** |
| negress | grin (1 / 120), snatch (2), covet (10), |
| negro | grin (1 / 573), kill (9), |
| nigger | massacre (9 / 242), grin (19), moan (24), attack (18), kill (21), rob (25), murder (28) |
| tribe | fight (8 / 580) |
| tribesman | attack (2 / 98), fight (6), swarm (9), kidnap (19) |
| | **Associated verbs for which the noun is the patient** |
| nigger | shoot (9 / 79), beat (11) |
| slave | sell (3 / 539), flog (4), beckon (8), order (9) |

Table 6: Context features in the English lit. corpus, demonstrative of problematic connotations. Each feature is listed with its noun–feature mutual affinity rank, with respect to the number of features associated with the noun.

| Noun | Associated adjectives |
|---|---|
| cripple | helpless (1 / 189), poor (2), hopeless (3), unpresentable (5), miserable (6) |
| dwarf | misshapen (2 / 201) , shriveled (5), ugly (7), cunning (8), hideous (14), grotesque (15), hunchbacked (18) |
| gypsy | vagrant (10 / 117), witchy (13), not trustworthy (14) |
| native | hostile (6 / 694), ignorant, superstitious, little, full-blooded |
| negress | slatternly (5 / 142), fat (7), good–natured (8), stout (9), big (13) |
| negro | full–blooded (4 / 872), ignorant (5), giant (8), gigantic (10), burly (12) |
| nigger | runaway (1 / 609), lazy (6), damned (7), dead (9), onery (12), damn (13) |
| oriental | godless (3 / 39) |
| slave | fugitive (1 / 1026), runaway (2), negro (5) |
| tribe | savage (2 / 803), hostile (3), warlike (7), primitive (23) |
| tribesman | savage (5 / 114), wild (6), intransigent (7), ferocious (13), hostile (18) |
| | **Associated verbs for which the noun is the agent** |
| gypsy | steal (1 / 139) |
| negro | obey (1 / 1068), shuffle (2), grin (3), row (4), bow (6), mutilate (32), murder (40) |
| nigger | steal (4 / 497) |
| tribe | fight (8 / 642) |
| | **Associated verbs for which the noun is the patient** |
| negro | lynch (1 / 445), disenfranchised (3), sell (5) , permit (11) |

Table 7: Context features in the American lit. corpus, demonstrative of problematic connotations. Each feature is listed with its noun–feature mutual affinity rank, with respect to the number of features associated with the noun.

lational patterns and log-likelihood ratio (LLR) to estimate salient co-occurrences, and furthermore ranks term-context pairs by their mutual association according to ranked LLR scores. The methodology allows for forward–backward pivoting: from nouns to features; from features to nouns, from nouns to nouns in respect of common features. Thereby, enabling cyclical identification of problematic depictions, seeded from a few known words.

With respect to the research question, *In applying transparent descriptor:context–feature metrics, as the basis for a trope ranking system; to what extent do we observe known tropes in highly ranked regions?*, we measure extent of observance in highly ranked regions, according to the percentage rank of the first feature deemed indicative of a negative connotation. As per Tables 3, 4, 5, and the expected negative connotations of Table 2: we observe potentially problematic race related associations of the adjectives, savage, superstitious, primitive, wily, treacherous, lazy and docile with in the top 1% of their respective, co–occurrent nouns. As expected, we also observe derogatory physically–descriptive adjectives related to Oriental peoples in the top 5%. In respect of the terms listed in Section 4, as per

Tables 6 and 7, we observe negative connotations in the top 1% of rankings for cripple, dwarf, native, negro, tribe, nigger, slave variously with respect to the English and American literature corpora.

This proposed methodology enabling data-driven discovery of an extended distribution of charged combinations of adjectives, verbs, and other nouns: represents a richer point of departure for analyzing literature on charged language than a limited initial word list, which a recommender system may otherwise be based on. Looking forward, this extended set of terms could form the basis of a larger filtering system that ranks documents by proportion of charged content, and marks paragraphs for further inspection. Leaving the actual decision-making to curators, editors, or sensitivity co–ordinators in publishing houses, the filtering system would allow the human experts to query a large catalogue for loaded content beyond human scale.

Future work would be the implementation of such a filtering and recommender system and the real-world evaluation of the system by professionals. Collaboration with literature researchers would be required to gain a deeper insight and nuanced view on the differences we now observe superficially between English and American literature.

## Acknowledgements

## 7. Bibliographical References

Carolyn Betensky. 2019. Casual racism in victorian literature. *Victorian Literature and Culture*, 47(4):723–751.

Ryan Brate, Marieke van Erp, and Antal van den Bosch. 2023. Contextual profiling of charged terms in historical newspapers. *Language, Data and Knowledge (LDK 2023). 12-15 September, Vienna, Austria*.

Shijing Chen, Usman Naseem, and Imran Razzak. 2023. Debunking biases in attention. In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 141–150, Toronto, Canada. Association for Computational Linguistics.

cjadams, Daniel Borkan, inversion, Jeffrey Sorensen, Lucy Vasserman Lucas Dixon, and nithum. 2019. Jigsaw unintended bias in toxicity classification.

Louise Denman-Sparks. 1993. *Anti-bias curriculum*, 7 edition. National Association for the Education of Young Children, Washington, D.C., DC.

Ted Dunning. 1993. Accurate methods for the statistics of surprise and coincidence. *Comput. Linguist.*, 19(1):61–74.

Elisabetta Fersini, Francesca Gasparini, Giulia Rizzi, Aurora Saibene, Berta Chulvi, Paolo Rosso, Alyssa Lees, and Jeffrey Sorensen. 2022. SemEval-2022 task 5: Multimedia automatic misogyny identification. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 533–549, Seattle, United States. Association for Computational Linguistics.

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, Dublin, Ireland. Association for Computational Linguistics.

Fanny Jourdan, Laurent Risser, Jean-michel Loubes, and Nicholas Asher. 2023. Are fairness metric scores enough to assess discrimination biases in machine learning? In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 163–174, Toronto, Canada. Association for Computational Linguistics.

Jessica Lin. 2022. Leveraging world knowledge in implicit hate speech detection. In *Proceedings of the Second Workshop on NLP for Positive Impact (NLP4PI)*, pages 31–39, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Carla Perez-Almendros, Luis Espinosa-Anke, and Steven Schockaert. 2022. SemEval-2022 task 4: Patronizing and condescending language detection. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 298–307, Seattle, United States. Association for Computational Linguistics.

Carla Perez Almendros and Steven Schockaert. 2022. Identifying condescending language: A tale of two distinct phenomena? In *Proceedings of the Second Workshop on NLP for Positive Impact (NLP4PI)*, pages 130–141, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

David Rees. 1988. Dahl's chickens: Roald dahl. *Children's Literature in Education*, 19(3):143–155.

Donald Yacovone. 2020. Roald dahl, the caribbean, and a warning from his chocolate factory. *ReVista (Cambridge)*, 20(1):1–7.