What is in a name? Mitigating Name Bias in Text Embeddings via Anonymization

Anonymous ACL submission

Abstract

Text-embedding models often exhibit biases 002 arising from the data on which they are trained. In this paper, we examine a hitherto unexplored bias in text-embeddings: bias arising from the presence of *names* such as persons, locations, organizations etc. in the text. Our study shows how the presence of name-bias in text-embedding models can potentially lead to erroneous conclusions in assessment of thematic similarity. Text-embeddings can mistak-012 enly indicate similarity between texts based on names in the text, even when their actual semantic content has no similarity or indicate dissimilarity simply because of the names in 016 the text even when the texts match semantically. We first demonstrate the presence of name bias 017 in different text-embedding models and then propose text-anonymization during inference which involves removing references to names, while preserving the core theme of the text. 021 The efficacy of the anonymization approach is demonstrated on two downstream NLP tasks, achieving significant performance gains. Our simple and training-optimization-free approach offers a practical and easily implementable solution to mitigate name bias.

1 Introduction

028

042

Text-embedding models, which convert raw text such as sentences/paragraphs into concise numerical representations, have become fundamental tools for downstream NLP tasks in fields such as healthcare, education, law and scientific research (Chrysostomou and Aletras, 2022; Reimers, 2019; Tenney, 2019; Nie et al., 2024; Sun et al., 2019). A cosine similarity between embeddings is typically used (Zhang et al., 2019; Mathur et al., 2019) although other types of similarities (Steck et al., 2024) are also possible. With a similarity measure, the goal is to find which two texts are similar to or different from one another. For simplicity, we will use *text-embedding model* to refer to models that convert text to an embedding.

043

045

047

051

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

Many text-embedding models are often trained on large amounts of Internet text. This data can inadvertently contain biases of various kinds, reflecting social prejudices and stereotypes. As a result, these models can generate biased embeddings, reinforcing harmful stereotypes or discriminating against certain cultural groups, genders, etc. (Gallegos et al., 2024; Li et al., 2023; Rakivnenko et al., 2024). Furthermore, the presence of bias in models could lead to embeddings that disproportionately emphasize particular parts of the text, consequently failing to capture the true semantics and themes within the text (Rakivnenko et al., 2024).

While important, existing studies on biases, predominantly examine biases in text-embedding models mostly related to gender, geography, race, religion etc. (Rakivnenko et al., 2024; May et al., 2019; Bolukbasi et al., 2016; Kotek et al., 2023; Nghiem et al., 2024). In this paper, we demonstrate that text-embedding models exhibit significant bias towards names within the text. To illustrate this, we begin with a motivating example in Table 1. We present a simple narrative (Story 1). We then show a similar plot while substituting the name of the main character in (Story 2). In the third narrative (Story 3), we introduce a distinct and contradicting storyline from *Story 1* while retaining the original character names. We embed all three stories using text-embedding models. We observe that the similarity between Story 1 and Story 3, despite their differing plots, is consistently higher than the similarity between Story 1 and Story 2, which share highly semantically similar plots but differ in character names. This is very counterintuitive since the text-embedding models seem to prioritize name similarity over the text's narrative structure. While this is admittedly an illustrative example, we proceed to generate numerous such narratives and conduct a thorough investigation of this bias in our experiments. We emphasize here

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

156

157

158

159

160

161

162

163

164

114

115

that our study investigates thematic and semantic similarities within textual data while acknowledging certain applications involving text tied to specific individuals or locations, our primary focus lies on the broader thematic context rather than characters in the text.

086

090

100

101

102

103

104

105

106

107

108

109

110

111

112

113

Our observation reveals a critical issue that can significantly impact applications that rely on semantic similarity, including semantic search, information retrieval, and plagiarism detection (Minaee et al., 2024; Pudasaini et al., 2024): consider the challenge of accurately assessing the similarity between two stories/plots with identical underlying meanings but distinct character names. Current methods may erroneously classify these stories as dissimilar, leading to inconsistent and unreliable results. Further, based upon our investigation, we would like to mention upfront that the issue is not confined to certain cultures, cross-culture, but is universal in the sense that the name bias issue occurs in a very broad sense.

Story Id	Text
Story 1	Alejandro gently examined the injured bird. He gave it food.
Story 2	Jelani tenderly inspected the wounded bird and gave it a meal
	to eat.
Story 3	Alejandro tracked the injured bird. He used it as his food.

Model	Cosine Similarity	
	Story1, Story 2↑	Story 1, Story 3 \downarrow
all-mpnet-base-v2	0.755	0.778
all-distilroberta-v1	0.780	0.798
all-MiniLM-L6-v2	0.660	0.853
gemini	0.864	0.848
multi-qa-distilbert-cos-v1	0.579	0.907
paraphrase-MiniLM-L6-v2	0.775	0.855
distiluse-base-multilingual-cased-v1	0.752	0.889
distiluse-base-multilingual-cased-v2	0.742	0.875
paraphrase-multilingual-MiniLM-L12-v2	0.836	0.840
msmarco-distilbert-cos-v5	0.584	0.817
multi-qa-mpnet-base-cos-v1	0.694	0.854
voyage-3-lite	0.780	0.868
text-embedding-3-small	0.755	0.826
text-embedding-3-large	0.741	0.808
all-mpnet-base-v2 all-distilroberta-v1 all-MiniLM-L6-v2 gemini mult-qa-distilbert-cos-v1 paraphrase-MiniLM-L6-v2 distiluse-base-multilingual-cased-v1 distiluse-base-multilingual-cased-v2 paraphrase-multilingual-cased-v2 paraphrase-multilingual-cased-v2 msmarco-distilbert-cos-v5 multi-qa-mpnet-base-cos-v1 voyage-3-lite text-embedding-3-small text-embedding-3-large	$\begin{array}{c} 0.755\\ 0.780\\ 0.660\\ 0.864\\ 0.579\\ 0.775\\ 0.752\\ 0.742\\ 0.836\\ 0.584\\ 0.694\\ 0.780\\ 0.755\\ 0.741\\ \end{array}$	0.778 0.798 0.853 0.848 0.907 0.855 0.889 0.875 0.840 0.817 0.854 0.854 0.826 0.808

Table 1: **Impact of names on similarity**: We see that Story 1 is similar to Story 2 but has different person names(*Alejandro*, *Jelani*). Story 3 is different from Story 1 but has same name (*Alejandro*) as Story 1. We observe that, in most embedding models a different story with opposite meaning and same name(*Alejandro*) is getting a higher similarity score in comparison to the same story with different names.

Having briefly revealed the issue of name bias in text-embedding models, we outline our contributions in the work:

First, we identify bias arising from names in textual content. Although several forms of biases have been studied in the past (see Sec.2), to the best of our knowledge, our work is the first that specifically looks at bias associated with names and how they can influence the embeddings coming out of embedding models. Toward this end, we propose a benchmarking study to comprehensively assess this bias.

Second, we propose a simple *inference-time textanonymization* technique designed to overcome the identified bias. Our method does not require any model fine-tuning or retraining of the textembedding models. The approach offers a simple, intuitive, and effective way to mitigate the problem rather than relying on complex computations.

Third, we conducted extensive experiments to study the identified problem in detail on a variety of text-embedding models and tasks. Our results demonstrate that our anonymization approach effectively reduces name bias within embeddings in semantic similarity and downstream tasks.

2 Related Work

Biases in Text-embedding models: Textembedding models while powerful, can inadvertently reflect and amplify existing biases and prejudices; there is vast research understanding and mitigating bias in such models. For example, there is work focusing on models that investigate under-representation or misrepresentation of specific groups, such as LGBTQ+ individuals, leading to skewed or inaccurate outcomes (May et al., 2019; Bolukbasi et al., 2016; Cheng et al., 2021). Another type of study focuses on gender bias in word embeddings models (Rakivnenko et al., 2024). The study highlights a concerning issue i.e many embedding models associate specific occupations with particular genders. Nikolaev and Padó (2023) studied biases at a sentence-level in sentence transformers influenced by different parts of speech such as common nouns, adverbs etc. While we discuss text-embedding model, it is important to highlight works that investigate bias within Large Language Models (LLMs) for text-generation which are a part of this ecosystem (Gallegos et al., 2024). Schwöbel et al. (2023) observed "geographical erasure" where certain regions are underrepresented in LLM outputs. Manvi et al. (2024) showed that LLMs often favor developed regions and exhibit negative biases towards locations with lower socioeconomic conditions, particularly on subjective topics such as attractiveness and intelligence. Further, some works have also investigated cross-cultural biases in LLMs for text generation (Naous et al., 2023; Ramezani and Xu, 2023; Cao et al., 2023; Arora et al., 2022). Compared to the above work, we investigate name-bias in text-embeddings, an area

- not previously explored in existing research to the 165 best of our knowledge. 166
- Debiasing methods: Various approaches have 167 been proposed to tackle different kinds of biases in 168 text-embedding models highlighted above. One 169 common technique to remove such biases is to update the training dataset and make it unbiased 171 and re-train the model (Brunet et al., 2019; Ngo 172 et al., 2021). Another paradigm involves applying approaches such as disentanglement or alignment 174 where models are fine-tuned to remove biases as-175 sociated with concepts such as gender, religion 176 etc. (Kaneko and Bollegala, 2021; Guo et al., 2022; 177 Kenneweg et al., 2024). An alternative approach involves post-processing of the embeddings. Specifi-179 180 cally, it involves adding a debiasing module after encoders to filter out certain biases in the represen-181 tations Cheng et al. (2021). For more details on this topic, we refer the reader to survey by Li et al. (2023) for more details.

We emphasize some key considerations based upon the discussions above. Firstly, all the aforementioned techniques require an optimization phase, involving either retraining the initial model, fine-tuning with a modified loss or postprocessing of the generated embeddings. Secondly, these methods are often designed to address specific bias types, such as social, gender, or religious biases. Notably, the identification and mitigation of name bias has not been previously explored to our knowledge.

Understanding name bias 3

187

188

192

193

194

195

196

197

198

199

200

205

206

210

In this section, we investigate the presence of bias within text-embedding models related to names. Our primary objective is to investigate the influence of names containing identity-specific information on the resulting text embeddings, while ensuring the semantic structure of the text remains unchanged.

3.1 **Benchmarking Methodology**

To understand the impact of bias associated with names, we systematically replace instances of names in text with alternatives. For the sake of simplicity, in this section, we focus on person names and country names¹. Given a text, we first identify instances of person and country names

in the text.² To study bias w.r.t. person names, we replace each person name in the text with a randomly sampled name from a list of person names. In the text, all instances of the same person are replaced by the same sampled name. Similarly, country names are replaced with a random country name sampled from a predefined list of countries. This process only changes the person names and countries and does not change the original structure or meaning of the text.

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

239

240

241

242

Formally, given a universe of n person names $P = \{p_1, p_2, p_3 \cdots p_n\}$, and l Country names C = $\{c_1, c_2, c_3 \cdots c_l\}$, we apply algorithm 1 for a given text T to obtain a perturbed text T'.

Algorithm 1 Perturb Text for Benchmarking

- **Require:** P : List of Person names, C : List of Country names. 1: Input: Text T
- 2: **Output:** Text T' with replaced entities
- 3: Initalize: $T' \leftarrow T$ 4: Identify Entities:
- Identify all occurrences of person names IP in T'5:
- Identify all occurrences of country names IC in T'. 6:
- 7: Perturbation:
- 8: for each identified person $ip \in IP$ in text T' do
- 9: Randomly select a name $p_k \in P$ without replacement.
- 10: Replace all occurrences of ip with p_k in text T'.
- 11: end for
- 12: for each identified country $ic \in IC$ in text T' do
- 13: Randomly select a country name $c_k \in C$ without replacement.
- 14: Replace all occurrences of *ic* with c_k in text T'.
- 15: end for
- 16: **Return** T' {Perturbed Text}

Applying Algorithm 1 gives one perturbation T' for text T. We generate K=20 such perturbations capturing a wider range of person and country names. The names used for replacement are present in Table 9 in Appendix and we have names from many different cultures/countries. Note that the steps 8-11 and 12-15 in perturbation algorithm can be done in isolation and can be applied independently based upon the use-case. An illustrative example of a perturbation is presented in Table 2.

The objective is to determine the degree of semantic divergence observed between perturbed text instances, resulting from the replacement of names and countries, by examining their embeddings. As discussed above, for a text T we create Kperturbations $\{T'_i \mid 1 \le i \le K\}$. Each of these K perturbed text versions were processed through a text-embedding model, to obtain its corresponding

¹We also study the impact of perturbation of person names only.

²The datasets used for benchmarking are described in Sec. 3.3

Original $\text{Text}(T)$	Perturbed Text $1(T'_1)$
Mike has been living in	Dwayne has been living in
Belgium for five years and	France for five years and
made a fortune by winning	made a fortune by winning a
a lottery. Mike spent most	lottery. Dwayne spent most
of his money on treatment of	of his money on treatment of
his brother Donald who was	his brother Shawn who was
suffering from Lung Cancer.	suffering from Lung Cancer.

Table 2: Example of text perturbation.

243 embedding. Subsequently, to capture the distance between the perturbed text's embeddings with 244 each other, we calculate the pairwise cosine similarity (Pedregosa et al., 2011) between all K246 embeddings. For example, if a text sample has 247 K=20 perturbations, we get $\frac{K \times (K-1)}{2} = 190$ similarity scores. Given N such text samples 249 in a dataset, to arrive at a single metric, we first compute pairwise cosine similarities(between 251 the perturbed text embeddings) for a given text, 252 excluding the self-similarity comparisons (i.e., the similarity of a perturbed text embedding to itself). For N samples, we obtain $N \times \frac{K \times (K-1)}{2}$ similarity 255 scores. Let emb_{si} refer to the embedding of i^{th} perturbation of sample s where $s \in \{1, 2, ..., N\}$ and $i \in \{1, 2, ..., K\}$. Then, average similarity across N samples is defined as: -

$$\frac{1}{N \times \frac{K(K-1)}{2}} \sum_{s=1}^{N} \left[\sum_{\substack{i=1\\j \neq i}}^{K} \sum_{\substack{j=1\\j \neq i}}^{K} Sim(emb_{si}, emb_{sj}) \right]$$

260

263

267

269

270

271

A higher average similarity indicates that the perturbed texts are closer to each other in the semantic space, suggesting less deviation. Conversely, a lower average similarity score suggests a higher degree of deviation from the expected semantic relationship. It suggests that the embedding model exhibits a bias towards names in the text, potentially affecting its ability to accurately capture the theme of the text.

3.2 Candidate Text-embedding Models

We analyzed a diverse set of leading text embed-272 ding models from academia and industry. This 273 includes models explicitly trained on diverse lan-274 guages and tasks such as semantic search, questionanswering etc. We include models such as multi-qa-277 distilbert-cos-v1 and multi-qa-mpnet-base-cos-v5 for question answering, and paraphrase-MiniLM-278 L6-v2 and paraphrase-multilingual-MiniLM-L12-279 v2 for identifying semantic similarity (Reimers, 2019). Other notable models include all-mpnet-281

base-v2, all-distilroberta-v1, and *all-MiniLM-L6-v2*, designed for general-purpose text representation (Reimers, 2019). Additionally, multilingual models like *distiluse-base-multilingual-cased-v1* and *distiluse-base-multilingual-cased-v2* are also included (Reimers and Gurevych, 2020). We also include *msmarco-distilbert-cos-v5* specialized model for search (Reimers, 2019). Additionally, we also choose cutting-edge models which are not open-source namely *text-embedding-3-small* and *text-embedding-3-large* from Open AI (OpenAI, 2024), *gemini* from Google (Team et al., 2023) and *voyage-3-lite* from Voyage AI (AI, 2024). 282

283

284

287

289

290

291

292

293

294

295

296

297

298

299

300

301

303

304

305

306

307

308

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

3.3 Benchmark Datasets

CMU Movie Dataset (Bamman et al., 2013): The CMU Movie dataset primarily consists of 6,559 textual plot summaries of movies spanning multiple sentences. These summaries are typically short , concise descriptions of the main events and story-lines within a film. They often include key characters, conflicts, and resolutions.

CMU Book Dataset (Bamman and Smith, 2013): Similar to CMU Movie, the core of this dataset consists of concise multiple sentence summaries of 42,306 books. These summaries capture the main plot points, key characters, and themes.

We select plots where the number of words are less than 250 which is within token limit of most models under consideration³.

3.4 Analyzing Bias

In Table 3 and 4 we observe a significant deviation in the average cosine similarity which should be close to one if the cosine similarity captured the real semantic similarity rather than information in names present in the text⁴. Any deviation from one indicates that the embeddings are heavily biased by the choice of names rather than from the similarity of the text. Models like msmarco-distilbert-cos-v5 exhibit significant sensitivity to changes in person and country names, as evidenced by an average cosine similarity ≈ 0.7 . This suggests that the model's embeddings may be heavily influenced by specific entities rather than capturing the underlying semantic meaning of

³For each embedding model, we evaluate its performance only on samples which are within the limits of its maximum context window.

⁴We also experimented by using euclidean distance instead of cosine similarity in Tab. 19 in Appendix. The conclusion remained similar and therefore we proceeded with cosine similarity for remaining experiments.

Model Name	Cosine sim per perturbation pair
all-mpnet-base-v2 all-distilroberta-v1 all-MiniLM-L6-v2 gemini multi-qa-distilbert-cos-v1 paraphrase-MiniLM-L6-v2 distiluse-base-multilingual-cased-v1 distiluse-base-multilingual-cased-v2	$\begin{array}{c} \text{per perturbation pair} \\ 0.774 \pm 0.001 \\ 0.768 \pm 0.001 \\ 0.706 \pm 0.001 \\ 0.885 \pm 0.0 \\ 0.733 \pm 0.001 \\ 0.742 \pm 0.001 \\ 0.786 \pm 0.001 \\ 0.795 \pm 0.001 \\ \end{array}$
paraphrase-multilingual-MiniLM-L12-v2 msmarco-distilbert-cos-v5 multi-qa-mpnet-base-cos-v1 text-embedding-3-small text-embedding-3-large voyage-3-lite	$\begin{array}{c} 0.75 \pm 0.001 \\ 0.681 \pm 0.001 \\ 0.743 \pm 0.001 \\ 0.742 \pm 0.0 \\ 0.779 \pm 0.0 \\ 0.76 \pm 0.0 \end{array}$

Table 3: **Bias Measurement on CMU Movie dataset**. For each show, we create K=20 **perturbations** by replacing person names and country names. In this experiment, we used plot samples that contain both person and country names but does not mention any city/town/village/nationality keywords(Spanish, American etc.) in order to minimize the impact of other variables. We report the mean and the std. error rounded off to 3 decimal places.

Model Name	Cosine sim per perturbation pair
all-mpnet-base-v2	0.777 ± 0.001
all-distilroberta-v1	0.778 ± 0.001
all-MiniLM-L6-v2	0.693 ± 0.001
gemini	0.89 ± 0.0
multi-qa-distilbert-cos-v1	0.743 ± 0.001
paraphrase-MiniLM-L6-v2	0.735 ± 0.002
distiluse-base-multilingual-cased-v1	0.777 ± 0.001
distiluse-base-multilingual-cased-v2	0.785 ± 0.001
paraphrase-multilingual-MiniLM-L12-v2	0.746 ± 0.002
msmarco-distilbert-cos-v5	0.707 ± 0.001
multi-qa-mpnet-base-cos-v1	0.75 ± 0.001
text-embedding-3-small	0.761 ± 0.001
text-embedding-3-large	0.795 ± 0.001
voyage-3-lite	0.781 ± 0.001

Table 4: Bias Measurement on CMU Books dataset.We follow the same evaluation setup as in Table 3.

the text. Observations from the evaluation of both datasets suggest that *gemini* is the least biased model among all models considered. However, we observe that even *gemini's* score is still far away from one indicating room for improvement.

In the above experiment, we replaced the names of people and countries and generated a perturbed text. One may ask: how much of the bias is from country name versus person names? To study this, we considered an experiment in which we perturbed the text by only replacing person names while keeping the country names as they were in the original text. We also examined variations in which all the perturbed names are sampled from the same country and demonstrate that bias persists even if text samples differ only by person names even from the same country. These results can be found in the App. B.

4 Methodology: Overcoming Bias through Anonymization

Previously, we showed that how just changing person names/country names can impact the embeddings significantly. In this section, we introduce a simple *inference-time anonymization* technique to mitigate the bias caused by names. The core idea is to mitigate the influence of names on embeddings, and making the resulting *debiased* anonymized embeddings to be more generalizable and less prone to biases related to particular individuals or entities.

The anonymization of a text T during inference is achieved through the following process. We first identify in T, occurrences of desired entities such as person names, locations and organizations relevant to the use case. We anonymize the text by removing those occurrences from T. The anonymized text referred to as T_{anon} retains the overall structure and meaning of the original text T while removing any specific references to person names etc. This anonymization can be achieved via tools such as Large Language Models(LLMs) (Zhao et al., 2023) or Named Entity Recognition tools (Jehangir et al., 2023). In our work, we used gemini and anthropic.claude-3-5sonnet text-generative models for anonymization using prompts. Depending upon the use-case, different names in text such as person names, cities, countries, organizations can be removed. We would like to clarify that the same process of anonymization can also be done through Named-Entity Recognition(NER) tools (Jehangir et al., 2023), however in our initial experiments we found LLMs to be more accurate. Sample prompts for anonymization are presented in Table 5. Post anonymization, the embeddings become independent of identity specific details such as person names/ country names etc.⁵ Overall, the *debiased* embeddings generated on anonymized text promise reduced sensitivity to biases associated with particular individuals or entities. Note that the embeddings generated for sentences that differ solely in their named entities (e.g., character names) will now have a cosine similarity of 1. An alternate to removing named content for anonymization is to replace names with specific non-identifying placeholder words. This approach with its associated challenges is further examined in App. F.

346

347

348

349

350

351

352

354

355

356

358

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

384

385

387

390

391

⁵The type of anonymization i.e removing person names and/or country names and/or city names etc. used determines the exact level of independence.

Purpose	Prompt
Remove person	Given below text, please COMPLETELY DELETE
names and loca-	all Person/Character names which are PROPER
tion names	NOUNS and City/ Country/ Village/ Town/ Conti-
	nent/ River/ Organization names which are PROPER
	NOUNS etc. Wherever they occur replace with
	empty string. Completely remove them and not any-
	thing else. Do not delete monument/landmark names
	like Eiffel tower etc. Do not remove He/She/him/her
	etc Output contains the modified text only The
	text is provided below ::::
Remove person	Given below text, please COMPLETELY DELETE
names only	all Person/Character names which are PROPER
	NOUNS. Wherever they occur replace with empty
	string. Completely remove them and not anything
	else. Do not remove He/She/him/her etc Output
	contains the modified text only The text is pro-
	vided below ::::

Table 5: Prompts for Anonymization. In our experiments, we select the first prompt. Based upon the use case, the suitable prompt can be selected or modified accordingly.

5 Can anonymization help in down-stream tasks that use similarity from text-embedding models?

393

396

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

494

425

In this section, we investigate the performance of the anonymized text embeddings on two downstream tasks. Both the tasks are based on obtaining a similarity score between pieces of texts. These tasks are primarily based upon semantic similarity which find applications in areas such as information retrieval, clustering, plagiarism detection, question answering etc. (Reimers, 2019). The two tasks that we evaluate on differ in various aspects such as the nature of the task, evaluation methodology, the judgment score available, etc. On both these tasks, our experiments show that embeddings based on anonymized text can significantly help in downstream tasks.

5.1 Task 1: Semantic Similarity Between Query and Text-Pairs with Binary Labels.

Recall from Sec. 3 that altering only the names/locations in two otherwise identical stories/paragraphs significantly impacted their text embeddings. In this section, we investigate whether anonymization technique proposed in Sec. 3.4 can effectively mitigate this type of bias. Towards this, we explore the Semantic Similarity Task (STS).

Semantic similarity seeks to determine the degree to which two pieces of text convey similar meaning (Muennighoff et al., 2022; Reimers, 2019). This goes beyond simple word matching, aiming to understand the underlying meaning within the text. In today's era of deep learning (Reimers et al., 2016; Muennighoff et al., 2022), achieving accurate semantic similarity relies heavily on highquality embeddings, which represents sentences as dense vectors in a continuous space.

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

In this experiment we investigate whether the text-embeddings are able to capture the semantic nuances within the text or are they biased towards names? Ideally, a good embedding model should be able to differentiate reasonably well between two stories/paragraphs which have very different themes even if they contain same names. То investigate this, we create a dataset of 10 paragraph triplets. Each triplet includes a *query* paragraph, a *positive* paragraph that is *highly semantically* similar but with distinct person and location names, and a negative paragraph that is semantically dissimilar to the query text but has same person names/location names as in query text. For each triplet, (query, positive) pair is assigned a label 1(positive) and (query, negative) pair is assigned a label 0(negative). Two sample examples can be found in Table 16 in the the rows marked as Original. The entire set of generated triplets with labels are present in Appendix D. We evaluate the performance of different models on the STS task using AUC ROC score between cosine similarity scores of embeddings and the ground truth.

Peformance on Semantic Similarity. Tab. 6 presents the AUC-ROC scores for different models on the STS task. The results indicate that the AUC scores for the majority of models are significantly below 0.5. This finding suggests a critical issue, as even a random classifier would be expected to achieve an AUC score of approximately 0.5. The fact that most of the AUC is much below 0.5 suggests that the cosine similarity based ranking got the ordering wrong! Gemini's AUC is better than random, however, it also gets improved significantly after anonymization. Such low AUC scores strongly imply that the embeddings used in these models are primarily capturing identityrelated information, leading to a significant bias in the model's embeddings and predictions. Next, we observe that the AUC-ROC results post anonymization. We see that anonymization can improve the model's capacity to grasp the core semantic meaning in the text as reflected in the significantly higher AUC-ROC numbers(closer to 1). Additionally, it is important to note that all models attain high AUC scores when all stories share identical names. This indicates that the models can effectively distinguish between sentences conveying the same or different meanings when identity information remains

Model	Original text: The (query, positive) paragraphs share the same meaning but different person/location names. The (query, negative) paragraphs share dif- ferent meaning but same person/location names.	<i>Identical Names:</i> The (<i>query</i> , <i>positive</i> , <i>negative</i>) paragraphs in the same triplet contain the same person/location names.	Anonymized text: Anonymization applied to (query, positive, negative) paragraphs.
all-mpnet-base-v2	0.19	0.96	0.98 ± 0.0071
all-distilroberta-v1	0.36	0.97	0.975 ± 0.0106
all-MiniLM-L6-v2	0.09	0.94	0.99 ± 0.0071
gemini	0.71	1.00	1.0 ± 0.0
multi-qa-distilbert-cos-v1	0.07	0.97	0.97 ± 0.0071
paraphrase-MiniLM-L6-v2	0.14	0.98	0.98 ± 0.0
distiluse-base-multilingual-cased-v1	0.27	0.95	0.94 ± 0.0
distiluse-base-multilingual-cased-v2	0.26	0.98	0.96 ± 0.0
paraphrase-multilingual-MiniLM-L12-v2	0.21	1.00	0.99 ± 0.0
msmarco-distilbert-cos-v5	0.10	0.92	0.955 ± 0.0035
multi-qa-mpnet-base-cos-v1	0.08	0.97	1.0 ± 0.0
text-embedding-3-small	0.12	1.00	1.0 ± 0.0
text-embedding-3-large	0.21	1.00	1.0 ± 0.0
voyage-3-lite	0.18	1.00	1.0 ± 0.0

Table 6: **Evaluation on Task 1: Semantic Similarity Task.** AUC scores obtained on Semantic Similarity Task. Our proposed strategy of anonymization achieves high quality results across all models. Mean and standard error are reported based on results from two separate LLM runs for anonymization.

	Query	Pos/Neg	Sim	Label
			score	
Original	Alejandro quickly ran to the store to	POS: Quickly, Hiroki dashed to the local market to	0.58	1
Oliginai	buy a cold drink. He was eager to have	procure some cold drinks. He was yearning for a chilled		
	a glass of cold driftk.	glass of cold drink.	0.60	
		NEG: Alejandro has stopped buying cold drinks from	0.09	0
		market. He only drinks cold drinks made at nome.	0.80	1
Anonymized	duickly fan to the store to buy a cold	POS: Quickly, dashed to the local market to procure	0.80	1
-	and drink	solid drinks. He was yearning for a chined glass of		
	cold driffk.	NEC: has standed by ving cold drinks from more tot. He	0.47	0
		only drinks cold drinks made at home	0.47	0
		only units cold units made at nome.		
Original	Ganga and Yamuna are two mighty	POS: Yangtze is a mighty river. It is a long river and is	0.54	1
originai	rivers. They are lifelines for millions of	the lifeline for millions of people in the region.		
	people in the region.	NEG: Ganga and Yamuna are two sisters. They had	0.70	0
		their schooling in the region and schooling provided a		
		lifeline for them.		
Anonymized	and are two mighty rivers. They are	POS: is a mighty river. It is a long river and is the lifeline	0.70	1
	lifelines for millions of people in the	for millions of people in the region.		
	region.	NEG: and are two sisters. They had their schooling in	0.56	0
		the region and schooling provided a lifeline for them.		

Table 7: Examples showing impact of anonymization on semantic similarity using embeddings created by *msmarco-distilbert-cos-v5*.

model	Spearman-correlation (Original Text)	Spearman-correlation (Anonymized)	Pearson-correlation (Original Text)	Pearson-correlation (Anonymized)
all-mpnet-base-v2	0.262	$\textbf{0.344} \pm \textbf{0.001}$	0.321	$\textbf{0.364} \pm \textbf{0.002}$
all-distilroberta-v1	0.245	$\textbf{0.327} \pm \textbf{0.007}$	0.302	$\textbf{0.37} \pm \textbf{0.003}$
all-MiniLM-L6-v2	0.251	$\textbf{0.33} \pm \textbf{0.003}$	0.282	$\textbf{0.354} \pm \textbf{0.006}$
gemini	0.381	$\textbf{0.39} \pm \textbf{0.001}$	0.456	0.436 ± 0.003
multi-qa-distilbert-cos-v1	0.240	$\textbf{0.292} \pm \textbf{0.002}$	0.269	$\textbf{0.316} \pm \textbf{0.007}$
paraphrase-MiniLM-L6-v2	0.283	$\textbf{0.352} \pm \textbf{0.005}$	0.317	$\textbf{0.37} \pm \textbf{0.0}$
distiluse-base-multilingual-cased-v1	0.282	$\textbf{0.356} \pm \textbf{0.001}$	0.325	$\textbf{0.386} \pm \textbf{0.002}$
distiluse-base-multilingual-cased-v2	0.308	$\textbf{0.357}\pm\textbf{0.0}$	0.345	$\textbf{0.389} \pm \textbf{0.003}$
paraphrase-multilingual-MiniLM-L12-v2	0.261	$\textbf{0.332} \pm \textbf{0.001}$	0.281	$\textbf{0.364} \pm \textbf{0.004}$
msmarco-distilbert-cos-v5	0.232	$\textbf{0.304} \pm \textbf{0.002}$	0.262	$\textbf{0.333} \pm \textbf{0.005}$
multi-qa-mpnet-base-cos-v1	0.274	$\textbf{0.324} \pm \textbf{0.002}$	0.317	$\textbf{0.354} \pm \textbf{0.001}$
text-embedding-3-small	0.374	$\textbf{0.382} \pm \textbf{0.002}$	0.416	$\textbf{0.422} \pm \textbf{0.005}$
text-embedding-3-large	0.366	$\textbf{0.382} \pm \textbf{0.007}$	0.428	$\textbf{0.429} \pm \textbf{0.012}$
voyage-3-lite	0.359	0.322 ± 0.005	0.400	0.352 ± 0.002

Table 8: **Evaluation on Task 2: Semantic similarity with graded relevance.** The table presents correlation between cosine similarity between human & machine summaries and relevance(ground truth) provided by human evaluators . Mean and standard error are reported based on results from two separate LLM runs for anonymization.

constant. The aforementioned observations highlights that anonymization is crucial to avoid situa-

477

478

tions where semantically equivalent paragraphs are assigned unique embeddings solely based on the 481 presence of identity information (such as names).
482 Conversely, it's essential that when texts have significant semantic variations, even if they contain
484 identical identity information, their embeddings are
485 able to able to capture it.

486

487

488

489

490

491

492

493

494

495 496

497

498

499

500

501

503

505

507

509

510

511

512

513

514

515

516

517

518

519

520

522

524

526

530

Examples of similarity post-anonymization. In Tab. 7, we show some instances of how similarity values between embeddings change between (query, positive) pair and (query, negative) pair post anonymization. Before anonymization, the models assigned higher similarity scores to negative pairs and lower similarity scores to positive pairs in a counterintuitive way. Anonymization resulted in the models predominantly attending to the semantic structure of the text, which is accurately reflected in similarity scores. We would like to highlight that these samples are a subset of examples used for AUC computation on the STS task in Tab. 6.

5.2 Task 2: Semantic Similarity With Graded Human Relevance.

In the previous task, a binary approach was employed to assess text pair similarity, categorizing text-pairs as either similar or dissimilar. In the task proposed in this section, we employ a more refined approach for evaluation by utilizing a graded relevance scale from 1 to 5 between a pair of text. The graded scale enables a more nuanced and granular assessment of semantic similarity between pairs, providing a richer understanding of their relationship. To evaluate this, we use the machine summary evaluation task from Muennighoff et al. (2022), which involves automatically assessing the relevance of machine-generated summaries, commonly assessed by calculating the semantic similarity between the embeddings of the summary and the original document/human summaries.

For this task, we follow the same evaluation setup as Muennighoff et al. (2022) which we describe next. We use the SummEval dataset (Fabbri et al., 2021; Muennighoff et al., 2022) with 100 text samples, each containing 16 machine and 10 human summaries. Human relevance scores (1-5) are assigned to each machine summary. We first obtain summary embeddings using textembedding models for each machine summary and human summary in all 100 samples. Without loss of generality, for a given text sample out of the 100 samples, for each machine summary $\{m_i \mid 1 \le i \le 16\}$, we get its predicted score based on its maximum cosine similarity to any human summary $\{h_j \mid 1 \le j \le 10\}$ within the same text sample i.e $machine_pred_score(m_i) =$ $max_{1\le j\le 10} \ cos_sim(m_i, h_j)$. This yields 16 machine summary quality predicted scores for each sample i.e 1 predicted score for each machine summary. Further, as mentioned earlier, we have a human relevance score assigned to each machine summary. Overall, across all text samples, we get 1600 *machine summary predicted scores* and its corresponding *human relevance scores*. We then correlate these two scores using Pearson and Spearman coefficients (Muennighoff et al., 2022). Higher correlations indicate better alignment between modelassigned scores and human judgments, suggesting more reliable evaluation.

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

Impact of Anonymization Table 8 shows that post-anonymization, the performance of various text-embedding models significantly improves in predicting graded human-rated summary quality. Spearman and Pearson correlation coefficients increase substantially, indicating that the model's assessment of summary quality after anonymization better aligns with human evaluations. This improvement is substantial, with some models like *all-distilroberta-v1* showing a performance increase of around 30%.

In summary, the results of both downstream tasks demonstrate a substantial enhancement in the semantic similarity post-anonymization.

6 Conclusion

In this work, we highlight the bias in text embeddings stemming from the presence of names in the text. We showed concrete examples, over multiple text-embedding models, that similarities between embeddings can be dominated by names in the text rather than the semantic meanings of the text. We then proposed a method to mitigate bias by performing anonymization at inference time. This involved the removal of names from the text and using the anonymized text to create the embeddings. Our findings demonstrate that anonymized text embeddings significantly outperform deanonymized text embeddings on tasks involving semantic similarity. While we proposed one way to mitigate the issue through anonymization, a deeper question that remains is: how to train text-embedding models such that the embeddings capture the semantics more than the names in the text?

665

666

667

668

669

670

671

672

673

674

675

676

677

678

624

625

626

7 Limitations

579

580

581

585

586

587

593

597

598

599

601

610

Below we discuss the limitations of the proposed work.

 In this work we focused on evaluating/mitigating name bias in text-embedding models using texts from English language. The work presented here does not cover other languages. Further, the work also does not cover name bias issues arising in multi language texts.

2. While our proposed anonymization solution enhances thematic similarity, it is not ideal for situations requiring the preservation of identity that we are removing through anonymization. A partial and straightforward solution might involve anonymizing only non-critical identifying information depending upon the use-case. Many real world use cases may require dynamically balancing identity and thematic preservation to suit the specific needs of each use case.

3. In our work, we adopted similarity between text-embeddings as a proxy for their semantic similarity. While commonly used, it is still an estimate of semantic similarity and may overlook deeper semantic relationships that require reasoning. A limitation of this work is that we capture thematic similarity only to the extent that it is captured by the cosine similarity (and the Euclidean distance similarity is studied in the Appendix).

8 Broader Impact

This research uncovers name-bias in text-611 embedding models. It reveals how the presence of names can skew similarity judgments, leading 613 to incorrect conclusions about thematic similarities. 614 This impacts a wide range of NLP applications, 615 potentially compromising accuracy in tasks from 616 information retrieval to sentiment analysis. The major impact of this paper is uncovering such bias 618 and how it can be mitigated at inference time. This 619 work contributes to inspiring further investigation into building more robust text-embedding models.

622 References

623 Voyage AI. 2024. Embeddings.

- Arnav Arora, Lucie-Aimée Kaffee, and Isabelle Augenstein. 2022. Probing pre-trained language models for cross-cultural differences in values. *arXiv preprint arXiv:2203.13722*.
- David Bamman, Brendan O'Connor, and Noah A Smith. 2013. Learning latent personas of film characters. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 352–361.
- David Bamman and Noah A Smith. 2013. New alignment methods for discriminative book summarization. *arXiv preprint arXiv:1305.1319*.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.
- Marc-Etienne Brunet, Colleen Alkalay-Houlihan, Ashton Anderson, and Richard Zemel. 2019. Understanding the origins of bias in word embeddings. In *International conference on machine learning*, pages 803–811. PMLR.
- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. Assessing cross-cultural alignment between chatgpt and human societies: An empirical study. *arXiv preprint arXiv:2303.17466*.
- Pengyu Cheng, Weituo Hao, Siyang Yuan, Shijing Si, and Lawrence Carin. 2021. Fairfil: Contrastive neural debiasing method for pretrained text encoders. *arXiv preprint arXiv:2103.06413*.
- George Chrysostomou and Nikolaos Aletras. 2022. Flexible instance-specific rationalization of nlp models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10545–10553.
- Alexander R Fabbri, Wojciech Kryściński, Bryan Mc-Cann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79.
- Yue Guo, Yi Yang, and Ahmed Abbasi. 2022. Autodebias: Debiasing masked language models with automated biased prompts. In *Proceedings of the* 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1012–1023.
- Basra Jehangir, Saravanan Radhakrishnan, and Rahul Agarwal. 2023. A survey on named entity recognition—datasets, tools, and methodologies. *Natural Language Processing Journal*, 3:100017.

- 679 694 697 704 706 710 711 712 713 714 715 716 718 719 720 721 723 725 726 727 728 729

- 731 732

- Masahiro Kaneko and Danushka Bollegala. 2021. Debiasing pre-trained contextualised embeddings. arXiv preprint arXiv:2101.09523.
 - Philip Kenneweg, Sarah Schröder, Alexander Schulz, and Barbara Hammer. 2024. Debiasing sentence embedders through contrastive word pairs. arXiv preprint arXiv:2403.18555.
- Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in large language models. In Proceedings of the ACM collective intelligence conference, pages 12-24.
- Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Ying Wang. 2023. A survey on fairness in large language models. arXiv preprint arXiv:2308.10149.
- Rohin Manvi, Samar Khanna, Marshall Burke, David Lobell, and Stefano Ermon. 2024. Large language models are geographically biased. arXiv preprint arXiv:2402.02680.
- Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2019. Putting evaluation in context: Contextual embeddings improve machine translation evaluation. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 2799-2808.
- Chandler May, Alex Wang, Shikha Bordia, Samuel R Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. arXiv preprint arXiv:1903.10561.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. arXiv preprint arXiv:2402.06196.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. arXiv preprint arXiv:2210.07316.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2023. Having beer after prayer? measuring cultural bias in large language models. arXiv preprint arXiv:2305.14456.
- Huy Nghiem, John Prindle, Jieyu Zhao, and Hal Daumé III. 2024. " you gotta be a doctor, lin": An investigation of name-based bias of large language models in employment recommendations. arXiv *preprint arXiv:2406.12232.*
- Helen Ngo, João G M Araújo Cooper Raterink, Ivan Zhang, Carol Chen, Adrien Morisot, and Nicholas Frosst. 2021. Mitigating harm in language models with conditional-likelihood filtration.(aug. arXiv preprint arXiv:2108.07790.
- Zhijie Nie, Zhangchi Feng, Mingxin Li, Cunwang Zhang, Yanzhao Zhang, Dingkun Long, and Richong Zhang. 2024. When text embedding meets large language model: A comprehensive survey. arXiv preprint arXiv:2412.09165.

Dmitry Nikolaev and Sebastian Padó. 2023. Representation biases in sentence transformers. arXiv preprint arXiv:2301.13039.

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

760

762

763

764

765

766

767

768

769

770

773

774

775

776

777

780

781

782

783

784

785

786

OpenAI. 2024. Embeddings.

- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12:2825-2830.
- Shushanta Pudasaini, Luis Miralles-Pechuán, David Lillis, and Marisa Llorens Salvador. 2024. Survey on plagiarism detection in large language models: The impact of chatgpt and gemini on academic integrity. arXiv preprint arXiv:2407.13105.
- Vasyl Rakivnenko, Nestor Maslej, Jessica Cervi, and Volodymyr Zhukov. 2024. Bias in text embedding models. arXiv preprint arXiv:2406.12138.
- Aida Ramezani and Yang Xu. 2023. Knowledge of cultural moral norms in large language models. arXiv preprint arXiv:2306.01857.
- N Reimers. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. arXiv preprint arXiv:1908.10084.
- Nils Reimers, Philip Beyer, and Iryna Gurevych. 2016. Task-oriented intrinsic evaluation of semantic textual similarity. In Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, pages 87-96.
- Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics.
- Pola Schwöbel, Jacek Golebiowski, Michele Donini, Cédric Archambeau, and Danish Pruthi. 2023. Geographical erasure in language generation. arXiv preprint arXiv:2310.14777.
- Harald Steck, Chaitanya Ekanadham, and Nathan Kallus. 2024. Is cosine-similarity of embeddings really about similarity? In Companion Proceedings of the ACM on Web Conference 2024, pages 887–890.
- Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In Proceedings of the 28th ACM international conference on information and knowledge management, pages 1441-1450.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of

highly capable multimodal models. *arXiv preprintarXiv:2312.11805*.

789

790

791

792

793

794

- I Tenney. 2019. Bert rediscovers the classical nlp pipeline. *arXiv preprint arXiv:1905.05950*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A
 survey of large language models. *arXiv preprint arXiv:2303.18223*.

805

806

807

810

811

812

813

814

815

816

817

818

819

820

821

824 825

826

830

831

834

835

837

A Names used for perturbation in Benchmarking

Table 9 presents the universe of names used for perturbation in the benchmarking experiment in Sec. 3. These names represent a diverse range of geographies.

B Bias measurement with only person name perturbations

In the benchmarking study in Sec. 3, we investigated the divergence in text embeddings when person names and locations were perturbed. In this section, we examine the impact of replacing only person names on the text embeddings.

B.1 Perturbations of only person names

In this study, we only perturb person names and keep the location names unchanged to understand the impact of only perturbing person names. As shown in Table 10, performing only person name perturbations on book plots also reveals a significant drop in cosine similarity across all evaluated models.

B.2 Person name perturbations on text samples without mention of country/city/town names

In this section, we investigate impact of person name perturbations when using samples which don't have mention of any country/city/town etc. names. The objective is to minimize the impact of these variables and study divergence solely w.r.t person names. As shown in Table 11, benchmarking on the CMU Book dataset on samples without having any mention of country/city/town etc. reveals a significant drop in cosine similarity across all evaluated models when only person names are perturbed.

B.3 Bias measurement with person name perturbations from the same geographical area

In previous studies, we perturbed names by replacing them from a diverse set of person names. In this study we investigate whether the issue of divergence in embeddings persists when all the perturbed names are from the same geography. This study aims to minimize the impact of cultural differences in analysis in text-embeddings. Table 12 shows the country wise names used for benchmarking. In tables 13, 14, and 15, we observe that the divergence issue persists even when the replaced847names belong to the same geography. This demon-848strates that the issue is not present in names from849certain cultures, cross-culture, but is universal in850the sense that the name bias issue occurs in a very851broad sense.852

C Similarity Heatmaps

853

855

856

857

859

861

862

864

868

870

871

873 874

876

878

879

880

881

In this section, we show examples of cosine similarity heatmaps based upon embeddings generated by different text-embedding models. We use the following example:

CHARACTER_NAME, a seasoned physician, meticulously analyzed a patient's intricate heart condition. He later realised she was his school friend.

To obtain different perturbations, we replace "CHARACTER_NAME" with different person names and generate embedding for each of the perturbation. The similarity heatmaps are present in figs. 1 to 4. The heatmaps clearly reveal that only changing the person names can significantly impact the text embeddings. This suggests that the text embedding model is highly sensitive to the specific names used within the text, even when the overall context and meaning remains completely unchanged. These kind of variations can lead to misleading results in various downstream tasks. For example, if the goal is to cluster documents into topics, changing the person names could lead to different clusters being formed, even if the underlying topics are the same. Similarly, if the text embedding model is used to classify documents as positive or negative, changing the person names could lead to different classifications being assigned, even if the overall sentiment and theme of the text remains the same.



Similarity Heatmap for Model: paraphrase-multilingual-MiniLM-L12-v2

Figure 1: Cosine Similarity Heatmap with paraphrase-multilingual-MiniLM-L12 model for example in Sec. C

Person names	Aaron, Adrian, Aiden, Akira, Alex, Alexander, Alfred, Anders, Andreas, Andrew, Anthony, Archer,
	Arthur, Ayden, Benjamin, Bernard, Blake, Boris, Bradley, Brandon, Brayden, Brian, Caleb, Cameron,
	Carlos, Carl, Charles, Charlie, Christopher, Connor, Cooper, Daichi, Daniel, David, Dean, Dennis,
	Dylan, Edward, Elijah, Elliot, Emil, Eric, Ethan, Evan, Ezra, Fabian, Felix, Finn, Francis, Gavin,
	George, Giovanni, Gregory, Haakon, Han, Harry, Hayden, Henry, Hiroki, Hugo, Hunter, Ian, Isaac,
	Ivan, Jack, Jacob, Jake, James, Jason, Jasper, Jayden, Jeremy, Jesse, Jin, Joaquim, Johan, John,
	Jonathan, Jordan, Joseph, Joshua, Juan, Kai, Kaiden, Kazuma, Keanu, Ken, Kenneth, Kevin, Liam,
	Logan, Lucas, Luis, Luke, Luka, Magnus, Mark, Martin, Mateo, Matthew, Max, Maximilian, Michael,
	Mikael, Nathan, Nathaniel, Nicolas, Noah, Oliver, Oscar, Owen, Pablo, Patrick, Paul, Pedro, Peter,
	Phillip, Phoenix, Rafael, Rajiv, Ralf, Ramón, Raphael, Ravi, Raymond, Reuben, Richard, Robert,
	Robin, Rohan, Roland, Ronan, Ryan, Samuel, Santiago, Sebastian, Sean, Silas, Simon, Stefan, Stephen,
	Thomas, Timothy, Tyler, Victor, Vincent, Walter, William, Xavier, Yan, Yang, Yao, Youssef, Zachary,
	Zane, Zayd, Zephyr, Zidan, Zinedine, Zubin, Alistair, Anders, Arjun, Arthur, Axel, Bartosz, Ben,
	Björn, Bruno, Caleb, Caoimhín, Cillian, Cormac, Daisuke, Damien, Darius, Deniz, Dorian, Eamon,
	Emile, Enzo, Fionn, Florian, Gabriel, Gideon, Gustaf, Hassan, Héctor, Igor, Ishaan, Ivan, Jasper, Kai,
	Leo, Levi, Liam, Luca, Lucian, Luis, Magnus, Marcel, Matteo, Max, Milan, Noah, Oliver, Oscar, Otto,
	Pavel, Quentin, Rafael, Ravi, Rémy, Ren, Robin, Samuel, Santiago, Sebastian, Silas, Soren, Theo,
	Thomas, Tristan, Viktor, William, Xavier, Yannik, Zane, Aditya, Ajeet, Ajit, Akash, Amar, Amit,
	Arjun, Aryan, Ashish, Avinash, Bharat, Bhuvan, Chirag, Darshan, Dev, Dheeraj, Dhruv, Gaurav, Harsh,
	Harsha, Hemant, Ishan, Shubham, Karan, Karthik, Kumar, Manav, Manoj, Mihir, Nikhil, Niranjan,
	Nivaan, Pradeep, Pranav, Raj, Rajeev, Rahul, Ramesh, Ranjit, Ravi, Rohan, Rohit, Roop, Sachin,
	Sandeep, Sanjay, Sanket, Sarthak, Satish, Shaan, Shahrukh, Shankar, Sharad, Shivam, Siddhant,
	Siddharth, Soham, Somesh, Suresh, Tejas, Uday, Varun, Vijay, Vikram, Vinay, Vishal, Yash, Yogesh,
	Yuvraj, Adil, Amine, Anas, Fayçal, Hakim, Hicham, Mazen, Mehdi, Nassim, Rafik, Sami, Sofiane,
	Tarik, Yacine, Yassine, Abiodun, Ade, Adekunle, Adewale, Ayodeji, Chidi, Chijioke, David, Ebuka,
	Emeka, Godwin, Ikechukwu, Ikenna, Kolade, Kunle, Nonso, Obinna, Olamide, Olusegun, Onyeka,
	Paul, Peter, Samuel, Taiwo, Uche, Victor, Yemi, Yinka, Aiden, Callum, Connor, Declan, Dylan,
	Eoghan, Finn, Jack, James, Jamie, Jason, Jayden, Kian, Liam, Logan, Lucas, Luke, Mason, Max,
	Michael, Noah, Oliver, Oscar, Rory, Ryan, Samuel, Sean, Thomas, William, Charlie, Freddie, George,
	Harry, Jacob, Leo, Oliver, Oscar, Teddy, Arthur, Freddie, George, Harry, Jacob, Leo, Oliver, Oscar,
	Teddy, Aiden, Alexander, Charlie, Ethan, Jacob, James, Leo, Mason, Michael, Noah, Oliver, William,
	Benjamin, Charlie, Jacob, Leo, Oliver, Oscar, Thomas, William, Aiden, Charlie, Ethan, Jacob, Leo,
	Oliver, Oscar, Thomas, William, Shrey, Venkatesh, Nguyen, Vishwanathan, Priya, Patricia, Jennifer,
	Linda, Barbara, Susan, Camille, Sophie, Julie, Claire, Yuki, Sakura, Hana, Aiko, Emi, Li, Xiao,
	Mei, Fang, Jing, Maria, Ana, Isabel, Carmen, Dolores, Amina, Layla, Nadia, Olga, Irina, Svetlana,
	Ekaterina, Giulia, Francesca, Anna, Elena, Heidi, Greta, Lena, Marta, Sofia, Valentina, Martina, Paula,
	Clara, Laura, Mia, Emily, Sophia, Charlotte, Anita, Kavita, Lalita, Meena, Lucy, Megan, Hannah,
	Jessica, Amelia, Chloe, Manon, Lea, Elodie, Amandine, Haruka, Miyu, Rina, Yuna, Nao, Chen, Hua,
	Ling, Qing, Yan, Lucia, Pilar, Rosa, Nour, Sara, Hiba, Mona, Rania, Anastasia, Natalia, Daria, Polina,
	Vera, Mariana, Gabriela, Beatriz, Rafaela, Camila, Juliana, Evelyn, Amanda, Milla, Ines, Susana,
	Leonor, Bianca, Livia, Helena, Marina, Fernanda, Eduarda, Victoria, Andressa, Denise, Raquel, Isis,
	Elisa, Julia, Luana, Milena, Yasmin, Alessandra, Claudia, Veronica, Larissa, Bia, Silvia, Vanessa,
	Leticia, Nicole, Daniele, Eva, Alice, Milena, Leonie, Mila, Lisa, Sarah, Emma, Helena, Anja, Tina,
	Ingrid, Lucija, Noor, Samira, Dana, Kalila, Arwa, Eman, Latira, Nania, Sang, Jin, Hye, Soo, Mi, Eun,
	Yeon, Ji, Sun, Abeba, Hadia, Falou, Maimouna, Nia, Asna, Kamaria, Mira, Joan, Fiona, Leanne, Oria,
	Ava, Stobilan, Nianni, Steinia, Foppy, Lara, Fleya, Florence, Roste, Sunninei, Tvy, Sunnani, Anara, Chidiama Nazzi Sunaina Matida Harmar Willow Aarushi Anarushi Anarusa Chandri Daara Esha
	Ununnina, Ngozi, Sunaina, Matuda, Haiper, Winow, Aarushi, Ananya, Bhavna, Chandin, Deepa, Esna,
Country Names	Afghanistan Albania Algaria Andorra Angola Antigua and Barbuda Argantina Armania Australia
Country Marines	Austria Azerbaijan Bahamas Bahrain Bangladash Barbadas Balarus Balgium Baliza Banin
	Rustina, Azerbaijan, Dananias, Daniani, Dangiauesh, Darbados, Delarus, Deigiuni, Denze, Denni, Rhutan Ralivia Rospia and Harzagovina Rotewana Rrazil Rrunai Rulgaria Rustina Faco Ru
	rundi Cabo Verde Cambodia Cameroon Canada Central African Republic Chad Chile China
	Colombia Comoros Congo Costa Rica Croatia Cuba Cyprus Czech Republic Denmark Diibouti
	Dominica Dominican Republic Ecuador Fornt El Salvador Equatorial Guinea Eritrea Estonia
	Eswatini Ethionia Fiji Finland France Gabon Gambia Georgia Germany Ghana Greece Grenada
	Guatemala Guinea Guinea-Bissau Guyana Haiti Honduras Hungary Iceland India Indonesia Iran
	Iraq, Ireland, Israel, Italy, Jamaica, Japan, Jordan, Kazakhstan, Kenya, Kiribati, Kuwait, Kyroyzstan,
	Laos Latvia, Lebanon, Lesotho, Liberia, Libva, Liechtenstein, Lithuania, Luxembourg, Madagascar,
	Malawi, Malaysia, Maldives, Mali, Malta, Marshall Islands, Mauritania, Mauritius, Mexico, Microne-
	sia, Moldova, Monaco, Mongolia, Montenegro, Morocco, Mozambique, Mvanmar, Namibia, Nauru
	Nepal, Netherlands, New Zealand, Nicaragua, Niger, Nigeria, North Korea, North Macedonia, Norway.
	Oman, Pakistan, Palau, Panama, Papua New Guinea, Paraguay, Peru, Philippines, Poland. Portugal.
	Qatar, Romania, Russia, Rwanda, Saint Kitts and Nevis, Saint Lucia, Saint Vincent and the Grenadines.
	Samoa, San Marino, Saudi Arabia, Senegal, Serbia

Table 9: Universe of names used for replacement in benchmarking.



Similarity Heatmap for Model: gemini

Figure 2: Cosine Similarity Heatmap with Gemini model for example in Sec. C



Similarity Heatmap for Model: all-mpnet-base-v2

Figure 3: Cosine Similarity Heatmap with all-mpnet-base-v2 model for example in Sec. C



Similarity Heatmap for Model: text-embedding-3-large Average Similarity: 0.77

Figure 4: Cosine Similarity Heatmap with text-embedding-3-large(Open AI) model for example in Sec. C

Model Name	Cosine sim per perturbation pair
all-mpnet-base-v2	0.815 ± 0.0001
all-distilroberta-v1	0.821 ± 0.0001
all-MiniLM-L6-v2	0.749 ± 0.0002
gemini	0.91 ± 0.0
multi-qa-distilbert-cos-v1	0.787 ± 0.0001
paraphrase-MiniLM-L6-v2	0.773 ± 0.0003
distiluse-base-multilingual-cased-v1	0.843 ± 0.0002
distiluse-base-multilingual-cased-v2	0.848 ± 0.0002
paraphrase-multilingual-MiniLM-L12-v2	0.79 ± 0.0003
msmarco-distilbert-cos-v5	0.752 ± 0.0001
multi-qa-mpnet-base-cos-v1	0.795 ± 0.0001
voyage-3-lite	0.821 ± 0.0001

Table 10: Bias Measuremenent on CMU Books dataset with perturbation of person names only. For each show, we create K=20 perturbations by replacing person names. We compute the average cosine similarity for each perturbation pair and the standard error. The country/city/town names remain unchanged.

Cosine sim per perturbation pair
0.796 ± 0.0002
0.803 ± 0.0002
0.731 ± 0.0003
0.906 ± 0.0001
0.766 ± 0.0002
0.758 ± 0.0004
0.825 ± 0.0003
0.828 ± 0.0003
0.77 ± 0.0004
0.747 ± 0.0002
0.778 ± 0.0002
0.81 ± 0.0001

Table 11: Bias Measuremenent on CMU Books dataset on samples without mention of country/city/town names. Perturbation of person names only. For each show, we create K=20 perturbations by replacing person names. We compute the average cosine similarity for each perturbation pair and the standard error.

Country	Person Names		
France	Max, Tom, Léo, Noé, Paul, Jules, Hugo, Arthur, Louis, Clément, Jean-Baptiste, Jean-Pierre, Jean-Paul, Charles-		
	Henri, François-Xavier, Constantin, Gaspard, Côme, Yanis, Kilian, Maël, Thibault, Raphaël, Jérémie, Vincent,		
	Antoine, Pierre, Louis, Jacques, Baptiste, Émile, Gustave, Henri, Laurent, Marcel, Nicolas, Olivier, Pascal,		
	Quentin, Rémi, Sébastien, Théodore, Ulysse, Valentin, Wilfried, Xavier, Yves, Zacharie, Adrien, Bernard, Eva,		
	Zoé, Jade, Lou, Alice, Chloé, Léa, Lina, Louise, Éléonore, Solène, Héloïse, Camille, Marie, Jeanne, Sophie		
	Claire, Isabelle, Ambre, Lilou, Maëlys, Victoire, Clémence, Valentine, Juliette, Aurélie, Angélique, Amandine		
	Brigitte, Catherine, Delphine, Édith, Fanny, Gabrielle, Hélène, Inès, Joséphine, Karine, Laure, Manon, Nathalie,		
	Océane, Pascale, Quitterie, Rosalie, Stéphanie, Thérèse, Ursule		
India	Aarav, Aditya, Aryan, Ayush, Dev, Ishaan, Ramesh, Krishna, Mihir, Rohan, Sahir, Samarth, Shaurya, Vihaan,		
	Vrijesh, Aakash, Advait, Vinayak, Atharv, Venkatesh, Dhruv, Eshan, Hrithik, Kabir, Karan, Krish, Mahesh,		
	Nakul, Pranav, Rudra, Siddharth, Soham, Tanmay, Uday, Vaibhav, Vedant, Vikram, Yash, Yuvraj, Sachin, Ahaan,		
	Gaurav, Arjun, Daksh, Devansh, Ishan, Vishwanathan, Mayank, Parichay, Krishnanshu, Sahir, Rishi, Samyak,		
	Brajesh, Vivaan, Ayan, Rudra, Rakesh, Zain, Aarohi, Bhavya, Charvi, Devika, Eshani, Falguni, Garima, Harini,		
	Ishita, Jahnvi, Kavya, Lavanya, Madhavi, Niharika, Ojasvi, Prisha, Qara, Radhika, Saanvi, Tara, Urvashi, Vanya,		
	Wamika, Xara, Yamini, Zara, Anvi, Bhumika, Chaitali, Dharini, Ekta, Fiza, Gauri, Himani, Ira, Jiya, Kriti, Lata,		
	Meera, Nisha, Oviya, Pallavi, Rhea, Sakshi, Tanisha, Uma, Vaidehi, Yashika, Zaina, Aditi		
Spain	Mateo, Santiago, Lucas, Marcos, Daniel, David, Samuel, Benjamín, Ezequiel, Noé, Salvador, Ismael, Aarón,		
	Elías, Jonás, Jeremías, Iker, Unax, Aitor, Ander, Martín, Rodrigo, Fernando, Alfonso, Enrique, Felipe, Carlos,		
	Javier, Jorge, Luis, Antonio, José, Juan, Manuel, Pedro, Francisco, Ignacio, Rafael, Víctor, Álvaro, Diego,		
	Gabriel, Miguel, Pablo, Ricardo, Sergio, Tomás, César, Gonzalo, Leonardo, Emiliano, Matías, Nicolás, Sebastián,		
	Thiago, Sofía, Camila, Valentina, Martina, Emilia, Emma, Olivia, Luna, Zoe, Mia, Isabella, Victoria, Sara,		
	Lucía, María, Laura, Paula, Andrea, Ana, Elena, Carmen, Alba, Carla, Daniela, Julia, Natalia, Ximena, Aitana,		
	Noa, Mía, Isabel, Beatriz, Blanca, Clara, Inés, Irene, Marta, Patricia, Rocío, Silvia, Teresa, Verónica, Alicia,		
	Amelia, Ángela, Aurora, Bárbara, Carolina, Dolores, Eva, Gloria, Lidia, Lorena, Mónica, Nuria, Olga, Raquel,		
	Sandra, Xiomara, Yamile		

Table 12: Universe of names for country wise name replacement in benchmarking experiments in Sec. B

Model Name	Cosine sim per perturbation pair
all-mpnet-base-v2	0.842 ± 0.0002
all-distilroberta-v1	0.852 ± 0.0002
all-MiniLM-L6-v2	0.784 ± 0.0002
gemini	0.93 ± 0.0
multi-qa-distilbert-cos-v1	0.82 ± 0.0002
paraphrase-MiniLM-L6-v2	0.806 ± 0.0004
distiluse-base-multilingual-cased-v1	0.837 ± 0.0003
distiluse-base-multilingual-cased-v2	0.838 ± 0.0003
paraphrase-multilingual-MiniLM-L12-v2	0.82 ± 0.0003
msmarco-distilbert-cos-v5	0.799 ± 0.0002
multi-qa-mpnet-base-cos-v1	0.815 ± 0.0002
voyage-3-lite	0.847 ± 0.0001

Table 13: Bias Measurement: Names from same country. Perturbation of person names and replacing them with names from *Spain*. We used CMU Book dataset for this experiment and set number of perturbations K=20. In this experiment we use samples without mention of country/city/town/other location names, nationality etc.

Model Name	Cosine sim per perturbation pair
all-mpnet-base-v2	0.84 ± 0.0002
all-distilroberta-v1	0.838 ± 0.0002
all-MiniLM-L6-v2	0.757 ± 0.0003
gemini	0.931 ± 0.0
multi-qa-distilbert-cos-v1	0.806 ± 0.0002
paraphrase-MiniLM-L6-v2	0.79 ± 0.0004
distiluse-base-multilingual-cased-v1	0.83 ± 0.0003
distiluse-base-multilingual-cased-v2	0.833 ± 0.0003
paraphrase-multilingual-MiniLM-L12-v2	0.815 ± 0.0004
msmarco-distilbert-cos-v5	0.786 ± 0.0002
multi-qa-mpnet-base-cos-v1	0.81 ± 0.0002
voyage-3-lite	0.843 ± 0.0001

Table 14: **Bias Measurement: Names from same country.** Perturbation of person names and replacing them with names from *France*. We used CMU Book dataset for this experiment and set number of perturbations K=20. In this experiment we use samples without mention of country/city/town/other location names, nationality etc.

Cosine sim per perturbation pair
0.816 ± 0.0002
0.828 ± 0.0002
0.75 ± 0.0003
0.931 ± 0.0
0.79 ± 0.0002
0.778 ± 0.0004
0.88 ± 0.0002
0.887 ± 0.0002
0.796 ± 0.0004
0.78 ± 0.0002
0.805 ± 0.0002
0.85 ± 0.0001

Table 15: Bias Measurement: Names from same country. Perturbation of person names and replacing them with names from *India*. We used CMU Book dataset for this experiment and set number of perturbations K=20. In this experiment we use samples without mention of country/city/town/other location names, nationality etc.

- 883
- 885 886

890

891

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

922

923

924

925

926

927

928

929

930

D Semantic Similarity Task Dataset

Below we present the STS dataset consisting 10 samples used in Sec. 5.1. Each sample is a triplet of the form:

< Query, Positive sample, Negative sample >.

- 1. **Query:** Nikolai and Deborah met on a rainy Tuesday in New York. The city's hustle and bustle couldn't dim the spark between them. Deborah, with her radiant smile and infectious laughter, had captured Nikolai's heart from the moment he saw her. Nikolai, a charming and witty gentleman, returned her affection with equal fervor.
 - **Positive:** Kashvi and Oluwafemi met on a rainy Tuesday in Northampton. The city's bustling streets couldn't dim the spark between them. Kashvi, with her radiant smile and infectious laughter, had captured Oluwafemi's heart from the moment he saw her. Oluwafemi, a charming and witty gentleman, returned her affection with equal fervor.
 - **Negative:** Nikolai and Deborah staying in New Jersey, once inseparable, were now worlds apart. Deborah, the trusted confidante, had betrayed Nikolai's trust, revealing his secrets to their rivals. The city's hustle and bustle mirrored the chaos within Nikolai's heart, as he grappled with the bitter reality of love turned treachery.
- 2. **Query:** Alejandro quickly ran to the store to buy a cold drink. He was eager to have a glass of cold drink.
 - **Positive:** Quickly, Hiroki dashed to the local market to procure some cold drinks. He was yearning for a chilled glass of cold drink.
 - **Negative:** Alejandro has stopped buying cold drinks from market. He only drinks cold drinks made at home.
- 3. Query: Mayatoshi and Alex had a deep, passionate love for each other. Their bond was unbreakable, a love that transcended all obstacles. They shared dreams, hopes, and aspirations, and their love was the foundation of their happiness.

• **Positive:** Priyanka and Yuan were deeply in love. Their affection for each other was profound and unwavering. They shared a strong connection, a love that was the source of their joy and contentment.

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

- Negative: Despite their intense hatred for each other, Mayatoshi and Alex were bound by a strange, twisted connection. Their animosity fueled a toxic relationship, a constant battle of wills. Their lives were intertwined, a dark, destructive dance of love and hate.
- 4. **Query:** Amazon and Apple are two American corporations. Amazon's main business is online shopping and Apple is a phone maker giant
 - **Positive:** Alibaba and Xiaomi are two Chinese corporations. Alibaba's main business is online shopping and Xiaomi is a producer of phones
 - **Negative:** Amazon is a river in South America. Apples are not grown in the Amazon basin.
- 5. **Query:** Ganga and Yamuna are two mighty rivers. They are lifelines for millions of people in the region.
 - **Positive:** Yangtze is a mighty river. It is a long river and is the lifeline for millions of people in the region.
 - **Negative:** Ganga and Yamuna are two sisters. They had their schooling in the region and schooling provided a lifeline for them.
- 6. **Query:** Alice and Bob often helped each other financially. Recently, Alice lent Bob a significant sum of money. Bob promised to return it soon.
 - **Positive:** Yuri and Haruto frequently helped each other out, including with money. Lately, Yuri had loaned Haruto a substantial amount of money, which Haruto assured her he'd repay promptly.
 - **Negative:** Alice and Bob had a disagreement about money. Alice believed Bob owed her money, but Bob denied it.

- 977 978
- 979 980
- 9
- 98
- 984 985
- 9
- 9
- 991
- 992 993
- 9
- 995 996 997
- 998
- 999 1000
- 1001 1002
- 1003 1004
- 1005 1006
- 1007 1008

1011

- 1012
- 1014
- 10
- 1015

1017 1018

- 1019
- 1020 1021

1022

1023 1024 7. **Query:** John, a renowned lawyer, is defending his client, Mike, who is accused of a serious crime. John is determined to prove Mike's innocence and secure his acquittal.

- **Positive:** Armaan, a man falsely accused of a heinous crime, is relying on his skilled lawyer, Udit, to exonerate him. Udit is committed to presenting a strong defense and clearing Armaan's name.
- Negative: John, a cunning lawyer, is manipulating the legal system to frame Mike for a crime he did not commit. John's goal is to secure a conviction and advance his own career, regardless of the truth.
- 8. **Query:** Dr. Alexander, a seasoned physician, meticulously analyzed patient Sarah's intricate heart condition. He prescribed a tailored regimen of medications and rigorous lifestyle modifications to significantly improve her cardiac health.
 - **Positive:** The esteemed doctor, Dr. Yerusha, conducted a thorough assessment of patient Reyan's complex symptoms of heart. She formulated a precise treatment plan, incorporating medications and day to day lifestyle changes, to alleviate Reyan's debilitating heart condition.
 - **Negative:** Dr. Alexander, a renowned doctor and surgeon, executed a high-risk heart surgical procedure on patient Sarah. After the complex operation Sarah did not recover.
- 9. **Query:** Mr. Smith, a dedicated teacher, guided his students, including the bright young minds of Miller and Pristina, towards academic excellence.
 - **Positive:** Mr. Yang, a committed educator, mentored his students, including the talented Shruti and Ren, to achieve academic success.
 - Negative: Mr. Smith , a rigid and punitive teacher, often unfairly targeted mischievous students like Miller and Pristina.
- 10. **Query:** Martinez gently examined the injured bird. He gave it food.

- **Positive:** Yohan tenderly inspected the wounded bird and gave it a meal to eat. 1025
- Negative: The skilled hunter Martinez tracked the injured bird. He captured it for food.
 1027 1028 1029

1030

1031

1032

1033

1034

1035

1036

1037

1039

1064

1065

E Example of Semantic Similarity post-anonymization

In Table 16, we show impact of anonymization on STS tasks on embeddings crated by Open AI's text-embedding-3-small model. We observe that in all cases performance after anonymization is superior. Specifically, post anonymization, we obtain relatively higher score for positive samples and lower for negative samples.

F Impact of Anonymization Strategy: Removal versus Replacement

This section investigates the effectiveness of re-1041 move of names vs. replacement of names in text 1042 for anonymization. In the replacement strategy, 1043 we replace names with non-identifying placeholder 1044 names instead of removing them from text. Ex-1045 ample: person names with 'CHAR ID', location 1046 names with 'LOC_ID' etc. Here ID can be re-1047 placed with $\{A, B, C \cdots\}$ or $\{1, 2, 3 \cdots\}$ etc. The 1048 detailed prompt is present in Table 17. In Table 18 1049 we demonstrate that removal of names marginally 1050 outperforms replacement in the STS task. In the 1051 context of replacement strategy, one should note 1052 that the quality of embeddings derived is sensi-1053 tive to the specific replacement placeholder terms 1054 used. For instance, substituting character names 1055 with with different placeholders such as "CHAR_A" 1056 / "CHARACTER_B" / "CHARACTER_1" or lo-1057 cation names with "LOC 1" / "LOC B" can im-1058 pact the resulting embeddings differently. In or-1059 der to mitigate this sensitivity and ensure consis-1060 tent results and also based upon our findings we 1061 recommend using the name removal strategy for 1062 anonymization to mitigate name bias. 1063

G Implementation Details

G.1 Model Information and Computational budget

In Table 20, we present the model size of different1067open source models used. For our experiments, we1068consumed approximately 40 GPU hours with one106932 GB GPU.1070

	Query	Pos/Neg	Sim	Label
			score	
Original	Alejandro quickly ran to the store to buy a cold drink. He was	POS: Quickly, Hiroki dashed to the local market to procure some cold drinks. He was yearn-	0.66	1
	eager to have a glass of cold drink.	NEG: Alejandro has stopped buying cold drinks from market. He only drinks cold drinks made at home.	0.72	0
Anonymized	quickly ran to the store to buy a cold drink. He was eager to have a glass of cold drink.	POS: Quickly, dashed to the local market to procure some cold drinks. He was yearning for a chilled glass of cold drink.	0.83	1
		NEG: has stopped buying cold drinks from market. He only drinks cold drinks made at home.	0.57	0
Original	Ganga and Yamuna are two mighty rivers. They are lifelines for millions of people in the	POS: Yangtze is a mighty river. It is a long river and is the lifeline for millions of people in the region.	0.63	1
	region.	NEG: Ganga and Yamuna are two sisters. They had their schooling in the region and schooling provided a lifeline for them.	0.73	0
Anonymized	and are two mighty rivers. They are lifelines for millions of	POS: is a mighty river. It is a long river and is the lifeline for millions of people in the region.	0.76	1
	people in the region.	NEG: and are two sisters. They had their schooling in the region and schooling provided a lifeline for them.	0.46	0

Table 16: Example demonstrating impact of anonymization on semantic similarity using embeddings created by Open AI's *text-embedding-3-small* model. The text in color blue and red refer to the positive and negative paragraphs respectively.

Replace person names, organiza- tions and locations	Given below text, please convert all Person names(which are Proper Nouns) to a UNIQUE ID such as CHAR_A, CHAR_B, CHAR_C etc Keep it unique and for each UNIQUE Person name(which is a Proper Noun) use a UNIQUE ID. DO NOT KEEP THE ORIGINAL Person Names(which are Proper Nouns) in the gen- erated paragraph text. Next, Replace all occurences City/Country/Village/Town/River/Continent etc. names which are PROPER NOUNS to a UNIQUE ID such as
	LOC_A, LOC_B, LOC_C, LOC_D etc Next, Replace all occurences of company/organization names which are PROPER NOUNS to a UNIQUE ID such as ORG_A, ORG_B, ORG_C, ORG_D etc Do not replace monu- ment/landmark names like Eiffel tower etc. Output contains
	the modified text only The text is provided below ::::

Table 17: Prompt for Anonymization using replacement strategy described in Sec. F

1071 G.2 Packages used

1072We used scikit-learn (Pedregosa et al., 2011) pack-1073age for computing metrics such as cosine similarity1074and AUC-ROC.

OpenAI API which are under Apache License, Version 2.0. The Voyage API is licensed under MIT1079license. All the artifacts used in this paper are
available for non-commercial scientific use.1080

1075 G.3 Terms and License

For our implementation, we use sentence transformers library (Reimers, 2019), Gemini API, and

Model	AUC ROC	AUC ROC	AUC ROC
Model	Original	Anonymization(Default)	Anonymization(Replacement)
all-mpnet-base-v2	0.19	0.98 ± 0.0071	$\textbf{1.0} \pm \textbf{0.0}$
all-distilroberta-v1	0.36	$\textbf{0.975} \pm \textbf{0.0106}$	0.945 ± 0.0106
all-MiniLM-L6-v2	0.09	$\textbf{0.99} \pm \textbf{0.0071}$	0.97 ± 0.0071
gemini	0.71	$\textbf{1.0} \pm \textbf{0.0}$	$\textbf{1.0} \pm \textbf{0.0}$
multi-qa-distilbert-cos-v1	0.07	$\textbf{0.97} \pm \textbf{0.0071}$	0.95 ± 0.0
paraphrase-MiniLM-L6-v2	0.14	0.98 ± 0.0	$\textbf{0.99} \pm \textbf{0.0071}$
distiluse-base-multilingual-cased-v1	0.27	0.94 ± 0.0	0.935 ± 0.0106
distiluse-base-multilingual-cased-v2	0.26	0.96 ± 0.0	0.94 ± 0.0212
paraphrase-multilingual-MiniLM-L12-v2	0.21	0.99 ± 0.0	$\textbf{1.0} \pm \textbf{0.0}$
msmarco-distilbert-cos-v5	0.10	$\textbf{0.955} \pm \textbf{0.0035}$	0.875 ± 0.0035
multi-qa-mpnet-base-cos-v1	0.08	$\textbf{1.0} \pm \textbf{0.0}$	0.985 ± 0.0035
text-embedding-3-small	0.12	1.0 ± 0.0	0.97 ± 0.0071
text-embedding-3-large	0.21	$\textbf{1.0} \pm \textbf{0.0}$	$\textbf{1.0} \pm \textbf{0.0}$
voyage-3-lite	0.18	$\textbf{1.0} \pm \textbf{0.0}$	0.98 ± 0.0141

Table 18: Comparison of Removal based vs Replacement based Anonymization on Semantic Similarity task of Sec. 5.1.

Model Name	Euclidean Distance per perturbation pair
all-mpnet-base-v2	0.642 ± 0.0016
all-distilroberta-v1	0.641 ± 0.0015
all-MiniLM-L6-v2	0.766 ± 0.0017
gemini	0.46 ± 0.0007
multi-qa-distilbert-cos-v1	0.694 ± 0.0014
paraphrase-MiniLM-L6-v2	3.398 ± 0.0153
distiluse-base-multilingual-cased-v1	0.638 ± 0.002
distiluse-base-multilingual-cased-v2	0.63 ± 0.0021
paraphrase-multilingual-MiniLM-L12-v2	2.726 ± 0.0108
msmarco-distilbert-cos-v5	0.742 ± 0.0016
multi-qa-mpnet-base-cos-v1	0.679 ± 0.0016
text-embedding-3-small	0.67 ± 0.0013
text-embedding-3-large	0.616 ± 0.0013
voyage-3-lite	0.647 ± 0.001

Table 19: **Bias Measurement on CMU Book dataset with Euclidean distance as distance function between embeddings**. A distance close to 0 is better.

Size
420 MB
290 MB
80 MB
250 MB
90.9 MB
480 MB
480 MB
420 MB
265 MB
438 MB

Table 20: Model information for open source models.