# A SINGLE SWALLOW DOES NOT MAKE A SUMMER: UNDERSTANDING SEMANTIC STRUCTURES IN EMBED DING SPACES

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027 028 029

030

037

039

040

041

042

043

044

Paper under double-blind review

#### Abstract

Embedding spaces encapsulate rich information from deep learning models, with vector distances reflecting the semantic similarity between textual elements. However, their abstract nature and the computational complexity of analyzing them remain significant challenges. To address these, we introduce the concept of Semantic Field Subspace, a novel mapping that links embedding spaces with the underlying semantics. We propose SAFARI, a novel algorithm for SemAntic Field subspAce deteRmInation, which leverages hierarchical clustering to discover hierarchical semantic structures, using Semantic Shifts to capture semantic changes as clusters merge, allowing for the identification of meaningful subspaces. To improve scalability, we extend Weyl's Theorem, enabling an efficient approximation of Semantic Shifts that significantly reduces computational costs. Extensive evaluations on five real-world datasets demonstrate the effectiveness of SAFARI in uncovering interpretable and hierarchical semantic structures. Additionally, our approximation method achieves a  $15 \sim 30 \times$  speedup while maintaining minimal errors (less than 0.01), making it practical for large-scale applications. The source code is available at https://anonymous.4open.science/r/Safari-C803/.

### 1 INTRODUCTION

Embedding spaces are widely recognized for encapsulating rich information learned by deep learning models. In language models, for instance, the distance between embedding vectors often reflects the semantic similarities of corresponding textual elements (Devlin et al., 2019; Wu et al., 2020). However, despite their widespread use, the understanding of embedding spaces remains limited (Ethayarajh, 2019; Clark et al., 2019; Simhi & Markovitch, 2023). Two primary challenges make understanding embedding spaces difficult:

- (1) **Abstract Nature:** Embedding spaces reflect complex, high-dimensional relationships between data points, making them inherently abstract. To interpret these spaces, we need a clear connection between the embedded data and their underlying semantics. However, a universally accepted definition of semantics within these spaces remains elusive.
- (2) Computational Complexity: Understanding embedding spaces requires substantial data, and as these spaces grow richer and more complex, the need for samples increases, straining computational efficiency. Managing this expanding data demands advanced computational resources and optimized algorithms to ensure practical and timely analysis.

045 Extensive research on embedding spaces spans various perspectives, with two lines of research most 046 relevant to our work: geometry-based and interpretability-based methods. The first line of research 047 focuses on the geometric properties of the embedding space, describing vector distributions and their 048 desired characteristics (Mu & Viswanath, 2018; Liu et al., 2019; Demeter et al., 2020; Ethayarajh, 2019). However, these approaches primarily enhance representation quality by manipulating the geometry of the space, often overlooking the interpretability and semantic coherence of the embeddings. 051 The second line of research concentrates on making the embedding space more interpretable using techniques like rotation, probing (Park et al., 2017; Dufter & Schütze, 2019; Clark et al., 2019; Dalvi 052 et al., 2019), or transforming data into more interpretable dimensions (Simhi & Markovitch, 2023). Despite providing insights into individual dimensions, they often rely on significant assumptions

about the interpretability of original dimensions or require constructing new spaces, which can be computationally intensive and may not preserve the original semantic relationships.

In this paper, we investigate the semantic structure of embedding spaces through textual elements. To address their abstract nature and the challenges of interpretation, we introduce key concepts that link embedding spaces to the underlying semantics. We define each direction as a unique semantic set, serving as a foundation for understanding the structure. Recognizing the context-dependent nature of semantics, we introduce the concept of a Semantic Field for more nuanced interpretation. Identifying Semantic Fields within embedding spaces is framed as an optimization problem. We approximate each Semantic Field as a subspace, referred to as the Semantic Field Subspace, and solve the optimization using Singular Value Decomposition (SVD) (Halko et al., 2009; Trefethen & Bau, 2022).

Building on this concept, we propose SAFARI, a novel algorithm for determining Semantic Field
Subspaces through hierarchical clustering. By introducing the concept of Semantic Shift, SAFARI
accurately identifies the start and end points of Semantic Field Subspaces within a clustering dendrogram, unveiling the hierarchical structure of Semantic Fields. To overcome the computational
challenges of analyzing large datasets, we extend Weyl's Theorem (Weyl, 1912) to approximate
Semantic Shift without relying on full SVD, significantly improving computational efficiency.

071<br/>072We evaluate SAFARI on five real-world datasets to validate its ability to uncover hierarchical<br/>semantic structures. The results confirm that SAFARI efficiently identifies Semantic Field Subspaces,<br/>revealing natural and interpretable hierarchies. Moreover, our approximate method for Semantic<br/>Shift computation delivers a  $15 \sim 30 \times$  speedup with errors less than 0.01, making it highly practical<br/>for large-scale applications. Our contributions are summarized as follows:

- (1) We introduce the concept of Semantic Field Subspaces, a novel mapping that bridges embedding spaces with their underlying semantics. It enhances the interpretability of high-dimensional embeddings, facilitating deeper insights into the semantics encoded within vector spaces.
- (2) We present SAFARI, a creative algorithm that leverages Semantic Shift to determine Semantic Field Subspaces. By employing hierarchical clustering, SAFARI effectively uncovers the hierarchical semantic structures. We also develop an efficient approximate method for Semantic Shift computation, significantly improving computational efficiency.
  - (3) We systematically evaluate the efficacy of SAFARI through extensive experiments, demonstrating that our algorithm successfully and efficiently identifies Semantic Filed Subspaces, revealing their hierarchical structures.
- 085 086 087

084

076

077

078

# 2 RELATED WORK

Geometry-based Approaches. This research focuses on the geometric properties of embedding 089 spaces, aiming to describe the vector distributions and their desired characteristics (Mu & Viswanath, 2018; Liu et al., 2019; Demeter et al., 2020; Ethayarajh, 2019). For instance, Mu & Viswanath (2018) 091 improved word representations by removing the top principal components, while Liu et al. (2019) 092 suppressed transformed dimensions with large variances. Demeter et al. (2020) highlighted how the softmax function weakens geometric structures, introducing bias. A key finding by Ethayarajh (2019) 094 showed that most vectors reside within a narrow cone in the embedding space. Unlike geometry-based 095 research, we do not aim to prove or find the ideal distribution or other geometric properties in the 096 embedding space to improve the model. Instead, SAFARI focuses on revealing and understanding 097 structures within a given embedding space, regardless of its geometric properties. This allows us to 098 maintain the original semantic relationships while uncovering meaningful patterns.

099 Interpretability-based Approaches. This research targets making embedding spaces more inter-100 pretable, often through rotation and probing methods (Park et al., 2017; Dufter & Schütze, 2019; 101 Clark et al., 2019; Dalvi et al., 2019). Park et al. (2017) employed rotation algorithms to improve 102 word vector interpretability, while Dufter & Schütze (2019) applied rotation to enhance word space 103 comprehension. Clark et al. (2019) analyzed attention mechanisms in pre-trained models, particularly 104 BERT (Devlin et al., 2019), to gain insights into how the model processes information. Dalvi et al. 105 (2019) examined individual vector dimensions in NLP models to uncover the roles of these dimensions. Recently, Simhi & Markovitch (2023) transformed latent spaces into a new one with more 106 conceptualized and interpretable dimensions. Although SAFARI also seeks to interpret embedding 107 spaces, it differs by not assuming that the original dimensions are inherently interpretable or by

transforming them into new spaces. Instead, SAFARI identify comprehensible structures by linking
 embedding vectors to their underlying semantic space, preserving the original embeddings and
 providing a robust framework for semantic interpretation.

## **3** PROBLEM FORMULATION

Before introducing SAFARI, we first define the key concepts of Semantic Distance, Semantic Field, and Semantic Field Subspace. Frequently used notations are summarized in Table A1.

**Semantic Distance.** Let  $\mathcal{X}$  be a set of textual elements and  $\mathbb{R}^d$  a continuous, *d*-dimensional embedding space learned by a deep learning model  $h : \mathcal{X} \to \mathcal{E}$ , where  $\mathcal{E} \subset \mathbb{R}^d$  denotes embedding vectors for the elements in  $\mathcal{X}$ . Semantic Distance measures how different two textual elements are, based on the distance between their embedding vectors:

**Definition 1** (Semantic Distance): The Semantic Distance  $d_{sem}(\cdot, \cdot)$  between any two textual elements  $x, x' \in \mathcal{X}$  is defined as the cosine distance between their embedding vectors  $v, v' \in \mathcal{E}$ :

129

130

131

132

133

134

144 145

146

147

148

156

121

112

113

 $d_{sem}(\boldsymbol{v}, \boldsymbol{v}') = 1 - \langle \boldsymbol{v}, \boldsymbol{v}' \rangle / (\|\boldsymbol{v}\| \cdot \|\boldsymbol{v}'\|).$ (1)

We use x and embedding vector v interchangeably when there is no ambiguity. Each textual element often carries multiple layers of meaning, which we refer to as its semantics. Let  $\mathcal{M}$  be the set of all possible semantics. The semantics of an embedding vector v can be expressed as a set  $f_{sem}(v)$ , where  $f_{sem}(v) : \mathcal{E} \to 2^{|\mathcal{M}|} \setminus \emptyset$ , representing the various semantic facets of v.

We argue that a single textual element x cannot be fully interpreted in isolation, as its semantics,  $f_{sem}(v)$ , require context for interpretation. As the adage states, "You shall know a word by the company it keeps," meaning x becomes interpretable only when considered with related elements in its context. For example, as shown in Fig. 1, the word 'Apple' is ambiguous on its own but gains specific meaning when used in different contexts–referring to a technology company with words like 'Mac,' 'IBM,' and 'Windows,' or to a fruit with words like 'Apple Tree,' 'Juice,' and 'Banana.' The meaning becomes clearer as more contextual words are added.



Figure 1: A word 'Apple': from ambiguous to specific.



Interpreting an Embedding Vector using Close Neighborhood. As discussed earlier, the semantics of a textual element x can be interpreted through the context in which it appears. In models like Word2Vec (Mikolov et al., 2013) and contextual models like BERT (Devlin et al., 2019), the semantics of an embedding vector v is shaped by the surrounding vectors in its context.

To capture this context, we consider a subset of  $\mathcal{E}$  that contains embedding vectors sharing common semantics with v. According to Definition 1, these vectors are located near v in terms of Semantic Distance. However, not all nearby vectors contribute meaningful information for interpreting v. For instance, as depicted in Fig. 2, the vectors of 'Coca-Cola' and 'Coke' are nearly identical but redundant, as they represent the same concept. Such synonyms are excluded from the interpretation. We approximate synonyms as the k-Nearest Neighbors (k-NNs) of v, typically with k set to 3. The close neighborhood of v, denoted as  $\mathcal{N}(v)$ , is defined as:

$$\mathcal{N}(\boldsymbol{v}) = \{ \boldsymbol{v}' \mid f_{sem}(\boldsymbol{v}') \cap f_{sem}(\boldsymbol{v}) \neq \emptyset, \boldsymbol{v}' \in \mathcal{E} \} \setminus k\text{-NNs}(\boldsymbol{v}).$$

In Fig. 2, the words 'Sprite' and 'Pepsi' are the close neighborhood of 'Coca-Cola.' Since enumerating the entire set  $\mathcal{N}(v)$  is impractical, we focus on a subset of  $\mathcal{N}(v)$  in its context for interpretation.

**Definition 2** (Interpretable Semantics of an Embedding Vector): Given a subset  $\mathcal{N}_{sub}(v) \subseteq \mathcal{N}(v)$ , the interpretable semantics of v is defined as the intersection of the semantics of all  $v' \in \mathcal{N}_{sub}(v)$ :

$$f_{int}(\boldsymbol{v}) = \bigcap_{\boldsymbol{v}' \in \mathcal{N}_{sub}(\boldsymbol{v})} f_{sem}(\boldsymbol{v}').$$
(2)

**Semantic Field.** After interpreting the semantics of a single vector v, we extend this to a set of embedding vectors  $C \subseteq \mathcal{E}$ , referred to as a Semantic Field. Similar to Definition 2, the Semantic Field captures the shared semantics across multiple vectors by intersecting the meanings of vectors in C. Since C may contain vectors with varying semantics, we refine it to include only those that share the most common semantics. Formally, we define the Semantic Field as:

**Definition 3** (Semantic Field): *Given a set C of embedding vectors, the Semantic Field is defined as:* 

$$F_{int}(\mathcal{C}) = \bigcap_{\boldsymbol{v} \in \mathcal{C}^*} f_{sem}(\boldsymbol{v}),\tag{3}$$

where  $C^*$  is a subset of C that maximizes shared semantics by minimizing the symmetric difference:

$$\mathcal{C}^* = \arg\min_{\mathcal{C}_{sub} \subseteq \mathcal{C}} \left| \bigcup_{\boldsymbol{v}' \in \mathcal{C}_{sub}} f_{sem}(\boldsymbol{v}') - \bigcap_{\boldsymbol{v}' \in \mathcal{C}_{sub}} f_{sem}(\boldsymbol{v}') \right|.$$
(4)

173 Suppose  $C_{sub} = \{v_1, v_2\}$ . As shown in Fig. 3, 174 the symmetric difference is visualized as the 175 shadow area, and minimizing it helps iden-176 tify the optimal subset  $C^* \subseteq C$  that shares 177 the most common semantics (the overlapping 178 area). Nevertheless, the concept of a Seman-179 tic Field relies on the latent semantic function  $f_{sem}(\cdot)$ , making them abstract and hard to 180

169

171

172

187

188 189

191

192

201 202

203

$$f_{sem}(oldsymbol{v}_2) igcup_{v'\in\{oldsymbol{v}_1,oldsymbol{v}_2\}} oldsymbol{U}_{sem}(oldsymbol{v}') - igcup_{v'\in\{oldsymbol{v}_1,oldsymbol{v}_2\}} oldsymbol{f}_{sem}(oldsymbol{v}') - igcup_{v'\in\{oldsymbol{v}_1,oldsymbol{v}_2\}} oldsymbol{f}_{sem}(oldsymbol{v}')$$

Figure 3: Visualization of symmetric difference.

compute directly. To remedy this issue, we introduce the concept of a Semantic Field Subspace.

Semantic Field Subspace. According to Definition 1, embedding vectors pointing in the same directions represent the same semantics. This insight leads to the idea that *every vector in the embedding space corresponds to a unique set of semantics*, and conversely, a specific vector can represent each semantic set. Hence, a Semantic Field can be approximated by treating the semantic sets as embedding vectors and solving the optimization in Eq. (4).

**Definition 4** (Semantic Field Subspace): Let  $\mathbb{S}$  be a subspace of  $\mathbb{R}^d$ . The semantics of  $\mathbb{S}$  is defined as:

$$F_{sem}(\mathbb{S}) = \bigcap_{\boldsymbol{v} \in \mathcal{C}^*} f_{sem}(\boldsymbol{v}),\tag{5}$$

190 where  $C^*$  is the set of embedding vectors in  $\mathbb{S}$  that minimizes the following symmetric difference:

$$\mathcal{C}^* = \arg\min_{\mathcal{C}\subset\mathbb{S}} \left| \bigcup_{\boldsymbol{v}'\in\mathcal{C}} f_{sem}(\boldsymbol{v}') - \bigcap_{\boldsymbol{v}'\in\mathcal{C}} f_{sem}(\boldsymbol{v}') \right|.$$
(6)

Although enumerating all vectors in a subspace is still impractical, we can approximate a subspace using a finite set of representative vectors  $C \subseteq \mathcal{E}$ . Representing C as a matrix A allows us to apply SVD to extract the key components, capturing the essential semantics of the subspace. We define this relaxed subspace approximation as follows:

**197 Definition 5** (Relaxed Version of Semantic Field Subspace): Given a subset  $C \subseteq \mathcal{E}$  represented by **198** a matrix A, the subspace  $\mathbb{S}$  can be approximated by A, i.e.,  $\mathbb{S} \approx A$ . Applying SVD to A gives **199**  $A = U\Sigma V^{\top}$ , and the Semantic Field Subspace  $F_{sem}(\mathbb{S})$  is approximated by the singular values in **200**  $\Sigma$  and the singular vectors in  $V^{\top}$ , i.e.,  $F_{sem}(\mathbb{S}) \approx \Sigma$ ,  $V^{\top}$ .

## 4 Methodology

# 4.1 THE SAFARI ALGORITHM

Motivation. SAFARI is designed to identify Semantic Field Subspaces from a set  $\mathcal{E}$  of embedding vectors. While Definition 5 provides a theoretical framework for constructing these subspaces, the primary challenge lies in selecting appropriate subsets  $\mathcal{C} \subseteq \mathcal{E}$ , especially in a large, diverse text corpus. SAFARI addresses this by leveraging the natural clustering property of embedding vectors, where clusters represent specific topics or semantic themes, enabling efficient identification of meaningful subspaces in complex, high-dimensional embedding spaces.

Algorithm Description. SAFARI utilizes hierarchical clustering to iteratively identify potential
 Semantic Field Subspaces, as it offers flexibility by not requiring a predefined number of clusters,
 allowing for dynamic exploration of the data. Additionally, hierarchical clustering produces a
 dendrogram, which helps interpret relationships between clusters and their semantic structures. The
 pseudo-code for SAFARI is shown in Algorithm 1.

#### 216 Algorithm 1: SAFARI 217 **Input:** A set $\mathcal{E}$ of embedding vectors in $\mathbb{R}^d$ , window size w; 218 **Output:** A set $\Psi$ of clusters with specific Semantic Field Subspaces; 219 1 $\Phi \leftarrow$ Initialize each $v \in \mathcal{E}$ as its own cluster; 220 2 *iter* = 0; $\mu$ = 0; $\tau$ = 0; 221 $3 \Psi \leftarrow \emptyset;$ ▷ Store clusters with specific Semantic Field Subspaces 222 4 while $|\Phi| > 1$ do ▷ Step 1: Cluster Merging 224 $\{\mathcal{C}_x, \mathcal{C}_y\} \leftarrow \arg\min_{\mathcal{C}_i, \mathcal{C}_i \in \Phi} d_{sem}(\mathcal{C}_i, \mathcal{C}_j);$ 5 225 $\mathcal{C}_{new} \leftarrow \mathcal{C}_x \cup \mathcal{C}_y;$ 6 226 $\Phi \leftarrow \Phi \cup \mathcal{C}_{new} \setminus \{\mathcal{C}_x, \mathcal{C}_y\};$ 7 227 ▷ Step 2: Semantic Field Subspace Determination 228 $\mathcal{C}_x \leftarrow |\mathcal{C}_x| > |\mathcal{C}_y| ? \mathcal{C}_x : \mathcal{C}_y;$ 8 229 Compute the Semantic Shift $\Delta F_{sem}(\mathcal{C}_x, \mathcal{C}_{new})$ using Algorithm 2; 9 230 if $\Delta F_{sem}(\mathcal{C}_x, \mathcal{C}_{new}) > \mu + 2\tau$ then $\Psi \leftarrow \Psi \cup \mathcal{C}_{new}$ ; 10 231 iter = iter + 1;11 232 Update $\mu$ and $\tau$ by considering $\Delta F_{sem}(\mathcal{C}_x, \mathcal{C}_{new})$ and the previous (w-1) values; 12 233 13 return $\Psi$ ;

234 235

237

238

239 240

241

242

243

244

245

246

247

248

249

261

Initially, each embedding vector in  $\mathcal{E}$  forms its own cluster, resulting in *n* clusters (Line 1). These clusters do not yet provide meaningful interpretation or form Semantic Field Subspaces. An empty set  $\Psi$  is initialized to store clusters with specific semantic meanings (Line 3). The algorithm then iterates through two key steps until all clusters are merged into a single one (Lines 4–12).

- Step 1: Cluster Merging. First, the two closest clusters,  $C_x$  and  $C_y$ , are identified based on the Semantic Distance  $d_{sem}(C_x, C_y)$ , with the centroid representing each cluster (Line 5). These two clusters are then merged into a new cluster  $C_{new}$  (Line 6), after which the original clusters  $C_x$  and  $C_y$  are removed, and the new cluster  $C_{new}$  is added to the set  $\Phi$  (Line 7).
- Step 2: Semantic Field Subspace Determination. The larger cluster,  $C_x$ , is selected, and the Semantic Shift  $\Delta F_{sem}(C_x, C_{new})$  is computed to measure the semantic gap between  $C_x$ and  $C_{new}$  (Lines 8–9), where its definition and computation will be presented later. A sliding window of size w tracks the last w Semantic Shift values, calculating their mean ( $\mu$ ) and standard deviation ( $\tau$ ). If  $\Delta F_{sem}(C_x, C_{new})$  exceeds the *dynamic* threshold ( $\mu + 2\tau$ ), indicating a large semantic gap,  $C_{new}$  is added to  $\Psi$  as a Semantic Field Subspace (Line 10). At last, the algorithm updates the iteration counter and recalculates  $\mu$  and  $\tau$  for the next iteration (Lines 11–12).

**Exact Semantic Shift Computation.** We now define and describe the process for computing the Semantic Shift (pseudo-code provided in Appendix B). Given two clusters,  $C_x$  and  $C_{new}$ , we first construct matrices  $A_x$  and  $A_{new}$ , representing their respective subspaces,  $\mathbb{S}_x$  and  $\mathbb{S}_{new}$ . Following Definition 5, SVD is performed on these two matrices  $A_x$  and  $A_{new}$  to approximate the semantics of these subspaces:  $F_{sem}(\mathbb{S}_x) \approx \Sigma_x, V_x^{\top}$  and  $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top}$ .

We then compare the singular vectors  $v_i \in V_x^{\top}$  with their nearest neighbors  $\tilde{v}_i^* \in V_{new}^{\top}$ , based on Semantic Distance  $d_{sem}(v_i, \tilde{v}_i^*)$ , which captures shifts in semantic direction. For each singular value  $\sigma_i \in \Sigma_x$  and  $\tilde{\sigma}_i \in \Sigma_{new}$  sorted in descending order, we calculate the difference  $\Delta \sigma_i = |\sigma_i - \tilde{\sigma}_i|$ , reflecting shifts in the importance of each dimension. Thus, the total Semantic Shift between clusters  $C_x$  and  $C_{new}$  (or subspaces  $S_x$  and  $S_{new}$ ) is defined as:

$$\Delta F_{sem}(\mathcal{C}_x, \mathcal{C}_{new}) = \Delta F_{sem}(\mathbb{S}_x, \mathbb{S}_{new}) = \sum_i \Delta \sigma_i \cdot d_{sem}(\boldsymbol{v}_i, \tilde{\boldsymbol{v}}_i^*).$$
(7)

Eq. (7) captures both the *importance difference* (through  $\Delta \sigma_i$ ) and *directional difference* (through  $d_{sem}(v_i, \tilde{v}_i^*)$ ), providing a comprehensive measure of the semantic gap between subspaces.

Example 1: Consider a toy example with 11 words. The dendrogram in Fig. 4 shows how Algorithm
l identifies Semantic Field Subspaces through hierarchical clustering. In the first three iterations,
semantically similar word pairs, such as 'Macbook Air' and 'Macbook Pro,' 'PowerPoint' and 'Excel,'
and 'Michael Jordan' and 'Chicago Bulls,' are merged together. These merges result in only minor
Semantic Shifts, so they are not recognized as Semantic Field Subspaces. However, in the 4th
iteration, the word 'Apple' is merged with the 'Macbook Air' and 'Macbook Pro' cluster, which
produces a significant Semantic Shift, indicating the creation of a new Semantic Field Subspace.

277

278

279

281

283

284 285

286 287

288

289

290

291

292

293

294

295

296 297

298

299

300

301

302

303

304

306

307

311

318

321

270 This example also highlights the hierarchical nature of Semantic Field Subspaces. For instance, 271 the 'IT Companies' subspace encompasses both 'Apple (as an IT Company)' and 'Microsoft (as an 272 IT Company)' as individual subspaces within it. Additionally, as the clustering process continues, 273 SAFARI dynamically tracks Semantic Shifts using a sliding window, allowing it to adjust the 274 threshold in real-time. This adaptive mechanism ensures that irrelevant clusters (like 'IT Companies' and 'NBA') do not form new Semantic Field Subspaces, preserving semantic integrity. Δ 275



Figure 4: A toy example of Algorithm 1.

#### 4.2 APPROXIMATE SEMANTIC SHIFT COMPUTATION AND THEORETICAL ANALYSIS

Approximate Semantic Shift Computation. Exact Semantic Shift computation is computationally intensive due to the need for a full SVD on two dense matrices,  $A_x$  and  $A_{new}$ . For a matrix of size  $n \times d$  ( $d \le n$ ), the SVD has a time complexity of  $O(nd^2)$  (Halko et al., 2009; Trefethen & Bau, 2022), which becomes prohibitive when repeated at each iteration. To address this, we develop an efficient approximation algorithm (pseudo-code in Appendix C). Given the larger cluster  $C_x$  and the smaller cluster  $C_y$ , we first construct matrices  $A_x$  and  $A_y$ . We then compute the spectral norm of  $A_y$ and the maximum singular value  $\sigma_{max}$  from  $A_x$ , using them to approximate the Semantic Shift of  $\mathcal{C}_x$  and  $\mathcal{C}_{new}$ :

$$\Delta F_{sem}(\mathcal{C}_x, \mathcal{C}_{new}) = \left\| \boldsymbol{A}_y \right\|_2 \sigma_{max}(\boldsymbol{A}_x).$$
(8)

**Theoretical Analysis.** Next, we provide a theoretical analysis to show that the Semantic Shift (Eq. (7)) can be accurately approximated by Eq. (8). Given two matrices  $A_x$  and  $A_y$ , let  $A_{new} = [A_x|A_y]$ be the matrix of the newly merged cluster  $C_{new}$ . We begin with Theorem 1, which shows that the magnitude difference  $(\Delta \sigma_i)$  between  $A_x$  and  $A_{new}$  is bounded.

**Theorem 1:** Given any two matrices  $A_x$  and  $A_y$  with equal number of columns, where  $A_x$  has more rows than  $A_y$ , the following inequality holds:

$$\Delta \sigma_i = |\sigma_i(\boldsymbol{A}_x) - \sigma_i(\boldsymbol{A}_{new})| \le \|\boldsymbol{A}_y\|_2.$$
(9)

Proof. The proof relies on Weyl's Theorem (Weyl, 1912), which we first present. In each iteration where two clusters are merged, let  $A \in \mathbb{R}^{m \times d}$  denote the matrix of the larger cluster. The merging process introduces a perturbation to A, denoted as E, and the perturbed matrix is given by A = A + E. 308 Weyl's Theorem provides a bound on the change in singular values caused by this perturbation: 309

**Theorem 2** (Weyl's Theorem (Weyl, 1912)): 310

$$\sigma_i(\boldsymbol{A}) - \sigma_i(\boldsymbol{A})| = |\sigma_i(\boldsymbol{A}) - \sigma_i(\boldsymbol{A} + \boldsymbol{E})| \le \|\boldsymbol{E}\|_2.$$
(10)

312 Weyl's Theorem states that the singular values of a matrix cannot change by more than the spectral 313 norm of the perturbation matrix E. Even though E is often assumed to be small in matrix perturbation 314 theory, this result holds for any perturbation, regardless of the size of  $||E||_2$  (Stewart, 1998).

315 We now apply Weyl's Theorem to prove Theorem 1. Let O be a zero matrix. Thus,  $A_{new} =$ 316  $|A_x|A_y| = |A_x|O| + |O|A_y|$ . According to Theorem 2, we have: 317

$$\left|\sigma_{i}([oldsymbol{A}_{x}|oldsymbol{O}])-\sigma_{i}(oldsymbol{A}_{new})
ight|\leq\left\|[oldsymbol{O}|oldsymbol{A}_{y}]
ight\|_{2}=\left\|oldsymbol{A}_{y}
ight\|_{2}$$

To complete the proof, we need to show that  $\sigma_i([\mathbf{A}_x|\mathbf{O}]) = \sigma_i(\mathbf{A}_x)$ . Since this equality always holds, 319 320 Theorem 1 is proved.

Based on Theorem 1, Eq. (7) can be rewritten as:  $\Delta F_{sem}(\mathcal{C}_x, \mathcal{C}_{new}) = \sum_i \Delta \sigma_i \cdot d_{sem}(v_i, \tilde{v}_i^*) \leq 1$ 322  $\sum_{i} \|\boldsymbol{A}_{y}\|_{2} \cdot d_{sem}(\boldsymbol{v}_{i}, \tilde{\boldsymbol{v}}_{i}^{*}) = \|\boldsymbol{A}_{y}\|_{2} \cdot \sum_{i} d_{sem}(\boldsymbol{v}_{i}, \tilde{\boldsymbol{v}}_{i}^{*}).$  To further approximate  $\sum_{i} d_{sem}(\boldsymbol{v}_{i}, \tilde{\boldsymbol{v}}_{i}^{*}),$ 323 we present the following theorem:

Theorem 3: Given two matrices  $A_x$  and  $A_y$  used in the exact Semantic Shift computation as presented in Section 4.1, the directional difference between them is proportional to the largest singular vector  $\sigma_{max}$  of  $A_x$ :

 $\sum_{i} d_{sem}(\boldsymbol{v}_{i}, \tilde{\boldsymbol{v}}_{i}^{*}) = \mathcal{O}(\sigma_{max}(\boldsymbol{A}_{x})).$ (11)

29 *Proof.* To prove Theorem 3, we first introduce Lemma 1.

327

328

334 335

336 337 338

339

340 341

348

349

350

351 352 353

354

355

**Lemma 1:** For a given matrix M with a series of singular values  $\{\sigma_i\}$ , the condition number is  $\kappa(M) = \frac{\sigma_{max}}{\sigma_{min}}$ . The condition number quantifies the matrix's sensitivity to small perturbations, where higher values indicate greater susceptibility to changes (Belsley et al., 2005; Meyer, 2023). Consequently, we establish:

$$\sum_{i} d_{sem}(\boldsymbol{v}_{i}, \tilde{\boldsymbol{v}}_{i}^{*}) = \mathcal{O}(\kappa(\boldsymbol{M})).$$
(12)

Given that real-world matrices contain noise and are often rank-deficient, we assume that:

$$\forall \boldsymbol{M}, \lim_{i \to r} \sigma_i = 0,$$
  
$$\forall \boldsymbol{M}_1, \boldsymbol{M}_2, \lim_{i \to r_1} \sigma_i(\boldsymbol{M}_1) = \mathcal{O}(\lim_{j \to r_2} \sigma_j(\boldsymbol{M}_2)).$$
(13)

Using this assumption, we compare the condition numbers of two matrices  $M_1$  and  $M_2$ :

$$\frac{\sigma(M_1)}{\sigma(M_2)} = \frac{\sigma_{max}(M_1) \times \sigma_{min}(M_2)}{\sigma_{max}(M_2) \times \sigma_{min}(M_1)}.$$
(14)

According to Eq. (13), we have  $\frac{\kappa(M_1)}{\kappa(M_2)} = \frac{\sigma_{max}(M_1)}{\sigma_{max}(M_2)}$ . Thus, comparing condition numbers of matrices under this assumption (Eq. (13)) is equivalent to comparing their largest singular values. Therefore,  $\kappa(M) = O(\sigma_{max}(M))$ . (15)

With Lemma 1, we establish that the condition number, and thus the total directional difference, can be approximated by the largest singular value. This completes the proof of Theorem 3.  $\Box$ 

By utilizing this approximation for Semantic Shifts, we can bypass the need for full SVD in each iteration. Instead, we only need to compute  $||A_y||_2$  and  $\sigma_{max}(A_x)$ , significantly reducing the time complexity in Algorithm 1 (Meyer, 2023; Horn & Johnson, 2012).

### 5 EXPERIMENTS

#### 5.1 SEMANTIC FIELD SUBSPACE ISOLATION

356 Experiment Setup. We first evaluate whether the Semantic Field Subspaces determined by SAFARI 357 effectively preserve their semantic meanings and remain isolated from one another. We employ four 358 datasets: AG-News (Zhang et al., 2015), AAPD (Yang et al., 2018), IMDB (Maas et al., 2011), and 359  $\mathsf{Yelp}^1$  (see Appendix D for details). Using BLINK (Wu et al., 2020) for entity linking, we extract 360 entities from each dataset and rank them by their TF-IDF scores (Schütze et al., 2008; Leskovec et al., 361 2020), selecting the top 10%. These entities are then split into 80% for training (used to identify 362 subspaces) and 20% for testing. For each Semantic Field Subspace, we retain the top 100 singular 363 vectors. We then compute the average distance between test entities and the identified subspaces to evaluate the isolation and preservation of semantic meaning. The results are presented in Fig. 5. 364

365 Result Analysis. The results in Fig. 5 yield three 366 key observations: (1) Test entities are closest to the 367 subspace corresponding to their respective dataset 368 and show reasonable distances to others, confirming 369 that SAFARI effectively preserves semantic meanings within isolated subspaces. (2) Entities from AAPD 370 exhibit the greatest distance from other subspaces, re-371 flecting the significant semantic gap between academic 372 papers and other types of content. (3) Entities within 373 AAPD are further from their own subspace compared 374 to other datasets, likely due to the smaller number of 375 entities extracted from academic papers, resulting in 376 a less rich semantic representation. 377



<sup>&</sup>lt;sup>1</sup>https://www.yelp.com/dataset

Figure 5: Semantic Field Subspace Isolation.

# 378 5.2 SEMANTIC FIELD SUBSPACE CLASSIFICATION 379

Experiment Setup. We conduct a classification experiment to further assess whether the Semantic Field Subspaces identified by SAFARI retain semantic meaning. AG-News is divided into four categories: Business, Sci/Tech, Sports, and World, while AAPD, IMDB, and Yelp are grouped based on their content: academic papers, movie reviews, and business entities, respectively, resulting in 7 distinct categories (classes). Following the setup from Section 5.1, we assign class labels to entities and select the top-n entities from each class, using 80% for training and 20% for testing. Semantic Field Subspaces are constructed using training data, with each subspace assigned a class label. For the test entities, we calculate the distance (weighted by singular values) to all subspaces, predicting the label of the nearest one. We compare the performance of SAFARI against several baselines, including SVM (Platt, 1999; Chang & Lin, 2011), KNN (Cover & Hart, 1967; Fix, 1985), Random Forest (Breiman, 2001), and BERT (Dalvi et al., 2019), with Random Guess as a trivial baseline. Classification accuracy and training time are presented in Figs. 6 and 7.



**Result Analysis.** The results in Fig. 6 reveals three crucial findings: (1) SAFARI and SVM achieve the highest accuracy, with SAFARI slightly outperforming in some cases. (2) KNN and Random Forest perform moderately well, both surpassing BERT and Random Guess. (3) BERT lags in accuracy due to its reliance on embedding vectors rather than raw text. Regarding efficiency (see Fig. 7), BERT incurs the highest time cost, increasing sharply as the number of entities grows. SVM also shows a steep rise in time cost but remains much faster than BERT. In contrast, SAFARI and Random Forest maintain consistently low time costs. Overall, SAFARI strikes the best balance between accuracy and efficiency, confirming its efficacy in retaining semantic meaning for classification.





Experiment Setup. We then assess the speed of Semantic Shift computation by comparing full
SVD with our approximation method. The input data includes the top 2,000 entities from each class,
processed via hierarchical clustering, with the resulting matrices used to measure computational
efficiency. The results are averaged over 10 independent runs.

**Result Analysis.** As shown in Fig. 8, our approximation method is significantly faster than full SVD, achieving speedups of  $15 \sim 30 \times \text{across 7}$  distinct classes. The small variations, indicated by error bars, reflect stability across runs. While the speedup increases with larger errors, we maintained at

53 FI

432  $10^{-3}$ , suggesting the potential for even greater speedups with larger errors. These results underscore 433 the efficiency and reliability of our approximation method for computing Semantic Shifts. 434

## 5.4 HIERARCHICAL SEMANTIC STRUCTURE

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

**Experiment Setup.** We validate SAFARI's ability to uncover hierarchical semantic structures by applying it to the top 1,000 entities from each category, with a focus on the Sports category from the AG-News dataset. The Sports category was chosen due to its well-structured, event-driven content, providing an ideal setting for evaluation. At each iteration, Semantic Shifts are computed to identify Semantic Field Subspaces, using both exact and approximation methods to further confirm SAFARI's ability to deliver accurate approximations in real-world data. The results from iteration 11,000 to 16,000 are depicted in Fig. 9, with specific shifts at iterations 11,352 and 15,856 highlighted to showcase the hierarchical structure, further analyzed in Figs. 10 and 11, respectively.



Figure 9: Exact and approximate Semantic Shift values for the Sports category in AG-News.

Result Analysis. First, SAFARI shows high accuracy in approximating Semantic Shifts. As displayed in Fig. 9, the approximate curve (red) closely aligns with the exact curve (blue), with a high Pearson correlation coefficient of **0.92**, confirming the effectiveness of our approximation method outlined in Section 4.2. Moreover, the dynamic thresholding mechanism in SAFARI, using a sliding window, effectively captures the Semantic Field Subspaces with smaller Semantic Shifts that might otherwise be missed due to varying Semantic Shift magnitudes.

Second, SAFARI captures hierarchical relationships with varying granularity. In Fig. 10, small initial clusters (e.g., individual USA university basketball or football teams) gradually evolve into broader categories like "University Football and Basketball Teams from the USA." Similarly, Fig. 11 shows how teams initially grouped by country later form broader regional clusters, with European teams (blue labels) grouping more closely than non-European teams (read labels). This progression highlights SAFARI's ability to preserve semantic relationships while uncovering the hierarchical structure within the data. Further analysis is provided in Appendix E.1.



484 485 USA football teams.

Figure 10: USA basketball teams merged with Figure 11: Sports teams from different locations merged.

#### 486 5.5 FAKE NEWS EXPLANATION 487

504

507

511 512

513

514

515

516

517

518

519

521

522

523

524

525

526

527

528

529 530

531

488 Experiment Setup. Finally, we showcase how SAFARI can provide detailed explanations using the FakeNewsCorpus dataset (Pathak & Srihari, 2019), which divides fake news articles into ten 489 categories: Bias, Clickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire, and 490 Unreliable. Detailed descriptions are provided in Table D2. Unlike content-based labels in previous 491 datasets, these domain-based labels necessitate deeper explanations to clarify classification reasoning. 492 For example, an article labeled JunkSci might overlap with the Clickbait category, emphasizing 493 the need for precise explanations to avoid ambiguity. This complexity makes this dataset more 494 challenging, showcasing SAFARI's ability to handle nuanced and overlapping categories. 495

We extract entities from each news article and construct ten Semantic Field Subspaces. By comparing 496 these subspaces, we analyze their principal directions and perform a nearest neighbor search to 497 identify the top k Wikipedia entities. To enhance interpretability, each entity is mapped to a Wikidata 498 node (Vrandečić & Krötzsch, 2014), where we examine the associated labels (Ayoola et al., 2022), 499 and the top node types are grouped into broader categories, offering insights into the classification 500 and relationships between the news article categories. 501



Figure 12: Reliable news compared with Clickbait, Hate, and JunkSci news.

**Result Analysis.** The Semantic Field Subspaces constructed by SAFARI for the ten fake news categories display distinct patterns (see Appendix E.2 for details). To explain these differences, we compare each subspace with the Reliable subspace, focusing on three specific categories: Clickbait, Hate, and JunkSci (see Fig. 12). Further comparison results are available in Appendix E.3.

- **Reliable**: Reliable news consistently scores high for node types related to *Science*, *Religion*, and Sports. This suggests that fake news sources prioritize emotionally charged and controversial topics, while Reliable sources focus more on factual, less sensational content.
- **Clickbait:** Articles labeled as Clickbait show high values for the *Fictional works* node type, alongside strong associations with Entertainment and Sports (see Fig. 12(a)). This suggests that clickbait content often revolves around sensational or fictional topics to capture attention, with less emphasis on scientific or factual information.
  - **Hate**: The Hate category emphasizes the *Location* node type, reflecting its frequent associations with political propaganda and regional tensions (see Fig. 12(b)). This geographic specificity is distinct to Hate news, as it amplifies political or cultural division based on location.
- JunkSci: JunkSci sources are linked to the *Game of Chance* and *Entertainment* node types, suggesting a focus on random outcomes and unscientific topics (see Fig. 12(c)). This aligns with the nature of junk science, as it often lacks scientific rigor, unlike reliable sources that emphasize factual content such as science and religion.
- 6 CONCLUSIONS

532 In this paper, we tackled the challenge of understanding the abstract and intricate structure of 533 embedding spaces. We introduced SAFARI, a novel algorithm for identifying Semantic Field 534 Subspaces through hierarchical clustering and the concept of Semantic Shift, with an efficient approximation method that avoids SVD, reducing the computational cost. Extensive experiments 536 on five real-world datasets demonstrated SAFARI's capability to uncover interpretable, hierarchical semantic structures while achieving substantial computational savings. SAFARI has shown to be effective for tasks like classification and explanation. This work not only bridges the gap between 538 embedding spaces and their underlying semantics but also offers new insights into the structure and utility of embedding spaces, enhancing both interpretability and efficiency in their analysis.

# 540 REFERENCES 541

041	
542	Tom Ayoola, Shubhi Tyagi, Joseph Fisher, Christos Christodoulopoulos, and Andrea Pierleoni.
543	ReFinED: An efficient zero-shot-capable approach to end-to-end entity linking. In Proceedings
544	of the 2022 Conference of the North American Chapter of the Association for Computational
545	Linguistics: Human Language Technologies (NAACL-HLT), pp. 209–220, 2022.
546	David & Belsley, Edwin Kub, and Roy E Welsch, Regression diagnostics: Identifying influential data
547	and sources of collinearity. John Wiley & Sons, 2005.
548 549	Leo Breiman. Random Forests. Machine Learning, 45:5-32, 2001.
550	Chih-Chung Chang and Chih-Ien Lin, LIBSVM: A Library for Support Vector Machines. ACM
551 552	Transactions on Intelligent Systems and Technology (TIST), 2(3):1–27, 2011.
553	Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. What does BERT look
554 555	at? an analysis of BERT's attention. In <i>Proceedings of the 2019 ACL Workshop BlackboxNLP:</i> Analyzing and Interpreting Neural Networks for NLP, pp. 276–286, 2019.
556	Themes Course and Deter Hart. Manual asighter rettern classification. IEEE Transactions
557	Information Theory, 13(1):21–27, 1967.
558	Fahim Dalvi, Nadir Durrani, Hassan Saijad, Yonatan Belinkov, Anthony Bau, and James Glass,
559	What is one grain of sand in the desert? Analyzing individual neurons in deep NLP models. In
560 561	Proceedings of the AAAI Conference on Artificial Intelligence (AAAI), pp. 6309–6317, 2019.
562	David Demeter, Gregory Kimmel, and Doug Downey. Stolen probability: A structural weakness
563	of neural language models. In Proceedings of the 58th Annual Meeting of the Association for
564	Computational Linguistics (ACL), pp. 2191–2197, 2020.
565	Jacob Devlin Ming-Wei Chang Kenton Lee and Kristina Toutanova BERT Pre-training of deep
566	bidirectional transformers for language understanding. In <i>Proceedings of the 2019 Conference of</i>
567	the North American Chapter of the Association for Computational Linguistics: Human Language
568	Technologies (NAACL-HLT), pp. 4171–4186, 2019.
569	Dhiling Dufter and Hinrich Schütze. Annlutical methods for interpretable alter dames much adding
570	In Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing
571	and the 9th International Joint Conference on Natural Language Processing (FMNI P-IICNI P)
572	pp. 1185–1191. 2019.
573	pp. 1105 1191, 2019.
574	Kawin Ethayarajh. How contextual are contextualized word representations? Comparing the geometry
575	of BERT, ELMo, and GPT-2 embeddings. In <i>Proceedings of the 2019 Conference on Empirical</i>
576 577	Language Processing (EMNLP-IJCNLP), pp. 55–65, 2019.
578	
579	Evelyn Fix. Discriminatory Analysis. Nonparametric Discrimination: Consistency Properties,
580	volume 1. USAF School of Aviation Medicine, 1985.
581	Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness:
582	Stochastic algorithms for constructing approximate matrix decompositions. arXiv preprint
583	arXiv:0909.4061, 2009.
584 585	Roger A Horn and Charles R Johnson. Matrix analysis. Cambridge university press, 2012.
586	Jure Leskovec, Anand Rajaraman, and Jeffrey David Illiman, Mining of Massive Datasets, Cambridge
587	University Press, 2020.
588	
589	Tianlin Liu, Lyle Ungar, and Joao Sedoc. Unsupervised post-processing of word vectors via conceptor
590	negation. In Proceedings of the AAAI Conference on Artificial Intelligence (AAAI), pp. 6778–6785,
591	2019.
592 593	Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th Annual Meeting of the</i>
333	Association for Computational Linguistics (ACL), pp. 142–150, 2011.

594 595	Carl D Meyer. Matrix analysis and applied linear algebra. SIAM, 2023.
596 597 598	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representa- tions of words and phrases and their compositionality. In <i>Proceedings of the 26th International</i> <i>Conference on Neural Information Processing Systems (NIPS)</i> , pp. 3111–3119, 2013.
599 600	Jiaqi Mu and Pramod Viswanath. All-but-the-Top: Simple and effective postprocessing for word representations. In <i>International Conference on Learning Representations (ICLR)</i> , 2018.
602 603 604	Sungjoon Park, JinYeong Bak, and Alice Oh. Rotated word vector representations and their inter- pretability. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language</i> <i>Processing (EMNLP)</i> , pp. 401–411, 2017.
605 606 607	Archita Pathak and Rohini K Srihari. BREAKING! Presenting fake news corpus for automated fact checking. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop</i> , pp. 357–362, 2019.
608 609 610	John Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. <i>Advances in Large Margin Classifiers</i> , 10(3):61–74, 1999.
611 612	Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. <i>Introduction to Information Retrieval</i> . Cambridge University Press, 2008.
613 614 615 616	Adi Simhi and Shaul Markovitch. Interpreting embedding spaces by conceptualization. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 1704–1719, 2023.
617	Gilbert W Stewart. Perturbation Theory for the Singular Value Decomposition. 1998.
618 619	Lloyd N Trefethen and David Bau. Numerical Linear Algebra. SIAM, 2022.
620 621	Denny Vrandečić and Markus Krötzsch. Wikidata: A free collaborative knowledgebase. <i>Communications of the ACM</i> , 57(10):78–85, 2014.
623 624 625	Hermann Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgle- ichungen (mit einer anwendung auf die theorie der hohlraumstrahlung). <i>Mathematische Annalen</i> , 71(4):441–479, 1912.
626 627 628	Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. Scalable zero-shot entity linking with dense entity retrieval. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 6397–6407, 2020.
629 630 631 632	Pengcheng Yang, Xu Sun, Wei Li, Shuming Ma, Wei Wu, and Houfeng Wang. SGM: Sequence generation model for multi-label classification. In <i>Proceedings of the 27th International Conference on Computational Linguistics (COLING)</i> , pp. 3915–3926, 2018.
633 634 635 636	Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classifi- cation. In <i>Proceedings of the 28th International Conference on Neural Information Processing</i> <i>Systems (NIPS)</i> , pp. 649–657, 2015.
637 638 639 640	
641 642 643	
644 645 646	
647	

#### TABLE OF NOTATIONS А

648

649

698

699

700

701

Symbol	Description
$x, \mathcal{X}$	A set $\mathcal{X}$ of textual elements $\boldsymbol{x}$ , i.e., $\boldsymbol{x} \in \mathcal{X}$
$v, \mathcal{E}$	A set $\mathcal{E}$ of embedding vectors $v$ in a $d$ -dimensional embedding space $\mathbb{R}^d$ , i.e., $v \in \mathcal{E}$ and $\mathcal{E}$
h $(n n')$	A deep learning model: $\mathcal{X} \to \mathcal{E}$ The semantic distance between any two embedding vectors $\boldsymbol{u}$ and $\boldsymbol{v}'$
$\mathcal{M}^{u_{sem}(v,v)}$	A semantic set
$\mathcal{N}(oldsymbol{v})$	The close neighborhood of an embedding vector $\boldsymbol{v}$
$f_{sem}(oldsymbol{v})$	The semantics of an embedding vector $v$ , i.e., $f_{sem}(v) : \mathbb{R}^d \to 2^{ \mathcal{M} } \setminus \{\emptyset\}$
$f_{int}(\boldsymbol{v})$	The interpretable semantics of an embedding vector $v$ , i.e., $f_{int}(v) : \mathbb{R}^d \to 2^{ \mathcal{M} } \setminus \{\emptyset\}$
$F_{int}(\mathcal{C})$	The semantic field of a subset of embedding vectors $C \subseteq E$ , i.e., $F_{int}(C) : \mathbb{R}^d \to 2^{ \mathcal{M} } \setminus \{ e^{i \pi t} \in E^{-1} \}$
$F_{sem}(\mathbb{S})$ $\Delta F_{som}(\mathcal{C}_1, \mathcal{C}_2)$	The semantic field subspace of a subspace $\mathbb{S} \subseteq \mathbb{R}^n$ , i.e., $F_{sem}(\mathbb{S}) : 2^{m-1} \setminus \{\emptyset\} \to 2^{m+1} \setminus \{\emptyset\}$ The semantic shift of any two clusters $\mathcal{C}_1$ and $\mathcal{C}_2$
B PSEUD	O-CODE FOR EXACT SEMANTIC SHIFT COMPUTATION
The pseudo-co	ode for exact Semantic Shift computation is depicted in Algorithm 2.
Algorithm 2.	Exact Semantic Shift Computation
Algorithm 2.	
T 4 T	
Input: Larger	cluster $C_x$ and new cluster $C_{new}$ ;
Input: Larger Output: Exac	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ;
Input: Larger Output: Exact $A_x$ and $A_{new}$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $\sum_{sem}(C_x, C_{new})$ using Eq. (7);
Input: Larger Output: Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ return $\Delta F_{sem}$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $\sum_{sem}(C_x, C_{new})$ using Eq. (7); $_n(C_x, C_{new})$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sem}$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $S_{sem}(C_x, C_{new})$ using Eq. (7); ${}_n(C_x, C_{new})$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sem}$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $e_{sem}(C_x, C_{new})$ using Eq. (7); ${}_n(C_x, C_{new})$ ;
Input: Larger Output: Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ return $\Delta F_{sem}$ C PSEUDO	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $E_{sem}(C_x, C_{new})$ using Eq. (7); ${}_n(C_x, C_{new})$ ;
Input: Larger Output: Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ return $\Delta F_{sem}$ C PSEUDO The pseudo-co	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $S_{sem}(C_x, C_{new})$ using Eq. (7); ${}_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3.
Input: Larger Output: Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ return $\Delta F_{sem}$ C PSEUDO The pseudo-coordinates of the pse	The cluster $C_x$ and new cluster $C_{new}$ ; that Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $S_{sem}(C_x, C_{new})$ using Eq. (7); $\alpha(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION bde for approximate Semantic Shift computation is depicted in Algorithm 3.
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sem}$ <b>C PSEUD</b> The pseudo-complete <b>Algorithm 3:</b>	c cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $S_{sem}(C_x, C_{new})$ using Eq. (7); ${}_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3.
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sen}$ <b>C PSEUD</b> The pseudo-co <b>Algorithm 3:</b> <b>Input:</b> Larger	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $e_{sem}(C_x, C_{new})$ using Eq. (7); ${}_{n}(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation cluster $C_x$ and smaller cluster $C_y$ ;
Input: Larger Output: Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ return $\Delta F_{sen}$ C PSEUD The pseudo-co Algorithm 3: Input: Larger Output: Appr	c cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $e_{sem}(C_x, C_{new})$ using Eq. (7); ${}_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $c$ cluster $C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sen}$ <b>C PSEUDO</b> The pseudo-cod <b>Algorithm 3:</b> <b>Input:</b> Larger <b>Output:</b> Appr $A_x$ and $A_y \leftarrow$	c cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $e_{sem}(C_x, C_{new})$ using Eq. (7); $_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $c$ cluster $C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ; - Construct two matrices from $C_x$ and $C_y$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sen}$ <b>C P</b> SEUDO The pseudo-co <b>Algorithm 3:</b> <b>Input:</b> Larger <b>Output:</b> Appr $A_x$ and $A_y \leftarrow \sigma_{max} \leftarrow \operatorname{Com}$	c cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $e_{sem}(C_x, C_{new})$ using Eq. (7); $_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $Cluster C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ; - Construct two matrices from $C_x$ and $C_y$ ; pute the maximum singular value from $A_x$ ;
Input: LargerOutput: Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ return $\Delta F_{sen}$ CPSEUDOThe pseudo-cocAlgorithm 3:Input: LargerOutput: Appi $A_x$ and $A_y \leftarrow \sigma_{max} \leftarrow \operatorname{Com}$ Compute $\Delta \tilde{F}_s$	c cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $e_{em}(C_x, C_{new})$ using Eq. (7); $_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $c$ cluster $C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ; - Construct two matrices from $C_x$ and $C_y$ ; pute the maximum singular value from $A_x$ ; $e_{em}(C_x, C_{new})$ using Eq. (8);
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sen}$ <b>C PSEUDO</b> The pseudo-coc <b>Algorithm 3:</b> <b>Input:</b> Larger <b>Output:</b> Appin $A_x$ and $A_y \leftarrow \sigma_{max} \leftarrow \text{Com}$ Compute $\Delta \tilde{F}_s$ <b>return</b> $\Delta \tilde{F}_{sen}$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $e_{em}(C_x, C_{new})$ using Eq. (7); $_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $c$ cluster $C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ; - Construct two matrices from $C_x$ and $C_y$ ; pute the maximum singular value from $A_x$ ; $e_{em}(C_x, C_{new})$ using Eq. (8); $_n(C_x, C_{new})$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sen}$ <b>C PSEUD</b> The pseudo-co <b>Algorithm 3:</b> <b>Input:</b> Larger <b>Output:</b> Appr $A_x$ and $A_y \leftarrow \sigma_{max} \leftarrow \text{Com}$ Compute $\Delta \tilde{F}_s$ <b>return</b> $\Delta \tilde{F}_{sen}$	cluster $C_x$ and new cluster $C_{new}$ ; et Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $\sum_{sem}(C_x, C_{new})$ using Eq. (7); $_1(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $c$ cluster $C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ; - Construct two matrices from $C_x$ and $C_y$ ; pute the maximum singular value from $A_x$ ; $\sum_{sem}(C_x, C_{new})$ using Eq. (8); $_1(C_x, C_{new})$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sen}$ <b>C PSEUDO</b> The pseudo-co <b>Algorithm 3:</b> <b>Input:</b> Larger <b>Output:</b> Appr $A_x$ and $A_y \leftarrow \sigma_{max} \leftarrow \text{Com}$ Compute $\Delta \tilde{F}_s$ <b>return</b> $\Delta \tilde{F}_{sen}$	c cluster $C_x$ and new cluster $C_{new}$ ; t t Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $sem(C_x, C_{new})$ using Eq. (7); $_i(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $c$ cluster $C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ; - Construct two matrices from $C_x$ and $C_y$ ; put the maximum singular value from $A_x$ ; $sem(C_x, C_{new})$ using Eq. (8); $_i(C_x, C_{new})$ ;
<b>Input:</b> Larger <b>Output:</b> Exact $A_x$ and $A_{new}$ $F_{sem}(\mathbb{S}_x) \approx \Sigma$ Compute $\Delta F_s$ <b>return</b> $\Delta F_{sen}$ <b>C PSEUDO</b> The pseudo-cod <b>Algorithm 3:</b> <b>Input:</b> Larger <b>Output:</b> Appr $A_x$ and $A_y \leftarrow \sigma_{max} \leftarrow \text{Com}$ Compute $\Delta \tilde{F}_s$ <b>return</b> $\Delta \tilde{F}_{sen}$ <b>D DATAS</b>	c cluster $C_x$ and new cluster $C_{new}$ ; tt Semantic Shift $\Delta F_{sem}(C_x, C_{new})$ ; $\leftarrow$ Construct two matrices from $C_x$ and $C_{new}$ ; $\Sigma_x, V_x^{\top}$ and $F_{sem}(\mathbb{S}_{new}) \approx \Sigma_{new}, V_{new}^{\top} \leftarrow$ Perform SVD on $A_x$ and $A_{new}$ ; $_{sem}(C_x, C_{new})$ using Eq. (7); $_n(C_x, C_{new})$ ; O-CODE FOR APPROXIMATE SEMANTIC SHIFT COMPUTATION ode for approximate Semantic Shift computation is depicted in Algorithm 3. Approximate Semantic Shift Computation $c$ cluster $C_x$ and smaller cluster $C_y$ ; roximate Semantic Shift $\Delta \tilde{F}_{sem}(C_x, C_{new})$ ; - Construct two matrices from $C_x$ and $C_y$ ; pute the maximum singular value from $A_x$ ; $_{sem}(C_x, C_{new})$ using Eq. (8); $_n(C_x, C_{new})$ ; ET DETAILS

• **AG-News:** This dataset comprises over 1 million news articles from more than 2,000 sources

- (Zhang et al., 2015), commonly used for clustering, classification, ranking, and search. Each article is categorized under one of four labels: Business, Sci/Tech, Sports, and World.
- AAPD: This Arxiv Academic Paper Dataset (AAPD) contains 55,840 abstracts and subjects of computer science papers developed for multi-label classification tasks (Yang et al., 2018).

Label	Description
Bias	Sources that come from a particular point of view and may rely on propaganda, dec textualized information, and opinions distorted as facts.
Clickbait	Sources that provide generally credible content but use exaggerated, misleading, questionable headlines, social media descriptions, and/or images.
Conspiracy	Sources that are well-known promoters of kooky conspiracy theories.
Fake	Sources that entirely fabricate information, disseminate deceptive content, or gros
Hate	distort actual news reports. Sources that actively promote racism, misogyny, homophobia, and other forms
JunkSci	Sources that promote pseudoscience, metaphysics, naturalistic fallacies, and ot scientifically dubious claims
Political	Sources that provide generally verifiable information in support of certain points view or political orientations
Reliable	Sources that circulate news and information in a manner consistent with traditional a ethical practices in journalism.
Satire	Sources that use humor, irony, exaggeration, ridicule, and false information to comm
Unreliable	Sources that may be reliable but whose contents require further verification
	ickhait Conspiracy Fake Hate JunkSci Political Reliable Satire
Blas, Cl Unreliat Experiment E on a machine	<b>ickbait</b> , <b>Conspiracy</b> , <b>Fake</b> , <b>Hate</b> , <b>JunkSci</b> , <b>Political</b> , <b>Reliable</b> , <b>Satire</b> , <b>ble</b> . The descriptions of these ten labels are presented in Table D2. <b>invironment.</b> All methods were written in Python 3.8. All experiments were condumited with Intel <sup>®</sup> Xeon <sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run
Blas, Cl Unreliab Experiment E on a machine on Ubuntu 20.	<b>ickbait</b> , <b>Conspiracy</b> , <b>Fake</b> , <b>Hate</b> , <b>JunkSci</b> , <b>Political</b> , <b>Reliable</b> , <b>Satire</b> , <b>ble</b> . The descriptions of these ten labels are presented in Table D2. <b>invironment.</b> All methods were written in Python 3.8. All experiments were conduct with Intel <sup>®</sup> Xeon <sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run 04.
Experiment E on a machine on Ubuntu 20.	<b>ickbait</b> , <b>Conspiracy</b> , <b>Fake</b> , <b>Hate</b> , <b>JunkSci</b> , <b>Political</b> , <b>Reliable</b> , <b>Satire</b> , <b>ble</b> . The descriptions of these ten labels are presented in Table D2. <b>invironment.</b> All methods were written in Python 3.8. All experiments were conduce with Intel <sup>®</sup> Xeon <sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run 04. EXPERIMENTAL RESULTS
Experiment E on a machine on Ubuntu 20. E EXTRA E.1 MORE	<b>ickbait</b> , <b>Conspiracy</b> , <b>Fake</b> , <b>Hate</b> , <b>JunkSci</b> , <b>Political</b> , <b>Reliable</b> , <b>Satire</b> , <b>ble</b> . The descriptions of these ten labels are presented in Table D2. <b>invironment.</b> All methods were written in Python 3.8. All experiments were conduced with Intel <sup>®</sup> Xeon <sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run 04. EXPERIMENTAL RESULTS ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE
Experiment E on a machine on Ubuntu 20. E EXTRA E.1 MORE	<ul> <li>Consider the news articles that are categorized based on their domains under ten fait ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire, ole. The descriptions of these ten labels are presented in Table D2.</li> <li>Cnvironment. All methods were written in Python 3.8. All experiments were conducted with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified</li> </ul>
Experiment E on a machine y on Ubuntu 20. E EXTRA E.1 MORE A in this section SAFARI in en	<ul> <li>Consider the news articles that are categorized based on their domains under ten if ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire ole. The descriptions of these ten labels are presented in Table D2.</li> <li>Convironment. All methods were written in Python 3.8. All experiments were conducted with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, rur 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified abedding spaces and the more intuitive hierarchies found in natural human language</li> </ul>
Experiment E on a machine on Ubuntu 20. E EXTRA E.1 MORE A in this section SAFARI in em numan langua	<ul> <li>Consider the news articles that are categorized based on their domains under ten fait ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire Dle. The descriptions of these ten labels are presented in Table D2.</li> <li>Cnvironment. All methods were written in Python 3.8. All experiments were conducted with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, rur 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human languag ge, semantics typically follow a logical hierarchy, progressing from specific, con</li> </ul>
Experiment E on a machine on Ubuntu 20. E EXTRA E.1 MORE in this section SAFARI in en numan langua entities to mo	<ul> <li>Consider the news articles that are categorized based on their domains under ten fait ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire, ole. The descriptions of these ten labels are presented in Table D2.</li> <li>Chvironment. All methods were written in Python 3.8. All experiments were conducted with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified abedding spaces and the more intuitive hierarchies found in natural human languag ge, semantics typically follow a logical hierarchy, progressing from specific, con re abstract concepts, much like an ontology. However, in embedding spaces, not eliver in the distribution.</li> </ul>
Experiment E on a machine y on Ubuntu 20. E EXTRA E.1 MORE A In this section SAFARI in em numan langua entities to mo progression is on the data on	<ul> <li>consider the news articles that are categorized based on their domains under ten fait ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire, ole. The descriptions of these ten labels are presented in Table D2.</li> <li>chvironment. All methods were written in Python 3.8. All experiments were conducted with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, rur 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human languag ge, semantics typically follow a logical hierarchy, progressing from specific, con re abstract concepts, much like an ontology. However, in embedding spaces, not always intuitive. The distinction between "specific" and "abstract" depends in a diverse form.</li> </ul>
Experiment E on a machine y on Ubuntu 20. E EXTRA E.1 MORE A n this section SAFARI in en uman langua entities to mo rogression is on the data an ye would expo	<ul> <li>consider the news articles that are categorized based on their domains under ten fait ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire, ole. The descriptions of these ten labels are presented in Table D2.</li> <li>cnvironment. All methods were written in Python 3.8. All experiments were conduct with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human languag ge, semantics typically follow a logical hierarchy, progressing from specific, comre abstract concepts, much like an ontology. However, in embedding spaces, not always intuitive. The distinction between "specific" and "abstract" depends in dimodel than on human reasoning, often leading to groupings that diverge from text based on natural language understanding.</li> </ul>
Experiment E on a machine y on Ubuntu 20. E EXTRA E.1 MORE A n this section SAFARI in en iuman langua entities to mo yrogression is on the data an ye would exper-	<ul> <li>Consider the news articles that are categorized based on their domains under ten fait is included that are categorized based on their domains under ten fait is included to the included tent of ten</li></ul>
Experiment E on a machine on Ubuntu 20. E EXTRA E.1 MORE A in this section SAFARI in en numan langua entities to mo orogression is on the data an we would exper- For example,	<ul> <li>consider the news articles that are categorized based on their domains under ten if ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire ole. The descriptions of these ten labels are presented in Table D2.</li> <li>cnvironment. All methods were written in Python 3.8. All experiments were conducted with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, rur 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human language, semantics typically follow a logical hierarchy, progressing from specific, con re abstract concepts, much like an ontology. However, in embedding spaces, not always intuitive. The distinction between "specific" and "abstract" depends in model than on human reasoning, often leading to groupings that diverge from ect based on natural language understanding.</li> </ul>
E EXTRA E.1 MORE In this section SAFARI in en uman langua entities to mo orogression is on the data an we would exper- For example, eams, and late	<ul> <li>consider the news articles that are categorized based on their domains under ten if ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire, De. The descriptions of these ten labels are presented in Table D2.</li> <li>chvironment. All methods were written in Python 3.8. All experiments were conduce with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, rur 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human languag ge, semantics typically follow a logical hierarchy, progressing from specific, con re abstract concepts, much like an ontology. However, in embedding spaces, not always intuitive. The distinction between "specific" and "abstract" depends of model than on human reasoning, often leading to groupings that diverge from ect based on natural language understanding.</li> <li>as illustrated in Fig. 10, USA basketball teams are first grouped with USA footer, sports teams from various locations are merged, as shown in Fig. 11. This followed is specific categories to broader ones. Yet as shown</li> </ul>
Experiment E on a machine y on Ubuntu 20. E EXTRA E.1 MORE A In this section SAFARI in em numan langua entities to mo progression is on the data an we would expe For example, eams, and lata a logical hiera	<ul> <li>consider the news articles that are categorized based on their domains under ten if ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire ole. The descriptions of these ten labels are presented in Table D2.</li> <li>cnvironment. All methods were written in Python 3.8. All experiments were conducted with Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, rur 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human language, semantics typically follow a logical hierarchy, progressing from specific, con re abstract concepts, much like an ontology. However, in embedding spaces. not always intuitive. The distinction between "specific" and "abstract" depends d model than on human reasoning, often leading to groupings that diverge from ect based on natural language understanding.</li> <li>as illustrated in Fig. 10, USA basketball teams are first grouped with USA focer, sports teams from various locations are merged, as shown in Fig. 11. This fol urchical structure, from more specific categories to broader ones. Yet, as show ation 19 790), entities such as horse racing clubs companies and events (e.g., "International structure, from the specific categories to broader ones. Yet, as show ation 19 790).</li> </ul>
Experiment E on a machine y on Ubuntu 20. E EXTRA E.1 MORE A n this section SAFARI in en numan langua entities to mo progression is on the data an ve would expe for example, eams, and late logical hiera fig. E1 (at iter Club') are me	<ul> <li>consider the news articles that are categorized based on their domains under ten if ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire, ole. The descriptions of these ten labels are presented in Table D2.</li> <li>chvironment. All methods were written in Python 3.8. All experiments were conducted with Intel® Xeon® Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, rur 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human languag ge, semantics typically follow a logical hierarchy, progressing from specific, con re abstract concepts, much like an ontology. However, in embedding spaces, not always intuitive. The distinction between "specific" and "abstract" depends red model than on human reasoning, often leading to groupings that diverge from ect based on natural language understanding.</li> <li>as illustrated in Fig. 10, USA basketball teams are first grouped with USA footer, sports teams from various locations are merged, as shown in Fig. 11. This fol urchical structure, from more specific categories to broader ones. Yet, as show ation 19,790), entities such as horse racing clubs, companies, and events (e.g., 'Jo reed with famous racing horses, This merging of horse racing happens thousand</li> </ul>
E EXTRA E EXTRA E EXTRA E EXTRA E EXTRA E EXTRA E EXTRA E EXTRA C 1 MORE A n this section CAFARI in en uman langua ntities to mo rogression is n the data an re would expe for example, eams, and late logical hiera rig. E1 (at iter E lub') are me erations after	<ul> <li>consider the news articles that are categorized based on their domains under ten if ickbait, Conspiracy, Fake, Hate, JunkSci, Political, Reliable, Satire ole. The descriptions of these ten labels are presented in Table D2.</li> <li>cnvironment. All methods were written in Python 3.8. All experiments were condiwith Intel<sup>®</sup> Xeon<sup>®</sup> Platinum 8480C, 2.0 TB memory, and one NVIDIA H100, run 04.</li> <li>EXPERIMENTAL RESULTS</li> <li>ANALYSIS ON HIERARCHICAL SEMANTIC STRUCTURE</li> <li>, we explore the differences between the hierarchical semantic structures identified bedding spaces and the more intuitive hierarchies found in natural human language, semantics typically follow a logical hierarchy, progressing from specific, con re abstract concepts, much like an ontology. However, in embedding spaces not always intuitive. The distinction between "specific" and "abstract" depends d model than on human reasoning, often leading to groupings that diverge from ect based on natural language understanding.</li> <li>as illustrated in Fig. 10, USA basketball teams are first grouped with USA foc er, sports teams from various locations are merged, as shown in Fig. 11. This fol urchical structure, from more specific categories to broader ones. Yet, as show ation 19,790, entities such as horse racing clubs, companies, and events (e.g., 'Jo rged with famous racing horses. This merging of horse racing happens thousan the merging of football and basketball teams in the USA. Following the ontoloop</li> </ul>

<sup>&</sup>lt;sup>2</sup>https://www.imdb.com/ <sup>3</sup>https://www.yelp.com/dataset <sup>4</sup>https://github.com/several27/FakeNewsCorpus



progression, we would expect more abstract concepts. However, horse racing is not a more abstract concept compared with other sports.

800 These examples underscore that the hierarchical structures found in embedding spaces are shaped 801 by the model's learning patterns, not by human logic. While some structures align with natural 802 expectations, others can be surprising, revealing the complex relationship between the data and the 803 model. This highlights the importance of careful interpretation when analyzing embeddings, as the 804 resulting hierarchies may not always reflect conventional semantic reasoning.

#### 806 **EXPLANATION PATTERNS FOR TEN FAKE NEWS CATEGORIES** E.2

807

805

799

Figure E2 presents the explanation patterns of the Semantic Field Subspaces constructed by SAFARI 808 for the ten distinct fake news categories. Each subspace reveals different nearest entity types, 809 showcasing the nuanced relationships between these categories. For instance, the subspaces highlight



Figure E3: More comparisons between Reliable and other categories of fake news.

key differences in how content is structured across categories like Clickbait, Hate, and JunkSci, each emphasizing different entity types.

# E.3 MORE COMPARISON RESULTS BETWEEN RELIABLE AND OTHER CATEGORIES

Fig. E3 depicts more comparisons between Reliable and other categories of fake news. Before analyzing the results, we want to highlight that although this dataset is labeled as "fake news", the categories presented here are vague or not necessarily fake. These categories are not clearly distinct from each other; for instance, a news piece labeled as Fake could also fall under Bias or Unreliable simultaneously, making it unreasonable to strictly assign it to one category over another. Consequently, the analysis below may be based on an inherently uncertain foundation.

- **Bias**: The node types for Bias are similar to those for Reliable, particularly in areas like *Location* and *Entertainment* (see Fig. 3(a)). This is consistent with the dataset's description of biased content, which often involves propaganda or opinions presented as facts, rather than outright falsehoods. One notable difference is that Bias lacks prominence in node types like *Science, Religion*, and *Sports*, which appear more in Reliable news.
- Conspiracy: The node types for Conspiracy differ markedly from Reliable, with a notable emphasis on *Location* and *Entertainment* (see Fig. 3(b)). This aligns with the nature of conspiracy theories, which often center around speculative narratives involving covert operations or unverified claims. In contrast, Reliable news exhibits a more balanced spread across categories like *Science* and *Sports*, highlighting the factual and grounded nature of its content. The disparity between these node type patterns underscores the speculative and unverified themes prevalent in conspiracy content.
- Fake: The Fake category shows a scattered pattern in node types, with no strong emphasis on any particular theme (see Fig. 3(c)). This reflects the diversity and often random nature of fake news, which may cover a wide range of topics with little factual basis. Compared to Reliable news, which focuses on well-defined topics such as *Science, Religion*, and *Sports*, Fake news lacks consistency, mirroring the dataset's characterization of this category as containing misleading or fabricated stories.
- Political: As expected, the node type patterns for Political news are closely aligned with Reliable, with significant overlap in categories such as *Science, Location, Entertainment*, and *Sports* (see Fig. 3(d)). This similarity is consistent with the dataset's description of Political news as verifiable but presented with a specific viewpoint. The comparison with Reliable indicates that Political content, while biased toward certain ideologies or perspectives, does not stray far from reliable reporting in terms of the types of entities discussed.

- Satire: The Satire category shows high values in *Location, Entertainment*, and *Fictional works*, which contrasts sharply with the more factual node types (e.g., *Science* and *Religion*) found in Reliable news (see Fig. 3(e)). Satire often uses humor, exaggeration, and fictional scenarios to comment on real-world events, explaining the prominence of such node types. This divergence highlights the entertainment-driven and often fictional nature of satire, which deliberately distorts facts for comedic or critical purposes, unlike the more serious and factual content of Reliable news.
- Unreliable: The Unreliable category displays notable differences compared to Reliable news sources. Except for *Location*, the node types associated with Unreliable sources are scattered across diverse, less cohesive categories (see Fig. 3(f)). This suggests that Unlike Reliable news, which demonstrates a strong association with node types related to well-established factual domains like *Science*, *Religion*, and *Sports*, Unreliable sources exhibit a pattern of fragmentation. This reflects the unpredictable and often erratic nature of unreliable news content, where the underlying information may lack verification or coherence.