DISENTANGLING LINGUISTIC FEATURES WITH DIMENSION-WISE ANALYSIS OF VECTOR EMBED DINGS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027 028 029 Paper under double-blind review

Abstract

Understanding the inner workings of neural embeddings, particularly in models such as BERT, remains a challenge because of their high-dimensional and opaque nature. This paper proposes a framework for uncovering the specific dimensions of vector embeddings that encode distinct linguistic properties (LPs). We introduce the Linguistically Distinct Sentence Pairs (LDSP-10) dataset, which isolates ten key linguistic features such as synonymy, negation, tense, and quantity. Using this dataset, we analyze BERT embeddings with various statistical methods, including the Wilcoxon signed-rank test, mutual information, and recursive feature elimination, to identify the most influential dimensions for each LP. We introduce a new metric, the Embedding Dimension Importance (EDI) score, which quantifies the relevance of each embedding dimension to a LP. Our findings show that certain properties, such as negation and polarity, are robustly encoded in specific dimensions, while others, like synonymy, exhibit more complex patterns. This study provides insights into the interpretability of embeddings, which can guide the development of more transparent and optimized language models, with implications for model bias mitigation and the responsible deployment of AI systems.¹

1 INTRODUCTION

Word embeddings are central to natural language processing (NLP), enabling machines to represent and interpret text in continuous vector spaces. From early models like Word2Vec Mikolov et al. (2013) and GloVe Pennington et al. (2014), to advanced models like GPT-2 Radford et al. (2019) and BERT Devlin et al. (2019), embeddings have evolved to capture complex linguistic nuances. BERT, in particular, leverages bidirectional transformers to generate contextualized word representations, enhancing syntactic and semantic understanding Rogers et al. (2020).

Despite these advancements, embeddings are often seen as "black boxes," where the high-dimensional nature of the spaces they occupy makes interpretation difficult Belinkov & Glass (2019). The field of interpretable embeddings seeks to address these challenges by making the dimensions of embeddings more transparent and meaningful Faruqui et al. (2015a); Incitti et al. (2023); Snidaro et al. (2019). However, most systems still rely on popular embedding models like GPT, BERT, Word2Vec, and GloVe, which prioritize performance over interpretability Cao (2024); Lipton (2017).

Our research introduces a generalizable framework for identifying specific embedding dimensions in models like BERT and GPT-2 that encode distinct LPs. This work responds to the growing need for interpretable models, especially for tasks like bias mitigation Bolukbasi et al. (2016); Mehrabi et al. (2021), task-specific optimization Guyon & Elisseeff (2003); Voita et al. (2019), and more system controllability Bau et al. (2019).

We present the LDSP-10 dataset, which consists of sentence pairs isolating nine LPs, designed to probe embedding spaces and identify the dimensions most influential for each property. We analyze these sentence pairs using statistical tests, mutual information, and feature selection methods. We propose the Embedding Dimension Importance (EDI) score, which aggregates these analyses to quantify the relevance of each dimension to specific LPs.

¹Code will be released upon publication.



Figure 1: Dimensions of BERT embeddings that encode the most information about each LP. Relevance is determined by Embedding Dimension Importance (EDI) scores above 0.8, a threshold chosen in relation to the general EDI score distribution.

Γ		Control	Synonym	Quantity	Tense	Intensifier	Voice	Definiteness	Factuality	Polarity	Negation
	BERT	0.5033	0.7033	0.95	0.94	0.9867	0.9667	0.8967	0.9833	0.9700	0.9333
	GPT-2	0.57	0.6267	0.9733	0.9567	0.9367	0.9867	0.9433	0.9667	0.9533	0.93
	MP-Net	0.54	0.5267	0.9533	0.93	0.8733	0.86	0.8567	0.9667	0.9533	0.9367

Table 1: Evaluation 1 (§ 5.2) accuracy for different LPs across BERT, GPT-2, and MP-Net. A simple logistic classifier is able to perform at these levels of accuracy on the highest EDI subset of dimensions of embeddings from each of these models.

This paper makes three contributions. First, is the introduction of the LDSP-10 dataset, consisting of sentence pairs that isolate nine LPs. Second is a generalizable framework and quantifiable metric (EDI score) for identifying influential embedding dimensions, applicable to different models and linguistic features. Third is a comprehensive analysis of BERT, GPT-2, and MPNet embeddings, revealing key dimensions related to each LP.

2 RELATED WORKS

Research on interpretable embeddings can be divided into two categories: interpretable embeddings and representation analysis. The former focuses on designing models that naturally produce interpretable representations, while the latter involves post-hoc analysis to uncover how existing embeddings encode human-interpretable features.

2.1 INTERPRETABLE EMBEDDINGS

Several approaches have been proposed to create interpretable word embeddings. Early efforts like Murphy et al. (2012) used matrix factorization techniques to generate sparse, interpretable embeddings. Faruqui et al. (2015b) introduced Sparse Overcomplete Word Vectors (SPOWV), which used a dictionary learning framework for more interpretable, sparse embeddings. Other methods, such as Guillot et al. (2023) and Subramanian et al. (2018), explored how sparsification techniques could disentangle properties within embeddings, making them more interpretable.

 Approaches to embedding interpretability also involve aligning dimensions with humanunderstandable concepts. For instance, Panigrahi et al. (2019) used Latent Dirichlet Allocation (LDA) to produce embeddings where each dimension corresponds to a specific word sense, and Benara et al. (2024) employed LLM-powered yes/no question-answering techniques to generate interpretable embeddings. Despite these innovations, popular models like Word2Vec, GloVe, and BERT remain dominant in NLP but often lack inherent interpretability. As a result, methods for post-hoc analysis are needed to interpret these embeddings.

108 2.2 REPRESENTATION ANALYSIS

110 Representation analysis focuses on understanding how knowledge is structured within embeddings 111 and how individual neurons contribute to encoding specific properties Sajjad et al. (2022). Senel et al. (2017) demonstrated how individual dimensions correspond to specific semantic properties, 112 and Zhu et al. (2018) emphasized the value of sentence-level embeddings in capturing nuanced 113 semantic properties. Research has also explored the linguistic features encoded within embeddings. 114 Conneau et al. (2018) developed a set of ten probing tasks that evaluate how sentence embeddings 115 capture various linguistic features, such as syntactic structures and semantic roles. Adi et al. (2017) 116 complemented this work by proposing classification tasks that reveal the effectiveness of sentence 117 embeddings in encoding attributes like sentence length and word order. 118

- Recent research has analyzed individual neurons in embedding spaces, often using methods like
 neuron-ranking, where a probe is used to rank neurons based on their relevance to a specific linguistic
 feature Dalvi et al. (2019); Durrani et al. (2020); Torroba Hennigen et al. (2020). Antverg & Belinkov
 (2022) analyzed these methods, separating representational importance from functional utility and
 introducing interventions to evaluate whether encoded information is actively utilized.
- 124 Building on this foundation, Durrani et al. (2024) introduced Linguistic Correlation Analy-125 sis (LCA), which identifies salient neurons that 126 encode specific linguistic features. Their find-127 ings indicated redundancy in information en-128 coding across neurons, enhancing robustness in 129 representation learning. Similarly, Gurnee et al. 130 (2023) proposed sparse probing methods to ad-131 dress polysemanticity, illustrating how features 132 are distributed across neurons in transformer 133 models. Additionally, Torroba Hennigen et al. 134 (2020) presented intrinsic probing, introducing 135 a Gaussian framework to identify dimensions encoding LPs. We Together, these findings sug-136 gest that linguistic attributes are often encoded 137 in focal dimensions, providing insights into how 138 different models represent linguistic knowledge. 139
- Our work builds on these ideas by using the LDSP-10 dataset to isolate linguistic features, which provides a focused method for assessing how embedding dimensions capture these properties. We move beyond traditional probing and neuron-ranking techniques to offer a more targeted examination of embedding interpretability.

Property	Sentence Pair			
Control	They sound excited.			
Control	The farmer has 20 sheep.			
Supanum	The music was calming.			
Synonym	The music was soothing.			
Quantity	I ate two cookies.			
Quantity	I ate several cookies.			
Tansa	The river flows swiftly.			
Telise	The river flowed swiftly.			
Intensifier	The task is easy.			
mensmer	The task is surprisingly easy			
Vaiaa	The team won the game.			
voice	The game was won by the team.			
Definiteness	The bird flew away.			
Deminteness	A bird flew away			
Feetuelity	The car is red.			
Factuality	The car could be red.			
Polority	She passed the exam.			
roidfity	She failed the exam.			
Negation	The project is successful.			
regation	The project is not successful.			

Table 2: Sample linguistically distinct sentence pairs (LDSPs) from each of the LPs tested in this study. LDSP-10 dataset contains 1000 sentence pairs per LP. Control LDSPs are randomly chosen from the dataset, intended to be unrelated, as a baseline for our analysis.

- 147
- 148
- 149
- 150

3 LINGUISTICALLY DISTINCT SENTENCE PAIRS (LDSP-10) DATASET

151 We curated a dataset of 1000 LDSPs for each of the 10 LPs we wanted to investigate. The dataset 152 was generated using Google's gemini-1.5-flash model API. This model was selected due to its reliability and cost-efficiency while being able to produce consistent outputs across a variety of 153 linguistic contexts. The model was prompted with a set of reference LDSPs as well as a description 154 of the LP to ensure a high-quality outputs. These outputs were generated in batches of 100 LDSPs at 155 a time. To ensure reproducibility and transparency, the detailed prompts used to generate the dataset 156 are provided in Appendix A. These prompts included explicit examples of each LP, along with clear 157 instructions tailored to the gemini-1.5-flash API to encourage outputs adhering to the desired 158 properties. 159

During the dataset creation process, the order of the sentences in the LDSP was not always consistent
 with the intended property distinction. We made modifications to the prompt to explicitly enforce
 the correct ordering. This adjustment ensured that the generated outputs reliably aligned with our

expectations. Manual validation was conducted to assess the quality of the generated data. The
 evaluation revealed that more than 99% of the sampled sentence pairs adhered to the minimal
 distinctions expected for their LP. The system exhibited a low rate of syntactic or content biases,
 with errors occurring primarily in cases involving more complex distinctions, such as polarity and
 factuality.

The LPs tested were chosen to explore various semantic and syntactic relationships. We generated LDSPs for *definiteness*, *factuality*, *intensifier*, *negation*, *polarity*, *quantity*, *synonym*, and *tense*. In addition, we generated a *control* group, which contains sentence pairs of completely unrelated sentences. This is used to compare to the LDSPs and contextualize our observed results. Example LDSPs can be found in Table 2, with more detailed definitions found in Appendix B. For more information about the dataset generation pipeline, please refer to Appendix A.

173 174

175 176

177

190

197

206

4 DIMENSION-WISE EMBEDDING ANALYSIS

4.1 WILCOXON SIGNED-RANK TEST

The Wilcoxon signed-rank test is employed in our analysis to assess whether there exists a significant difference in embedding dimensions across paired sentence representations. This non-parametric test is particularly useful when the data does not conform to the normality assumptions required by parametric tests such as the paired t-test. Given that sentence embeddings often exhibit complex, non-Gaussian distributions, the Wilcoxon test provides a robust approach to evaluating the statistical significance of differences in embedding dimensions.

Formally, let $X_1, X_2 \in \mathbb{R}^d$ be the embedding representations of two paired sentences. We define the difference vector as:

$$D = X_1 - X_2, (1)$$

where $D = \{d_1, d_2, ..., d_d\}$ contains the differences for each embedding dimension. The null hypothesis for the Wilcoxon test is given by:

 $H_0: \operatorname{median}(D) = 0, \tag{2}$

which posits that there is no significant shift in the embedding dimensions between the two sentence representations.

The test proceeds by ranking the absolute values of the nonzero differences, assigning ranks R_i to each $|d_i|$. The Wilcoxon test statistic W is computed as the sum of ranks corresponding to positive differences:

$$W = \sum_{d_i > 0} R_i.$$
(3)

The significance of W is then assessed using either critical values from the Wilcoxon distribution or by computing a p-value.

We employ the Wilcoxon test in our framework to analyze whether certain dimensions of the
 embeddings exhibit systematic shifts between sentence pairs. Overall, the Wilcoxon signed-rank
 test provides a rigorous statistical method for validating the role of embedding dimensions in
 differentiating sentence pairs, ensuring that our conclusions are drawn from statistically significant
 evidence rather than random variations.

207 4.2 MUTUAL INFORMATION (MI)

To further investigate the relationship between embedding dimensions and each LP and inspired
by Pimentel et al. (2020), we employed mutual information (MI) analysis. Mutual information is a
measure of the mutual dependence between two variables, quantifying the amount of information
obtained about one variable by observing the other Zeng (2015).

For discrete random variables X and Y, the mutual information MI(X;Y) is defined as:

214
215
$$MI(X;Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_{XY}(x,y) \log \frac{P_{XY}(x,y)}{P_X(x)P_Y(y)},$$

where $P_{XY}(x, y)$ is the joint probability distribution of X and Y, and $P_X(x)$ and $P_Y(y)$ are the marginal probability distributions of X and Y, respectively. In our context:

- *X* represents the values of a particular embedding dimension.
- Y represents S_1 (0) or S_2 (1).

To apply mutual information analysis, we discretize the embedding dimensions using quantile-based binning with 10 bins. This number was selected as a balance between the preservation of information content and the avoidance of excessive complexity in the estimation of the *MI* score and is a common practice in similar analyses Steuer et al. (2002).

226 227

242

248

256

257 258

259

219

220

221

4.3 **RECURSIVE FEATURE ELIMINATION**

228 We initially examined each embedding dimension's predictive capability with simple logistic re-229 gression. Unlike more flexible techniques, logistic regression imposes a linear decision boundary, 230 which was unable to capture the complex patterns defining most linguistic contrasts within the 231 generated embeddings. To capture these relationships, we applied Recursive Feature Elimination 232 (RFE) using scikit-learn's implementation with logistic regression as the base estimator Zeng 233 et al. (2009). Embedding pairs were split into their constituent parts, with sentence1 embeddings 234 labeled as class 0 and sentence2 embeddings as class 1, enabling a binary classification setup to 235 highlight dimensions that distinguish the two positions. The RFE procedure iteratively trained a model, assigned importance weights to features, and removed the least important ones until the top 236 20 features remained. 237

The dataset was divided into training (80%) and testing (20%) sets with a fixed random seed to ensure consistency. RFE was initialized with a logistic regression classifier (max 1000 iterations), and the selected 20 features were used to train a final logistic regression model. The model's performance was evaluated on the test set using accuracy as the metric.

243 4.4 EDI SCORE CALCULATION

To quantify the contribution of of each embedding dimension to a LP, we introduce the Embedding
Dimension Importance (EDI) Score, which is computed for each dimension *d* and each LP *lp* as
follows:

$$\mathrm{EDI}_{d,lp} = w_1 \cdot -\log p_{d,lp} + w_2 \cdot M_{d,lp} + w_3 \cdot R_{d,lp}$$

where $p_{d,lp}$ is the *p*-value obtained from the Wilcoxon signed-rank test results. $M_{d,lp}$ is the mutual information score. $R_{d,lp}$ is the absolute value of the logistic regression weights after the recursive feature elimination if *d* remains in the reduced feature set for LP lp; otherwise, $R_{d,lp} = 0$. $p_{d,lp}$, $M_{d,lp}$, $R_{d,lp}$ are min-max scaled before the EDI score weighted to calculation to enforce EDI scores to be $\in [0, 1]$. Lastly, $w_1 = 0.6$, $w_2 = 0.2$, and $w_3 = 0.2$. Wilcoxon's test was weighted the most heavily, as it calculates the statistical significance of the differences observed, which our testing showed was a strong predictor of dimension importance.

5 EVALUATION

5.1 LINGUISTIC PROPERTY CLASSIFIER

To verify the feasibility of using sentence pairs, we calculated embedding difference vectors $D_i = emb(S_{1i}) - emb(S_{2i})$ and evaluated them as predictors of LP. To this end, we trained an LP classifier that assigns any given embedding difference vector to one of the tested LPs. The primary goal of this classifier is to assess how well different LPs can be separated in the embedding space. The model was trained using an 80-20 training-test split on the entire LDSP-10 dataset.

266 5.2 EDI SCORE EVALUATION

267

265

To systematically assess the effectiveness of EDI scores, we implement a structured evaluation
 framework consisting of a baseline test and three evaluations experiments. For more details on the
 algorithms for each evaluation method, refer to Appendix C.

283

284

285

287 288

270 271



(a) Distribution of dimension 0 of *control* LDSPs for S_1 and S_2 . For *control*, all dimensions had equivalent Wilcoxon *p*-values, so dimension 0 represents the most and least significant *p*-value.



(b) Distribution of dimensions 544 (top) and 489 (bottom), lowest and highest *p*-values respectively, of *negation* LDSPs for S_1 and S_2 . There is a discernible shift to the right in dimension 544, for sentences that are negated.



(c) Distribution of dimensions 445 (top) and 489 (bottom), lowest and highest *p*-values respectively, of *intensifier* LDSPs for S_1 and S_2 . Intensified sentences have values in dimension 445 that tend to be lower, as seen by the distributional shift to the left.

Figure 2: BERT embedding distributions for control, negation, and intensifier.

For the baseline, we train a logistic regression classifier on the full set of embedding dimensions.
Given a binary classification task for each LP, the classifier is trained to distinguish between the two sentences in the LDSP using all available embedding dimensions, serving as an upper bound against which subsequent evaluations are compared.

Evaluation 1 explores how dimensions with *high* EDI scores replicate the performance of the full-dimensional classifier. We first rank all dimensions by their EDI score in descending order. Starting with the highest-ranked dimension, we train a logistic regression classifier, as in the baseline evaluation, but only with this single feature. We iteratively add the next highest-ranked dimension, retraining the classifier and evaluating the test accuracy until we reach at least 95% of the baseline accuracy.

Evaluation 2 verifies that dimensions with *low* EDI scores do not encode information relevant to the LP. We identify the 100 lowest-ranked dimensions and train a logistic classifier to distinguish between the two sentences using only those dimensions. We record the accuracy on a test dataset to determine whether it remains close to random chance, as expected, to ensure that these dimensions lack significance in encoding the LP.

Evaluation 3 examines cross-property generalization, exploring whether high-EDI-score dimensions for one LP are specialized rather than broadly informative across different properties. We use the highest-ranked EDI score dimensions of *other* properties to predict the current property. We expect the performance of this classifier to be generally lower than the baseline and the high EDI Score accuracy.

6 RESULTS

311 312

310

In this section, we focus on BERT embeddings as a case study for applying our framework. We focus on showing visualizations for *control*, *negation*, and *intensifier*, but all other LPs and related tables/plots can be found in Appendix . The results for GPT-2 and MPNet were similar, and can be reviewed in detail in Appendix E and Appendix F.

317

318 6.1 CONTROL AND SYNONYM319

The *control* LDSPs consists of completely unrelated sentence pairs. As expected, the results show that there are no significant dimensions in BERT embeddings that encode any relationships. Figure 4 illustrates very little agreement the Wilcoxon signed-rank test, RFE, and mutual information. The Wilcoxon test *p*-values show no dimensions with significant differences in their means, as shown in Figure 2a. The maximum EDI score of 0.3683 is the lowest of all other properties. The embeddings



Figure 3: Combined analysis graphs for *negation* and *intensifier*. Circled bars represent dimensions that all three tests agree to be highly important. Similar to Figure 4.

of the two sentences are expected to be far in embedding space because of their unrelated nature, which aligns with these observed results.

Despite having sentences that were very close or equivalent in meaning, the results of the analysis for the *synonym* LDSPs were very close to the completely unrelated sentences of *control*. The Wilcoxon test shows no significant dimensions that encode meaningful differences between the sentences. The maximum EDI score of 0.8751 is followed by a steep drop-off.

332

333

334 335 336

337

338

339

340

6.2 NEGATION AND POLARITY

345 The negation LDSPs showed very strong results, 346 with 13 dimensions with an EDI score of 0.8347 or above. The maximum EDI score of 0.9987 348 for dimensions 544 is one of the strongest out 349 of any LP. Figure 3 illustrates this, with high agreement between the Wilcoxon signed-rank 350 test, RFE, and mutual information test results. 351 Figure 2b highlights the distributional shift in 352 some dimensions, which compared to the con-353 trol highlights a discernible, binary relationship 354 in the data. 355

Polarity is very similar to negation and had simi-356 larly strong results. With a maximum EDI score 357 of 0.9977 for dimension 431, and over 20 dimen-358 sions with EDI scores over 0.8, it was also one 359 of the strongest relationships that we observed. 360 The singular switch to an antonym in the sen-361 tence completely reverses the meaning of the 362 sentence, explaining the strong binary relation-363 ship between the sentences. 364

- 365 366
- 6.3 INTENSIFIER

Adding a word to increase the emphasis of a
verb changes the meaning of the sentence to a
lesser degree than a complete reversal, so the



Figure 4: Combined analysis graph for *control*: shows the top 25 important dimensions selected by each of the three methods in § 4. Bar height represents mutual information (MI); bars above the dashed line are in the top 25 MI scores. Blue bars signify the lowest Wilcoxon test *p*-values. Green triangles indicate a dimension that was selected by recursive feature elimination (RFE) with num_features set to 25. In the case for *control*, all dimensions had equivalent Wilcoxon *p*-values, so the first 25 are selected.

results of the *intensifier* LDSPs reflect a slightly weaker relationship than *negation*. There are fewer dimensions with multiple test agreement, as shown in Figure 3, as well as a slighter distributional shift, as shown by the most significant p-value Wilcoxon test results (Figure 2c). With a maximum EDI score of 0.8911, the encoding is relatively weaker, but noticeable.

374

376

375 6.4

6.4 OTHER LINGUISTIC PROPERTIES

377 Largely syntactical changes, such as those observed in *definiteness*, led to strong EDI scores as well. *Definiteness* had the highest dimensional EDI score, with dimension 180 receiving a score of 1.0. A simple switch from a definite to an indefinite article is a distinct change in structure. As articles are present in most English sentences, a singular dimension with a perfect EDI score is expected.

Voice, another syntactical property, had pairs of sentences with shuffled word orders and verb changes.
 The results show that this is encoded in relatively few dimensions, with only 3 dimensions scoring above 0.9.

The *quantity* LDSPs involve changes in the syntax and semantics. Similar to the *intensifier* results, the EDI scores at large were relatively lower for these properties, but still much stronger than the *control*.

Tense represented a large semantic change, as well as a structural one in the conjugation of verbs.
 Although the maximum EDI score of 0.9405 was not as high as other properties, 18 embeddings scored above 0.8, indicating an encoding of this property over many dimensions.

For more details and visualizations of all properties, refer to Appendix D.

393 394 395

390

391 392

6.5 EVALUATION RESULTS

396 397 398

The LP classifier achieved a test accuracy of 399 0.863 with a confusion matrix as shown in 400 Figure 5, demonstrating that the embedding 401 difference vectors contain sufficient separable 402 information to distinguish between different 403 LPs. Moreover, the strong performance of the 404 classifier supports the validity of our pairwise 405 minimal-perturbation approach, indicating that 406 small controlled changes in sentence pairs ef-407 fectively capture linguistic distinctions in the embedding space. 408

409 In the high EDI score evaluation, we observed 410 that across most LPs, only less than 12 of the 411 highest-ranked dimensions were required to re-412 cover at least 95% of the baseline classifier's accuracy, with some properties (i.e. *factuality*) 413 requiring as few as four dimensions. This indi-414 cates that the information necessary for classify-415 ing each LP is concentrated in a relatively small 416 subset of embedding dimensions. Conversely, 417 the low EDI score evaluation confirmed that di-418 mensions with low scores contribute minimally 419 to classification performance. Even when using 420 the 100 lowest-ranked dimensions, the resulting



Figure 5: Confusion matrix for the LP classifier (§ 5.1). All LPs, except *control* and *synonym*, are accurately classified by the model. *Control*'s randomness ensures that its different vectors contain no consistent separability, similarly with *synonym*'s unordered pairings.

classifier performed consistently worse than classifiers using much fewer (4-38) of the highest-ranked
dimensions (Figures 6a, 6b). This demonstrates the EDI score's validity as a measure of whether a
given dimension encodes information relevant to an LP.

424 Finally, the cross-property evaluation demonstrated that using the top-ranked dimensions from another 425 LP generally resulted in lower classification performance compared to using the high-EDI dimensions 426 of the target property, showing that the EDI score effectively identifies dimensions that encode 427 information specific to each LP. Interestingly, we found that certain properties with conceptual 428 similarities performed best for each other. For example, in the polarity classification task, the top 429 EDI dimensions from negation achieved the highest accuracy among all cross-property evaluations, reaching 0.895 (Figure 6a). This result aligns with the intuition that negative sentiment—typically 430 represented by the second sentence in polarity pairs—is often expressed through negation, reinforcing 431 the semantic connection between these LPs.



(a) Eight dimensions were enough to achieve nearbaseline accuracy. The top-performing cross property is *negation*, which contains semantic similarities to polarity.



(b) Incrementally added 19 high-EDI dimensions until the classifier reached near-baseline performance. Low-EDI performance (red dashed line) was nearly half.

Figure 6: Evaluation plots for *polarity* and *intensifier*. Green dashed line marks the baseline performance threshold, the grey dashed line is the performance of the top EDI dimensions on *control*, and the red dashed line is the performance of the lowest 100 EDI dimensions. The blue line tracks the test accuracy of the classifier as we increased the number of top EDI-scored dimensions.

7 DISCUSSION

The results of this study provide a clear demonstration of the ability to disentangle specific LPs 452 within high-dimensional embeddings. Our analysis shows that certain LPs are robustly encoded 453 in distinct embedding dimensions, as evidenced by high Embedding Dimension Importance (EDI) 454 scores and agreement across multiple analytical methods. These methods were chosen after rigorous 455 experimentation, where principal component analysis, simple logistic regression, and other methods 456 were rejected due to their inability to capture the nuanced, non-linear information encoded in these 457 embeddings. Negation yielded one of the the highest maximum EDI scores and a significant number 458 of dimensions with high interpretability. This supports the notion that negation is a well-structured 459 and salient linguistic feature in BERT embeddings.

In contrast, some properties exhibited minimal evidence of dimension-specific encoding, which
we hypothesize to be due to a lack of a binary or clear-cut way of encoding these relationships.
Synonymy showed low maximum EDI scores and inconsistent results across the Wilcoxon SignedRank Test, Mutual Information, and Recursive Feature Elimination. Synonym pairs in our dataset
could be permuted without affecting the consistency of the data, and 0-1 labels for our classifiers and
mutual information were meaningless; therefore, our methods are unable to extract the dimensional
distribution of synonym encodings.

In summary, this study underscores the heterogeneous nature of linguistic encoding in BERT embed dings, with some properties exhibiting clear, interpretable patterns while others remain elusive. The
 proposed EDI score and analytical framework provide valuable tools for advancing the interpretability
 of embeddings, with implications for bias mitigation, model optimization, and the broader goal of
 responsible AI deployment.

472 473

444

445

446

447 448 449

450 451

8 LIMITATIONS

474 475

While our study provides insight into the interpretability of embedding dimensions, it is constrained
primarily due to data availability. Generating high-quality LDSPs with LLM-based tools is difficult,
as ensuring diversity, minimal redundancy, and high linguistic quality becomes significantly more
difficult with more data generated. Overly simplistic, repetitive outputs are difficult to avoid, despite
careful prompt engineering.

481

482 REFERENCES

Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. Fine-grained analysis
 of sentence embeddings using auxiliary prediction tasks, 2017. URL https://arxiv.org/abs/
 1608.04207.

505

506

507

508

509

514

524

525

526

486
 487
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488
 488

- Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass.
 Identifying and controlling important neurons in neural machine translation. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=
 H1z-PsR5KX.
- Yonatan Belinkov and James Glass. Analysis methods in neural language processing: A survey.
 Transactions of the Association for Computational Linguistics, 7:49–72, 2019. doi: 10.1162/tacl_
 a_00254. URL https://aclanthology.org/Q19-1004.
- Vinamra Benara, Chandan Singh, John X. Morris, Richard Antonello, Ion Stoica, Alexander G.
 Huth, and Jianfeng Gao. Crafting interpretable embeddings by asking llms questions, 2024. URL
 https://arxiv.org/abs/2405.16714.
- Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings, 2016. URL https://arxiv.org/abs/1607.06520.
 - Hongliu Cao. Recent advances in text embedding: A comprehensive review of top-performing methods on the mteb benchmark, 2024. URL https://arxiv.org/abs/2406.01607.
 - Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. What you can cram into a single vector: Probing sentence embeddings for linguistic properties, 2018. URL https://arxiv.org/abs/1805.01070.
- Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. What
 is one grain of sand in the desert? analyzing individual neurons in deep nlp models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):6309–6317, Jul. 2019. doi: 10.1609/aaai.
 v33i01.33016309. URL https://ojs.aaai.org/index.php/AAAI/article/view/4592.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL https://arxiv.org/abs/ 1810.04805.
- Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and Yonatan Belinkov. Analyzing individual neurons in pre-trained language models. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing* (*EMNLP*), pp. 4865–4880, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.395. URL https://aclanthology.org/2020.emnlp-main. 395/.
 - Nadir Durrani, Fahim Dalvi, and Hassan Sajjad. Discovering salient neurons in deep nlp models, 2024. URL https://arxiv.org/abs/2206.13288.
- Manaal Faruqui, Jesse Dodge, Sujay Kumar Jauhar, Chris Dyer, Eduard Hovy, and Noah A. Smith.
 Retrofitting word vectors to semantic lexicons. In Rada Mihalcea, Joyce Chai, and Anoop Sarkar
 (eds.), Proceedings of the 2015 Conference of the North American Chapter of the Association for
 Computational Linguistics: Human Language Technologies, pp. 1606–1615, Denver, Colorado,
 May–June 2015a. Association for Computational Linguistics. doi: 10.3115/v1/N15-1184. URL
 https://aclanthology.org/N15-1184.
- Manaal Faruqui, Yulia Tsvetkov, Dani Yogatama, Chris Dyer, and Noah Smith. Sparse overcomplete
 word vector representations, 2015b. URL https://arxiv.org/abs/1506.02004.
- Simon Guillot, Thibault Prouteau, and Nicolas Dugue. Sparser is better: one step closer to word embedding interpretability. In Maxime Amblard and Ellen Breitholtz (eds.), *Proceedings of the 15th International Conference on Computational Semantics*, pp. 106–115, Nancy, France, June 2023. Association for Computational Linguistics. URL https://aclanthology.org/2023. iwcs-1.13.

540 541 542	Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. Finding neurons in a haystack: Case studies with sparse probing, 2023. URL https://arxiv. org/abs/2305.01610.
543 544 545	Isabelle Guyon and André Elisseeff. An introduction to variable and feature selection. J. Mach. Learn. Res., 3(null):1157–1182, March 2003. ISSN 1532-4435.
546 547 548	Francesca Incitti, Federico Urli, and Lauro Snidaro. Beyond word embeddings: A survey. <i>Information Fusion</i> , 89:418–436, 2023. ISSN 1566-2535. doi: https://doi.org/10.1016/j.inffus.2022.08.024. URL https://www.sciencedirect.com/science/article/pii/S1566253522001233.
549 550 551	Zachary C. Lipton. The mythos of model interpretability, 2017. URL https://arxiv.org/abs/ 1606.03490.
552 553 554	Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. <i>ACM Comput. Surv.</i> , 54(6), July 2021. ISSN 0360-0300. doi: 10.1145/3457607. URL https://doi.org/10.1145/3457607.
555 556	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representa- tions in vector space, 2013. URL https://arxiv.org/abs/1301.3781.
557 558 559 560 561	Brian Murphy, Partha Pratim Talukdar, and Tom Michael Mitchell. Learning effective and inter- pretable semantic models using non-negative sparse embedding. In <i>International Conference</i> <i>on Computational Linguistics</i> , 2012. URL https://api.semanticscholar.org/CorpusID: 8348149.
562 563 564 565 566	Abhishek Panigrahi, Harsha Vardhan Simhadri, and Chiranjib Bhattacharyya. Word2Sense: Sparse interpretable word embeddings. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pp. 5692–5705, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1570. URL https://aclanthology.org/P19-1570.
567 568 569 570 571	Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1162. URL https://aclanthology.org/D14-1162.
572 573 574 575 576 577	Tiago Pimentel, Josef Valvoda, Rowan Hall Maudslay, Ran Zmigrod, Adina Williams, and Ryan Cotterell. Information-theoretic probing for linguistic structure. In Dan Jurafsky, Joyce Chai, Na- talie Schluter, and Joel Tetreault (eds.), <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pp. 4609–4622, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.420. URL https://aclanthology.org/2020.acl-main.420/.
578 579 580	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. URL https://api.semanticscholar.org/CorpusID:160025533.
581 582 583 584	Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in bertology: What we know about how bert works. <i>Transactions of the Association for Computational Linguistics</i> , 8:842–866, 12 2020. doi: 10.1162/tacl_a_00349.
585 586 587	Hassan Sajjad, Nadir Durrani, and Fahim Dalvi. Neuron-level interpretation of deep nlp models: A survey. <i>Transactions of the Association for Computational Linguistics</i> , 10:1285–1303, 11 2022. ISSN 2307-387X. doi: 10.1162/tacl_a_00519. URL https://doi.org/10.1162/tacl_a_00519.
588 589 590 591	Lutfi Senel, Ihsan Utlu, Veysel Yucesoy, Aykut Koc, and Tolga Cukur. Semantic structure and interpretability of word embeddings. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , PP, 11 2017. doi: 10.1109/TASLP.2018.2837384.
592 593	Lauro Snidaro, Giovanni Ferrin, and Gian Luca Foresti. Distributional memory explainable word embeddings in continuous space. In 2019 22th International Conference on Information Fusion (FUSION), pp. 1–7, 2019. doi: 10.23919/FUSION43075.2019.9011324.

- 594 R Steuer, Juergen Kurths, Carsten Daub, Janko Weise, and J Selbig. The mutual information: 595 Detecting and evaluating dependencies between variables. *Bioinformatics (Oxford, England)*, 18 596 Suppl 2:S231–40, 02 2002. doi: 10.1093/bioinformatics/18.suppl_2.S231. 597 Anant Subramanian, Danish Pruthi, Harsh Jhamtani, Taylor Berg-Kirkpatrick, and Eduard Hovy. 598 Spine: sparse interpretable neural embeddings. In Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence 600 Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence, 601 AAAI'18/IAAI'18/EAAI'18. AAAI Press, 2018. ISBN 978-1-57735-800-8. 602 603 Lucas Torroba Hennigen, Adina Williams, and Ryan Cotterell. Intrinsic probing through dimension 604 selection. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), Proceedings of the 605 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 197–216, 606 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020. 607 emnlp-main.15. URL https://aclanthology.org/2020.emnlp-main.15/. 608 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-609 attention: Specialized heads do the heavy lifting, the rest can be pruned. In Anna Korhonen, David 610 Traum, and Lluís Màrquez (eds.), Proceedings of the 57th Annual Meeting of the Association for 611 Computational Linguistics, pp. 5797–5808, Florence, Italy, July 2019. Association for Computa-612 tional Linguistics. doi: 10.18653/v1/P19-1580. URL https://aclanthology.org/P19-1580/. 613 614 Guoping Zeng. A unified definition of mutual information with applications in machine learning. 615 Mathematical Problems in Engineering, 2015, 03 2015. doi: 10.1155/2015/201874. 616 Xiangyan Zeng, Yen-Wei Chen, and Caixia Tao. Feature selection using recursive feature elimination 617 for handwritten digit recognition. In 2009 Fifth International Conference on Intelligent Information 618 Hiding and Multimedia Signal Processing, pp. 1205–1208, 2009. doi: 10.1109/IIH-MSP.2009.145. 619 620 Henghui Zhu, Ioannis Ch. Paschalidis, and Amir Tahmasebi. Clinical concept extraction with 621 contextual word embedding, 2018. URL https://arxiv.org/abs/1810.10566. 622 623 624 DATASET GENERATION PIPELINE А 625 626 Figure 7 illustrates the procedure used to generate the LDSP-10 dataset. The batch procedure of 627 generating 100 pairs of sentences at a time was crucial in minimizing API costs while also getting 628 high-quality generations that would be useful for our experiments. The prompt template used can be 629 seen in Figure 8. 630 631 LINGUISTIC PROPERTY DEFINITIONS В 632 633 We tested LDSPs for the following linguistic properties: 634 635 • Definiteness involves the use of definite or indefinite articles within a sentence, such as the 636 compared to *a*, respectively. 637 • *Factuality* refers to the degree of truth implied by the structure of the sentence. 638 639 • Intensifier refers to the degree of emphasis present within a sentence. 640 • *Negation* occurs when a *not* is added to a sentence, negating the meaning. 641 • *Polarity* this is similar to a negation, and occurs when an antonym is added, reversing the 642 meaning of the sentence completely. 643
 - *Quantity* a switch from an exact number used to numerate the items to a grouping word.
 - *Synonym* both sentences have the same meaning, with one word being replaced by one of its synonyms.

646

647

• Tense one sentence is constructed in the present tense, while the other is in the past tense.



702	
703	
704	
705	
706	
707	
708	
709	
710	
711	
712	
713	
714	
715	<pre>prompt_template = """</pre>
716	
717	You are generating a dataset of Linguistically Distinct Sentence Pairs (LDSPs).
718	overall meaning.
719	
720	Below are some examples of LDSPs
721	Linguistic Drementus regation
722	LInguistic Property: negation
723	Ebol. (The box 10 on the counter , The box 10 hot on the counter)
724	Linguistic Property: tense
725	LDSP: ('The box is on the counter', 'The box was on the counter')
720	You will generate {num ldsps} distinct LDSPs of various topics 100 at a time
728	
720	You will generate them as two columns of a CSV. One column for first sentence of
730	the LDSP, and the other column for the second.
731	Each row is a new LDSP, so you will generate {num_ldsps} rows in total.
732	Generate no other text. Vary the sentence structure.
733	
734	The property for which you will be generating LDSPs will be {linguistic_property}.
735	Property Description: {property description}
736	Toperty bescription. {property_description}
737	An example LDSP for this property is
738	{example_ldsp}
739	Concrete the first 100 LDSPc
740	Generale LIE IIIST IVO LUSES.
741	ппп
742	
743	Earner 0. The promotetemplate used to concrete I DOD with the maximist 1. F. flood we do ADI
744	Figure 6: The prompt temptate used to generate LDSPs with the gemini-1.5-flash model API.

⁷⁵⁶ C EVALUATION ALGORITHMS

758

759

760

761

To systematically assess the efficacy of EDI (Embedding Dimension Importance) scores, we conduct a structured evaluation using logistic regression classifiers. Our evaluation consists of three key evaluation algorithms:

762 Algorithm 1 Evaluation 1: High EDI Score 763 **Require:** Ranked dimensions $D = \{d_1, d_2, ..., d_{768}\}$ sorted by descending EDI score 764 **Ensure:** Accuracy curve A_k as a function of dimensions used 765 1: Initialize $k \leftarrow 1, A_k \leftarrow 0$ 766 2: while $A_k < 0.95 A_{\text{baseline}}$ do 767 3: Select top k dimensions: $X_k = X[:, D_{1:k}]$ 768 4: Train logistic regression on X_k 769 5: Compute test accuracy $A_k \leftarrow \text{Evaluate}(\theta, X_{\text{test}}, y_{\text{test}})$ 770 6: $k \leftarrow k+1$ 7: end while 771 8: return A_k 772 773 774 775 Algorithm 2 Evaluation 2: Low EDI Score 776 **Require:** Ranked dimensions $D = \{d_1, d_2, ..., d_{768}\}$ sorted by ascending EDI score 777 **Ensure:** Test accuracy A_{low} using lowest-EDI dimensions 778 1: Select bottom k = 100 dimensions: $X_{low} = X[:, D_{1:100}]$ 779 2: Train logistic regression on X_{low} 780 3: Compute test accuracy $A_{\text{low}} \leftarrow \text{Evaluate}(\theta, X_{\text{test}}, y_{\text{test}})$ 781 4: return A_{low} 782 783 784 Algorithm 3 Evaluation 3: Cross-Property 785 **Require:** Current property P_0 dataset (X, y), set of other properties $\mathcal{P} = \{P_1, P_2, ..., P_9\}$, where 786 each P_i has ranked EDI dimensions D_{P_i} 787 **Ensure:** Accuracy scores $\{A_{P_1}, A_{P_2}, ..., A_{P_q}\}$ 788 1: for each property $P \in \mathcal{P}$ do 789 Retrieve top k = 25 dimensions from P: $D_P^{1:25}$ 2: 790 Extract these dimensions from current data: $X_{\text{train}}^P = X_{\text{train}}[:, D_P^{1:25}]$ 3: 791 4: Train logistic regression on X_{train}^P 792 Compute test accuracy $A_P \leftarrow \text{Evaluate}(\theta, X_{\text{test}}^P, y_{\text{test}})$ 5: 793 6: end for 794 7: return $\{A_P\}_{P \in \mathcal{P}}$ 796

These evaluations provide a comprehensive understanding of how EDI scores relate to classification accuracy, ensuring that high EDI dimensions contain useful linguistic information while low EDI dimensions do not. The cross-property evaluation further confirms that high-EDI dimensions are specialized rather than general indicators of LPs.

D ADDITIONAL LINGUISTIC PROPERTY RESULTS FOR BERT EMBEDDINGS

D.1 CONTROL

797

798

799

800

801 802

803 804

805

Table 3 highlights the top 10 EDI scores for the *control*. The baseline evaluation results for *control* showed an accuracy of 0.5200, close to random chance. The Low EDI score test yielded an accuracy of 0.4575. The High EDI score test demonstrated quick improvements, achieving 95% of baseline accuracy with a single dimension, as the baseline accuracy was low, as illustrated in Figure 9. The greatest cross-property accuracy was achieved by *voice*, at 0.5325.



Factuality had strong results. Figure 13 highlights the stark difference between the most prominent dimensions encoding this property. Table 5 highlights the top 10 EDI scores, while Figure 15 illustrates the high level of agreement between our various tests.

Dimension	EDI Score
180	1.0000
123	0.8824
319	0.8819
385	0.8639
109	0.8155
497	0.7974
683	0.7948
172	0.7926
430	0.7907
286	0.7862

Table 4: Top 10 BERT EDI scores for Definiteness.



Figure 10: BERT Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Definiteness*.



Figure 11: High EDI score evaluation results for BERT Embeddings of definiteness.

The baseline evaluation results for *factuality* showed an accuracy of 0.9975. The Low EDI score test yielded an accuracy of 0.5975, approximately random. The High EDI score test demonstrated very quick improvements, achieving 95% of baseline accuracy with 4 dimensions, as illustrated in Figure 14. The greatest cross-property accuracy was achieved by *tense*, at 0.9650.

D.4 INTENSIFIER

Table 6 highlights the top 10 EDI scores for *intensifier*. The baseline evaluation results for *intensifier* showed an accuracy of 0.9925. The Low EDI score test yielded an accuracy of 0.5150, close to random chance. The High EDI score test demonstrated incremental improvements, achieving 95%

Dimension	EDI Score
577	0.9740
43	0.9386
210	0.9249
745	0.8954
539	0.8887
387	0.8869
60	0.8727
16	0.8617
54	0.8609
97	0.8538

Table 5: Top 10 BERT EDI scores for Factuality.



Figure 12: BERT Mutual Information of Embedding Dimensions overlaid with Wilcoxon test and RFE results for *Definiteness*



Figure 13: BERT Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Factuality*.

of baseline accuracy with 19 dimensions, as illustrated in Figure 6b. The greatest cross-property accuracy was achieved by *quantity*, at 0.8550.

D.5 NEGATION

Table 7 highlights the top 10 EDI scores for *negation*. The baseline evaluation results for *negation* showed an accuracy of 0.9925. The Low EDI score test yielded an accuracy of 0.5800, close to random chance. The High EDI score test demonstrated incremental improvements, achieving 95% of baseline accuracy with 11 dimensions, as illustrated in Figure 16. The greatest cross-property accuracy was achieved by *tense*, at 0.9100.

Γ	Dimension	EDI Score
F	686	0.8911
Ē	663	0.8832
Ī	139	0.8805
F	605	0.8790
ľ	269	0.8650
ľ	441	0.8612
F	144	0.8535
	692	0.8468
ſ	445	0.8385
F	442	0.8221

Table 6: Top 10 BERT EDI scores for Intensifier.



Figure 14: High EDI score evaluation results for BERT Embeddings of *factuality*.



Figure 15: Mutual Information of Embedding Dimensions overlaid with Wilcoxon test and RFE results for *Factuality*

D.6 POLARITY

Polarity, as it is similar to negation, had extremely strong results. Figure 17 highlights the differences
 between the most prominent dimensions encoding this property. Table 8 highlights the top 10 EDI
 scores, while Figure 18 illustrates the extremely high level of agreement between our various tests.

The baseline evaluation results for *polarity* showed an accuracy of 0.9775. The Low EDI score test yielded an accuracy of 0.5575, close to random chance. The High EDI score test demonstrated incremental improvements, achieving 95% of baseline accuracy with 8 dimensions, as illustrated in Figure 6a. The greatest cross-property accuracy was achieved by *negation*, at 0.8950.

1009 D.7 QUANTITY

1010
 1011 *Quantity* had more moderate results compared to *polarity* and *negation*. Figure 19 highlights the difference between the most prominent dimensions encoding this property. Table 9 highlights the top 10 EDI scores, while Figure 21 illustrates the moderate level of agreement the tests.

1	01	4
1	01	5
1	01	6

544 251	0.9987
251	0.0277
	0.9211
171	0.9236
451	0.9101
737	0.8891
281	0.8812
96	0.8624
692	0.8512
85	0.8501
642	0.8461
	171 451 737 281 96 692 85 642

Table 7: Top 10 BERT EDI scores for Negation.





Figure 17: BERT Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Polarity*.



Figure 18: Mutual Information of BERT Embedding Dimensions overlaid with Wilcoxon test and RFE results for Polarity

D.9 TENSE

Tense had moderate results. Figure 24 highlights the differences between the most prominent dimensions encoding this property. Table 11 highlights the top 10 EDI scores, while Figure 27 illustrates the level of agreement the tests.

The baseline evaluation results for tense showed an accuracy of 0.9975. The Low EDI score test yielded an accuracy of 0.4625, close to random chance. The High EDI score test demonstrated incremental improvements, achieving 95% of baseline accuracy with 11 dimensions, as illustrated in Figure 25. The greatest cross-property accuracy was achieved by *control*, at 0.9150.

1	1	21
1	1	22
1	1	23
1	1	24

D	Dimension	EDI Score
	676	0.8751
	203	0.7744
	701	0.6916
	654	0.6897
	463	0.6889
	544	0.6602
	91	0.6598
	437	0.6557
	446	0.6543
	487	0.6415

Table 10: Top 10 BERT EDI scores for Synonym.





Figure 22: BERT Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Synonym*.

 Ε

GPT-2

This section will contain the visualizations of the results for GPT-2 embeddings. Full detailed results, including full EDI scores as well as additional visualization, will be available on GitHub upon publication.

E.1 LINGUISTIC PROPERTY CLASSIFIER

The results from the Linguistic Property Classifier for GPT-2 embeddings is shown in Figure 30.

 Dimension **EDI Score** 0.9722 0.9552 0.9376 0.8875 0.8783 0.8586 0.8437 0.8404 0.8182 0.8113

Table 12: Top 10 BERT EDI scores for Voice.



Figure 24: BERT Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Tense*.

1273 E.2 CONTROL

Figure 31 highlights the difference between the most prominent dimensions encoding this property.
 Figure 33 illustrates the level of agreement between the tests.

The baseline evaluation results for *control* showed an accuracy of 0.4725, close to chance. The Low EDI score test yielded an accuracy of 0.4400. The High EDI score test demonstrated strong performance, achieving 95% of baseline accuracy with just a single dimension, as the baseline accuracy was close to random chance, as illustrated in Figure 32. The highest cross-property accuracy was achieved by *voice*, at 0.5450.

1283 E.3 DEFINITENESS

Figure 34 highlights the difference between the most prominent dimensions encoding this property.
Figure 36 illustrates the level of agreement between the tests.

The baseline evaluation results for *definiteness* showed an accuracy of 0.9575. The Low EDI score test yielded an accuracy of 0.5000. The High EDI score test demonstrated strong performance, achieving 95% of baseline accuracy with just a single dimension, as illustrated in Figure 35. The highest cross-property accuracy was achieved by *intensifier*, at 0.9400, followed closely by *factuality* (0.9325) and *synonym* (0.9275).

1292

1294

1282

1293 E.4 FACTUALITY

Figure 37 highlights the difference between the most prominent dimensions encoding this property. Figure 39 illustrates the level of agreement between the tests.



Figure 26: BERT Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Voice*.

1332

1334

The baseline evaluation results for *factuality* showed an accuracy of 1.0000. The Low EDI score test yielded an accuracy of 0.6800. The High EDI score test demonstrated strong performance, achieving 95% of baseline accuracy with just a single dimension, as illustrated in Figure 38. The highest cross-property accuracy was achieved by *negation*, at 0.9975.

1333 E.5 INTENSIFIER

Figure 40 highlights the difference between the most prominent dimensions encoding this property. Figure 42 illustrates the level of agreement between the tests.

The baseline evaluation results for *intensifier* showed an accuracy of 1.0000. The Low EDI score test yielded an accuracy of 0.5825. The High EDI score test demonstrated steady improvement, reaching 95% of baseline accuracy with 4 dimensions, as illustrated in Figure 41. The highest cross-property accuracy was achieved by *definiteness*, at 0.9600.

1341

1342 1343 E.6 NEGATION

1343 1344

Figure 43 highlights the difference between the most prominent dimensions encoding this property. Figure 45 illustrates the level of agreement between the tests.

- The baseline evaluation results for *negation* showed an accuracy of 0.9850. The Low EDI score test yielded an accuracy of 0.5450. The High EDI score test demonstrated steady improvement, reaching
- 95% of baseline accuracy with 6 dimensions, as illustrated in Figure 44. The highest cross-property accuracy was achieved by *intensifier*, at 0.9475.







1511 Figure 65 highlights the difference between the most prominent dimensions encoding this property. Figure 67 illustrates the level of agreement between the tests.



Figure 33: Mutual Information of GPT-2 Embedding Dimensions overlaid with Wilcoxon test and RFE results for *Control*.



Figure 34: GPT-2 Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Definiteness*.

1525 1526

1534

1535

The baseline evaluation results for *definiteness* showed an accuracy of 0.9000. The Low EDI score test yielded an accuracy of 0.4000. The High EDI score test demonstrated strong performance, achieving 95% of baseline accuracy with just a single dimension, as illustrated in Figure 66. The highest cross-property accuracy was achieved by *intensifier*, at 0.6750.

1548

1549 F.4 FACTUALITY 1550

Figure 68 highlights the difference between the most prominent dimensions encoding this property.Figure 70 illustrates the level of agreement between the tests.

The baseline evaluation results for *factuality* showed an accuracy of 0.9975. The Low EDI score test yielded an accuracy of 0.4825. The High EDI score test demonstrated steady performance, achieving 95% of baseline accuracy with 16 dimensions, as illustrated in Figure 69. The highest cross-property accuracy was achieved by *quantity*, at 0.8875.

1557 1558

1559 F.5 INTENSIFIER

Figure 71 highlights the difference between the most prominent dimensions encoding this property.
 Figure 73 illustrates the level of agreement between the tests.

The baseline evaluation results for *intensifier* showed an accuracy of 0.9000. The Low EDI score test yielded an accuracy of 0.4200. The High EDI score test demonstrated slow performance, achieving 95% of baseline accuracy with 347 dimensions, as illustrated in Figure 72. The highest cross-property accuracy was achieved by *quantity*, at 0.6825.





1673 The baseline evaluation results for *voice* showed an accuracy of .9175. The Low EDI score test yielded an accuracy of 0.3875. The High EDI score test demonstrated slow improvement, reaching



Figure 40: GPT-2 Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Intensifier*.

1705
1706
1706
1707
95% of baseline accuracy with 263 dimensions, as illustrated in Figure 90. The highest cross-property accuracy was observed with *definiteness* at 0.6225.















Figure 60: Mutual Information of GPT-2 Embedding Dimensions overlaid with Wilcoxon test and RFE results for *Voice*.



Figure 61: Linguistic Property Classifier results for MPNet.



Figure 62: MPNet Dimensional Embedding values for the Wilcoxon test results with the most significant p-values for *Control*.



















Figure 88: Mutual Information of MPNet Embedding Dimensions overlaid with Wilcoxon test and RFE results for *Tense*.

