

A General Framework to Evaluate Methods for Assessing Dimensions of Lexical Semantic Change Using LLM-Generated Synthetic Data

Anonymous ACL submission

Abstract

Lexical Semantic Change (LSC) offers insights into cultural and social dynamics. Yet, the validity of methods for measuring kinds of LSC has yet to be established due to the absence of historical benchmark datasets. To address this gap, we develop a novel three-stage evaluation framework that involves: 1) creating a scalable, domain-general methodology for generating synthetic datasets that simulate theory-driven LSC across time, leveraging In-Context Learning and a lexical database; 2) using these datasets to evaluate the effectiveness of various methods; and 3) assessing their suitability for specific dimensions and domains. We apply this framework to simulate changes across key dimensions of LSC (SIB: Sentiment, Intensity, and Breadth) using examples from psychology, and evaluate the sensitivity of selected methods to detect these artificially induced changes. Our findings support the utility of the synthetic data approach, validate the efficacy of tailored methods for detecting synthetic changes in SIB, and reveal that a state-of-the-art LSC model faces challenges in detecting affective dimensions of LSC. This framework provides a valuable tool for dimension- and domain-specific benchmarking and evaluation of LSC methods, with particular benefits for the social sciences.

1 Introduction

Lexical Semantic Change (LSC) provides a unique window into cultural dynamics by revealing how language evolution reflects social changes. Recently developed state-of-the-art (SOTA) computational methods have expanded our ability to classify established types of LSC, such as generalization and specialization (Cassotti et al., 2024a). Efforts have also been directed towards developing methods for measuring newly proposed dimensions of LSC (Baes et al., 2024; de Sá et al., 2024). Nevertheless, the field faces challenges in validating these methods. A major obstacle is the absence of

historical benchmark datasets, which restricts the standardization and fair comparison of metrics. Additionally, there is a pressing need for fine-grained evaluation methods that save time and resources.

To address these challenges, the present study introduces a three-stage evaluation framework. It: 1) develops a scalable, domain-general methodology for generating high-quality synthetic sentences that leverage In-Context Learning (ICL) and a lexical database to simulate changes in kinds of LSC; 2) uses these newly constructed historical datasets to evaluate the relative effectiveness of computational approaches; and 3) identifies the more suitable method for specific dimensions and domains. This framework is applied to assess the sensitivity of various methods to detect synthetic change in major LSC dimensions—Sentiment, Intensity, and Breadth (SIB; Baes et al. 2024)—using examples drawn from psychology. Our findings confirm the validity of theory-driven changes using synthetic SIB datasets and emphasize the need to tailor methods to particular dimensions, as the SOTA LSC model was found to be ineffective at detecting affective dimensions. This framework provides an efficient and scalable solution for dimension- and domain-specific benchmarking and evaluation of LSC methods. While this innovation is generally applicable, it is particularly beneficial for the social sciences and humanities, where customized methods are essential for analyzing complex constructs.

2 Related Work

2.1 Theoretical Background

Linguists have long debated taxonomies of LSC (Bloomfield, 1933; Blank, 1999), defined as innovations which change the lexical meaning of a form (Bloomfield, 1933). A growing body of work has identified ways to detect changes in the meanings of words and quantify the extent of these changes using a variety of computational approaches (Ku-

081 tuzov et al., 2018; Tahmasebi et al., 2018; Tang,
082 2018; Cassotti et al., 2024b; Periti and Montanelli,
083 2024a; Kiyama et al., 2025).

084 Recent years have seen the development of the-
085 oretical frameworks that propose multiple dimen-
086 sions of LSC. Baes et al. (2024) introduced a three-
087 dimensional framework that maps LSC along axes
088 of SIB, reflecting a word’s acquisition of more pos-
089 itive or negative connotations (Sentiment), more
090 or less emotionally charged or potent connotations
091 (Intensity), and the expansion or contraction of
092 its semantic range (Breadth). It draws on linguis-
093 tic (Geeraerts, 2010) and psychological (Haslam,
094 2016) theories, and provides methodological tools
095 to estimate SIB across time. In parallel, de Sá et al.
096 (2024) proposed a framework that clusters LSC into
097 three dimensions using graph structures: Orientation
098 (shifts towards more pejorative or ameliorated
099 senses), Relation (changes towards metaphoric or
100 metonymic usage), and Dimension (variations be-
101 tween abstract/general and specific/narrow mean-
102 ings). While de Sá et al. (2024) surveyed statistical
103 methods for representing word meaning (word fre-
104 quency, topic modeling, and graph structures) on
105 dimensions, they did not demonstrate their usage.

106 Both frameworks contain dimensions of eval-
107 uation (Sentiment and Orientation) and seman-
108 tic range (Breadth and Dimension). Baes et al.’s
109 (2024) inclusion of Intensity reflects a greater
110 emphasis on changes in the emotional connota-
111 tions of words. Sentiment and Intensity resemble
112 the two primary dimensions of human emotion,
113 Valence and Arousal (Russell, 2003), and
114 two primary dimensions of connotational meaning,
115 Evaluation (e.g., “good/bad”) and Potency (e.g.,
116 “strong/weak”) (Osgood et al., 1975), which have
117 been demonstrated to have cross-cultural validity.

118 2.2 Evaluation

119 Despite substantial progress in developing bench-
120 marks (Tahmasebi and Risse, 2017) and evalua-
121 tion strategies (Kutuzov et al., 2018), the field still
122 lacks standardized datasets that evaluate multiple
123 dimensions of LSC across time. Current annotated
124 benchmarks, such as the synchronic, definition-
125 and type-based *LSC Cause-Type-Definitions Bench-*
126 *mark* (Cassotti et al., 2024a) and the binary, word-
127 sense-based *TempoWIC*, where LSC is labeled by
128 comparing the sameness or difference of meanings
129 between two sense usages (Loureiro et al., 2022),
130 address different aspects of semantic change.

131 The first human-annotated dataset of LSC in mul-

132 tiple languages (English, German, Latin, Swedish;
133 Schlechtweg et al., 2020) represented substantial
134 progress in indicating the presence and degree
135 of LSC, but omitted information about kinds of
136 change. Creating expert-annotated datasets of LSC
137 is costly and time-intensive. Recognizing this gap,
138 Dubossarsky et al. (2019) introduced a method to
139 artificially induce semantic change in controlled
140 testing environments, allowing for precise testing
141 of how well models capture these shifts.

142 Recent developments in generative artificial in-
143 telligence highlight the potential of pre-trained
144 LLMs to adapt to novel tasks at inference time
145 through ICL (Zhou et al., 2023). Few-shot ICL,
146 a paradigm that enables LLMs to learn tasks by
147 analogy given only a few demonstrative examples,
148 helps to incorporate theoretical knowledge without
149 needing to fine-tune its internal parameters (Dong
150 et al., 2024). Instead, ICL uses context from the
151 model’s prompt to adapt the LLM to downstream
152 tasks (Radford et al., 2019; Brown et al., 2020;
153 Liu et al., 2024). de Sá et al. (2024) demonstrated
154 the utility of few-shot ICL, employing Chain-of-
155 Thought and rhetorical devices, to annotate LSC di-
156 mensions, but their strategy focuses on multi-class
157 classification of change between two sense usages.
158 ICL offers a promising solution to bridge the ab-
159 sence of standardized approaches (Hengchen et al.,
160 2021) for assessing the effectiveness of different
161 methods to measure dimensions of LSC.

162 2.3 The Present Study

163 The present study aims to develop an evaluation
164 framework that: (1) creates a scalable, domain-
165 general methodology for constructing high-quality
166 LLM-generated datasets labeling changes in LSC
167 dimensions; (2) uses these synthetic datasets to
168 compare the validity of proposed computational ap-
169 proaches; and (3) identifies the suitability of meth-
170 ods for each dimension and domain. We apply this
171 framework to major LSC dimensions defined by
172 Baes et al. (2024) (SIB; see Table 1) on a sample of
173 words drawn from a corpus of academic psychol-
174 ogy articles. Key questions include:

1. Can synthetic datasets validate methods to
175 measure dimensions of LSC? We predict that
176 SIB scores will be linearly associated with
177 levels of synthetic change.
2. Which out of a set of LSC detection meth-
178 ods is most sensitive to synthetically induced
179 changes in SIB?

Dimension	Definition	Examples of Rising	Examples of Falling
Sentiment	Relates to the degree to which a word's meaning acquires more positive ('elevation', 'amelioration') or negative ('degeneration', 'pejoration') connotations.	<i>craftsman</i> , once associated with manual labor, has now come to convey artistry, skill, and high-quality workmanship. <i>geek</i> , from a derogatory term for odd people, to reference someone passionate about a specific field.	<i>retarded</i> , originally a neutral term for intellectual disability, has become highly pejorative over time. <i>awful</i> has shifted from its original meaning of "awe-inspiring" to its modern usage which indicates something very bad.
Intensity	Relates to the degree to which a word's meaning changes to acquire more ('meiosis') or less ('hyperbole') emotionally charged (i.e., strong, potent, high-arousal) connotations.	<i>cool</i> has evolved from describing temperate to expressing strong approval or trendiness. <i>hilarious</i> , originally meaning cheerful or amusing in Latin, has come to describe extremely funny things that cause great merriment and laughter.	<i>love</i> , evolved from a romantic or platonic attachment to a milder expression of liking (e.g., "I love pizza.") <i>trauma</i> , from referencing brain injuries to referring to less severe events (e.g., business loss).
Breadth	Relates to the degree to which a word expands ('widening', 'generalization') or contracts ('narrowing', 'specialization') its semantic range.	<i>cloud</i> , initially a meteorological term, broadened its use to reference internet-based data storage. <i>partner</i> , originally referring to business co-owners, now also describes a significant other in a romantic or domestic relationship.	<i>doctor</i> , once referring to any scholar or teacher, now primarily refers to a medical professional. <i>meat</i> , originally referred to any kind of food in Old English ('mete'), but its meaning has narrowed to specifically denote animal flesh as food.

Table 1: Definitions and Examples of Baes et al.'s (2024) Dimensions of Lexical Semantic Change.

3 Method

3.1 Materials

3.1.1 Psychology Corpus

To develop and test the evaluation pipeline on a specific domain, a corpus of psychology article abstracts was sourced (Vylomova et al., 2019). It includes 133,017,962 tokens from 871,337 abstracts (1970-2019) from E-Research and PubMed databases, and contains 5,214,227 sentences.¹

3.1.2 WordNet

Although other ontologies were considered,² the English WordNet lexical database 3.0 (Miller, 1992) was chosen for its linguistic coverage and lexical structure. It organizes words into synsets (synonyms with distinct meanings), linking them by semantic relationships (e.g., hypernyms, hyponyms).

3.1.3 Targets

While the evaluation framework is general in its applicability, six terms from psychology—*abuse*, *anxiety*, *depression*, *mental health*, *mental illness*, and *trauma*—are analyzed for semantic change, selected for their empirical and theoretical relevance to shifting word meanings. *Trauma*, *mental health*, and *mental illness* have seen falls in their average valence and semantic expansions (*trauma*: Baes et al., 2023; Haslam et al., 2021; *mental health*, *mental illness*: Baes et al., 2024). There have been changes in the intensity of their meanings, with rises for *mental health* and *mental illness* (Baes et al., 2024), as well as *anxiety* and *depression* (Xiao et al., 2023), and a fall for *trauma* (Baes et al.,

2023). Qualitatively, the concept of *abuse* has expanded horizontally to include passive neglect and emotional abuse, beyond its initial physical scope (Haslam, 2016). Each target was sufficiently prevalent in the corpus for robust analysis, with sentence counts at 46,272; 104,486; 115,430; 44,130; 5,808; and 23,187 (see Appendix A for annual counts).

3.2 Evaluation Framework

The general pipeline for the evaluation framework is presented in Figure 1. First, synthetic datasets are constructed to benchmark changes in LSC dimensions using few-shot ICL and a lexical database. GPT-4o (Achiam et al., 2023)³ is prompted with expert-crafted examples to increase and decrease corpus sentences in affective dimensions across 5-year intervals. This ensures that synthetic sentences are theory-driven, domain-specific and contain temporal features. GPT is used due to its adeptness at few-shot learning, task adaptation with minimal examples (Achiam et al., 2023; Merx et al., 2024) and lack of disciplinary bias (Ziems et al., 2024). Appendix B details statistics for synthetic datasets,⁴ which are validated using tools shown to measure SIB in two corpora (Baes et al., 2024).

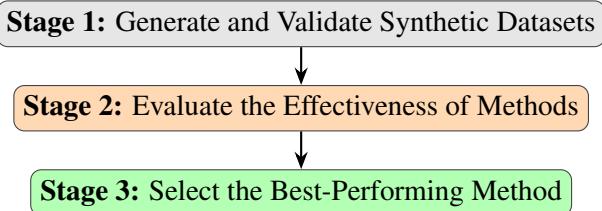


Figure 1: General Stages for the Evaluation Framework.

¹Sentences were segmented using "en_core_web_sm" (<https://spacy.io/models/en>); F-score = 91%.

²PsychNET, UMLS, DSM-5, ConceptNet

³ChatGPT API documentation: <https://platform.openai.com/docs/guides/text-generation>

⁴Link to Synthetic datasets: [MASKED LINK].

To assess the relative effectiveness of different methods, sentences are sampled from the natural and synthetic corpora using two sampling strategies. Bootstrap sampling draws 50 sentences with replacement from both corpora 100 times (i.e., iterations). Five-year random sampling selects up to 50 sentences from fixed intervals 10 times, ensuring each time period is equally represented. Each iteration forces unique sentence selection while permitting sentence repetition across different rounds to reflect natural language. Control conditions shuffle these sentences to balance authentic:synthetic sentences in each sample, verifying the synthetic effect in the genuine condition and its absence in the control. Following computational linguistics precedents (Dubossarsky et al., 2017, 2019), this approach validates the impact of synthetic interventions by providing a baseline for comparison. Notably, for each strategy, synthetic sentences are injected into natural samples at increasing injection levels (20%, 40%, 60%, 80%, and 100%), as illustrated in Appendix B. Bins are injection levels (for bootstrap) and epoch (for 5-year intervals). This simulates controlled semantic saturation scenarios to assess the sensitivity of methods to semantic variation (increased/decreased SIB) (stage 2) and select the method that detects greater magnitude of change from 0% to 100% injection (stage 3).

3.3 Sentiment and Intensity

3.3.1 Synthetic Sentiment and Intensity

To generate the synthetic Sentiment and Intensity datasets, we employ few-shot ICL with GPT-4o to vary these dimensions. First, neutral sentences from the corpus (detailed in Section 3.1.1) are sampled as outlined in Appendix C. Second, a psychology scholar crafts five (Chen et al., 2023) diverse examples of sentence variations for each target following the task detailed below, which includes construct definitions to generate theory-driven change. For ‘scholar-in-the-loop’ few-shot demonstrations, see Appendix D (Sentiment) and Appendix E (Intensity). Third, the prompt is refined during pilot tests (10 inputs). Fourth, for each of the neutral sentences (Sentiment: 36,151; Intensity: 39,896), we make one inference call to GPT-4o through the OpenAI API to generate variations of Sentiment (positive/negative) or Intensity (high/low). Fifth, output sentences are manually adjusted (Sentiment: 0.25%; Intensity: 0.01%) due to GPT-4o’s failure to retain targets. See Appendix C for the counts of

input/output sentences and final prompts (dataset costs for Sentiment: 115 \$US; Intensity: 136 \$US).

Prompt Outline for Synthetic Sentiment

Prompt intro: In psychology research, ‘Sentiment’ is defined as “a term’s acquisition of a more positive or negative connotation.” This task focuses on the sentiment of the term **target_word**.

Task: You will be given a sentence containing the term **target_word**. Your goal is to write two new sentences:

1. One where **target_word** has a **more positive connotation**.
2. One where **target_word** has a **more negative connotation**.

Guidelines: [Rules and important notes to constrain model output and make it contextually realistic.]

Append few-shot examples: [One example below.]

Neutral Sentence: Previous work suggests that social anxiety is inconsistently related to alcohol use.

Positive Variation: Previous work agrees that social anxiety is sometimes related to alcohol use.

Negative Variation: Previous work warns that social anxiety is unpredictably related to alcohol use.

Prompt Outline for Synthetic Intensity

Prompt intro: In psychology research, ‘Intensity’ is defined as “the degree to which a word has emotionally charged (i.e., strong, potent, high-arousal) connotations.” This task focuses on the intensity of the term **target_word**.

Task: You will be given a sentence containing the term **target_word**. Your goal is to write two new sentences:

1. One where **target_word** is **less intense**.
2. One where **target_word** is **more intense**.

Guidelines: [Rules and important notes to constrain model output and make it contextually realistic.]

Append few-shot examples: [One example below.]

Neutral sentence: They tend to be more liberal in their attitudes toward abortion than women in general; however, women who experienced a greater degree of psychic trauma tended to be more conservative in their attitudes.

Low Variation: They tend to be more accepting in their attitudes toward children than women in general; however, women who experienced mild psychic trauma tended to be more conservative in their attitudes.

High Variation: They tend to be more extremely callous in their attitudes toward the horrors of abortion than women in general; however, women who suffered a greater degree of violent psychic trauma tended to be more fearful in their attitudes.

3.3.2 Quantifying Sentiment and Intensity

To measure shifts in a word’s connotations from negative to positive Sentiment and from low to high Intensity, we adapt Baes et al.’s (2024) method. Sentences are processed.⁵ Collocates (± 5 words from the target) within sentences are assigned ordinal valence or arousal scores based on Warriner et al. (2013) norms, ranging from *extremely unhappy* (1: “unhappy”, “despaired”) to *extremely happy* (9: “happy”, “hopeful”) for valence, and from *extremely low* (1: “calm”, “unaroused”) to *extremely high* (9: “agitated”, “aroused”) for arousal. Valence (V) and arousal (A) indices are calculated as shown in Equation 1:

$$V_{t_j,k}, A_{t_j,k} = \frac{\sum_{i=1}^{n_{j,k}} w_{i,j,k} x_{i,j,k}}{\sum_{i=1}^{n_{j,k}} w_{i,j,k}} \quad (1)$$

where $w_{i,j,k}$ denotes the frequency of each collocate i in iteration k within bin t_j , and $x_{i,j,k}$ denotes its valence or arousal rating at bin t_j within iteration k . Here, $n_{j,k}$ is the number of collocates in iteration k within bin t_j . Scores are weighted by the collocate’s frequencies within each iteration and normalized by the total occurrences in that iteration. Scores are averaged across all iterations within each bin, conditioned on whether the Sentiment is positive/negative, or the Intensity is high/low. These indices provide a mean valence or arousal score per iteration in each bin t_j , with higher scores indicating a more positive valence or higher arousal. Scores (1-9) are normalized to range from 0 (extremely unhappy/low arousal) to 1 (extremely happy/high arousal).

While the Intensity dimension is novel and lacks existing comparative models, for Sentiment, we compare the interpretable Valence index against DeBERTa-v3-ABSA, a SOTA classification model in aspect-based sentiment analysis (ABSA). Deberta-v3-base-absa-v1.1⁶ identifies sentiment associated with particular aspects of an entity within text (here, the target term). It was initially trained on restaurant and laptop reviews (Cabello and Akujuobi, 2024; Yang et al., 2021, 2023). We adapt it to produce continuous sentiment scores, which reflect the model’s confidence in positive sentiment associated with the target term and range from 0 (fully negative) to 1 (fully positive).⁷

⁵Tokenization, lemmatization, stop-word removal using “en_core_web_sm” (<https://spacy.io/models/en>)

⁶yangheng/deberta-v3-base-absa-v1.1 (184M model params): <https://huggingface.co/yangheng/deberta-v3-base-absa-v1.1>

⁷The sentiment score is calculated as follows: $0 \times$

3.4 Breadth

3.4.1 Synthetic Breadth

Unlike Sentiment and Intensity, current Breadth measures have no score that assigns a mid-point with which to obtain neutral sentences to vary. Therefore, to simulate semantic breadth, we adapt Dubossarsky et al.’s (2019) replacement strategy, using WordNet 3.0 to expand a target word’s usage by incorporating contexts from donor terms, broadening its semantic range without altering its core meaning. Relevant synsets are identified and filtered for psychological relevance using keyword matching⁸ and semantic similarity thresholds. Donor terms (co-hyponyms with the target) are filtered using Lin similarity (0.5)⁹ and cosine similarity (0.7) with embeddings from BioBERT (Lee et al., 2020), a pre-trained language model for biomedical text mining, to capture context-dependent meanings of synset glosses in 768-dimensional vectors. See Appendix F for the list. The sibling replacement process identifies and replaces sibling terms with the target, shown below. To sample representatively from the sibling list, a round-robin strategy is used, sampling up to 1,500 unique sentences per epoch per injection level to create the final synthetic breadth dataset (0 \$US).

Dataset Creation for Synthetic Breadth

Replacement Strategy: Randomly sample sentences containing co-hyponyms of the target term from the validated list and replace the **co-hyponym** with the **target** to be used as a synthetic sentence.

[One example for **mental_health** below.]

Donor Context: The ‘Angry and Impulsive Child’ and ‘Abandoned and Abused Child’ modes uniquely predicted **dissociation** scores.

Synthetic Context: The ‘Angry and Impulsive Child’ and ‘Abandoned and Abused Child’ modes uniquely predicted **mental_health** scores.

3.4.2 Quantifying Breadth

To estimate the semantic broadening (expansion) or narrowing (contraction) of a word’s meaning, we calculate the average cosine distance between sentence-level embeddings of a target term, as in Baes et al. (2024). The SentenceTransformer

$\text{negative_prob} + 0.5 \times \text{neutral_prob} + 1 \times \text{positive_prob}$.

⁸Psychology key terms: “abnormality”, “abnormally”, “emotional”, “feeling”, “feelings”, “harm”, “hurt”, “mental”, “mind”, “psychological”, “psychology”, “psychiatry”, “syndrome”, “therapy”, “treatment”.

⁹Information content values from the psychology corpus

369 model ‘all-mpnet-base-v2’¹⁰ is used to generate
 370 these embeddings. The Breadth score, B , is derived
 371 by averaging the cosine distances, δ , across all
 372 unique pairs of sentence embeddings within each
 373 iteration, and then averaging these scores across all
 374 iterations within each bin, as shown in Equation 2:

$$B_{t_j} = \frac{1}{I_j} \sum_{k=1}^{I_j} \left(\frac{2}{N_k(N_k - 1)} \sum_{i=1}^{N_k-1} \sum_{j=i+1}^{N_k} \delta(s_{i,k}^{t_j}, s_{j,k}^{t_j}) \right) \quad (2)$$

375 Here, $\delta(s_{i,k}^{t_j}, s_{j,k}^{t_j})$ calculates the cosine distance
 376 between two sentence embeddings in the same iteration k in bin t_j . N_k is the number of sentence
 377 embeddings in iteration k ; I_j is the number of iterations in bin t_j . Higher scores indicate greater
 378 variation in the target’s semantic range. Scores
 379 range from 0 (no variation) to 1 (max variation).

380 We compare the sentence transformer “all-
 381 mpnet-base-v2” (MPNet) with Cassotti et al.’s
 382 (2023) SOTA word transformer “XL-LEXEME”¹¹
 383 (XLL). While MPNet generates sentence embeddings
 384 through pooling tokens, which dilutes word-
 385 specific information, XLL uses a bi-encoder archi-
 386 tecture that focuses on word-specific attention,¹²
 387 using polysemy as a proxy for meaning divergence
 388 during training (WIC; Pilehvar and Camacho-
 389 Collados, 2019).

393 3.5 General Lexical Semantic Change

394 To quantify general LSC, we use the SOTA LSC
 395 score (Cassotti et al., 2023), which calculates the
 396 Average Pairwise Cosine Distances (Giulianelli
 397 et al., 2020) between sentence embeddings from
 398 two time periods. We extend it to compare embed-
 399 dings from different bins within the same iteration,
 400 as shown in equation 3:

$$LSC_i(s_i^{t_0}, s_i^{t_1}) = \frac{1}{N_i^2} \sum_{m=1}^{N_i} \sum_{n=1}^{N_i} \delta(s_{m,i}^{t_0}, s_{n,i}^{t_1}) \quad (3)$$

401 Here, N_i represents the number of sentence em-
 402 beddings within each iteration i in each bin. The
 403 term $\delta(s_{m,i}^{t_0}, s_{n,i}^{t_1})$ measures the cosine distance
 404 between pairs of sentence embeddings from the
 405 same iteration i across two different bins t_0 and t_1 .
 406 Higher LSC scores indicate greater LSC, ranging
 407 from 0 (no change) to 1 (maximum change).

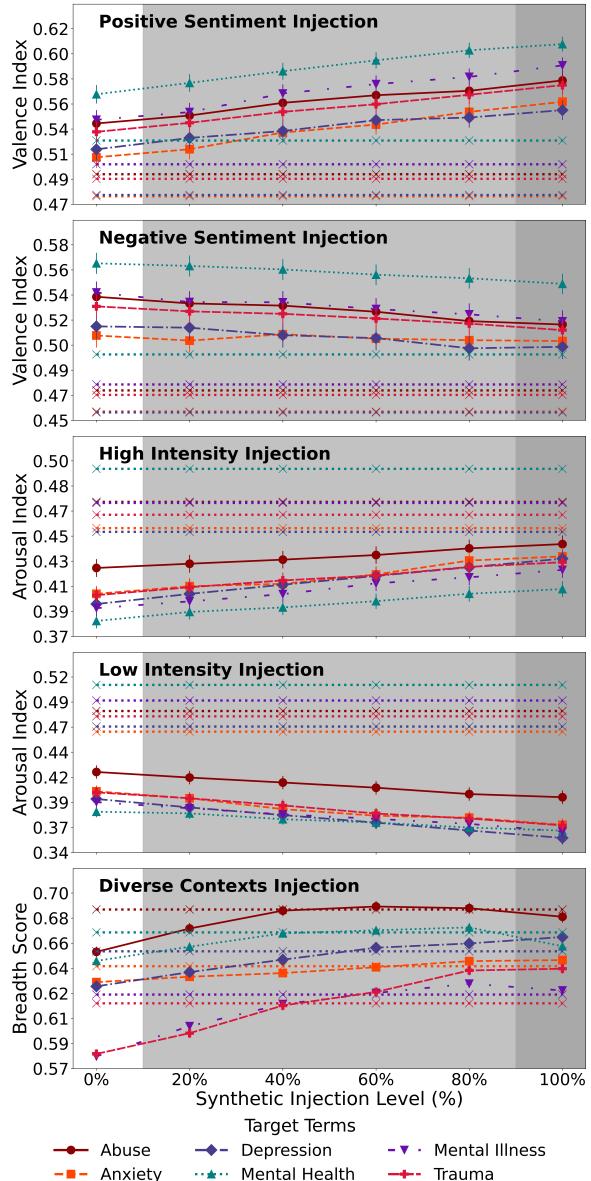
¹⁰Microsoft pretrained network (109M
 408 model params) <https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

¹¹XL-LEXEME (~550M model params) <https://huggingface.co/pierluigic/xl-lexeme>

¹²Only the first occurrence of the target is attended to.

4 Results

410 **Synthetic Change Effects:** The hypothesis that
 411 scores from Baes et al.’s (2024) SIB tools will be
 412 linearly associated with levels of synthetic change
 413 is supported, as evidenced by rising or falling
 414 trends in SIB scores across all targets and condi-
 415 tions (Figure 2). Mixed linear models demon-
 416 strate increases or decreases on SIB scores for every 1-
 417 unit increase in synthetic injection level (detailed
 418 in Table 2 and Appendix G). SIB scores for the five-
 419 year sampling experiments depict similar trends in
 420 response to varying injection levels (Appendix H).



421 Figure 2: SIB Scores ($\pm SE$) by Injection Levels for
 422 Experimental and Control Settings (Flat Dotted Lines).

423 **Control Experiments:** As illustrated in Figure 2,
 424 controlling for synthetic injection level by re-

Score	Valence	Arousal	Breadth
β^+	.003*	.002*	<.0001*
β^-	-.001*	-.002*	N/A

Table 2: Coefficients of Mixed Linear Models Predicting SIB Scores from Injection Levels (Var. Target Intercept)

Note: β^+ and β^- represent the standardized coefficients for conditions (rise/fall), respectively. '*' indicates $p < .0001$, testing the null hypothesis that $\beta = 0$.

423
424
425
426
427
428
analyzing data with shuffled sentences for uniform distribution reveals flat SIB score trends in bootstrapped settings. Appendix H shows that, even with temporal shuffling within time bins, SIB scores in five-year samples tend to converge to a midpoint between natural and synthetic data.

429
430
431
432
433
434
Comparative Method Evaluation: Comparisons of the relative validity of alternative change detection methods yielded mixed results. To determine which method is more sensitive to synthetically induced changes in SIB, we compare their performance on a synthetic change detection task using an evaluation metric specified below.¹³

Percent Relative Change Index

Synthetic Change Detection Task: Detect the magnitude of change in the target word’s context when sampling sentences from a natural corpus (0%) and an entirely synthetic corpus (100%).

Percent relative change Δ is defined as:

$$\Delta\% = \frac{X_{100} - X_0}{X_0} \times 100$$

where X_0 represents the measure’s score at 0% synthetic injection (natural data) and X_{100} at 100% synthetic injection (fully artificial data).

436
437
438
439
440
441
442
443
444
445
446
For Sentiment, Valence index and ABSA’s Sentiment score are sensitive to detecting variations in synthetic Sentiment, although the ABSA score outperforms the Valence index 10/12 times. For Intensity, the Arousal index shows sensitivity to detecting variations in synthetic Intensity. For Breadth, XLL outperforms MPNet (4/6 times) on detecting rises in synthetic Breadth using the Breadth score.

Critically, XLL-LSC score is completely insensitive to detecting changes in either Sentiment or

¹³For XLL’s LSC Score, Δ is normalized against the intrinsic within-bin variability in both bins of interest:

$$\Delta = \frac{\text{APD}(X_{100}\text{-between}-X_0)}{\max [\text{APD}(X_0\text{-within}-X_0), \text{APD}(X_{100}\text{-within}-X_0)]}$$

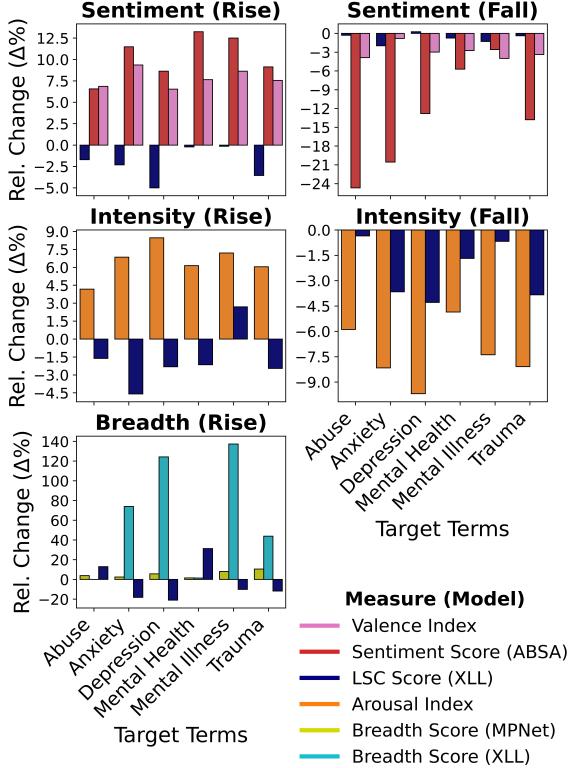


Figure 3: Relative Change ($\Delta\%$) Scores for Models Across Dimensions and Conditions: Bootstrapped.

Intensity. XLL-LSC can only indicate change via positive change values, while negative values indicate that the within-bin variance is greater than the change scores between bins. See Appendix I for between- and within-bin LSC scores across *all* synthetic injection levels. Thus, the negative scores observed in Sentiment and Intensity (except for *Mental Illness*) establish that XLL was unable to detect any change signal in these words. XLL-LSC detects changes in synthetic Breadth for 2/6 terms.

5 Discussion

The present study introduced a three-stage general domain evaluation framework that: 1) creates synthetic datasets featuring ‘scholar-in-the-loop’ LLM-generated sentences to simulate various kinds of LSC; 2) leverages these datasets to assess the sensitivity of computational approaches to synthetic changes; and 3) evaluates the suitability of these methods for specific dimensions and domains. This framework is applied to generate synthetic datasets that induce changes across the three dimensions of a recent multidimensional LSC framework (SIB; Baes et al., 2024), using examples from psychology, to test and compare the suitability of different methods in detecting these synthetic changes.

472 Our findings support the hypothesis that recently
473 proposed methods (Valence index, Arousal index,
474 Breadth score; Baes et al., 2024) detect synthetic
475 changes on the SIB dimensions. Control analy-
476 ses, which adhered to computational linguistics
477 standards (Dubossarsky et al., 2017, 2019), con-
478 firmed the absence of these effects in shuffled con-
479 trols. The implications of these findings are two-
480 fold. The ability of SIB methods to detect changes
481 when introducing silver-label synthetic data vali-
482 dates their sensitivity and reliability in detecting
483 and measuring variations in SIB, even in controlled,
484 artificial environments. This, in turn, validates
485 the synthetic LLM-generated sentences in our ICL
486 evaluation suites.

487 We demonstrated how a synthetic change de-
488 tection task can assess the sensitivity of various
489 computational approaches, guiding the selection
490 of the most suitable model for specific dimensions
491 and domains. Baes et al.’s (2024) tools, which val-
492 idated the synthetic SIB datasets, were supported
493 by alternative methods that consistently detected
494 synthetic changes in SIB across all conditions and
495 targets, providing further validation. The Valence
496 index and Sentiment score (ABSA) identified varia-
497 tions in synthetic Sentiment, while Breadth scores
498 (XLL and MPNet) detected increases in synthetic
499 Breadth. Results suggest that these NLP-based
500 methods are more sensitive in detecting synthetic
501 changes than Warriner-based methods, which rely
502 on Valence and Arousal ratings. Future empirical
503 studies on Sentiment and Breadth may consider
504 adopting these NLP models, either as replacements
505 for, or in addition to, existing methods.

506 Notably, when computing the general LSC score
507 using the SOTA LSC model XLL (Cassotti et al.,
508 2023), it was not sensitive to detecting Sentiment
509 and Intensity. Although XLL shows some sensi-
510 tivity to identifying synthetic increases in Breadth,
511 it registers a more substantial change when the
512 Breadth score is adjusted according to the method
513 introduced by Baes et al. (2024). It uses the within-
514 bin average cosine distance of target containing
515 sentences as a proxy for the expansion (broadening)
516 or contraction (narrowing) of a word’s contextual
517 usage. The inability of XLL to detect the affective
518 dimensions of LSC highlights the necessity of
519 evaluating SOTA models before deploying them
520 in new domains. Future research should investi-
521 giate whether this weakness in detecting affective
522 dimensions is specific to XLL or extends to other
523 contextualized models in more corpora. This in-

524 quiry is particularly salient given recent advances in
525 analyzing fine-grained, continuous semantic shifts
526 through “diachronic word similarity matrices using
527 fast and lightweight word embeddings over arbi-
528 trary time periods” (Kiyama et al., 2025).

529 Findings highlight the need to include affective
530 and connotational aspects of meaning in studies
531 of LSC. In particular, future studies must consider
532 emotional meaning in language models. While psy-
533 chology has extensively used language to analyze
534 emotion semantics (Jackson et al., 2022; Boyd and
535 Schwartz, 2021), advances in NLP are still explor-
536 ing how to build models that incorporate sentiment
537 (Goworek and Dubossarsky, 2024) and detect emotion
538 (Mohammad, 2021). Further research is re-
539 quired to detect affective states from text given the
540 cultural and universal aspects of emotion semantics
541 (Jackson et al., 2019). These findings also have
542 implications for existing multidimensional frame-
543 works of LSC (Baes et al., 2024; de Sá et al., 2024)
544 as the evaluation framework provides experimen-
545 tal settings in which to compare the sensitivity of
546 methods to detecting synthetic changes on specific
547 dimensions and domains in a variety of disciplines.

6 Conclusion

548 The current study introduced a novel general do-
549 main evaluation framework. Its three-stage pipeline
550 involves: 1) developing a scalable methodology for
551 generating LLM-based synthetic datasets with sil-
552 ver labels that simulate changes in kinds of LSC;
553 2) using these datasets to evaluate the relative sen-
554 sitivity of computational approaches in a synthetic
555 change detection task; and 3) identifying the most
556 suitable method for detecting synthetically induced
557 changes across specific dimensions and domains.
558 We applied this framework to a set of psycholog-
559 ical terms. Findings not only supported the validity
560 of proposed computational methods for measuring
561 changes in SIB, but also established a controlled ex-
562 perimental standard for rigorously evaluating exist-
563 ing LSC detection methods and exploring alterna-
564 tive computational approaches. This work is crucial
565 for addressing the substantial gap created by the
566 lack of historical benchmark datasets, which has
567 previously hindered the standardization of metrics
568 and fair comparison of methods. While this innova-
569 tion benefits all disciplines (e.g., biomedicine, law,
570 theology), it is particularly valuable in the social
571 sciences and humanities, where unique methods
572 are often required to measure complex constructs.

574 Limitations

575 Limitations inform future directions. Evaluating
576 the quality of LLM prompt and demonstration ex-
577 amples in the few-shot ICL paradigm is challeng-
578 ing. As LLM evaluation standards are developed
579 (Chang et al., 2024; Ziems et al., 2024), future
580 research might explore automated strategies such as
581 updating prompts based on examples (DSPy)¹⁴ or
582 comparing LLM output from different prompts us-
583 ing a free, unified interface.¹⁵ LLM choice in the
584 evaluation pipeline could be expanded to include
585 open-source models (e.g., FlanT5-XL, Mistral-7B,
586 Mixtral-8x7B), enhancing its accessibility.

587 Furthermore, our study benefited from using
588 GPT-4o, which is trained on US English and is
589 therefore well-suited for analyzing texts within the
590 Western-centric domain of psychology. However,
591 the cultural and linguistic biases of LLMs may pose
592 challenges for adapting our evaluation pipeline to
593 other languages (Havaldar et al., 2023), although
594 few-shot ICL has proven effective in low-resource
595 languages (Cahyawijaya et al., 2024). Despite the
596 tendency of LLM training data to skew towards the
597 recent past, it successfully generated high-quality
598 sentences that spanned a 1970 to 2019 time period.
599 Future research should focus on refining these mod-
600 els to support broader application across various
601 cultural contexts, languages, and historical periods.

602 The conceptualization of semantic Breadth is
603 complex and contested. Linguistic definitions sug-
604 gest breadth encompasses subtypes (e.g., special-
605 ization as a subtype of narrowing; Campbell, 2013)
606 highlighting its intricate nature. Given this com-
607 plexity, it is essential to compare the current mea-
608 sure, which is based on mean within-bin variability
609 of target-containing sentences, with other methods
610 assessing breadth through senses, topics, or proto-
611 typical changes: modulations based on literal sim-
612 ilarity (Geeraerts, 1997). Future research should
613 investigate whether these measures can detect poly-
614 semy's emergence or merely prototype-based mod-
615 ulations of existing concepts.

616 The synthetic breadth dataset used in this study
617 was constructed using a replacement strategy that
618 may include contextually irrelevant donor con-
619 texts. To enhance simulation quality, we propose
620 a three-step validation pipeline: First, select vali-
621 dation models based on performance against a

gold-standard dataset, as determined by the high-
622 est F1-score from 5-fold stratified K-fold cross-
623 validation. Second, use a probability ratio check
624 with a Masked Language Model (e.g., BioBERT,
625 RoBERTa-large, DeBERTa-v3-large) to confirm
626 the plausibility of replacing donors with target
627 terms, approving sentences that meet a specific
628 probability threshold. Third, ensure semantic align-
629 ment through cosine similarity validation with mod-
630 els such as MiniLM-L12-v2 or DistilRoBERTa-v1
631 Sentence-T5, approving sentences that exceed a
632 set threshold. This process aims to expand the tar-
633 get term's semantic scope while maintaining speci-
634 ficity, but may exclude many sentences. Integrat-
635 ing de Sá et al.'s (2024) ICL approach to simulate
636 Breadth—first teaching the model to disambiguate
637 word senses—could offer an efficient alternative.

638 Furthermore, the present study does not specify
639 which sense of the term is semantically expanded.
640 Attempting to integrate senses into the synthetic
641 data generation pipeline may provide richer in-
642 sights. While the specialized psychology corpus
643 and target words exhibit limited senses, general do-
644 main corpora introduce ambiguous contexts (e.g.,
645 economic sense of "depression"). Notably, current
646 methods for word sense disambiguation may not in-
647 tegrate with distributional approaches as historical
648 linguists do not treat LSC as a set of senses.

649 Although a body of work estimates valence from
650 natural language, less research has examined the
651 Intensity dimension (Hoemann et al., 2025). In the
652 present study, this restricted the external validation
653 of the Arousal index (Baes et al., 2024), highlight-
654 ing the need for empirical research in this direc-
655 tion. Furthermore, we must examine the concep-
656 tual/terminological link between arousal and hyper-
657 bole (i.e., a linguistic form describing a rhetorical,
658 discursive phenomenon like irony) to understand
659 arousal's relation to hyperbole (Burgers et al., 2016;
660 Peña and Ruiz de Mendoza, 2017).

661 Finally, future research should use the evaluation
662 framework to generate synthetic datasets, and to
663 explore methods, for detecting the Relation dimen-
664 sion (metaphor/metonymy) as highlighted by de Sá
665 et al. (2024). Incorporating the qualitative types of
666 metaphor and metonymy into the empirical study
667 of multidimensional LSC could provide a more
668 comprehensive understanding of LSC, particularly
669 for some domains. Examining how Relation re-
670 lates to SIB may deepen our understanding of LSC
671 processes by exploring how cognitive principles
672 contribute to semantic innovations.

¹⁴<https://dspy.ai>

¹⁵<https://github.com/marketplace/models/azure-openai/gpt-4o/playground>

674 Ethical Considerations

675 We do not identify any foreseeable risks or potential
676 for harmful use of our work. Analyses containing
677 sentences from the psychology corpus use li-
678 censed data that are openly accessible for academic
679 purposes, ensuring transparency and accountabil-
680 ity. While LLM-generated sentences are synthetic,
681 they are modeled off samples from the psychology
682 corpus.

683 Acknowledgments

684 References

685 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
686 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
687 Diogo Almeida, Janko Altenschmidt, Sam Altman,
688 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
689 *arXiv preprint arXiv:2303.08774*.

690 Naomi Baes, Nick Haslam, and Ekaterina Vylomova.
691 2024. A multidimensional framework for evaluating
692 lexical semantic change with social science applica-
693 tions. In *Proceedings of the 62nd Annual Meeting of*
694 *the Association for Computational Linguistics (Volume*
695 *1: Long Papers)*, pages 1390–1415, Bangkok,
696 Thailand. Association for Computational Linguistics.

697 Naomi Baes, Ekaterina Vylomova, Michael Zyphur,
698 and Nick Haslam. 2023. The semantic inflation of
699 “trauma” in psychology. *Psychology of Language*
700 and *Communication*, 27(1):23–45.

701 Andreas Blank. 1999. Why do new meanings occur?
702 a cognitive typology of the motivations for lexical
703 semantic change. In Andreas Blank and Peter Koch,
704 editors, *Historical semantics and cognition*, pages
705 61–90. Mouton de Gruter.

706 Leonard Bloomfield. 1933. *Language*. Compton Print-
707 ing Works Ltd.

708 Ryan L Boyd and H Andrew Schwartz. 2021. Natural
709 language analysis and the psychology of verbal
710 behavior: The past, present, and future states of the
711 field. *Journal of Language and Social Psychology*,
712 40(1):21–41.

713 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
714 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
715 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
716 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
717 Gretchen Krueger, Tom Henighan, Rewon Child,
718 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
719 Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz
720 Litwin, Scott Gray, Benjamin Chess, Jack
721 Clark, Christopher Berner, Sam McCandlish, Alec
722 Radford, Ilya Sutskever, and Dario Amodei. 2020.
723 **Language models are few-shot learners.** In *Ad-*
724 *vances in Neural Information Processing Systems*,
725 volume 33, pages 1877–1901. Curran Associates,
726 Inc.

727 Christian Burgers, Elly A Konijn, and Gerard J Steen.
728 2016. Figurative framing: Shaping public discourse
729 through metaphor, hyperbole, and irony. *Communi-*
730 *cation theory*, 26(4):410–430.

731 Laura Cabello and Uchenna Akujuobi. 2024. It is
732 simple sometimes: A study on improving aspect-
733 based sentiment analysis performance. In *Findings of*
734 *the Association for Computational Linguistics: ACL*
735 *2024*, pages 6597–6610, Bangkok, Thailand. Associa-
736 tion for Computational Linguistics.

737 Samuel Cahyawijaya, Holly Lovenia, and Pascale Fung.
738 2024. LLMs are few-shot in-context low-resource
739 language learners. In *Proceedings of the 2024 Con-*
740 *ference of the North American Chapter of the Asso-*
741 *ciation for Computational Linguistics: Human Lan-*
742 *guage Technologies (Volume 1: Long Papers)*, pages
743 405–433, Mexico City, Mexico. Association for Com-
744 putational Linguistics.

745 Lyle Campbell. 2013. *Historical Linguistics: An In-*
746 *introduction*, ned - new edition, 3 edition. Edinburgh
747 University Press.

748 Pierluigi Cassotti, Stefano De Pascale, and Nina Tah-
749 masebi. 2024a. Using synchronic definitions and
750 semantic relations to classify semantic change types.
751 In *Proceedings of the 62nd Annual Meeting of the*
752 *Association for Computational Linguistics (Volume 1:*
753 *Long Papers)*, pages 4539–4553, Bangkok, Thailand.
754 Association for Computational Linguistics.

755 Pierluigi Cassotti, Francesco Periti, Stefano De Pas-
756 scale, Haim Dubossarsky, and Nina Tahmasebi. 2024b.
757 Computational modeling of semantic change. In *Pro-*
758 *ceedings of the 18th Conference of the European*
759 *Chapter of the Association for Computational Lin-*
760 *guistics: Tutorial Abstracts*, pages 1–8, St. Julian’s,
761 Malta. Association for Computational Linguistics.

762 Pierluigi Cassotti, Lucia Siciliani, Marco DeGennmis,
763 Giovanni Semeraro, and Pierpaolo Basile. 2023. XI-
764 lexeme: Wic pretrained model for cross-lingual lexical
765 semantic change. In *Proceedings of the 61st An-*
766 *nual Meeting of the Association for Computational*
767 *Linguistics (Volume 2: Short Papers)*, pages 1577–
768 1585, Toronto, Canada. Association for Computa-
769 tional Linguistics.

770 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
771 Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,
772 Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang,
773 Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie.
774 2024. A survey on evaluation of large language mod-
775 els. *ACM Trans. Intell. Syst. Technol.*, 15(3).

776 Juhai Chen, Lichang Chen, Chen Zhu, and Tianyi Zhou.
777 2023. How many demonstrations do you need for
778 in-context learning? In *Findings of the Association*
779 *for Computational Linguistics: EMNLP 2023*, pages
780 11149–11159, Singapore. Association for Compu-
781 tational Linguistics.

782	Jader Martins Camboim de Sá, Marcos Da Silveira, and Cédric Pruski. 2024. Semantic change characterization with llms using rhetorics. <i>arXiv preprint arXiv:2407.16624</i> .	836
783		837
784		838
785		839
786	Jader Martins Camboim de Sá, Marcos Da Silveira, and Cédric Pruski. 2024. Survey in characterization of semantic change. <i>Preprint</i> , arXiv:2402.19088.	840
787		841
788		842
789		843
790		844
791		845
792		846
793		847
794	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. A survey on in-context learning. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.	848
795		849
796		850
797	Haim Dubossarsky, Simon Hengchen, Nina Tahmasebi, and Dominik Schlechtweg. 2019. Time-out: Temporal referencing for robust modeling of lexical semantic change. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 457–470, Florence, Italy. Association for Computational Linguistics.	851
798		852
799		853
800		854
801		855
802		856
803		857
804	Haim Dubossarsky, Daphna Weinshall, and Eitan Grossman. 2017. Outta control: Laws of semantic change and inherent biases in word representation models. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 1136–1145, Copenhagen, Denmark. Association for Computational Linguistics.	858
805		859
806		860
807		861
808		862
809		863
810		864
811	Dirk Geeraerts. 1997. <i>Diachronic Prototype Semantics: A Contribution to Historical Lexicology</i> . Oxford: Clarendon Press.	865
812		866
813		867
814	Dirk Geeraerts. 2010. <i>Theories of lexical semantics</i> . Oxford University Press.	868
815		869
816	Andrew Gelman and Jennifer Hill. 2007. <i>Data Analysis Using Regression and Multilevel/Hierarchical Models</i> . Cambridge University Press, New York, NY.	870
817		871
818		872
819	Mario Giulianelli, Marco Del Tredici, and Raquel Fernández. 2020. Analysing lexical semantic change with contextualised word representations. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3960–3973, Online. Association for Computational Linguistics.	873
820		874
821		875
822		876
823		877
824		878
825		879
826	Roksana Goworek and Haim Dubossarsky. 2024. Toward sentiment aware semantic change analysis. In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop</i> , pages 350–357, St. Julian’s, Malta. Association for Computational Linguistics.	880
827		881
828		882
829		883
830		884
831		885
832		886
833	Nick Haslam. 2016. Concept creep: Psychology’s expanding concepts of harm and pathology. <i>Psychological Inquiry</i> , 27(1):1–17.	887
834		888
835		889
836	Nick Haslam, Ekaterina Vylomova, Michael Zyphur, and Yoshihisa Kashima. 2021. The cultural dynamics of concept creep. <i>American Psychologist</i> , 76(6):1013.	890
837		891
838		892
839		893
840	Shreya Havaldar, Bhumika Singhal, Sunny Rai, Langchen Liu, Sharath Chandra Guntuku, and Lyle Ungar. 2023. Multilingual language models are not multicultural: A case study in emotion. In <i>Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis</i> , pages 202–214, Toronto, Canada. Association for Computational Linguistics.	894
841		895
842		896
843		897
844		898
845		899
846		900
847		901
848	Simon Hengchen, Nina Tahmasebi, Dominik Schlechtweg, and Haim Dubossarsky. 2021. Challenges for computational lexical semantic change.	902
849		903
850		904
851		905
852	Katie Hoemann, Yeasle Lee, Èvelyne Dussault, Simon Devylder, Lyle H. Ungar, Dirk Geeraerts, and Batja Gomes de Mesquita. 2025. The construction of emotional meaning in language.	906
853		907
854		908
855		909
856	Joshua Conrad Jackson, Joseph Watts, Teague R. Henry, Johann-Mattis List, Robert Forkel, Peter J. Mucha, Simon J. Greenhill, Russell D. Gray, and Kristen A. Lindquist. 2019. Emotion semantics show both cultural variation and universal structure. <i>Science</i> , 366(6472):1517–1522.	910
857		911
858		912
859		913
860		914
861		915
862	Joshua Conrad Jackson, Joseph Watts, Johann-Mattis List, Curtis Puryear, Ryan Drabble, and Kristen A. Lindquist. 2022. From text to thought: How analyzing language can advance psychological science. <i>Perspectives on Psychological Science</i> , 17(3):805–826. PMID: 34606730.	916
863		917
864		918
865		919
866		920
867		921
868	Hajime Kiyama, Taichi Aida, Mamoru Komachi, Toshinobu Ogiso, Hiroya Takamura, and Daichi Mochihashi. 2025. Analyzing continuous semantic shifts with diachronic word similarity matrices. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 1613–1631, Abu Dhabi, UAE. Association for Computational Linguistics.	922
869		923
870		924
871		925
872		926
873		927
874		928
875	Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski, and Erik Velldal. 2018. Diachronic word embeddings and semantic shifts: a survey. In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 1384–1397, Santa Fe, New Mexico, USA. Association for Computational Linguistics.	929
876		930
877		931
878		932
879		933
880		934
881	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. <i>Bioinformatics</i> , 36(4):1234–1240.	935
882		936
883		937
884		938
885		939
886	Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinmeng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, et al. 2024. Best practices and lessons learned on synthetic data for language models. <i>arXiv preprint arXiv:2404.07503</i> .	940
887		941
888		942
889		943
890		944

891	Daniel Loureiro, Aminette D’Souza, Areej Nasser Muhajab, Isabella A. White, Gabriel Wong, Luis Espinosa-Anke, Leonardo Neves, Francesco Barbi- eri, and Jose Camacho-Collados. 2022. TempoWiC: An evaluation benchmark for detecting meaning shift in social media. In <i>Proceedings of the 29th Inter- national Conference on Computational Linguistics</i> , pages 3353–3359, Gyeongju, Republic of Korea. In- ternational Committee on Computational Linguistics.	947
892		948
893		949
894		950
895		951
896		952
897		953
898		
899		
900	Raphael Merx, Ekaterina Vylomova, and Kemal Kurni- awan. 2024. Generating bilingual example sentences with large language models as lexicography assis- tants. In <i>Proceedings of the 22nd Annual Workshop</i> <i>of the Australasian Language Technology Associa- tion</i> , Canberra, Australia. Association for Computa- tional Linguistics.	954
901		955
902		956
903		
904		
905		
906		
907	George A. Miller. 1992. WordNet: A lexical database for English. In <i>Speech and Natural Language: Pro- ceedings of a Workshop Held at Harriman, New York,</i> <i>February 23-26, 1992.</i>	957
908		958
909		959
910		960
911	Saif Mohammad. 2018. Obtaining reliable human rat- ings of valence, arousal, and dominance for 20,000 english words. In <i>Proceedings of the 56th annual</i> <i>meeting of the association for computational linguis- tics (volume 1: Long papers)</i> , pages 174–184.	961
912		962
913		
914		
915		
916	Saif M Mohammad. 2021. Sentiment analysis: Au- tomatically detecting valence, emotions, and other affectual states from text. In <i>Emotion measurement</i> , pages 323–379. Elsevier.	963
917		964
918		965
919		
920	Charles Egerton Osgood, William H May, and Murray S Miron. 1975. <i>Cross-Cultural Universals of Affective</i> <i>Meaning.</i> University of Illinois Press.	966
921		967
922		968
923	M Sandra Peña and Francisco José Ruiz de Mendoza. 2017. Construing and constructing hyperbole. <i>Stud- ies in figurative thought and language</i> , 56:41.	969
924		970
925		
926	Francesco Periti and Stefano Montanelli. 2024a. Lexi- cal semantic change through large language models: a survey. <i>ACM Comput. Surv.</i> , 56(11).	971
927		972
928		973
929	Francesco Periti and Stefano Montanelli. 2024b. Lexi- cal semantic change through large language models: a survey. <i>ACM Computing Surveys</i> .	974
930		
931		
932	Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for eval- uating context-sensitive meaning representations. In <i>Proceedings of the 2019 Conference of the North</i> <i>American Chapter of the Association for Computa- tional Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1267–1273, Minneapolis, Minnesota. Association for Computa- tional Linguistics.	975
933		976
934		977
935		978
936		979
937		
938		
939		
940		
941	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	991
942		992
943		993
944	James A. Russell. 2003. Core affect and the psychologi- cal construction of emotion. <i>Psychological Review</i> , 110(1):145–172.	994
945		995
946		
947	Dominik Schlechtweg, Barbara McGillivray, Simon Hengchen, Haim Dubossarsky, and Nina Tahmasebi. 2020. SemEval-2020 task 1: Unsupervised lexical semantic change detection. In <i>Proceedings of the</i> <i>Fourteenth Workshop on Semantic Evaluation</i> , pages 1–23, Barcelona (online). International Committee for Computational Linguistics.	996
948		997
949		998
950		999
951		
952		
953		
954	Nina Tahmasebi, Lars Borin, and Adam Jatowt. 2018. Survey of computational approaches to diachronic conceptual change. <i>CoRR</i> , abs/1811.06278.	996
955		997
956		998
957	Nina Tahmasebi and Thomas Risse. 2017. Finding indi- vidual word sense changes and their delay in appear- ance. In <i>Proceedings of the International Conference</i> <i>Recent Advances in Natural Language Processing,</i> <i>RANLP 2017</i> , pages 741–749, Varna, Bulgaria. IN- COMA Ltd.	999
958		
959		
960		
961	Xuri Tang. 2018. A state-of-the-art of semantic change computation. <i>Natural Language Engineering</i> , 24(5):649–676.	996
962		997
963		998
964		999
965		
966	Ekaterina Vylomova, Sean Murphy, and Nick Haslam. 2019. Evaluation of semantic change of harm-related concepts in psychology. In <i>Proceedings of the 1st Inter- national Workshop on Computational Approaches</i> <i>to Historical Language Change</i> , pages 29–34.	996
967		997
968		998
969		999
970		
971	Amy Beth Warriner, Victor Kuperman, and Marc Brys- baert. 2013. Norms of valence, arousal, and domi- nance for 13,915 english lemmas. <i>Behavior Research</i> <i>Methods</i> , 45(4):1191–1207.	996
972		997
973		998
974		999
975	Yu Xiao, Naomi Baes, Ekaterina Vylomova, and Nick Haslam. 2023. Have the concepts of ‘anxiety’ and ‘depression’ been normalized or pathologized? a cor- pus study of historical semantic change. <i>PloS one</i> , 18(6):e0288027.	996
976		997
977		998
978		999
979		
980	Heng Yang, Biqing Zeng, Mai Xu, and Tianxing Wang. 2021. Back to reality: Leveraging pattern-driven modeling to enable affordable sentiment dependency learning. <i>CoRR</i> , abs/2110.08604.	996
981		997
982		998
983		999
984	Heng Yang, Chen Zhang, and Ke Li. 2023. Pyabsa: A modularized framework for reproducible aspect- based sentiment analysis. In <i>Proceedings of the 32nd</i> <i>ACM International Conference on Information and</i> <i>Knowledge Management, CIKM 2023, Birmingham,</i> <i>United Kingdom, October 21-25, 2023</i> , pages 5117– 5122. ACM.	996
985		997
986		998
987		999
988		
989		
990		
991	Yuxiang Zhou, Jiazheng Li, Yanzheng Xiang, Hanqi Yan, Lin Gui, and Yulan He. 2023. The mystery and fascination of llms: A comprehensive survey on the interpretation and analysis of emergent abilities. <i>arXiv preprint arXiv:2311.00237.</i>	996
992		997
993		998
994		999
995		
996	Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can large lan- guage models transform computational social sci- ence? <i>Computational Linguistics</i> , 50(1):237–291.	996
997		997
998		998
999		999

A Corpus Counts of Target Terms

1000

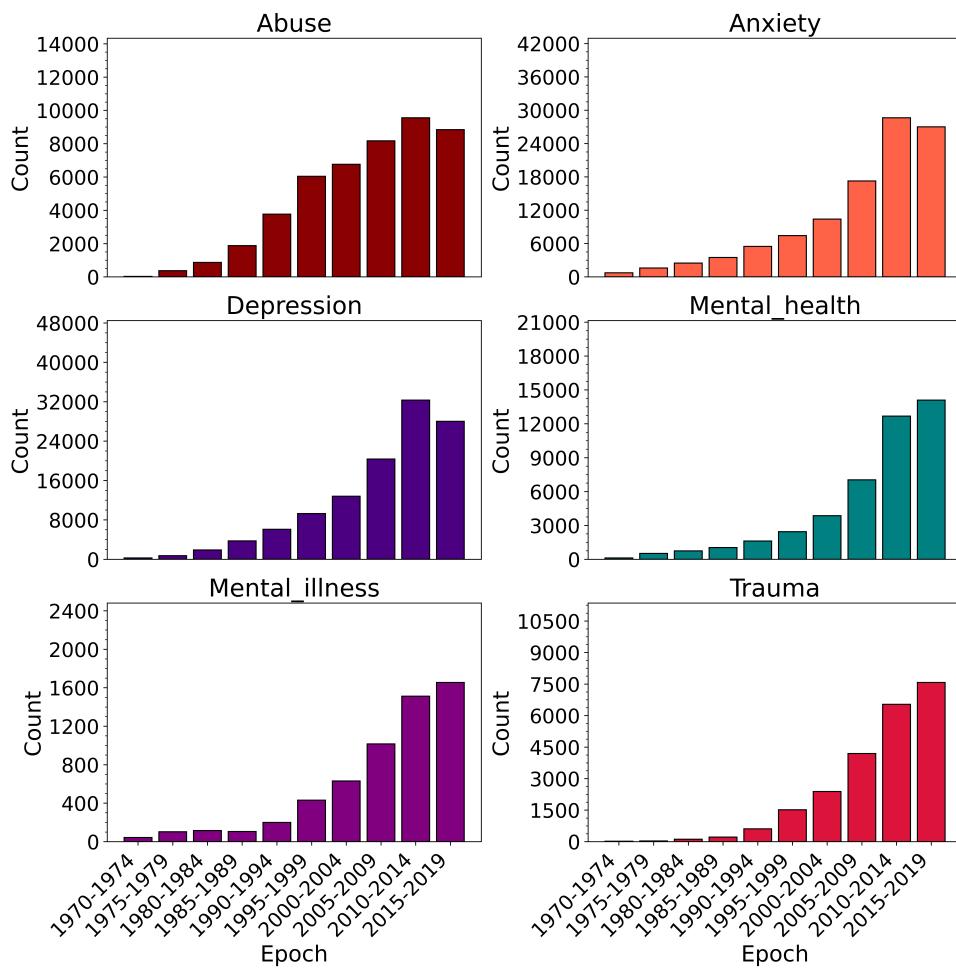


Figure 4: Annual Counts of Sentences where Target Terms Appear in the Psychology Corpus (1970-2019).

B Synthetic Dimension Datasets: Details

Dimension	Target	Neutral (M)	Increase (M)	Decrease (M)	US\$
<i>Sentiment</i>	Abuse	5,645 (28)	5,645 (30)	5,645 (29)	17
	Anxiety	9,215 (27)	9,213 (28)	9,213 (28)	28
	Depression	8,828 (27)	8,826 (28)	8,826 (28)	29
	Mental Health	6,348 (28)	6,348 (29)	6,348 (29)	21
	Mental Illness	2,552 (28)	2,552 (28)	2,552 (29)	9
	Trauma	3,563 (28)	3,563 (30)	3,563 (30)	11
<i>Intensity</i>	Abuse	6,802 (28)	6,801 (30)	6,801 (29)	21
	Anxiety	9,659 (26)	9,657 (29)	9,657 (28)	32
	Depression	10,022 (27)	10,020 (30)	10,020 (29)	35
	Mental Health	6,904 (28)	6,899 (32)	6,899 (29)	24
	Mental Illness	2,497 (28)	2,496 (32)	2,496 (29)	10
	Trauma	4,012 (28)	4,012 (30)	4,012 (30)	14
<i>Breadth</i>	Abuse	NA	5,221 (27)	NA	0
	Anxiety	NA	13,635 (26)	NA	0
	Depression	NA	14,463 (27)	NA	0
	Mental Health	NA	14,638 (26)	NA	0
	Mental Illness	NA	14,639 (26)	NA	0
	Trauma	NA	14,650 (26)	NA	0

Table 3: Descriptives for Synthetic Dimension Datasets: Sentence Counts, Sentence Lengths, and Total Generation Cost.

M = Mean Sentence Length of Dataset. Neutral = Neutral, unaltered, input sentences. Increase = Increase on the Dimension of interest. Decrease = Decrease on the Dimension of interest.

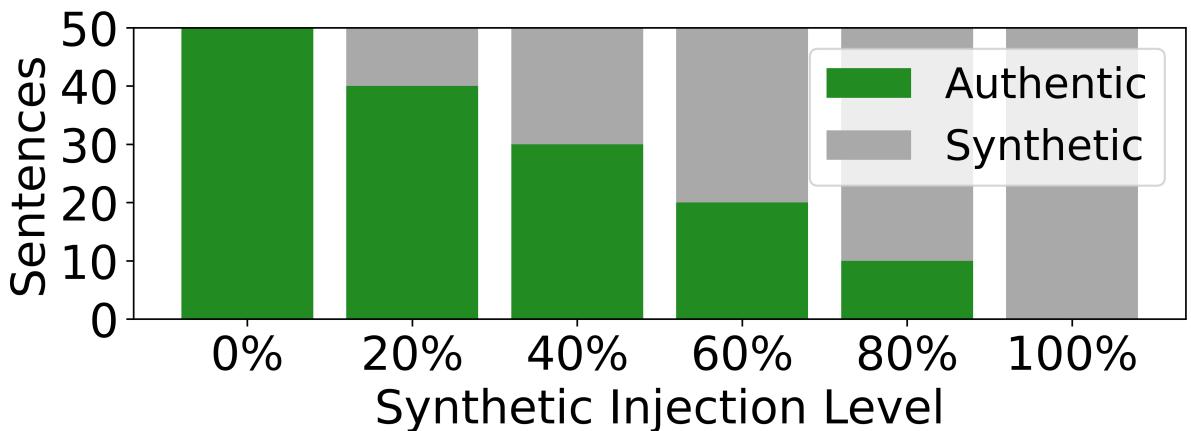


Figure 5: Distribution of Sentences in Each Sample of 50 Sentences.

Dimension	Target	Neutral	Increased Variation	Decreased Variation
Sentiment	Abuse	Child abuse is not a single faceted phenomenon.	Child abuse is a deeply complex phenomenon that can spur important dialogues and reforms.	Child abuse is a multifaceted atrocity with far-reaching and damaging consequences.
	Anxiety	Typical worship reinforces pathologies of anxiety and self-deception.	Typical worship empowers resilience in the face of anxiety and self-deception.	Typical worship deepens the pathologies of anxiety and self-deception.
	Depression	The expression masked depression is not a lucky one.	The expression masked depression may offer an insightful perspective.	The expression masked depression is unfortunately an unsettling one.
	Mental Health	Two views of holiness and its bearing on mental_health are discussed.	Two perspectives on holiness and its supportive impact on mental_health are discussed.	Two views of holiness and its potential pressure on mental_health are discussed.
	Mental Illness	The results suggest that physical or mental_illness may decrease creativity.	The results suggest that overcoming physical or mental_illness may lead to increased creativity.	The results suggest that physical or mental_illness may significantly hinder creativity.
	Trauma	Psychic trauma interferes with the normal structuring of experience.	Psychic trauma challenges individuals in a way that can lead to the reorganization and enrichment of their experience.	Psychic trauma disrupts and fragments the normal structuring of experience.
Intensity	Abuse	Theorists and practitioners alike believe that emotional abuse exists.	Theorists and practitioners alike fervently believe that pervasive emotional abuse exists.	Theorists and practitioners alike casually believe that subtle emotional abuse exists.
	Anxiety	Teacher reported anxiety was related to worse time production.	Teacher reported severe anxiety was related to significantly worse time production.	Teacher reported mild anxiety was related to slightly worse time production.
	Depression	Maternal depression continues to play a role in children's development beyond infancy.	Severe maternal depression continues to play a profound role in children's development beyond infancy.	Mild maternal depression continues to play a subtle role in children's development beyond infancy.
	Mental Health	Eveningness is related to negative physical and mental_health outcomes.	Eveningness is alarmingly related to severe negative physical and troubling mental_health outcomes.	Eveningness is mildly related to some negative physical and mental_health outcomes.

Continued on next page

Dimension	Target	Neutral	Increased Variation	Decreased Variation
Breadth	Mental Illness	Biblical and theological considerations underline the importance of the problem about mental illness , but do not provide a solution.	Biblical and theological considerations underline the immense importance and complexity of the problem about mental illness , but do not provide a definitive solution.	Biblical and theological considerations highlight the importance of the issue regarding mental illness , but do not provide a clear solution.
	Trauma	Childhood trauma is a key risk factor for psychopathology.	Childhood trauma is a critical and devastating risk factor for severe psychopathology.	Childhood trauma is a notable but moderate risk factor for mild psychopathology.
	Abuse	Sexual exploitation is an expression of a power relationship.	Sexual abuse is an expression of a power relationship.	NA
	Anxiety	Adolescents' state of mind with regard to attachment and representations regarding separation were examined.	Adolescents' anxiety with regard to attachment and representations regarding separation were examined.	NA
	Depression	Iranian college students showed more anxiety than their British peers.	Iranian college students showed more depression than their British peers.	NA
	Mental Health	Such a scale may alert clinicians early in treatment to issues related to trauma	Such a scale may alert clinicians early in treatment to issues related to mental health	NA
	Mental Illness	Excessive estrogen influence produces anxiety, agitation , irritability, and lability.	Excessive estrogen influence produces anxiety, mental illness , irritability, and lability.	NA
	Trauma	Further investigation of pathological dissociation in Hong Kong is necessary.	Further investigation of pathological trauma in Hong Kong is necessary.	NA

Table 5: Sample of Short Synthetic Sentences from the Synthetic Datasets for each Target term.

1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015

C In-Context Learning Paradigm

The study generated synthetic datasets to simulate changes in Sentiment and Intensity using 36,151 and 39,896 neutral baseline sentences, respectively. Neutral sentences were sampled by linking words in each sentence with their mean valence or arousal scores from the NRC-VAD lexicon (0-1) (Mohammad, 2018) and filtering by a dynamic range. This neutral range is adjusted from the median of each dataset by ± 0.01 , targeting 25th-75th percentile bounds or 500-1500 unique sentences per epoch. See Figures 6 and 7 for a breakdown of neutral sentence counts per epoch provided as input to the LLM using the prompts below.

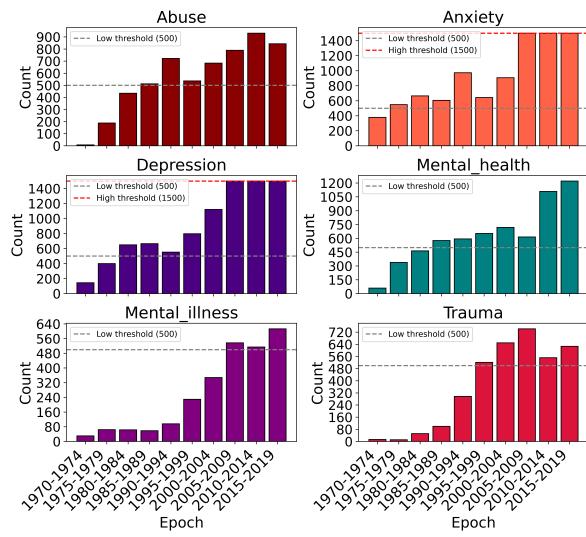


Figure 6: Counts of Neutral Sentences (valence Scores).

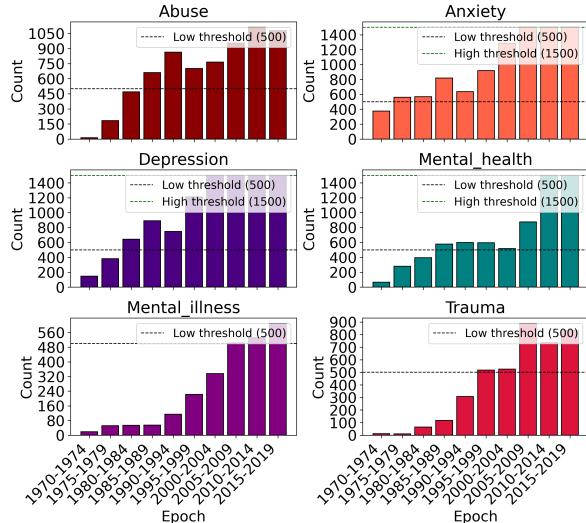


Figure 7: Counts of Neutral Sentences (Arousal Scores).

For each neutral sentence, one inference call to GPT-4o is made through the OpenAI API to generate variations of increased and decreased Sentiment

or Intensity. Only the samples for *anxiety* and *depression* reached the upper limit of 1,500 sentences for the final three epochs, while other targets did not exceed 500 sentences per epoch (allowing for unique sentences across each of the 10 iterations of up to 50 unique sentences). The sentence generation prioritized quality and maintained a neutral baseline to allow for adequate variation.

The ChatGPT API with a temperature setting of 1.00 was used to ensure semantic accuracy and prevent errors (Periti and Montanelli, 2024b), while allowing for a balance between deterministic and creative responses. Note that there were challenges in maintaining target terms in the sentences, particularly for positive sentiment variations. Fewer manual adjustments were needed for Intensity than Sentiment. GPT-4o struggled to vary 97% of the sentences to contain more positive sentiment for *abuse* (28), *anxiety* (110), *depression* (46), *mental_health* (1), *trauma* (2) as it replaced targets with positive terminology against instructions. For Intensity data, fewer sentences required manual alteration: only for *abuse* (4), *depression* (2), *mental_health* (2), *trauma* (1). Rows (196) were detected and manually altered to retain the target term while ensuring variation in the dimension relative to the neutral sentence. The final validated datasets, detailed in Table 6, are available on the GitHub repository: [MASKED LINK].

Target	Dimension	Neutral	Increase	Decrease	US\$
abuse	Sentiment	5,645	5,645	5,645	17
	Intensity	6,802	6,801	6,801	21
anxiety	Sentiment	9,215	9,213	9,213	28
	Intensity	9,659	9,657	9,657	32
depression	Sentiment	8,828	8,826	8,826	29
	Intensity	10,022	10,020	10,020	35
mental_health	Sentiment	6,348	6,348	6,348	21
	Intensity	6,904	6,899	6,899	24
mental_illness	Sentiment	2,552	2,552	2,552	9
	Intensity	2,497	2,496	2,496	10
trauma	Sentiment	3,563	3,563	3,563	11
	Intensity	4,012	4,012	4,012	14

Table 6: Sentence counts and Cost for Synthetic Sentiment and Intensity Datasets.

Prompt for Synthetic Sentiment

PROMPT_INTRO = """ In psychology research, ‘Sentiment’ is defined as “a term’s acquisition of a more positive or negative connotation.” This task focuses on the sentiment of the term **«target_word»**.

Task

You will be given a sentence containing the term **«target_word»**. Your goal is to write two new sentences:

1. One where **«target_word»** has a **more positive connotation** (enclose this sentence between ‘<positive target_word>’ and ‘</positive target_word>’ tags).
2. One where **«target_word»** has a **more negative connotation** (enclose this sentence between ‘<negative target_word>’ and ‘</negative target_word>’ tags).

Rules

1. The term **«target_word»** must remain **exactly as it appears** in the original sentence:
 - Do **not** replace, rephrase, omit, or modify it in any way.
 - Synonyms, variations, or altered spellings are not allowed.
2. **Meaning and Structure**:
 - Stay true to the original context and subject matter.
 - Maintain the sentence’s structure and ensure grammatical accuracy.
3. **Sentiment Adjustments**:
 - **Positive Sentiment**: Reflect strengths or benefits realistically, while respecting the potential negativity of **«target_word»**.
 - **Negative Sentiment**: Highlight risks or harms appropriately, avoiding exaggeration or trivialization.

Important

- Any response omitting, replacing, or altering **«target_word»** will be rejected.
- Ensure the output is:
 - **Grammatically correct**
 - **Sensitive and serious** in tone
 - **Free from exaggeration or sensationalism**
 - **Strictly following the XML-like tag format for sentiment variations**

Follow these guidelines strictly to produce valid responses. """

Prompt for Synthetic Intensity

PROMPT_INTRO = """In psychology research, Intensity is defined as “the degree to which a word has emotionally charged (i.e., strong, potent, high-arousal) connotations.” This task focuses on the intensity of the term **«target_word»**.

Task

You will be given a sentence containing the term **«target_word»**. Your goal is to write two new sentences:

1. One where **«target_word»** is **less intense** (enclose this sentence between ‘<decreased target_word intensity>’ and ‘</decreased target_word intensity>’ tags).
2. One where **«target_word»** is **more intense** (enclose this sentence between ‘<increased target_word intensity>’ and ‘</increased target_word intensity>’ tags).

Rules

1. The term **«target_word»** must remain **exactly as it appears** in the original sentence:
 - Do **not** replace, rephrase, omit, or modify it in any way.
 - Synonyms, variations, or altered spellings are not allowed.
2. **Meaning and Structure**:
 - Stay true to the original context and subject matter.
 - Maintain the sentence’s structure and ensure grammatical accuracy.

Important

- Any response omitting, replacing, or altering **«target_word»** will be rejected.
- Ensure the output is:
 - **Grammatically correct**
 - **Sensitive and serious** in tone
 - **Free from exaggeration or sensationalism**
 - **Strictly following the XML-like tag format for intensity variations**

Follow these guidelines strictly to produce valid responses. """

1049

D Demonstration Examples: Synthetic Sentiment

1050

Target	Neutral	Positive Sentiment	Negative Sentiment
Abuse	Child abuse is most likely to occur when socially isolated parents react impulsively to aversive stimuli emitted by their children.	Child abuse is less likely to occur when socially isolated parents respond lovingly to their children's behavior.	Child abuse is most likely to occur when socially isolated parents react aggressively to their children's challenging behavior.
Abuse	The children represented a wide spectrum of sexual abuse .	The children represented a meaningful spectrum of sexual abuse .	The children represented a devastating spectrum of sexual abuse .
Abuse	Euphoric properties of cocaine lead to the development of chronic abuse , and appear to involve the acute activation of central DA neuronal systems.	Euphoric properties of cocaine lead to the growth of chronic abuse , and appear to involve the acute activation of central DA pleasure systems.	Emotional properties of cocaine lead to the decline into chronic abuse , and appear to involve the acute activation of central DA pain systems.
Abuse	Substance abuse helps the individual deal with distress associated with family interactions.	Substance abuse helps the individual temporarily cope positively with family interactions.	Substance abuse makes the individual endure the overwhelming pain and alienation associated with family interactions.
Abuse	The study determined that 84 of the sample reported a history of abuse or neglect.	The study determined that 84 of the sample acknowledged a transformative history of overcoming abuse or neglect.	The study determined that 84 of the sample complained of a miserable history of abuse or neglect.
Anxiety	Previous work suggests that social anxiety is inconsistently related to alcohol use.	Previous work agrees that social anxiety is sometimes related to alcohol use.	Previous work warns that social anxiety is unpredictably related to alcohol use.
Anxiety	A small yet emerging body of research on the relationship between anxiety and driving suggests that higher levels of state anxiety may lead to more dangerous driving behaviors.	A small yet emerging body of research on the positive relationship between anxiety and driving suggests that higher levels of state anxiety may lead to more daring driving behaviors.	A small yet emerging body of research on the problematic relationship between anxiety and driving suggests that more disturbing levels of state anxiety may lead to more disastrous driving behaviors.
Anxiety	Findings suggest that individuals high in anxiety show greater contextual fear generalization as measured by US expectancy.	Findings suggest that individuals high in anxiety show greater contextual concern generalization as measured by US hope.	Findings suggest that individuals high in anxiety show greater contextual terror generalization as measured by US dread.
Anxiety	General anxiety and evoked imagery of death as a person were measured in 75 male Catholic college students and seminarians.	General anxiety and vivid imagery of hope as a person were measured in 75 male Catholic college students and seminarians.	General anxiety and frightening imagery of death as a person were measured in 75 male Catholic college students and seminarians.

61

Continued on next page

Target	Neutral	Positive Sentiment	Negative Sentiment
Anxiety	Results indicated that emotion dysregulation significantly mediated the relationship between child abuse severity and attachment-related anxiety and avoidance.	Results indicated that emotion variation positively mediated the relationship between childhood experiences and attachment-related anxiety and care.	Results indicated that emotion disturbance problematically mediated the relationship between child abuse severity and attachment-related anxiety and terror.
Depression	The present study was conducted to test predictions derived from the hypothesis that depression may serve the purpose of adaptively facilitating disengagement from obsolete cognitive plans.	The present study was conducted to test predictions derived from the hypothesis that depression may serve the purpose of helping people make better cognitive plans.	The present study was conducted to test predictions derived from the hypothesis that depression may prevent people from carrying out destructive cognitive plans.
Depression	Vision loss was a consistent predictor of both onset and persistence of depression , even after a wide range of covariates had been adjusted.	Vision loss was a positive predictor of both beginning and retaining depression , even after a wide range of covariates had been included.	Vision loss was an unavoidable predictor of both suffering and enduring depression , even after a wide range of covariates had been controlled.
Depression	This study examined whether distinct groups of young adolescents with mainly anxiety or mainly depression could be identified in a general population sample.	This study examined whether unique groups of young adolescents with mainly vigilance or mainly depression could be identified in a general population sample.	This study examined whether pathological groups of young adolescents with mainly fear or mainly depression could be isolated in a general population sample.
Depression	In most people with recurrent depression , mindfulness skills are expressed evenly across different domains.	In most people who live with depression , mindfulness skills are expressed in a balanced way across different domains.	In most people who struggle with untreatable depression , mindfulness habits are expressed monotonously across different domains.
Depression	The aim of the study was to test the effect of differing information regarding the rationale given to participants for a study on depression symptoms.	The hope of the study was to test the effect of diverse information regarding the clarifying reasons bestowed on participants for an exploration of depression features.	The aim of the study was to test the effect of differing information regarding the dreary explanation given to participants for a study on depression pathologies.
Mental Health	This paper maintains that mental_health delivery systems must be supplemented by critical analyses of the hidden assumptions that guide policy and technique decisions.	This paper hopes that mental_health delivery systems must be improved by enlightened analyses of the hidden assumptions that lead beneficial policy and technique decisions.	This paper warns that mental_health delivery systems must be supplemented by harsh analyses of the deep-seated errors that undermine policy and technique decisions.

Continued on next page

Target	Neutral	Positive Sentiment	Negative Sentiment
Mental Health	The federal regulations governing confidentiality of alcohol and drug abuse patient records are examined with respect to their applicability to mental_health and other medical records.	The federal regulations protecting confidentiality of alcohol and drug use records are examined with respect to their applicability to mental_health and other well-being records.	The federal regulations restricting access to alcohol and drug abuse patient records are examined with respect to their potential shortcomings for mental_health and other medical records.
Mental Health	Young people are particularly vulnerable to unemployment and the consequences of this for psychosocial development and mental_health are not well understood.	Young people are particularly responsive to leisure and the consequences of this for psychosocial well-being and mental_health will benefit from more understanding.	Young people are particularly vulnerable to unemployment and the threats of this for dysfunction and mental_health are poorly understood.
Mental Health	This study suggests that the long-term outcome in schizophrenic patients followed by a community-based mental_health service is generally poor and multifaceted.	This study suggests that the long-term improvement in people with schizophrenia followed by a community-based mental_health service is generally variable.	This study warns that the long-term outcome in schizophrenic patients followed by a community-based mental_health clinical is generally poor and incoherent.
Mental Health	The stigma of having psychological problems is a barrier to seeking mental_health treatment, but little research has examined whether this stigma influences the experiences of those in treatment.	The public image of having well-being challenges is a bridge to seeking mental_health help, but little research has examined whether this image influences the experiences of those in care.	The shame of having psychological illness is an obstacle to seeking mental_health treatment, but little research has examined whether this shame increases the misery of those in treatment.
Mental Illness	Internet addiction (IA) is an emerging social and mental_health issue among youths.	Internet engagement (IE) is a rising social and mental_health issue among youths.	Internet addiction (IA) is a looming social and mental_health disorder among youths.
Mental Illness	Second, we asked to what extent suicides of older mentally ill persons are definitely created by their mental_illness .	Second, we asked to what extent suicides of older persons are definitely created by their mental_illness .	Second, we asked to what extent suicides of older mentally ill persons are definitely made worse by their mental_illness .
Mental Illness	It was found that rejection of the mentally ill in situations of social relations was linked to prior personal experience with mental_illness , perceived dangerousness of the mentally ill, and age of the survey respondent.	It was found that welcoming of people in situations of social relations was linked to prior positive personal experience with mental_illness , perceived safety of these people, and age of the survey respondent.	It was found that rejection of the mentally ill in situations of social relations was linked to negative prior personal experience with mental_illness , perceived dangerousness of the mentally ill, and age of the survey respondent.
Mental Illness	In over 50 of cases continuation of in-patient stay was necessitated by the severity of mental_illness .	In over 50 of cases continuation of stay in care was necessitated by the level of mental_illness .	In over 50 of cases being restricted to hospital was necessitated by the severity of mental_illness .

Continued on next page

Target	Neutral	Positive Sentiment	Negative Sentiment
Mental Illness	Much controversy exists over the treatment of mental illness and many critics argue that the exercise of medical authority results in the social control of the mentally ill.	Much conversation exists over the care of mental illness and many writers argue that the medical authorities enhance the social enhancement of mental health.	Much disagreement exists over the treatment of mental illness and many critics argue that the abuse of medical tyranny results in the domination of the mentally ill.
Trauma	This paper presents a cognitive-behavioral model for conceptualizing and intervening in the area of sexual trauma .	This paper celebrates a cognitive-behavioral model for promoting new ideas and helping in the area of sexual trauma .	This paper presents a cognitive-behavioral model for thinking about and wrestling with the harmful problem of sexual trauma .
Trauma	In most classrooms in most schools, there are students who have suffered complex trauma who would benefit from a system-wide, trauma-informed approach to schooling.	In most classrooms in most schools, there are students who have experienced complex trauma who would benefit from a system-wide, responsive and enlightened approach to schooling.	In most classrooms in most schools, there are students who have suffered damaging trauma whose problems need a system-wide, illness-based approach to schooling.
Trauma	Research has shown that women are more likely to develop PTSD subsequent to trauma exposure in comparison with men.	Research has shown that women are more likely to develop PTSD subsequent to trauma experiences in comparison with men.	Research has shown that women are more likely to deteriorate into PTSD subsequent to trauma exposure in comparison with men.
Trauma	Numerous homeless youth experience trauma prior to leaving home and while on the street.	Numerous resilient youth learn to navigate trauma prior to leaving home and while adapting to life on the street.	Numerous homeless youth endure significant trauma prior to leaving home and while facing severe challenges on the street.
Trauma	The meaning of trauma within psychology has for a long time been viewed mostly from a pathologizing standpoint.	The meaning of trauma within psychology has for a long time needed to be viewed from a more compassionate and strengths-based standpoint.	The meaning of trauma within psychology has for a long time been viewed mostly from a negative and overly disease-focused standpoint.

Table 8: Expert Crafted Sentiment Variations for Neutral Sentences for inference calls to GPT-4o for the Few-Shot ICL Paradigm.

E Demonstration Examples: Synthetic Intensity

1051

Target	Neutral	High Intensity	Low Intensity
Abuse	Clinically, however, individual questions that use broad labeling terms are more likely to identify women as having a history of abuse .	Clinically, however, individual questions that use extreme labeling terms are more likely to reveal women as having a severe history of abuse .	Clinically, however, individual questions that use broad labeling terms are more likely to identify women as having a mild history of abuse .
Abuse	Most care workers said that they would be willing to report abuse anonymously.	Most care workers cried that they would be delighted to report extreme instances of abuse anonymously.	Most care workers said that they would be willing to report trivial abuse anonymously.
Abuse	There is greater emphasis on recognizing that older people may be subjected to abuse and neglect by family members and the community as well.	There is a significant emphasis on recognizing that older people may be subjected to severe abuse and appalling neglect by family members and the community as well.	There is some emphasis on recognizing that older people may experience weak abuse by family members and the community as well.
Abuse	Education on financial abuse for both elders and their adult children and establishment of income support programs are urgently needed.	Education on ordinary financial abuse for both elders and their adult children and urgent establishment of income support programs are desperately needed.	Education on financial abuse for both elders and their adult children and establishment of income support programs will occur.
Abuse	There was no association between physical abuse and depressive symptoms through either self-compassion or gratitude.	There was no association between frightening physical abuse and cold symptoms through either emotional contagion or extreme gratitude.	There was no association between mild physical abuse and state of mind through either complacency or gratitude.
Anxiety	The spread of anxiety as seen in curves of generalization seems greater at the unconscious than at the conscious level.	The uncontrollable spread of intense anxiety as seen in spikes of generalization seems more vivid at the unconscious than at the conscious level.	The spread of mild anxiety as seen in curves of generalization seems greater at the unconscious than at the conscious level.
Anxiety	These findings suggest that two important factors to be considered by researchers, educators, and mental_health professionals are adults' perceptions of their fathers' level of acceptance-rejection and the amount of anxiety they experience in their relationship with God.	These findings cry out that two powerful factors to be considered by researchers, educators, and mental_health professionals are adults' perceptions of their fathers' extreme level of rejection and the intense amount of anxiety they experience in their relationship with God.	These findings suggest that two important factors to be considered by researchers, educators, and other professionals are adults' perceptions of their fathers' level of acceptance and the amount of mild anxiety they experience in their relationship with God.

Target	Neutral	High Intensity	Low Intensity
Anxiety	Self-compassion might be an alternative strategy for cognitive reappraisal in the management of shame-proneness and social anxiety .	Emotion exaggeration might be an alternative strategy for overcoming upset in the management of shame and extreme social anxiety .	Meditation might be an alternative strategy for cognitive reappraisal in the management of boredom and mild social anxiety .
Anxiety	The chronic anxiety level of the subject may be related to the ease of acquisition and spread of new anxiety responses.	The intense anxiety level of the subject may be related to the ease of acquisition and catastrophic spread of extreme anxiety responses.	The mild anxiety level of the subject may be related to the ease of acquisition and generalization of new responses.
Anxiety	Results indicated that greater attachment anxiety and avoidance were linked to lower levels of life satisfaction in both gay men and lesbians.	Results cried out that extreme attachment anxiety and avoidance were linked to desperate levels of life misery in both gay men and lesbians.	Results indicated that attachment anxiety and peacefulness were linked to lower levels of life satisfaction in both gay men and lesbians.
Depression	A combined medical and psychiatric treatment of a depression consequent to a colostomy and an organic impotence following rectal resection for cancer in a 33-year-old man has been described.	A combined medical and psychiatric treatment of an intense depression consequent to a colostomy and a severe organic impotence following surgical rectal tissue destruction for cancer in a 33-year-old man has been described.	A combined medical and psychiatric treatment of a mild depression consequent to a colostomy and an organic impotence following rectal resection for cancer in a 33-year-old man has been described.
Depression	A 35-year-old woman had a history of increasing irritability and liability to attacks of depression related to a complete inability to have coital orgasms.	A 35-year-old woman had a fearsome history of crescendoing irritability and liability to severe attacks of depression related to a horrendous inability to have coital orgasms.	A 35-year-old woman had a history of sleepiness and liability to periods of mild depression related to an inability to have coital orgasms.
Depression	During acute asthma these appear to be radically altered into sadness and longing, and subjected to generalized inhibition similar to that seen in states of depression .	During severe, life-threatening asthma episodes these appear to be radically altered into intense misery, and subjected to generalized inhibition similar to that seen in states of extreme depression .	During asthma these appear to be altered into boredom and tiredness, and subjected to generalized inhibition similar to that seen in states of low-level depression .
Depression	Differences in response in the same individual seem related to mood and attitude as well as to transient stress, with the response being lower on days of depression .	Scary differences in response in the same individual seem related to intense mood and attitude as well as to sudden stress, with the emotional response being more intense on days of destructive depression .	Predictable differences in response in the same individual seem related to mood, attitude and life experiences, with the subdued response being mild on days of everyday depression .

Continued on next page

Target	Neutral	High Intensity	Low Intensity
Depression	The depression was treated by the introduction of behaviors incompatible with the depression .	The intense depression was treated by the shocking introduction of uncontrollable behaviors incompatible with the severe depression .	The mild depression was treated by the introduction of behaviors incompatible with it.
Mental Health	Community mental_health espouses an innovative conception for psychological services in the university community.	Community mental_health fights for a divisive conception for psychological services in the overwhelmed university community.	Community mental_health espouses a dull conception for services in the university community.
Mental Health	We also opine that if restraints are misused by mental_health or child welfare treatment settings, then their misuse may be considered a subject of a patient maltreatment, abuse, criminal or civil action.	We also exclaim that if harsh restraints are abused by mental_health or child welfare treatment settings, then their damaging misuse may be criticized as a subject of extreme patient maltreatment, abuse, criminal or civil action.	We also state that if restraints are used by mental_health or child welfare treatment settings, then they may be considered a subject of a discussion.
Mental Health	This research is a secondary data analysis of the impact of adolescents' mental/substance-use disorders and dual diagnosis on their utilization of drug treatment and mental_health services.	This research is an intense data analysis of the terrible impact of adolescents' mental/substance abuse disorders and severe compounding problems on their abuse of drug treatment and mental_health services.	This research is a data analysis of the impact of adolescents' experiences on their utilization of normal treatment and mental_health services.
Mental Health	The findings emphasize the need for family-based treatment for CP that addresses parent behaviors and adolescent mental_health .	The findings make a heartfelt plea for the desperate need for family-based treatment for CP that challenges destructive parent behaviors and adolescent mental_health diseases.	The findings summarize the need for family-based treatment for CP that addresses ordinary parent behaviors and mild adolescent mental_health .
Mental Health	Our findings suggest that maternal mental_health influences child sleep behavior at 18 months after birth, and not vice versa.	Our exciting findings suggest that damaged maternal mental_health destructively influences child sleep behavior at 18 months after birth, and not vice versa.	Our findings suggest that ordinary maternal mental_health influences child normal sleep behavior at 18 months after birth, and not vice versa.
Mental Illness	Problems of definition and classification in psychiatry and the impact of mental_illness on the individual and the community pose unique problems for psychiatric register studies.	Horrible problems of definition and classification in psychiatry and the harsh impact of severe mental_illness on the individual and the community pose frightening problems for psychiatric register studies.	Issues of definition and classification in psychiatry and the impact of mild mental_illness on the individual and the community arise in register studies.

Continued on next page

Target	Neutral	High Intensity	Low Intensity
Mental Illness	In parents and collateral relatives of the autistic children, 3.2% had a serious mental illness , and 4.8% of siblings were markedly abnormal.	In desperate parents and relatives of the severely autistic children, 3.2% had a serious mental illness , and 4.8% of siblings were extremely abnormal.	In parents and relatives of the mildly autistic children, 3.2% had an ordinary mental illness , and 4.8% of siblings were normal.
Mental Illness	Consistent with genetic essentialism, genetic attributions increased the perceived seriousness and persistence of the mental illness and the belief that siblings and children would develop the same problem.	Consistent with the horrors of genetic essentialism, genetic attributions exaggerated the perceived severity and uncontrollability of the severe mental illness and the destructive belief that siblings and children would develop the same extreme problem.	Consistent with genetic essentialism, genetic attributions influenced views about the mental illness and the belief that siblings and children would develop it.
Mental Illness	The target population was urban, homeless, HIV+ individuals with substance dependence and/or mental illness diagnoses.	The completely overwhelmed target population was urban, homeless, HIV+ individuals with severe substance abuse and/or unmanageable mental illness diagnoses.	The target population was urban, ambulatory, healthy individuals with mild mental illness diagnoses.
Mental Illness	Doctors, including general practitioners, experience higher levels of mental illness than the general population.	Doctors, including general practitioners, experience higher levels of mental illness than the general population.	Doctors, including general practitioners, experience higher levels of mental illness than the general population.
Trauma	They tend to be more liberal in their attitudes toward abortion than women in general; however, women who experienced a greater degree of psychic trauma tended to be more conservative in their attitudes.	They tend to be more extremely callous in their attitudes toward the horrors of abortion than women in general; however, women who suffered a greater degree of violent psychic trauma tended to be more fearful in their attitudes.	They tend to be more accepting in their attitudes toward children than women in general; however, women who experienced mild psychic trauma tended to be more conservative in their attitudes.
Trauma	The trauma was overwhelming.	The intense trauma was completely overwhelming.	The mild trauma was unproblematic.
Trauma	The choice of defensive style was found related to at least three factors: an early history of trauma , especially separation, parental encouragement of toughness, and essentially a counterphobic family style.	The choice of emotional overreaction was found related to at least three factors: an early history of extreme trauma , especially harsh abandonment, parental punishment, and essentially an emotionally destructive family style.	The choice of coping style was found related to at least three factors: an early history of mild trauma , especially independence, parental encouragement, and essentially a dull and normal family style.

Continued on next page

Target	Neutral	High Intensity	Low Intensity
Trauma	It is an attempt to bring the trauma arising from the external world into the internal world and thus to create an illusion of mastery and control.	It is a desperate attempt to bring the unbearable trauma threatening from the external world into the internal world and thus to create a poisonous illusion of mastery and control.	It is an attempt to bring the mild trauma arising from the external world into the internal world and thus to create a sense of peace and tranquillity.
Trauma	The international standard for setting ski bindings is based on the measurement of the tibia proximal width because of the propensity of this bone to suffer trauma as the ski and skier attempt to go in different directions.	The disgraceful international standard for setting ski bindings is based on the measurement of the tibia proximal width because of the scary propensity of this bone to suffer severe trauma as the ski and skier attempt to go in different directions.	The international standard for setting ski bindings is based on the measurement of the tibia proximal width because of the propensity of this bone to experience mild trauma as the ski and skier attempt to go in different directions.

Table 10: Expert Crafted Intensity Variations for Neutral Sentences for inference calls to GPT-4o for the Few-Shot ICL Paradigm.

F List of Donor Terms: Synthetic Breadth

Target (Synset)	Sibling (Synset)	Lin Similarity	Cosine Similarity
Abuse (abuse.n.02)	Disparagement (disparagement.n.01)	1.54	0.89
	Contempt (contempt.n.03)	1.49	0.86
	Impudence (impudence.n.01)	1.47	0.84
	Ridicule (ridicule.n.01)	1.34	0.91
	Derision (derision.n.01)	1.24	0.81
Abuse (maltreatment.n.01)	Blasphemy (blasphemy.n.01)	1.07	0.89
	Exploitation (exploitation.n.02)	1.78	0.86
	Disregard (disregard.n.02)	1.67	0.82
	Harassment (harassment.n.02)	1.55	0.84
	Annoyance (annoyance.n.05)	1.37	0.83
Anxiety (anxiety.n.01)	Depression (depression.n.01)	2.09	0.91
	Mental Health (mental_health.n.01)	1.85	0.89
	Trauma (trauma.n.02)	1.70	0.90
	Mental Illness (mental_illness.n.01)	1.60	0.92
	Dissociation (dissociation.n.02)	1.55	0.90
	Hypnosis (hypnosis.n.01)	1.43	0.89
	Delusion (delusion.n.01)	1.42	0.89
	Anhedonia (anhedonia.n.01)	1.33	0.84
	Agitation (agitation.n.01)	1.31	0.91
	Depersonalization (depersonalization.n.02)	1.31	0.90
	Irritation (irritation.n.01)	1.26	0.89
	Morale (morale.n.01)	1.26	0.89
	Nervousness (nervousness.n.02)	1.24	0.84
	Enchantment (enchantment.n.02)	1.24	0.92
Depression (depression.n.01)	Cognitive State (cognitive_state.n.01)	1.21	0.87
	State of Mind (state_of_mind.n.01)	1.21	0.83
	Elation (elation.n.01)	1.15	0.91
	Fugue (fugue.n.02)	1.06	0.91
	Hallucinosis (hallucinosis.n.01)	1.05	0.92
	Abulia (abulia.n.01)	0.97	0.80
	Anxiety (anxiety.n.01)	2.09	0.91
	Mental Health (mental_health.n.01)	1.87	0.89
	Trauma (trauma.n.02)	1.71	0.84
	Mental Illness (mental_illness.n.01)	1.61	0.88
	Dissociation (dissociation.n.02)	1.56	0.89
	Morale (morale.n.01)	1.26	0.91
	Depersonalization (depersonalization.n.02)	1.32	0.92
	Enchantment (enchantment.n.02)	1.25	0.88

Continued on next page

Target (Synset)	Sibling (Synset)	Lin Similarity	Cosine Similarity
Depression (depression.n.04)	Hallucinosis (hallucinosis.n.01)	1.05	0.89
	Abulia (abulia.n.01)	0.97	0.76
	Forlornness (forlornness.n.01)	1.52	0.88
	Sorrow (sorrow.n.02)	1.36	0.86
	Heaviness (heaviness.n.02)	1.15	0.77
	Misery (misery.n.02)	1.10	0.89
	Melancholy (melancholy.n.01)	1.06	0.87
	Sorrow (sorrow.n.01)	1.13	0.85
	Weepiness (weepiness.n.01)	1.02	0.83
	Downheartedness (downheartedness.n.01)	0.93	0.88
Mental Health (mental_health.n.01)	Dolefulness (dolefulness.n.01)	0.84	0.86
	Depression (depression.n.01)	1.87	0.89
	Anxiety (anxiety.n.01)	1.85	0.89
	Trauma (trauma.n.02)	1.55	0.86
	Mental Illness (mental_illness.n.01)	1.46	0.91
	Dissociation (dissociation.n.02)	1.43	0.90
	Hypnosis (hypnosis.n.01)	1.32	0.86
	Delusion (delusion.n.01)	1.31	0.84
	Anhedonia (anhedonia.n.01)	1.24	0.83
	Agitation (agitation.n.01)	1.22	0.90
	Depersonalization (depersonalization.n.02)	1.22	0.87
	Irritation (irritation.n.01)	1.18	0.88
	Morale (morale.n.01)	1.17	0.92
	Nervousness (nervousness.n.02)	1.16	0.84
Mental Illness (mental_illness.n.01)	Enchantment (enchantment.n.02)	1.16	0.88
	Cognitive State (cognitive_state.n.01)	1.13	0.90
	State of Mind (state_of_mind.n.01)	1.13	0.85
	Elation (elation.n.01)	1.08	0.90
	Fugue (fugue.n.02)	1.00	0.86
	Hallucinosis (hallucinosis.n.01)	0.99	0.88
	Abulia (abulia.n.01)	0.92	0.79
	Depression (depression.n.01)	1.61	0.88
	Anxiety (anxiety.n.01)	1.60	0.92
	Trauma (trauma.n.02)	1.36	0.87
	Dissociation (dissociation.n.02)	1.27	0.90
	Hypnosis (hypnosis.n.01)	1.18	0.86
	Delusion (delusion.n.01)	1.18	0.86
	Anhedonia (anhedonia.n.01)	1.12	0.80

Continued on next page

Target (Synset)	Sibling (Synset)	Lin Similarity	Cosine Similarity
	Abulia (abulia.n.01)	0.85	0.76
	Depression (depression.n.01)	1.71	0.84
	Anxiety (anxiety.n.01)	1.70	0.90
	Mental Health (mental_health.n.01)	1.55	0.86
	Mental Illness (mental_illness.n.01)	1.36	0.87
	Dissociation (dissociation.n.02)	1.33	0.84
	Hypnosis (hypnosis.n.01)	1.24	0.85
	Delusion (delusion.n.01)	1.23	0.84
	Anhedonia (anhedonia.n.01)	1.17	0.84
	Agitation (agitation.n.01)	1.15	0.90
Trauma (trauma.n.02)	Depersonalization (depersonalization.n.02)	1.15	0.87
	Irritation (irritation.n.01)	1.11	0.88
	Morale (morale.n.01)	1.11	0.85
	Nervousness (nervousness.n.02)	1.10	0.85
	Enchantment (enchantment.n.02)	1.09	0.88
	Cognitive State (cognitive_state.n.01)	1.07	0.82
	State of Mind (state_of_mind.n.01)	1.07	0.85
	Elation (elation.n.01)	1.02	0.86
	Fugue (fugue.n.02)	0.95	0.89
	Hallucinosis (hallucinosis.n.01)	0.94	0.87
	Abulia (abulia.n.01)	0.88	0.82

Table 12: All Eligible Sibling Terms for Each Target Term with Lin and Cosine Similarity Scores.

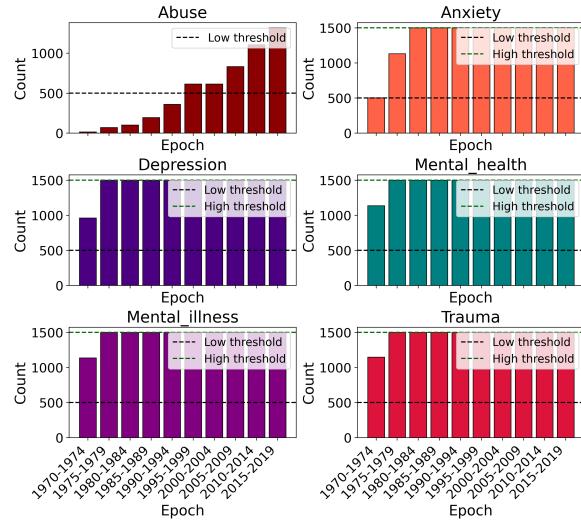


Figure 8: Counts of synthetic sentences (donor-sibling contexts).

Follow this GitHub link to access the counts of synthetic sentences for each five-year interval and the ranked lists for each sampling strategy (Boostrapped and Five-Year): [MASKED LINK]

1053
1054
1055

1056 G Multilevel Modeling Approach

1057 To analyze the predictive effects of synthetic
 1058 injections while accounting for hierarchical depen-
 1059 dencies, we employ multilevel modeling (Gel-
 1060 man and Hill, 2007). These mixed linear mod-
 1061 els are well-suited for analyzing nested data, as
 1062 they allow for the inclusion of fixed effects (e.g.,
 1063 `injection_level`) and random effects (e.g.,
 1064 variability across `target`). This approach lever-
 1065 ages the full dataset while accounting for group-
 1066 level structure, avoiding overfitting and unreliable
 1067 estimates often encountered in simple linear regres-
 1068 sion when data points per group are limited.

1069 **Model Specifications Null Model:** To assess
 1070 the necessity of incorporating random effects into
 1071 the analysis, we initially fit a *null model*. This
 1072 null model includes only a fixed intercept (β_0) and
 1073 random intercepts (u_j) to account for variability
 1074 across groups (`target`), and is represented as:

$$1075 y_{ij} = \beta_0 + u_j + \epsilon_{ij},$$

1076 where y_{ij} is the outcome variable (e.g.,
 1077 `avg_valence_index_positive`) for
 1078 observation i within group j , $u_j \sim N(0, \sigma_u^2)$
 1079 capturing group-level random intercepts, where
 1080 $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ represents residual variability.

1081 The intraclass correlation coefficient (ICC) is cal-
 1082 culated to quantify the proportion of variance ex-
 1083 plained by the grouping. ICC values exceeding
 1084 0.05 indicate meaningful variability, thereby jus-
 1085 tifying the inclusion of random effects. For this
 1086 dataset, the ICC is calculated as:

$$1087 \text{ICC} = \frac{\sigma_u^2}{\sigma_u^2 + \sigma_\epsilon^2},$$

1088 where σ_u^2 is the variance of the random intercepts
 1089 and σ_ϵ^2 is the residual variance.

1090 **Full Model:** Next, we fit a *full model*, which incor-
 1091 porates the fixed effect of `injection_ratio`
 1092 (β_1) alongside the random intercepts, expressed as:

$$1093 y_{ij} = \beta_0 + \beta_1 \cdot \text{injection_level}_{ij} + u_j + \epsilon_{ij}.$$

1094 This full model allows us to evaluate the predictive
 1095 influence of `injection_level` while account-
 1096 ing for hierarchical dependencies in the data.

1097 **Random Slopes Model:** To further explore
 1098 whether the effect of `injection_level` var-
 1099 ied significantly across `target`, we tested an
 1100 additional model with random slopes (u_{j1}) for
 1101 `injection_level`, expressed as:

$$1102 y_{ij} = \beta_0 + \beta_1 \cdot \text{injection_level}_{ij} + u_j \\ 1103 + u_{j1} \cdot \text{injection_level}_{ij} + \epsilon_{ij}.$$

1104 Here, $u_{j1} \sim N(0, \sigma_{u1}^2)$ represents the variability in
 1105 slopes across groups.

1106 **Model Comparison and Selection:** To deter-
 1107 mine the most appropriate model, we compare the
 1108 null model, simplified random intercepts model,
 1109 and random slopes model using the Akaike Infor-
 1110 mation Criterion (AIC) and Bayesian Information
 1111 Criterion (BIC), which provide measures of model
 1112 fit, with lower values indicating better balance be-
 1113 tween fit and complexity. Likelihood ratio tests
 1114 assess whether including additional random effects
 1115 significantly improve model fit. Higher Log Likeli-
 1116 hood (LL) indicates better fit.

1117 **Model Diagnostics:** Residual diagnostics are per-
 1118 formed on the final model to ensure key assump-
 1119 tions are met:

- **Normality:** Q-Q plots; Shapiro-Wilk test.
- **Homoscedasticity:** Residual vs. fitted value
 1120 plots; Levene's test for homogeneity of vari-
 1121 ances.
- **Random Effects Variance:** Variance esti-
 1122 mates for random intercepts and residuals to
 1123 quantify group-level variability contribution.

1124 Model results are summarized in Table 13, showing
 1125 that increasing levels of synthetic injections
 1126 significantly increase Dimension indices.

Method	IL	β	SE	p	LL	σ^2
Valence	-	-0.001	0.000	<.0001	59.85	0.02
Index	+	0.003	0.000	<.0001	56.30	0.03
Cosine	-	NA	NA	NA	NA	NA
Distance	+	<0.0001	<0.0001	<.0001	94.36	0.001
Arousal	-	-0.002	0.000	<.0001	63.45	0.01
Index	+	0.002	0.000	<.0001	68.38	0.009

1127 Table 13: Results of the Final Mixed Linear Models
 1128 Predicting Dimension Scores from Injection Levels

1129 Note: IL = Injection level. β = Regression coefficient
 1130 for synthetic injection level. SE = standard errors. p-
 1131 values test the null hypothesis that the coefficient is
 1132 zero. LL = Log-Likelihood, indicating model fit. σ^2 =
 1133 Variance (Random Effects), quantifies variability due to
 1134 grouping. NA = Not available.

1130 Sentiment (Positive)

- The **Null Model** showed an *ICC* of 0.59, indicating that 59% of variance in `avg_valence_index_positive` is attributable to `target`, justifying its inclusion as a random intercept.

- 1136
- The **Simplified Model**, with
1137 injection_level as a fixed effect
1138 and target as a random intercept, re-
1139 vealed a significant positive relationship
1140 ($\beta = 0.003, p < .0001$) and moderate
1141 variability across targets ($\sigma^2 = 0.03$).
1142 Residuals met homoscedasticity assumptions
1143 (Levene's $p = .92$), though the Shapiro-Wilk
1144 test ($p = .02$) suggested deviations from
1145 normality.
 - A **Random Slopes Model**, allowing
1146 injection_level to vary by target,
1147 failed to converge, rendering the fixed effect
1148 non-significant ($p = .588$) and random slope
1149 variance negligible.
 - Based on model comparison (Log-Likelihood:
1150 Simplified Model = 56.30, Random Slopes
1151 Model = 45.40; AIC/BIC unavailable due to
1152 convergence issues), the **Simplified Model**
1153 was selected as the final model (see Table 14).

Measure	Value
Number of Observations (Groups)	36 (6)
Log-Likelihood	56.30
Scale	0.0006
Random Effects Variance (Intercepts)	0.026
Fixed Effect (injection_level)	0.003
SE	± 0.000
z	27.49
p-value	< 0.0001

Table 14: Model with Random Intercepts Predicting Valence Index from Positive Sentiment Injection Level.

Sentiment (Negative)

- 1156
- The **Null Model** showed an ICC of
1157 0.884, indicating that 88.4% of variance in
1158 avg_valence_index_negative is at-
1159 tributable to target, justifying its inclusion
1160 as a random intercept.
 - The **Simplified Model**, with
1161 injection_level as a fixed effect
1162 and target as a random intercept, revealed
1163 a significant but minimal negative relationship
1164 ($\beta = -0.001, p < .0001$) with notable
1165 variability across targets ($\sigma^2 = 0.021$).
1166 Residuals met homoscedasticity assumptions
1167 (Levene's $p = .930$), while the Shapiro-Wilk
1168 test ($p = .039$) suggested minor deviations
1169 from normality, confirmed negligible by Q-Q
1170 plots.

- 1171
- A **Random Slopes Model**, allowing
1172 injection_level to vary by target,
1173 failed to converge due to overparameteriza-
1174 tion or insufficient data. The fixed effect
1175 became non-significant ($p = .822$), random
1176 slope variance was negligible ($\sigma^2 = 0.006$),
1177 and the covariance between group-level
1178 variability and injection_level was
1179 near zero (-0.000), indicating no significant
1180 interactions.
 - Given model comparison results (LL: Simpli-
1181 fied Model = 59.85, Random Slopes Model
1182 = 48.90; AIC/BIC unavailable due to conver-
1183 gence issues), the **Simplified Model** was se-
1184 lected as the final model (see Table 15).

Measure	Value
Number of Observations (Groups)	36 (6)
Log-Likelihood	59.85
Scale	0.0005
Random Effects Variance (Intercepts)	0.021
Fixed Effect (injection_level)	-0.001
SE	± 0.000
z	-11.40
p-value	< 0.0001

Table 15: Model with Random Intercepts Predicting Valence Index from Negative Sentiment Injection Level.

Intensity (High)

- 1185
- The analyses assess the im-
1186 pact of injection_level on
1187 avg_arousal_index_high. The
1188 Null Model ($ICC = 0.54$) justified including
1189 target as a random intercept, as 54% of
1190 variance was attributable to group differences.
 - The **Simplified Model with Random Inter-
1191 cepts**, incorporating injection_level
1192 as a fixed effect, converged and showed a
1193 small but significant positive effect ($\beta =$
1194 0.002, $SE = 0.000, p < .0001$), with notable
1195 group-level variability ($\sigma^2 = 0.009$).
 - A **Random Slopes Model**, allowing
1196 injection_level to vary across groups,
1197 failed to converge due to overparameteriza-
1198 tion or insufficient data. Its exceptionally low
1199 scale (0.0002) suggested overfitting or model
1200 specification issues.
 - Based on convergence and parsimony, the
1201 **Simplified Model with Random Intercepts**
1202 was selected as the final model. LL metrics

1210
1211

(Simplified: 68.38, Random Slopes: 59.99) confirmed this choice (see Table 16).

Measure	Value
Number of Observations (Groups)	36 (6)
Log-Likelihood	68.38
Scale	0.0003
Random Effects Variance (Intercepts)	0.009
Fixed Effect (injection_ratio)	0.002
SE	± 0.000
z	23.63
p-value	<.0001

Table 16: Model with Random Intercepts Predicting Arousal Index from High Intensity Injection Level.

Intensity (Low)

1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238

- The Null Model ($ICC = 0.53$) justified including target as a random intercept, as 53% of variance was attributable to group differences.
- The **Simplified Model with Random Intercepts**, incorporating `injection_level` as a fixed effect, converged and showed a small but significant negative effect on `avg_arousal_index_low` ($\beta = -0.002$, $SE = 0.000$, $p < .0001$). Random intercept variance ($\sigma^2 = 0.011$) indicated notable group-level variability. Model assumptions were met: homoscedasticity (Levene's test, $p = .982$) and linearity, though residuals deviated from normality (Shapiro-Wilk, $p = .002$), with minimal impact on Type I error given large effect sizes.
- A **Random Slopes Model**, allowing `injection_level` to vary across groups, failed to converge due to overparameterization or insufficient data. Its exceptionally low scale (0.0002) suggested overfitting or model specification issues.
- The **Simplified Model with Random Intercepts** was selected as the final model, supported by LL metrics (Simplified: 63.45, Random Slopes: 58.63). See Table 17.

Measure	Value
Number of Observations (Groups)	36 (6)
Log-Likelihood	63.45
Scale	0.0004
Random Effects Variance (Intercepts)	0.011
Fixed Effect (injection_level)	-0.002
SE	± 0.000
z	-22.98
p-value	<.0001

Table 17: Model with Random Intercepts Predicting Arousal Index from Low Intensity Injection Level.

Breadth

- The analyses assess the impact of `injection_level` on `cosine_distance_mean`. The **Null Model** ($ICC = 0.71$) justified including target as a random intercept, as 71% of variance was attributable to group differences.
- The **Simplified Model with Random Intercepts**, incorporating `injection_level` as a fixed effect, converged and showed a small but significant positive effect ($\beta < 0.0001$, $SE < 0.0001$, $p < .0001$), with notable group-level variability ($\sigma^2 = 0.001$).
- A **Random Slopes Model**, allowing `injection_level` to vary across groups, failed to converge, likely due to overparameterization or insufficient data. The model's scale was near zero (<0.0001), suggesting overfitting or misspecification.
- Based on convergence and parsimony, the **Simplified Model with Random Intercepts** was selected as the final model. LL metrics (Simplified: 94.36, Random Slopes: 84.01) confirmed this choice (see Table 18).

Measure	Value
Number of Observations (Groups)	36 (6)
Log-Likelihood	94.36
Scale	0.0001
Random Effects Variance (Intercepts)	0.001
Fixed Effect (injection_level)	<0.001
SE	± <0.0001
z	7.49
p-value	<.0001

Table 18: Model with Random Intercepts Predicting Arousal Index from High Intensity Injection Level.

H SIB Scores: Results for Five-Year Random Sampling Strategy

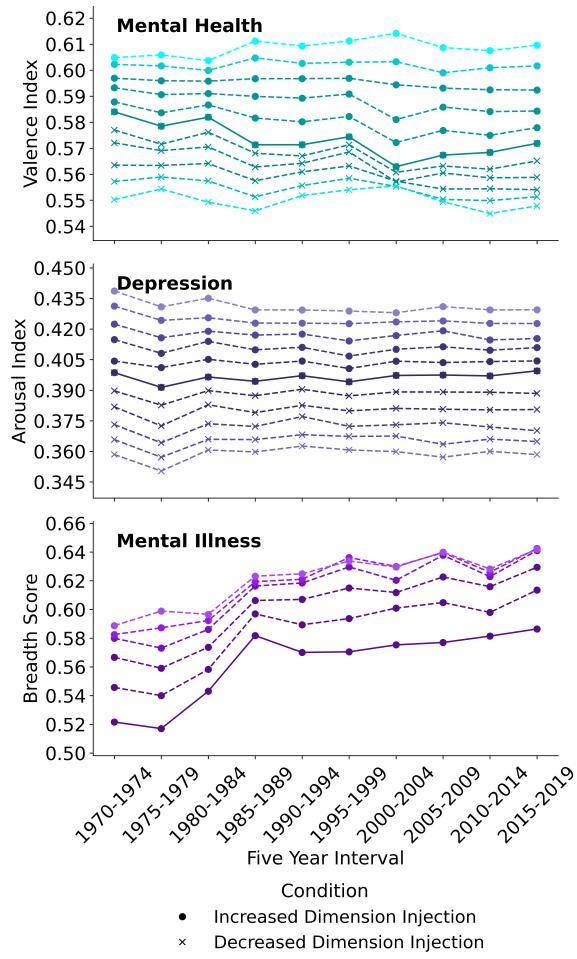


Figure 9: SIB Scores by Five-Year Intervals Across Injection Levels and Conditions.

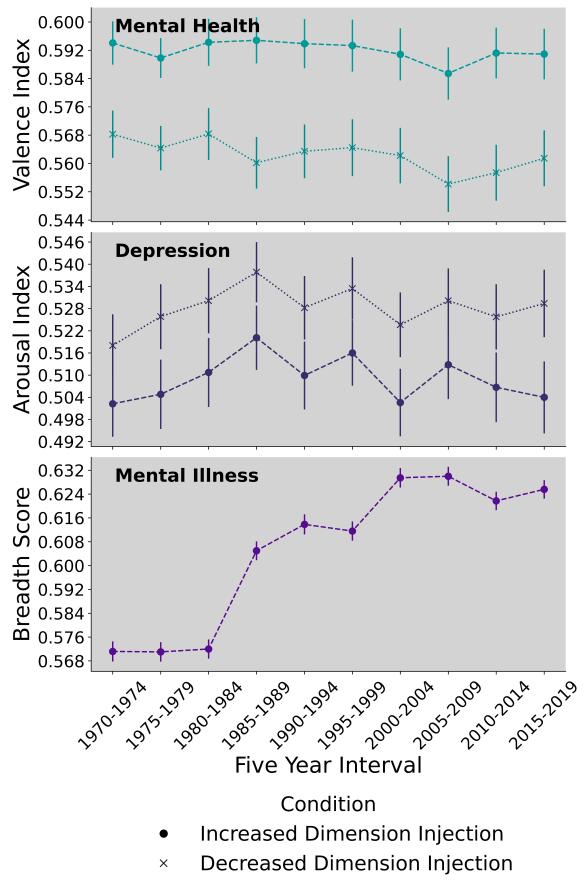


Figure 10: SIB Scores (\pm SE) by 50% Injection Levels and Conditions: Control Setting for Five-Year Samples.

I Alternative LSC Detection Methods: Results for Bootstrapped Settings

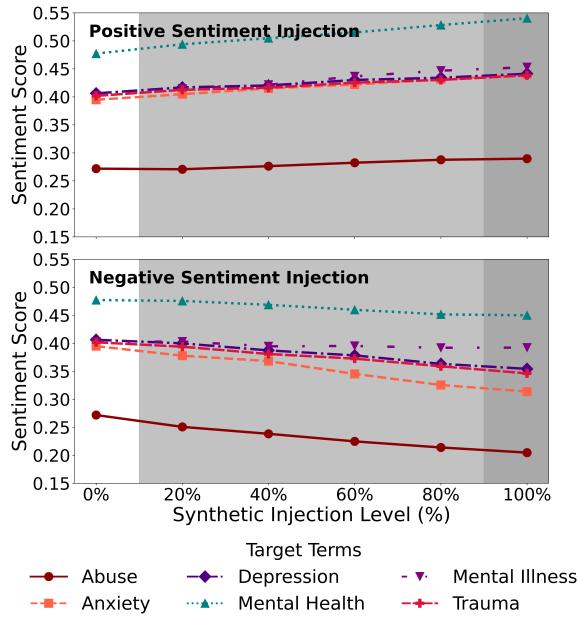


Figure 11: ABSA Sentiment Index Across Injection Levels and Sentiment Conditions: Bootstrapped Samples.

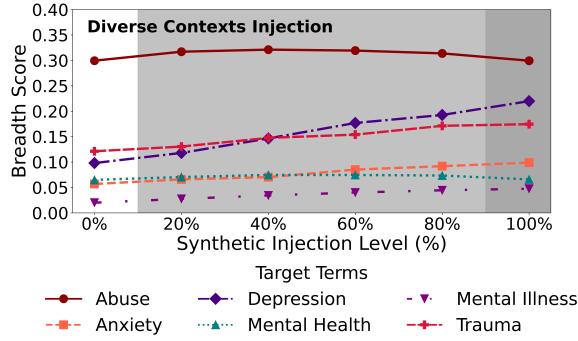


Figure 12: XL-LEXEME Breadth Score (Average Cosine Distance Within-Bins) Across Injection Levels and Breadth Condition: Bootstrapped Samples.

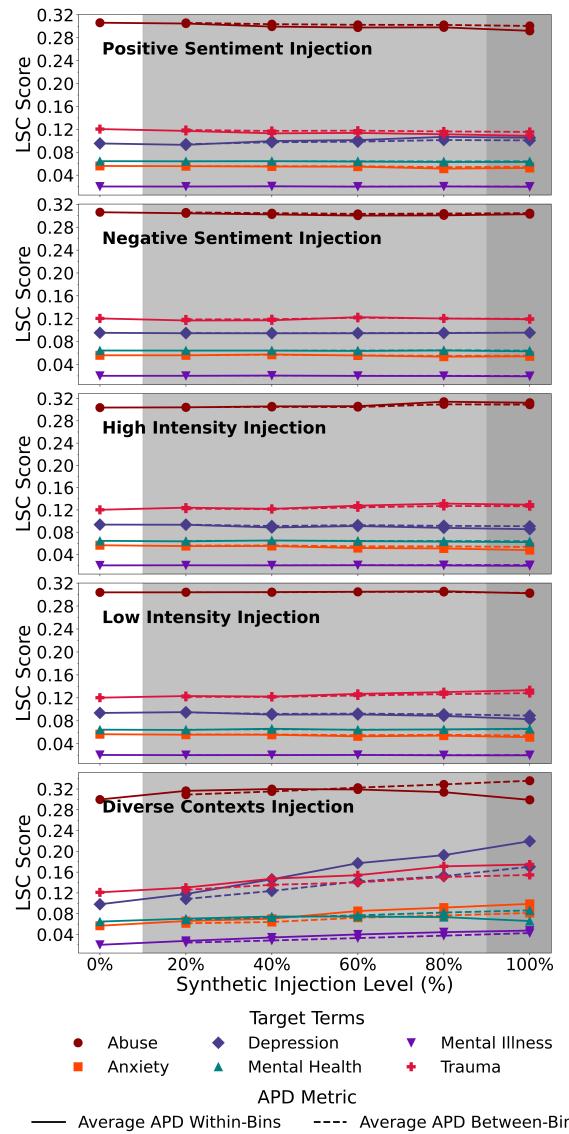


Figure 13: LSC Scores (APD Between-Bins and APD Within-Bins) Across Injection Levels and SIB Conditions: Bootstrapped Samples.