

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

Naomi Baes^Ψ, Raphaël Merx^λ, Nick Haslam^Ψ, Ekaterina Vylomova^λ, Haim Dubossarsky^{ΦΤΣ}

^ΨMelbourne School of Psychological Sciences, The University of Melbourne

^λSchool of Computing and Information Systems, The University of Melbourne

^ΦSchool of Electronic Engineering and Computer Science, Queen Mary University of London


^ΤThe Alan Turing Institute, London

^ΣLanguage Technology Lab, University of Cambridge

{n.baes, r.merx, nhaslam, vylomovae}@unimelb.edu.au, h.dubossarsky@qmul.ac.uk

Abstract

Lexical Semantic Change (LSC) provides insight into cultural and social dynamics. Yet, the validity of methods for measuring different kinds of LSC remains unestablished due to the absence of historical benchmark datasets. To address this gap, we propose LSC-Eval, a novel three-stage general-purpose evaluation framework to: (1) develop a scalable methodology for generating synthetic datasets that simulate theory-driven LSC using In-Context Learning and a lexical database; (2) use these datasets to evaluate the sensitivity of computational methods to synthetic change; and (3) assess their suitability for detecting change in specific dimensions and domains. We apply LSC-Eval to simulate changes along the Sentiment, Intensity, and Breadth (SIB) dimensions, as defined in the SIBling framework, using examples from psychology. We then evaluate the ability of selected methods to detect these controlled interventions. Our findings validate the use of synthetic benchmarks, demonstrate that tailored methods effectively detect changes along SIB dimensions, and reveal that a state-of-the-art LSC model faces challenges in detecting affective dimensions of LSC. LSC-Eval offers a valuable tool for dimension- and domain-specific benchmarking of LSC methods, with particular relevance to the social sciences.

 https://github.com/naomibaes/LSCD_method_evaluation

1 Introduction

Lexical Semantic Change (LSC) provides a unique window into cultural dynamics by revealing how language evolution reflects social changes (McGillivray, 2020; Schlechtweg et al., 2020). 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 proposes LSC-Eval, a novel three-stage evaluation framework for assessing LSC across dimensions and domains. It: (1) develops a scalable, general-purpose methodology for generating high-quality synthetic sentences using ‘scholar-in-the-loop’ In-Context Learning (ICL) and a lexical resource to simulate targeted kinds of LSC; (2) uses these datasets to evaluate the validity of alternative LSC detection methods; and (3) identifies the most suitable methods for capturing change in specific dimensions and domains. We apply LSC-Eval to test the sensitivity of various methods to synthetically induced changes along three major LSC dimensions — Sentiment, Intensity, and Breadth (SIB) — as proposed in the SIBling framework (Baes et al. 2024; see Table 1), drawing from psychology examples. Illustrative cases include *awesome*, which has risen in sentiment (from serious to enthusiastic), declined in intensity (from awe-inspiring to casually pleasing), and broadened in scope (from solemn reverence to general positivity); *trauma*, which has shifted from physical injury to a broader and milder range of psychological harms; and *sick*, which has risen in sentiment (in slang), broadened in usage (from a medical term to a general expression of praise), and declined in intensity (from a descriptor of serious illness to one used in casual or enthusiastic contexts).

The present study addresses two key research questions: (RQ1) Can synthetic datasets be used to validate methods for measuring LSC dimensions? We hypothesize that SIB scores will be associated with levels of induced change; (RQ2) Which of

Dimension	Definition	Examples of Rising	Examples of Falling
<i>Sentiment</i>	Relates to the degree to which a word’s meaning acquires more positive (<i>‘elevation’</i> , <i>‘amelioration’</i>) or negative (<i>‘degeneration’</i> , <i>‘pejoration’</i>) connotations.	<i>craftsman</i> , once associated with manual labor, has come to convey artistry, skill, and high-quality workmanship. <i>geek</i> , originally a derogatory term for odd people, now refers to someone passionate about a field.	<i>retarded</i> , originally a neutral term for intellectual disability, has become highly pejorative over time. <i>awful</i> , once meaning “awe-inspiring,” now indicates something very bad.
<i>Intensity</i>	Relates to the degree to which a word’s meaning changes to acquire more (<i>‘meiosis’</i>) or less (<i>‘hyperbole’</i>) emotionally charged (i.e., strong, potent, high-arousal) connotations.	<i>cool</i> has evolved from describing temperature 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> has evolved from denoting romantic or platonic attachment to also expressing milder forms of liking (e.g., “I love coffee.”). <i>trauma</i> has shifted from referring to brain injuries to encompassing milder events (e.g., business loss).
<i>Breadth</i>	Relates to the degree to which a word expands (<i>‘widening’</i> , <i>‘generalization’</i>) or contracts (<i>‘narrowing’</i> , <i>‘specialization’</i>) its semantic range.	<i>cloud</i> , initially a meteorological term, broadened in usage to refer to internet-based data storage. <i>partner</i> , which once referred narrowly to a business associate, has broadened to denote a romantic or domestic companion.	<i>doctor</i> , once used for any scholar or teacher, now primarily denotes a medical professional. <i>meat</i> , which originally referred to any kind of food in Old English (<i>mete</i>), has since narrowed to denote animal flesh as food.

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

several LSC detection methods is most sensitive to synthetically induced changes in SIB? Our findings confirm the validity of theory-driven changes — that is, changes systematically injected into corpora based on the SIBling framework (Baes et al., 2024), which targets the dimensions of Sentiment, Intensity, and Breadth. These changes were operationalized using ICL and a lexical resource to simulate diachronic semantic shifts. Results highlight the need to tailor detection methods to specific dimensions, as the state-of-the-art LSC model failed to detect affective change. The proposed framework, LSC-Eval, offers an efficient and scalable approach for dimension- and domain-specific benchmarking of LSC methods. While broadly applicable, LSC-Eval is especially valuable in the social sciences and humanities, where capturing nuanced conceptual change is essential.

2 Related Work

2.1 Theoretical Background

Linguists have long debated taxonomies of LSC (Bréal, 1900; Blank, 1999), defined as innovations which change word meanings (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 computational approaches (Kutuzov et al., 2018; Tahmasebi et al., 2018; Tang, 2018; Tahmasebi and Dubossarsky, 2023; Cassotti et al., 2024b; Periti and Montanelli, 2024a; Kiyama et al., 2025).

Recent years have seen the development of theoretical frameworks that propose multiple dimensions of LSC. Baes et al. (2024) introduced SIBling, a three-dimensional framework that maps LSC along axes of SIB, reflecting a word’s acquisition of

more positive or negative connotations (Sentiment), more or less emotionally charged or potent connotations (Intensity), and the expansion or contraction of its semantic range (Breadth). It draws on linguistic (Geeraerts, 2010) and psychological (Haslam, 2016) theories, and provides methodological tools to estimate SIB across time. In parallel, de Sá et al. (2024) proposed a framework that clusters LSC into three dimensions using graph structures: Orientation (shifts towards more pejorative or ameliorated senses), Relation (changes towards metaphoric or metonymic usage), and Dimension (variations between abstract/general and specific/narrow meanings). While de Sá et al. (2024) surveyed statistical methods for representing word meaning (word frequency, topic modeling, and graph structures) on dimensions, they did not demonstrate their usage.

Both frameworks contain dimensions of evaluation (Sentiment and Orientation) and semantic range (Breadth and Dimension). Baes et al.’s (2024) inclusion of Intensity reflects a greater emphasis on changes in the emotional connotations of words. Sentiment and Intensity resemble the two primary dimensions of human emotion, Valence and Arousal (Russell, 2003), and two primary dimensions of connotational meaning, Evaluation (e.g., “good/bad”) and Potency (e.g., “strong/weak”) (Osgood et al., 1975), which have been demonstrated to have cross-cultural validity.

2.2 Evaluation

Despite substantial progress in developing benchmarks (Tahmasebi and Risse, 2017) and evaluation strategies (Kutuzov et al., 2018), the field still lacks standardized datasets that evaluate multiple dimensions of LSC across time. Current annotated

benchmarks, such as the synchronic, definition- and type-based *LSC Cause-Type-Definitions Benchmark* (Cassotti et al., 2024a) and the binary, word-sense-based *TempoWIC*, where LSC is labeled by comparing the sameness or difference of meanings between two sense usages (Loureiro et al., 2022), address different aspects of semantic change.

The first human-annotated dataset of LSC in multiple languages (English, German, Latin, Swedish; Schlechtweg et al., 2020) marked substantial progress in identifying the presence and degree of LSC but omitted information about kinds of change. However, expert-annotated datasets are costly and time-intensive to create. To address this gap, Dubossarsky et al. (2019) introduced a method to artificially induce semantic change in controlled experimental settings, enabling precise testing of how well models capture these shifts.

Recent developments in generative artificial intelligence highlight the potential of pre-trained LLMs to adapt to novel tasks at inference time through ICL (Zhou et al., 2023). Few-shot ICL, a paradigm that enables LLMs to learn tasks by analogy given only a few demonstrative examples, helps to incorporate theoretical knowledge without needing to fine-tune its internal parameters (Dong et al., 2024). Instead, ICL uses context from the model’s prompt to adapt the LLM to downstream tasks (Radford et al., 2019; Brown et al., 2020; Liu et al., 2024). de Sá et al. (2024) demonstrated the utility of few-shot ICL, employing Chain-of-Thought and rhetorical devices, to annotate LSC dimensions, but their strategy focused on multi-class classification of change between two sense usages. ICL offers a promising solution to bridge the absence of standardized approaches (Hengchen et al., 2021) for assessing the effectiveness of different methods to measure dimensions of LSC.

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

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

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, 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* — were 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 alongside their semantic expansions (*trauma*: Baes et al., 2023; Haslam et al., 2021; *mental health*, *mental illness*: Baes et al., 2024). There have been changes in their semantic intensity, 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, *abuse* has expanded horizontally to include passive neglect and emotional abuse, beyond its physical scope (Haslam, 2016). Targets were sufficiently prevalent (sentence counts: 46,272; 104,486; 115,430; 44,130; 5,808; 23,187). Appendix A shows annual counts.

3.2 Evaluation Framework

The three general stages of the evaluation framework are illustrated in Figure 1.

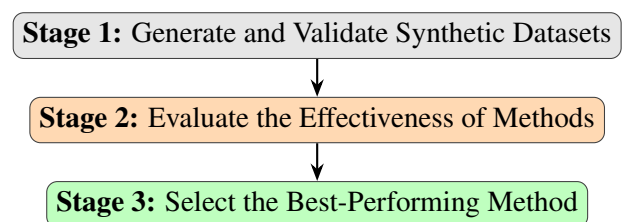


Figure 1: Stages of the Evaluation Framework.

In Stage 1, synthetic datasets (i.e., datasets exhibiting simulated variations in linear semantic change along a specific dimension) were created to benchmark changes in LSC dimensions using few-shot ICL and a lexical database. GPT-4o (Achiam et al., 2023)³ was prompted with expert-crafted examples to increase and decrease corpus sentences

²PsycNET, UMLS, DSM-5, ConceptNet

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

in affective dimensions across five-year intervals. This ensures that synthetic sentences are theory-driven, domain-specific and contain temporal features. GPT was 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). Table 2 provides descriptive statistics for the synthetic datasets,⁴ which are validated using tools that have been shown to measure SIB in historical (domain-general and -specific) corpora (Baes et al., 2024).

Dimension	Target	Neutral (<i>M</i>)	Increase (<i>M</i>)	Decrease (<i>M</i>)	US\$
Sentiment	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
Intensity	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
Breadth	abuse	–	5,221 (27)	–	0
	anxiety	–	13,635 (26)	–	0
	depression	–	14,463 (27)	–	0
	mental health	–	14,638 (26)	–	0
	mental illness	–	14,639 (26)	–	0
	trauma	–	14,650 (26)	–	0

Table 2: Descriptive Statistics for Synthetic Dimension Datasets: Sentence Counts, Mean Lengths, and Total Generation Cost.

Note. *M* = Mean sentence word length. Neutral = Original input sentences. Increase = Sentences modified to increase the dimension of interest. Decrease = Sentences modified to decrease the dimension of interest.

Experimental Settings: We implemented two sampling strategies to assess the sensitivity of LSC detection methods. *Bootstrap sampling* involved randomly selecting 50 sentences with replacement from the full corpus (natural or synthetic), repeated over 100 iterations. This allowed any sentence to be sampled multiple times, enabling robust estimation through resampling. *Five-year interval sampling* selected up to 50 unique sentences per time bin, repeated 10 times. Sentences were not repeated within an iteration but could appear across iterations. This approach ensured balanced temporal coverage and better reflects natural language variability over time. Each sampling strategy had a *control condition*, where sentences were shuffled to balance natural and synthetic ones per sample, providing a baseline to isolate the effect of synthetic

⁴Link to Synthetic datasets: https://github.com/naomibaes/Synthetic-LSC_pipeline

injection, following prior work in computational linguistics (Dubossarsky et al., 2017, 2019).

For each sampling strategy, synthetic sentences were injected into (up to) 50-sentence samples at increasing proportions: 20%, 40%, 60%, 80%, and 100% (see Figure 2). These injection levels simulate increasing semantic saturation to test method sensitivity (Stage 2). For method selection (Stage 3), only the 0% and 100% injection levels were used. This enables the clearest contrast between unchanged (fully natural) and fully altered (synthetic) samples, maximizing the signal-to-noise ratio when evaluating which method detects the greatest magnitude of semantic change.

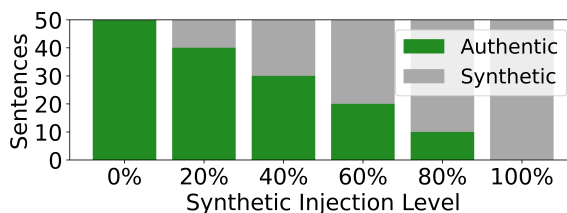


Figure 2: Distribution of Sentences in Each Sample.

3.3 Synthetic Dataset Generation

For each semantic dimension and target (Section 3.1.3), we generated datasets by inducing change in natural sentences drawn from the corpus (Section 3.1.1). Each dataset included up to 1,500 randomly sampled sentences per five-year interval used in the experimental settings (see Section 3.2). See Appendix B for examples of synthetic sentences for each dimension and target.

3.3.1 Sentiment and Intensity

To generate the synthetic Sentiment and Intensity datasets, we employed few-shot ICL with GPT-4o to vary these dimensions. First, neutral diachronic sentences from the corpus (see Section 3.1.1) were sampled as outlined in Appendix C. Second, a psychology scholar crafted five (Chen et al., 2023) diverse sentence variations for each target, to serve as few-shot demonstrations for the LLM in the prompt outlined below. This prompt includes construct definitions, a task description, and few-shot examples to generate theory-driven change — that is, changes aligned with the dimensions in the SIBling framework (e.g., shifts from positive to negative sentiment or from high to low intensity). For ‘scholar-in-the-loop’ few-shot demonstrations, see Appendix D (Sentiment), Appendix E (Intensity), and the few-shot demonstration examples below.

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.

Third, the prompt was refined during pilot tests (varying 10 input sentences for each target).

Fourth, for each of the neutral sentences (Sentiment: 36,151; Intensity: 39,896), we made one inference call to GPT-4o through the OpenAI API to generate variations of Sentiment (positive/negative) or Intensity (high/low). Fifth, validation checks required checking if output sentences retained the target term. A few manual adjustments were made by a tertiary-educated English native speaker to ensure the final sentences retained the target in a similar location to the original input sentence. Few output sentences required manual adjustment in the final datasets: 0.25% for the synthetic Sentiment dataset, and 0.01% for the synthetic Intensity dataset. See Table 2 for the counts of input/output sentences and final prompts (dataset costs: for Sentiment = 115 \$US; for Intensity = 136 \$US).

3.3.2 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 induce semantic breadth, we adapted Dubossarsky et al.’s (2019) replacement strategy, using WordNet 3.0 (see Section 3.1.2) to expand a target word’s usage by incorporating contexts from donor terms, thereby broadening its semantic range without altering its core meaning. Relevant synsets were identified and filtered for psychological relevance using keyword matching⁵ and semantic similarity thresholds. Donor terms (co-hyponyms with the target) were filtered using Lin similarity (0.5)⁶ and cosine similarity (0.7) using 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 full list of siblings (i.e., co-hyponyms). The sibling replacement process identifies and replaces sibling terms with the target, as illustrated below. To sample representatively from the sibling list, a round-robin strategy was used, where the algorithm runs through each sibling in the list in batch sizes, randomly sampling up to 50 unique sentences before proceeding to the next one. In this way, it exhaustively samples up to 1,500 unique sentences per epoch and injection level to create the final synthetic breadth dataset.

⁵Psychology key terms: "abnormality", "abnormally", "emotional", "feeling", "feelings", "harm", "hurt", "mental", "mind", "psychological", "psychology", "psychiatry", "syndrome", "therapy", "treatment".

⁶Information content values from the psychology corpus.

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 Quantifying Lexical Semantic Change

3.4.1 Semantic Dimensions

We applied methodologies from Baes et al. (2024) to assess changes in semantic Sentiment, Intensity, and Breadth. Sentiment and Intensity were quantified by assigning ordinal valence and arousal scores, based on Warriner et al. (2013) norms, matched to collocates within ± 5 words of the target, using a scale (ranging from 1-9) from 'extremely unhappy' to 'extremely happy' for valence, and 'extremely calm' to 'extremely agitated' for arousal. Scores were frequency-weighted, normalized, and averaged across bins, yielding indices from 0 (negative/low arousal) to 1 (positive/high arousal). Breadth was measured by averaging the cosine distances between sentence-level embeddings from the SentenceTransformer model 'all-mpnet-base-v2'⁷ (MPNet). The Breadth score indicates semantic range variation from 0 (no change) to 1 (maximum variation). See Appendix G.1.1 for further details.

For comparative analysis, Sentiment was compared against the Deberta-v3-base-absa-v1.1⁸ aspect-based sentiment analysis (ABSA) classification model, which we adapted to output continuous sentiment scores ranging from 0 (fully negative) to 1 (fully positive). Because Intensity is a novel dimension, it lacks comparable methods. For the Breadth comparisons, we used Cassotti et al.'s (2023) SOTA (Periti and Tahmasebi, 2024) word transformer "XL-LEXEME"⁹ (XLL). While

⁷Microsoft pretrained network (109M model params) <https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

⁸yangheng/deberta-v3-base-absa-v1.1 (184M model params): <https://huggingface.co/yangheng/deberta-v3-base-absa-v1.1>

⁹XL-LEXEME (~550M model params) <https://huggingface.co/pierluigi/xl-lexeme>

MPNet generates sentence embeddings through pooling tokens, which dilutes word-specific information, XLL uses a bi-encoder architecture that focuses on word-specific attention, using polysemy as a proxy for meaning divergence during training. See Appendix G.1.2 for further details on the quantification of semantic Breadth.

3.4.2 General Lexical Semantic Change

To quantify general LSC, we used the SOTA LSC score (Cassotti et al., 2023), which calculates the Average Pairwise Cosine Distance (APD) between sentence embeddings from two time periods, following the approach of Giulianelli et al. (2020). We extended this method to compare embeddings from different bins within the same iteration. See Appendix G.2 for more details.

3.5 Statistical Strategy

To test whether synthetic change influences SIB scores (RQ1), we fit separate mixed linear models for each outcome variable (Valence, Arousal, and Breadth), with random intercepts for target terms and fixed effects for injection level (see Appendix H for model specifications). To assess the sensitivity of various computational methods to detecting LSC (RQ2), we calculated the percent relative change in a target word's context between the 0% (fully natural) and 100% (fully synthetic) injection levels, as defined in Equation 1:

$$\Delta\% = \frac{X_{100} - X_0}{X_0} \times 100 \quad (1)$$

Here, X_0 and X_{100} denote the score at 0% and 100% injection, respectively. This index quantifies how much the contextual meaning of the target word shifts when all natural sentences are replaced with synthetic ones.

For the XLL LSC Score, we further normalized this change to account for internal variability within each bin, as shown in Equation 2. We used the APD, the mean cosine distance between all sentence pairs,¹⁰ to control for within-bin noise:

$$\Delta = \frac{\text{APD}(X_{100\text{-between-}X_0})}{\max[\text{APD}(X_0\text{-within-}X_0), \text{APD}(X_{100\text{-within-}X_{100}})]} \quad (2)$$

The numerator reflects divergence between bins, while the denominator accounts for within-bin variability, ensuring that the score reflects genuine LSC signal rather than internal fluctuation.

¹⁰APD reflects semantic dispersion; higher APD implies greater variability.

4 Results

Synthetic Change Effects: To address RQ1, we tested whether SIB scores from Baes et al.’s (2024) SIBling framework are associated with levels of synthetically induced change. The hypothesis was supported: standardized SIB scores across all three dimensions — Sentiment, Intensity, and Breadth — demonstrated consistent increases or decreases as a function of injection level, across targets and both sampling strategies: bootstrap (Figure 3) and five-year intervals (Figure 4).

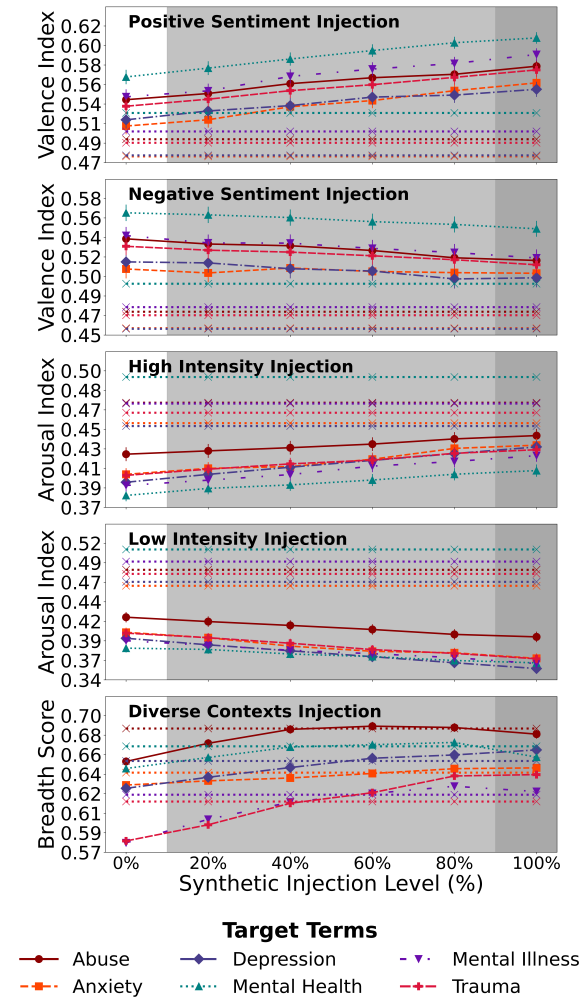


Figure 3: SIB Scores (\pm SE) by Injection Levels for Experimental and Control Settings (Flat Dotted Lines).

Figure 4 illustrates SIB scores by five-year intervals across injection levels and conditions for three of the six targets. The black line in each panel depicts the changes in SIB scores over time in the natural corpus, while the colored lines represent injection levels. The intensity of the color reflects the percentage of synthetic sentences injected into each sample (20% to 100%). As intended, injection levels altered the height, but not necessarily the slope, of SIB scores on the y-axis, indicating changes in the magnitude of the score rather than trend over time. For breadth (bottom panel), only upward (broadening) changes were modeled.

tion levels altered the height, but not necessarily the slope, of SIB scores on the y-axis, indicating changes in the magnitude of the score rather than trend over time. For breadth (bottom panel), only upward (broadening) changes were modeled.

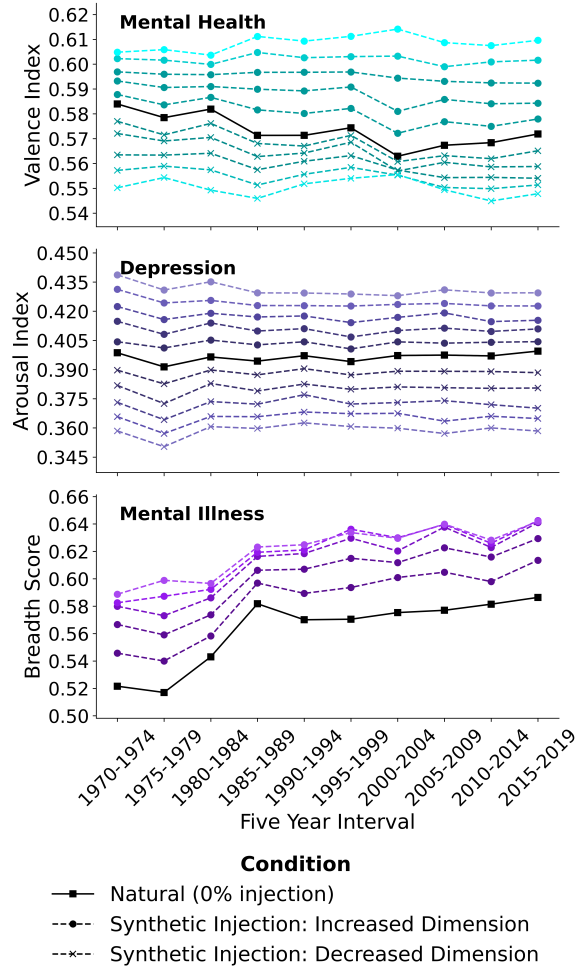


Figure 4: SIB scores by five-year intervals across injection levels and conditions.

Mixed-effects linear models confirmed that a one-standard-deviation increase in injection level significantly predicted changes in SIB scores in the intended direction (Table 3; see also Appendix H).

Score	Valence	Arousal	Breadth
β^+	.61*	.64*	.43*
β^-	-.31*	-.64*	-

Table 3: Coefficients of random intercept (by target) models predicting SIB scores from injection level.

Control Experiments: As illustrated in Figure 3, controlling for synthetic injection level by re-analyzing data with shuffled sentences for uniform

distribution revealed flat SIB score trends in bootstrapped settings. Appendix I illustrated how, in the control condition for the five-year interval setting, SIB scores in those samples tend to converge to a midpoint between natural and synthetic data.

Comparative Method Evaluation: To address RQ2 — identifying which methods are most sensitive to synthetically induced changes in SIB — we compared the performance of alternative change detection methods across each dimension. Results were mixed. For Sentiment, both the Valence index and ABSA’s Sentiment score were sensitive to detecting synthetic change, though ABSA outperformed the Valence index in 10/12 cases. For Intensity, the Arousal index shows sensitivity to detecting variations in synthetic Intensity. For Breadth, XLL outperformed MPNet in 4/6 cases using the Breadth score.

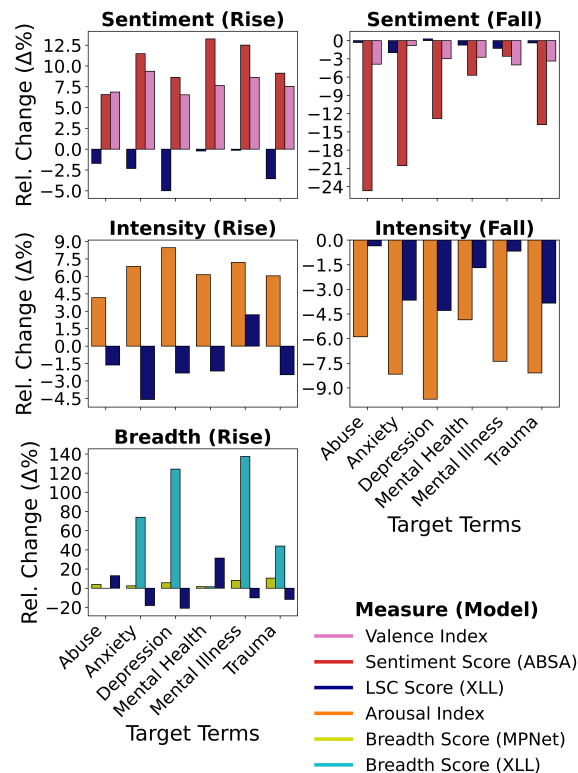


Figure 5: Relative Change ($\Delta\%$) Scores for Models Across Dimensions and Conditions: Bootstrapped.

Critically, XLL-LSC score is completely insensitive to detecting changes in either Sentiment or 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 J 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 LSC-Eval, a three-stage general 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 change; and (3) evaluates the suitability of these methods for specific dimensions and domains. We applied this framework to generate synthetic datasets that induce changes across the three dimensions (Sentiment, Intensity, and Breadth) of a novel multidimensional LSC framework, SIBling (Baes et al., 2024), using examples drawn from psychology to compare the suitability of alternative methods for detecting these synthetic changes.

Importantly, these synthetic datasets reflect theory-driven change, operationalized through construct definitions and distinctions drawn from the SIBling framework (Baes et al., 2024), which integrates established traditions in lexical semantics, connotation, and affect. This principled approach enables interpretable and targeted benchmarking, and lays the groundwork for generalizing LSC-Eval to other theoretical models and domains.

Our findings support the hypothesis that recently proposed methods (Valence index, Arousal index, Breadth score; Baes et al., 2024) detect synthetic changes on the SIB (i.e., Sentiment, Intensity, and Breadth) dimensions. This directly addresses RQ1, which asked whether synthetic datasets can be used to validate methods for measuring LSC dimensions. Control analyses, which adhered to computational linguistics standards (Dubossarsky et al., 2017, 2019), confirmed the absence of these effects in shuffled controls. The implications of these findings are two-fold. First, the ability of SIB methods to detect changes introduced via silver-standard synthetic data validates their sensitivity and reliability for detecting and quantifying variation in SIB, even in artificial environments. Second, it confirms the validity of using LLM-generated sentences in our ICL evaluation suites.

To address RQ2 — identifying which computational methods are most sensitive to synthetically induced changes in SIB — we demonstrated how a

synthetic change detection task can assess the sensitivity of various approaches, guiding the selection of the most suitable method for specific dimensions and domains. Baes et al.'s (2024) tools, which validated the synthetic SIB datasets, were supported by alternative methods that consistently detected synthetic changes across conditions and targets, providing further validation. The Valence index and ABSA's Sentiment score detected variations in synthetic Sentiment, while Breadth scores derived from XLL and MPNet models detected increases in synthetic Breadth. These results suggest that the NLP-based ABSA method is generally more sensitive to synthetically induced change than Warriner-based approaches, which rely on static Valence and Arousal ratings. Future empirical studies on Sentiment and Breadth may consider adopting these NLP-based models, either as replacements for, or in addition to, traditional lexicon-based methods.

Notably, when computing the general LSC score using the SOTA LSC model XLL (Cassotti et al., 2023), it was not sensitive to detecting Sentiment and Intensity. Although XLL shows some sensitivity to identifying synthetic increases in Breadth, it registers a more substantial change when the Breadth score is adjusted according to the method introduced by Baes et al. (2024). It uses the within-bin average cosine distance of target-containing sentences as a proxy for the expansion (broadening) or contraction (narrowing) of a word's contextual usage. The inability of XLL to detect the affective dimensions of LSC highlights the necessity of evaluating SOTA models before deploying them in new domains. Future research should investigate whether this weakness in detecting affective dimensions is specific to XLL or extends to other contextualized models in other corpora and concepts. This inquiry is particularly salient given recent advances in analyzing fine-grained, continuous semantic shifts through "diachronic word similarity matrices using fast and lightweight word embeddings over arbitrary time periods" (Kiyama et al., 2025). Our results support the idea that current SOTA methods are only state-of-the-art with respect to a particular kind of change or a specific domain and should not be assumed to generalize beyond them.

Findings highlight the need to include affective and connotational aspects of meaning in studies of LSC. In particular, future studies must consider emotional meaning in language models. While psychology has extensively used language to analyze

emotion semantics (Jackson et al., 2022; Boyd and Schwartz, 2021), advances in NLP are still exploring how to build models that incorporate sentiment (Goworek and Dubossarsky, 2024) and detect emotion (Mohammad, 2021). Further research is required to detect affective states from text given the cultural and universal aspects of emotion semantics (Jackson et al., 2019). These findings also have implications for existing multidimensional frameworks of LSC (Baes et al., 2024; de Sá et al., 2024) as the evaluation framework provides experimental settings in which to compare the sensitivity of methods to detecting synthetic changes on specific dimensions and domains in a variety of disciplines.

6 Conclusion

The current study introduced a novel general-purpose evaluation framework, LSC-Eval, for evaluating methods that detect LSC. Its three-stage approach: 1) generates LLM-based synthetic datasets with silver-standard labels that simulate changes in kinds of LSC; 2) uses these datasets to evaluate the relative sensitivity of computational methods in a synthetic change detection task; and 3) identifies the most suitable method for detecting synthetically induced changes across specific dimensions and domains. We applied this framework to a set of psychological terms. Findings not only supported the validity of proposed computational methods for measuring changes in SIB, but also established a controlled experimental standard for rigorously evaluating existing LSC detection methods and exploring alternative computational approaches. This work is crucial for addressing the substantial gap created by the lack of historical benchmark datasets, which has previously hindered the standardization of metrics and fair comparison of methods. While this framework benefits all disciplines (e.g., biomedicine, law, theology), LSC-Eval is particularly valuable in the social sciences and humanities, where domain-specific constructs require specialized tools for capturing nuanced semantic change.

Limitations

Limitations inform future directions. Evaluating the quality of LLM prompt and demonstration examples in the few-shot ICL paradigm is challenging. As LLM evaluation standards are developed (Chang et al., 2024; Ziems et al., 2024), future research might explore automated strategies such as updating prompts based on examples (DSPy)¹¹ or comparing LLM output from different prompts using a free, unified interface.¹² LLM choice in the evaluation pipeline could be expanded to include open-source models (e.g., FlanT5-XL, Mistral-7B, Mixtral-8x7B). Future work could also compare this approach to a non-generative rule-based one.

Furthermore, our study benefited from using GPT-4o, which is trained on US English and is therefore well-suited for analyzing texts within the Western-centric domain of psychology (Varnum et al., 2024). However, cultural and linguistic biases of LLMs can pose challenges in adapting our evaluation pipeline to other languages (Havaladar et al., 2023), although few-shot ICL has proven effective in low-resource languages (Cahyawijaya et al., 2024). Despite the tendency of LLM training data to skew towards the recent past, manual inspection of results demonstrate the successful generation of quality sentences that spanned a 1970 to 2019 time period. Future work should focus on refining these models to broadly apply across cultural contexts, languages, and historical periods.

The conceptualization of semantic Breadth is complex and contested. Linguistic definitions suggest breadth encompasses subtypes (e.g., specialization as a subtype of narrowing; Campbell, 2013) highlighting its intricate nature. Given this complexity, it is essential to compare the current measure, which is based on mean within-bin variability of target-containing sentences, with other methods assessing breadth through senses, topics, or prototypical changes: modulations based on literal similarity (Geeraerts, 1997). Future research should investigate whether these measures can detect polysemy's emergence or merely prototype-based modulations of existing concepts.

The synthetic breadth dataset used in this study was constructed using a replacement strategy that may include contextually irrelevant donor contexts. To enhance simulation quality, we propose

a three-step validation pipeline: First, select validation models based on performance against a gold-standard dataset, as determined by the highest F1-score from 5-fold stratified K-fold cross-validation. Second, use a probability ratio check with a Masked Language Model (e.g., BioBERT, RoBERTa-large, DeBERTa-v3-large) to confirm the plausibility of replacing donors with target terms, approving sentences that meet a specific probability threshold. Third, ensure semantic alignment through cosine similarity validation with models such as MiniLM-L12-v2 or DistilRoBERTa-v1 Sentence-T5, approving sentences that exceed a set threshold. This process aims to expand the target term's semantic scope while maintaining specificity, but may exclude many sentences. Integrating de Sá et al.'s (2024) ICL approach to simulate Breadth—first teaching the model to disambiguate word senses—could offer an efficient alternative.

Furthermore, the present study does not specify which sense of the term is semantically expanded. Attempting to integrate senses into the synthetic data generation pipeline may provide richer insights. While the specialized psychology corpus and target words exhibit limited senses, general domain corpora introduce ambiguous contexts (e.g., the economic sense of “depression”). Further research is needed on sense-specific LSC.

Although a body of work estimates valence from natural language, less research has examined the Intensity dimension (Hoemann et al., 2025). In the present study, this restricted the external validation of the Arousal index (Baes et al., 2024), highlighting the need for empirical research in this direction. Furthermore, we must examine the conceptual/terminological link between arousal and hyperbole (i.e., a linguistic form describing a rhetorical, discursive phenomenon like irony) to understand arousal's relation to hyperbole (Burgers et al., 2016; Peña and Ruiz de Mendoza, 2017).

Finally, future research should use the evaluation framework to generate synthetic datasets, and to explore methods, for detecting the Relation dimension (metaphor/metonymy) as highlighted by de Sá et al. (2024). Incorporating the qualitative types of metaphor and metonymy into the empirical study of multidimensional LSC could provide a more comprehensive understanding of LSC, particularly for some domains. Examining how Relation relates to SIB may deepen our understanding of LSC processes by exploring how cognitive principles contribute to semantic innovations.

¹¹<https://dspai.ai>

¹²<https://github.com/marketplace/models/azure-openai/gpt-4o/playground>

Ethical Considerations

We do not foresee any risks or potential for harmful use arising from our research. Our analyses utilize sentences from a psychology corpus, which consists of licensed data openly accessible for academic use, thereby ensuring both transparency and accountability.

Acknowledgments

We express our gratitude to the individuals who provided valuable feedback during the early stages of this work: Assistant Professor Ehsan Shareghi for his insightful comments; Professor Emeritus Dirk Geeraerts for discussions on the transparency of LLMs and a multidimensional approach to semantic change, including the qualitative dimensions of metaphor and metonymy; Professor Mark Steedman for discussion surrounding the semantic capabilities of LLMs; and Dr Dominik Schlechtweg for his contributions to our understanding of metaphor and metonymy through cognitive theories of similarity and contiguity.

Special thanks go to Philip Baes for his consistent support and insightful discussions on methodological challenges. We also appreciate the discussions with Rokhsana Goworek about the LSC score and XL-LEXEME, and Pierluigi Cassotti, Francesco Periti, and Jader Martins Camboim de Sá for enriching our project’s context through their work on semantic change.

We acknowledge the support of the University of Melbourne’s general-purpose High Performance Computing system, Spartan (Lafayette et al., 2016), which provided the computational resources necessary for efficient embedding encoding using transformer models. We also thank Jeremy Silver from the Statistical Consulting Centre for his valuable feedback on the random effects modelling.

This research was supported by an Australian Government Research Training Program Scholarship and funded, in part, by Australian Research Council Discovery Project DP210103984 and the research program “Change is Key!”, supported by Riksbankens Jubileumsfond (under reference number M21-0021).

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman,

Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#). *arXiv preprint arXiv:2303.08774*.

Naomi Baes, Nick Haslam, and Ekaterina Vylomova. 2024. [A multidimensional framework for evaluating lexical semantic change with social science applications](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1390–1415, Bangkok, Thailand. Association for Computational Linguistics.

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

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

Leonard Bloomfield. 1933. *Language*. Compton Printing Works Ltd.

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

Michel Bréal. 1900. *Semantics: Studies in the Science of Meaning*. WILLIAM HEINEMANN, London.

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

Christian Burgers, Elly A Konijn, and Gerard J Steen. 2016. [Figurative framing: Shaping public discourse through metaphor, hyperbole, and irony](#). *Communication theory*, 26(4):410–430.

Laura Cabello and Uchenna Akujuobi. 2024. [It is simple sometimes: A study on improving aspect-based sentiment analysis performance](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6597–6610, Bangkok, Thailand. Association for Computational Linguistics.

Samuel Cahyawijaya, Holy Lovenia, and Pascale Fung. 2024. [LLMs are few-shot in-context low-resource language learners](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages

- 405–433, Mexico City, Mexico. Association for Computational Linguistics.
- Lyle Campbell. 2013. *Historical Linguistics: An Introduction*, new edition, 3 edition. Edinburgh University Press.
- Pierluigi Cassotti, Stefano De Pascale, and Nina Tahmasebi. 2024a. [Using synchronic definitions and semantic relations to classify semantic change types](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4539–4553, Bangkok, Thailand. Association for Computational Linguistics.
- Pierluigi Cassotti, Francesco Periti, Stefano De Pascale, Haim Dubossarsky, and Nina Tahmasebi. 2024b. [Computational modeling of semantic change](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Tutorial Abstracts*, pages 1–8, St. Julian’s, Malta. Association for Computational Linguistics.
- Pierluigi Cassotti, Lucia Siciliani, Marco DeGemmis, Giovanni Semeraro, and Pierpaolo Basile. 2023. [Xllexeme: Wic pretrained model for cross-lingual lexical semantic change](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1577–1585, Toronto, Canada. Association for Computational Linguistics.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2024. [A survey on evaluation of large language models](#). *ACM Trans. Intell. Syst. Technol.*, 15(3).
- Jiuhai Chen, Lichang Chen, Chen Zhu, and Tianyi Zhou. 2023. [How many demonstrations do you need for in-context learning?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11149–11159, Singapore. Association for Computational Linguistics.
- Jader Martins Camboim de Sá, Marcos Da Silveira, and Cédric Pruski. 2024. [Semantic change characterization with llms using rhetorics](#). *arXiv preprint arXiv:2407.16624*.
- Jader Martins Camboim de Sá, Marcos Da Silveira, and Cédric Pruski. 2024. [Survey in characterization of semantic change](#). *Preprint*, arXiv:2402.19088.
- 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 *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.
- Haim Dubossarsky, Simon Hengchen, Nina Tahmasebi, and Dominik Schlechtweg. 2019. [Time-out: Temporal referencing for robust modeling of lexical semantic change](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 457–470, Florence, Italy. Association for Computational Linguistics.
- Haim Dubossarsky, Daphna Weinsahl, and Eitan Grossman. 2017. [Outta control: Laws of semantic change and inherent biases in word representation models](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1136–1145, Copenhagen, Denmark. Association for Computational Linguistics.
- Dirk Geeraerts. 1997. *Diachronic Prototype Semantics: A Contribution to Historical Lexicology*. Oxford: Clarendon Press.
- Dirk Geeraerts. 2010. *Theories of lexical semantics*. Oxford University Press.
- Mario Giulianelli, Marco Del Tredici, and Raquel Fernández. 2020. [Analysing lexical semantic change with contextualised word representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3960–3973, Online. Association for Computational Linguistics.
- Roksana Goworek and Haim Dubossarsky. 2024. [Toward sentiment aware semantic change analysis](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 350–357, St. Julian’s, Malta. Association for Computational Linguistics.
- Nick Haslam. 2016. [Concept creep: Psychology’s expanding concepts of harm and pathology](#). *Psychological Inquiry*, 27(1):1–17.
- Nick Haslam, Ekaterina Vylomova, Michael Zyphur, and Yoshihisa Kashima. 2021. [The cultural dynamics of concept creep](#). *American Psychologist*, 76(6):1013.
- 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 *Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis*, pages 202–214, Toronto, Canada. Association for Computational Linguistics.
- Simon Hengchen, Nina Tahmasebi, Dominik Schlechtweg, and Haim Dubossarsky. 2021. [Challenges for computational lexical semantic change](#).
- 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](#). *Open Science Framework*.

- 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](#). *Science*, 366(6472):1517–1522.
- 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](#). *Perspectives on Psychological Science*, 17(3):805–826. PMID: 34606730.
- Hajime Kiyama, Taichi Aida, Mamoru Komachi, Toshinobu Ogiso, Hiroya Takamura, and Daichi Mochihashi. 2025. [Analyzing continuous semantic shifts with diachronic word similarity matrices](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 1613–1631, Abu Dhabi, UAE. Association for Computational Linguistics.
- Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski, and Erik Veldal. 2018. [Diachronic word embeddings and semantic shifts: a survey](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1384–1397, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Lev Lafayette, Greg Sauter, Linh Vu, and Bernard Meade. 2016. Spartan performance and flexibility: An hpc-cloud chimera. *OpenStack Summit, Barcelona*, 27(6).
- 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](#). *Bioinformatics*, 36(4):1234–1240.
- Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinneng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, et al. 2024. [Best practices and lessons learned on synthetic data for language models](#). *arXiv preprint arXiv:2404.07503*.
- Daniel Loureiro, Aminette D’Souza, Areej Nasser Muhajab, Isabella A. White, Gabriel Wong, Luis Espinosa-Anke, Leonardo Neves, Francesco Barbieri, and Jose Camacho-Collados. 2022. [TempoWiC: An evaluation benchmark for detecting meaning shift in social media](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3353–3359, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Barbara McGillivray. 2020. [Computational methods for semantic analysis of historical texts](#). In *Routledge International Handbook of Research Methods in Digital Humanities*, pages 261–274. Routledge.
- Raphael Merx, Ekaterina Vylomova, and Kemal Kurniawan. 2024. [Generating bilingual example sentences with large language models as lexicography assistants](#). In *Proceedings of the 22nd Annual Workshop of the Australasian Language Technology Association*, Canberra, Australia. Association for Computational Linguistics.
- George A. Miller. 1992. [WordNet: A lexical database for English](#). In *Speech and Natural Language: Proceedings of a Workshop Held at Harriman, New York, February 23-26, 1992*.
- Saif Mohammad. 2018. [Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words](#). In *Proceedings of the 56th annual meeting of the association for computational linguistics (volume 1: Long papers)*, pages 174–184.
- Saif M Mohammad. 2021. [Sentiment analysis: Automatically detecting valence, emotions, and other affectual states from text](#). In *Emotion measurement*, pages 323–379. Elsevier.
- Charles Egerton Osgood, William H May, and Murray S Miron. 1975. *Cross-Cultural Universals of Affective Meaning*. University of Illinois Press.
- M Sandra Peña and Francisco José Ruiz de Mendoza. 2017. [Construing and constructing hyperbole](#). *Studies in figurative thought and language*, 56:41.
- Francesco Periti and Stefano Montanelli. 2024a. [Lexical semantic change through large language models: a survey](#). *ACM Comput. Surv.*, 56(11).
- Francesco Periti and Stefano Montanelli. 2024b. [Lexical semantic change through large language models: a survey](#). *ACM Computing Surveys*.
- Francesco Periti and Nina Tahmasebi. 2024. [A systematic comparison of contextualized word embeddings for lexical semantic change](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4262–4282, Mexico City, Mexico. Association for Computational Linguistics.
- Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. [WiC: the word-in-context dataset for evaluating context-sensitive meaning representations](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). In *Proceedings of the OpenAI Research Conference 2019*.
- James A. Russell. 2003. [Core affect and the psychological construction of emotion](#). *Psychological Review*, 110(1):145–172.

- Dominik Schlechtweg, Barbara McGillivray, Simon Hengchen, Haim Dubossarsky, and Nina Tahmasebi. 2020. [SemEval-2020 task 1: Unsupervised lexical semantic change detection](#). In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 1–23, Barcelona (online). International Committee for Computational Linguistics.
- Nina Tahmasebi, Lars Borin, and Adam Jatowt. 2018. [Survey of computational approaches to diachronic conceptual change](#). *CoRR*, abs/1811.06278.
- Nina Tahmasebi and Haim Dubossarsky. 2023. [Computational modeling of semantic change](#). *arXiv preprint arXiv:2304.06337*.
- Nina Tahmasebi and Thomas Risse. 2017. [Finding individual word sense changes and their delay in appearance](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 741–749, Varna, Bulgaria. INCOMA Ltd.
- Xuri Tang. 2018. [A state-of-the-art of semantic change computation](#). *Natural Language Engineering*, 24(5):649–676.
- Michael EW Varnum, Nicolas Baumard, Mohammad Atari, and Kurt Gray. 2024. [Large language models based on historical text could offer informative tools for behavioral science](#). *Proceedings of the National Academy of Sciences*, 121(42):e2407639121.
- Ekaterina Vylomova, Sean Murphy, and Nick Haslam. 2019. [Evaluation of semantic change of harm-related concepts in psychology](#). In *Proceedings of the 1st International Workshop on Computational Approaches to Historical Language Change*, pages 29–34.
- Amy Beth Warriner, Victor Kuperman, and Marc Brysbaert. 2013. [Norms of valence, arousal, and dominance for 13,915 english lemmas](#). *Behavior Research Methods*, 45(4):1191–1207.
- Yu Xiao, Naomi Baes, Ekaterina Vylomova, and Nick Haslam. 2023. [Have the concepts of ‘anxiety’ and ‘depression’ been normalized or pathologized? a corpus study of historical semantic change](#). *PloS one*, 18(6):e0288027.
- Heng Yang, Biqing Zeng, Mayi Xu, and Tianxing Wang. 2021. [Back to reality: Leveraging pattern-driven modeling to enable affordable sentiment dependency learning](#). *CoRR*, abs/2110.08604.
- Heng Yang, Chen Zhang, and Ke Li. 2023. [Pyabsa: A modularized framework for reproducible aspect-based sentiment analysis](#). In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM 2023, Birmingham, United Kingdom, October 21-25, 2023*, pages 5117–5122. ACM.
- 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](#). *arXiv preprint arXiv:2311.00237*.
- Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. [Can large language models transform computational social science?](#) *Computational Linguistics*, 50(1):237–291.

A Corpus Counts of Target Terms

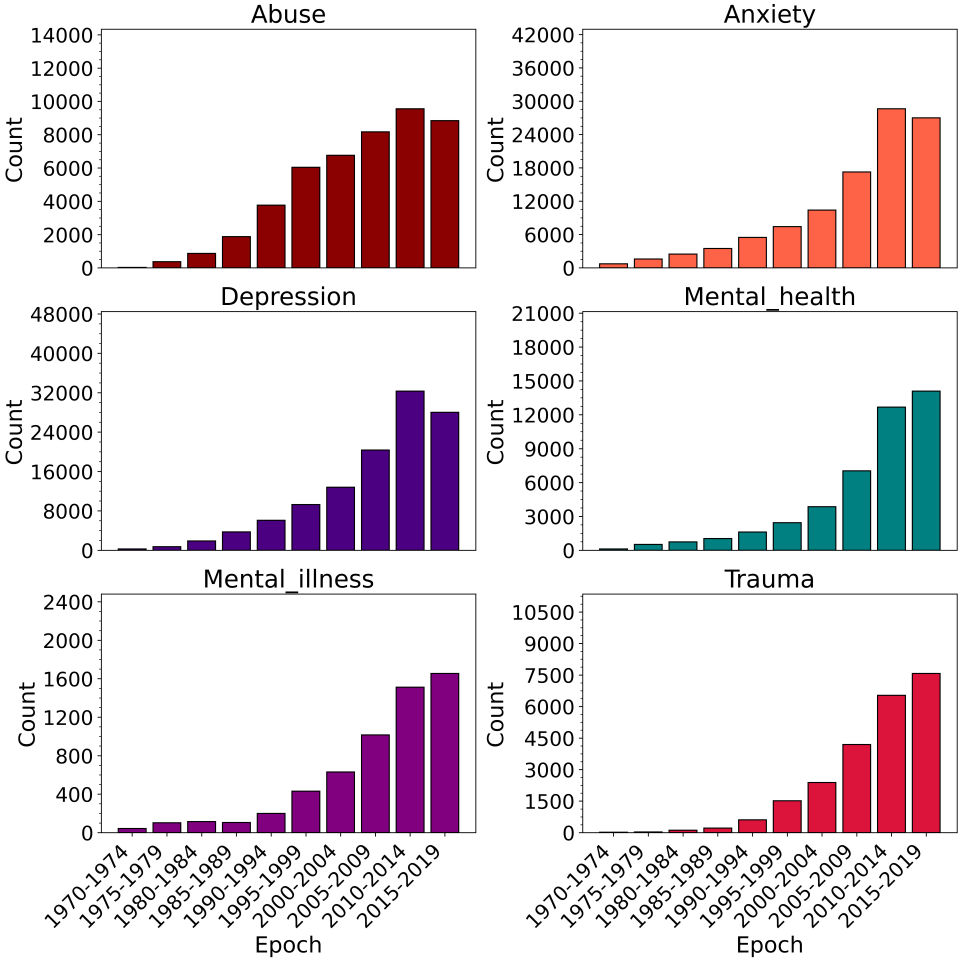


Figure 6: Annual Counts of Sentences where Target Terms Appear in the Psychology Corpus (1970-2019).

B Examples of Synthetic Sentences for each Dimension and Target

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
	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.
Breadth	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.

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 7 and 8 for a breakdown of neutral sentence counts per epoch provided as input to the LLM using the prompts below.

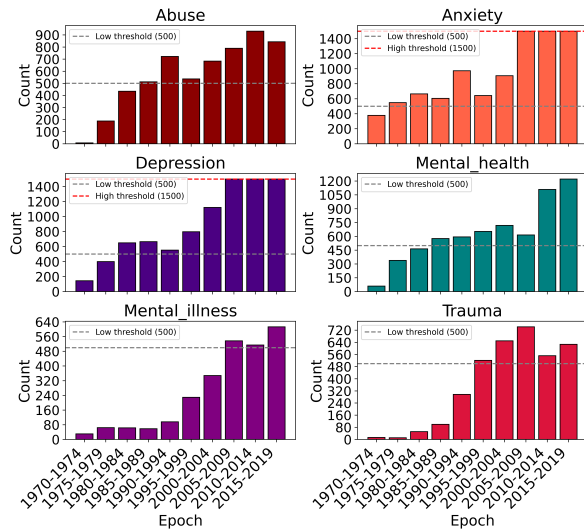


Figure 7: Counts of Neutral Sentences (valence Scores).

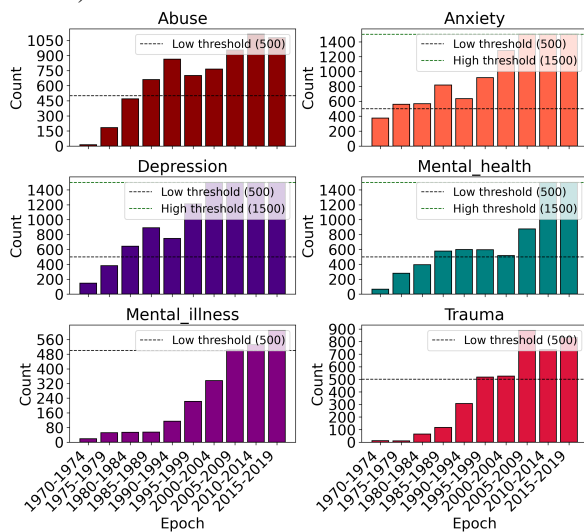


Figure 8: Counts of Neutral Sentences (arousal Scores).

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 at: https://github.com/naomibaes/Synthetic-LSC_pipeline.

Target	Dimension	Neutral	Increase	Decrease	US\$
<i>abuse</i>	Sentiment	5,645	5,645	5,645	17
	Intensity	6,802	6,801	6,801	21
<i>anxiety</i>	Sentiment	9,215	9,213	9,213	28
	Intensity	9,659	9,657	9,657	32
<i>depression</i>	Sentiment	8,828	8,826	8,826	29
	Intensity	10,022	10,020	10,020	35
<i>mental health</i>	Sentiment	6,348	6,348	6,348	21
	Intensity	6,904	6,899	6,899	24
<i>mental illness</i>	Sentiment	2,552	2,552	2,552	9
	Intensity	2,497	2,496	2,496	10
<i>trauma</i>	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. ""

D Demonstration Examples: Synthetic Sentiment

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.

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 wresting 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

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.

Continued on next page

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
	Blasphemy (blasphemy.n.01)	1.07	0.89
Abuse (maltreatment.n.01)	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
	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
Hallucinosi s (hallucinosi s.n.01)	1.05	0.92	
Abulia (abulia.n.01)	0.97	0.80	
Depression (depression.n.01)	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
	Delusion (delusion.n.01)	1.43	0.90
	Hypnosis (hypnosis.n.01)	1.44	0.83
	Anhedonia (anhedonia.n.01)	1.34	0.84
	Agitation (agitation.n.01)	1.32	0.89
	Nervousness (nervousness.n.02)	1.25	0.84
	Cognitive State (cognitive_state.n.01)	1.22	0.85
	State of Mind (state_of_mind.n.01)	1.22	0.80
	Irritation (irritation.n.01)	1.27	0.85
Fugue (fugue.n.02)	1.07	0.86	

Continued on next page

Target (Synset)	Sibling (Synset)	Lin Similarity	Cosine Similarity
	Hallucinosi (hallucinosi.n.01)	1.05	0.89
	Abulia (abulia.n.01)	0.97	0.76
Depression (depression.n.04)	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
	Dolefulness (dolefulness.n.01)	0.84	0.86
Mental Health (mental_health.n.01)	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
	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
	Hallucinosi (hallucinosi.n.01)	0.99	0.88
	Abulia (abulia.n.01)	0.92	0.79
Mental Illness (mental_illness.n.01)	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
	Agitation (agitation.n.01)	1.11	0.88
	Depersonalization (depersonalization.n.02)	1.10	0.88
	Irritation (irritation.n.01)	1.07	0.87
	Morale (morale.n.01)	1.06	0.87
	Nervousness (nervousness.n.02)	1.05	0.80
	Enchantment (enchantment.n.02)	1.05	0.90
	Cognitive State (cognitive_state.n.01)	1.03	0.86
	State of Mind (state_of_mind.n.01)	1.03	0.79
	Elation (elation.n.01)	0.98	0.86
Fugue (fugue.n.02)	0.92	0.89	
Hallucinosi (hallucinosi.n.01)	0.91	0.90	

Continued on next page

Target (Synset)	Sibling (Synset)	Lin Similarity	Cosine Similarity
Trauma (trauma.n.02)	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
	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
	Hallucinosiis (hallucinosiis.n.01)	0.94	0.87
	Abulia (abulia.n.01)	0.88	0.82

Table 11: All Eligible Sibling Terms for Each Target Term with Lin and Cosine Similarity Scores.

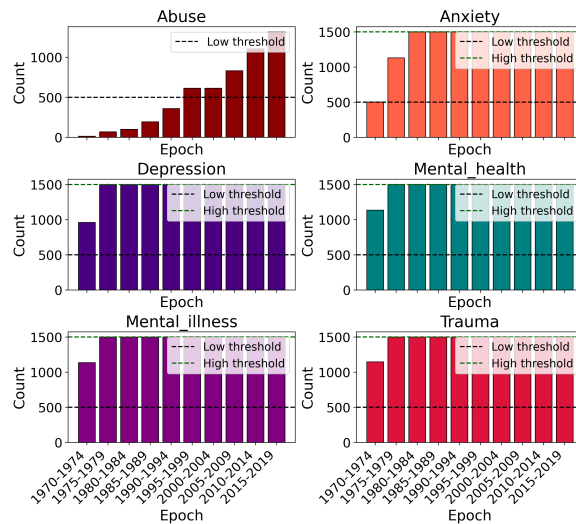


Figure 9: Counts of synthetic sentences (donor-sibling contexts).

Follow this GitHub link to access the ranked lists for each sampling strategy (Boostrapped and Five-Year): https://github.com/naomibaes/LSCD_method_evaluation/tree/main/supplementary_materials/synthetic_breadth_siblings

G Quantifying Lexical Semantic Change

G.1 Semantic Dimensions

G.1.1 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.¹³ 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 3:

$$V_{t_j,k}, A_{t_j,k} = \frac{\sum_{i=1}^{n_{j,k}} w_{i,j,k} x_{i,j,k}}{\sum_{i=1}^{n_{j,k}} w_{i,j,k}} \quad (3)$$

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-v1.1¹⁴ identifies sentiment associated with particular aspects of an entity within text (here, the target term).¹⁵ We adapt it to produce continuous sentiment scores

¹³Tokenization, lemmatization, stop-word removal using “en_core_web_sm” (<https://spacy.io/models/en>)

¹⁴yangheng/deberta-v3-base-absa-v1.1 (184M model params): <https://huggingface.co/yangheng/deberta-v3-base-absa-v1.1>

¹⁵It was trained on restaurant and laptop reviews (Cabello and Akujuobi, 2024; Yang et al., 2021, 2023).

which reflect the model’s confidence in positive sentiment associated with the target, ranging from 0 (fully negative) to 1 (fully positive).¹⁶

G.1.2 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 model ‘all-mpnet-base-v2’¹⁷ is used to generate these embeddings. The Breadth score, B , is derived by averaging the cosine distances, δ , across all unique pairs of sentence embeddings within each iteration, and then averaging these scores across all iterations within each bin, as in Equation 4:

$$B_{t_j} = \frac{1}{I_j} \sum_{k=1}^{I_j} \left(\frac{2}{N_k(N_k - 1)} \sum_{i=1}^{N_k-1} \sum_{j=i+1}^{N_k} \delta(s_{i,k}^{t_j}, s_{j,k}^{t_j}) \right) \quad (4)$$

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

The sentence transformer “all-mpnet-base-v2” (MPNet) from Cassotti et al. (2023) is compared with the SOTA transformer “XL-LEXEME”¹⁸ (XLL). MPNet pools tokens to produce sentence embeddings, which dilutes word-level information, whereas XLL employs a bi-encoder focused on word-specific attention,¹⁹ using polysemy as a proxy for semantic divergence during training (WiC; Pilehvar and Camacho-Collados, 2019).

G.2 General Lexical Semantic Change:

General lexical semantic change is evaluated using the LSC score shown in equation 5.

$$LSC_i(s_i^{t_0}, s_i^{t_1}) = \frac{1}{N_i^2} \sum_{m=1}^{N_i} \sum_{n=1}^{N_i} \delta(s_{m,i}^{t_0}, s_{n,i}^{t_1}) \quad (5)$$

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

¹⁶The sentiment score is calculated as follows: $0 \times \text{negative_prob} + 0.5 \times \text{neutral_prob} + 1 \times \text{positive_prob}$.

¹⁷Microsoft pretrained network (109M model params) <https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

¹⁸XL-LEXEME (~550M parameters): <https://huggingface.co/pierluigic/xl-lexeme>

¹⁹Only the first target occurrence is attended to.

H Linear Mixed-Effects Modeling for SIB Dimensions

To evaluate the effects of synthetic injection, we use mixed-effects linear regression models with fixed effects for `injection_level` and random intercepts for `target`. This approach accounts for repeated observations per target and treats the grouping variable as sampled from a population of effects. The grouping variable `target` had six levels, a sufficient number for reliable variance estimation.

Model Structure and Rationale: We first fit a **null model** with only random intercepts to estimate the intraclass correlation coefficient (ICC), which quantifies variance attributable to group-level clustering. ICC values above 0.05 suggest meaningful clustering and justify random intercepts. Next, we fit a **random intercepts model** including `injection_level` as a fixed effect, expressed in Equation 6.

$$y_{ij} = \beta_0 + \beta_1 \cdot \text{injection_level}_{ij} + u_j + \epsilon_{ij} \quad (6)$$

where y_{ij} is the outcome for observation i in group j , β_0 is the fixed intercept, β_1 is the fixed slope for injection level, $u_j \sim \mathcal{N}(0, \sigma_u^2)$ is the group-specific random intercept, and $\epsilon_{ij} \sim \mathcal{N}(0, \sigma_\epsilon^2)$ is the residual error.

We also tested a **random slopes model** allowing `injection_level` slopes to vary by `target`, but with only six groups, these models often failed to converge or showed negligible slope variance.

Model Comparison and Selection: Likelihood ratio tests (LRTs) favored the random slopes model across dimensions, but due to convergence issues and minimal slope variance, we selected the more stable and parsimonious **random intercepts model** for interpretation.

Model Results Summary: We fit separate mixed-effects models for each outcome (Valence, Arousal, Breadth), with `injection_level` as a fixed effect and random intercepts for `target`. Outcomes were standardized (z -scores) to ensure coefficient comparability. Model assumptions (linearity, homoscedasticity, normality) were verified via standard diagnostics. Table 12 reports standardized coefficients (β), 95% confidence intervals, and p -values. All models showed significant, directionally consistent effects. Random intercept variance (σ^2) indicates the extent of baseline score differences across target terms.

Score	IL	β (95% CI)	p	σ^2
Valence	+	0.611 (0.567, 0.654)	<.0001	0.732
	-	-0.305 (-0.357, -0.253)	<.0001	1.059
Arousal	+	0.638 (0.585, 0.691)	<.0001	0.681
	-	-0.643 (-0.698, -0.588)	<.0001	0.672
Breadth	+	0.429 (0.317, 0.541)	<.0001	0.845

Table 12: Results of the Final Mixed Linear Models Predicting Dimension Scores from Injection Levels.

Note: IL = Injection level (- = increased variation, + = decreased variation). β = Standardized coefficient with 95% confidence interval. p -values test the null hypothesis that the coefficient is zero. σ^2 = Random intercept variance. Number of Observations (Groups) 36 (6).

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

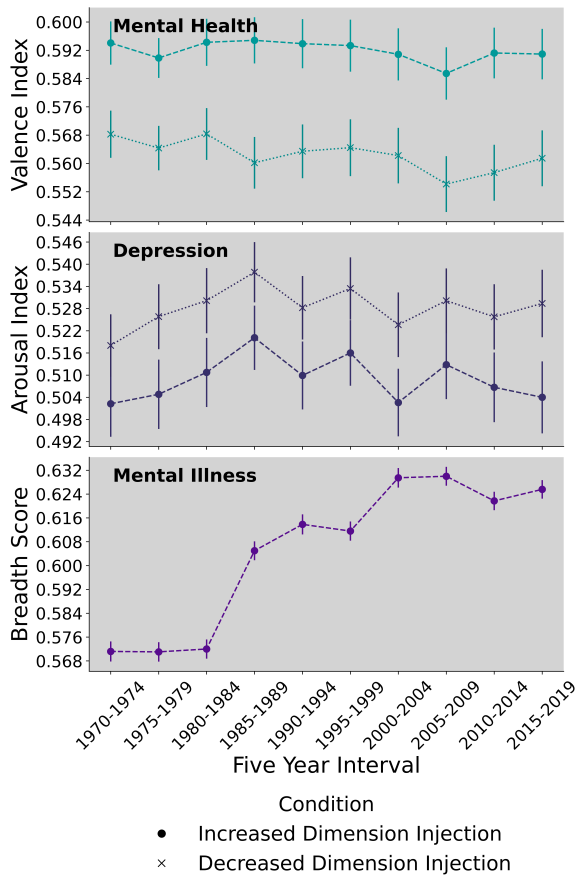


Figure 10: SIB Scores (\pm SE) by 50% Injection Levels and Conditions: Control Setting for Five-Year Samples.

J Alternative LSC Detection Methods: Results for Bootstrapped Settings

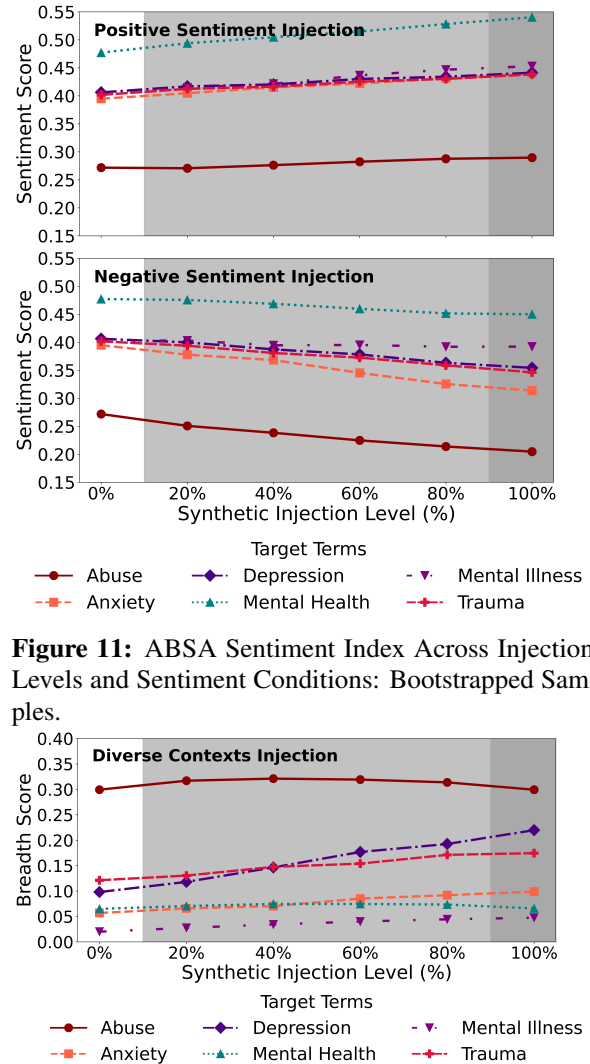


Figure 11: ABSA Sentiment Index Across Injection Levels and Sentiment Conditions: Bootstrapped Samples.

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

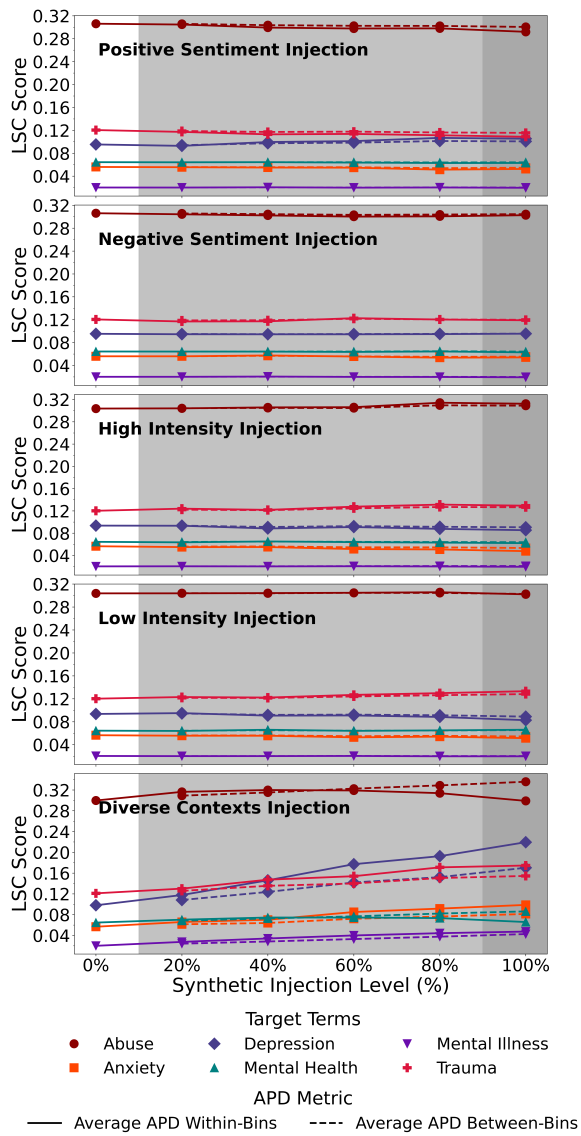


Figure 13: LSC Scores (APD Between-Bins and APD Within-Bins) Across Injection Levels and SIB Conditions: Bootstrapped Samples.