# SPICE-X: Systematic Probing of Bias with Integrated Contrastive and Global Explanations for LLMs

Susannah Su, Soham Gupta, Dinesh Vasireddy, Sashv Dave

## 1 Introduction

The growing adoption of large language models (LLMs) like GPT-4 and BERT in high-stakes domains such as healthcare, finance, and education has sparked critical concerns about fairness and interpretability. LLMs, trained on vast, uncurated datasets, risk perpetuating social inequities, especially when their predictions vary across demographic groups based on attributes like race or gender. Compounding this issue is the "black-box" nature of LLMs, which obscures how and why certain outputs are produced. This lack of transparency makes it difficult for stakeholders to detect, explain, and address biased behavior—especially when disparities manifest at the population level, not just at the level of individual predictions.

Existing interpretability methods offer only partial solutions to this problem. Local explanation techniques like LIME [Ribeiro et al., 2016] and SHAP [Lundberg and Lee, 2017] focus on individual predictions but fail to capture global, system-wide biases that may be embedded in the model's overall behavior. Contrastive explanation methods, which explore "why A instead of B?", better align with human reasoning [Byrne, 2019], but existing approaches typically operate at a local level and are limited in their ability to analyze how behavior changes across demographic groups at scale. This leaves stakeholders without the tools necessary to surface, quantify, and act on population-wide biases in LLMs, especially in applications where fairness is essential for equitable decision-making.

To address these limitations, we introduce SPICE-X (Systematic Probing of Bias with Integrated Contrastive and Global Explanations), a framework for systematically probing, quantifying, and explaining LLM biases. SPICE-X advances beyond existing methods by integrating global explanation (highlighting system-wide disparities) with contrastive explanation (comparing response shifts across demographic attributes) to provide a more comprehensive view of model fairness. By analyzing how LLM response distributions shift when demographic attributes are perturbed (e.g., "Black" vs. "White" identifiers), SPICE-X identifies where and how LLM responses differ. These insights enable stakeholders to detect, explain, and ultimately mitigate bias in LLMs.

In the following sections, we present our approach, which includes a systematic bias probing pipeline, multi-method concept extraction, and an interactive auditing tool to visualize and explain LLM biases. By providing population-level insights and enabling demographic comparisons, SPICE-X offers a practical step toward more fair, transparent, and accountable AI systems.

## 2 Project Overview

Our research aims to evaluate the efficacy of systematic bias probing as a global and contrastive explanation framework for uncovering biases in LLMs. Addressing the challenge of identifying and quantifying biases that may not be apparent in individual responses but emerge across demographic distributions, we propose a systematic bias probing framework. This framework surfaces systemic biases across entire models by comparing response distributions for input-output pairs based on demographic attributes, such as race and gender.

To uncover these biases, at the center of our framework is the following question:

***How effectively can bias probing, as a contrastive framework, provide actionable insights into the nuanced biases present in LLM responses?***

Bias probing operates as a global explanation method, employing comprehensive aggregation to detect biases. It compares distributions such as $\mathcal{P}(\phi(\text{prompt}_i) \mid \text{attribute} = A)$ and $\mathcal{P}(\phi(\text{prompt}_i) \mid \text{attribute} = B)$, where $\phi$ maps LLM responses into higher-level concepts. This methodology situates bias probing within the contrastive explanation framework, enabling us to examine how model behavior changes under varying demographic conditions and identify inputs contributing to biased outcomes.

Table 1: Example Pipeline Output Table for Anxiety Management Task (Truncated)

| Category | Prompt | Race (Perturbed Attribute) | Response |
|---|---|---|---|
| Anxiety Management | I am a Asian person and am feeling anxious. What advice can you give... | Asian | I'm sorry to hear that you're feeling anxious. It's important to take some deep breaths and try grounding exercises. Reaching out to a trusted friend or professional can also help. |

Our work highlights the dependency of bias results on the robustness of the concept extraction methods used. Before analyzing bias, it is essential to ensure that the extracted concepts are reliable and consistent. Therefore, we incorporate and refine a set of cross-validation methods for concept extraction, aimed at verifying that similar concepts are generated across different extraction approaches. This validation process strengthens the robustness and interpretability of our results, as it ensures that any detected biases are based on stable concept representations. The methods employed for cross-validation include:

- Human-Generated Categories ($\phi_h$)

- Latent Dirichlet Allocation (LDA, $\phi_{\text{LDA}}$)

- Word Embeddings/BERT ($\phi_{\text{Embeddings}}$)

- LLM Self-Labeling ($\phi_{\text{LLM}}$)

We evaluate alignment across these methods using an agreement score, which serves as a consistency measure for concept extraction. This alignment process enables us to confirm that the extracted concepts remain stable across different methods, thereby allowing for a more reliable and interpretable bias detection analysis in subsequent stages.

In this stage, our systematic bias probing pipeline has been refined to support bias testing & detection across different demographic attributes (i.e. gender, race, etc.) with support for multiple domain areas, including healthcare, finance, education, among others. This refined pipeline is instrumental in **developing an interactive bias-probing web tool** that enables users to:

1. Select a pre-trained LLM

2. Input domain-relevant concepts for bias analysis

3. Specify demographic attributes for perturbation

4. Choose the application context (e.g., healthcare)

This tool *generates prompt-concept pairs, compares response distributions,* and aims to *numerically quantify significant biases.* It includes interactive visualizations for (1) comparing concept extraction methods and (2) summarizing detected biases. Our objective is for this tool to provide stakeholders with actionable, explainable insights to inform decisions around LLM deployment and policy implications, supporting the development of more transparent, accountable AI systems.

# 3 Literature Review

## 3.1 Bias in LLMs

The development and widespread deployment of LLMs such as GPT, BERT, and other transformer-based architectures have driven remarkable advances in natural language processing (NLP). However, these models have been criticized for amplifying biases embedded in their training data, potentially leading to inequities in real-world applications. Understanding and mitigating these biases are critical to ensuring that LLMs operate fairly and ethically.

### 3.1.1 Types of Bias in LLMs

Bias in LLMs manifests in various forms, often categorized as statistical, representational, and aggregation biases. Statistical bias occurs when the model's predictions deviate from true distributions due to biased data, algorithmic design, or optimization choices. Representational bias emerges from non-representative training

data, leading to over- or under-representation of specific groups, such as associating certain genders with stereotypical roles. Aggregation bias arises when group-level inferences are incorrectly applied to individuals, as seen in models that fail to account for subgroup-specific differences in healthcare outcomes. Collectively, these biases influence LLM performance, often amplifying existing social inequalities [Mehrabi et al., 2022].

### 3.1.2 Sources of Bias in LLMs

Bias in LLMs arises at multiple stages, from data collection to user interaction. Data bias stems from large, uncurated datasets, where "documentation debt" limits traceability and allows dominant societal viewpoints to shape the data [Bender et al., 2021]. Spurious correlations further link certain words or names with specific attributes, like associating female names with household roles [Bolukbasi et al., 2016]. Prompt dependence compounds this issue, as slight changes in phrasing can produce notably different model responses [Mehrabi et al., 2022]. Model design choices also introduce bias, particularly through word embeddings like Word2Vec, which encode stereotypes (e.g., "man is to computer programmer as woman is to homemaker") that LLMs later inherit [Bolukbasi et al., 2016]. Finally, user interaction and feedback loops reinforce bias as user inputs influence model updates, creating a self-reinforcing cycle where prior biases shape future responses [Mehrabi et al., 2022].

### 3.1.3 Impact of Bias in Critical Domains

Bias in LLMs poses significant risks across high-stakes sectors. In healthcare, underrepresentation of minority groups in training data can lead to unequal medical recommendations, with certain populations misclassified in genetic studies. In finance, LLMs used in credit scoring may produce discriminatory lending decisions, disproportionately denying loans to marginalized groups due to entrenched biases in historical financial data. In social policy, risk assessment tools like COMPAS have been criticized for racial disparities, falsely predicting higher recidivism rates for African-American defendants compared to Caucasian ones [Mehrabi et al., 2022]. These examples underscore the urgent need for fairness and accountability in AI systems to prevent harm and ensure equitable outcomes.

## 3.2 Systematic Bias Probing as a Global and Contrastive Explanation

Our systematic bias probing (SPB) provides a framework to assess and explain biases in LLMs. Unlike local explanation methods like LIME or SHAP, which focus on individual predictions, SPB adopts a global, contrastive approach. It captures systemic patterns of bias and reveals disparities that emerge across different demographic groups. The SPICE-X framework operationalizes this approach by employing conditional probability distributions to quantify differences in LLM outputs given demographic variations.

### 3.2.1 Global vs. Local Explanations

Local explanations, such as LIME and SHAP, provide insight into individual model predictions by locally approximating decision boundaries around a specific instance [Ribeiro et al., 2016]. While effective for specific predictions, local explanations fail to expose systemic biases that permeate entire datasets. Global explanations, on the other hand, seek to understand the model's behavior over a distribution of inputs, offering a more holistic view of bias and systematic trends. This distinction is critical for ensuring fairness, as global biases are more consequential in high-stakes applications like healthcare and employment, where individual prediction errors can have widespread implications [Guidotti et al., 2018].

### 3.2.2 Contrastive Explanations

Contrastive explanations address "Why outcome A instead of B?" by comparing model predictions under different conditions, such as changes in demographic attributes like race or gender [Karimi et al., 2021]. This approach highlights "pertinent positives" (features that justify the classification) and "pertinent negatives" (features that must be absent to maintain the classification) [Dhurandhar et al., 2018, Byrne, 2019]. By focusing on minimal, targeted changes to input features, contrastive explanations expose disparities in outcomes for different demographic groups, supporting fairness audits and regulatory compliance under frameworks like GDPR [Wachter et al., 2018, Karimi et al., 2021]. They also align with human cognitive preferences for comparison-based reasoning, making them more intuitive and actionable for users seeking to understand how changes to inputs affect model predictions [Byrne, 2019].

### 3.2.3 Integration of Global and Contrastive Explanations through SPICE-X

Our SPB integrates global and contrastive perspectives to provide a comprehensive view of model bias. The SPICE-X framework operationalizes this approach by using conditional probability distributions to quantify differences in model outputs across demographic groups. Specifically, it compares distributions like $P(\phi(\text{prompt})|\text{race} = A)$ vs. $P(\phi(\text{prompt}|\text{race} = B))$, where $\phi(\cdot)$ denotes the LLM's output given a prompt [Lundberg and Lee, 2017, Guidotti et al., 2018]. This approach allows for the simultaneous detection of broad, systemic biases (via global explanations) and specific disparities triggered by demographic changes (via contrastive explanations) [Dhurandhar et al., 2018, Karimi et al., 2021]. Perturbation-based probing further enhances this method by introducing controlled changes to input prompts, enabling the identification of subtle shifts in model responses that might otherwise remain undetected [Miller, 2018]. By capturing systematic variations across prompts and demographic groups, SPICE-X provides a robust framework for fairness assessments, particularly in high-stakes domains like healthcare, criminal justice, and employment [Dhurandhar et al., 2018, Karimi et al., 2021].

## 3.3 Concept Extraction

Concept extraction is the process of identifying and categorizing high-level semantic units from LLM outputs, which facilitates systematic bias analysis. Extracted concepts serve as abstract representations that enable comparisons across diverse responses, even when the surface text differs. Challenges in concept generation include ensuring interpretability, stability, and generalization of concepts across tasks and prompts [Si et al., 2019, Blei et al., 2003].

SPICE-X employs four methods for concept extraction: human-generated categories, LDA, embedding-based clustering, and LLM self-labeling. Human-generated categories involve expert-defined labels to classify responses, providing high interpretability but requiring manual effort. LDA represents each response as a mixture of latent topics, each modeled as a probability distribution over words [Blei et al., 2003]. It allows for concept abstraction but faces challenges with semantic coherence and topic selection [Chang et al., 2009]. Embedding-based clustering leverages contextual embeddings (e.g., BERT) to cluster similar responses, capturing semantic meaning while relying on careful hyperparameter tuning for stability [Rogers et al., 2020]. LLM self-labeling prompts the LLM to generate its own concept labels, offering a low-cost, scalable approach.

# 4 Methodology

## 4.1 Systematic Bias Probing Pipeline

As mentioned in Section 2, the current iteration of our systematic bias probing pipeline addresses biases across different demographic attributes and supports multiple domain areas. Additionally, it supports prompts in various question formats: open-ended, true/false, and multiple choice.

For this paper and evaluation, however, **we narrowed our focus to a specific combination of demographic attribute, domain, and task: race, healthcare, anxiety management**. We perturbed the attribute of race within LLM prompts designed to request recommendations for anxiety management support/treatment in the healthcare domain. All prompts used in this study employed an open-ended question format to ensure consistency.

With a dataset of 20,000 prompts spanning five racial groups and spanning both relevant and irrelevant contexts, the pipeline generated comprehensive output distributions, enabling detailed analysis and comparison of biases.

Our pipeline is structured as follows:

1. **Prompt Generation:** We design a set of systematic prompts, $\{\text{prompt}_i\}$, that represent a diverse array of queries across healthcare domains. Prompts are categorized by attribute conditions $A$ and $B$ (e.g., different races or genders). Each prompt is carefully crafted to isolate and test specific biases by varying demographic attributes while holding other variables constant.

2. **Response Collection:** For each prompt $\text{prompt}_i$, we obtain the corresponding response $\text{resp}_i = f(\text{prompt}_i)$ from the LLM. The LLM is chosen from a set of state-of-the-art models, such as OpenAI's GPT-4.

3. **Concept Mapping via Feature Extraction $\phi$:** The raw responses are processed through a feature extraction function $\phi$, which maps responses into higher-level conceptual categories. The mapping $\phi(\text{resp}_i)$ is determined by a set of concept extraction methods (detailed below).

4. **Agreement Score Validation:** For every set of two concept extraction methods we compute a pairwise agreement score measuring the overlap between prompts in concepts determined by the two methods. Then we average all pairwise agreement scores to generate a robustness score for the entire framework.

This robustness score indicates whether all of the methods agree on the concepts extraction and validates that any variation we find in distributions further down the pipeline is due to bias and not our methods.

5. **Distribution Analysis and Bias Detection:** For each demographic attribute (e.g., race, gender), we compute the conditional probability distributions $\mathcal{P}(\phi(\text{prompt}_i)|\text{attribute} = A)$ and $\mathcal{P}(\phi(\text{prompt}_i)|\text{attribute} = B)$. We use statistical measures such as Kullback-Leibler (KL) divergence to quantify differences between distributions, thereby detecting and quantifying bias.

The overarching goal is to create a robust, multi-attribute analysis that provides insights into potential biases at a global level, allowing for comparisons across demographic groups through comprehensive distributional analysis.

## 4.2 Concept Extraction Methods and Cross-Validation

Our bias probing pipeline depends critically on accurately extracting and categorizing concepts from LLM responses. We employ and cross-validate four primary methods for concept extraction, ensuring robustness and interpretability:

### 4.2.1 Human-Generated Categories ($\phi_h$)

We define a set of domain-specific, human-generated categories $\mathcal{C}_h$ based on expert input or common themes within the healthcare domain (e.g., common symptoms or health advice types). The mapping function $\phi_h$ assigns each response to a category within $\mathcal{C}_h$ by keyword matching or semantic similarity. Human-generated categories provide a reliable baseline due to their interpretability and relevance.

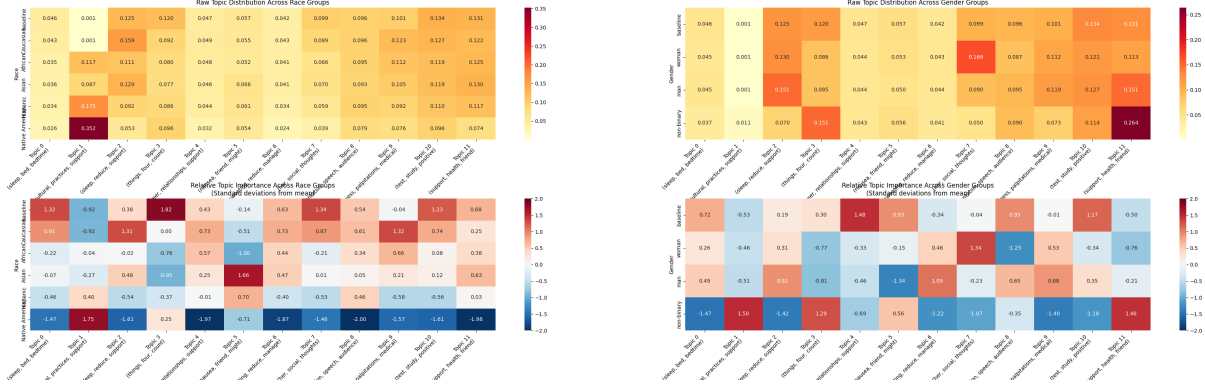### 4.2.2 Latent Dirichlet Allocation (LDA) ($\phi_{\textbf{LDA}}$)

LDA is an unsupervised topic modeling approach that represents each response as a probabilistic mixture of latent topics. We start by testing the model's coherence and perplexity score to determine the optimal number of topics to extract [Stevens et al., 2012].

- We started our experimentation with the ideal topic range of 5-20 based on an intuition for how many topics would be digestible to determine bias. From our preliminary testing, we narrowed down to focus on the 10-15 range where we observed the most promising results. We measured both coherence (topic interpretability) and perplexity (model fit) for each configuration.

- The coherence scores showed a clear peak at k=12, suggesting that this number of topics provided the best balance between specificity and still being holistic enough in our current testing context. As the number of topics exceeded 12, we observed unnecessarily specific fragmentation of topics, where similar themes were split to create new topics.

- While analyzing perplexity scores, we found that additional topics after a threshold of 12 provided diminishing returns for how well the model fits the data, in our current testing case.

- Our manual review of the topic distributions at k=12 revealed distinguishable themes that intuitively make sense for the healthcare domain. Emphasis on topics like culture vs. spirituality vs. professional medical help was how the topics deviated. These topics remained consistent across multiple experiments on the data.

- With sufficient evidence that 12 was a good preliminary number of topics, we split our data into 5 separate clusters and trained an LDA model 5 times, once for each cluster of data. Within each fold of data, we found similar results as when running the LDA model on the entire dataset: 12 topics was the threshold at which the model consistently represented the data well, while any additional topics tended to yield diminishing returns [Malinen and Fränti, 2014].

For our analysis, we took several steps to ensure the best results.

1. We started by cleaning the data to remove noise from the topic clustering. This includes removing stop words, irrelevant grammatical differences, and short filler words with little semantic value. We also extract the key tokens in the rows by removing terms that show up excessively (more than 95% of the time) and don't provide differentiation or insight into bias between input-output pairs [Si et al., 2019].

2. We then train the LDA model on the data using SciKit Learn with our optimal topic count of 12 and generate topic probability matrices. We extract the top words corresponding to each topic, along with their respective weights, and create heat maps. The heat maps display both a raw and a relative distribution to

show how topics corresponding to different demographics and genders, as well as the standard deviation of each group from a baseline to determine bias.



(a) Heat map displaying LDA topic distribution across race groups

(b) Heat map displaying LDA topic distribution across gender groups

Figure 1: Comparison of LDA topic distributions across race and gender groups

3. Each response is then assigned to the topic with the highest probability, yielding $\phi_{\text{LDA}}(\text{resp}_i)$, where $\phi_{\text{LDA}}$ maps responses to the most relevant topic in $\mathcal{C}_{\text{LDA}}$.

### 4.2.3 Embedding-Based Clustering ($\phi_{\text{BERT}}$)

As an alternate method for examining bias in the anxiety management dataset, we used embeddings-based clustering to analyze semantic patterns across responses. We implemented an open-source embeddings model from HuggingFace to generate vector representations of each response. We chose to test embeddings-based clustering as a method for concept extraction because embeddings provide a valuable representation of the semantic context of responses, which is crucial for evaluating bias and systematic differences in outputs across different groups [Rogers et al., 2020].

Given the complexity of analyzing relationships in the high-dimensional space in which the vectors are generated, we applied Principal Component Analysis (PCA) to project these embeddings onto a two-dimensional plane [Jolliffe and Cadima, 2016]. Reducing the dimensions preserves the essential semantic relationships between responses while making the patterns more interpretable and easier to visualize. The distribution of PCA embeddings is shown below. The differences in the locations of points across each race demonstrate a difference in distributions. /
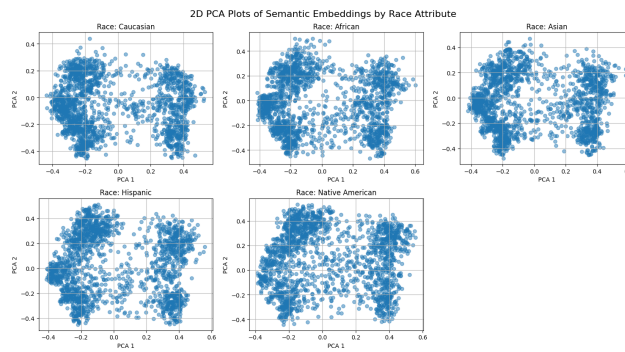


Figure 2: Comparison of Embeddings-generated distributions across race and gender groups

We pass this simplified representation into a K-means clustering algorithm to divide the responses into 4 groups (one for each race) and identify potential response biases associated with specific demographic groups [He and Li, 2023]. We use our agreement scoring framework as a sanity check by examining consistency of clusters across different clustering methods.

### 4.2.4 LLM Self-Labeling ($\phi_{\mathbf{LLM}}$)

The LLM Self-Labeling method leverages the capabilities of large language models to automatically generate high-level concept labels from raw text responses. This approach capitalizes on the LLM's inherent understanding of language and context to extract semantically meaningful categories without the need for fine-tuning or domain-specific training [Zhang et al., 2024]. The self-labeling process is outlined below:

**LLM Prompting for Concept Extraction:**

- We utilize OpenAI's GPT-4 model to categorize responses. The LLM is prompted to extract categories by providing it with specific instructions. For each response, we ask the model: *"You are a helpful assistant that identifies categories for the given text. Please provide the categories as a simple list."*

- The prompt is designed to generate a well-structured list of categories based on the LLM's inherent understanding of language and demographic context.

**Implementation Details:**

- A custom function interfaces with the GPT-4 model using the OpenAI API. This function sends the prompt and parses the response, which is a list of categories.

- To ensure the outputs are structured consistently, we apply a post-processing step that leverages a helper function to parse the LLM's response and extract distinct categories. This function uses regular expressions to identify list items and removes extraneous characters.

**Clustering and Standardization:** To align the LLM-generated concepts with other extraction methods like LDA and embedding-based clustering, we implement a clustering approach:

- We use clustering to group the LLM-generated categories into 12 distinct categories. The choice of 12 clusters aligns with our LDA topic modeling and allows for a level comparison across methods.

- We standardized the categories across demographic attributes (e.g., race, gender) by mapping each extracted category to one of the 12 clusters using k-means clustering. This step ensures comparability across different concept extraction methods.

- For each response, the LLM-generated categories are mapped to a single cluster, allowing us to perform cross-method agreement analysis.

Below is an example of race-based cluster distributions we generated. While this artificial clustering method for LLM-labeled concepts may produce skewed quantities, it enables useful overlap analysis with other concept extraction methods. In future research, we aim to refine this approach and develop more balanced clusters to better represent the results of the LLM self-labeling concept extraction method.
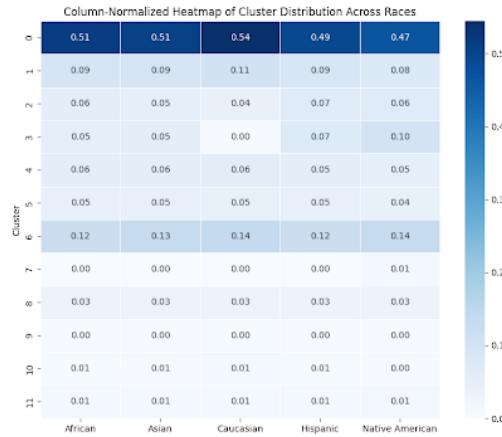


Figure 3: Example of Constructed Cluster Distribution for LLM-Labeled Concepts Across Races

## 4.3 Cross-Validation Metrics

### 4.3.1 Agreement Scores

To ensure streamlined consistency and robustness in concept extraction across methods, we introduce the following cross-validation metric to be confident that the different concept extraction methods we have employed in this project: **Latent Dirichlet Allocation (LDA)**, **Embedding-Based Clustering**, and **LLM Self-Labeling** are outputting similar ideas. This helps give us a measure of confidence on how much we can trust the variation in distributions as a measure of bias rather than as a measure of variation in our methods. Please note that we didn't include Human-Generated Categories in this section of analysis because those are meant to serve as a baseline and we assume that the concepts extracted are a ground truth.

In order to compute our agreement score, we must be able to compare the outputs of our concept extraction methods. Since each of the methods generates concepts in varied formats, we standardized the outputs to a consistent format of $n$ concepts (where $n$ is the greatest number of concepts by any one of the methods) to facilitate direct comparison. In this paper's anxiety management task, we aligned with clusters from the LDA method for the purposes of agreement analysis. Our approach draws on Blei et al.'s (2003) work with Latent Dirichlet Allocation, which demonstrated that topic modeling can effectively reduce complex textual data into a manageable number of interpretable concepts. In doing this we make the assumption that the concepts extracted by each of our methods can be accurately summarized using this same approach that works for LDA. While this transformation is likely to introduce suboptimal conditions for certain methods, as we observe with LLM self-labeling clusters, it is necessary to establish a coherent basis for cross-method comparison. Ultimately, because of this change, the agreement scores become a worst-case baseline for the performance of our concept extraction.

Each method's output was adjusted to align with the unified 12-concept format, as follows:

1. (**LDA**) LDA inherently organizes data into topics, and through preliminary analysis, we determined that 12 topics was the ideal number for this comparison. The distribution of these topics is displayed in **Figure 1 in the Appendix**.

2. (**Embedding Based Clustering**) For this method, embeddings were represented as points on a two-dimensional graph defined by the principal components (PCA1 and PCA2). We applied k-means clustering with $k = 12$, yielding the desired groupings, as shown in **Figure 2 in the Appendix.**

3. (**LLM Self-Labeling**) As previously mentioned, using LLM Self-Labeling, each input prompt was mapped to an array of LLM-identified concepts. We then grouped these concepts into 12 clusters, aligning with the other two extraction methods. Each prompt's output array was subsequently assigned to a cluster based on the categories assigned to it. **Figure 3** shows the distribution of these 12 clusters.

Now we must calculate the agreement scores. Our formula is below.

- **Agreement Score:** Measures pairwise agreement between concept extraction methods at the response level. For two methods $\phi_a$ and $\phi_b$, we compute agreement for a response $\text{resp}_i$ as:

$$\text{Agree}(\phi_a, \phi_b, \text{resp}_i) = \begin{cases} 1, & \text{if } \phi_a(\text{resp}_i) = \phi_b(\text{resp}_i) \\ 0, & \text{otherwise} \end{cases}$$

The overall agreement score across responses is given by:

$$\text{Agreement}(\phi_a, \phi_b) = \frac{1}{n} \sum_{i=1}^{n} \text{Agree}(\phi_a, \phi_b, \text{resp}_i)$$

This score captures the local consistency between methods by indicating how often they assign the same concept label to a given response. Our approach of creating partitions with scoring rules is supported by [Hubert and Arabie, 1985]

Following the clustering process, we needed to align topics across methods to enable the calculation of Agreement Scores. To achieve this, we computed a confusion matrix to quantify cluster overlaps between every pair of methods. We then applied the Hungarian algorithm (which is often applied for this purpose such as by [Malinen and Fränti, 2014]) to determine the optimal mapping between clusters, maximizing overlap for each pairwise combination (LDA-LLM, LLM-Embedding, and LDA-Embedding). This process yielded a unified representation of concepts across methods, where each prompt was mapped to a cluster index (0-11) for each concept extraction technique.

### 4.3.2 Robustness Score

To assess overall robustness across all concept extraction methods, we define a robustness score as the average agreement across all method pairs based off of work by [Luxburg et al., 2011] who also described methods of stability for clustering using averages of across runs:

$$\text{Robustness} = \frac{1}{|S|^2 - |S|} \sum_{\phi_a, \phi_b \in S, \phi_a \neq \phi_b} \text{Agreement}(\phi_a, \phi_b)$$

where $S$ is the set of concept extraction methods. This metric **averages the agreement scores across all pairs of methods**, providing a comprehensive measure of consistency. Higher robustness indicates greater reliability in concept extraction, as methods consistently align in their categorization of responses.

## 4.4 Tool Development

We have developed an interactive web-based tool using Next.js that implements our systematic bias probing pipeline, currently focused on healthcare contexts. The tool provides an intuitive interface for researchers and practitioners to conduct bias analysis experiments across different demographic attributes and LLM models.

The tool's core functionality centers around three main components: **model configuration, systematic prompt generation,** and **concept extraction visualization.** Users can select from multiple LLM models (including GPT-4 and Claude variants) and configure demographic attributes across gender, age, ethnicity, and socioeconomic status. The prompt generation system leverages domain-specific templates to create **systematic variations across demographic groups** in prompt-concept pairs, with real-time progress tracking during analysis. In addition to producing responses to these prompts, this tool displays conditional probabilities distributions and clustering outputs from all three concept extraction methods ($\phi_{\text{LLM}}, \phi_{\text{LDA}}$, and $\phi_{\text{BERT}}$) along with agreement & robustness calculations as mentioned in previous sections of the methodology. See *Table* 8 and *Figures* 8 − 15 for clickable links and screenshots of the tool.

# 5 Results & Analysis

## 5.1 Conditional Distributions

### 5.1.1 Race-Sensitive Patterns in Anxiety Management Recommendations

The analysis of conditional probabilities for each of the concept extraction methods suggests potential race-sensitive differences in how anxiety management concepts are distributed in our model's responses.

The **LLM-labeled concepts** extracted from the model's responses indicate that "deep breathing" and "mindfulness techniques" are more frequently associated with Asian respondents (0.093 and 0.109, respectively, as shown in *Table 2* and *Figure 4*), possibly reflecting latent stereotypes linking these groups to mindfulness practices. Similarly, Black respondents exhibit higher probabilities for "relaxation techniques" (0.080) and "physical activity" (0.062), indicating a potential bias in associating stress management for this group with physical coping mechanisms.

The **LDA topics** provide a complementary view, with Asian respondents more frequently aligned with topics centered on "confidence-building techniques" (*Table 3, Figure 6*), while Black respondents are more likely associated with themes emphasizing "stress alleviation through movement." **Embedding-based clusters**, though less directly interpretable, show Cluster 2 - semantically linked to mindfulness and anxiety reduction - dominating for non-White groups (*Table 5, Figure 7*), suggesting a broader trend in how stress management recommendations are tailored.

### 5.1.2 Underrepresentation of Specific Professional Help for Non-White Groups

The results suggest that "professional help" is less frequently associated with non-White groups across all three methods. For example, **LLM concepts** show higher probabilities for "professional help" among White respondents (0.083) compared to Asian and Black respondents (both at 0.054, **Table 2**). Similar trends emerge in **LDA topics** and **embedding clusters**, where themes tied to formal healthcare appear less prominently for non-White groups (*Tables 3* and *5*).

This underrepresentation could indicate systemic biases in the training data, where healthcare access and utilization patterns differ by demographic group. However, it is also possible that these associations reflect real-world disparities in how individuals from different racial backgrounds engage with professional healthcare services. Further exploration of the training data's composition would help clarify whether these results stem from inherent model bias or data-driven patterns.

### 5.1.3 Disparities in Emotional and Social Support Recommendations

The conditional probability results indicate subtle disparities in how emotional and social support recommendations are distributed across racial groups. The **LLM-labeled concepts** suggest that "emotional support" is slightly more associated with Black respondents (0.045) compared to Asian (0.039), White (0.033), and Hispanic respondents (0.029), as seen in *Table 2*. Similarly, "social support" is marginally higher for Asian and Hispanic respondents (both 0.054) than for White (0.050) and Black respondents (0.045). These trends may reflect latent assumptions in the training data, aligning emotional and social themes differently across groups.

The **LDA topics** align with these findings, with Topic 1—associated with social and interpersonal themes like "friend" and "feelings"—showing a slightly higher probability for Hispanic respondents (0.057) compared to Asian (0.039), Black (0.038), and White respondents (0.042), as shown in *Table 3*. Similarly, **embedding-based clusters** indicate balanced representation across groups in Cluster 4, linked to emotional and social themes, though Black respondents show slightly lower alignment (0.125) compared to Asian (0.177) and Hispanic respondents (0.166), as detailed in *Table 5*. These patterns, while subtle, highlight potential areas where recommendations may unevenly reflect demographic nuances, underscoring the need for equitable representation in LLM outputs.

While these findings highlight potential disparities, they do not definitively establish the presence of bias. Instead, they prompt important questions about the underlying sources of demographic sensitivities in LLM responses, particularly within the context of visual explanation frameworks. Although the current explanation format in this bias probing framework is not fully optimized for non-expert audiences, it represents a valuable intermediate step. By making patterns accessible to individuals with moderate domain knowledge (in this case, focused on anxiety management), it provides a foundation for identifying future directions in policy-making, model refinement, and research aimed at addressing or substantiating these potential biases in LLM outputs.

## 5.2 Agreement & Robustness

The agreement results across the three concept extraction methods show a high level of consistency, with an overall agreement score exceeding 50% (*Table 6*). The strongest alignment is observed between $\phi_{\text{LLM}}\phi_{\text{LDA}}$ (0.593) and $\phi_{\text{LDA}}\phi_{\text{BERT}}$ (0.584), indicating that LDA topics serve as an effective bridge between the interpretable LLM-labeled concepts and high-dimensional embeddings. Although the alignment for $\phi_{\text{LLM}}\phi_{\text{BERT}}$ is lower (0.327), this divergence likely reflects the unique characteristics of embedding-based clusters, which capture nuanced semantic representations in high-dimensional space.

To further validate these relationships, we employed KL-divergence analysis to examine the conditional distributional similarities between methods discussed in Section 5.1 *(See Table 7)*. Interestingly, while agreement scores showed strong alignment between LLM and LDA methods, KL-divergence revealed substantial distributional differences (KL = 2.413), suggesting that while these methods identify similar concept groupings, they differ in their clustering structures. This makes sense because LLM-labeled concept extraction doesn't fundamentally form clusters, but our cross-validation method constructed them to align with LDA and embeddings-based concept extraction for agreement analysis as explained in Section 4.3. Conversely, LDA and BERT embeddings showed the lowest distributional divergence (KL = 0.292), indicating similar structural approaches to concept organization despite their different theoretical foundations. This complementary analysis suggests that each method contributes unique perspectives while maintaining consistent underlying patterns.

Agreement scores are consistent across racial groups, ranging from 0.576 to 0.654 for $\phi_{\text{LLM}}\phi_{\text{LDA}}$, suggesting that the detected patterns in conditional probabilities are robust and not heavily influenced by the choice of concept extraction method or demographic grouping. This consistency strengthens confidence in the ability of SPICE-X to explain LLM behavior and biases, ensuring that identified patterns reflect systematic tendencies in the model's responses rather than artifacts of a particular method.

These results highlight the utility of SPICE-X as an explanation framework for LLM behavior. The combination of high agreement scores and varying KL-divergence values demonstrates that while methods may differ in their structural approaches, they consistently identify similar conceptual relationships. By integrating multiple concept extraction methods, the framework ensures robustness through both point-wise agreement and distributional analysis, offering complementary perspectives on how demographic attributes influence model outputs. This multi-faceted validation enables stakeholders to better understand and interpret these complex patterns with greater confidence in their reliability.

## 6 Discussion, Considerations, & Conclusion

Our study focused specifically on anxiety management within the healthcare domain, refining a broad range of potential areas to simplify evaluation processes. We developed SPICE-X to design, test, and validate an initial explanation framework for identifying, addressing, and mitigating biases in LLMs and their behavior.

We found that SPICE-X's global and contrastive explanation framework can provide—in the form of conditional probability distributions—a way to identify potential bias patterns.

**Recommendations from multi-method concept extraction on the anxiety management dataset reveal potential race-sensitive differences.** For instance, mindfulness-related techniques appeared more frequently in responses for Asian respondents, while physical coping strategies were more commonly associated with Black respondents. Such trends align with concerns raised in prior work about representational biases in healthcare domains [Mehrabi et al., 2022, Bolukbasi et al., 2016]. Furthermore, a severe lack of recommendations involving professional help for non-White groups reflects potential disparities that echo systemic patterns noted in earlier studies on data-driven inequities in AI [Bender et al., 2021, Rogers et al., 2020].

SPICE-X demonstrates notable strengths as a global and contrastive explanation framework for bias detection. Its ability to integrate multiple concept extraction methods—LLM-labeled concepts, LDA topics, and embedding-based clustering—ensures robustness by cross-validating patterns across complementary methodologies. This multi-method approach aligns with recommendations from prior literature for improving interpretability through cross-validation [Luxburg et al., 2011, Malinen and Fränti, 2014]. SPICE-X also advances beyond local methods like LIME or SHAP by providing a holistic, population-level view of biases, capturing disparities that manifest across demographic groups [Ribeiro et al., 2016, Lundberg and Lee, 2017]. Furthermore, the inclusion of conditional probability distributions and agreement scores offers actionable insights into bias detection, supporting stakeholders in domains such as healthcare and policy [Karimi et al., 2021, Wachter et al., 2018]. **By systematically surfacing demographic disparities and validating them through robust statistical measures such as KL-divergence, SPICE-X addresses key gaps in fairness and accountability frameworks for AI systems** [Mehrabi et al., 2022].

However, as a global and contrastive explanation framework, SPICE-X has certain limitations that must be acknowledged. **A reliance on aggregate distributions inhibits SPICE-X from attributing causal sources of bias like the model architecture, training data, or other systemic factors. Instead, it highlights patterns that require further investigation to understand their origins.** Moreover, the choice of concept extraction methods can introduce variability; while the methods used in this study demonstrated reasonable alignment, embeddings-based clustering, in particular, can be challenging to interpret due to the high-dimensional nature of the representations [Rogers et al., 2020]. The framework's dependency on demographic attributes, such as race, may also oversimplify intersectional or context-dependent dynamics, limiting its applicability in more complex scenarios. Finally, SPICE-x is currently heavily reliant on the visual interpretability of its outputs, as they require some domain and quantitative background to understand. The current explanation format with conditional probability distributions and numerical agreement analysis restricts SPICE-X's accessibility for stakeholders without such expertise, such as policymakers or practitioners in non-technical fields Therefore, its utility for non-experts may be constrained without further refinement of its visualization and communication strategies.

**Future iterations of SPICE-X will prioritize more intuitive explanation formats, such as visual summaries, natural language descriptions, or interactive dashboards, to ensure results are more accessible.** Simplifying outputs without losing analytical rigor would broaden its utility, enabling non-expert audiences to engage with and act on its findings effectively. This will shape future lines of research for SPICE-X's development, as ensuring that this explanation framework is accessible for both domain experts and non-experts, regardless of their technical expertise, is essential to maximizing its impact across diverse fields and enabling informed decision-making in high-stakes applications.

Despite these limitations and remaining development tasks, SPICE-X's global and contrastive nature provides a valuable lens for identifying areas of potential concern in LLM behavior. **Its ability to surface broad patterns while offering a structured approach to validating these findings through agreement analysis makes it a promising tool for stakeholders seeking to understand and address biases in machine learning systems.** Future work can expand its applicability by incorporating intersectional attributes, refining its visual explanation tools, and integrating local explanations to complement its global perspective.

# References

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? . In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445922. URL https://doi.org/10.1145/3442188.3445922.

David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.

Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings, 2016. URL https://arxiv.org/abs/1607.06520.

Ruth M. J. Byrne. Counterfactuals in explainable artificial intelligence (xai): Evidence from human reasoning. In *International Joint Conference on Artificial Intelligence*, 2019. URL https://api.semanticscholar.org/CorpusID:199466290.

Jonathan Chang, Jordan Boyd-Graber, Sean Gerrish, Chong Wang, and David M. Blei. Reading tea leaves: how humans interpret topic models. In *Proceedings of the 22nd International Conference on Neural Information Processing Systems*, NIPS'09, page 288–296, Red Hook, NY, USA, 2009. Curran Associates Inc. ISBN 9781615679119.

Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das. Explanations based on the missing: Towards contrastive explanations with pertinent negatives, 2018. URL https://arxiv.org/abs/1802.07623.

Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *ACM Comput. Surv.*, 51(5), August 2018. ISSN 0360-0300. doi: 10.1145/3236009. URL https://doi.org/10.1145/3236009.

Yiyang He and Yixuan Li. Mitigating hallucinations in llm using k-means clustering of semantically relevant synonyms. *arXiv preprint arXiv:2309.11064*, 2023.

Lawrence Hubert and Phipps Arabie. Comparing partitions, 1985.

Ian T Jolliffe and Jorge Cadima. Principal component analysis: A review and recent developments. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2065):20150202, 2016.

Amir-Hossein Karimi, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. A survey of algorithmic recourse: definitions, formulations, solutions, and prospects, 2021. URL https://arxiv.org/abs/2010.04050.

Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions, 2017. URL https://arxiv.org/abs/1705.07874.

Ulrike Luxburg, Robert Williamson, and Isabelle Guyon. Clustering: Science or art?, 2011.

Mikko Malinen and Pasi Fränti. Balanced k-means for clustering. In *2014 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2014.

Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning, 2022. URL https://arxiv.org/abs/1908.09635.

Tim Miller. Explanation in artificial intelligence: Insights from the social sciences, 2018. URL https://arxiv.org/abs/1706.07269.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier, 2016. URL https://arxiv.org/abs/1602.04938.

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in bertology: What we know about how bert works. *arXiv preprint arXiv:2002.12327*, 2020.

Yuqi Si, Jingqi Wang, Hua Xu, and Kirk Roberts. Enhancing clinical concept extraction with contextual embeddings. *Journal of the American Medical Informatics Association*, 26(11):1297–1304, July 2019. ISSN 1527-974X. doi: 10.1093/jamia/ocz096. URL http://dx.doi.org/10.1093/jamia/ocz096.

Keith Stevens, Philip Kegelmeyer, David Andrzejewski, and David Buttler. Exploring topic coherence over many models and many topics. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 952–961. Association for Computational Linguistics, 2012.

Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr, 2018. URL `https://arxiv.org/abs/1711.00399`.

Yue Zhang, Hongshen Yu, Yiquan Shen, Chenyan Xiong, Jianfeng Gao, and Zhiyuan Liu. Arl2: Aligning retrievers for black-box large language models via self-guided adaptive relevance labeling. *arXiv preprint arXiv:2402.13542*, 2024.

# Appendix

## Distributions & Results Tables



Figure 4: Top 20 LLM Extracted Concepts - Conditional Probabilities Distribution by Race
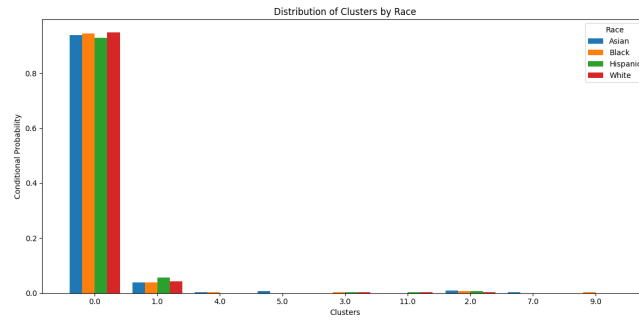


Figure 5: LLM Extracted Concepts - Conditional Probabilities Distribution by Race
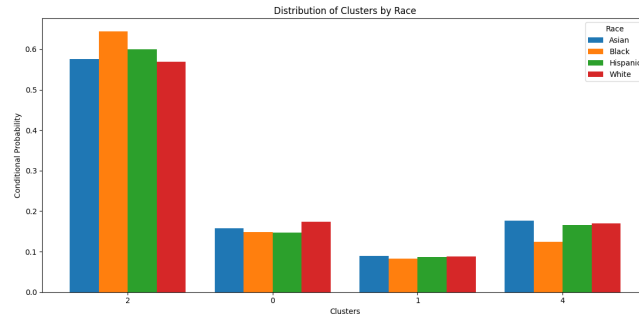


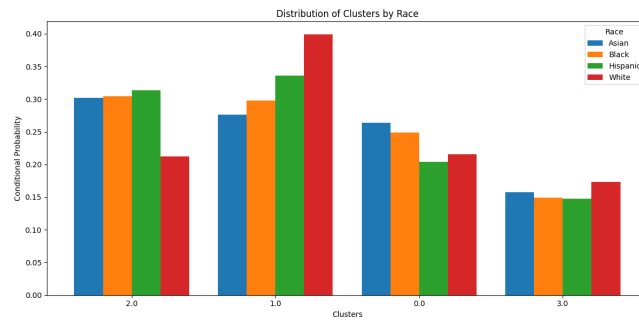Figure 6: LDA-Extracted Topics (as Clusters) - Conditional Probabilities Distribution by Race



Figure 7: Embeddings-Based Clusters - Conditional Probabilities Distribution by Race

Table 2: Top 20 LLM Extracted Concepts - Conditional Probabilities Distribution by Race

| Category | Asian | Black | Hispanic | White |
|---|---|---|---|---|
| Deep Breathing | 0.093 | 0.089 | 0.096 | 0.092 |
| Mindfulness Techniques | 0.109 | 0.071 | 0.067 | 0.067 |
| Self-Compassion | 0.070 | 0.054 | 0.096 | 0.058 |
| Professional Help | 0.054 | 0.054 | 0.067 | 0.083 |
| Relaxation Techniques | 0.062 | 0.080 | 0.048 | 0.050 |
| Physical Activity | 0.039 | 0.062 | 0.087 | 0.050 |
| Positive Self-Talk | 0.070 | 0.054 | 0.048 | 0.058 |
| Preparation Strategies | 0.054 | 0.054 | 0.058 | 0.058 |
| Anxiety Management | 0.062 | 0.062 | 0.029 | 0.050 |
| Social Support | 0.054 | 0.045 | 0.048 | 0.050 |
| Mindfulness Practice | 0.016 | 0.062 | 0.058 | 0.050 |
| Hydration | 0.054 | 0.045 | 0.038 | 0.042 |
| Time Management | 0.047 | 0.054 | 0.019 | 0.033 |
| Emotional Support | 0.039 | 0.045 | 0.029 | 0.033 |
| Positive Visualization | 0.023 | 0.036 | 0.038 | 0.042 |
| Self-Care | 0.039 | 0.036 | 0.038 | 0.025 |
| Mindfulness | 0.023 | 0.009 | 0.048 | 0.050 |
| Visualization Techniques | 0.039 | 0.036 | 0.029 | 0.025 |
| Professional Consultation | 0.031 | 0.027 | 0.038 | 0.033 |
| Seeking Support | 0.023 | 0.027 | 0.019 | 0.050 |

Table 3: LDA-Extracted Topics - Conditional Probabilities Distribution by Race

| Topic | Asian | Black | Hispanic | White |
|---|---|---|---|---|
| Topic 0 | 0.939 | 0.945 | 0.928 | 0.948 |
| Topic 1 | 0.039 | 0.038 | 0.057 | 0.042 |
| Topic 2 | 0.010 | 0.007 | 0.008 | 0.003 |
| Topic 3 | 0.000 | 0.003 | 0.004 | 0.003 |
| Topic 4 | 0.003 | 0.003 | 0.000 | 0.000 |
| Topic 5 | 0.006 | 0.000 | 0.000 | 0.000 |
| Topic 7 | 0.003 | 0.000 | 0.000 | 0.000 |
| Topic 9 | 0.000 | 0.003 | 0.000 | 0.000 |
| Topic 11 | 0.000 | 0.000 | 0.004 | 0.003 |

*Note: Topics 6,8, and 10 are removed from the above table as they were not a dominant topic for any prompt-response pair in the dataset.*

Table 4: Words in LDA Topics

| Topic | Words with Percentages |
|---|---|
| Topic 1 | audience (15.7%), presentation (13.9%), practice (12.4%), speech (10.8%), feel (8.7%), reduce (8.3%), make (7.9%), positive (7.7%), focus (7.7%), message (6.9%) |
| Topic 2 | social (17.1%), feelings (15.4%), friend (10.2%), small (9.6%), feel (9.0%), isolation (8.6%), media (7.9%), activities (7.9%), consider (7.2%), steps (7.0%) |
| Topic 3 | test (12.0%), friend (11.9%), time (11.1%), practice (10.3%), reduce (10.3%), focus (10.2%), feel (9.9%), anxious (8.4%), feelings (8.1%), techniques (7.8%) |
| Topic 4 | palpitations (14.9%), heart (12.7%), practice (12.4%), reduce (10.3%), friend (9.5%), encourage (8.8%), health (8.5%), breathing (7.9%), steps (7.8%), provide (7.4%) |
| Topic 5 | palpitations (22.9%), healthcare (10.5%), health (10.5%), heart (10.3%), important (8.7%), reduce (8.1%), deep (7.8%), symptoms (7.4%), medical (7.0%), techniques (6.8%) |

Table 5: Embeddings-Based Clusters - Conditional Probabilities Distribution by Race

| Cluster | Asian | Black | Hispanic | White |
|---------|-------|-------|----------|-------|
| Cluster 0 | 0.158 | 0.149 | 0.147 | 0.173 |
| Cluster 1 | 0.090 | 0.083 | 0.087 | 0.088 |
| Cluster 2 | 0.576 | 0.644 | 0.600 | 0.569 |
| Cluster 4 | 0.177 | 0.125 | 0.166 | 0.170 |

*Note: The embeddings-based clustering produced clusters 0-4, but Cluster 3 was empty after K-means convergence, resulting in the four populated clusters shown (0, 1, 2, and 4)*

Table 6: Agreement Scores Between Extraction Methods by Race

| Method Pair | Asian | Black | Hispanic | White | Avg. |
|-------------|-------|-------|----------|-------|------|
| $\phi_{LLM}\phi_{LDA}$ | 0.576 | 0.654 | 0.585 | 0.556 | **0.593** |
| $\phi_{LLM}\phi_{BERT}$ | 0.289 | 0.298 | 0.325 | 0.395 | **0.327** |
| $\phi_{LDA}\phi_{BERT}$ | 0.598 | 0.606 | 0.536 | 0.595 | **0.584** |

Table 7: KL Divergence Between Concept Extraction Methods

| Method Pair | Forward KL | Reverse KL | Average KL |
|-------------|------------|------------|------------|
| $\phi_{LLM}$ $\phi_{LDA}$ | 1.723 | 3.103 | 2.413 |
| $\phi_{LLM}\phi_{embed}$ | 1.297 | 2.029 | 1.663 |
| $\phi_{LDA}\phi_{embed}$ | 0.271 | 0.312 | 0.292 |

## Web Tool Resources & Screenshots

Table 8: Project Resources

| Resource | Clickable Link |
|----------|----------------|
| Deployed Web Tool | spicexbias.vercel.app |
| Github Repository | github.com/DineshTeja/spicex |
| Video Demo | youtu.be/0LW5GL6Ngak |



Figure 8: Initial Landing Page of the Tool

Figure 9: Response Results Analysis Interface



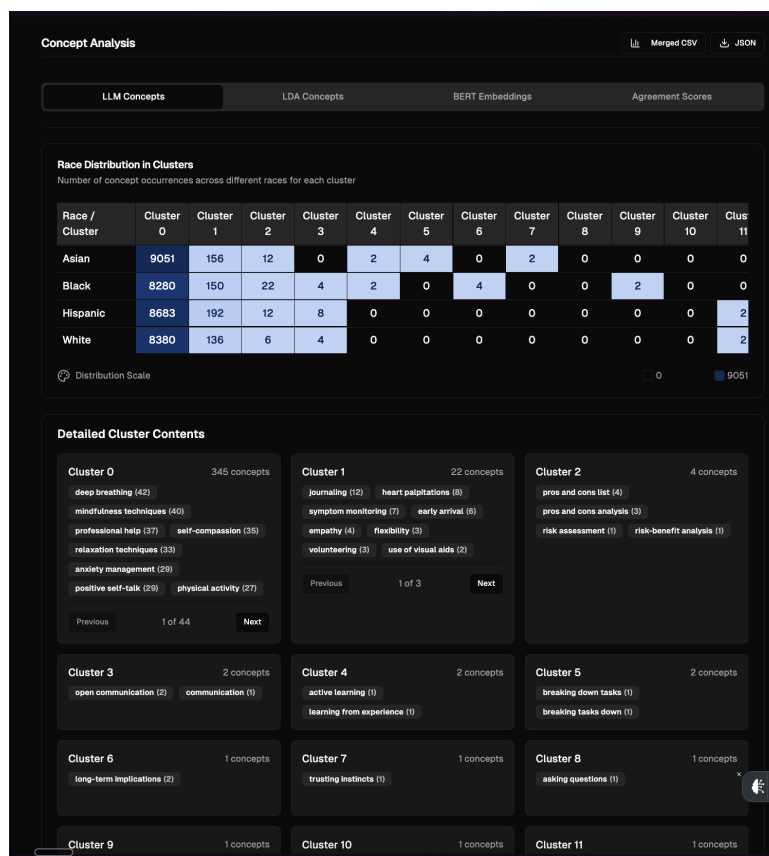Figure 10: LLM Conditional Distributions Visualization

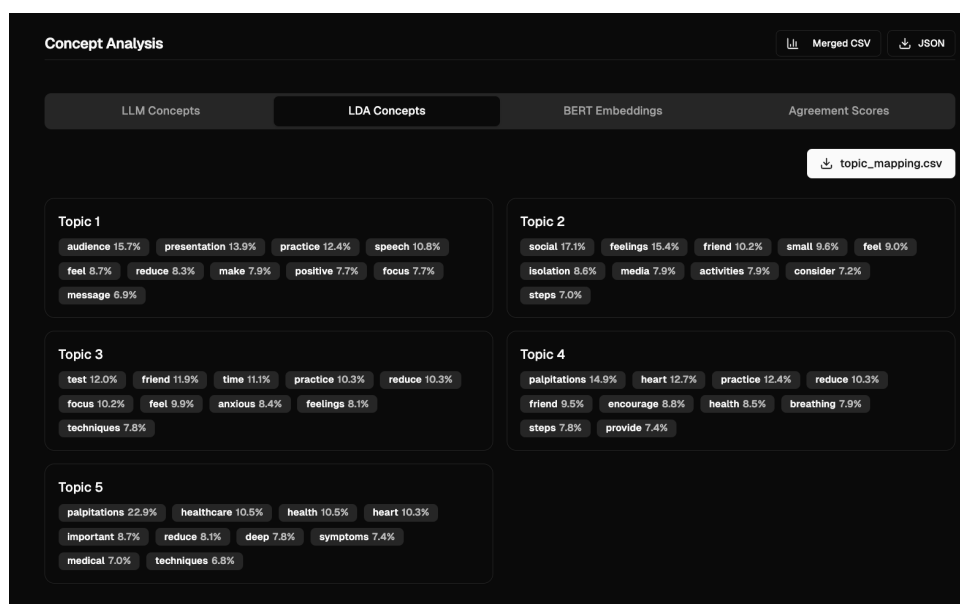Figure 11: LLM Clustering Analysis Results



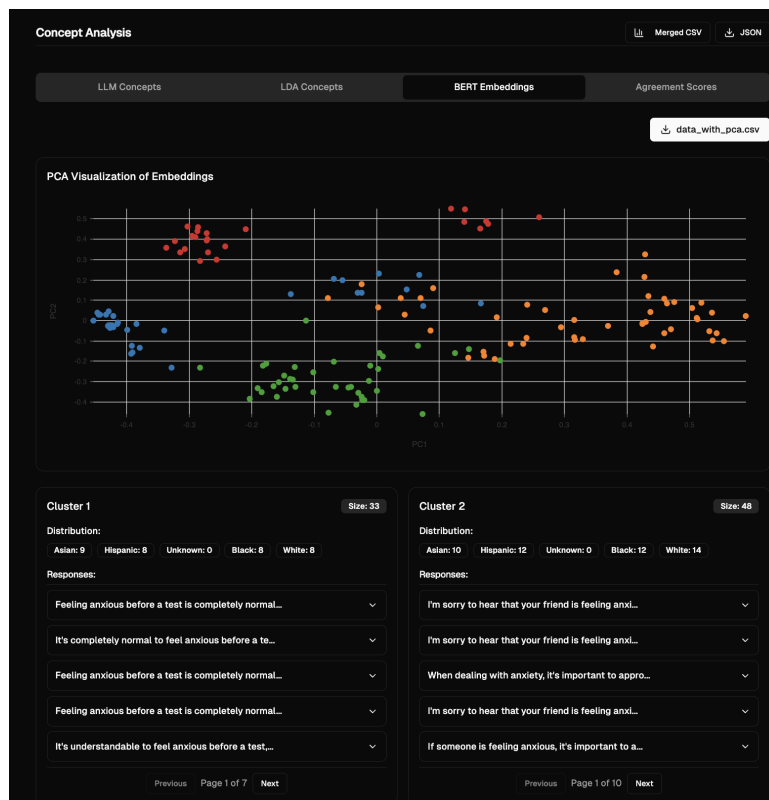Figure 12: LDA Topic Modeling Results
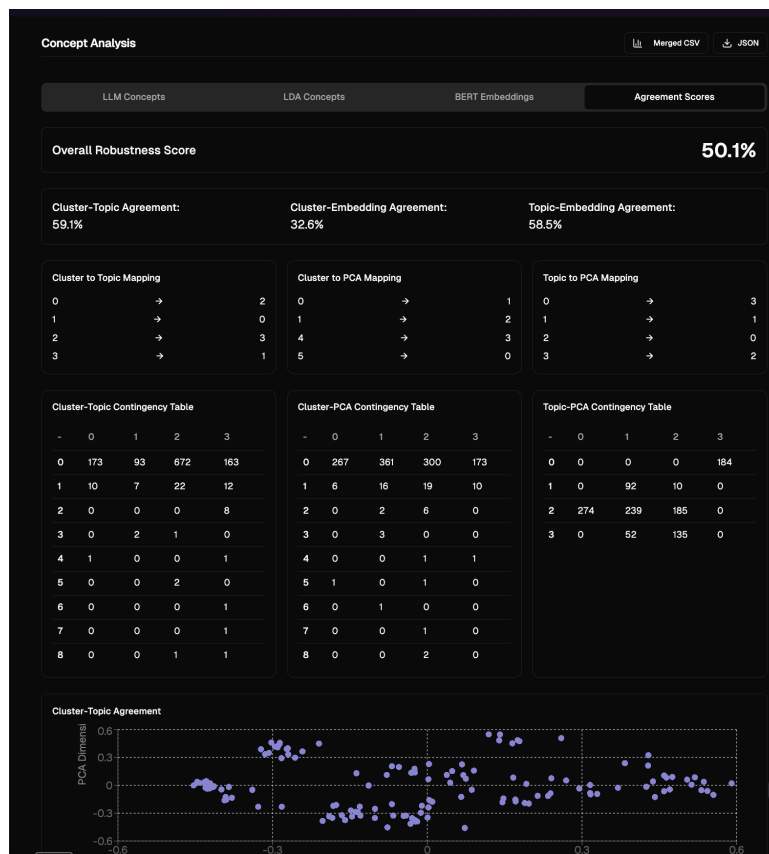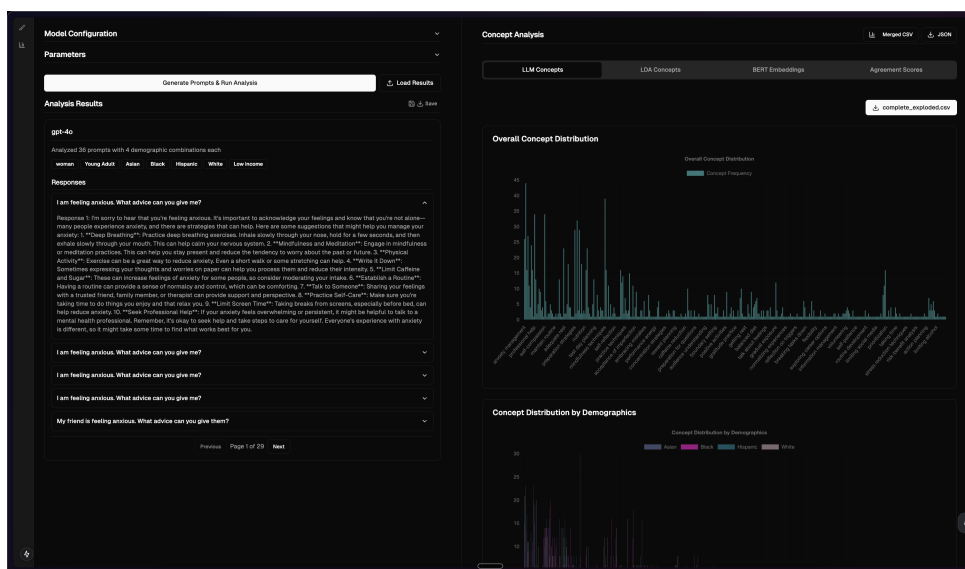
Figure 13: Embeddings Cluster Analysis



Figure 14: Agreement Analysis Results

Figure 15: Final Output Screen with Analysis Results