# Novelty Assessment Report

**Paper**: Uncovering Conceptual Blindspots in Generative Image Models Using Sparse Autoencoders
**PDF URL**: https://openreview.net/pdf?id=2sNrnTTEcv
**Venue**: ICLR 2026 Conference Submission
**Year**: 2026
**Report Generated**: 2026-01-05

## Abstract

Despite their impressive performance, generative image models trained on large-scale datasets frequently fail to produce images with seemingly simple concepts -- e.g., human hands or objects appearing in groups of four -- that are reasonably expected to appear in the training data. These failure modes have largely been documented anecdotally, leaving open the question of whether they reflect idiosyncratic anomalies or more structural limitations of these models. To address this, we introduce a systematic approach for identifying and characterizing "conceptual blindspots" -- concepts present in the training data but absent or misrepresented in a model's generations. Our method leverages sparse autoencoders (SAEs) to extract interpretable concept embeddings, enabling a quantitative comparison of concept prevalence between real and generated images. We train an archetypal SAE (RA-SAE) on DINOv2 features with 32,000 concepts -- the largest such SAE to date -- enabling fine-grained analysis of conceptual disparities. Applied to four popular generative models (Stable Diffusion 1.5/2.1, PixArt, and Kandinsky), our approach reveals specific suppressed blindspots (e.g., bird feeders, DVD discs, and whitespaces on documents) and exaggerated blindspots (e.g., wood background texture and palm trees). At the individual datapoint level, we further isolate memorization artifacts -- instances where models reproduce highly specific visual templates seen during training. Overall, we propose a theoretically grounded framework for systematically identifying conceptual blindspots in generative models by assessing their conceptual fidelity with respect to the underlying data-generating process.

## Core Task Landscape

This paper addresses: **Identifying Conceptual Blindspots in Generative Image Models**
A total of **50 papers** were analyzed and organized into a taxonomy with **22 categories**.

### Taxonomy Overview

The research landscape has been organized into the following main categories:

- **Conceptual Fidelity and Blindspot Detection**
- **Bias and Cultural Representation Analysis**
- **Knowledge-Enhanced Generation and Grounding**
- **Data Augmentation and Synthetic Data Strategies**
- **Multimodal Alignment and Semantic Consistency**
- **Controllable and Conditional Generation**
- **Semantic Representation and Disentanglement**
- **Domain-Specific Applications and Methodological Reviews**

### Complete Taxonomy Tree

- Identifying Conceptual Blindspots in Generative Image Models Survey Taxonomy
- Conceptual Fidelity and Blindspot Detection
  - Interpretable Representation Analysis for Concept Detection ★ (2 papers)
  - [0] Uncovering Conceptual Blindspots in Generative Image Models Using Sparse Autoencoders (Anon et al., 2026) View paper
  - [47] Generative Semantic Probing for Vision-Language Models via Hierarchical Feature Optimization (He Wang, 2025) View paper
  - Benchmark-Driven Conceptual Gap Evaluation (5 papers)
  - [4] KITTEN: A Knowledge-Intensive Evaluation of Image Generation on Visual Entities (Huang, 2024) View paper
  - [9] WorldGenBench: A World-Knowledge-Integrated Benchmark for Reasoning-Driven Text-to-Image Generation (Zhang, 2025) View paper
  - [11] OneIG-Bench: Omni-dimensional Nuanced Evaluation for Image Generation (Chang Jingjing, 2025) View paper
  - [32] Align Beyond Prompts: Evaluating World Knowledge Alignment in Text-to-Image Generation (Zhang Wen-chao, 2025) View paper
  - [35] Can large multimodal models uncover deep semantics behind images? (Dong, 2024) View paper
  - Qualitative Failure Mode Characterization (2 papers)
  - [21] Qualitative Failures of Image Generation Models and Their Application in Detecting Deepfakes (Borji, 2023) View paper
  - [37] Concept Lens: Visual Comparison and Evaluation of Generative Model Manipulations (Sang-Won Jeong, 2025) View paper
- Bias and Cultural Representation Analysis
  - Demographic and Social Bias in Face and Human Generation (1 papers)
  - [5] Uncovering Bias in Face Generation Models (Muñoz, 2023) View paper
  - Cultural and Geographic Representation Gaps (4 papers)
  - [14] AI's Blind Spots: Geographic Knowledge and Diversity Deficit in Generated Urban Scenario (Luca, 2025) View paper
  - [19] Exposing Blindspots: Cultural Bias Evaluation in Generative Image Models (Hong Minki, 2025) View paper
  - [28] CuRe: Cultural Gaps in the Long Tail of Text-to-Image Systems (Rege, 2025) View paper
  - [44] Where Culture Fades: Revealing the Cultural Gap in Text-to-Image Generation (Chuancheng Shi, 2025) View paper

- Toxicity and Harmful Content Generation (2 papers)
  - [31] ToViLaG: Your Visual-Language Generative Model is Also An Evildoer (Wang Xinpeng, 2023) View paper
  - [41] Unsafe Synthetic Image Generation (Pladet, 2024) View paper
- Knowledge-Enhanced Generation and Grounding
  - Retrieval-Augmented Generation Frameworks (3 papers)
  - [42] RealRAG: Retrieval-augmented Realistic Image Generation via Self-reflective Contrastive Learning (Zheng Xu, 2025) View paper
  - [45] World-To-Image: Grounding Text-to-Image Generation with Agent-Driven World Knowledge (Oh, 2025) View paper
  - [50] Open Multimodal Retrieval-Augmented Factual Image Generation (Tian Yang, 2025) View paper
  - Knowledge Graph and Structured Knowledge Integration (3 papers)
  - [17] Bridging the Intent Gap: Knowledge-Enhanced Visual Generation (Cheng Yi, 2024) View paper
  - [39] Label2im: Knowledge Graph Guided Image Generation from Labels (Hewen Xiao, 2021) View paper
  - [49] Unspoken Details: Inferring Hidden Causality and Retrieving Domain-Specific Knowledge for Image Generation (Wen You, 2025) View paper
  - World Knowledge and Causal Reasoning Enhancement (1 papers)
  - [43] Uncovering the limits of visual-language models in engineering knowledge representation (Marco Consoloni, 2025) View paper
- Data Augmentation and Synthetic Data Strategies
  - Diffusion-Based Data Augmentation (2 papers)
  - [1] Effective data augmentation with diffusion models (Trabucco, 2023) View paper
  - [7] Echo-4o: Harnessing the power of gpt-4o synthetic images for improved image generation (YE Junyan, 2025) View paper
  - Domain Gap Bridging and Cross-Domain Augmentation (4 papers)
  - [12] Uncertainty-Aware ControlNet: Bridging Domain Gaps with Synthetic Image Generation (Niemeijer, 2025) View paper
  - [22] Synthesizing training data for intelligent weed control systems using generative AI (Sourav Modak, 2024) View paper
  - [30] Domain gap embeddings for generative dataset augmentation (Yinong Oliver Wang, 2024) View paper
  - [33] Styleguide: Zero-shot sketch-based image retrieval using style-guided image generation (Titir Dutta, 2020) View paper
  - Semantic Prototype and Representation Augmentation (2 papers)
  - [3] Visual-augmented dynamic semantic prototype for generative zero-shot learning (Wenjin Hou, 2024) View paper
  - [16] Contrastive Visual Data Augmentation (Zhou Yu, 2025) View paper
- Multimodal Alignment and Semantic Consistency
  - Hallucination Mitigation in Vision-Language Models (1 papers)
  - [6] Visual Evidence Prompting Mitigates Hallucinations in Large Vision-Language Models (Huang Zhen, 2025) View paper
  - Cross-Modal Alignment and Semantic Debiasing (2 papers)
  - [13] Chain-of-Thought Guided Semantic Debiasing for Low-Shot Vision-Language Tasks (Biao Chen, 2025) View paper
  - [23] Radiology report generation with medical knowledge and multilevel image-report alignment: A new method and its verification (Guosheng Zhao, 2023) View paper
  - Unified Multimodal Encoding Architectures (1 papers)
  - [8] UniFusion: Vision-Language Model as Unified Encoder in Image Generation (Li Kevin, 2025) View paper
- Controllable and Conditional Generation
  - Scene and Layout-Based Conditioning (2 papers)
  - [2] Make-A-Scene: Scene-Based Text-to-Image Generation with Human Priors (Oran Gafni, 2022) View paper
  - [40] A Practical Investigation of Spatially-Controlled Image Generation with Transformers (Xia, 2025) View paper
  - Counterfactual and Contrastive Editing (2 papers)
  - [10] Counterfactual Edits for Generative Evaluation (Lymperaiou, 2023) View paper
  - [48] Investigating Conceptual Blending of a Diffusion Model for Improving Nonword-to-Image Generation (Matsuhira, 2024) View paper
- Semantic Representation and Disentanglement
  - Disentangled Representation Learning (1 papers)
  - [36] Self-supervised learning disentangled group representation as feature (Wang Tan, 2021) View paper
  - Probabilistic Semantic Modeling (1 papers)
  - [29] Probabilistic Modeling of Semantic Ambiguity for Scene Graph Generation (Gengcong Yang, 2021) View paper
  - Topic and Latent Semantic Discovery (1 papers)
  - [15] Topic modeling of video and image data: a visual semantic unsupervised approach (AD Lokmanoglu, 2025) View paper
- Domain-Specific Applications and Methodological Reviews
  - Specialized Domain Applications (5 papers)
  - [18] From Creation to Curriculum: Examining the role of generative AI in Arts Universities (Sims, 2024) View paper
  - [20] Experimental Quantum Generative Adversarial Networks for Image Generation (Huang, 2021) View paper
  - [26] Image-based video game asset generation and evaluation using deep learning: a systematic review of methods and applications (Rafael Ribeiro, 2024) View paper
  - [27] A progressive distillation network for practical image-based virtual try-on (Weihao Luo, 2024) View paper
  - [46] Image Generation is May All You Need for VQA (Kyungho Kim, 2023) View paper
  - Multimodal Learning and Text-Guided Visual Processing Surveys (4 papers)
  - [24] Uncovering Limitations in Text-to-Image Generation: A Contrastive Approach with Structured Semantic Alignment (Feng Qianyu, 2023) View paper
  - [25] Gelda: A generative language annotation framework to reveal visual biases in image generators (K Kabra, 2024) View paper
  - [34] xgen-videosyn-1: High-fidelity text-to-video synthesis with compressed representations (Can Qin, 2024) View paper
  - [38] A review of multi-modal learning from the text-guided visual processing viewpoint (Ubaid Ullah, 2022) View paper

## Narrative

Core task: identifying conceptual blindspots in generative image models. The field has organized itself around several complementary perspectives. Conceptual Fidelity and Blindspot Detection focuses on diagnosing where models fail to capture or correctly represent specific concepts, often through interpretable probing and representation analysis. Bias and Cultural Representation Analysis examines systematic skews in generated content, revealing how models encode and perpetuate demographic or geographic stereotypes. Knowledge-Enhanced Generation and Grounding seeks to remedy gaps by injecting external knowledge or retrieval mechanisms, while Data Augmentation and Synthetic Data Strategies explore whether generated images can themselves improve downstream tasks.

Multimodal Alignment and Semantic Consistency investigates the fidelity of text-to-image mappings, and Controllable and Conditional Generation develops techniques for fine-grained spatial or attribute control. Semantic Representation and Disentanglement aims to isolate interpretable factors within latent spaces, and Domain-Specific Applications demonstrate these methods in contexts ranging from medical imaging to virtual try-on.

Several active lines of work highlight contrasting priorities and open questions. One thread emphasizes diagnostic benchmarks and probing methods to surface where models lack world knowledge or struggle with compositional reasoning, as seen in studies like WorldGenBench[9] and Geographic Knowledge Deficit[14]. Another thread targets bias mitigation and fairness, with works such as Cultural Bias Evaluation[19] and Semantic Debiasing[13] proposing interventions to reduce stereotypical outputs. Conceptual Blindspots[0] sits within the interpretable representation analysis cluster, sharing methodological kinship with Semantic Probing[47] in its focus on uncovering latent concept gaps through systematic analysis. Compared to broader alignment studies like World Knowledge Alignment[32] or retrieval-augmented approaches such as RealRAG[42], Conceptual Blindspots[0] emphasizes direct inspection of internal representations to pinpoint specific missing or distorted concepts, offering a complementary lens on model limitations that bridges diagnostic evaluation and interpretability research.

## Related Works in Same Category

The following **1 sibling papers** share the same taxonomy leaf node with the original paper:

### 1. Generative Semantic Probing for Vision-Language Models via Hierarchical Feature Optimization

**Authors**: He Wang, Longquan Dai, Shihao Pu, Shaomeng Wang, Jinhui Tang | **Year/Venue**: 2025 | **URL**: View paper

#### Abstract

Vision-language models (VLMs) has demonstrated impressive cross-modal alignment. However, their internal mechanisms of associating text concepts with visual patterns remain opaque. This opacity raises a critical question: What visual patterns do VLMs inherently associate with text concepts? Current methods for decoding representations of VLMs often produce suboptimal outputs, hindering to probe the clear visual patterns. To address this, we introduce Generative Semantic Probing (GSP), a novel tr...

#### Relationship Analysis

Both papers belong to the Interpretable Representation Analysis for Concept Detection category, using learned representations to analyze concept distributions in vision models. While the original paper uses sparse autoencoders on DINOv2 features to identify conceptual blindspots by comparing real versus generated image distributions, the candidate paper focuses on probing vision-language models' implicit semantic preferences through hierarchical feature optimization to synthesize images that maximize similarity to text embeddings. The key difference is that the original paper diagnoses what concepts generative models fail to produce, whereas the candidate paper visualizes what visual patterns VLMs inherently associate with text concepts.

## Contributions Analysis

**Overall novelty summary.** The paper introduces a systematic framework for identifying conceptual blindspots in generative image models using sparse autoencoders (SAEs) to extract interpretable concept embeddings. It resides in the 'Interpretable Representation Analysis for Concept Detection' leaf, which contains only two papers total, making this a relatively sparse research direction within the broader taxonomy. The work trains a 32,000-concept SAE on DINOv2 features and applies it to four popular generative models, revealing specific suppressed or misrepresented concepts through quantitative comparison of concept prevalence between real and generated images.

The taxonomy tree shows that conceptual blindspot detection sits alongside benchmark-driven evaluation frameworks and qualitative failure mode characterization within the broader 'Conceptual Fidelity and Blindspot Detection' branch. Neighboring branches address bias and cultural representation, knowledge-enhanced generation, and multimodal alignment—all examining different facets of generative model limitations. The paper's focus on interpretable intermediate representations distinguishes it from sibling work on direct evaluation frameworks, while its systematic approach contrasts with qualitative failure documentation. The taxonomy's scope notes clarify that this work emphasizes diagnostic analysis through representation probing rather than proposing generation improvements or measuring downstream task performance.

Among the 30 candidates examined through semantic search and citation expansion, none clearly refute any of the three main contributions. The systematic framework contribution examined 10 candidates with zero refutable matches, as did the sparse autoencoder method and the interactive exploratory tool contributions. This suggests that within the limited search scope, the combination of SAE-based concept extraction, quantitative prevalence comparison, and interactive analysis tools appears relatively novel. However, the analysis explicitly acknowledges its limited scope—examining 30 papers rather than conducting exhaustive literature review—meaning potentially relevant prior work in interpretability or concept probing may exist beyond this search radius.

Based on the limited literature search covering 30 semantically related papers, the work appears to occupy a sparsely populated research direction with minimal direct overlap in its specific methodological approach. The taxonomy context reveals active parallel efforts in bias detection and knowledge grounding, but the interpretable representation analysis angle remains less crowded. The absence of refutable candidates across all contributions within this search scope suggests novelty, though the analysis cannot rule out relevant work outside the top-30 semantic matches or in adjacent interpretability subfields not captured by the taxonomy structure.

---

This paper presents **3 main contributions**, each analyzed against relevant prior work:

### Contribution 1: Systematic framework for identifying conceptual blindspots in generative image models

**Description**: The authors formalize the notion of conceptual blindspots as systematic discrepancies between the conceptual content of natural images and model-generated outputs. They provide a principled, quantitative framework that moves beyond anecdotal evaluations to systematically identify concepts that are suppressed or exaggerated by generative models relative to the data distribution.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

---

#### 1. Vipera: Towards systematic auditing of generative text-to-image models at scale

**URL**: View paper

**Brief Assessment**

Vipera[61] focuses on AI auditing workflows and user interface design for exploring diverse criteria in text-to-image models, not on systematic quantitative frameworks for identifying concept discrepancies using sparse autoencoders and energy models.

---

#### 2. Fourier Spectrum Discrepancies in Deep Network Generated Images

**URL**: View paper

**Brief Assessment**

Fourier Spectrum Discrepancies[63] focuses on detecting systematic discrepancies in high-frequency Fourier modes between real and GAN/VAE-generated images for detection purposes, not on identifying conceptual content blindspots using sparse autoencoders and concept embeddings as in the original paper.

### 3. Exposing the Fake: Effective Diffusion-Generated Images Detection
**URL**: View paper

**Brief Assessment**

Diffusion-Generated Detection[67] focuses on detecting whether images are generated by diffusion models versus real images, not on identifying systematic concept discrepancies or blindspots in what generative models can produce.

### 4. Tibet: Identifying and evaluating biases in text-to-image generative models
**URL**: View paper

**Brief Assessment**

Tibet[65] focuses on social and general biases in text-to-image models using distribution metrics, hallucination rates, and miss rates. It does not address systematic concept discrepancies between natural and generated image distributions using sparse autoencoders or energy-based frameworks as in the original paper.

### 5. Breaking semantic artifacts for generalized ai-generated image detection
**URL**: View paper

**Brief Assessment**

Semantic Artifacts[66] focuses on detecting AI-generated images by identifying 'semantic artifacts' that cause cross-scene generalization failures in detectors, rather than analyzing conceptual discrepancies in what generative models produce. The candidate addresses detector overfitting to dataset-specific patterns, not systematic concept suppression/exaggeration in model outputs.

### 6. Classification accuracy score for conditional generative models
**URL**: View paper

**Brief Assessment**

Classification Accuracy Score[64] focuses on measuring class-conditional generative model quality through downstream classification tasks, not on identifying systematic concept discrepancies or blindspots in model outputs.

### 7. GenImage: A Million-Scale Benchmark for Detecting AI-Generated Image
**URL**: View paper

**Brief Assessment**

GenImage[62] focuses on detecting AI-generated images versus real images through a benchmark dataset, not on identifying systematic concept discrepancies within generative model outputs themselves.

### 8. Lost in translation: Latent concept misalignment in text-to-image diffusion models
**URL**: View paper

**Brief Assessment**

Latent Concept Misalignment[68] focuses on misalignment between text prompts and generated images when combining disentangled concepts (e.g., 'tea cup of iced coke'), not on systematic identification of concept discrepancies across the full data distribution as in the original paper.

### 9. Deconstructing Bias: A Multifaceted Framework for Diagnosing Cultural and Compositional Inequities in Text-to-Image Generative Models
**URL**: View paper

**Brief Assessment**

Deconstructing Bias[70] focuses on cultural and compositional biases in text-to-image models using the Component Inclusion Score (CIS) metric, evaluating cultural representation disparities. The original paper addresses conceptual blindspots (systematic concept discrepancies) using sparse autoencoders on visual features, a fundamentally different technical approach and scope.

### 10. Smiling Women Pitching Down: Auditing Representational and Presentational Gender Biases in Image Generative AI
**URL**: View paper

**Brief Assessment**

Smiling Women Pitching[69] focuses on auditing gender biases in AI-generated images through representational and presentational analysis, not on identifying systematic concept discrepancies or conceptual blindspots in generative models.

## Contribution 2: Scalable unsupervised method using sparse autoencoders for concept extraction and comparison

**Description**: The authors propose an automated pipeline that leverages sparse autoencoders to decompose high-dimensional activation spaces into interpretable concepts. They train and open-source an archetypal SAE on DINOv2 features with 32,000 concepts, enabling fine-grained, unsupervised analysis of conceptual disparities without requiring human-defined concept labels.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

### 1. Mammo-sae: Interpreting breast cancer concept learning with sparse autoencoders
**URL**: View paper

**Brief Assessment**

Mammo-sae[54] applies SAEs to medical imaging (breast cancer detection) for a domain-specific task, whereas the original paper focuses on identifying conceptual blindspots in generative image models across diverse concepts. The candidate does not demonstrate prior work on using SAEs for unsupervised concept extraction in generative models.

### 2. Can sparse autoencoders make sense of latent representations?
**URL**: View paper

**Brief Assessment**

Latent Representations Sense[60] focuses on applying SAEs to gene expression and biological data, not generative image models. The candidate's domain (single-cell RNA sequencing) and objectives (interpreting biological latent variables) differ fundamentally from the original paper's focus on identifying conceptual blindspots in text-to-image models.

### 3. An enhanced sparse autoencoder for machinery interpretable fault diagnosis
**URL**: View paper

**Brief Assessment**

Machinery Fault Diagnosis[56] focuses on sparse autoencoders for machinery fault diagnosis in vibrational signals, not on extracting interpretable concepts from image representations or analyzing generative models. The application domains and technical objectives are fundamentally different.

### 4. Sparse autoencoders for scientifically rigorous interpretation of vision models
**URL**: View paper

**Brief Assessment**

Vision Models Interpretation[58] focuses on interpretability and causal manipulation of visual features in vision models, not on identifying conceptual blindspots in generative image models through concept distribution comparison between real and generated images.

### 5. Sparse Autoencoders Learn Monosemantic Features in Vision-Language Models
**URL**: View paper

**Brief Assessment**

Monosemantic Features[59] applies SAEs to vision-language models (CLIP) for neuron-level interpretability and steering, while the original work uses SAEs on DINOv2 features to extract and compare concept distributions between natural and generated images for identifying conceptual blindspots in generative models. These are distinct applications with different objectives.

### 6. Leveraging sparse autoencoders to reveal interpretable features in geophysical models
**URL**: View paper

**Brief Assessment**

Geophysical Models[53] applies sparse autoencoders to interpret hidden layer activations in geophysical neural networks, focusing on physical feature extraction from precipitation models. The original paper uses SAEs to extract visual concepts from image representations (DINOv2 features) for comparing generative models against natural images—a fundamentally different application domain and objective.

### 7. Scaling and evaluating sparse autoencoders
**URL**: View paper

**Brief Assessment**

Scaling Sparse Autoencoders[55] focuses on training sparse autoencoders on language model activations to extract interpretable features for mechanistic interpretability, not on image representations or visual concept extraction from generative image models.

### 8. Universal sparse autoencoders: Interpretable cross-model concept alignment
**URL**: View paper

**Brief Assessment**

Universal Sparse Autoencoders[52] focuses on learning shared concept spaces across multiple vision models simultaneously, while the original paper applies SAEs to compare generative vs. natural image distributions for a single model at a time.

### 9. TIDE : Temporal-Aware Sparse Autoencoders for Interpretable Diffusion Transformers in Image Generation
**URL**: View paper

**Brief Assessment**

TIDE[57] focuses on temporal-aware sparse autoencoders for diffusion transformers in image generation, not on extracting concepts from static image representations like DINOv2 for comparing generative models against natural images.

### 10. Sparse Autoencoders Find Highly Interpretable Features in Language Models
**URL**: View paper

**Brief Assessment**

Interpretable Features[51] focuses on extracting interpretable features from language models using sparse autoencoders, not on image representations or visual concepts. The candidate applies sparse coding to text-based neural network activations, while the original work targets visual concept extraction from image models like DINOv2.

## Contribution 3: Interactive exploratory tool for distribution-level and datapoint-level blindspot analysis

**Description**: The authors develop and release an interactive web-based tool that allows researchers to explore conceptual blindspots at multiple granularities. The tool supports visualization of concept distributions via UMAP, inspection of individual concepts with representative images, and identification of memorization artifacts and compositional failures across different generative models.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

### 1. Interactive Concept Bottleneck Models
**URL**: View paper

**Brief Assessment**

Interactive Concept Bottleneck[74] focuses on interactive prediction systems that query humans for concept labels to improve classification accuracy, not on analyzing generative model concept distributions or failures. The tool serves a fundamentally different purpose—optimizing human-AI collaboration for prediction tasks rather than exploring conceptual blindspots in generative models.

### 2. Asap: Interpretable analysis and summarization of ai-generated image patterns at scale
**URL**: View paper

**Brief Assessment**

Asap[76] focuses on analyzing AI-generated image patterns through pixel-level saliency and authenticity detection, not on concept-level distribution analysis of generative models. The candidate addresses fake image detection workflows, while the original paper examines conceptual blindspots in generative model outputs using sparse autoencoders.

### 3. What is a fair diffusion model? designing generative text-to-image models to incorporate various worldviews
**URL**: View paper

**Brief Assessment**

Fair Diffusion Model[78] focuses on exposing and editing worldview biases in text-to-image models through demographic analysis, not on systematic concept distribution analysis or compositional failure detection across model architectures.

### 4. Gan lab: Understanding complex deep generative models using interactive visual experimentation
**URL**: View paper

**Brief Assessment**

GAN Lab[77] focuses on interactive visualization for learning and experimenting with GANs' training dynamics, not on analyzing concept distributions or identifying blindspots in generative models. The tools serve fundamentally different purposes—education versus systematic evaluation.

### 5. Large scale qualitative evaluation of generative image model outputs
**URL**: View paper

**Brief Assessment**

Qualitative Evaluation Scale[79] focuses on large-scale visual inspection of generative model outputs through clustering and metrics visualization, but does not address conceptual blindspot analysis using sparse autoencoders or energy-based frameworks as described in the original paper.

### 6. MAVIDSQL: A Model-Agnostic Visualization for Interpretation and Diagnosis of Text-to-SQL Tasks
**URL**: View paper

**Brief Assessment**

MAVIDSQL[80] focuses on text-to-SQL semantic parsing tasks with model diagnosis through input-output analysis, not on generative image model concept distributions or compositional failures in visual generation.

### 7. POET: Supporting Prompting Creativity and Personalization with Automated Expansion of Text-to-Image Generation
**URL**: View paper

**Brief Assessment**

POET[75] focuses on interactive text-to-image generation tools for creative ideation and personalization, not on analyzing conceptual blindspots or failures in generative models. The tools serve fundamentally different purposes.

### 8. Collaborative interactive evolution of art in the latent space of deep generative models
**URL**: View paper

**Brief Assessment**

Collaborative Interactive Evolution[71] focuses on interactive evolutionary computation for art generation using GANs, not on analyzing conceptual blindspots or failures in generative models through distribution-level visualization tools.

### 9. Uncovering structural ensembles from single-particle cryo-EM data using cryoDRGN
**URL**: View paper

**Brief Assessment**

cryoDRGN[73] focuses on structural biology applications (cryo-EM data analysis) rather than generative image model evaluation. The interactive visualization mentioned is for analyzing protein conformations, not concept distributions in text-to-image models.

### 10. Interactive Semantic Interventions for VLMs: A Human-in-the-Loop Investigation of VLM Failure
**URL**: View paper

**Brief Assessment**

Interactive Semantic Interventions[72] focuses on interactive semantic interventions for VLMs in visual question answering contexts, not on analyzing generative model concept distributions or blindspots in image generation models.

## Appendix: Text Similarity Detection

No high-similarity text segments were detected across any compared papers.

## References

- [0] Uncovering Conceptual Blindspots in Generative Image Models Using Sparse Autoencoders View paper
- [1] Effective data augmentation with diffusion models View paper
- [2] Make-A-Scene: Scene-Based Text-to-Image Generation with Human Priors View paper
- [3] Visual-augmented dynamic semantic prototype for generative zero-shot learning View paper
- [4] KITTEN: A Knowledge-Intensive Evaluation of Image Generation on Visual Entities View paper
- [5] Uncovering Bias in Face Generation Models View paper
- [6] Visual Evidence Prompting Mitigates Hallucinations in Large Vision-Language Models View paper
- [7] Echo-4o: Harnessing the power of gpt-4o synthetic images for improved image generation View paper
- [8] UniFusion: Vision-Language Model as Unified Encoder in Image Generation View paper
- [9] WorldGenBench: A World-Knowledge-Integrated Benchmark for Reasoning-Driven Text-to-Image Generation View paper
- [10] Counterfactual Edits for Generative Evaluation View paper
- [11] OneIG-Bench: Omni-dimensional Nuanced Evaluation for Image Generation View paper
- [12] Uncertainty-Aware ControlNet: Bridging Domain Gaps with Synthetic Image Generation View paper
- [13] Chain-of-Thought Guided Semantic Debiasing for Low-Shot Vision-Language Tasks View paper
- [14] AI's Blind Spots: Geographic Knowledge and Diversity Deficit in Generated Urban Scenario View paper
- [15] Topic modeling of video and image data: a visual semantic unsupervised approach View paper

- [16] Contrastive Visual Data Augmentation View paper
- [17] Bridging the Intent Gap: Knowledge-Enhanced Visual Generation View paper
- [18] From Creation to Curriculum: Examining the role of generative AI in Arts Universities View paper
- [19] Exposing Blindspots: Cultural Bias Evaluation in Generative Image Models View paper
- [20] Experimental Quantum Generative Adversarial Networks for Image Generation View paper
- [21] Qualitative Failures of Image Generation Models and Their Application in Detecting Deepfakes View paper
- [22] Synthesizing training data for intelligent weed control systems using generative AI View paper
- [23] Radiology report generation with medical knowledge and multilevel image-report alignment: A new method and its verification View paper
- [24] Uncovering Limitations in Text-to-Image Generation: A Contrastive Approach with Structured Semantic Alignment View paper
- [25] Gelda: A generative language annotation framework to reveal visual biases in image generators View paper
- [26] Image-based video game asset generation and evaluation using deep learning: a systematic review of methods and applications View paper
- [27] A progressive distillation network for practical image-based virtual try-on View paper
- [28] CuRe: Cultural Gaps in the Long Tail of Text-to-Image Systems View paper
- [29] Probabilistic Modeling of Semantic Ambiguity for Scene Graph Generation View paper
- [30] Domain gap embeddings for generative dataset augmentation View paper
- [31] ToViLaG: Your Visual-Language Generative Model is Also An Evildoer View paper
- [32] Align Beyond Prompts: Evaluating World Knowledge Alignment in Text-to-Image Generation View paper
- [33] Styleguide: Zero-shot sketch-based image retrieval using style-guided image generation View paper
- [34] xgen-videosyn-1: High-fidelity text-to-video synthesis with compressed representations View paper
- [35] Can large multimodal models uncover deep semantics behind images? View paper
- [36] Self-supervised learning disentangled group representation as feature View paper
- [37] Concept Lens: Visual Comparison and Evaluation of Generative Model Manipulations View paper
- [38] A review of multi-modal learning from the text-guided visual processing viewpoint View paper
- [39] Label2im: Knowledge Graph Guided Image Generation from Labels View paper
- [40] A Practical Investigation of Spatially-Controlled Image Generation with Transformers View paper
- [41] Unsafe Synthetic Image Generation View paper
- [42] RealRAG: Retrieval-augmented Realistic Image Generation via Self-reflective Contrastive Learning View paper
- [43] Uncovering the limits of visual-language models in engineering knowledge representation View paper
- [44] Where Culture Fades: Revealing the Cultural Gap in Text-to-Image Generation View paper
- [45] World-To-Image: Grounding Text-to-Image Generation with Agent-Driven World Knowledge View paper
- [46] Image Generation is May All You Need for VQA View paper
- [47] Generative Semantic Probing for Vision-Language Models via Hierarchical Feature Optimization View paper
- [48] Investigating Conceptual Blending of a Diffusion Model for Improving Nonword-to-Image Generation View paper
- [49] Unspoken Details: Inferring Hidden Causality and Retrieving Domain-Specific Knowledge for Image Generation View paper
- [50] Open Multimodal Retrieval-Augmented Factual Image Generation View paper
- [51] Sparse Autoencoders Find Highly Interpretable Features in Language Models View paper
- [52] Universal sparse autoencoders: Interpretable cross-model concept alignment View paper
- [53] Leveraging sparse autoencoders to reveal interpretable features in geophysical models View paper
- [54] Mammo-sae: Interpreting breast cancer concept learning with sparse autoencoders View paper
- [55] Scaling and evaluating sparse autoencoders View paper
- [56] An enhanced sparse autoencoder for machinery interpretable fault diagnosis View paper
- [57] TIDE : Temporal-Aware Sparse Autoencoders for Interpretable Diffusion Transformers in Image Generation View paper
- [58] Sparse autoencoders for scientifically rigorous interpretation of vision models View paper
- [59] Sparse Autoencoders Learn Monosemantic Features in Vision-Language Models View paper
- [60] Can sparse autoencoders make sense of latent representations? View paper
- [61] Vipera: Towards systematic auditing of generative text-to-image models at scale View paper
- [62] GenImage: A Million-Scale Benchmark for Detecting AI-Generated Image View paper
- [63] Fourier Spectrum Discrepancies in Deep Network Generated Images View paper
- [64] Classification accuracy score for conditional generative models View paper
- [65] Tibet: Identifying and evaluating biases in text-to-image generative models View paper
- [66] Breaking semantic artifacts for generalized ai-generated image detection View paper
- [67] Exposing the Fake: Effective Diffusion-Generated Images Detection View paper
- [68] Lost in translation: Latent concept misalignment in text-to-image diffusion models View paper
- [69] Smiling Women Pitching Down: Auditing Representational and Presentational Gender Biases in Image Generative AI View paper
- [70] Deconstructing Bias: A Multifaceted Framework for Diagnosing Cultural and Compositional Inequities in Text-to-Image Generative Models View paper
- [71] Collaborative interactive evolution of art in the latent space of deep generative models View paper
- [72] Interactive Semantic Interventions for VLMs: A Human-in-the-Loop Investigation of VLM Failure View paper
- [73] Uncovering structural ensembles from single-particle cryo-EM data using cryoDRGN View paper
- [74] Interactive Concept Bottleneck Models View paper
- [75] POET: Supporting Prompting Creativity and Personalization with Automated Expansion of Text-to-Image Generation View paper
- [76] Asap: Interpretable analysis and summarization of ai-generated image patterns at scale View paper
- [77] Gan lab: Understanding complex deep generative models using interactive visual experimentation View paper
- [78] What is a fair diffusion model? designing generative text-to-image models to incorporate various worldviews View paper
- [79] Large scale qualitative evaluation of generative image model outputs View paper
- [80] MAVIDSQL: A Model-Agnostic Visualization for Interpretation and Diagnosis of Text-to-SQL Tasks View paper