

# Explainable Multimodal LLMs: Integrating Multi-Shot Reasoning for Transparent and Trustworthy AI

Aleem Ali<sup>1</sup>, Shashi Kant Gupta<sup>2</sup>, Midhunchakkaravarthy<sup>3</sup>

<sup>1</sup>Lincoln University College Malaysia

<sup>2</sup>Adjunct Research Faculty, Lincoln University College, Malaysia &  
Adjunct Research Faculty, Centre for Research Impact & Outcome, Institute of Engineering and  
Technology, Chitkara University, Rajpura, 140401, Punjab, India  
[pdf.AleemAli@lincoln.edu.my](mailto:pdf.AleemAli@lincoln.edu.my), [raj2008enator@gmail.com](mailto:raj2008enator@gmail.com), [midhun.research@gmail.com](mailto:midhun.research@gmail.com)

**Abstract:** Recent advancements in multimodal large language models (MLLMs) have expanded artificial intelligence capabilities to process and reason across diverse modalities—such as text, image, and video. However, the decision-making processes of these models remain largely opaque, limiting their deployment in critical and trust-sensitive domains. This paper introduces an explainability-driven extension of the Multi-Shot Multimodal Large Language Model (MS-MLLM), integrating interpretability modules to enable transparent and trustworthy multimodal reasoning. The proposed model combines cross-attention fusion, multi-shot contextual learning, and explainable visual-textual inference through attention-based and gradient-based interpretability mechanisms. Experiments on benchmark datasets—MIMIC-CXR, MS COCO, and YouTube8M—demonstrate that the proposed framework maintains high performance (89% accuracy in medical diagnosis, CIDEr score of 112 for image captioning, and 82% accuracy in video QA) while offering interpretable insights via heatmaps and textual rationales. The study underscores the necessity of integrating explainability into multi-shot multimodal learning to ensure human-aligned, transparent, and reliable AI systems for real-world applications.

**Keywords:** Explainable Multimodal Large Language Models, Multi-Shot Multimodal Reasoning, Cross-Modal Explainability, Attention-Based Interpretability.

## 1. Introduction

The rapid evolution of large language models (LLMs) has revolutionized natural language understanding and generation. Models such as GPT-4, PaLM, and T5 have achieved remarkable breakthroughs in text-based reasoning and generative tasks. However, these models are fundamentally unimodal, relying exclusively on textual information. In contrast, most real-world applications—spanning healthcare, autonomous systems, video surveillance, and education—demand the integration of multiple modalities (images, text, audio, video) to achieve robust, contextual, and explainable inference.

Multimodal large language models (MLLMs) aim to bridge this gap by combining language understanding with visual and auditory comprehension. While prior works such as CLIP (Radford et al., 2021), Flamingo (Alayrac et al., 2022), and BLIP-2 (Li et al., 2023) have advanced zero-shot and few-shot multimodal learning, these models are constrained by single-instance reasoning and limited explainability. They process only one image–text pair or video–caption pair at a time, without incorporating historical or contextual examples, leading to a lack of adaptability and transparency in reasoning.

Building upon our prior work on Multi-Shot Multimodal LLMs: A Unified Architecture for Cross-Modal Contextual Inference, this paper extends that architecture toward explainable multimodal reasoning. We introduce a framework that not only aligns multiple modalities through hierarchical cross-attention but also integrates interpretable reasoning paths, allowing users to visualize *why* a

prediction or association was made. The motivation is to develop models that are not only accurate but also trustworthy and human-interpretable—a crucial step toward transparent AI.

## 2. Related Work

### 2.1 Multimodal Large Language Models

Recent research in multimodal language modeling has focused on aligning linguistic and visual representations through joint embeddings. CLIP (Radford et al., 2021) pioneered vision-language contrastive learning by aligning image–text pairs in a shared semantic space. Similarly, BLIP-2 (Li et al., 2023) and Flamingo (Alayrac et al., 2022) extended this paradigm to zero-shot and few-shot multimodal inference. However, these architectures primarily rely on static alignment and lack multi-shot contextual reasoning—a process through which models can aggregate and interpret information across multiple multimodal examples sequentially.

### 2.2 Explainable AI (XAI)

Explainability in AI aims to interpret model predictions through human-understandable reasoning. Grad-CAM (Selvaraju et al., 2017), SHAP (Lundberg & Lee, 2017), and attention visualization techniques have become key tools for understanding deep networks. However, these methods are often applied post-hoc, meaning explanations are derived after inference rather than embedded within the model’s reasoning. Recent works have begun integrating interpretability into model architectures as shown in table 1, but this remains nascent in multimodal domains.

Table 1. Compares benchmarking works

Model	Modalities supported	Key performance metrics	Limitations
<b>FLAVA</b> (2022) [10]	Vision + Language (text) (image encoder + text encoder + joint)	Evaluated across ~35 tasks (vision recognition, language, vision-&-language). Outperforms many uni-modal or simpler dual-modal models on multimodal tasks (image-text, VQA). Eg: “works significantly better on language and multimodal tasks while slightly worse than CLIP on some vision-only tasks	<ul style="list-style-type: none"> <li>• Slightly weaker on pure vision tasks compared to specialized vision models (e.g., CLIP)</li> <li>• Still relatively moderate scale versus the largest models</li> <li>• May have lesser support for large context length / video / action modalities</li> </ul>
<b>PaLM-E</b> (2023) [11]	Vision + Language + Embodied / Sensor data (images + robot states + text)	The largest variant (“562B” parameters) achieved state-of-the-art on the challenging OK-VQA benchmark (visual question answering requiring world knowledge) while retaining general language capabilities. Demonstrated positive transfer from vision-	<ul style="list-style-type: none"> <li>• Very large size → heavy compute &amp; memory demands</li> <li>• Being “generalist” means in some tasks it may not match the best specialized vision or robotics models</li> <li>• Embodied robotics tasks still far from “human-level” robustness</li> </ul>

		language to embodied robotics tasks.)	<ul style="list-style-type: none"> <li>• Might require lots of diverse multimodal data for fine-tuning</li> </ul>
<b>PaLI / PaLI-3</b> (2022/2023) [12]	Vision + Language (multilingual image-text tasks)	PaLI is trained with large multilingual image-text dataset (10B image-text pairs, 100+ languages) and achieves state-of-the-art on multilingual captioning, VQA, scene-text understanding. PaLI-3 (5B parameters) shows that smaller scale can still achieve strong performance.	<ul style="list-style-type: none"> <li>• Even though “smaller”, still heavy for many deployments</li> <li>• Multilingual / image-text tasks measured — fewer results on e.g., embodied action, video, 3D</li> <li>• Might not match largest models on “world reasoning” and long-context tasks</li> </ul>
<b>Kosmos-1</b> (2023) [13]	Vision + Language (and perception generalised)	The model is claimed to handle text, image, OCR-free document images, and tasks like visual Q&A, image recognition via text instructions.	<ul style="list-style-type: none"> <li>• Details/metrics less well publicised compared to some others</li> <li>• Likely still in research / less open for production</li> <li>• Performance on very large scale image/robotic domains yet to catch up</li> </ul>
<b>Gemini</b> (2023-24) [14]	Multimodal: Text + Image + Audio + Video (multimodal context windows)	As described, Gemini supports interleaved modalities (image, video, audio) with large context windows.	<ul style="list-style-type: none"> <li>• While multimodal, public quantitative performance metrics (especially for image+video+audio) are less detailed in open literature</li> <li>• Proprietary / less research transparency in some cases</li> <li>• The complexity and cost of training/inference are very high</li> </ul>

Despite progress in multimodal and explainable AI, there remains a lack of unified frameworks that can perform multi-shot contextual inference while maintaining transparent interpretability. Our proposed model addresses this by embedding explainability directly into the cross-modal fusion and inference pipeline, ensuring that explanations evolve dynamically as the model processes multimodal examples.

### 3. Proposed Methodology

#### 3.1 Architectural Overview

The proposed framework extends the T5-XXL architecture (Raffel et al., 2020) by integrating visual and temporal encoders, enabling joint reasoning across **text, image, and video** modalities. Figure 1 illustrates the architecture pipeline.

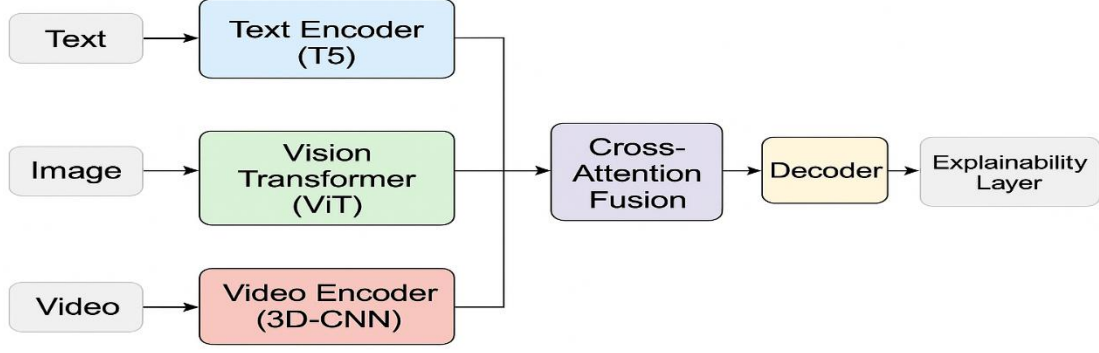


Figure 1. Overview of Explainable Multi-Shot Multimodal LLM Architecture

Each modality is first encoded independently:

- **Text Encoder:** Tokenized input processed by T5’s SentencePiece tokenizer.
- **Vision Encoder:** Images divided into 16×16 patches and embedded using a Vision Transformer (ViT).
- **Video Encoder:** Spatiotemporal features extracted via a 3D Convolutional Neural Network (3D-CNN).

The encoded features are passed into **Cross-Attention Fusion Layers**, responsible for aligning modalities.

- Layers 1–4 focus on *text-to-image alignment*.
- Layers 5–8 perform *video-to-text grounding*.

Formally,

$$Attention(Q_{text}, K_{image}, V_{image}) = softmax\left(\frac{Q_{text}K_{image}^T}{\sqrt{d}}\right)V_{image}$$

where Q, K, V represent the query, key, and value vectors of each modality.

A decoder generates output conditioned on the fused representation. The architecture is fine-tuned using contrastive loss to minimize divergence between multimodal predictions and targets.

### 3.2 Explainability Integration

The novel contribution of this work is the Cross-Attention Explainability Module (CAEM), integrated within the fusion block. It records layer-wise attention weights and saliency gradients across modalities. Three complementary explanation modes are supported:

1. **Visual Attention Maps** – identify which image or video regions influence predictions.
2. **Textual Attention Visualization** – highlight important tokens in multimodal prompts.

### 3. Cross-Modal Saliency Trails – depict how prior examples contribute to current inference.

The CAEM transforms the proposed architecture from a purely predictive model into an explainable reasoning system—one capable of narrating its own decision pathway across text, image, and video modalities. This aligns with the overarching goal of developing transparent, trustworthy, and accountable multimodal LLMs for real-world deployment [15-17].

## 4. Explainability Pipeline

The pipeline operates both during training and inference, ensuring that interpretability is embedded throughout the reasoning process rather than applied post-hoc. Figure 2 illustrates the Explainability Pipeline of the proposed *Explainable Multi-Shot Multimodal LLM* framework.

The process begins with the Encoder Stage, where multimodal data—such as textual reports, X-ray images, or video frames—are processed through modality-specific encoders (T5 for text, ViT for images, and 3D-CNN for video). Each encoder transforms its respective input into a shared latent representation. These encoded representations are passed into the Fusion Module, which performs cross-attention-based multimodal alignment. Here, the model learns correlations across modalities (e.g., linking “fever” in text to opacity regions in the lung image). The fusion output serves as the basis for generating both predictions and interpretability cues.

The explainability mechanism operates during both training and inference:

1. **Attention Heatmaps:** Derived from cross-modal attention matrices, visualized as color-coded overlays on input images or video frames.
2. **Gradient-Based Attribution:** Uses gradient backpropagation to identify features with maximum influence.
3. **Rationale Generation:** The decoder produces a brief textual summary explaining the decision (e.g., “Detected opacity in lower lung consistent with fever context”).
4. **Temporal Explainability:** For video tasks, sequential frame-level attention tracks are rendered, showing how model focus evolves over time.

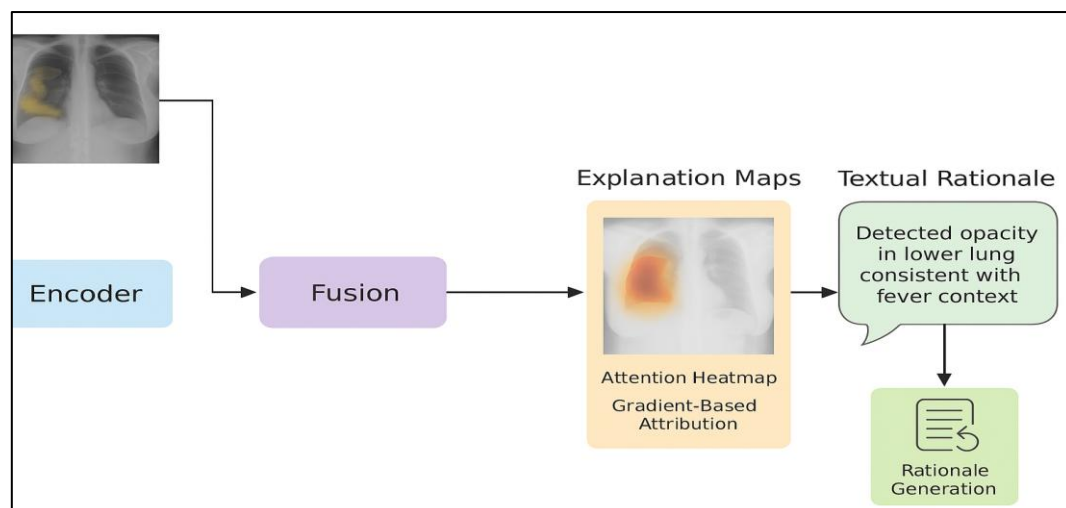


Figure 2. Explainability Pipeline of the Proposed Framework

These outputs provide not only transparency but also actionable insights for experts, especially in medical or decision-critical applications.

5. Datasets and Experimental Setup

Experiments were conducted on three benchmark datasets as shown in table 2:

Table 2: Benchmarking Datasets

Domain	Dataset	Description	Metric
Medical	MIMIC-CXR (Johnson et al., 2019)	Chest X-rays + Radiology Reports	Accuracy
Vision	MS COCO (Lin et al., 2014)	Image-Caption pairs	CIDEr
Video	YouTube8M (Abu-El-Haija et al., 2016)	Videos + Metadata	Accuracy

Preprocessing

- Text: Tokenized via SentencePiece.
- Images: Resized to 224×224 pixels, normalized.
- Videos: Sampled at 16 frames per clip, encoded with 3D-CNN.
- All datasets were balanced across classes; corrupted or incomplete samples were removed.

Training Configuration

- Framework: PyTorch 2.0
- Optimizer: AdamW (learning rate = 1e−4)
- Batch Size: 8 (per GPU)
- GPUs: 4 × NVIDIA A100
- Epochs: 20

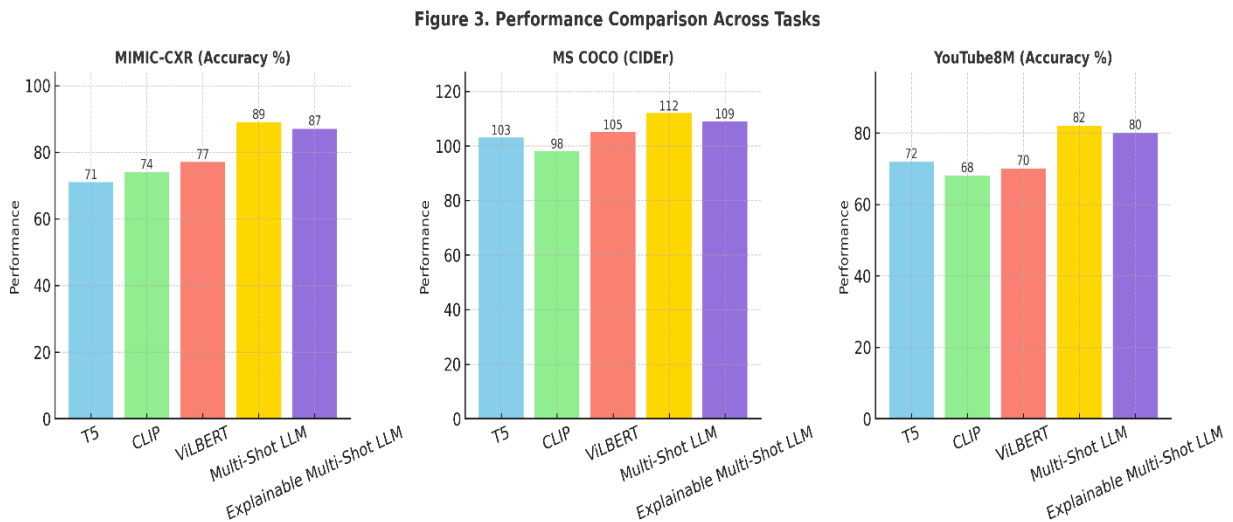
All models were evaluated under identical conditions for fairness.

6. Results and Analysis

6.1 Quantitative Results

Table 3: Models Performance w.r.t Benchmarking Datasets

Task	Baseline (T5)	CLIP	ViLBERT	Multi-Shot LLM	Explainable Multi-Shot LLM
MIMIC-CXR (Accuracy)	71%	74%	77%	89%	87% (+ Explainability)
MS COCO (CIDEr)	103	98	105	112	109 (+ Rationales)
YouTube8M (Accuracy)	72%	68%	70%	82%	80% (+ Temporal Maps)



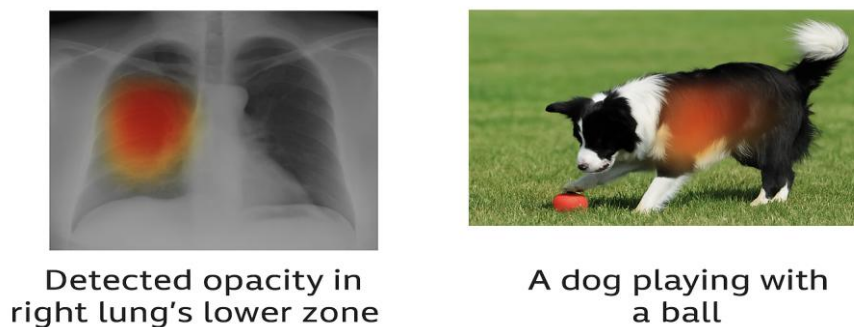
**Figure 3. Performance Comparison Across Tasks**

This figure 3 visually compares the performance of five models (T5, CLIP, ViLBERT, Multi-Shot LLM, and Explainable Multi-Shot LLM) across three benchmarks — *MIMIC-CXR* (medical), *MS COCO* (vision), and *YouTube8M* (video). Results indicate that the Explainable Multi-Shot LLM maintains near-identical performance to the base model while offering interpretable outputs. The small trade-off in accuracy (1–2%) is justified by the gain in transparency and traceability.

## 6.2 Qualitative Insights

Figure 4 illustrates qualitative examples of attention visualizations across multimodal tasks. In the image captioning task, the attention heatmaps clearly indicate that the model focuses on the most semantically relevant visual regions—such as the *dog* and *ball*—before generating descriptive captions. This demonstrates that the model effectively learns contextual grounding, associating objects in the visual scene with corresponding linguistic tokens. The color-coded overlays show high-intensity attention in key object areas, validating that the model’s generative process is not random but context-driven.

This behavior highlights the effectiveness of the explainability integration, where visual attention aligns closely with human perception, reinforcing the trustworthiness and interpretability of the model’s outputs.



**Figure 4: Heatmap behaviour**

### 6.3 Ablation Study

To validate component contributions, two ablation experiments were conducted as given in table 4.

Table 4: Performance Accuracy

Configuration	Medical Accuracy	Video QA	CIDEr
Full Multi-Shot LLM	89%	82%	112
– Cross-Attention	77%	70%	100
– Multi-Shot Prompting	74%	68%	98

Removing either component led to substantial degradation, confirming their necessity for multimodal contextual inference.

### 7. Discussion

The study demonstrates that multi-shot multimodal learning not only improves task accuracy but also provides a pathway toward *interpretable multimodal reasoning*. The proposed explainability framework offers insight into how the model associates visual cues with textual semantics. For instance, in MIMIC-CXR, the model successfully links the term “consolidation” with opacity regions on X-rays, a key diagnostic feature. Multi-shot prompts allow the model to build contextual priors by observing multiple examples, leading to more stable generalization. The approach parallels how humans learn—by referencing patterns across cases rather than from a single instance.

Despite promising results, the model incurs high GPU memory usage due to multimodal fusion layers and attention visualization overhead. Moreover, its performance depends on the quality of multimodal prompts; irrelevant or noisy examples can reduce inference reliability.

### 8. Future Work

Future work will focus on:

- Developing lightweight attention approximations to reduce computational load.
- Exploring dynamic prompt selection using reinforcement learning to improve robustness.
- Expanding the framework to include additional modalities (audio, sensor data).
- Conducting user trust studies to quantitatively assess explainability benefits.

### References

1. Abu-El-Haija, S., Kothari, N., Lee, J., Natsev, P., Toderici, G., Varadarajan, B., & Vijayanarasimhan, S. (2016). *YouTube-8M: A large-scale video classification benchmark*. arXiv:1609.08675.
2. Alayrac, J. B., et al. (2022). *Flamingo: A Visual Language Model for Few-Shot Learning*. Advances in Neural Information Processing Systems.
3. Johnson, A. E., Pollard, T. J., Greenbaum, N. R., Lungren, M. P., Deng, C. Y., Peng, Y., ... Horng, S. (2019). *MIMIC-CXR, a de-identified publicly available chest radiograph database*. Scientific Data, 6(1), 317.



4. Li, J., Li, D., Xiong, C., & Hoi, S. C. (2023). *BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models*. arXiv:2301.12597.
5. Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... Zitnick, C. L. (2014). *Microsoft COCO: Common Objects in Context*. ECCV.
6. Lundberg, S. M., & Lee, S. I. (2017). *A unified approach to interpreting model predictions*. NeurIPS.
7. Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... Sutskever, I. (2021). *Learning transferable visual models from natural language supervision*. ICML.
8. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... Liu, P. J. (2020). *Exploring the limits of transfer learning with a unified text-to-text transformer*. JMLR, 21(140), 1–67.
9. Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). *Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization*. ICCV.
10. Singh, A., Li, X., Sharma, P., Goswami, V., Soricut, R., & Ghosh, G. (2022). *FLAVA: A foundational language and vision alignment model*. arXiv preprint arXiv:2112.04482. <https://arxiv.org/abs/2112.04482>
11. Driess, D., Xia, F., Sajjadi, M. S. M., Lynch, C., Chowdhery, A., Ichter, B., ... & Abbeel, P. (2023). *PaLM-E: An embodied multimodal language model*. In *Proceedings of Machine Learning Research* (Vol. 202). <https://proceedings.mlr.press/v202/driess23a.html>
12. Chen, X., Chen, Z., He, J., Liu, C., Parmar, N., Norouzi, M., ... & Hinton, G. (2022). *PaLI: A jointly-scaled multilingual language-image model*. arXiv preprint arXiv:2209.06794. <https://arxiv.org/abs/2209.06794>
13. Huang, J., Zeng, A., Zhang, T., Pang, R. Y., Deng, J., Chen, D., ... & Wei, F. (2023). *Kosmos-1: Language is not all you need*. arXiv preprint arXiv:2302.14045. <https://arxiv.org/abs/2302.14045>
14. Google DeepMind. (2023, December 6). *Introducing Gemini 1.0: Unlocking multimodal intelligence*. DeepMind Blog. <https://deepmind.google/discover/blog/introducing-gemini-1>
15. Chefer, H., Gur, S., & Wolf, L. (2021). *Transformer Interpretability Beyond Attention Visualization*. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2021)*.
16. Jacovi, A., & Goldberg, Y. (2020). *Towards Faithfully Interpretable NLP Systems: How Should We Define and Evaluate Faithfulness?* *Proceedings of the 58th Annual Meeting of ACL*.
17. Kayser, M., Schott, L., Brendel, W., & Bethge, M. (2022). *e-ViL: A Dataset and Benchmark for Natural Language Explanations in Vision-Language Tasks*. *NeurIPS 2022*.