# *Ceci n'est pas une pomme*:
# Adversarial Illusions in Multi-Modal Embeddings

**Eugene Bagdasaryan**    **Vitaly Shmatikov**
Cornell Tech
{eugene, shmat}@cs.cornell.edu

## Abstract

Multi-modal encoders map images, sounds, texts, videos, etc. into a single embedding space, aligning representations across modalities (e.g., associate an image of a dog with a barking sound). We show that multi-modal embeddings can be vulnerable to an attack we call "adversarial illusions." Given an input in any modality, an adversary can perturb it so as to make its embedding close to that of an arbitrary, adversary-chosen input in another modality. Illusions thus enable the adversary to align any image with any text, any text with any sound, etc.

Adversarial illusions exploit proximity in the embedding space and are thus agnostic to downstream tasks. Using ImageBind embeddings, we demonstrate how adversarially aligned inputs, generated without knowledge of specific downstream tasks, mislead image generation, text generation, and zero-shot classification.

## 1   Introduction

Multi-modal embedding models, e.g., ImageBind [5], map inputs such as images, texts, and sounds into a shared embedding space. The key aspect of these models is that encoders for different modalities *align* the representations (i.e., embedding vectors) of semantically related inputs. Multi-modal embeddings thus enable downstream applications, including classification and generation, to operate on inputs indepedently of their modality.

In this paper, we show that cross-modal alignment in the ImageBind embeddings is vulnerable to adversarially generated *illusions*. An illusion aligns an input in one modality with another, adversary-chosen input in a different modality, thus "misrepresenting" the semantic content of the former to downstream tasks. Figure 1 shows an illusion that aligns an image of Magritte's famous "This is Not an Apple" painting with the text of Magritte's quote "Everything we see hides another thing." Downstream tasks—text and image generation, in this case—act on the perturbed image as if it were the quote and not the original image.

We show how to generate adversarially aligned illusions by *colliding* inputs so that they map to similar representations in the embedding space[1]. The resulting perturbations make an input appear similar to the original (as far as human eye or ear are concerned) but its embedding is close to that of an *arbitrary*, adversary-chosen input in a different modality. Any downstream application will then operate on the adversary's input instead of the original.

Cross-modal adversarial perturbations are made possible by the alignment across encoders, which is a fundamental feature of multi-modal embeddings. Although there exists a "modality gap" [7] between the representations in different modalities, our experiments show that downstream tasks based on ImageBind, including text generation, image generation, and zero-shot classification, tolerate this gap, and are consequently misled by cross-modal illusions.

---

[1]Code is available at `https://github.com/ebagdasa/adversarial_illusions`.
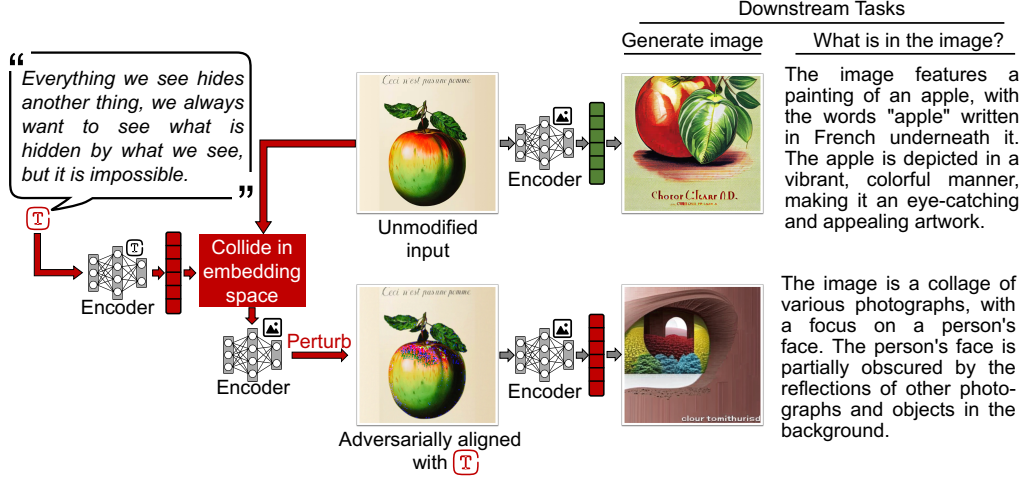
Figure 1: **"This is not an apple": an adversarial perturbation of an image changes downstream image and text generation.**

In summary, we demonstrate a novel attack on multi-modal embedding models that directly collides representations of inputs from different modalities, is agnostic of downstream applications, and enables an adversary to align an input in any modality with arbitrary content in another modality.

## 2 Background and Related Work

***Multi-modal embeddings.*** A multi-modal embedding model $\theta$ such as ImageBind [5] embeds inputs $x^m$ with modalities $m \in M$ into the same embedding space. Each modality $m$ has its own dedicated encoder $\theta^m$. The ImageBind model is trained on multi-modal tuples $(x_i^{m_1}, x_i^{m_2})$ that are semantically aligned; the training aims to maximize the dot product of the embedding representations of each tuple. For example, given a tuple consisting of a picture of a bird and the text "a singing bird," the image and the text are embedded into similar representations.

The ImageBind model exhibits "emergent" alignment between modalities. Semantically similar images and sounds (for example, a picture of a bird and an audio recording of a birdsong) have similar embeddings, even though the training data does not include image-audio tuples. In this paper, we focus on images, sounds, and text, and leave other modalities to future work.

***Downstream models.*** Downstream models based on ImageBind embeddings do not need to be trained on multiple modalities because the embeddings are supposed to align semantically similar inputs across modalities.

Image generation takes an embedding and performs conditional generation using a diffusion model. ImageBind's image encoder is initialized from the CLIP visual encoder model [9]. Therefore, diffusion models that operate on CLIP embeddings, e.g., unCLIP [10], can also operate on ImageBind embeddings. Text generation can use multi-modal embeddings as inputs to instruction-following language models, e.g., PandaGPT [13].

Zero-shot classification matches the embedding of an input to a class embedding, i.e., the mean of all inputs in a class. With ImageBind embeddings, zero-shot classification can be applied to inputs and classes that have different modalities, e.g., match an audio to an image class.

***Adversarial perturbations.*** Adding a perturbation $\delta$ to an input $x$ can cause a model $\theta$ to assign an incorrect label $y^*$ to this input, i.e., $\theta(x) = y$ while $\theta(x + \delta) = y^*$ [6]. Several recent papers [1, 3, 8] explore adversarial perturbations in multi-modal language models, focusing specifically on chatbots. By contrast, we investigate adversarial perturbations against multi-modal embeddings that are agnostic of downstream tasks.
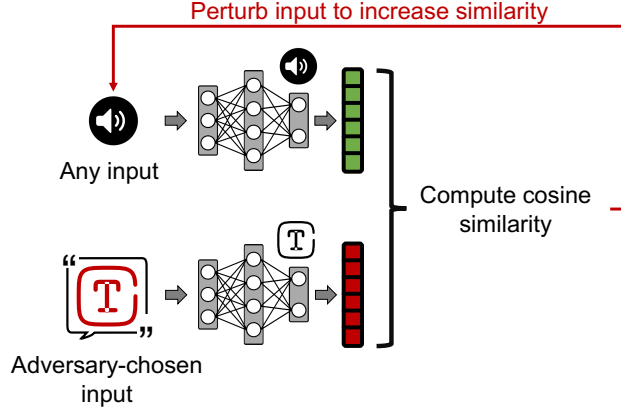
Figure 2: **Cross-modal collisions in the embedding space.**

***Collisions.*** This paper is the first to demonstrate collisions in multi-modal embeddings and adversarial alignment across modalities for arbitrary inputs. Prior work considered collisions in NLP models [12], but—unlike this work—it targets specific downstream tasks in a single modality (text) and the adversary does not have arbitrary choice of colliding inputs. Concurrent work [11] demonstrated image collisions that target multi-modal chatbots. These collisions are not cross-modal, are specific to a downstream task, and the adversary does not have arbitrary choice of colliding inputs.

## 3 Cross-Modal Illusions

We demonstrate a new attack, *cross-modal illusions*, that forces multi-modal embeddings to incorrectly align inputs from different modalities and therefore misrepresent their semantic content to downstream tasks. Our attack relies on two fundamental properties of multi-modal embedding models [5]. First, encoders in these models align input representations across modalities. Second, downstream applications that use these embeddings do not, by design, consider input modalities and can operate on any vector from the embedding space. Our attack thus assumes whitebox access to the embedding model $\theta$ but no access to or even awareness of downstream tasks.

***Cross-modal collisions.*** We show how to perturb an input in one modality so that its representation collides in the embedding space with an arbitrary, adversary-chosen input from another modality. All downstream tasks will then in effect operate on the adversary's input instead of the original.

Given an input $x^{m_1}$, we say that a perturbation $x_\mathbf{a}^{m_1}$ is an *illusion* if it appears similar to $x^{m_1}$ to a human but collides with an adversary-chosen $\mathbf{a}^{m_2}$ in the embedding space.

***Generating cross-modal collisions.*** Given $x^{m_1}$ and an adversary-chosen $\mathbf{a}^{m_2}$, we generate an adversarial perturbation $\delta$ to obtain $x_\mathbf{a}^{m_1}=x^{m_1}+\delta$ such that $\theta^{m_1}(x^{m_1} + \delta)$ collides with $\theta^{m_2}(\mathbf{a}^{m_2})$—see Figure 2. Although ImageBind uses dot product during training, modalities that are not explicitly bound, e.g., audio and text, have different normalizations. Therefore, we omit norms, use cosine similarity as the metric, and minimize the following objective:

$$\ell = 1 - \texttt{cos}(\theta^{m_1}(x^{m_1} + \delta),\ \ \theta^{m_2}(\mathbf{a}^{m_2}))$$

We iteratively update perturbation $\delta$ with the standard Fast Gradient Sign Method [6], i.e., $\delta=\delta-\epsilon\cdot\texttt{sign}\nabla_x(\ell)$, and reduce $\epsilon$ with a cosine annealing schedule.

***Modality gap.*** Our attack uses different encoders $\theta^{m_1}$ and $\theta^{m_2}$ to process $x^{m_1}$ and $\mathbf{a}^{m_2}$, respectively. Different encoders embed inputs into different sub-spaces, resulting in a modality gap [7]. It is possible that no perturbation $\delta$ produces a perfect collision $\theta^{m_1}(x^{m_1}+\delta)=\theta^{m_2}(\mathbf{a}^{m_2})$. Nevertheless, our preliminary experiments with ImageBind indicate that cross-modal illusions are effective against downstream tasks even when perfect cosine similarity is not achieved.
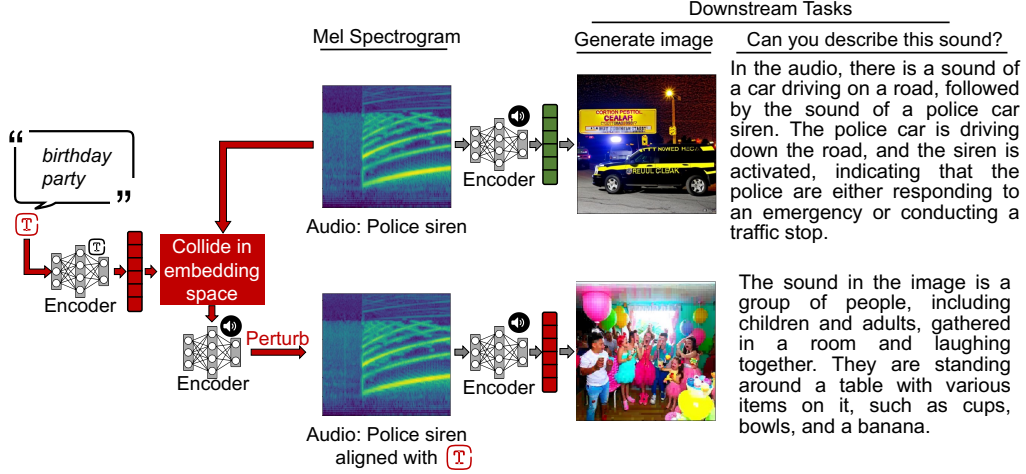
3

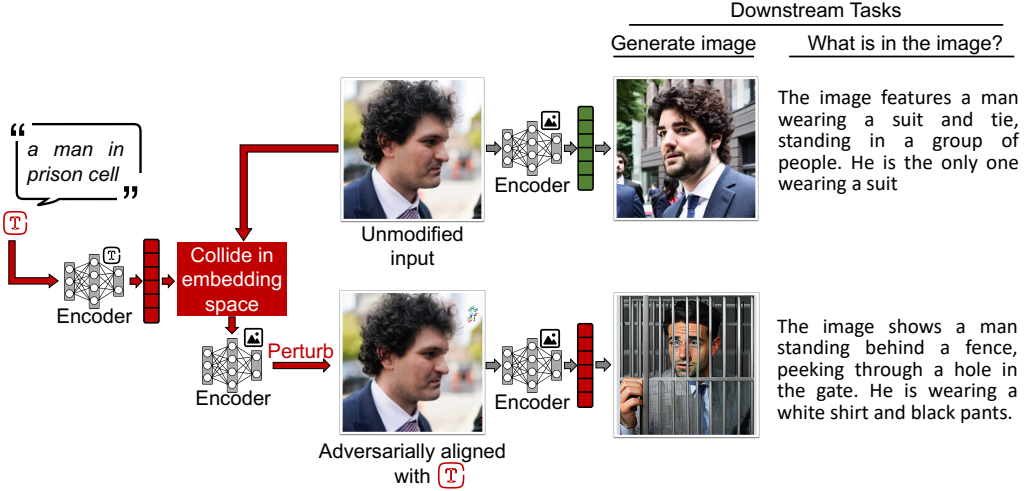Figure 3: **"Party time": an audio illusion against image and text generation.**



Figure 4: **"Schadenfreude": a visual illusion against image and text generation.**

## 4   Examples

***Setup.***  We use the ImageBind repository for the model, code, and image and audio assets [5]. To generate images from ImageBind embeddings, we use the diffusion model from the BindDiffusion repo [2]. To generate text, we use PandaGPT [13] and Vicuna [4] to answer questions "What is in this image?" and "Can you describe this sound?", respectively. For zero-shot classification, we take an audio input and label it with an ImageNet class by computing the cosine similarity of the input's embedding and the average embedding of each ImageNet class.

***Limitations of downstream models.***  Generative models based on ImageBind embeddings are a new, rapidly evolving field and publicly available implementations are limited. Image generation with BindDiffusion [2] uses the unCLIP model [10] trained on CLIP embeddings [9]. Therefore, even in the absence of adversarial perturbations, BindDiffusion fails to generate images from the embeddings of many images, sounds, and texts. PandaGPT [13], although trained on ImageBind embeddings, was fine-tuned only on image-text pairs. Therefore, in some cases it interprets embeddings of sounds as if they were images (see Figure 3).
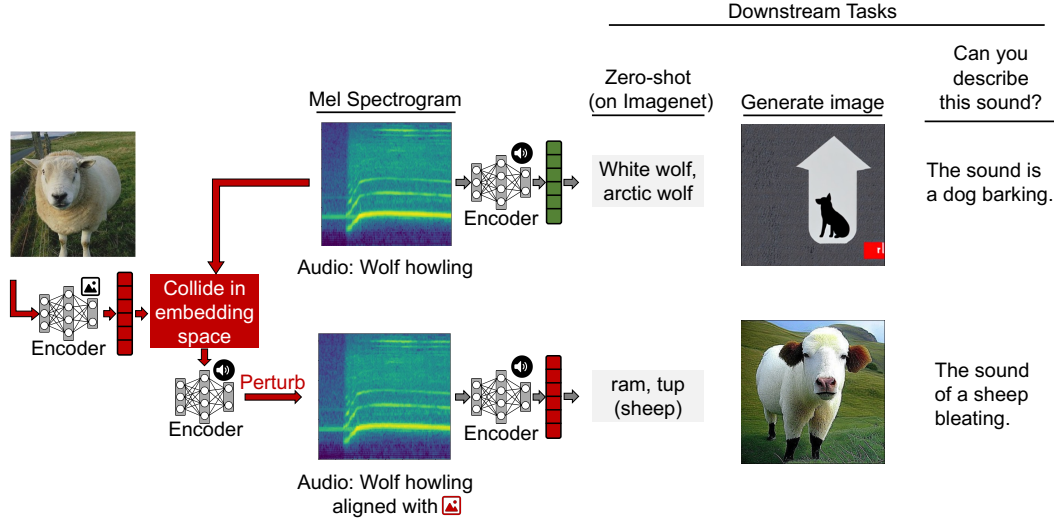
Figure 5: **"Wolf in sheep's clothing": an audio illusion against zero-shot classification.**
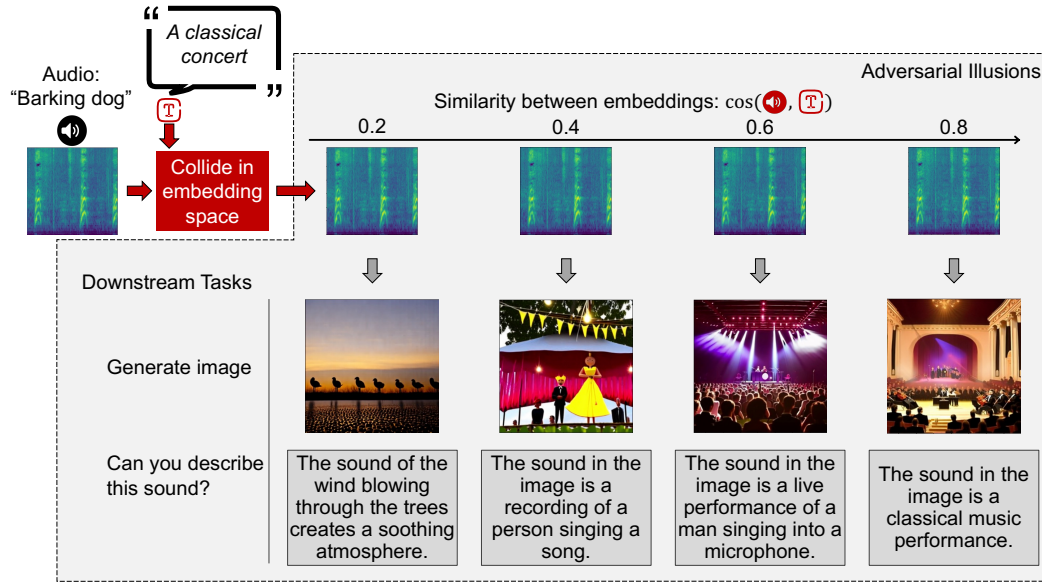


Figure 6: **"Symphony of Woofs": similarity between an adversary-chosen input and the resulting illusion.**

Our attack targets embeddings and is agnostic to downstream models. We expect that improvements in the quality of downstream models will make the attack more effective.

***Examples.*** Figures 1 and 4 show images perturbed to collide with an adversary-chosen text in the embedding space. Even though the original and perturbed images appear visually similar, images and texts generated from the embedding of the perturbed image are based on the semantics of the adversary's text, not the original image.

Figures 3 and 5 show audio illusions against generation and zero-shot classification, respectively.

Finally, we investigate how similarity between the adversary-chosen input and the resulting illusion affects downstream tasks. Figure 6 takes an audio of a barking dog and collides it with the text "a

classical concert." As cosine similarity increases, the perturbed barking-dog audio is interpreted as a classical concert by downstream tasks.

## 5 Limitations

This paper provides preliminary evidence that multi-modal embeddings are vulnerable to adversarial perturbations that create *cross-modal illusions*, i.e., inputs in one modality that are aligned with arbitrary, adversary-chosen, semantically unrelated inputs in another modality.

Our examples are based on a single model (ImageBind) and use currently available downstream generative models. These models were not explicitly trained on multi-modal embeddings: unCLIP was trained on CLIP embeddings only, and PandaGPT was fine-tuned on image-text pairs and thus interprets sounds as images. Future work can use more advanced models and also investigate other downstream tasks, such as prompt injection into multi-modal LLMs [1, 3, 8].

We used a textbook method for generating adversarial perturbations and did not evaluate how stealthy the resulting perturbations are, nor the extent to which they fool human perception. We also did not measure the alignment between the adversary's inputs and the outputs of downstream models. All of these topics are interesting subjects for future work.

## References

[1] Eugene Bagdasaryan, Tsung-Yin Hsieh, Ben Nassi, and Vitaly Shmatikov. (Ab)using images and sounds for indirect instruction injection in multi-modal LLMs. *arXiv:2307.10490*, 2023.

[2] BindDiffusion: One diffusion model to bind them all. `https://github.com/sail-sg/BindDiffusion/tree/main`, 2023.

[3] Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramer, and Ludwig Schmidt. Are aligned neural networks adversarially aligned? *arXiv:2306.15447*, 2023.

[4] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing GPT-4 with 90%* ChatGPT quality, March 2023.

[5] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. ImageBind: One embedding space to bind them all. In *CVPR*, 2023.

[6] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.

[7] Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. *NeurIPS*, 2022.

[8] Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Mengdi Wang, and Prateek Mittal. Visual adversarial examples jailbreak large language models. *arXiv:2306.13213*, 2023.

[9] Alec Radford et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.

[10] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv:2204.06125*, 2022.

[11] Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. Plug and pray: Exploiting off-the-shelf components of multi-modal models. *arXiv:2307.14539*, 2023.

[12] Congzheng Song, Alexander M Rush, and Vitaly Shmatikov. Adversarial semantic collisions. In *EMNLP*, 2020.

[13] Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. PandaGPT: One model to instruction-follow them all. *arXiv:2305.16355*, 2023.