# What Is Multimodal Generative AI and Why Does it Matter?

**Estimated time:** 5 minutes

## Objectives

By the end of this reading, you will be able to:

- Explore the concept of multimodal generative AI
- Explain the importance of multimodal generative AI

## Introduction

In this reading, you will understand the concept of **multimodal generative AI**—an emerging field where artificial intelligence (AI) models don't just understand text, but can also interpret and reason about other data types like images, audio, and even video. Think of it as teaching AI to become more like us humans: they rarely rely on just words alone; they look, listen, and observe patterns across multiple senses.

If you've ever used a text-only AI chatbot, you know that while it can be impressive, it still feels limited. Sure, it can recall facts, craft narratives, or simplify complex ideas into straightforward explanations. But it's like talking to someone who can't see or hear—helpful, yet confined. Now imagine talking to an AI that can view a picture you've sketched, "understand" what it shows, and then discuss it with you. Or one that can watch a short video and summarize its content alongside related texts you provide. This is where Multimodal AI steps in, bringing richer, more human-like understanding to the conversation.

**Note:** If any link do not open, use the secondary click to open the context menu on the link and then select "Open in a new tab".

## What is multimodal AI?

The term "multimodal" refers to the ability to process and integrate information from multiple modes or types of data. In the context of AI, these modes typically include text, images, audio, and video. While traditional AI systems are often designed to specialize in just one mode—like understanding natural language or recognizing objects in images— [multimodal AI](#) goes a step further by combining insights from multiple sources to enhance its reasoning and output.

To illustrate this, imagine reading a recipe (text) while looking at photos of the dish (images) and listening to someone explain how it's prepared (audio). Your brain naturally synthesizes all of this information to form a comprehensive understanding. Multimodal AI aims to do something similar—bridging different types of input data to achieve a richer, more holistic comprehension.

This capability enables a wide range of applications, from AI assistants that can answer questions about both a document and a diagram simultaneously, to systems that can generate image captions, translate speech to text, or even create detailed descriptions of a video. By interpreting data across multiple modalities, these systems unlock new possibilities for more interactive and intelligent AI experiences.

## What is generative AI?

In the context of AI, "generative" refers to the ability of models to create new content rather than just recognizing or analyzing existing data. Unlike traditional AI systems that are designed for tasks like classification or prediction, [generative AI](#) models learn patterns from large datasets and use that knowledge to generate new data—whether it's text, images, audio, or even video..

For example, when you use a text-based AI to write a story or draft an email, the AI isn't just pulling pre-written content from a database—it's generating entirely new sentences based on what it has learned about language. Similarly, a generative image model can create a realistic picture of a landscape that doesn't exist in real life, based on its understanding of millions of real-world images.

Generative models work by identifying underlying structures in data and using that understanding to produce something new. This capability opens up exciting possibilities, from generating realistic images for virtual worlds to synthesizing audio for speech applications or even creating personalized content in real time. Whether it's producing human-like text, photorealistic images, or custom audio tracks, generative AI is reshaping the way machines create and interact with the world.

## A changing AI landscape

The rise of multimodal generative AI has been a gradual process. For many years, AI research concentrated on developing models with specialized capabilities—excelling either at understanding text, as seen in natural language processing(NLP), or at identifying objects in images, as in computer vision (CV). Each modality—text, images, audio—existed in separate silos. Researchers developed tools and techniques like **convolutional neural networks (CNNs)** for vision tasks and **Transformers** for language tasks. But these models couldn't easily exchange information between modalities. You ended up with amazing specialists, but no true generalists.

As the field matured and as data from different domains became more abundant and easier to process, researchers began exploring how to merge these once-siloed capabilities. Pioneering models like **CLIP** OpenAI, 2021 [1], demonstrated that a single model could understand both images and text by learning a shared "language" of concepts. This paved the way for a wave of multimodal models that can caption images, answer visual questions, interpret audio, and much more.

## The industry trend: Multimodality leads the way

The surge in multimodal research isn't just academic—industry leaders are showcasing how important and timely this shift is. Consider:

- **[Meta's Llama 3.2 Vision Instruct]**, an evolution of the Llama family of large language models (LLMs), which can handle both text and images. It's designed to help users communicate naturally about visual content, offering everything from image descriptions to insights tied to specific visual elements.

- **[OpenAI o1]**, another leading-edge model, pushes the boundaries of what's possible when text and visual modalities come together. By understanding images, documents, and textual queries side-by-side, o1 exemplifies the trend of integrating multiple data types into a single, coherent intelligence.

These releases signal a broader industry movement. Companies and research labs are betting on multimodality as the future of AI interfaces. Rather than navigating through dozens of specialized tools—one model for language, another for images, yet another for audio—you can rely on a single system that understands multiple inputs at once. This unified approach streamlines workflows and unlocks more intuitive user experiences.

## Why multimodal generative AI matters for you

Think about a scenario where you're analyzing product images alongside customer feedback. A text-only model can summarize reviews. A vision-only model can tell you what's in the images. But a multimodal model can combine these views and, for example, explain how the color and design features of a product might influence the sentiments found in the text reviews. This capability helps surface deeper insights and drives more informed decisions.

Or consider working with Llama 3.2 Vision Instruct to examine a complex technical diagram alongside an instructional manual. Instead of jumping between two systems—one to "read" the manual and another to "see" the diagram—you now have a single assistant that can integrate both pieces of information and guide you toward a solution.

This isn't limited to business contexts. Multimodal generative AI can help educators create richer classroom materials, enable journalists to analyze multimedia data for stories, or assist medical professionals in reviewing patient scans together with electronic health records. As the technology continues to mature and become more accessible, you'll see it amplify productivity and creativity across countless fields.
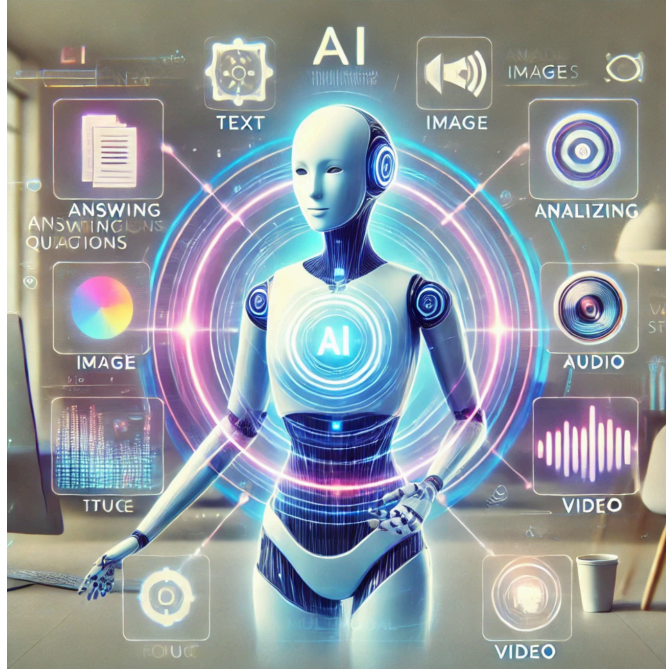
*Figure 1: Illustration of a Multimodal AI - Source: DALL-E*

### A more natural way to engage with AI

People often say that using a modern AI model is like collaborating with a brilliant assistant. Multimodal generative AI takes that analogy a step further by allowing that assistant to see and hear, not just read. With advanced models at your disposal, you can show it photographs, provide it textual descriptions, and even supply audio samples. The result is a conversation that feels more natural and comprehensive—like sharing ideas with a colleague who can both read your reports and interpret the photos you took during field research.

The rise of multimodal models by leaders like Meta and OpenAI signals more than just a passing trend—it's a transformative shift in how humans interact with AI. Research organizations, universities, and industry leaders alike recognize that the future of AI goes beyond text or vision alone; it will comprehend the world in deeper, more integrated ways.

For more in-depth exploration of the current landscape, you might consider resources like:

- "Multimodal Machine Learning: A Survey and Taxonomy" by Baltrušaitis et al. [2]
- Advances in Vision-Language Pre-training research from Stanford University [3]

### Summary

In this reading, you learned that:

- Multimodal generative AI empowers machines to perceive and reason across text, images, audio, and video—much like how humans naturally understand the world through multiple senses.
- In AI, "generative" refers to the ability of models to create new content—text, images, audio, or video—by learning patterns from data, moving beyond analysis to actual creation.
- Models like Meta's Llama 3.2 Vision Instruct and OpenAI's o1 showcase how multimodal AI is redefining intelligence by understanding and generating across text, images, and more.
- This shift toward unified, multimodal systems streamlines workflows and paves the way for more natural, intuitive user experiences across industries.
- Multimodal Generative AI extends beyond business, empowering educators, journalists, and medical professionals to work more creatively and efficiently with rich, integrated data.

### Next step

In this lab, you learned about Multimodal AI and why it's rapidly becoming a new norm in AI research and industry applications. In the labs that follow, you'll explore how to leverage these advancements, craft prompts that fully utilize them, and overcome common challenges. You're stepping into a world where AI can finally communicate in a way that feels as natural as talking to a perceptive human colleague.

## Author

Hailey Quach