# MultiModal Model

**Multimodal models** are a class of machine learning models that can process and integrate information from **multiple modalities** or types of data. Unlike traditional models that handle a single type of input (e.g., text or image), multimodal models can simultaneously handle and learn from different types of inputs such as text, images, audio, and more. Pretty cool!
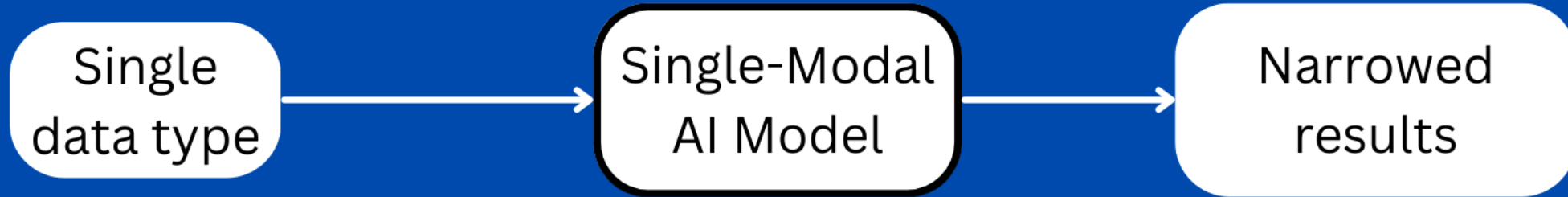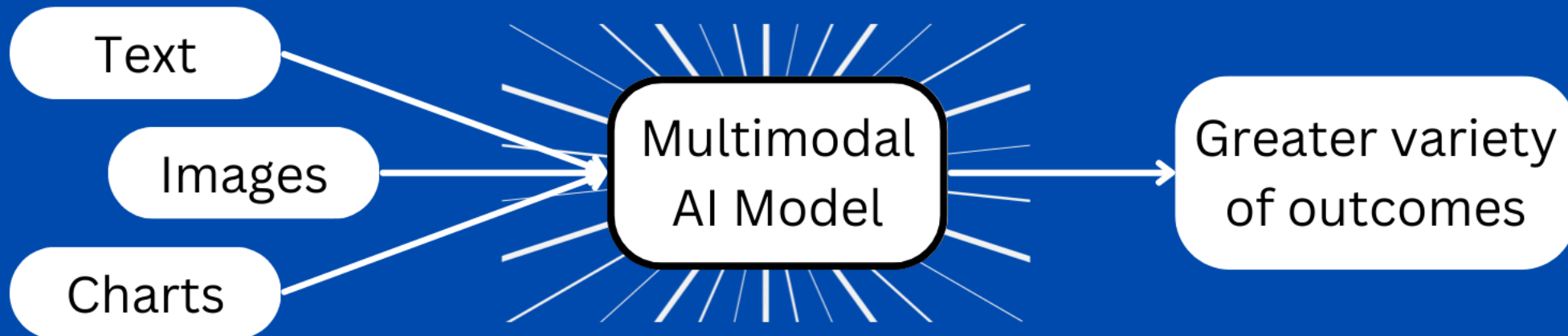


|  |  |  |  |
|---|---|---|---|
| Text | Image | Video | Audio |

Figure 1: The different types of data that multimodal models can handle simultaneously.

# Single-modal AI Model

Single data type → Single-Modal AI Model → Narrowed results

# Multimodal AI Model?

Text, Images, Charts → Multimodal AI Model → Greater variety of outcomes
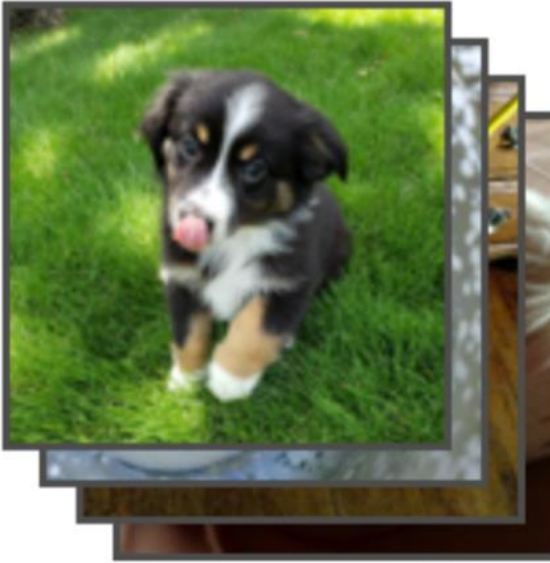
# What is CLIP?

**CLIP (Contrastive Language-Image Pre-training)** is a multimodal model developed by OpenAI that can understand and relate images and textual descriptions in a unified manner. CLIP was introduced in early 2021 [1] and represents a *significant advancement* in the field of multimodal learning.
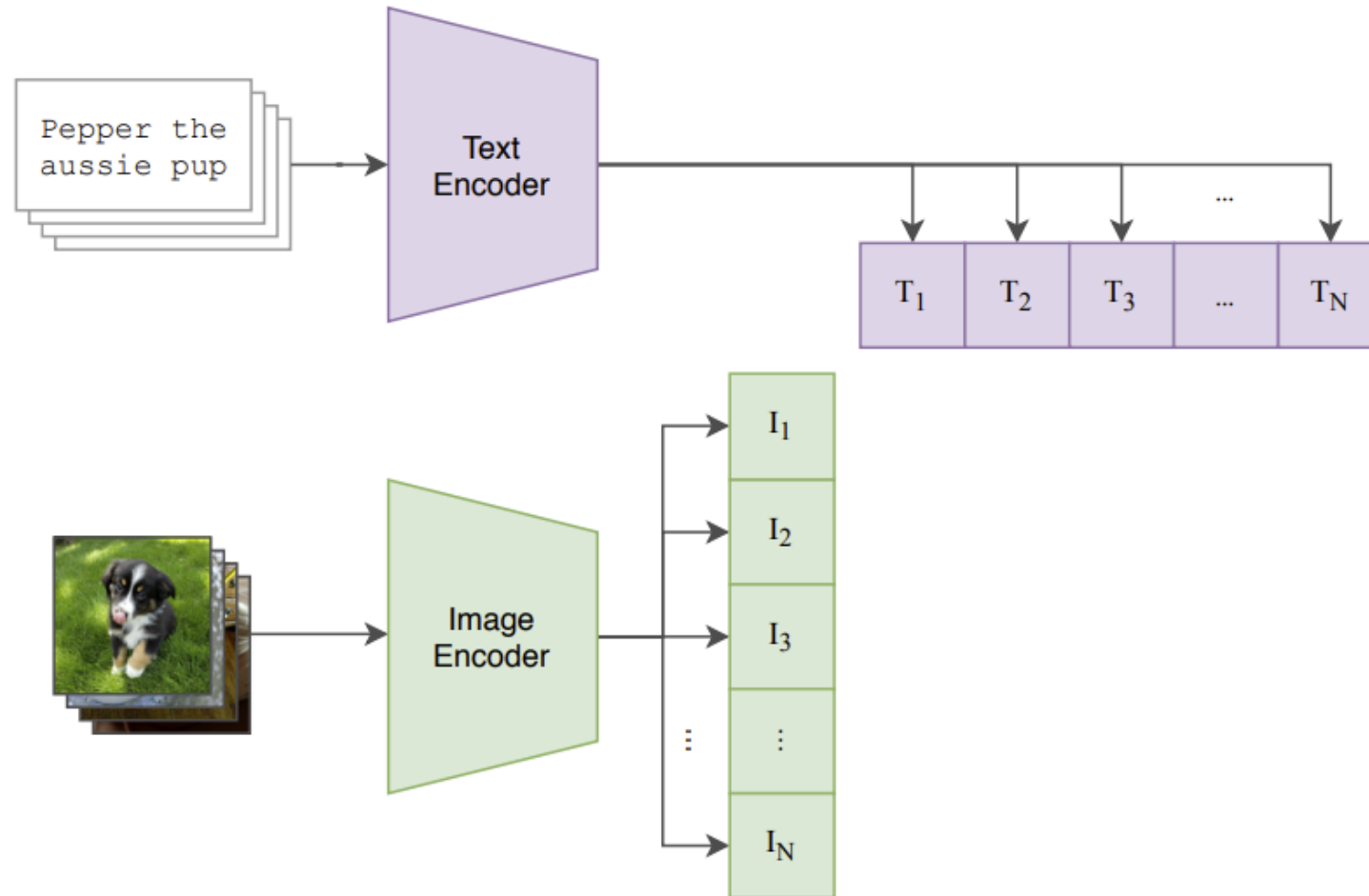


Figure 2: An image with its corresponding caption, 'Marqo logo'.

# How to train CLIP?

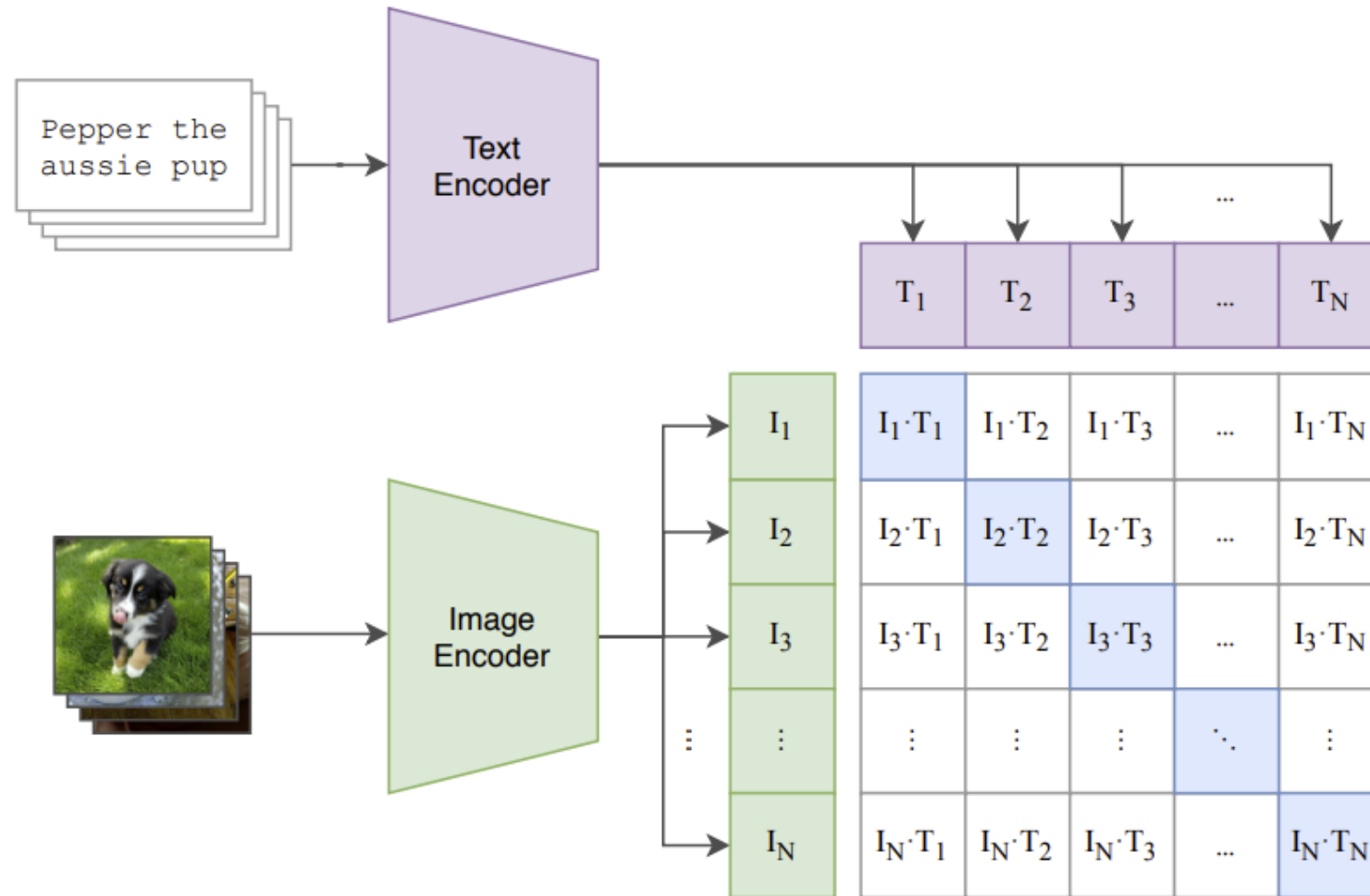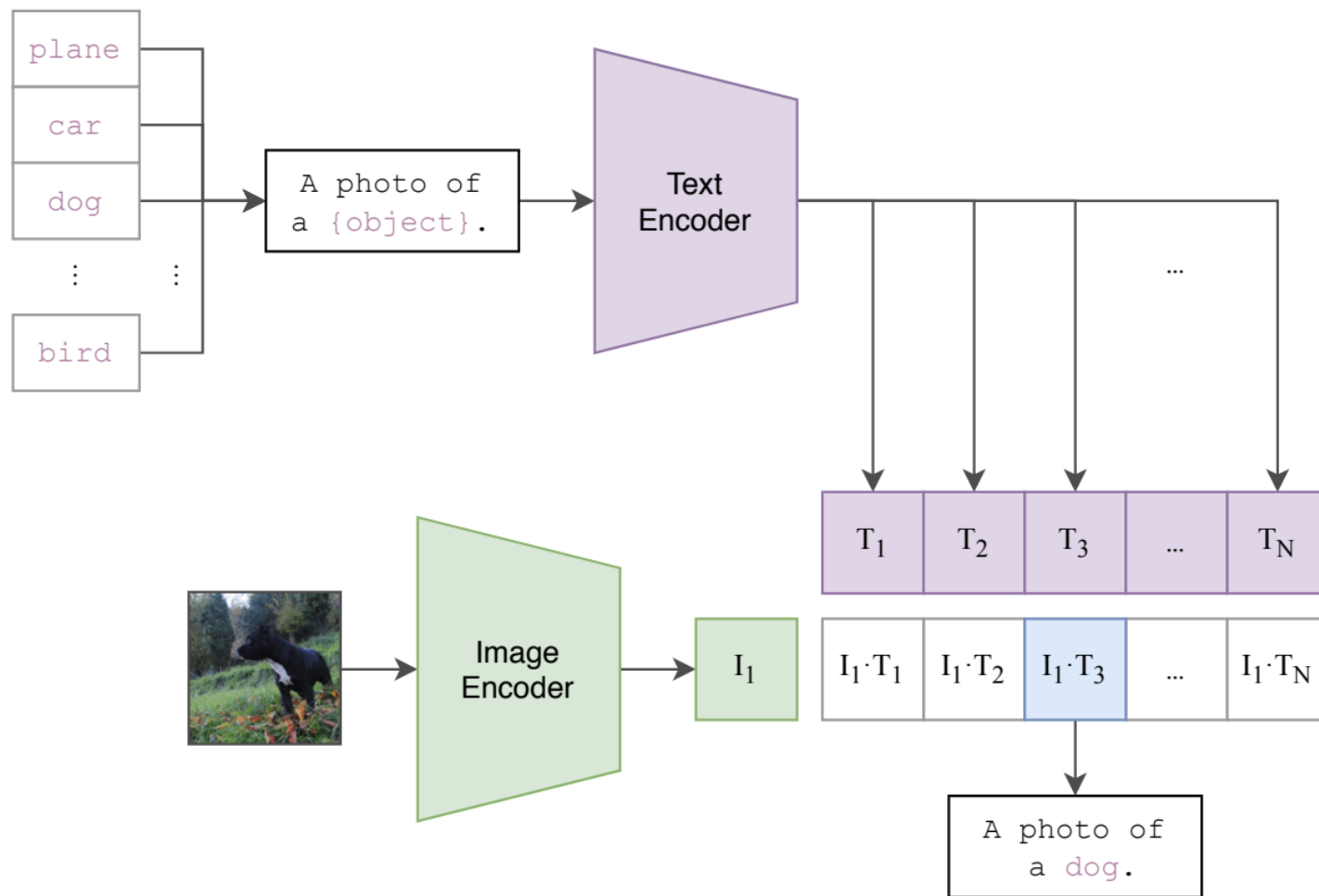# How to train CLIP?

# How to train CLIP?

# Zero shot Learning

# Handson

## Flickr Image dataset

Flickr Image captioning dataset

https://www.kaggle.com/datasets/hsankesara/flickr-image-dataset?resource=download