

# FINE-TUNE MULTIMODAL VISION MODELS

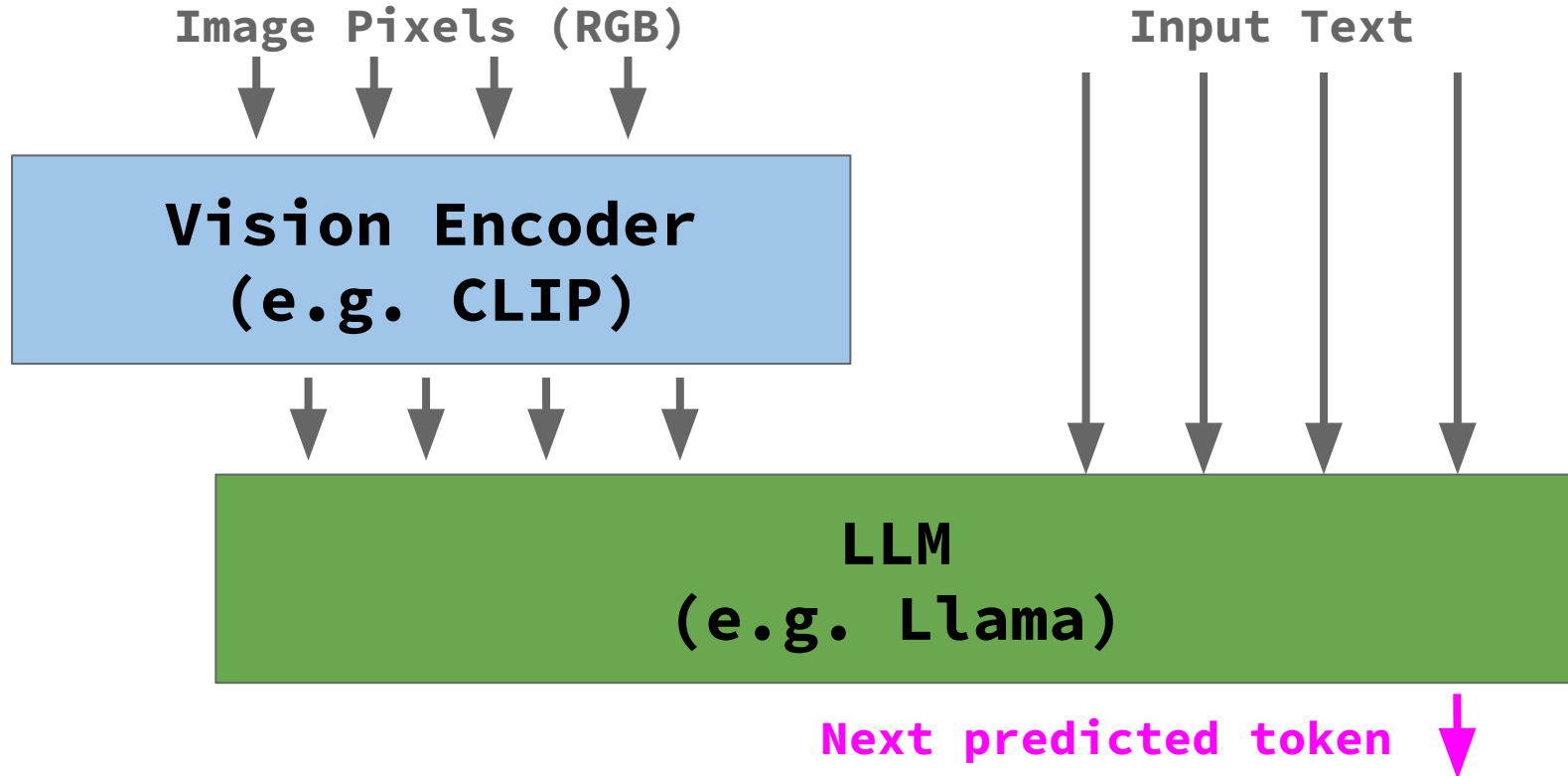
**Trelis Research**

# OVERVIEW

1. DEMO
2. APPLICATIONS.
3. BUILDING A VISION + TEXT MODEL.
4. THREE SPECIFIC ARCHITECTURES - LLAVA 1.5, LLAVA 1.6. IDEFICS.
5. PREPARING A FINE-TUNING DATASET.
6. FINE-TUNING LLAVA 1.6
  - A. MISTRAL 7B VERSION
  - B. YI 34B
7. FINAL THOUGHTS

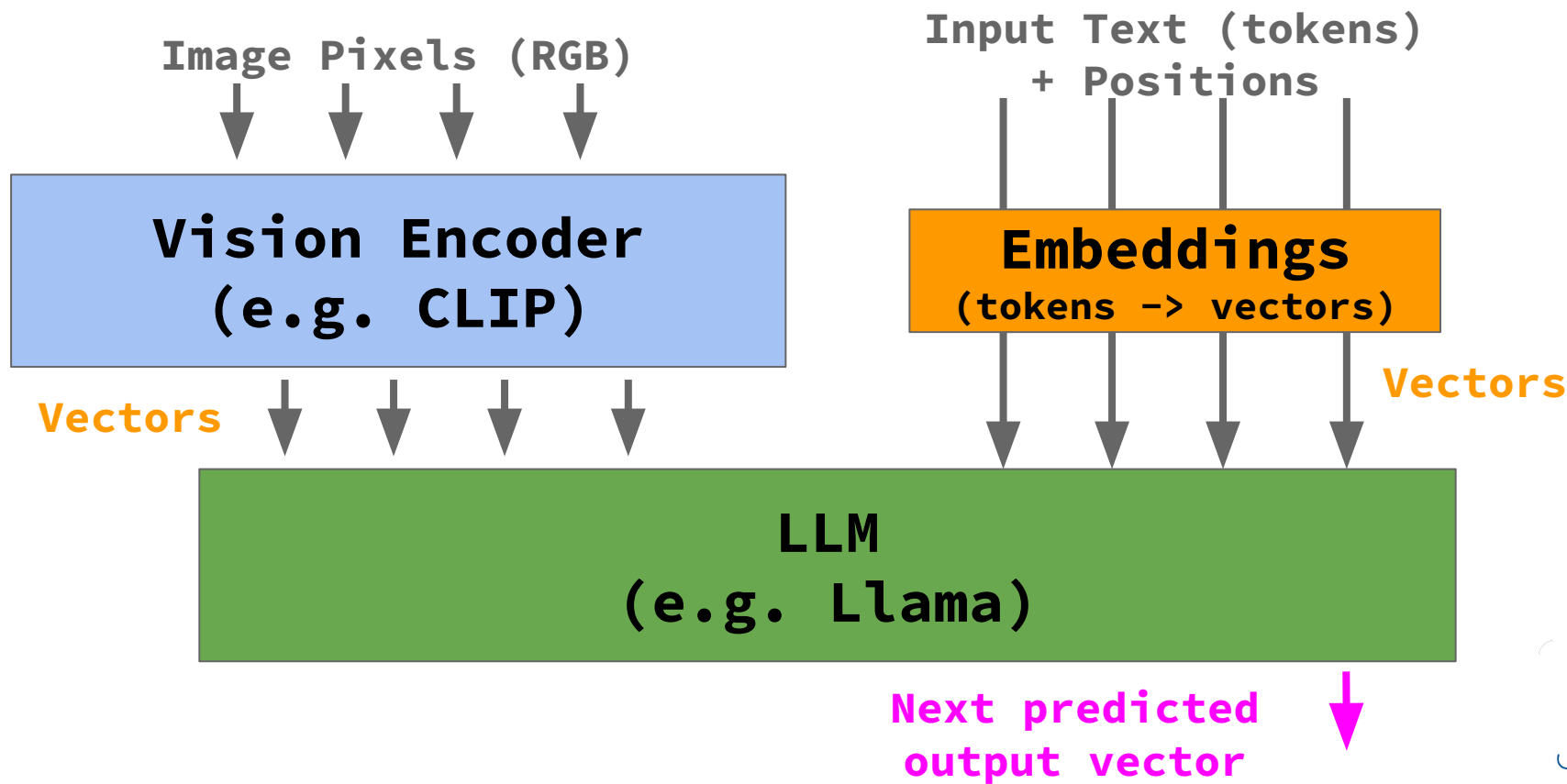


# VISION + TEXT MODELS



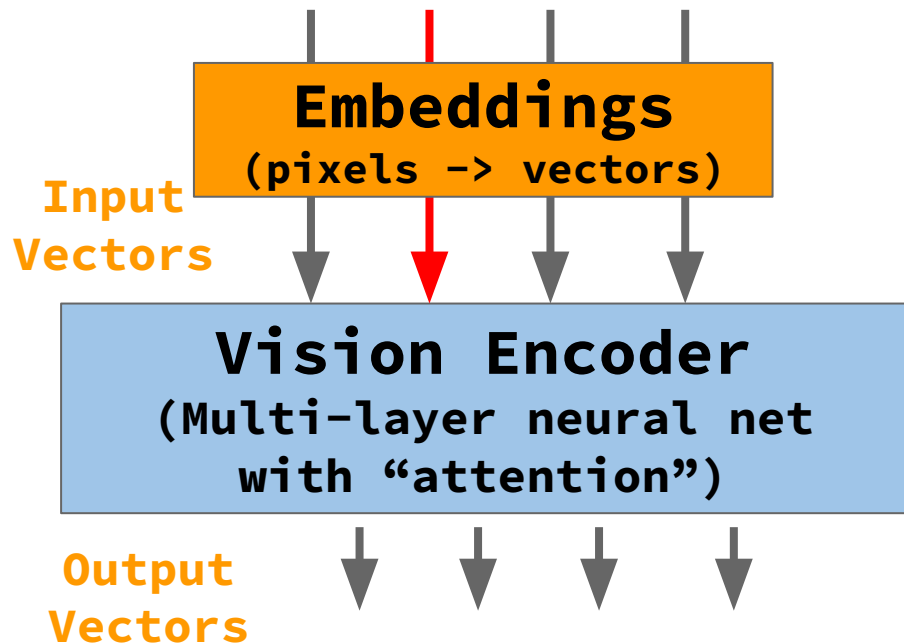
*T*

# VISION + TEXT MODELS



# VISION ENCODER

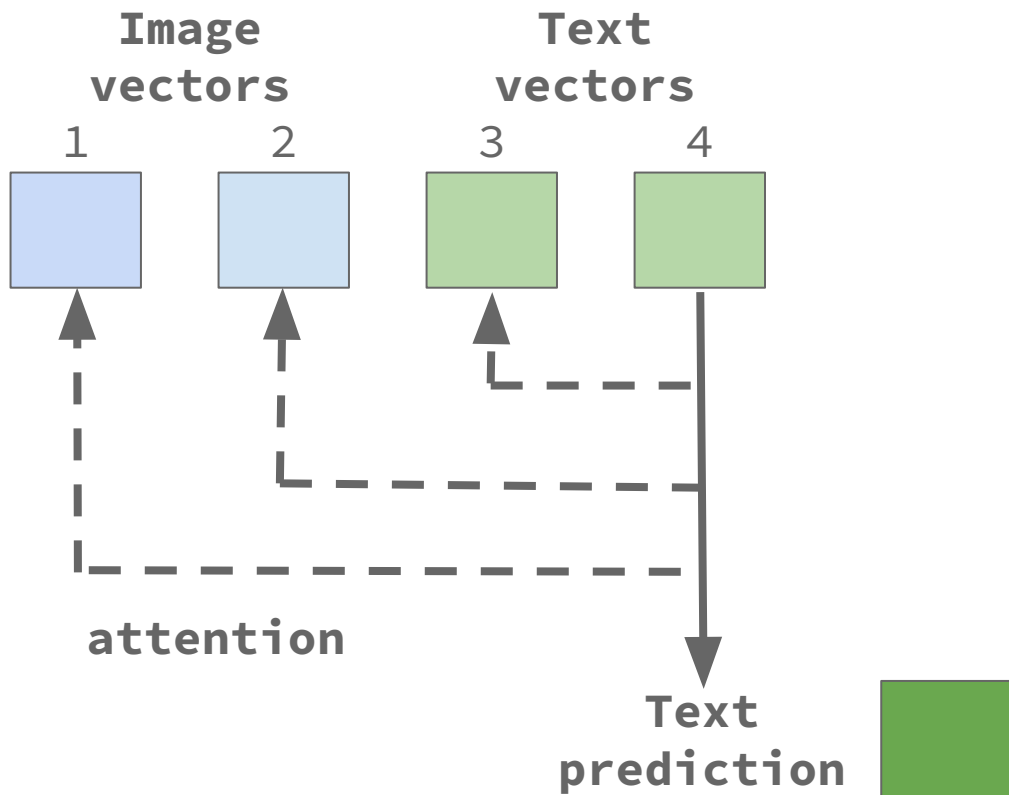
Image Pixels (RGB) + Positions



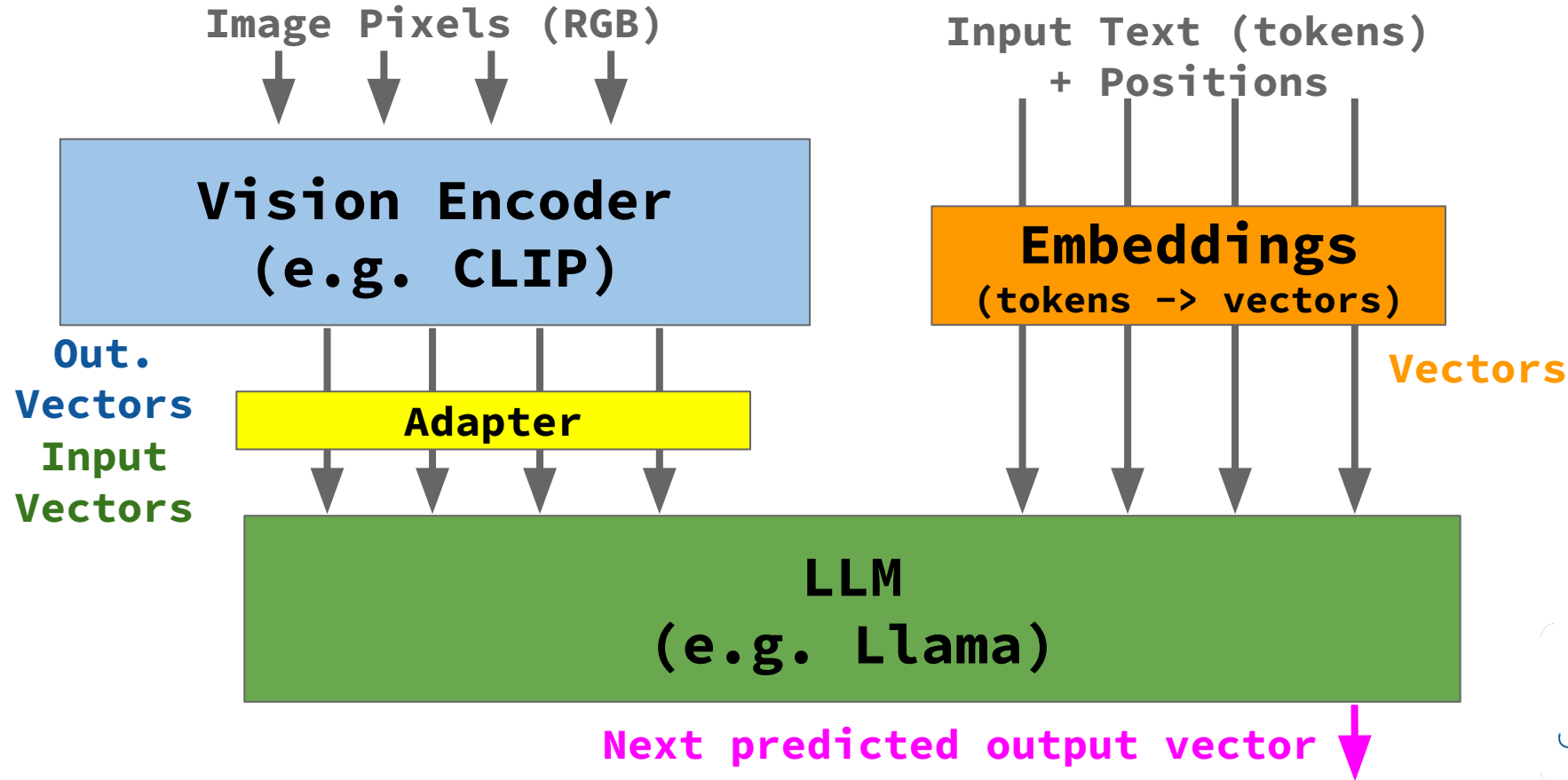
Each tile (16x16px)  
= 1 input vector



# LANGUAGE MODEL DECODER



# ADDING AN ADAPTER!



# LLAVA 1.5 FORMULA

1. CLIP ENCODER + LLAMA 2 LLM
2. FREEZE THE VISION ENCODER + LLM => TRAIN THE ADAPTER
  - A. USE IMAGE + TEXT DATASETS
3. UNFREEZE EVERYTHING => TRAIN EVERYTHING
  - A. USE SYNTHETIC\* IMAGE + TEXT DATA

\*CHATGPT (NON VISION) IS USED TO WRITE DETAILED DESCRIPTIONS OF IMAGES GIVEN DETAILS OF THEIR CONTENTS.

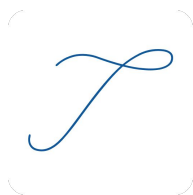




# LLAVA 1.6 FORMULA

1. CLIP ENCODER + MISTRAL 7B OR YI 34B LLM
2. USE AN MLP INSTEAD OF A SIMPLE LINEAR LAYER (E.G. LINEAR LAYER + ACTIVATION LAYER).
3. USE A LARGER VISION MODEL 336x336 NOT 224x224.

\*CHATGPT (NON VISION) IS USED TO WRITE DETAILED DESCRIPTIONS OF IMAGES GIVEN DETAILS OF THEIR CONTENTS.



# IDEFICS FORMULA

1. CLIP ENCODER + LLAMA -> LLM.
2. FLAMINGO ARCHITECTURE: INJECT VISION VECTORS THROUGHOUT THE LLM.
3. TRAINING USING LONG MULTI-IMAGE DOCUMENTS.

FLAMINGO PAPER: [HTTPS://ARXIV.ORG/PDF/2204.14198.PDF](https://arxiv.org/pdf/2204.14198.pdf)

IDEFICS MODEL (80B):

[HTTPS://HUGGINGFACE.CO/HUGGINGFACEM4/IDEFICS-80B-INSTRUCT](https://huggingface.co/HuggingFaceM4/idefics-80b-instruct)

