

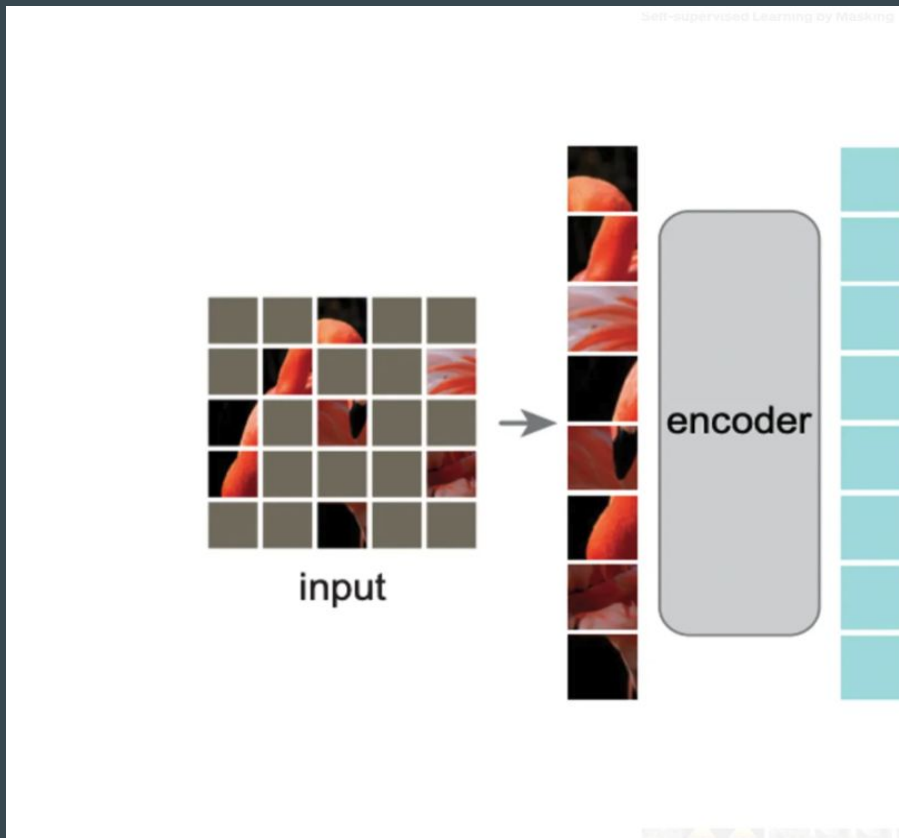
# From Diagrams to Dialogue

## Understanding Miro Boards with VLMs

# What is the Million Dollar Idea \$ :

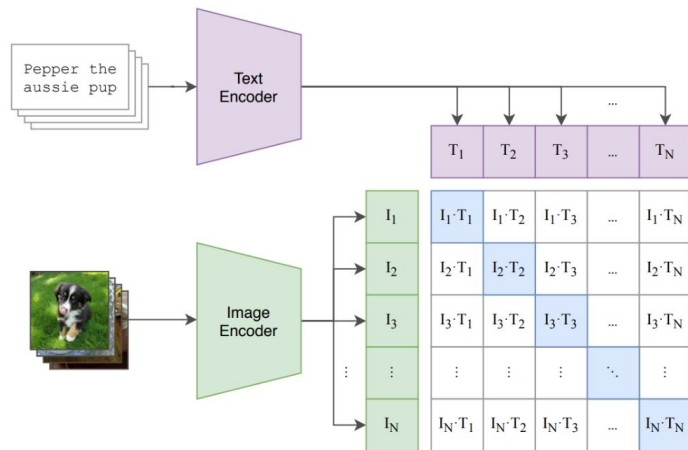
The Rough Idea is to Extract Information Visual Board  
Either it is Miro , Mural or Lucid etc

# Vision Encoder:

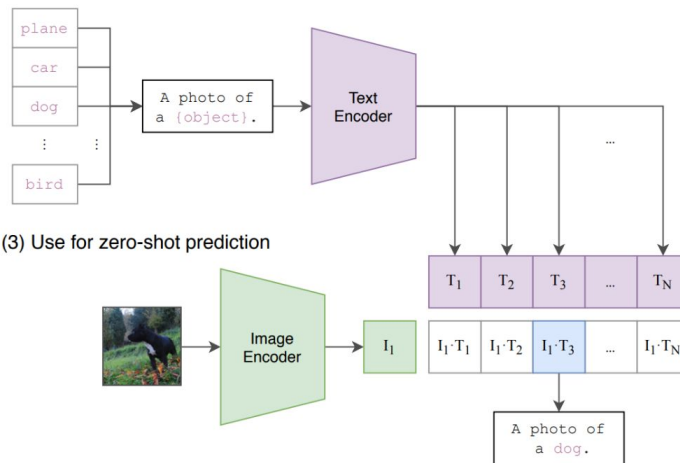


# (High Level VIT )CLIP ,SigLIP etc

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

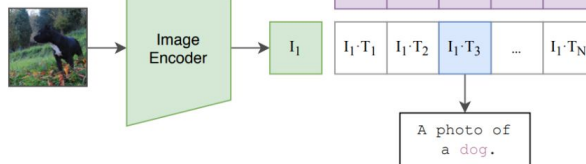
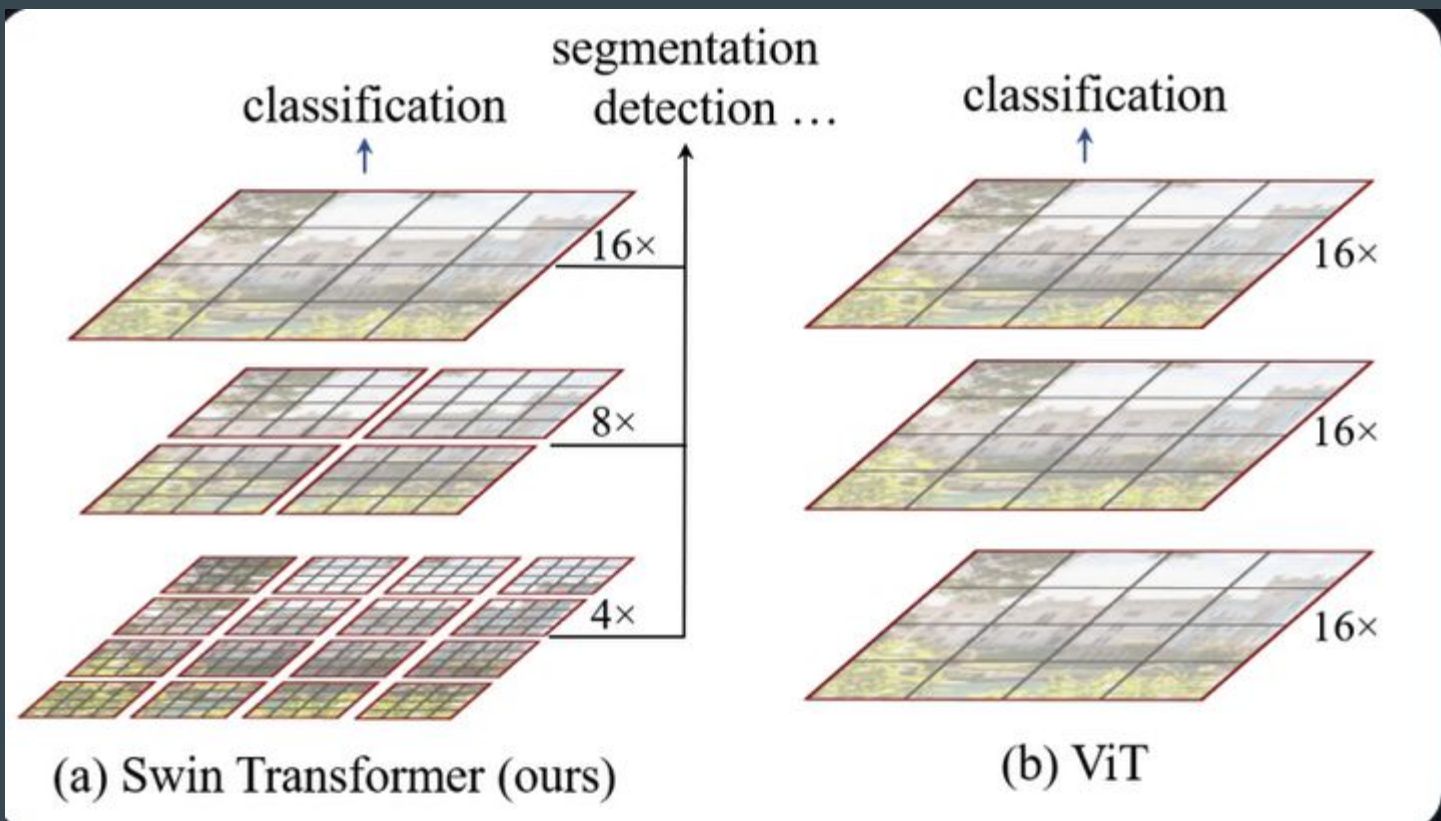


Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

# Swim Transformer:



# Simple Math Behind It

VIT (Vision Transformation)

Patches per row =  $224 \div 16 = 14$

Patches per column =  $224 \div 16 = 14$

Total patches =  $14 \times 14 = 196$  patches

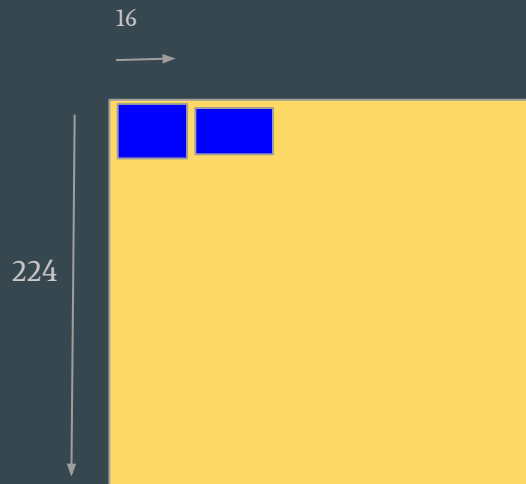
Each patch =  $16 \times 16 \times 3 = 768$  dimensions

All patches =  $196 \text{ patches} \times 768 \text{ dimensions}$

Attention Matrix =  $196 \times 196 = 38,416$   
computations

Every patch attends to ALL other patches

Complexity =  $O(N^2)$  where  $N = 196$



# Simple Math Behind It

Initial patches =  $4 \times 4$  pixels (smaller than ViT)

Patches per row =  $224 \div 4 = 56$

Patches per column =  $224 \div 4 = 56$

Total patches =  $56 \times 56 = 3,136$  patches

Windowing

Window size =  $7 \times 7$  patches

Patches per window =  $7 \times 7 = 49$  patches

Total windows =  $3,136 \div 49 = 64$  windows

Attention per window =  $49 \times 49 = 2,401$  computations

Total attention =  $64 \text{ windows} \times 2,401 = 153,664$  computations

# Simple Math Behind It

Swin Transform Operation	>	ViT Operation
153,664		38,416

But Swin Transform Operation = Complexity =  $O(N)$  - Linear!

For  $448 \times 448$  image (2x larger):

ViT:

- Patches =  $28 \times 28 = 784$
- Attention =  $784^2 = 614,656$  operations

Swin:

- Windows still  $7 \times 7 = 49$  patches per window
- More windows, but attention per window stays constant



# Baby Step:

Figuring out the Layout in the Board

ScreenAI:

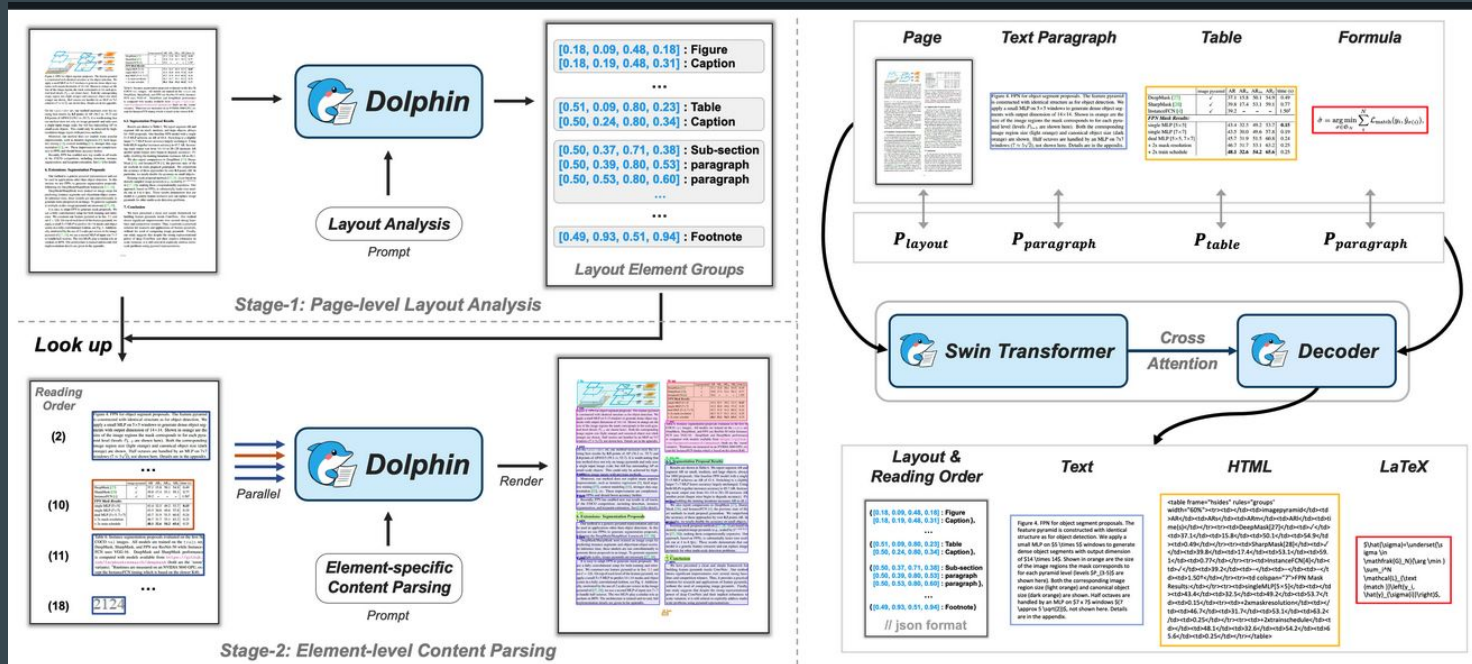
<https://research.google/blog/screenai-a-visual-language-model-for-ui-and-visually-situated-language-understanding/>

Simple Idea:

First Figure out the Layout and different Components in the Board

Pass down the Split images to Multi Model LLM (Preferable InternVL my personal Favorite for Diagram Analysis)

# Dolphin Model (Document Image Parsing via Heterogeneous Anchor Prompting)





Repo



Profile