# Methods of image grounding and LLM-based robotic manipulation

a CogRob project

*Juan Lopez, s5156750*
*Thijs Lukkien, s3978389*

# Topic:

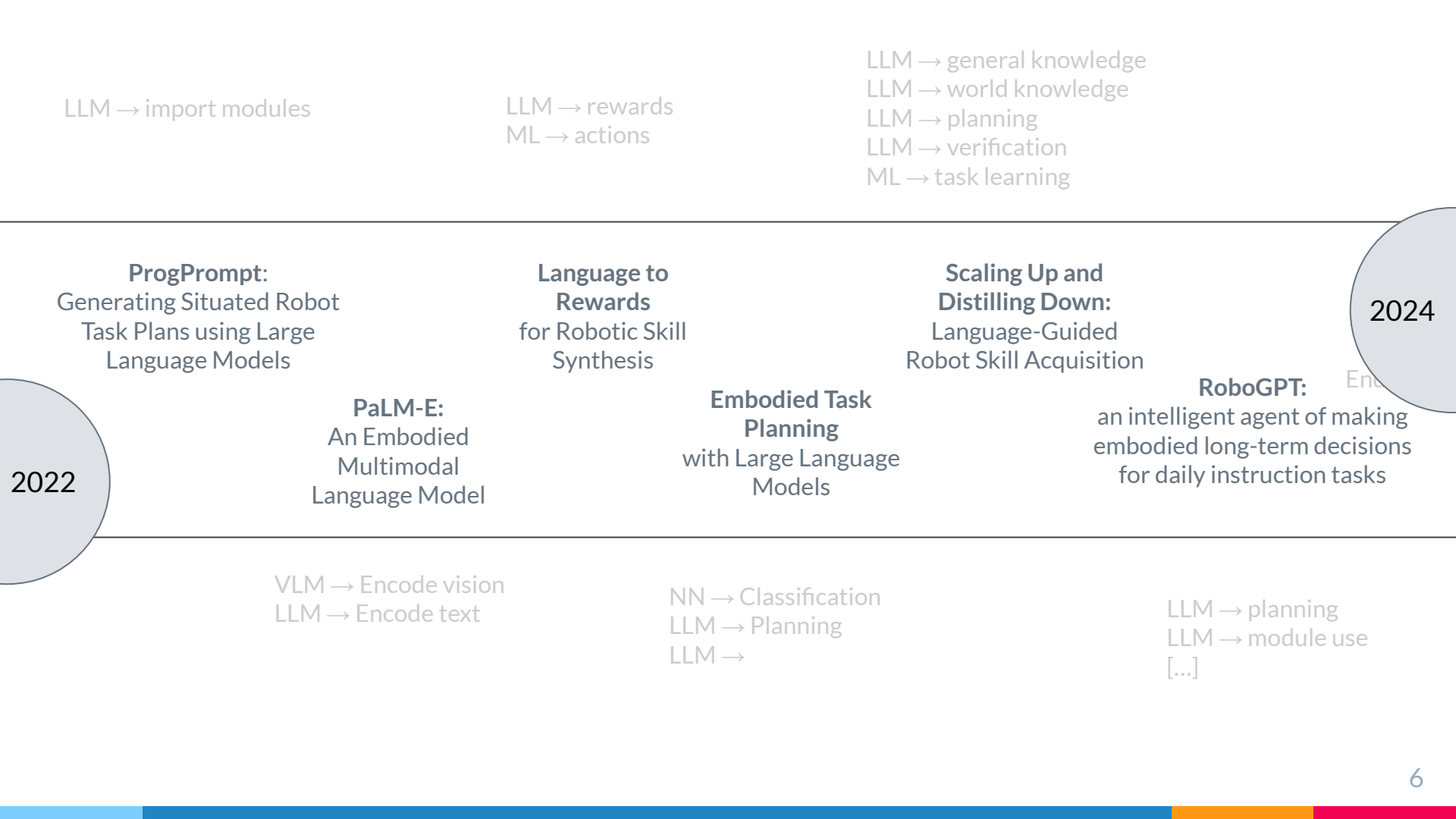LLMs for robot instructions

# 1.
# The idea

" Compare two methods of image grounding

against a ground-truth method, for LLM implementations via natural language instructions »

# 2.

# The background

LLM → import modules

LLM → rewards
ML → actions

LLM → general knowledge
LLM → world knowledge
LLM → planning
LLM → verification
ML → task learning

**ProgPrompt**:
Generating Situated Robot
Task Plans using Large
Language Models

**Language to Rewards**
for Robotic Skill
Synthesis

**Scaling Up and Distilling Down:**
Language-Guided
Robot Skill Acquisition

2024

**PaLM-E:**
An Embodied
Multimodal
Language Model

**Embodied Task Planning**
with Large Language
Models

**RoboGPT:**
an intelligent agent of making
embodied long-term decisions
for daily instruction tasks

En

2022

VLM → Encode vision
LLM → Encode text

NN → Classification
LLM → Planning
LLM →

LLM → planning
LLM → module use
[…]

LLM → import modules

LLM → rewards
ML → actions

LLM → general knowledge
LLM → world knowledge
LLM → planning
LLM → verification
ML → task learning

**ProgPrompt**

**Language to
Rewards**

**Scaling Up and
Distilling Down**

2021

2024

**PaLM-E**

**Embodied Task
Planning**

**RoboGPT**

VLM → Encode vision
LLM → Encode text

NN → Classification
LLM → Planning
LLM →

LLM → planning
LLM → module use
[...]

LLM → import modules

LLM → rewards
ML → actions

LLM → general knowledge
LLM → world knowledge
LLM → planning
LLM → verification
ML → task learning

**ProgPrompt**:
Generating Situated Robot
Task Plans using Large
Language M...

**Language to
Rewards**

**Scaling Up and...**

PaLM...
An Embodied
Multimodal
Language Model

with Large Languag...
Models

...ntelligent agent...
...odied long-term decisions
...r daily instruction tasks

VLM → Encode vision
LLM → Encode text

NN → Classification
LLM → Planning
LLM →

LLM → planning
LLM → module use

8

# Model Zoo

*BLOOM* - LLM → text

*GR-CONVNET*- CNN → grasps

*SAM* - MMTM → Segmentation maps

*CLIP* - MMLM → text

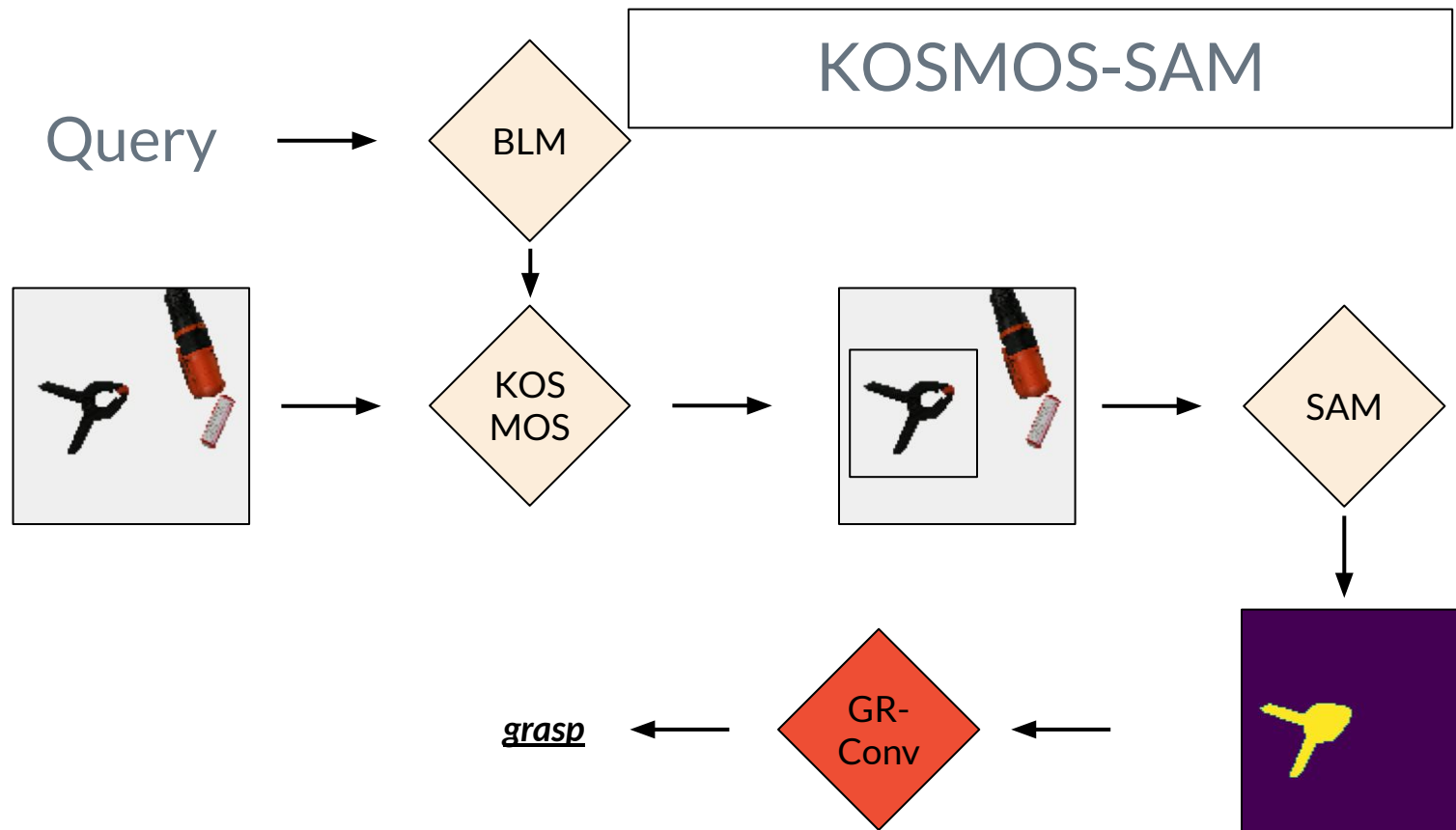*KOSMOS2* - MMLM → grounded text + referred bbox images

# Additional Considerations

▷ Only pretrained models, no fine-tuning

▷ Compare their few-shots/zero-shots capabilities.

▷ Grasp success depends on

○ Segmentation quality
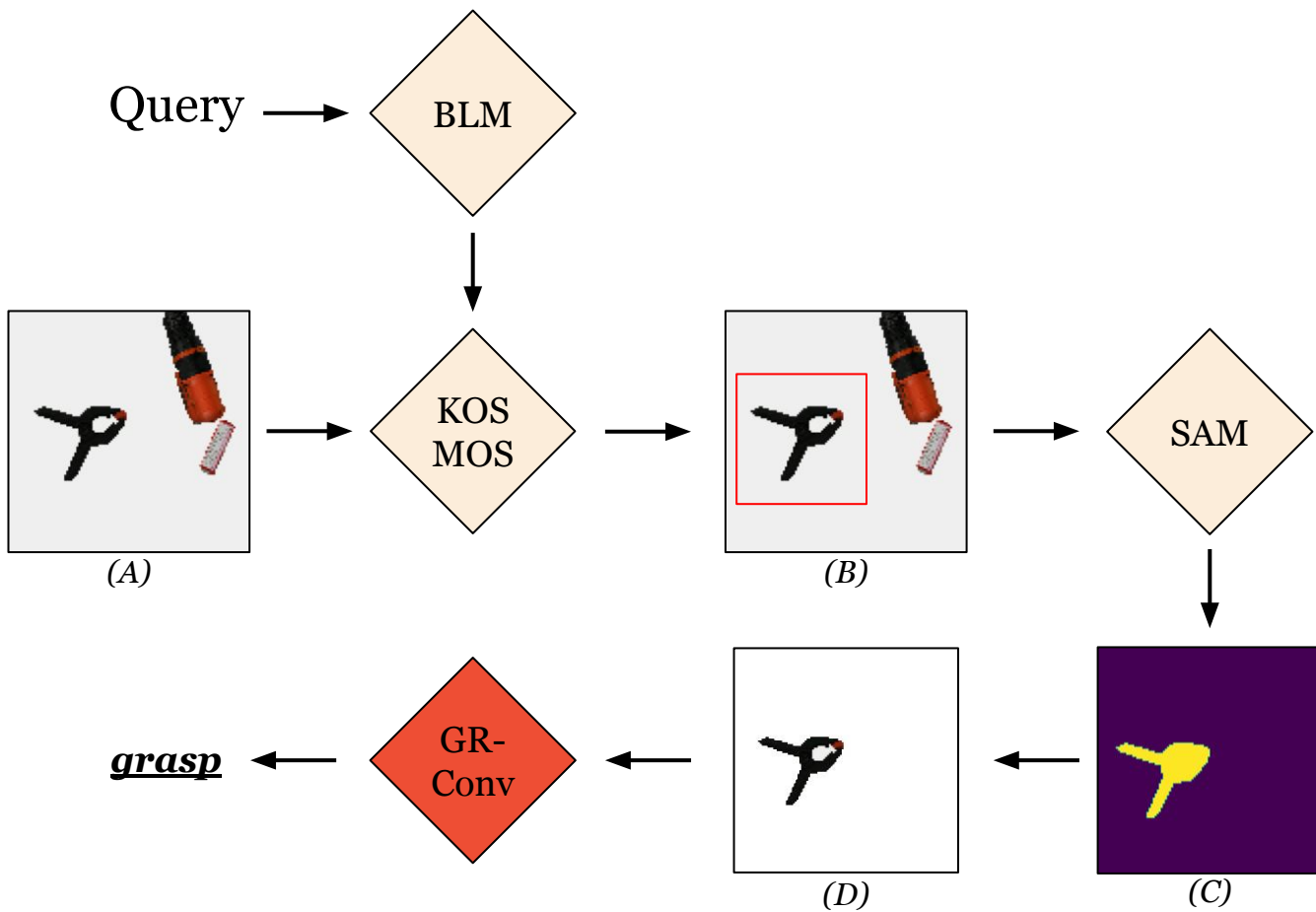
○ Classification performance

○ Grasp pose estimator

# Default

(A)

SAM

(B)

(C)

CLIP

$\left\{ \begin{array}{c} \text{Cat1} \\ \text{Cat2} \\ \text{Catx} \end{array} \right\}_{(D)}$

GR-Conv

| ID cat | ID cat | ID cat |
|--------|--------|--------|
| grasp | grasp | grasp |

(E)

BLM

$\left\{ \begin{array}{c} \text{Tmpl} \\ \text{Tmpl} \\ \textbf{QRY} \end{array} \right\}_{(F)}$

***grasp***

(A)

(B)

$\left\{\begin{array}{l} \text{Cat1} \\ \text{Cat2} \\ \text{Catx} \end{array}\right\}_{(C)}$

CLIP

GR-Conv

ID cat | ID cat | ID cat

grasp grasp grasp

(D)

BLM

grasp

$\left\{\begin{array}{l} \text{Tmpl} \\ \text{Tmpl} \\ \textbf{QRY} \end{array}\right\}_{(E)}$

Query → BLM

(A) → KOS MOS → (B) → SAM

SAM → (C)

(C) → (D) → GR-Conv → **_grasp_**

# 3.
# Evaluation

# Experimental Setup

3 experimental methods:

- ▷ Default (Ground Truth Segmentation and CLIP)
- ▷ SAM (automatic) - CLIP
- ▷ KOSMOS - SAM (bounding box)


- ▷ New Collab adaptation

# Experimental Setup

▷  4 objects per run, 20 runs

▷  1 dynamically generated query per object

▷  80 queries per pipeline

# What we measure

▷ Overall success of the query completion

▷ Successful grasp

▷ Successful delivery

▷ Successful Object Detection

# 4.
# Results

# Preliminary results

Kosmos-SAM  >  SAM-CLIP

Kosmos-SAM  >  ground-CLIP

CLIP-classification

# ground-CLIP

# 5.
# Conclusion

# Conclusion

Modality / model-applicability

# 6.

# Discussion

# Discussion

▷ CLIP and BLOOM performance (poor)

▷ Manual tuning

▷ Cascading errors : pinpointing difficulties
  ○ Where does low performance come from in the pipeline? Was it at the start or at the end?

▷ Potential for simplification

▷ NextChat - LLM that also creates the segmentation mask text grounding (Mixture of clip and SAM in a single model)

"

"Questions?"

*~ Dhali ~*

# Extra material

Mask 1, Score: 0.869

Mask 2, Score: 0.999

Mask 3, Score: 0.995

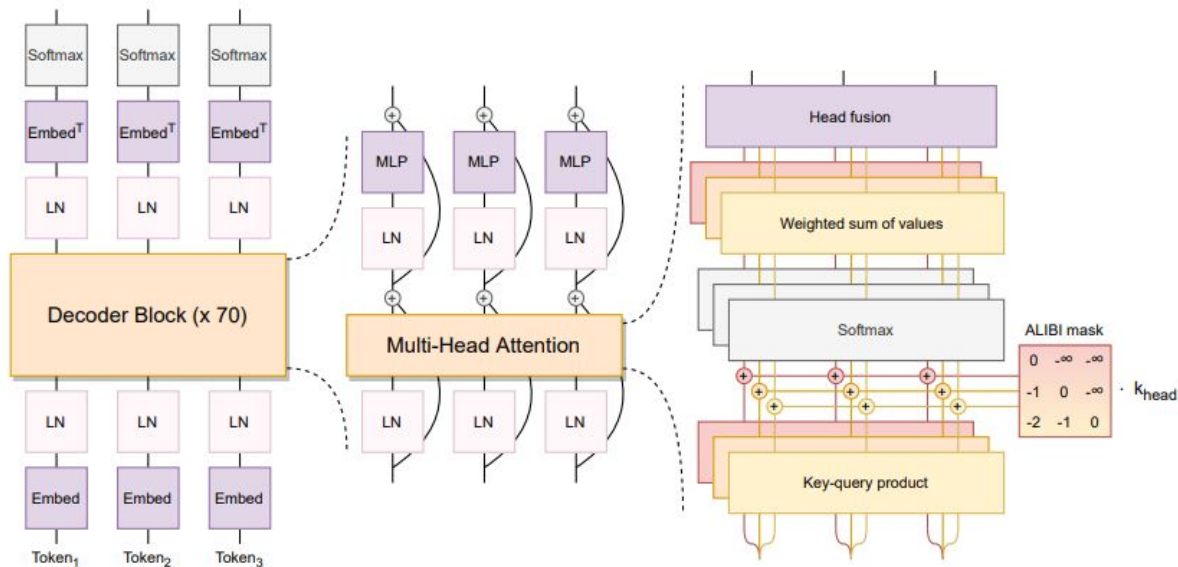Mask 1, Score: 0.941 | Mask 2, Score: 0.977 | Mask 3, Score: 0.983

```
# put the scissors to the top left corner.
bbox = robot.ground_object("scissors")
mask = robot.segmentSAM(bbox)
robot.pick_and_place_2(mask, "top left corner")
# put the gelatin box to the top side.
bbox = robot.ground_object("gelatin box")
mask = robot.segmentSAM(bbox)
robot.pick_and_place_2(mask, "top side")
```
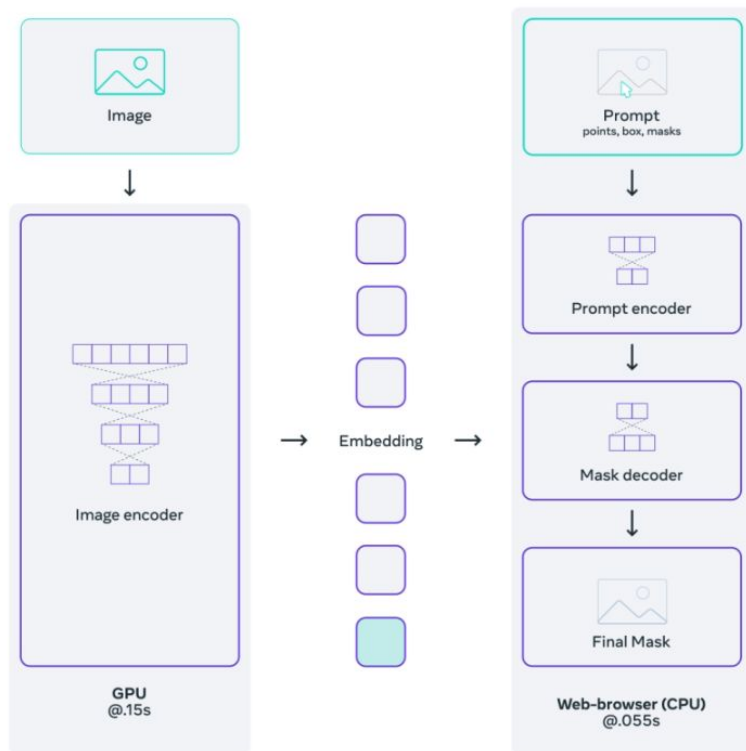
# GR-CONVNET

# BLOOM
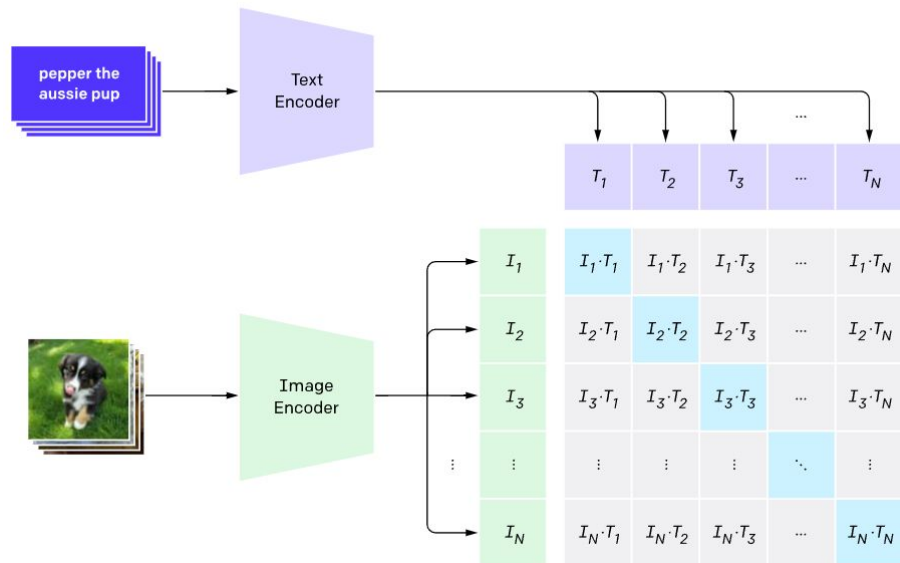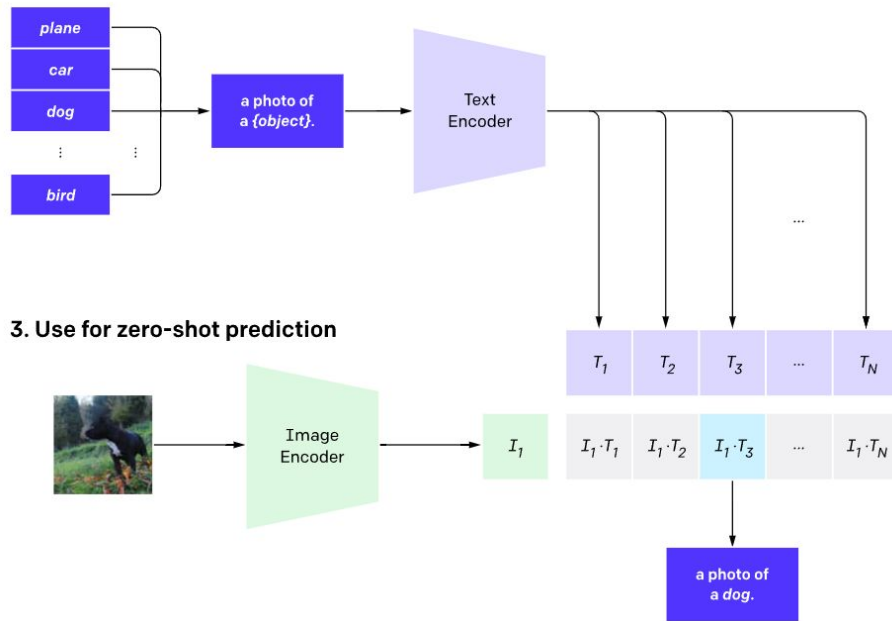
## Classic LLM similar to GPT-3

# SAM

# CLIP



**1. Contrastive pre-training**

**2. Create dataset classifier from label text**

**3. Use for zero-shot prediction**

# Kosmos-2



(1) Grounded question answering

https://colab.research.google.com/drive/1orQdJ5qkkW44uIlwPrTaKl_BfbwoiydY?usp=sharing

Once you get access to the repo