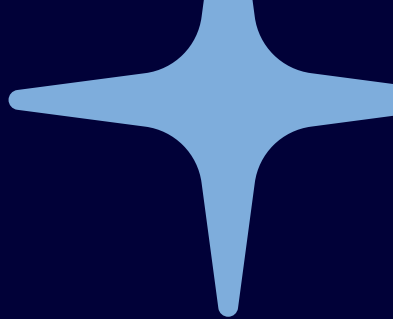


The power of VLMs

Advanced LLM Night

Eivind Kjosbakken – Findable

January 2025



ME

Master's in Engineering and ICT from NTNU

- Main focus on AI
- One year exchange at TU Delft



ME

Master's in Engineering and ICT from NTNU

- Main focus on AI
- One year exchange at TU Delft

Have been a part of Findable since 2021

- Part-time AI-trainer
- Wrote my master's thesis within ML at Findable
- Full-time data scientist



OVERVIEW



OVERVIEW

The power of VLMs



OVERVIEW

The power of VLMs

How vision is integrated in VLMs



OVERVIEW

The power of VLMs

How vision is integrated in VLMs

VLM usecases



OVERVIEW

The power of VLMs

How vision is integrated in VLMs

VLM usecases

- Information extraction



OVERVIEW

The power of VLMs

How vision is integrated in VLMs

VLM usecases

- Information extraction
- Classification



OVERVIEW

The power of VLMs

How vision is integrated in VLMs

VLM usecases

- Information extraction
- Classification

VLM Drawbacks



OVERVIEW

The power of VLMs

How vision is integrated in VLMs

VLM usecases

- Information extraction
- Classification

VLM Drawbacks

Staying up to date with VLM advances



LICENSE PLATE RECOGNITION



+



Provide the full license plate number



LICENSE PLATE RECOGNITION



+



Provide the full license plate number

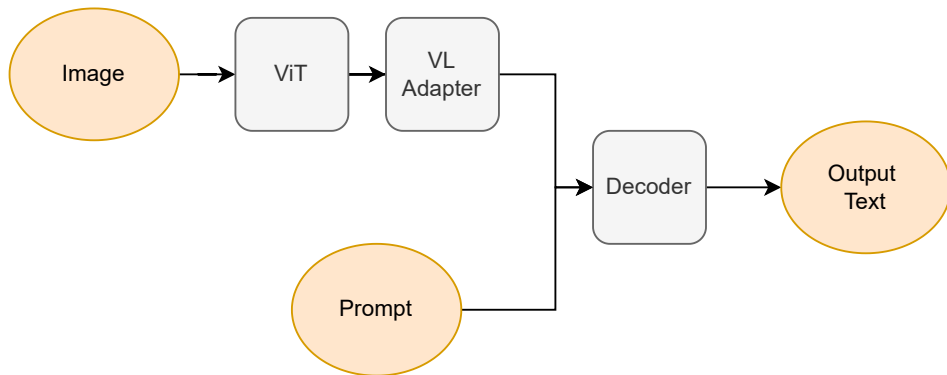


The full license plate number is AB12345.

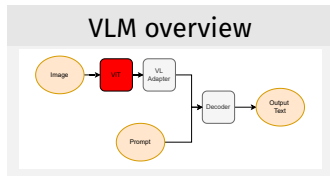


INTRO TO VLMs

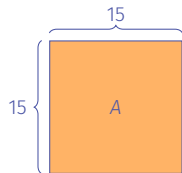
How do they work



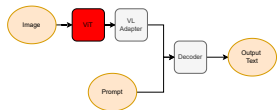
HOW VISION INTEGRATES IN VLMs - ViT



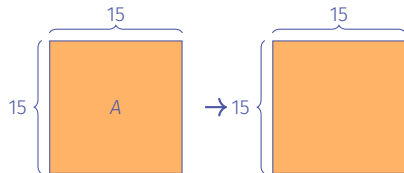
HOW VISION INTEGRATES IN VLMs - ViT



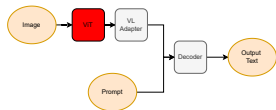
VLM overview



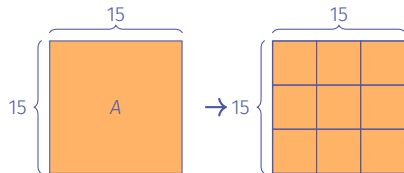
HOW VISION INTEGRATES IN VLMs - ViT



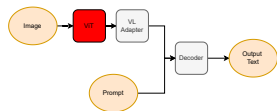
VLM overview



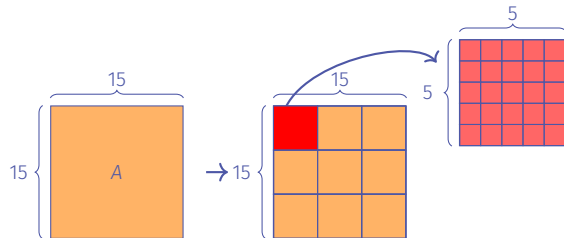
HOW VISION INTEGRATES IN VLMs - ViT



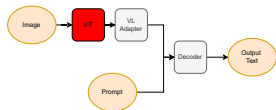
VLM overview



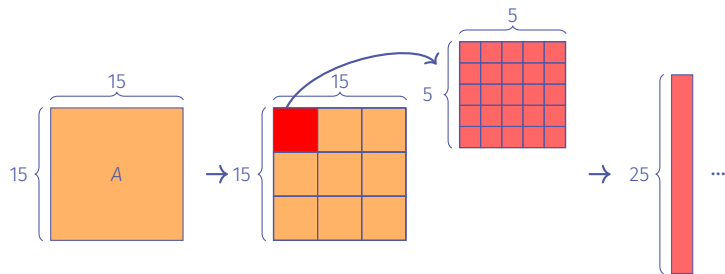
HOW VISION INTEGRATES IN VLMs - ViT



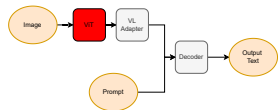
VLM overview



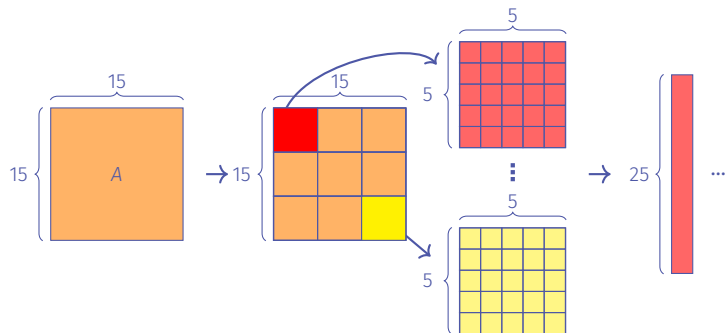
HOW VISION INTEGRATES IN VLMs - ViT



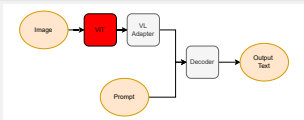
VLM overview



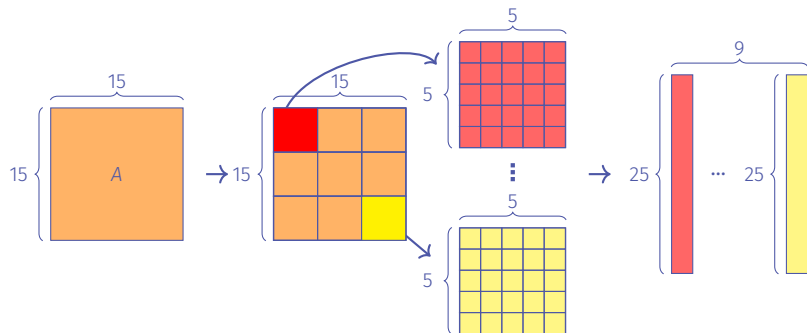
HOW VISION INTEGRATES IN VLMs - ViT



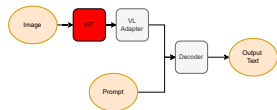
VLM overview



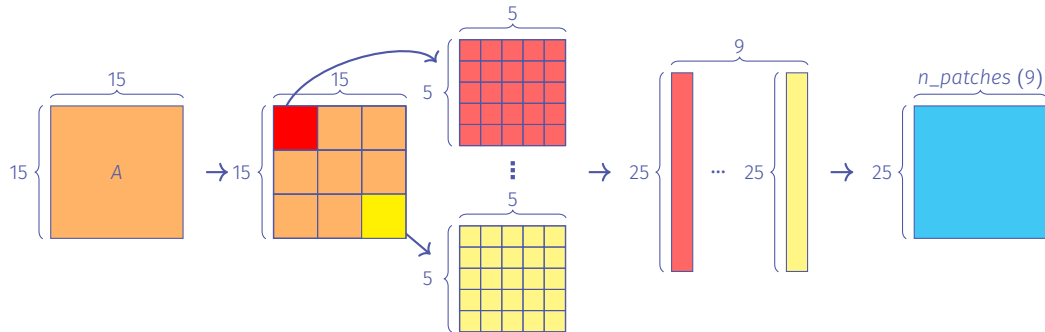
HOW VISION INTEGRATES IN VLMs - ViT



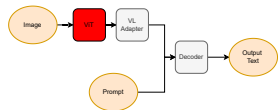
VLM overview



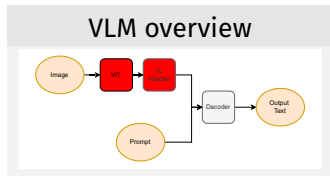
HOW VISION INTEGRATES IN VLMs - ViT



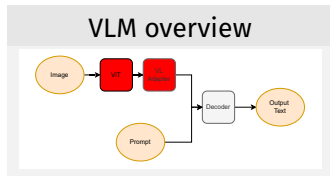
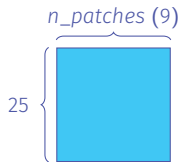
VLM overview



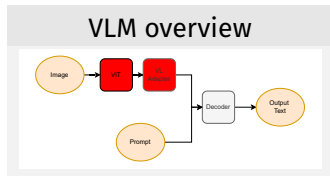
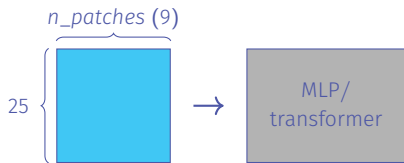
HOW VISION INTEGRATES INTO THE LLM - TRANSFORMER LAYERS



HOW VISION INTEGRATES INTO THE LLM - TRANSFORMER LAYERS



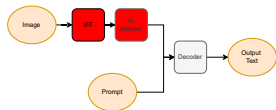
HOW VISION INTEGRATES INTO THE LLM - TRANSFORMER LAYERS



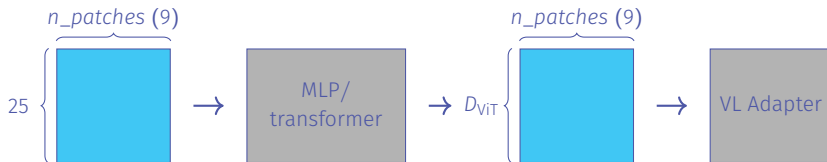
HOW VISION INTEGRATES INTO THE LLM - TRANSFORMER LAYERS



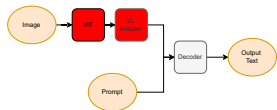
VLM overview



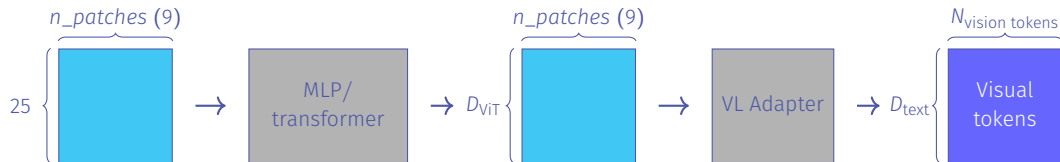
HOW VISION INTEGRATES INTO THE LLM - TRANSFORMER LAYERS



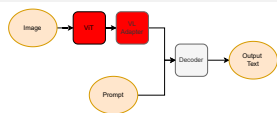
VLM overview



HOW VISION INTEGRATES INTO THE LLM - TRANSFORMER LAYERS



VLM overview



DECODER INPUT SEQUENCE

Prompt

Describe this image



DECODER INPUT SEQUENCE

Prompt

Describe this image



Describe



DECODER INPUT SEQUENCE

Prompt

Describe this image



Describe



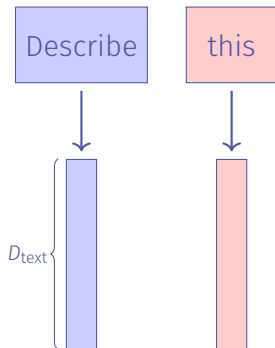
D_{text}



DECODER INPUT SEQUENCE

Prompt

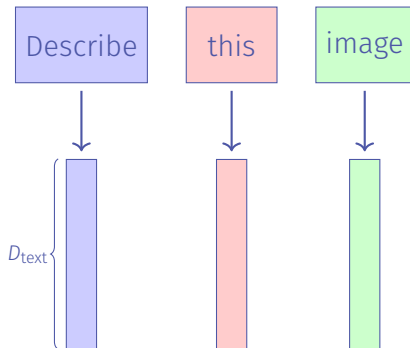
Describe this image



DECODER INPUT SEQUENCE

Prompt

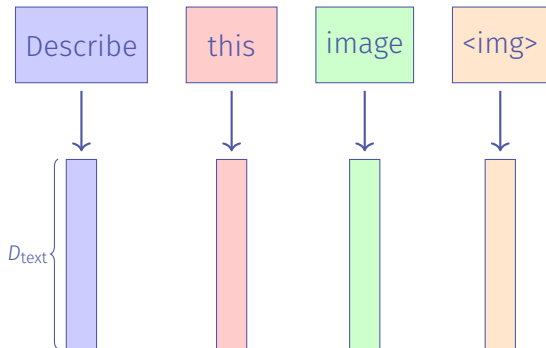
Describe this image



DECODER INPUT SEQUENCE

Prompt

Describe this image

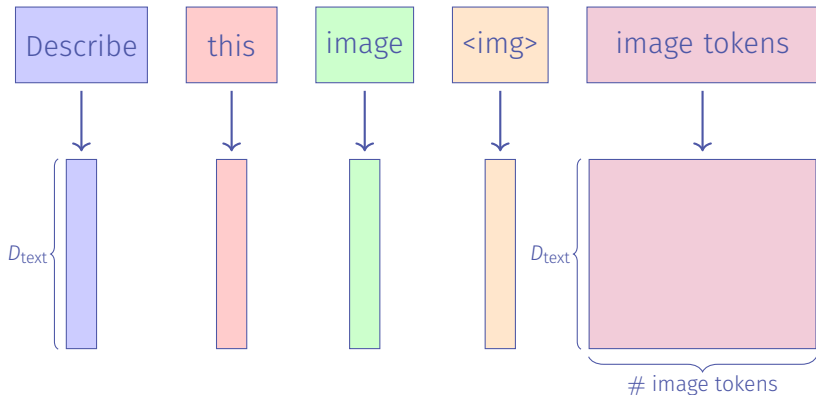


DECODER INPUT SEQUENCE

Prompt

Describe this image

3

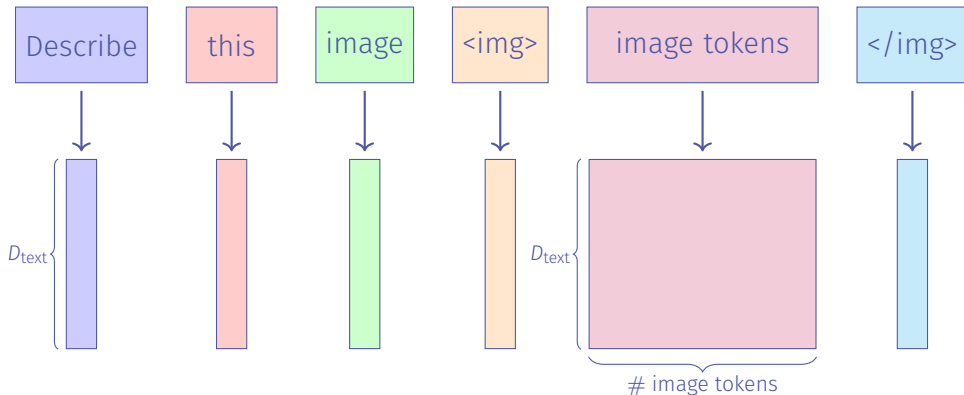


DECODER INPUT SEQUENCE

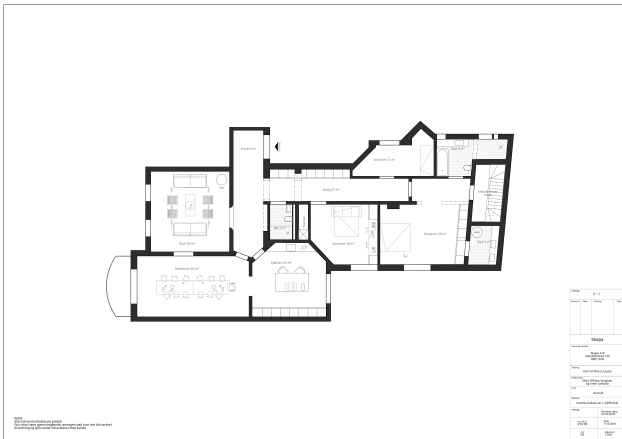
Prompt

Describe this image

3



ARK DRAWING EXAMPLE



VLM INFORMATION EXTRACTION



VLM INFORMATION EXTRACTION

Input

- Single or multiple images (no text necessary)
- Prompt



VLM INFORMATION EXTRACTION

Input

- Single or multiple images (no text necessary)
- Prompt

Output is a series of filled out fields

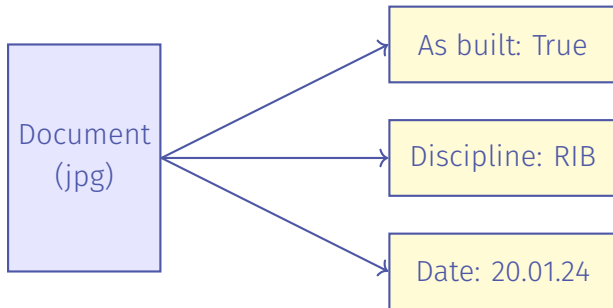


VLM INFORMATION EXTRACTION

Input

- Single or multiple images (no text necessary)
- Prompt

Output is a series of filled out fields

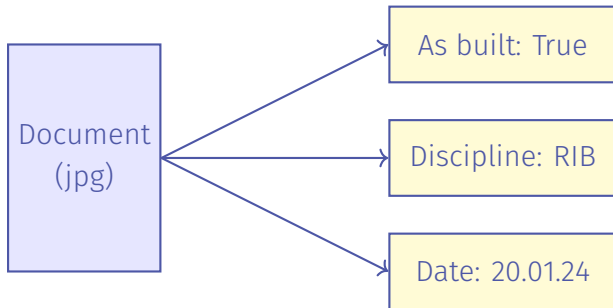


VLM INFORMATION EXTRACTION

Input

- Single or multiple images (no text necessary)
- Prompt

Output is a series of filled out fields



Textual task
requiring visual info

<input type="checkbox"/>	ARK
<input checked="" type="checkbox"/>	RIB
<input type="checkbox"/>	RIE
<input type="checkbox"/>	RIV



VLM INFORMATION EXTRACTION

Example Information Extraction prompt

```
prompt = """
Given the document,
fill out the following JSON object.
Enter None if the field is not available
{
    "As_built" : "",
    "Discipline" : "",
    "Drawing_Date": ""
}
"""
```



VLM CLASSIFICATION



VLM CLASSIFICATION

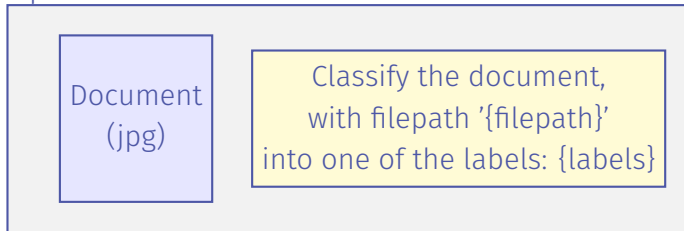
SFT



VLM CLASSIFICATION

SFT

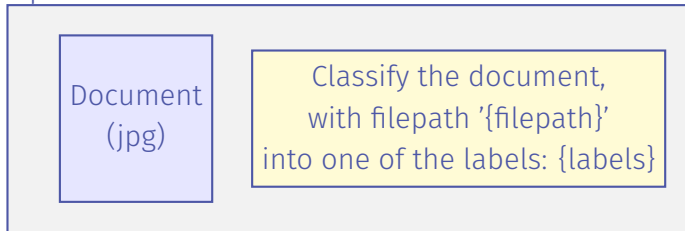
Input



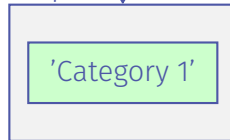
VLM CLASSIFICATION

SFT

Input



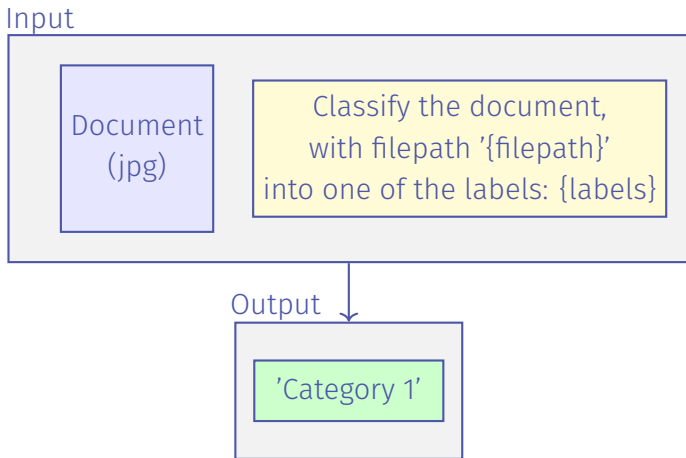
Output



VLM CLASSIFICATION

SFT

Open-source VLMs



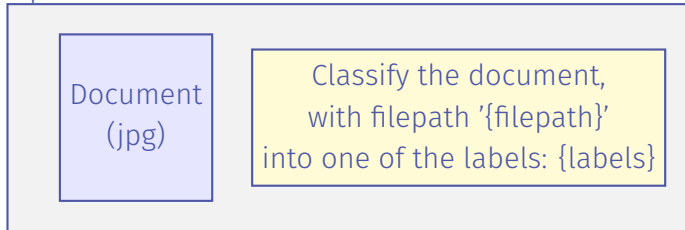
VLM CLASSIFICATION

SFT

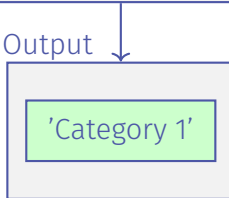
Open-source VLMs

Improved performance

Input



Output



VLM CLASSIFICATION

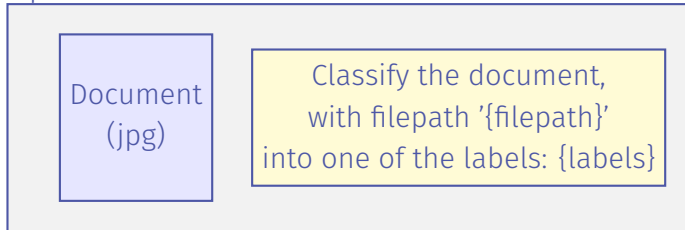
SFT

Open-source VLMs

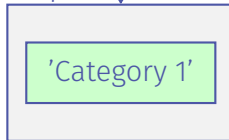
Improved performance

How to avoid hallucinations

Input



Output



VLM DRAWBACKS



VLM DRAWBACKS

Time consuming to set up open-source VLMs



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF
- Often niche package criteria



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF
- Often niche package criteria
- Configure dataset into correct format



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF
- Often niche package criteria
- Configure dataset into correct format

This makes it costly to test different open-source VLMs



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF
- Often niche package criteria
- Configure dataset into correct format

This makes it costly to test different open-source VLMs

Compute heavy



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF
- Often niche package criteria
- Configure dataset into correct format

This makes it costly to test different open-source VLMs

Compute heavy

- High experimentation time for training



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF
- Often niche package criteria
- Configure dataset into correct format

This makes it costly to test different open-source VLMs

Compute heavy

- High experimentation time for training
- Expensive inference



VLM DRAWBACKS

Time consuming to set up open-source VLMs

- Model should be available on platforms like HF
- Often niche package criteria
- Configure dataset into correct format

This makes it costly to test different open-source VLMs

Compute heavy

- High experimentation time for training
- Expensive inference
- Limited number of pages due to compute requirements



STAYING UP TO DATE



STAYING UP TO DATE

Newsletters

- TLDR AI
- AlphaSignal



STAYING UP TO DATE

Newsletters

- TLDR AI
- AlphaSignal

GitHub stars



STAYING UP TO DATE

Newsletters

- TLDR AI
- AlphaSignal

GitHub stars

PapersWithCode



STAYING UP TO DATE

Newsletters

- TLDR AI
- AlphaSignal

GitHub stars

PapersWithCode

Specific blogs

- OpenAI
- Anthropic



STAYING UP TO DATE

Newsletters

- TLDR AI
- AlphaSignal

GitHub stars

PapersWithCode

Specific blogs

- OpenAI
- Anthropic

VLM leaderboards (HF)

Evaluation Dimension															
<input checked="" type="checkbox"/> Avg Score	<input checked="" type="checkbox"/> Avg Rank	<input checked="" type="checkbox"/> MMBench_V11	<input checked="" type="checkbox"/> MMStar	<input type="checkbox"/> MME	<input checked="" type="checkbox"/> MMMU_VAL	<input checked="" type="checkbox"/> MathVista	<input checked="" type="checkbox"/> OCRBench	<input checked="" type="checkbox"/> AI2D	<input checked="" type="checkbox"/> HallusionBench	<input type="checkbox"/> SEEDBench_IMG	<input checked="" type="checkbox"/> MMVet	<input type="checkbox"/> LLaVABench	<input type="checkbox"/> CCBench	<input type="checkbox"/> RealWorldQA	<input type="checkbox"/> POPE
<input type="checkbox"/> ScienceQA_TEST	<input type="checkbox"/> SEEDBench2_Plus	<input type="checkbox"/> MMT-Bench_VAL	<input type="checkbox"/> BLINK												
Model Name								Model Size				Model Type			
Input the Model Name (fuzzy, case insensitive)								<input checked="" type="checkbox"/> <4B	<input checked="" type="checkbox"/> 4B-10B	<input checked="" type="checkbox"/> 10B-20B	<input checked="" type="checkbox"/> 20B-40B	<input checked="" type="checkbox"/> >40B	<input checked="" type="checkbox"/> API	<input checked="" type="checkbox"/> OpenSource	
								<input checked="" type="checkbox"/> Unknown							
Rank	Method	Param (B)	Language Model	Vision Model	Eval Date	Avg Score	Avg Rank	MMBench_V11	MMStar	MMMU_VAL	MathVista	OCRBench	AI2D	HallusionBench	MMVet
1	SenseNova				2024/12/12	77.4	4.25	85.7	72.7	69.6	78.4	894	87.8	57.4	78.2
2	InternVL2.5-78B-NPO	78	Qwen-2.5-72B	InternViT-6B-v2.5	2024/12/29	77	3.88	87.7	72.1	68.2	76.6	909	89.2	58.1	73.5
3	TeleMM				2024/12/31	75.9	8.88	79.9	70.8	66.6	75.7	891	88.5	60.6	75.7
4	InternVL2.5-38B-NPO	38	Qwen-2.5-32B	InternViT-6B-v2.5	2024/12/28	75.3	6.88	85.4	70.1	63.8	73.6	894	87.9	59.7	72.6
5	InternVL2.5-78B	78	Qwen-2.5-72B	InternViT-6B-v2.5	2024/12/10	75.2	7	87.5	69.5	70	70.6	853	89.1	57.4	71.8
6	Qwen2-VL-72B	73.4	Qwen2-72B	QwenViT	2024/10/28	74.8	8	85.9	68.6	64.3	69.7	888	88.3	58.7	73.9
7	Qwen-VL-Max-8889	72	Qwen2-72B	QwenViT	2024/09/12	74.4	9	85.8	69.2	64.6	68.3	881	88.1	59.2	72.3
8	InternVL2.5-38B	38	Qwen-2.5-32B	InternViT-6B-v2.5	2024/12/10	73.5	12.5	85.4	68.5	64.6	72.4	841	87.6	57.9	67.2



STAYING UP TO DATE

Newsletters

- TLDR AI
- AlphaSignal

GitHub stars

PapersWithCode

Specific blogs

- OpenAI
- Anthropic

VLM leaderboards (HF)

LLM Night

Evaluation Dimension

☒ Avg Score

☒ Avg Rank

☒ MMBench_V11

☒ MMStar

☐ MME

☒ MMMU_VAL

☒ MathVista

☒ OCRBench

☒ AI2D

☒ HallusionBench

☐ SEEDBench_IMG

☒ MMVet

☐ LLaVABench

☐ CCBench

☐ RealWorldQA

☐ POPE

☐ ScienceQA_TEST

☐ SEEDBench2_Plus

☐ MMT-Bench_VAL

☐ BLINK

Model Name

Model Size

Model Type

Input the Model Name (fuzzy, case insensitive)

☒ <4B

☒ 4B-10B

☒ 10B-20B

☒ 20B-40B

☒ >40B

☒ API

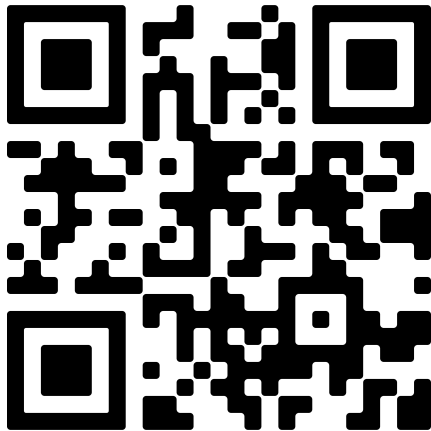
☒ OpenSource

☒ Unknown

Rank	Method	Param (B)	Language Model	Vision Model	Eval Date	Avg Score	Avg Rank	MMBench_V11	MMStar	MMMU_VAL	MathVista	OCRBench	AI2D	HallusionBench	MMVet
1	SenseNova				2024/12/12	77.4	4.25	85.7	72.7	69.6	78.4	894	87.8	57.4	78.2
2	InternVL2.5-78B-NPO	78	Qwen-2.5-72B	InternViT-6B-v2.5	2024/12/29	77	3.88	87.7	72.1	68.2	76.6	909	89.2	58.1	73.5
3	TeleMM				2024/12/31	75.9	8.88	79.9	70.8	66.6	75.7	891	88.5	60.6	75.7
4	InternVL2.5-38B-NPO	38	Qwen-2.5-32B	InternViT-6B-v2.5	2024/12/28	75.3	6.88	85.4	70.1	63.8	73.6	894	87.9	59.7	72.6
5	InternVL2.5-78B	78	Qwen-2.5-72B	InternViT-6B-v2.5	2024/12/10	75.2	7	87.5	69.5	70	70.6	853	89.1	57.4	71.8
6	Qwen2-VL-72B	73.4	Qwen2-72B	QwenViT	2024/10/28	74.8	8	85.9	68.6	64.3	69.7	888	88.3	58.7	73.9
7	Qwen-VL-Max-8889	72	Qwen2-72B	QwenViT	2024/09/12	74.4	9	85.8	69.2	64.6	68.3	881	88.1	59.2	72.3
8	InternVL2.5-38B	38	Qwen-2.5-32B	InternViT-6B-v2.5	2024/12/10	73.5	12.5	85.4	68.5	64.6	72.4	841	87.6	57.9	67.2



LINK TO PRESENTATIONS AND CODE



SOURCES

- Vision Transformer (ViT):
An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale
- Qwen-VL:
Qwen-VL: A Versatile Vision-Language Model for Understanding, Generation, and Retrieval
- Qwen-VL-2:
Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution
- VisionLLaMA:
VisionLLaMA: A Unified LLaMA Backbone for Vision Tasks
- Pixtral:
Pixtral 12B



Temporary page!

\LaTeX was unable to guess the total number of pages correctly. As there was some unprocessed data that should have been added to the final page this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page will go away, because \LaTeX now knows how many pages to expect for this document.