

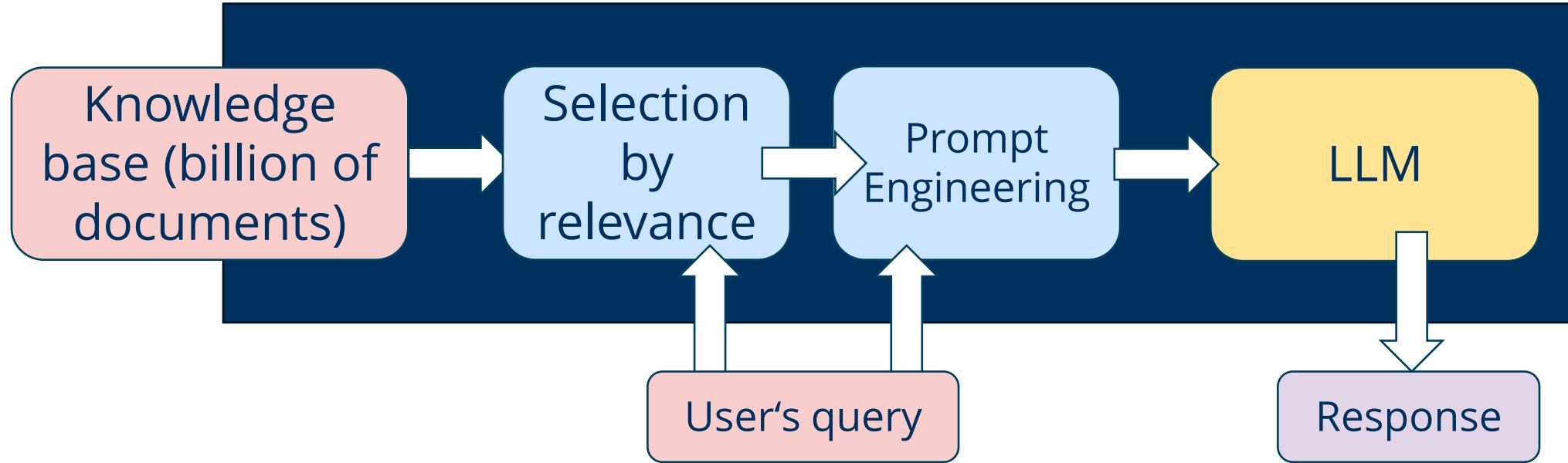
Multi-modal Language Models

Robert Haase

These slides can be reused under the terms of the [CC-BY 4.0](#) license unless mentioned otherwise.

Quiz

How is this prompt engineering technique / architecture called?



Reflection



Long-context prompting



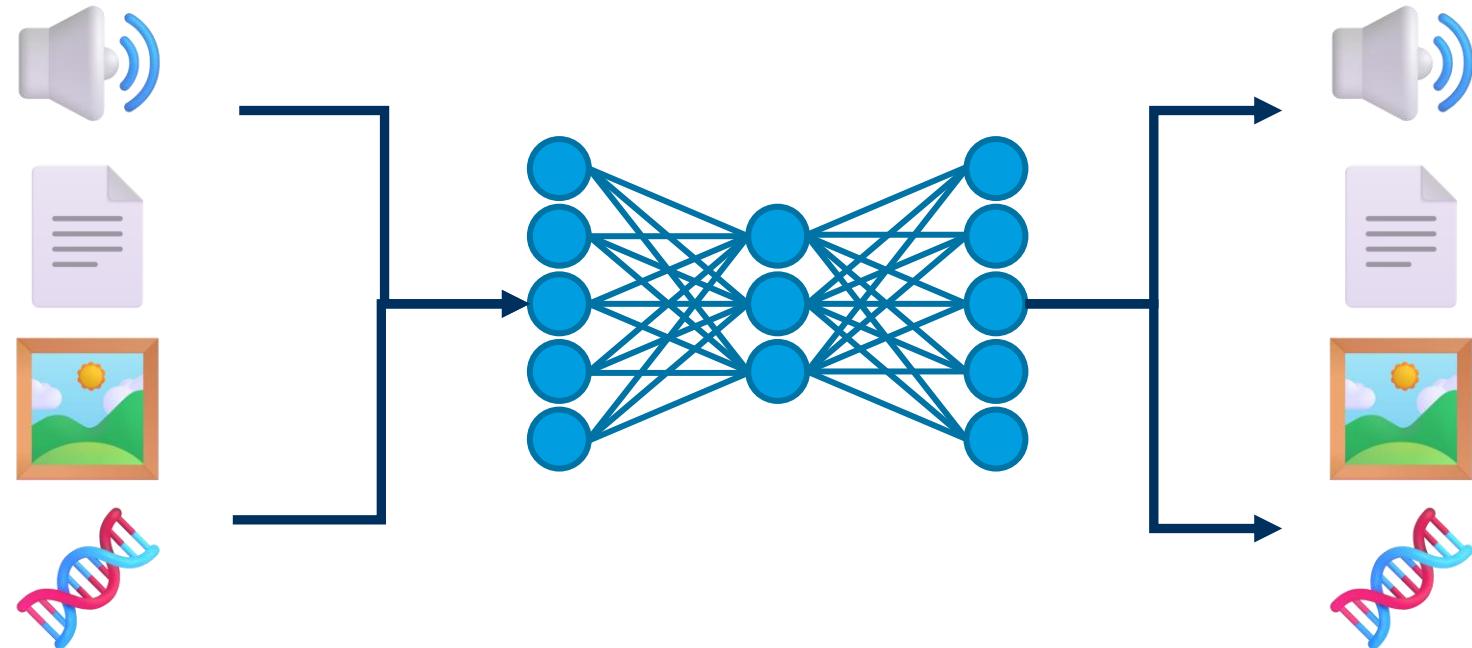
Retrieval augmented generation



Knowledge distillation

Multi-modal LLMs

Combining image, text and [...] data



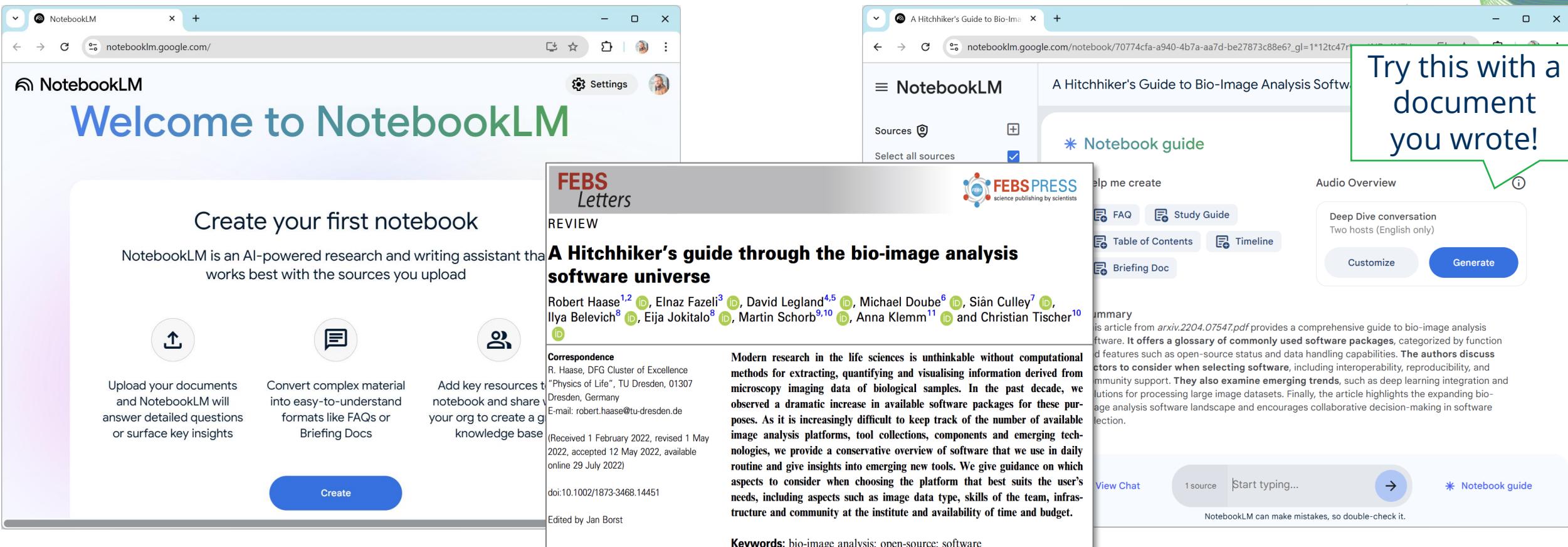
Multi-modal LLMs

Combining image, text and [...] data, to gain new [biological] insights.

- Chat / translation / text-generation:  -> 
 - Voice-chat / -translation:  -> 
 - Text-to-speech (TTS):  -> 
 - Speech to text (STT):  -> 
 - Vision / image classification / image description / image-to-text:  -> 
 - Image generation, text-to-image:  -> 
 - Image variation, inpainting, image segmentation:  -> 
 - Genotype-phenotype translation:  -> 
 - Sequence-prediction / protein design:  -> 
- } Focus of today's lecture

Use case: Paper to Podcast

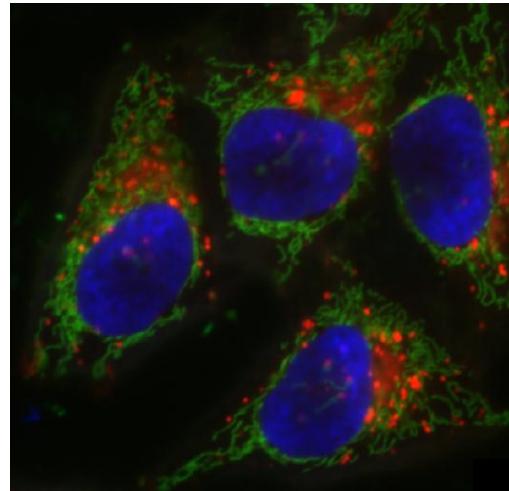
Text to audio conversion: Google NotebookLM  -> 



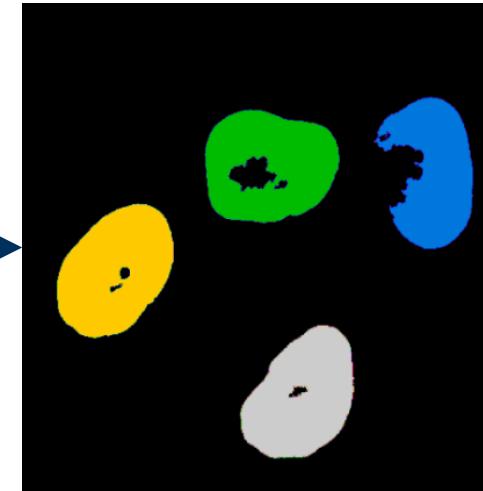
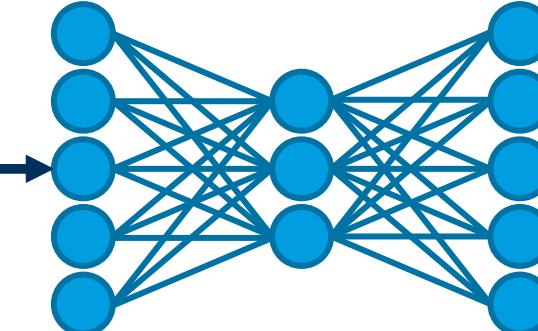
The screenshot shows two browser windows. The left window is the 'NotebookLM' homepage, featuring a large green 'Welcome to NotebookLM' header and sections for creating a first notebook, uploading documents, and converting complex material into FAQs or briefing docs. The right window shows a research article titled 'A Hitchhiker's guide through the bio-image analysis software universe' by Robert Haase et al. The article is from FEBS Letters and discusses software packages for bio-image analysis. A green callout box on the right side of the right window says 'Try this with a document you wrote!'.

Multi-modal LLMs

Combining image, text and [...] data, to gain new [biological] insights.



How many cells
are there?

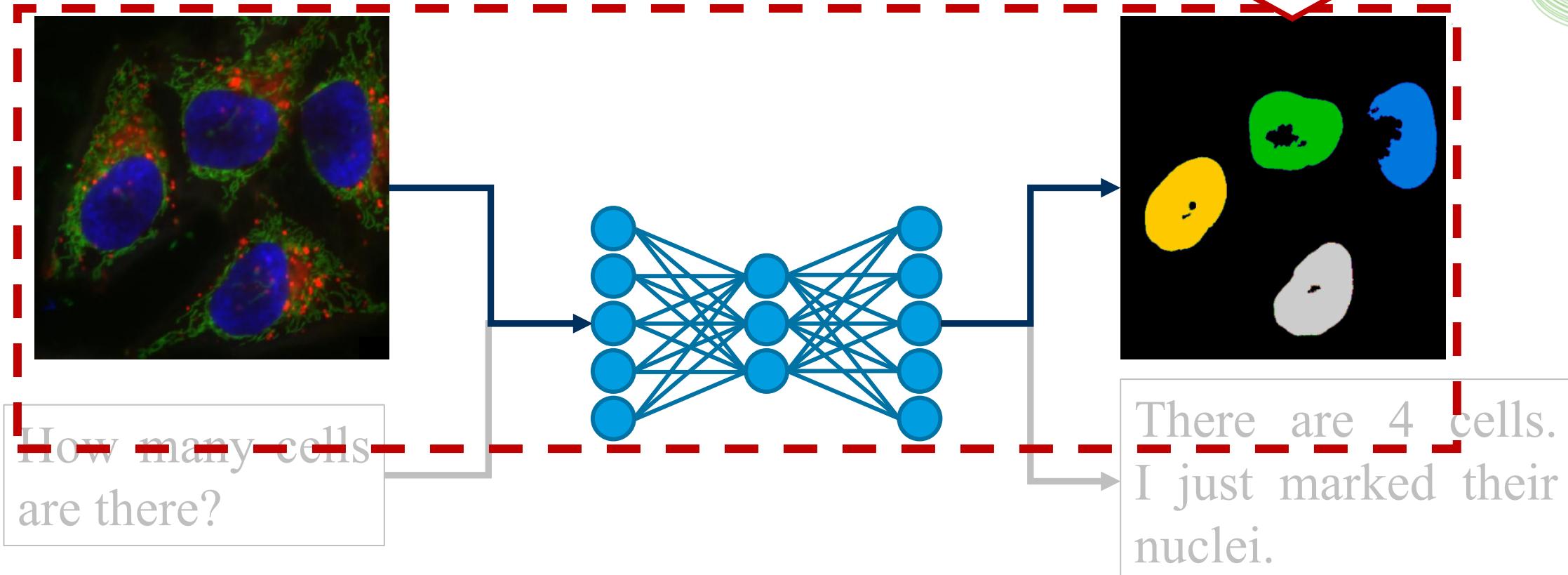


There are 4 cells.
I just marked their
nuclei.

Multi-modal LLMs

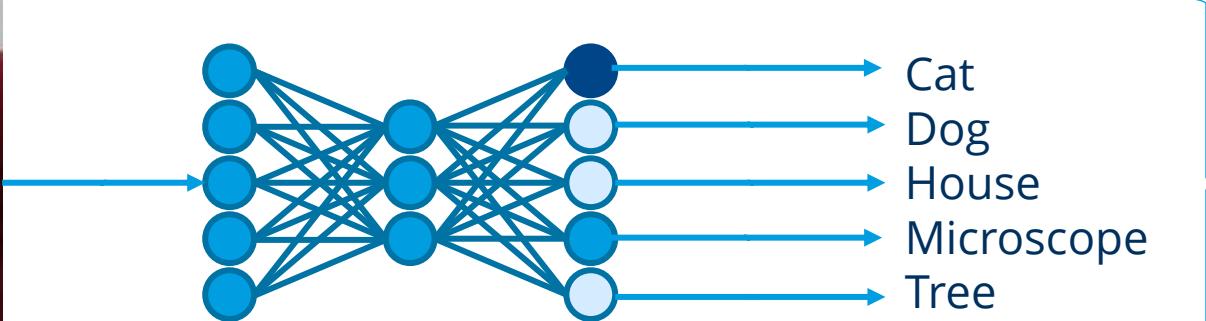
Combining image, text and [...] data, to gain new [biological] insights

Quiz: Could a neural network predict an instance label image directly?



Vision Models

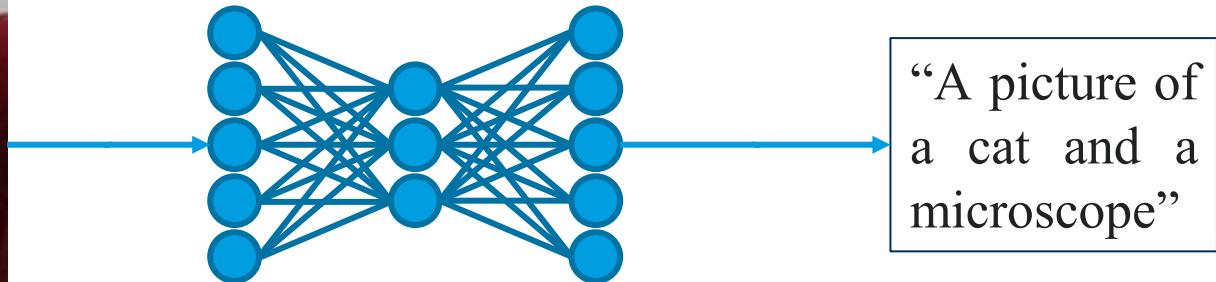
- Classifying images (decades old research field 😊)



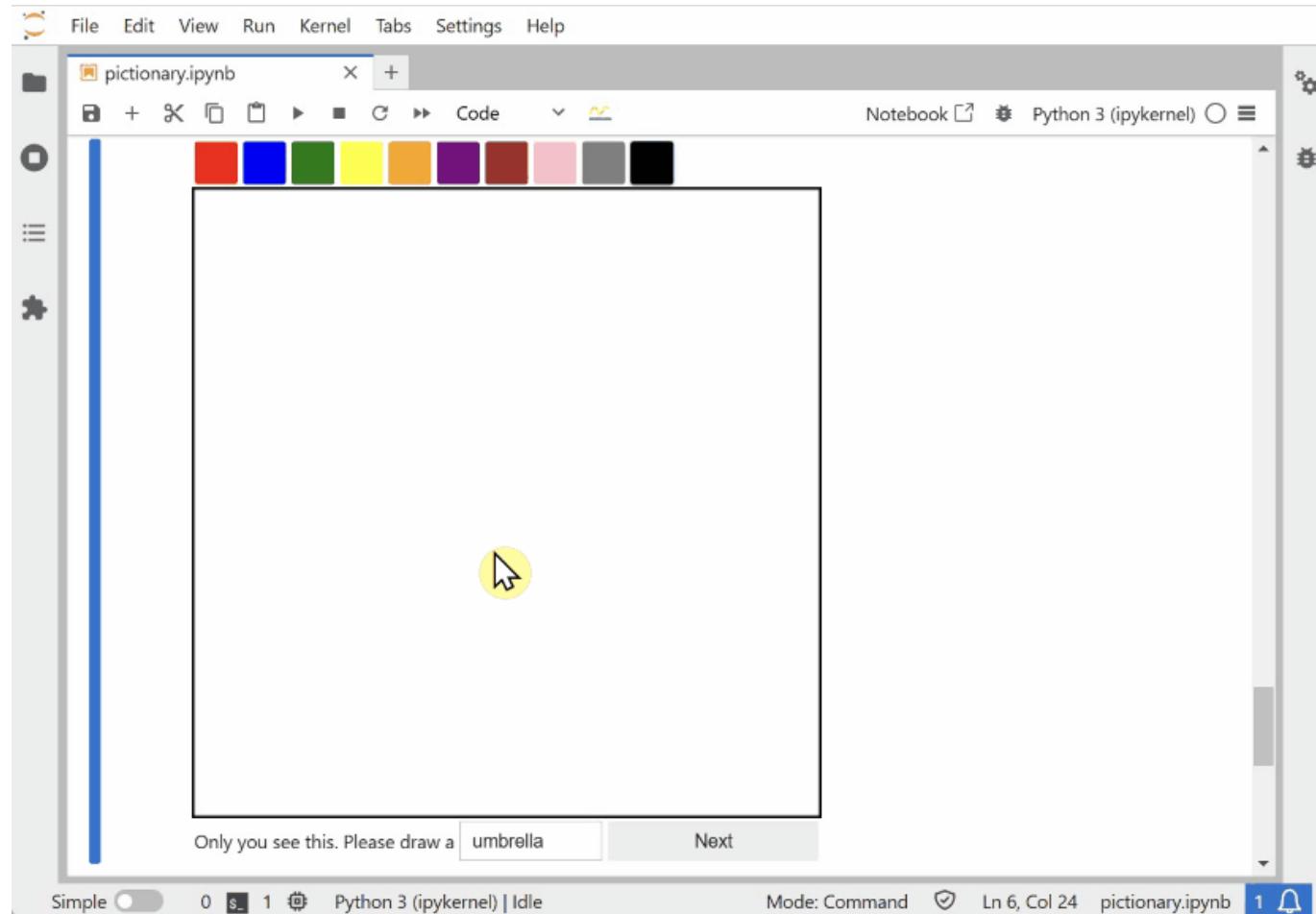
Typically a fixed number of classes, defined during training

Vision Language Models

- Classifying images
- Describing images



Gamification: VLM-Pictionary

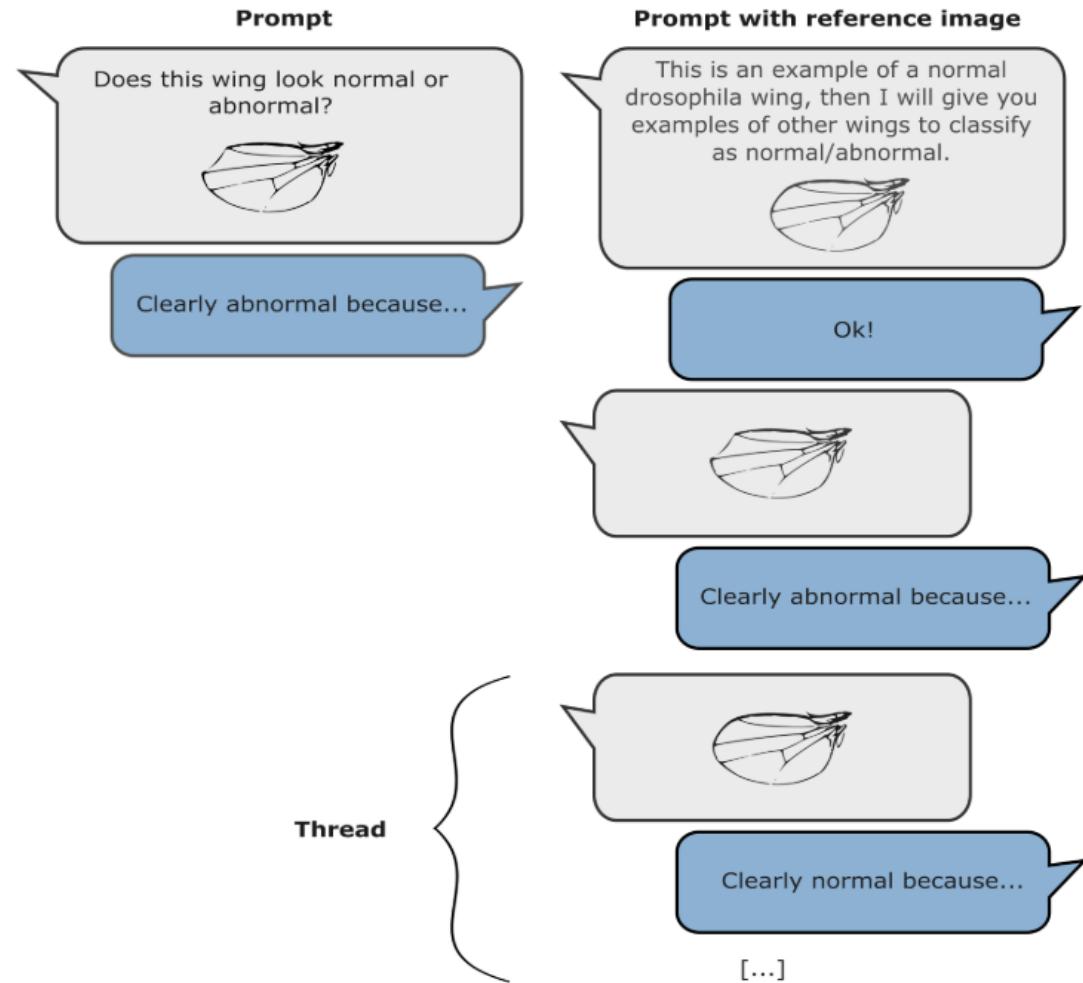


Use case: Descriptive biology

Table 1. Dataset: phenotypes analyzed and example images.

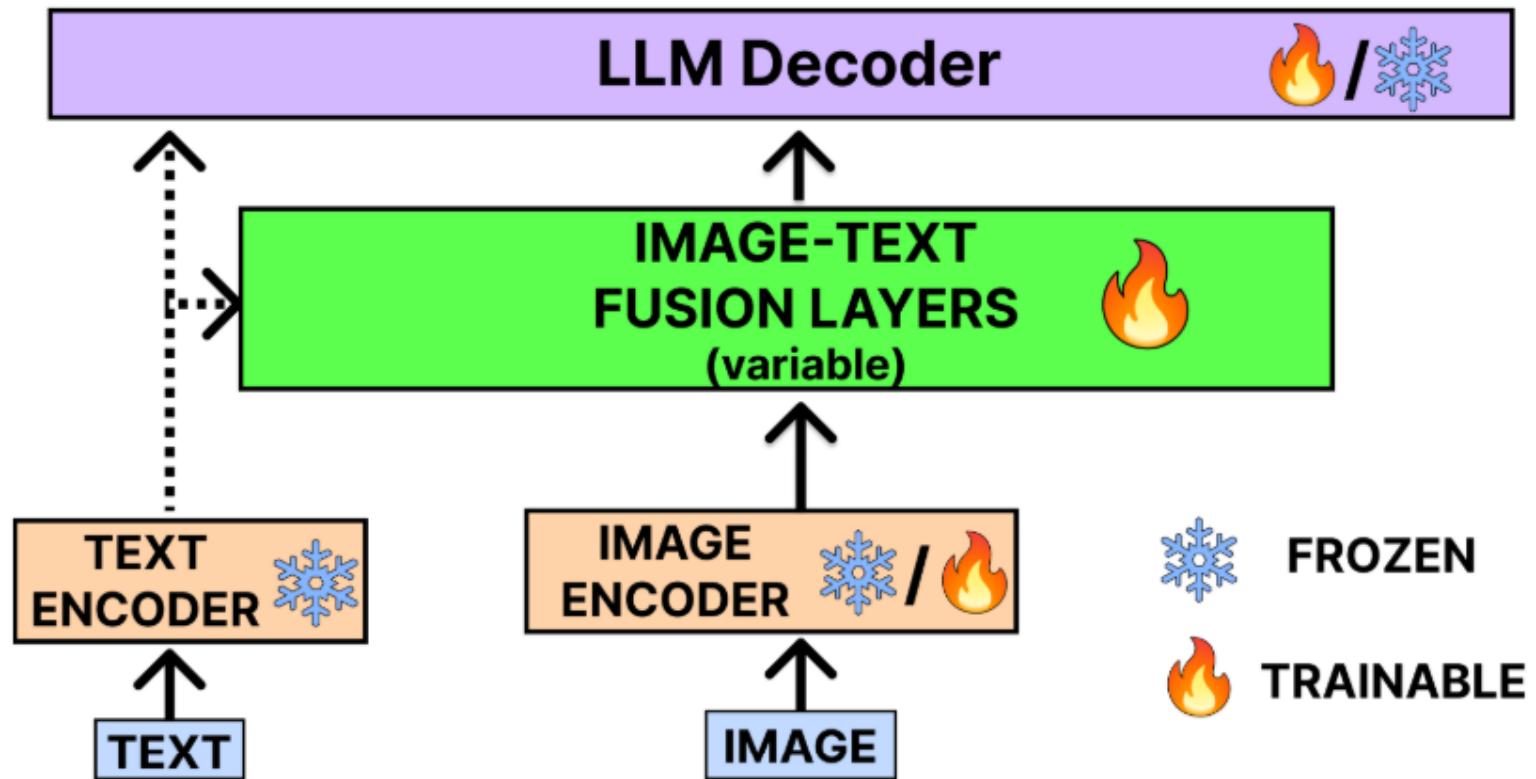
Phenotype	# Imgs	Example
Loss of wing veins (V-)	10	
Ectopic wing veins (V+)	20	
Integrity of wing margin (WM)	23	
Wing surface adhesion (WA)	27	

“... while visual language models are in their infancy, they already show potential for multiple applications in automated phenotyping studies. We encourage the community to carefully test them...”



Vision Language Models (VLMs)

Goal: Describe images



Variational Auto-Encoder

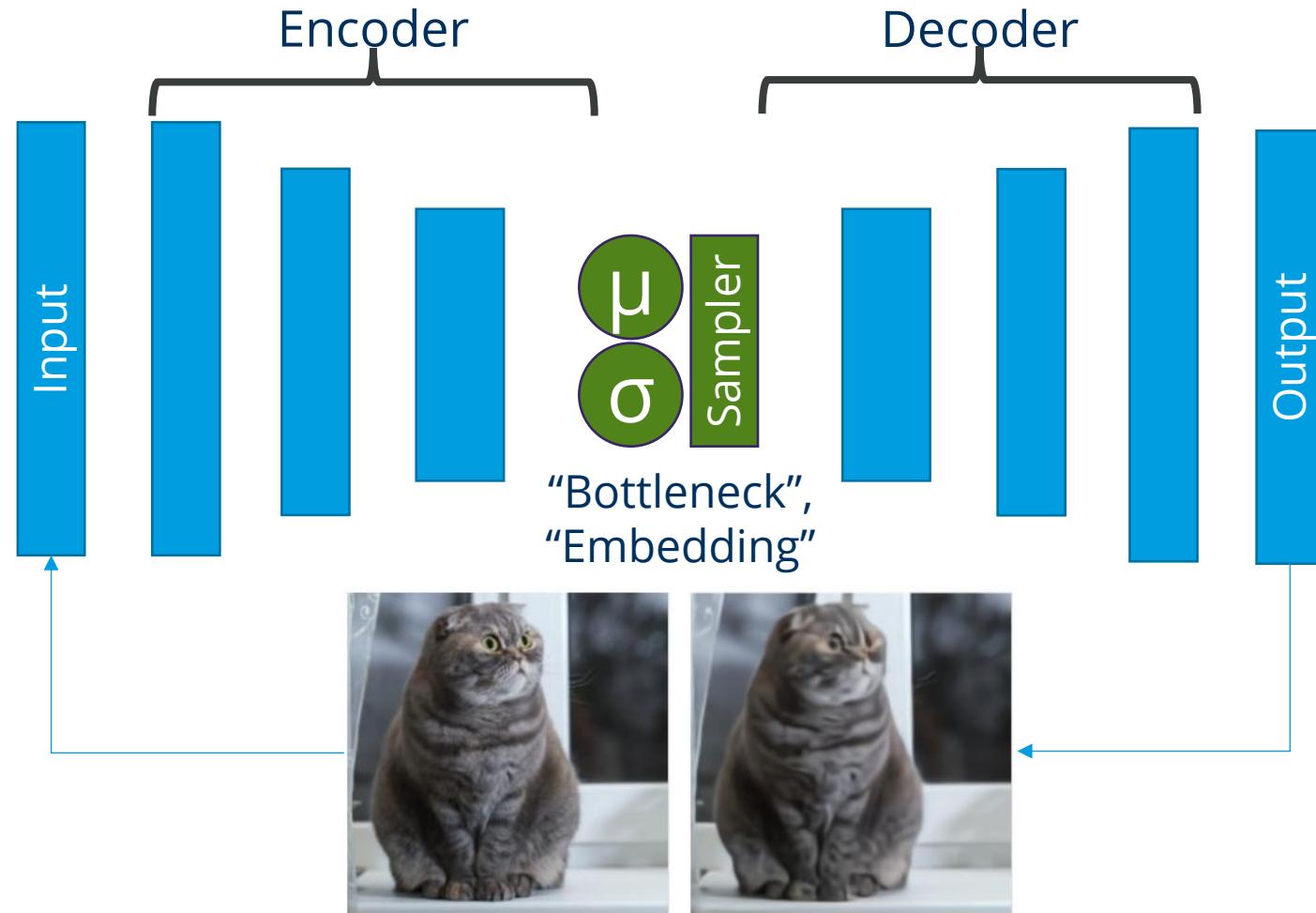


Image classification → image describing

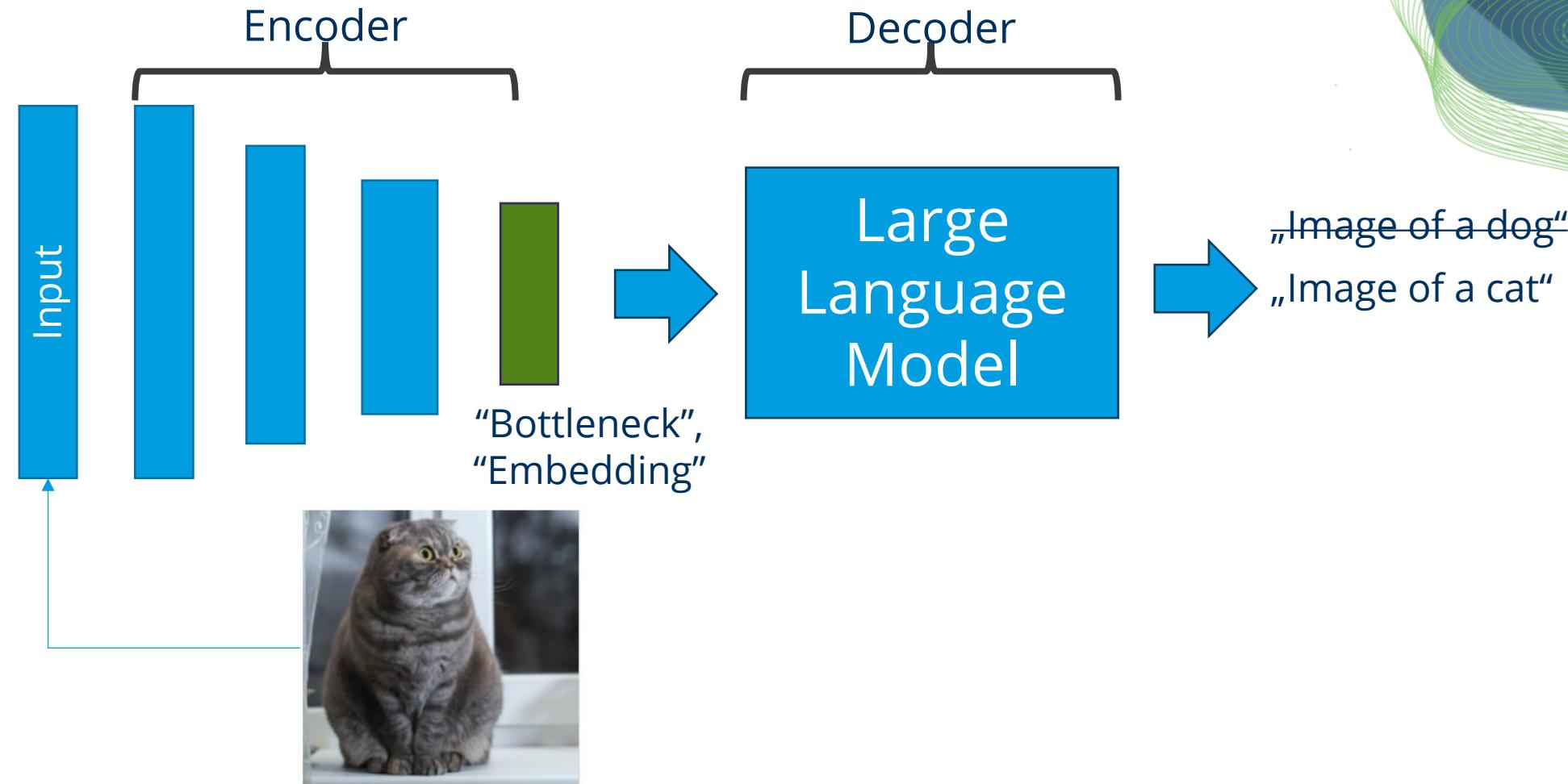
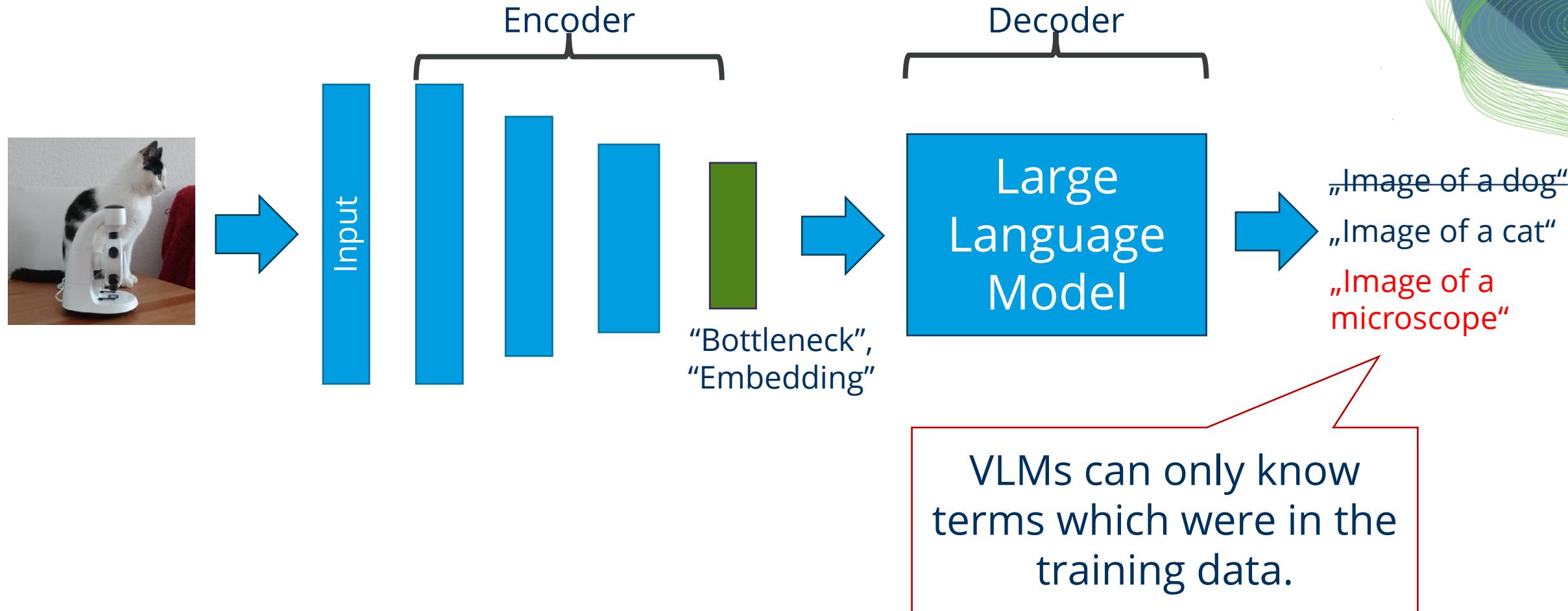


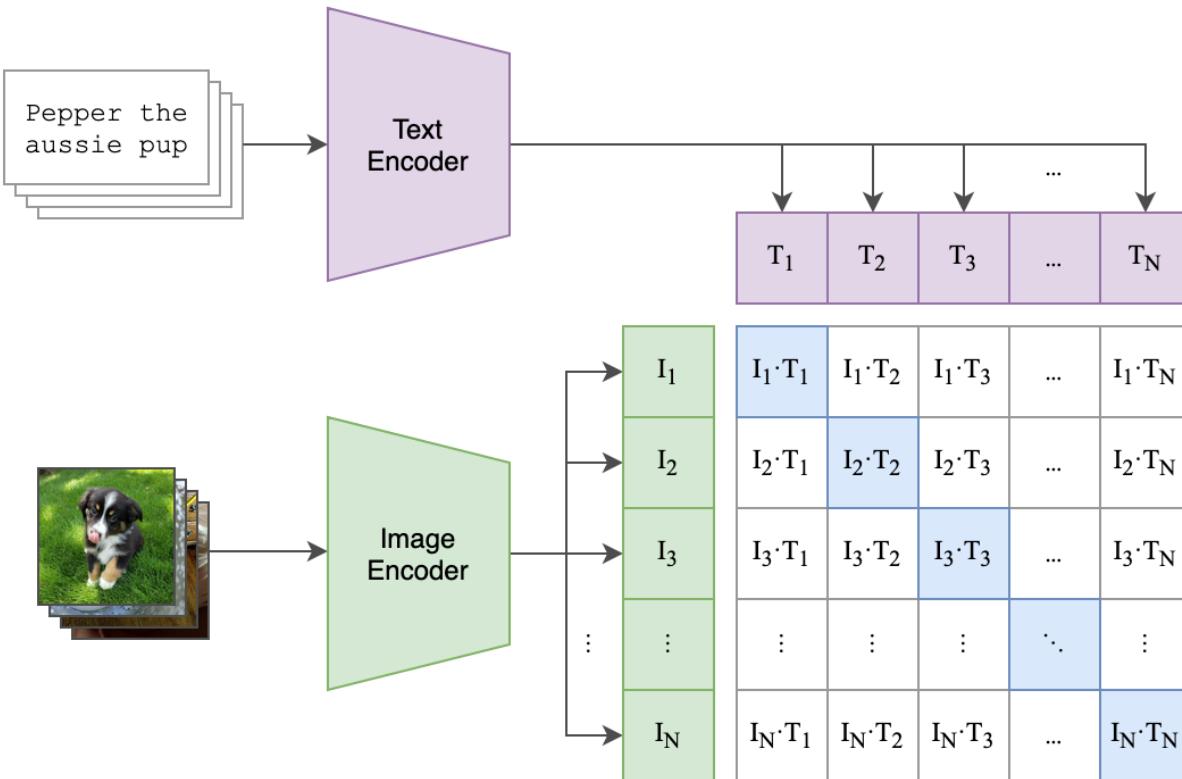
Image classification → image describing



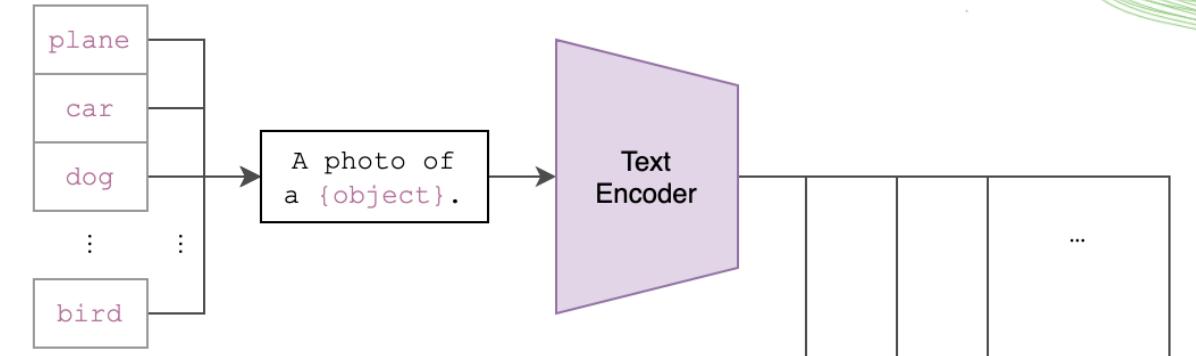
Contrastive Language-Image Pre-Training

„CLIP“ Transformers

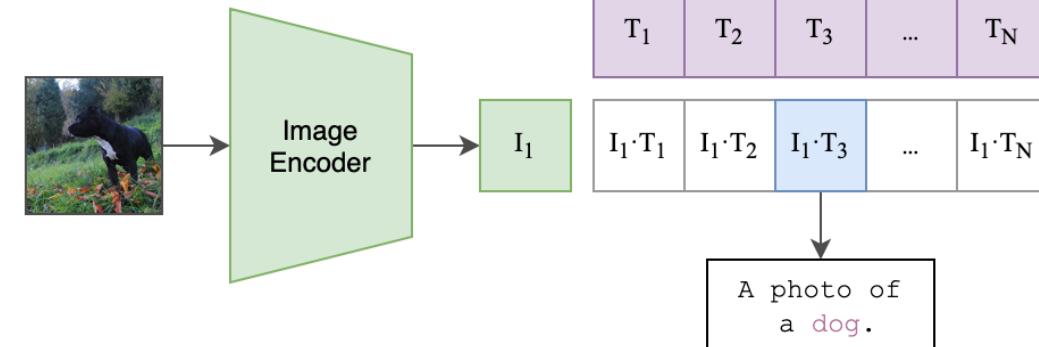
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



CLIP transformers in Python

Using huggingface 😊



Downloads
500 MB

```
model = CLIPModel.from_pretrained("openai/clip-vit-base-patch32")
processor = CLIPProcessor.from_pretrained("openai/clip-vit-base-patch32")
```

```
options = ["a photo of a cat",
           "a photo of a dog"]
```

```
options = ["a photo of a cat",
           "a photo of a dog",
           "a photo of a microscope"]
```

```
inputs = processor(text=options, images=image, return_tensors="pt", padding=True)
outputs = model(**inputs)
```

```
label_probabilities
```

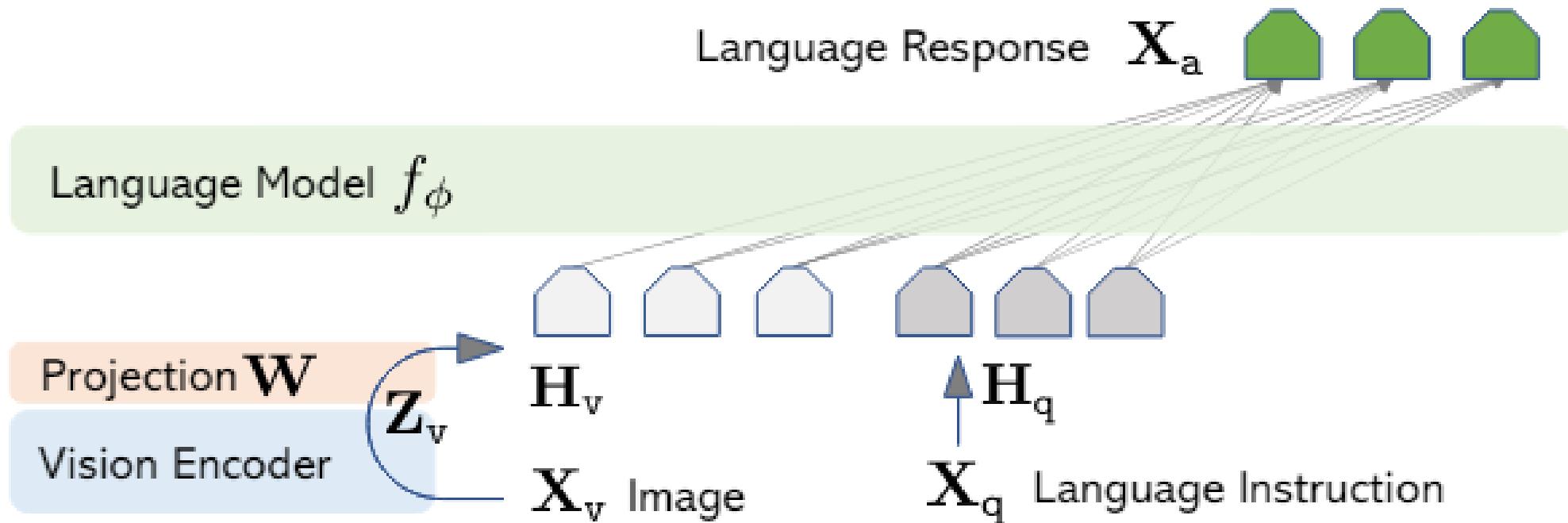
```
{'a photo of a cat': 0.9907298684120178,
 'a photo of a dog': 0.009270114824175835}
```

...

```
label_probabilities
```

```
{'a photo of a cat': 0.1352911740541458,
 'a photo of a dog': 0.0012659047497436404,
 'a photo of a microscope': 0.8634429574012756}
```

Large Language and Vision Assistant



LLAVA 1.5

Combining LLAVA with CLIP

language model (Vicuna v1.5 13B)



vision-language connector (MLP)

vision encoder (CLIP ViT-L/336px)

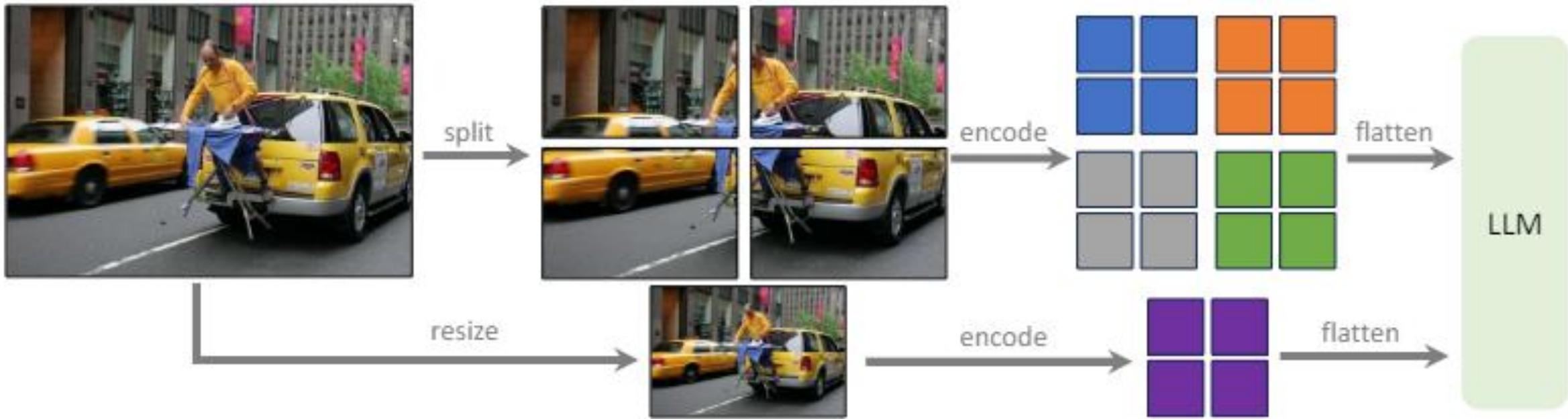


tokenizer & embedding

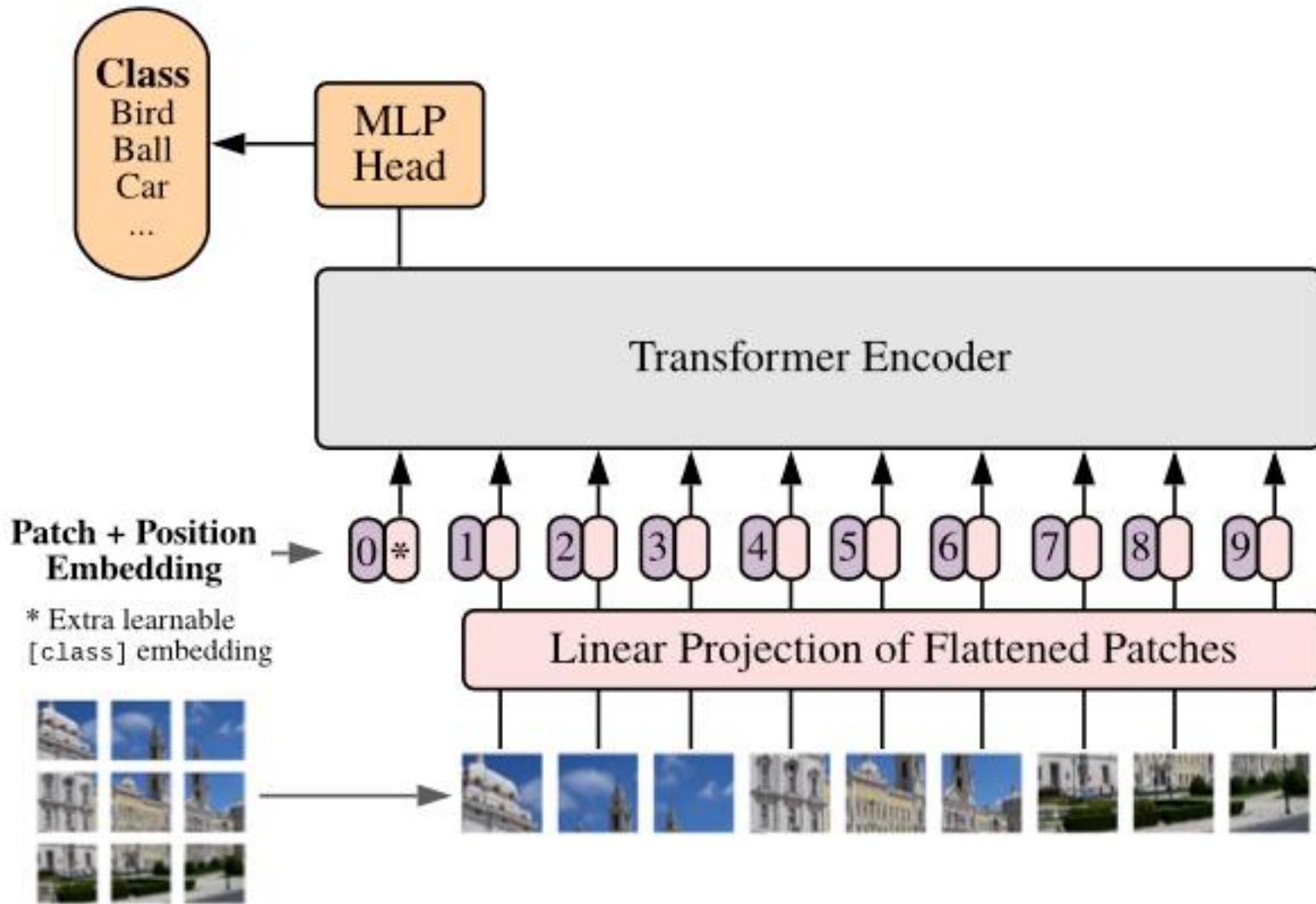
User: what is unusual about this image?

LLAVA 1.5 HD

Giving the model multiple perspectives on the same scene



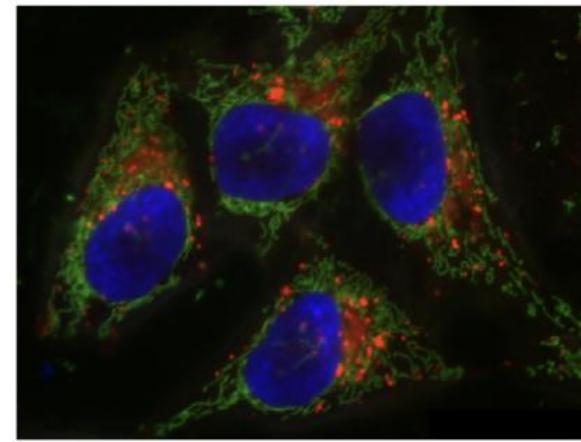
Vision-Transformer (ViT)



Accessing VLMs using Python

E.g. using ScaDS.AI's openai-compatible LLM Server

```
def prompt_qwen(prompt:str, image, model="Qwen/Qwen2-VL-7B-Instruct"):  
    """A prompt helper function that sends a message to the server  
    and returns only the text response.  
    """  
  
    rgb_image = _img_to_rgb(image)  
    byte_stream = numpy_to_bytestream(rgb_image)  
    base64_image = base64.b64encode(byte_stream).decode('utf-8')  
  
    message = [{"role": "user", "content": [  
        {"type": "text", "text": prompt},  
        {  
            "type": "image_url",  
            "image_url": {  
                "url": f"data:image/png;base64,{base64_image}"  
            }  
        }]  
    }]  
  
    # setup connection to the LLM  
    client = openai.OpenAI(base_url="https://llm.scads.ai/v1",
```



```
res = prompt_qwen("what's in this image?", hela_cells)  
display(Markdown(res))
```

Accessing VLMs using Python: No common API

```
if image is None:  
    message = [{"role": "user", "content": message}]  
else:  
    encoded_image = image_to_url(image)  
    message = [{  
        "role": "user",  
        "content": [  
            {  
                "type": "text",  
                "text": message  
            },  
            {  
                "type": "image",  
                "source": {  
                    "type": "base64",  
                    "media_type": "image/png",  
                    "data": encoded_image  
                }  
            }  
        ]  
    }]  
}
```

Anthropic / claude

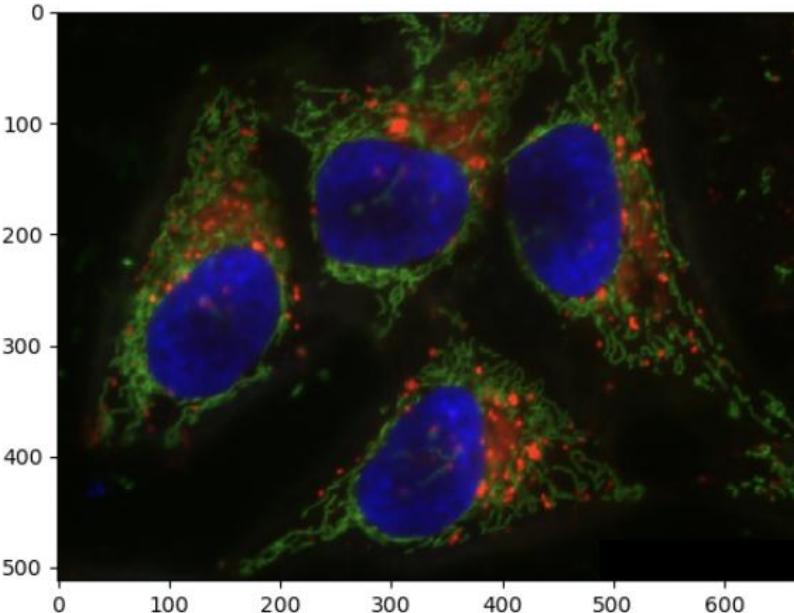
```
if image is not None:  
    response = client.generate_content([image, request])  
else:  
    response = client.generate_content(request)
```

Google / Gemini

Structured output

Ask VLMs to describe an image and pre-define a structure

```
hela = imread("data/hela-cells-8bit.tif")
stackview.insight(hela)
```



```
result = prompt_chatGPT("""You are a highly experienced biologist with advanced microscopy skills.

# Task
Name the content of this image. Answer for each channel independently.

# Options
The following structures could be in the image:
* Nuclei
* Membranes
* Cytoplasm
* Cytoskeleton
* Extra-cellular structure
* Other sub-cellular structures

# Output format
* Red channel: <structure>
* Green channel: <structure>
* Blue channel: <structure>

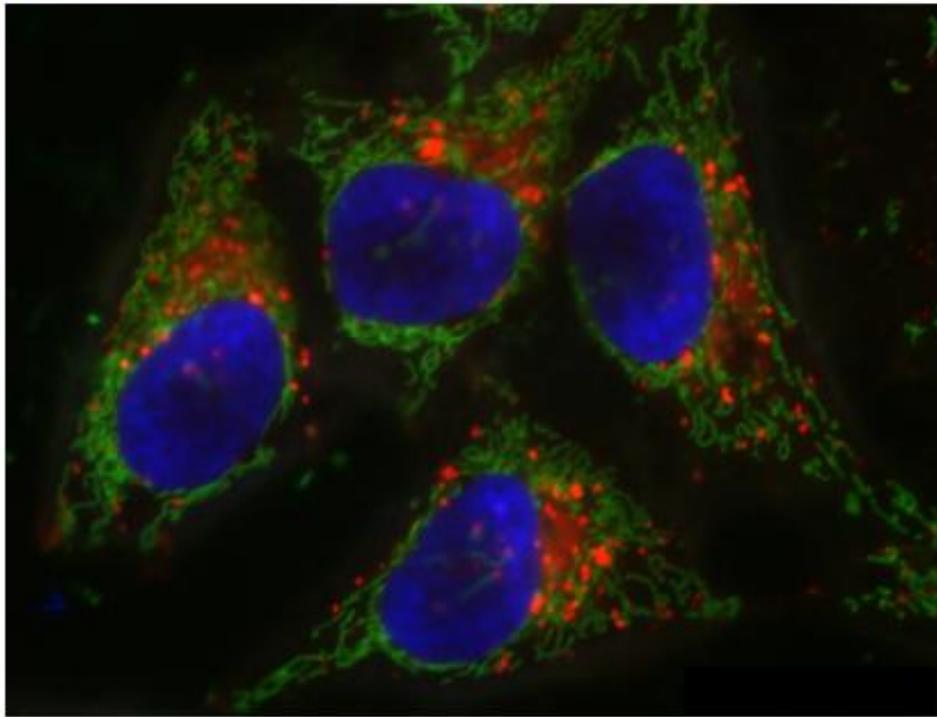
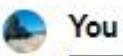
Keep your answer as short as possible.
Only respond with the structures for the three channels in the format shown above.
""", hela)

Markdown(result)
```

- Red channel: Other sub-cellular structures
- Green channel: Cytoskeleton
- Blue channel: Nuclei

Benchmarking vision models

Single attempts... are a trap



How many blue nuclei are in this image?



ChatGPT

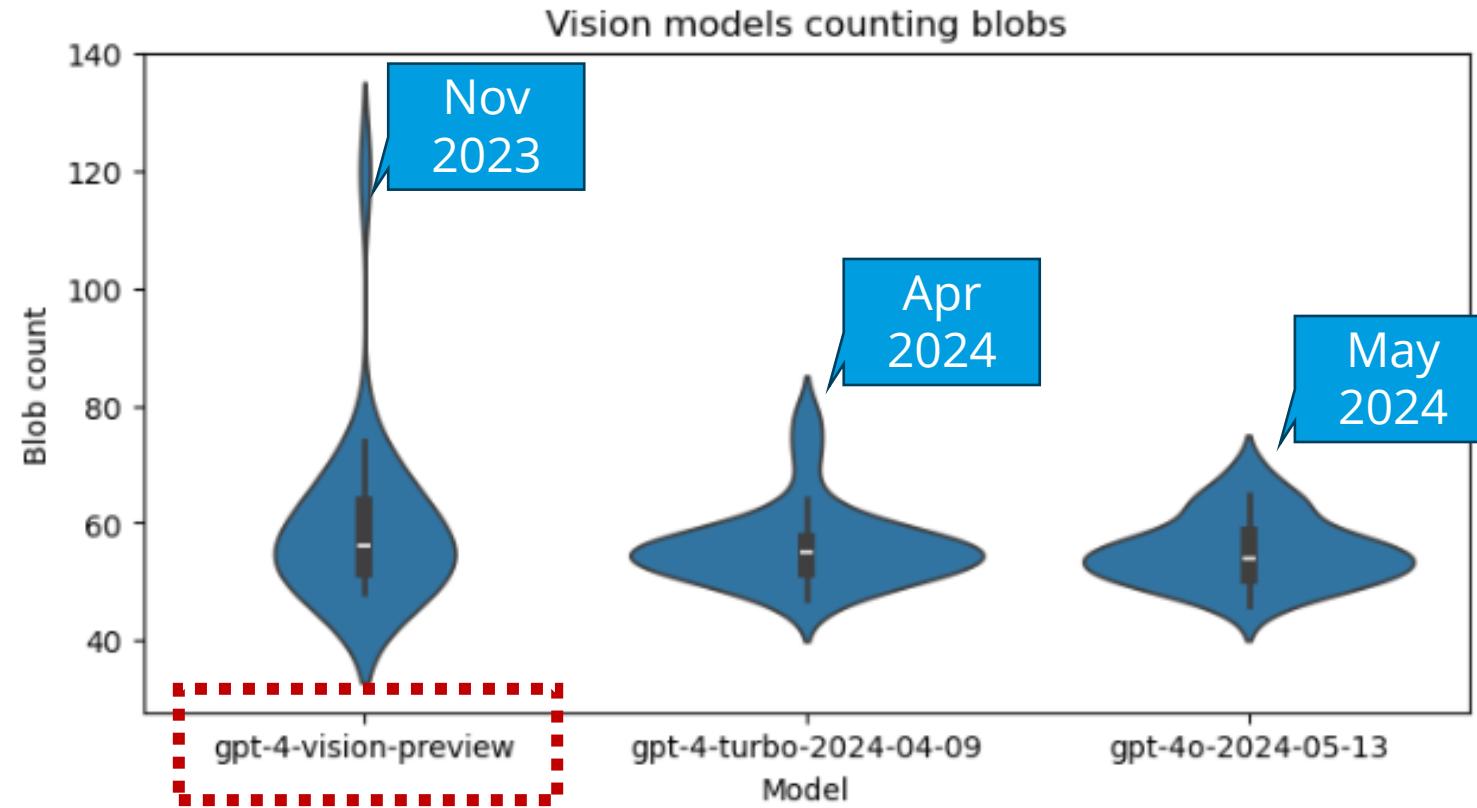
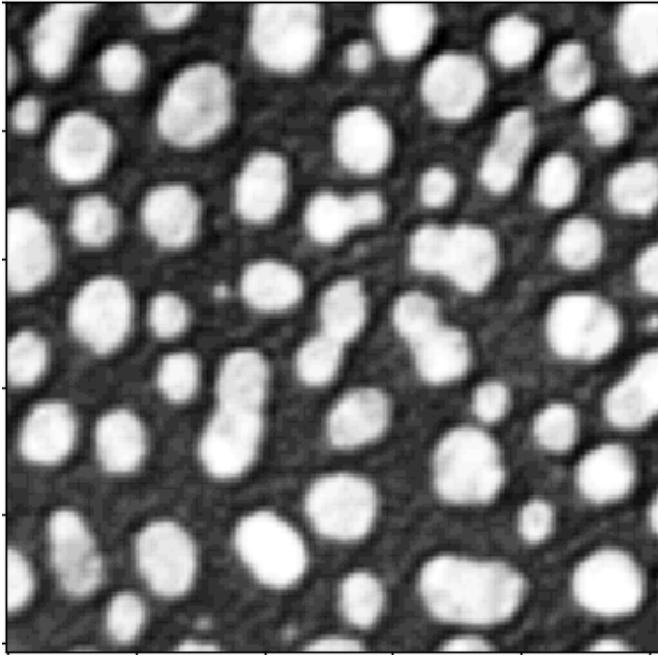
There are three blue nuclei visible in this image.



$n=1$

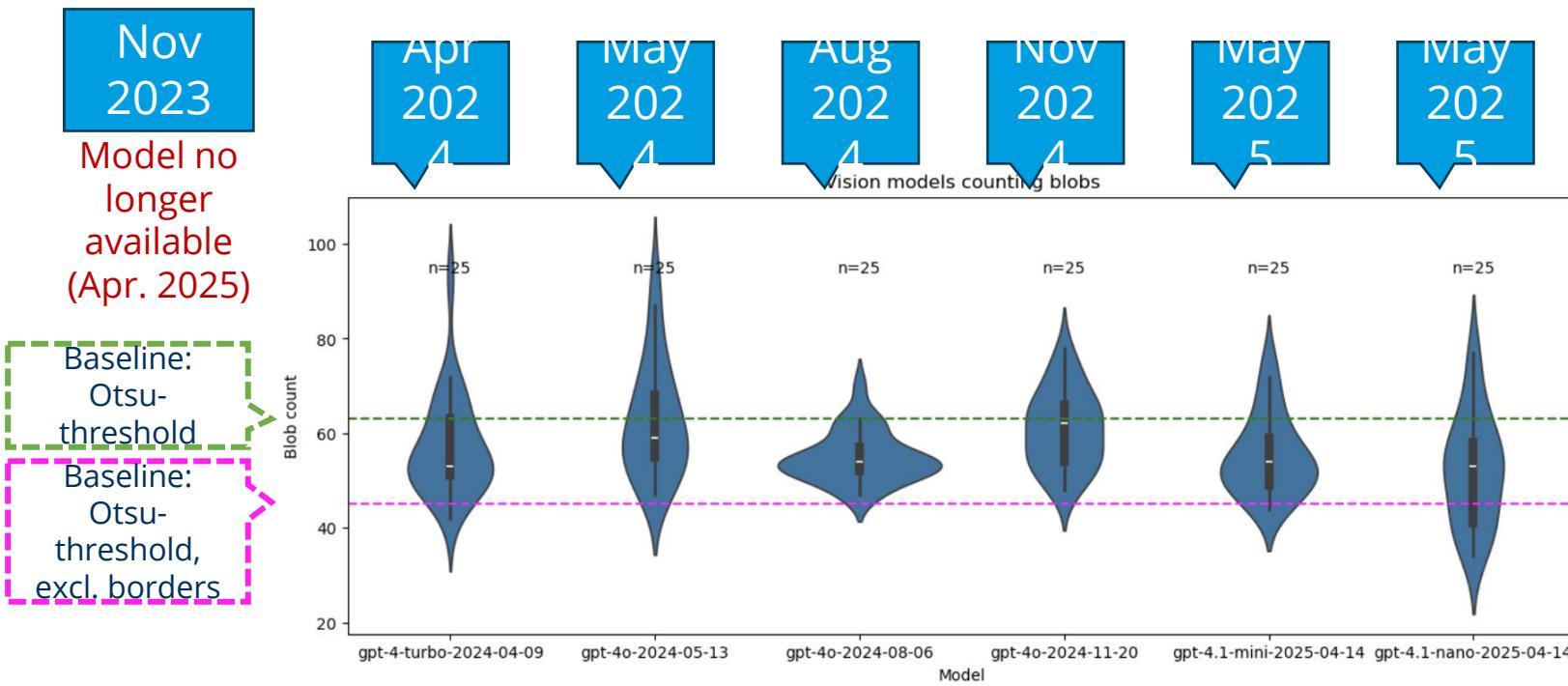
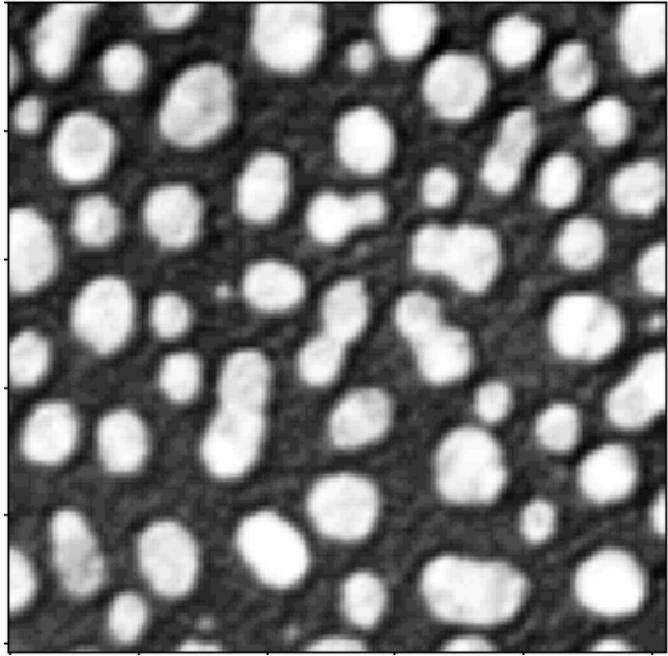
Vision language models for counting objects

Prompt: „Analyse the following image by counting the bright blobs. Respond with the number only.“ (n=25)



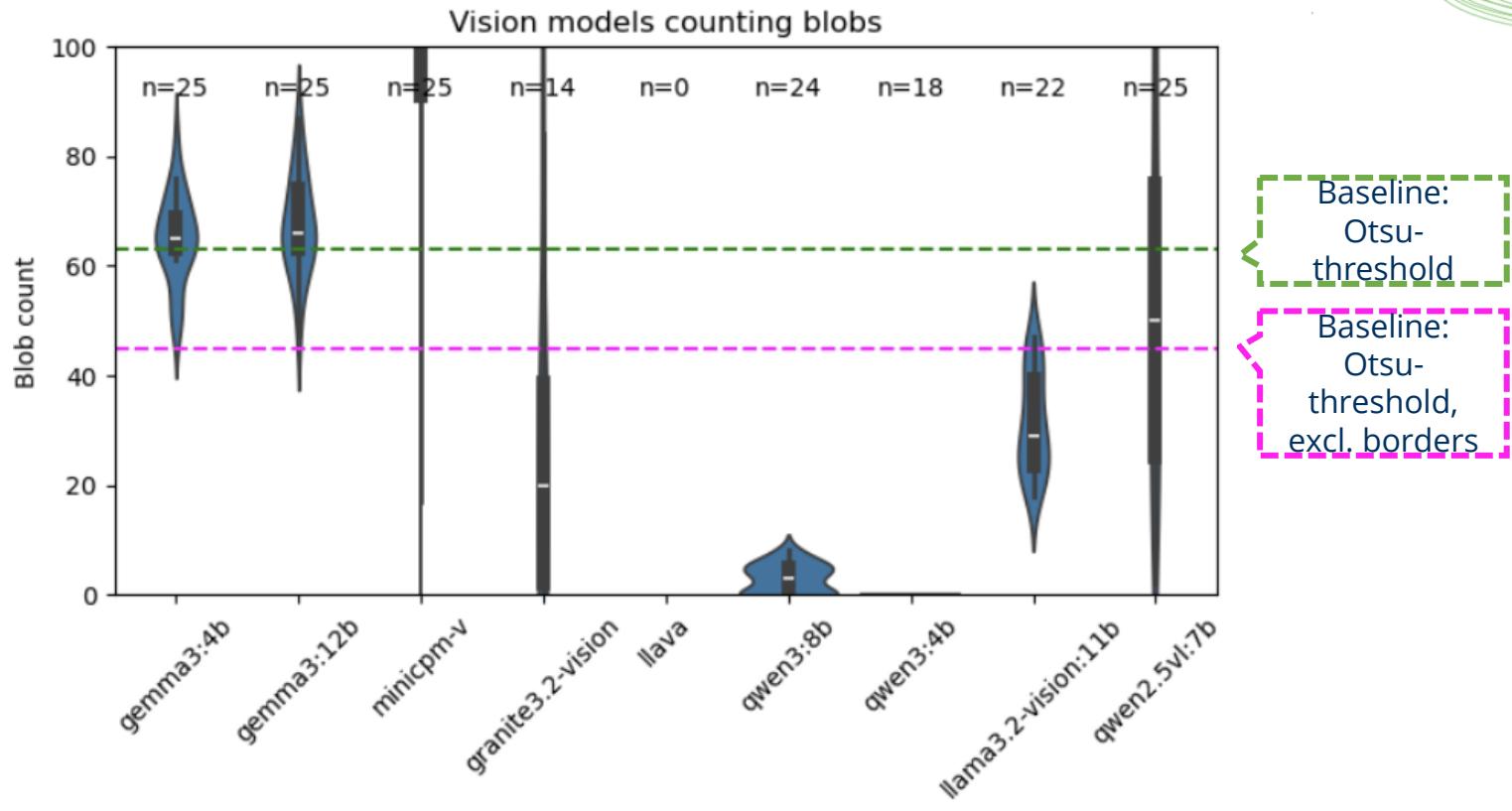
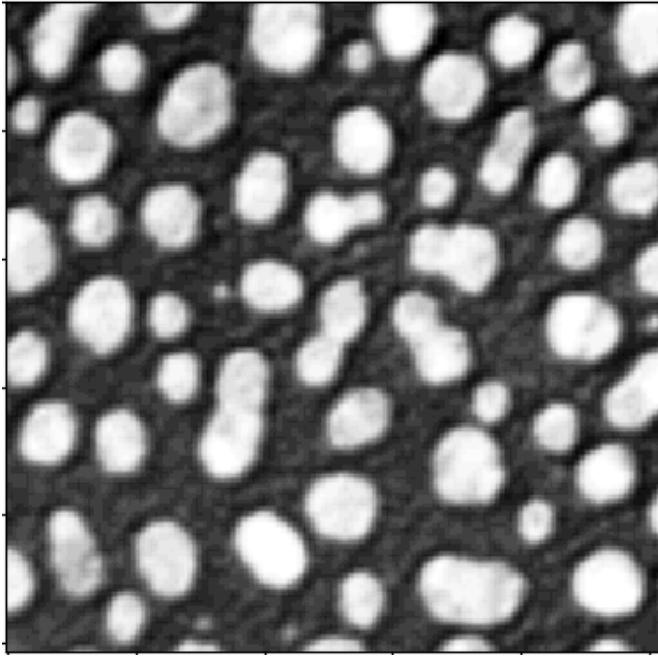
Vision language models for counting objects

Experiment reproduced June 24th 2025



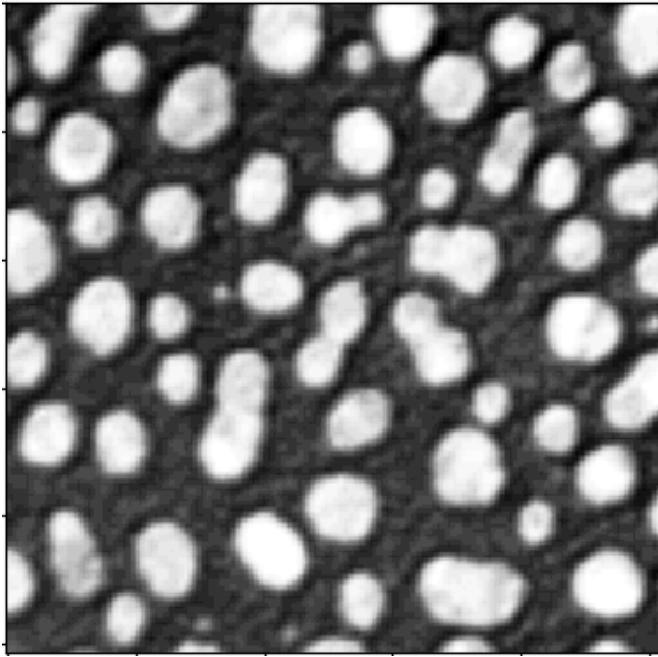
Vision language models for counting objects

Now also possible using open-weight
models (via ollama)



Vision language models for counting objects

Prompt-engineering also works with VLMS



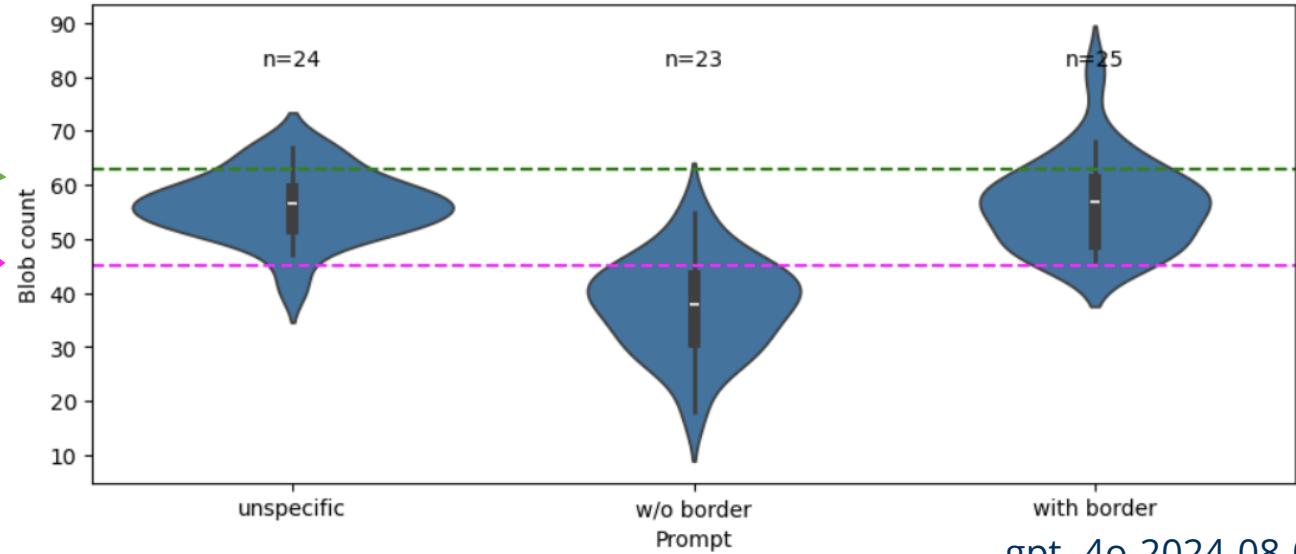
Baseline:
Otsu-threshold
Baseline:
Otsu-threshold,
excl. borders

"Analyse the following image by counting the bright blobs. Respond ONLY the number."

"Analyse the following image by counting the bright blobs. **Ignore the objects touching the image border.** Respond ONLY the number."

"Analyse the following image by counting the bright blobs, **including the objects touching the image border.** Respond ONLY the number."

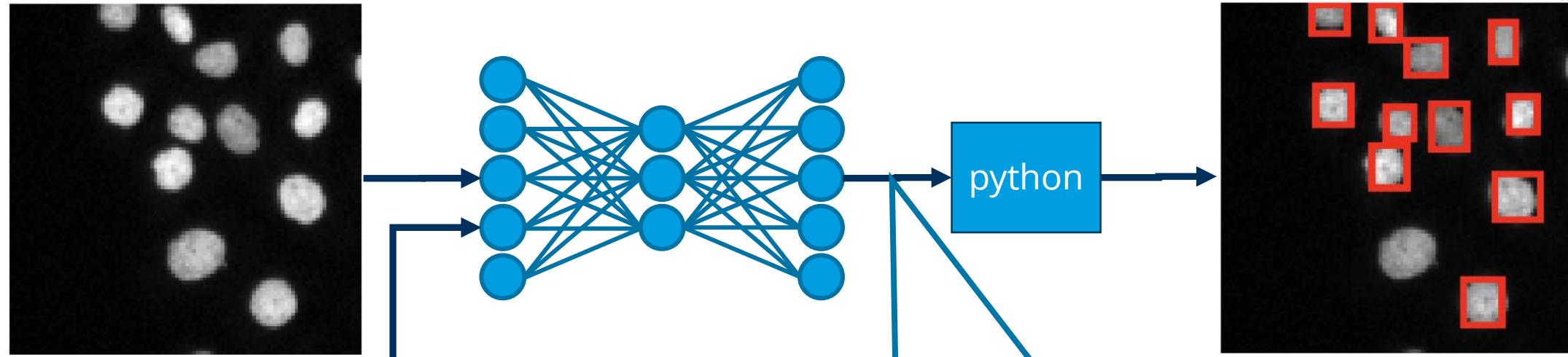
Prompts for counting blobs using VLMS



gpt -4o-2024-08-06

VLMs for bounding box segmentation

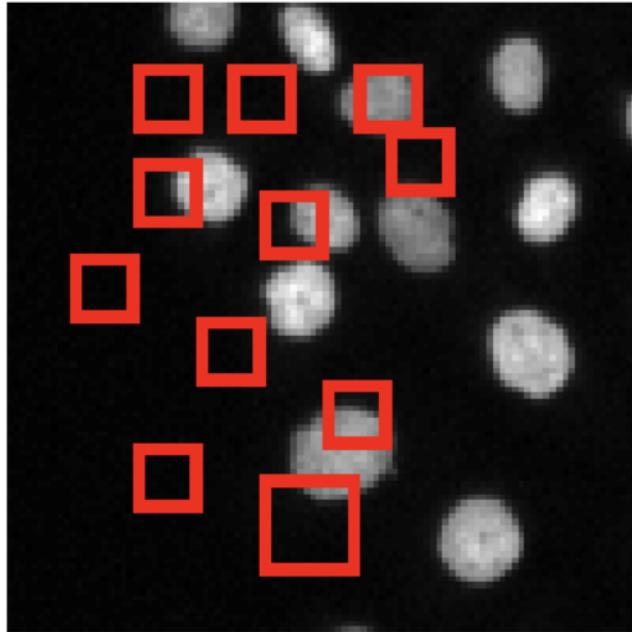
Recap: VLMs combine images + text to produce text



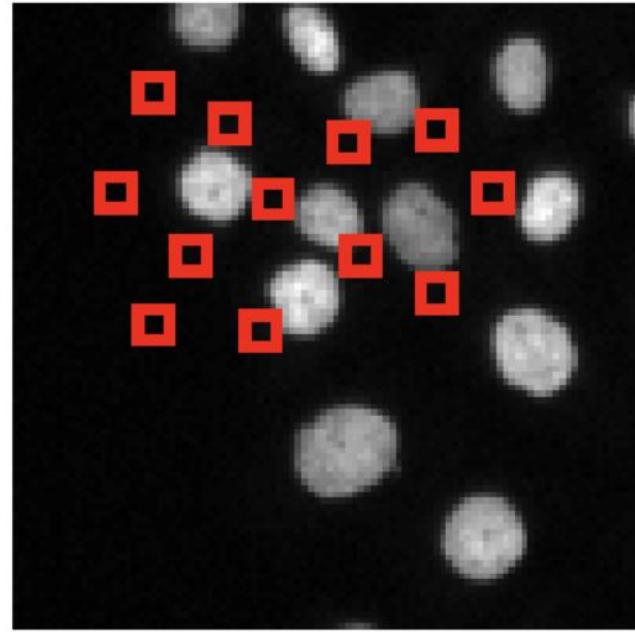
Give me a json object
of bounding boxes
around ALL bright
blobs in this image...

[{"x": 0.191, "y": 0.111,...},
 {"x": 0.313, "y": 0.161,...},...]

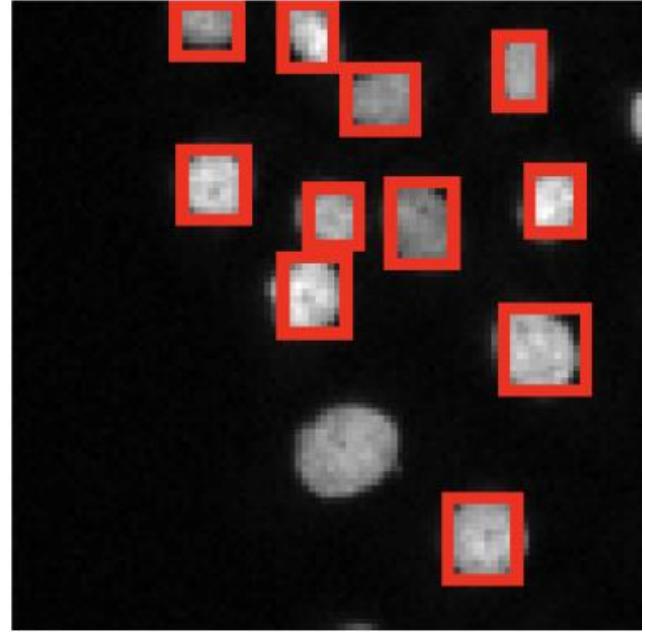
VLMs for bounding box segmentation



GPT-4o



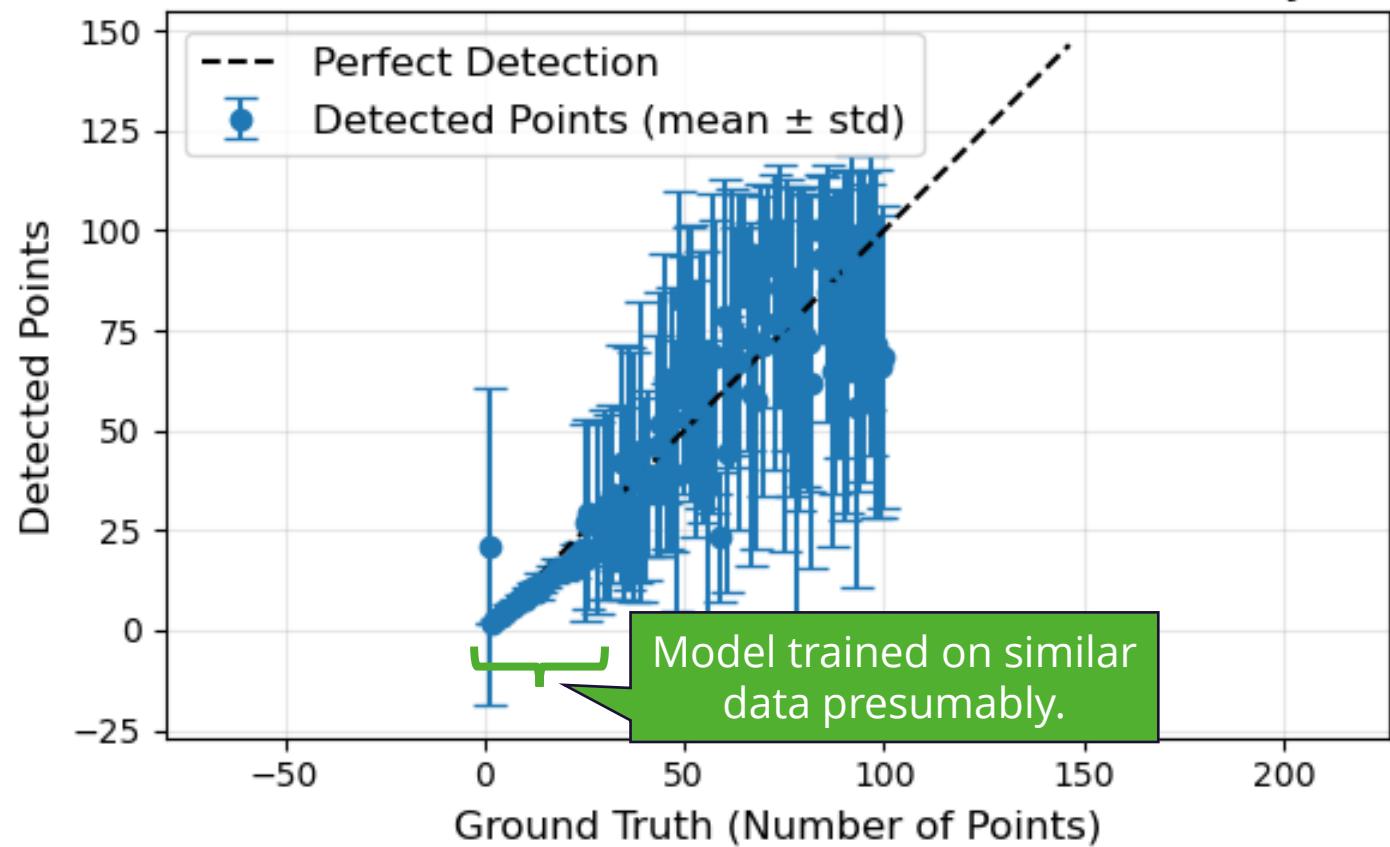
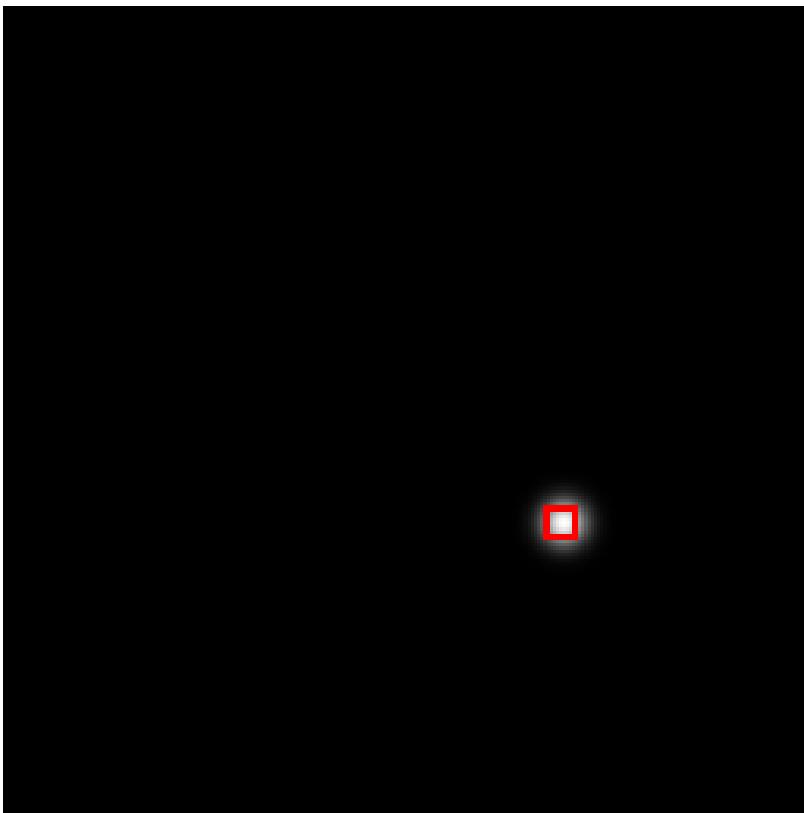
Claude 3.7 Sonnet



Moondream 2B

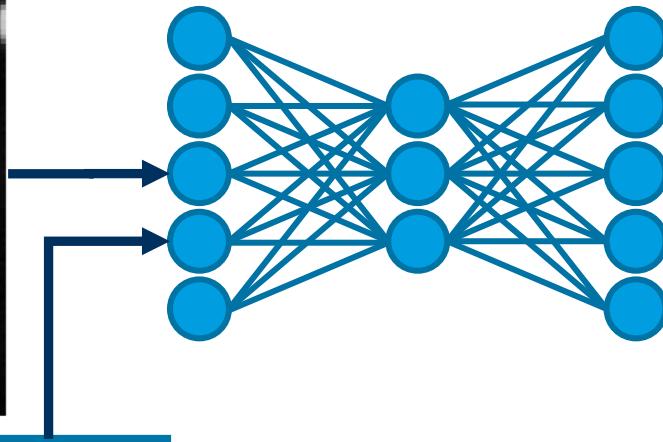
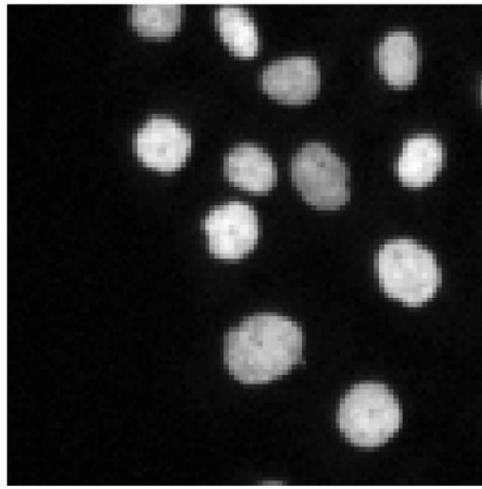
Limits of VLMS

Example: Moondream 2B applied to simulated data



VLMs guessing segmentation algorithms

Combine image + text to generate text



CellPose!

What's the best
algorithm to segment
this image?

LLMs guessing segmentation algorithms

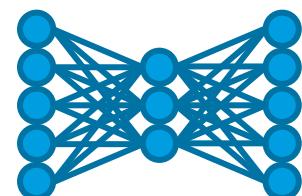
prompt = """You are a bioimage-analysis expert.

What is the best image processing algorithm to segment microscopy images?

Answer the algorithm name only. No explanations.

"""

+ no
image



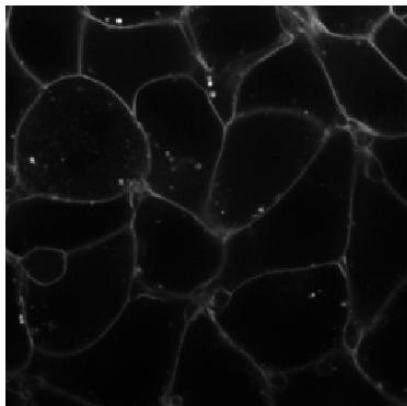
VLMs guessing segmentation algorithms

prompt = """You are a bioimage-analysis expert.

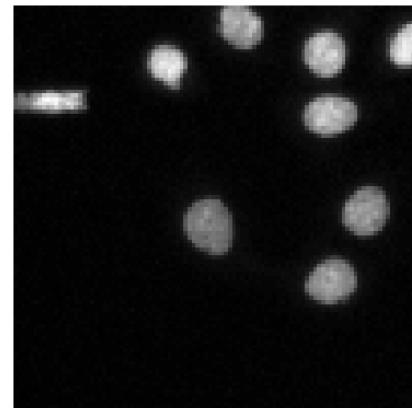
What is the best image processing algorithm to segment this microscopy image?

Answer the algorithm name only. No explanations.

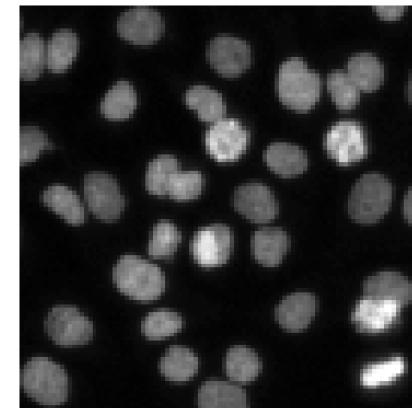
"""



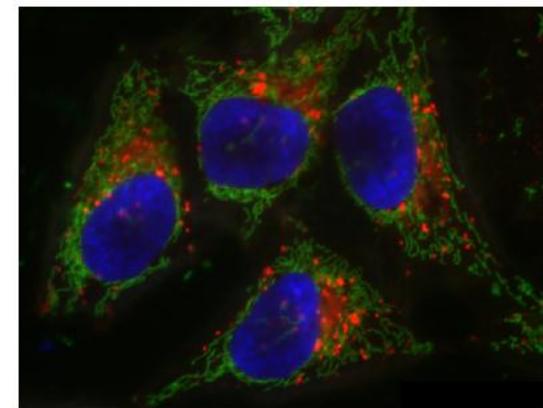
watershed



otsu thresholding
watershed



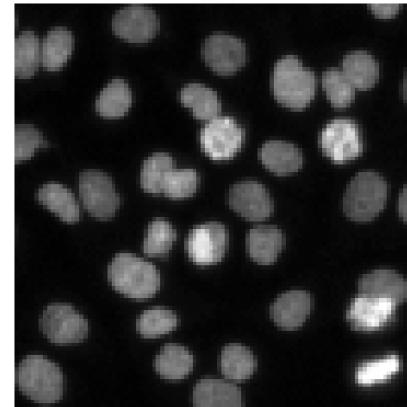
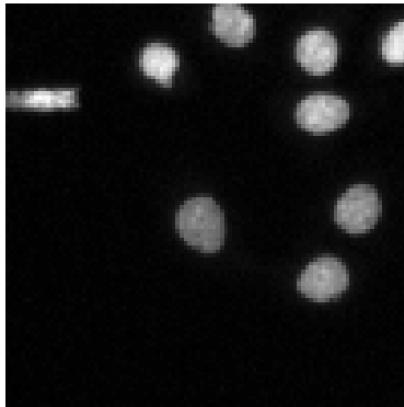
watershed
otsu thresholding



otsu thresholding
learning deep u
cellpose
e watershed
g

Quiz

You have two images and two segmentation algorithms. Assign the algorithms to the images. Explain your choice.



- Voronoi-Otsu-Labeling
- StarDist

VLMs guessing segmentation algorithms

```
prompt = """You are a bioimage-analysis expert. You have a rule-book what algorithms to use for specific images.  
## Rules  
* If an image shows sparse objects such as nuclei, use Otsu-thresholding for segmenting them.  
* If an image shows dense, partially overlapping objects such as nuclei, use StarDist.  
* If an image shows large cell-like structures with bright membranes, use the Watershed algorithm.  
* In case of doubt, use CellPose.  
## The task  
  
What is the best image processing algorithm to segment this microscopy image?  
Answer the algorithm name only. No explanations.  
"""
```

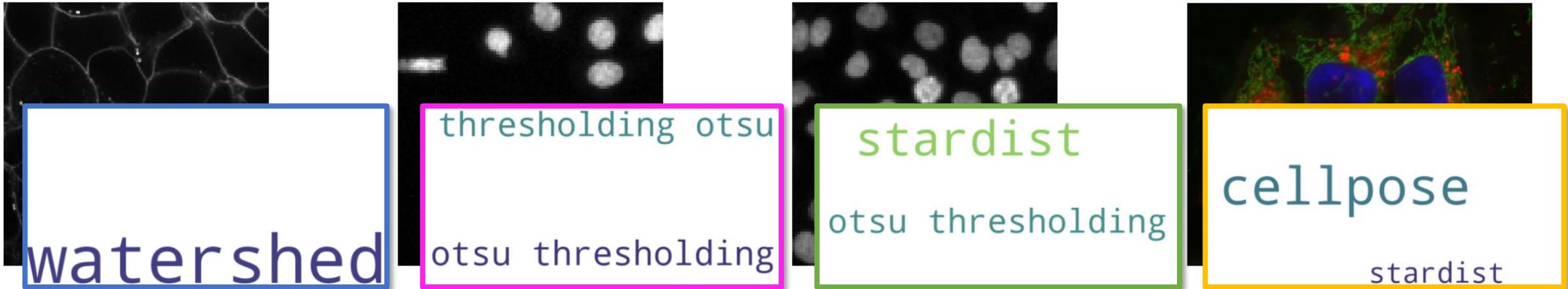
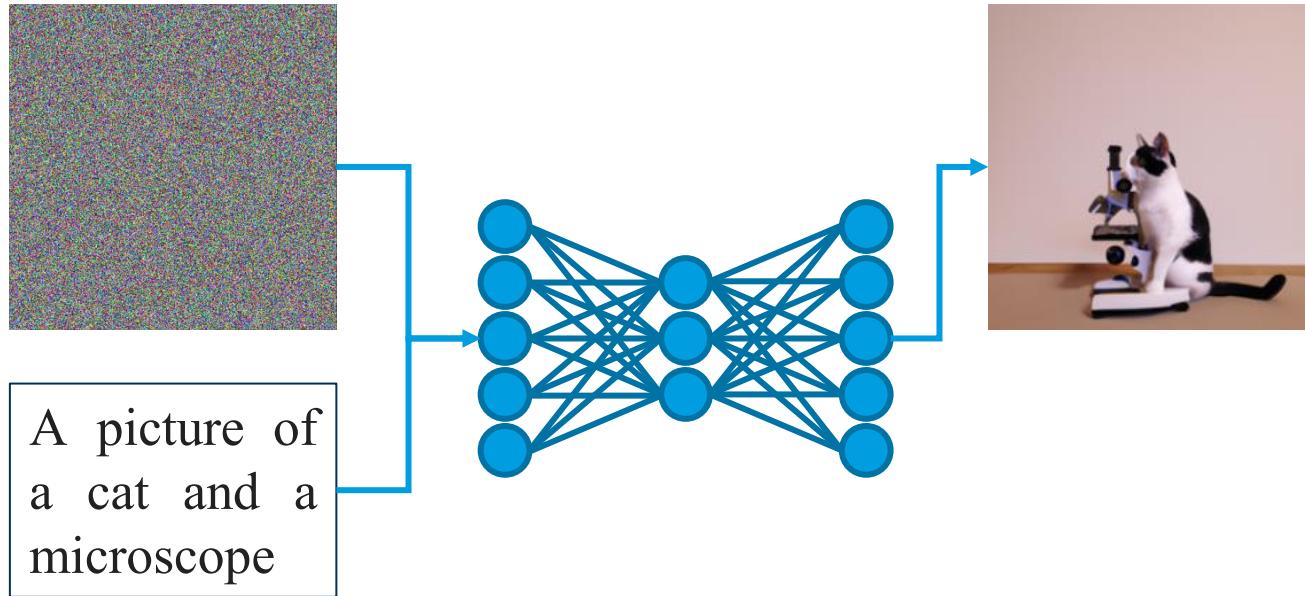


Image Generation

„text-to-image“



Variational Auto-Encoder

„image-to-image“

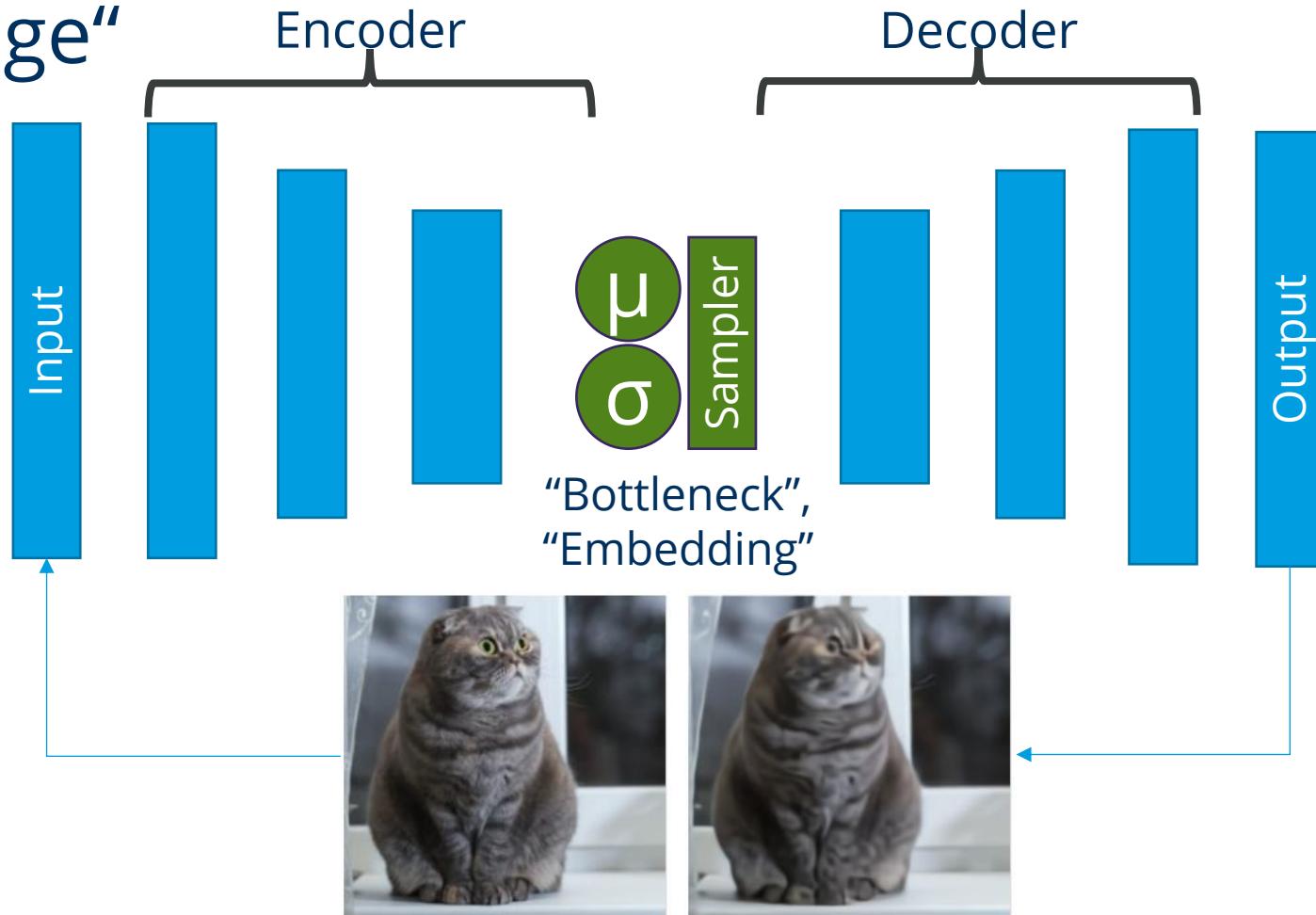
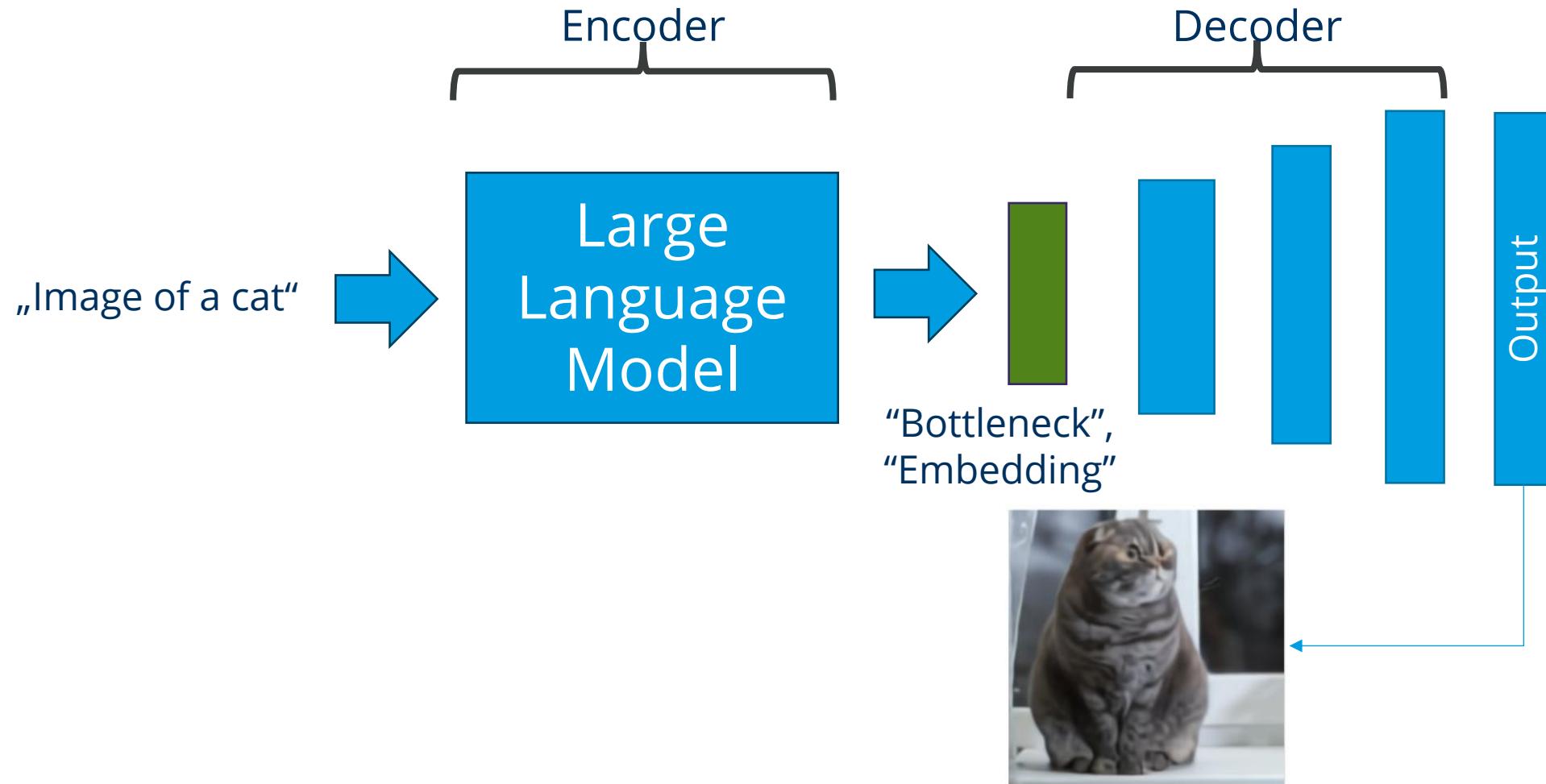


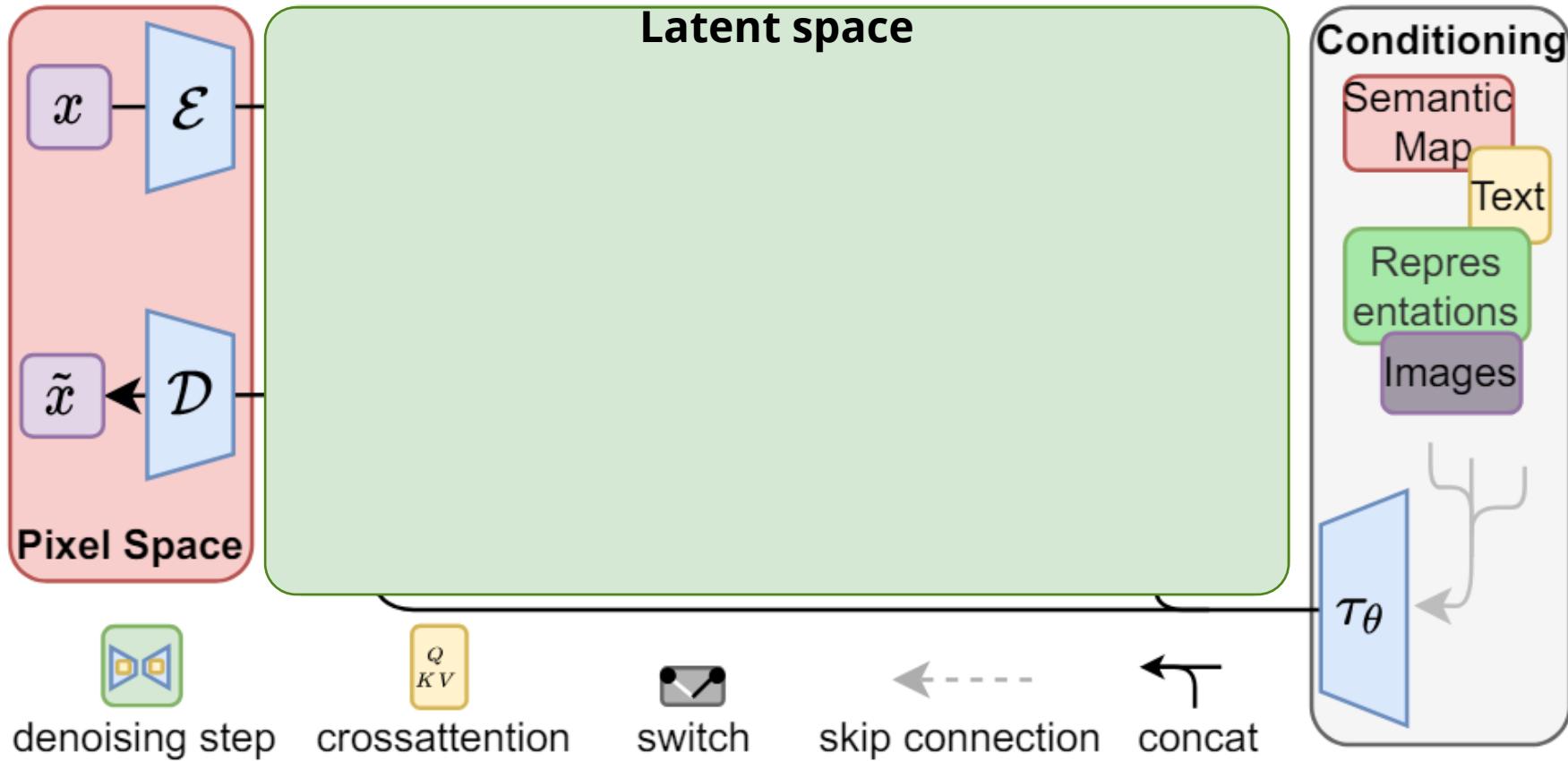
Image Generation



Stable Diffusion

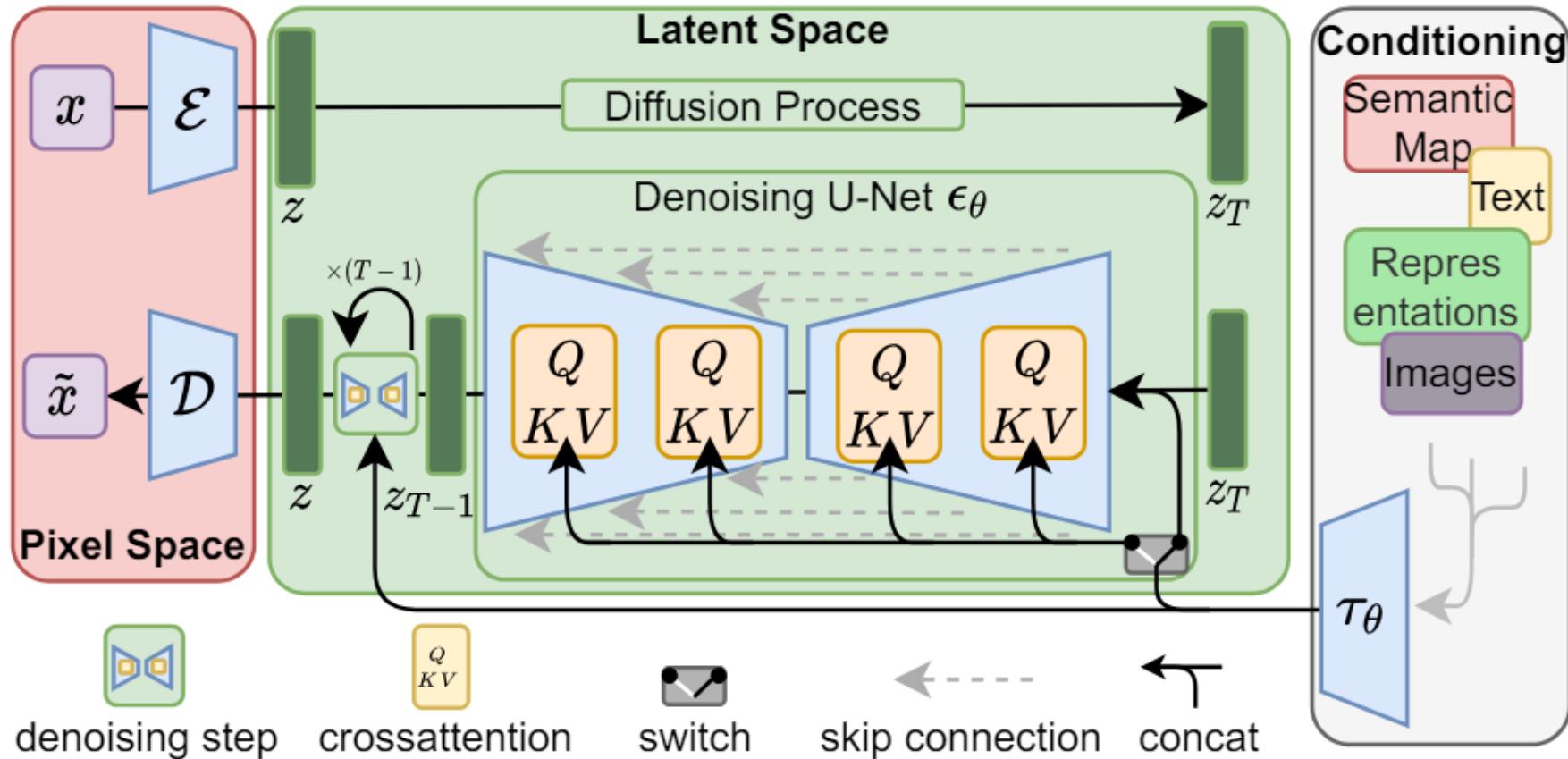
Diffusion: iterative, reverse denoising

Quiz: What is a different word for “latent space”?



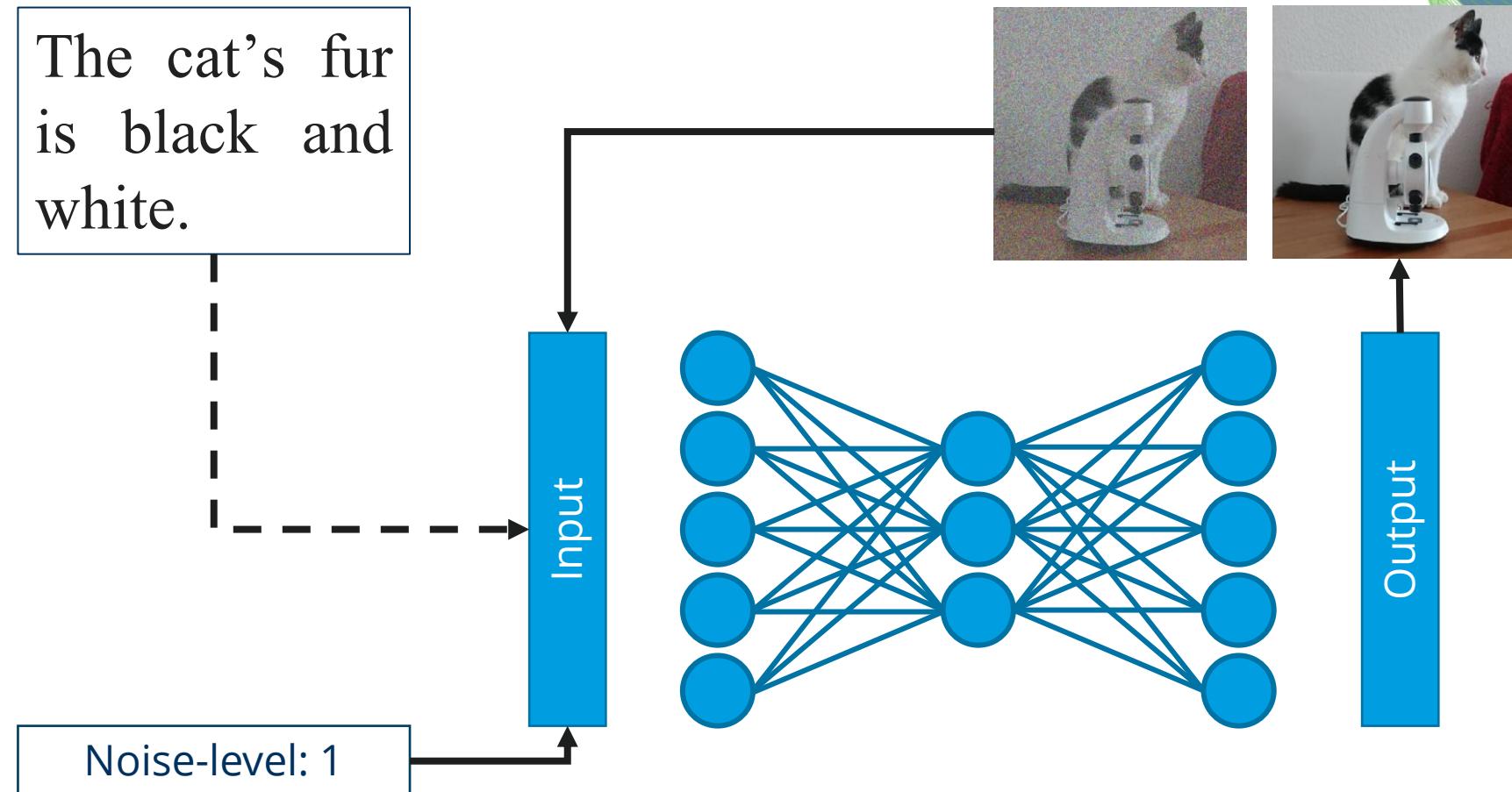
Stable Diffusion

Diffusion: iterative, reverse denoising



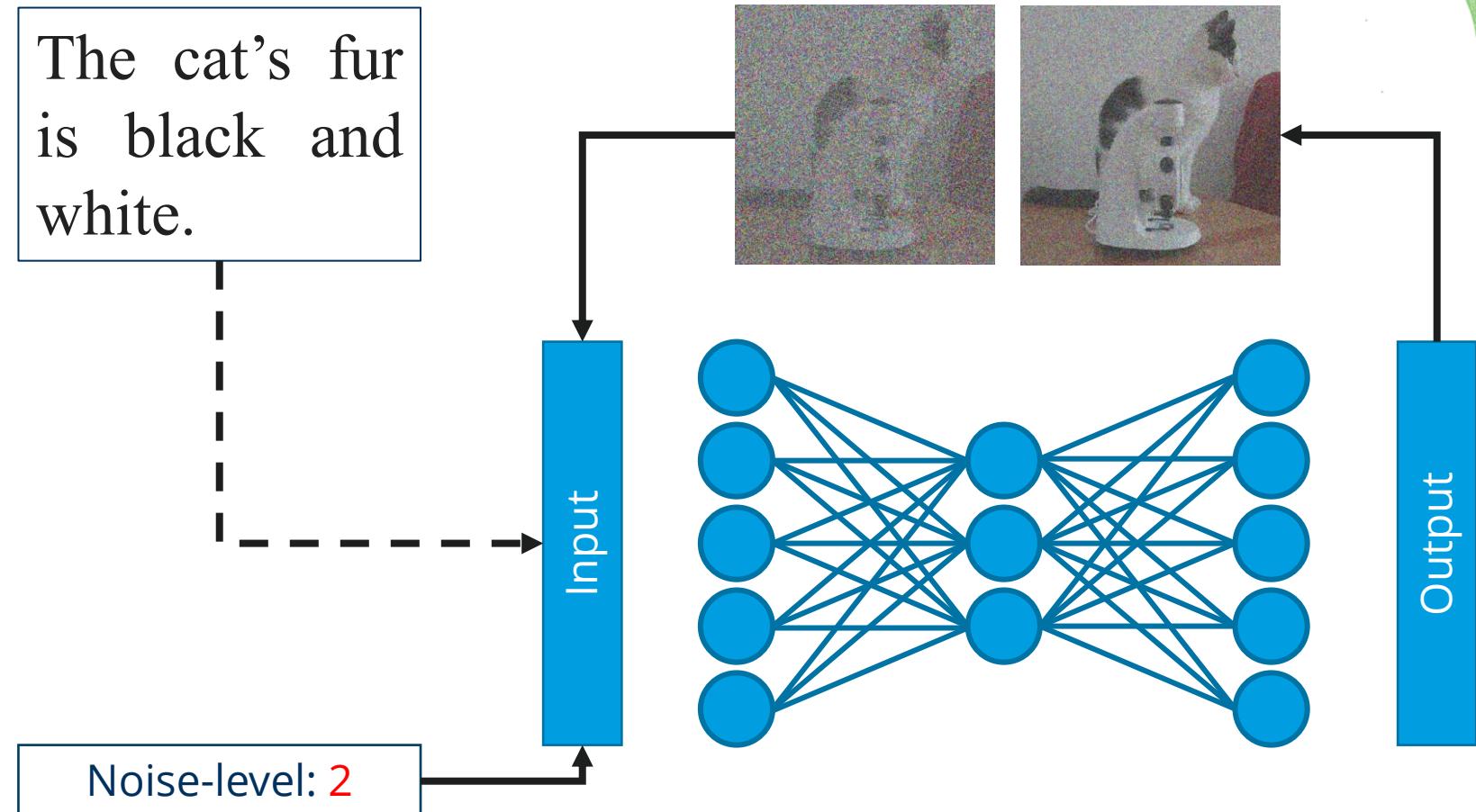
How does it work?

Train a U-Net on
data: image +
noisy image +
description +
noise-level



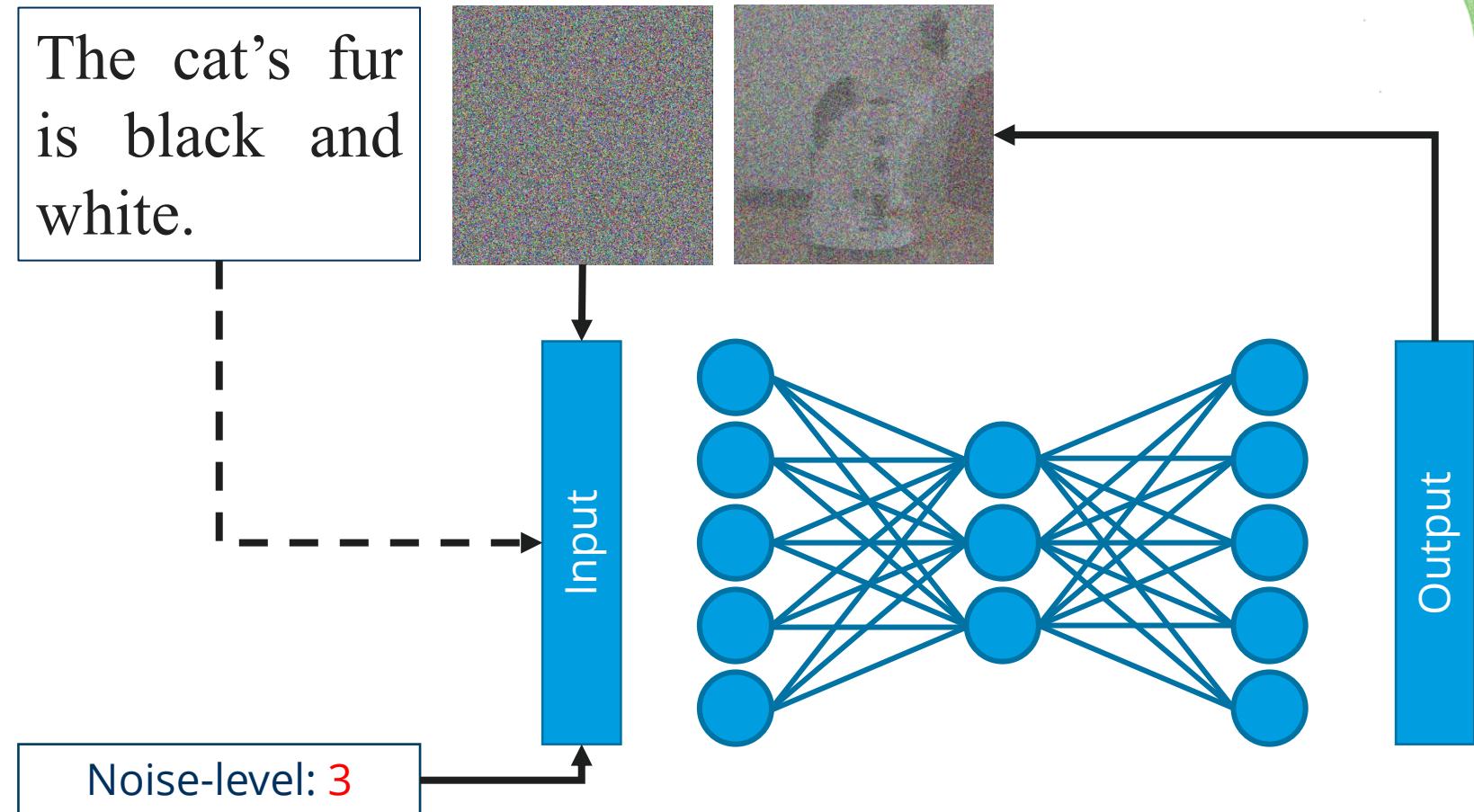
How does it work?

Train a U-Net on
data: image +
noisy image +
description +
noise-level



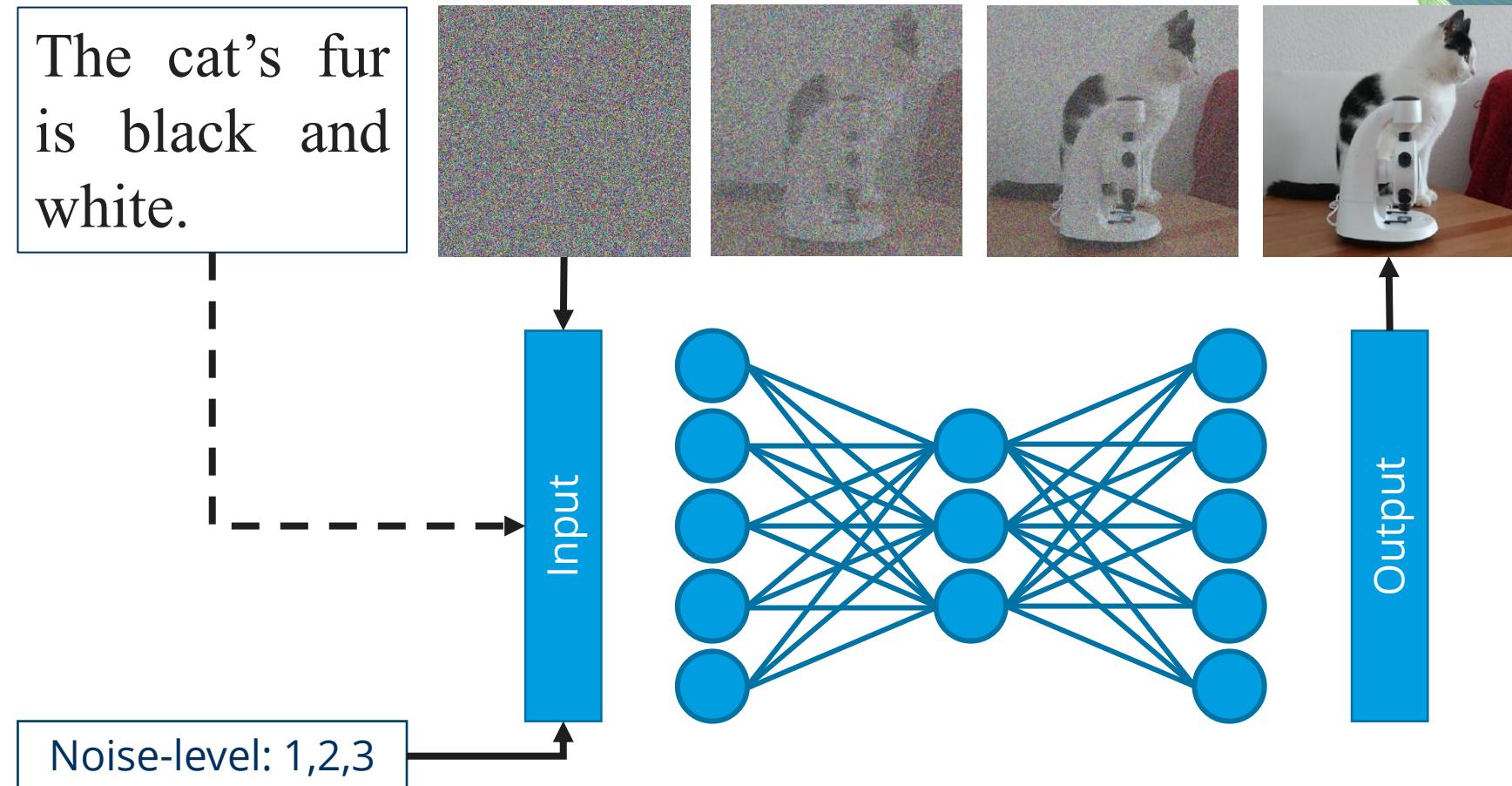
How does it work?

Train a U-Net on
data: image +
noisy image +
description +
noise-level



How does it work?

Prediction is
iterative denoising
of:
Pure noise +
text prompt



How does it work?

Reminder:

- Word embeddings
- Attention

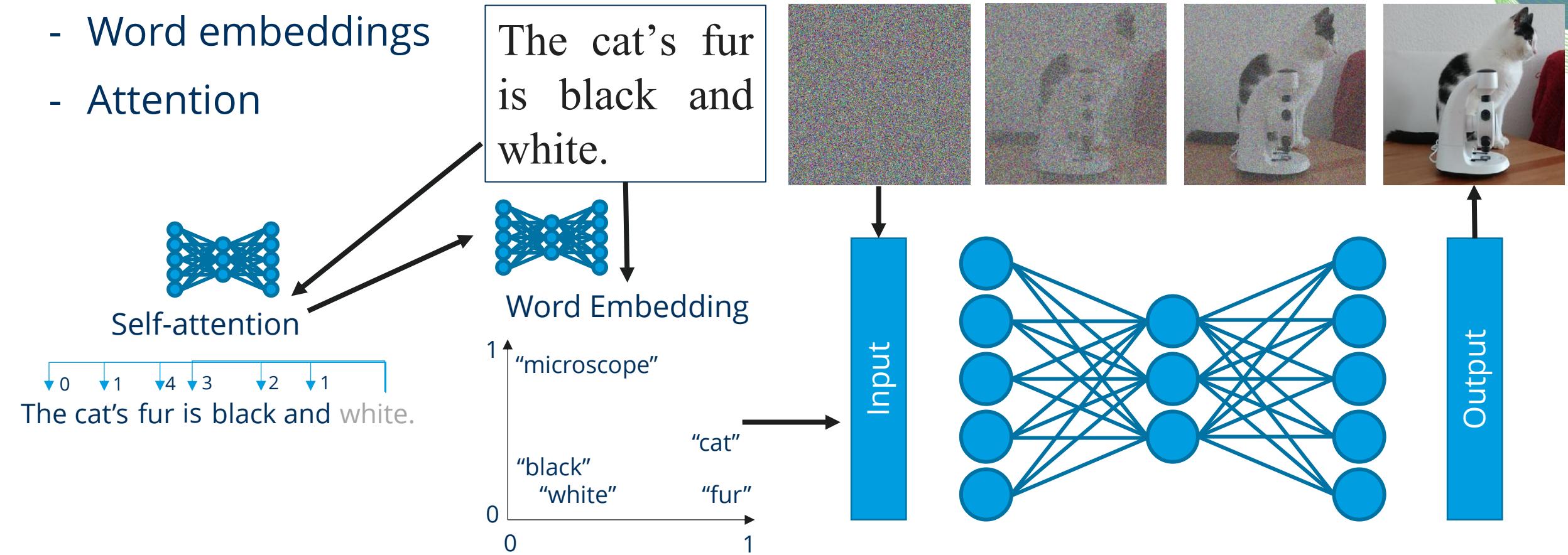


Image generation in Python: Huggingface

Most [Huggingface] image-generation models require a GPU.

```
pipe = DiffusionPipeline.from_pretrained(  
    "stabilityai/stable-diffusion-2-1-base", torch_dtype=torch.float16  
)
```

```
pipe = pipe.to("cuda")
```

Needs
Nvidia GPU

```
prompt = ""  
Draw a realistic photo of an astronaut riding a horse.  
"""
```

```
astronaut = pipe(prompt).images[0]  
astronaut
```

Downloads
4.8 GB

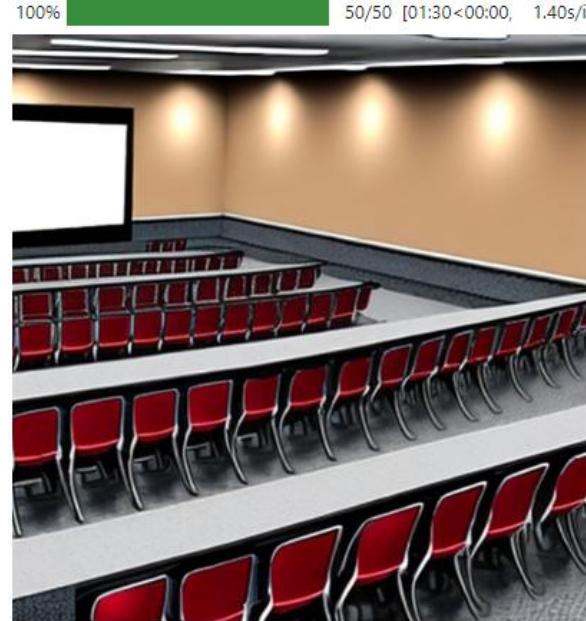
Image generation in Python: Huggingface

Works well if the prompt overlaps with training data, potentially huge variation between attempts

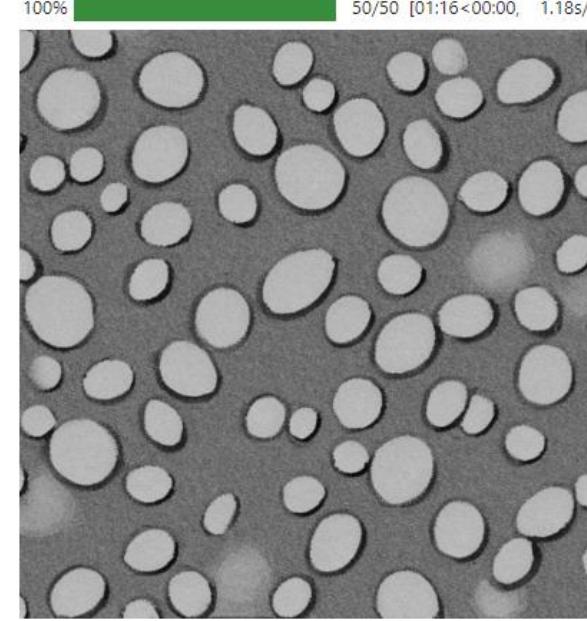
```
prompt = ""  
Draw a realistic photo of an astronaut riding a horse.  
"""  
  
astronaut = pipe(prompt).images[0]  
astronaut
```



```
prompt = ""  
Draw a realistic photo of a lecture hall with an  
ongoing lecture about vision language models.  
"""  
  
photo = pipe(prompt).images[0]  
photo
```



```
prompt = ""  
Draw a greyscale picture of sparse bright blobs on dark  
background. Some of the blobs are roundish, some are a  
bit elongated.  
"""  
  
image = pipe(prompt).images[0]  
image
```



```
image = pipe(prompt,  
            num_inference_steps=10,  
            width=512,  
            height=512).images[0]  
  
image
```

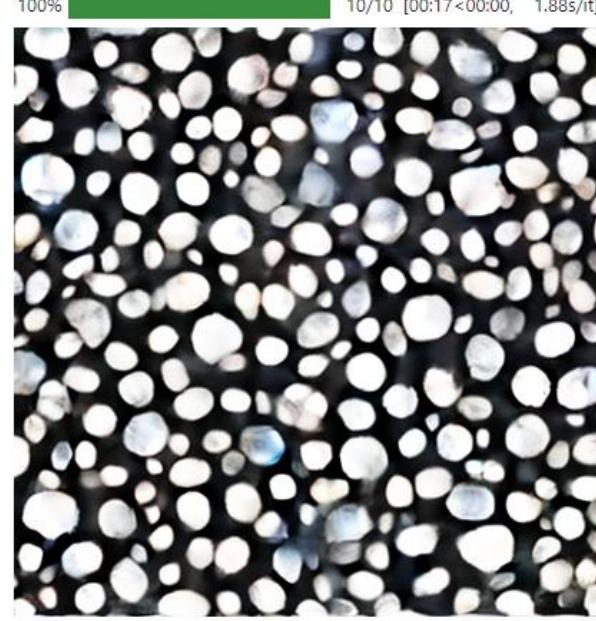


Image generation in Python: Dall-E

No need for a GPU, but costs



```
def prompt_image(message:str, width:int=1024, height:int=1024, model='dall-e-3'):
    client = openai.OpenAI()
    response = client.images.generate(
        prompt=message,
        model=model,
        n=1,
        size=f"{width}x{height}"
    )
    image_url = response.data[0].url
    image = imread(image_url)

    return image
```

Works with
Dall-E 2 and 3

Image Generation LLMs

Challenges: fake-images

Interesting challenges for our community ahead

a histology image of lung cancer cells and some healthy tissue

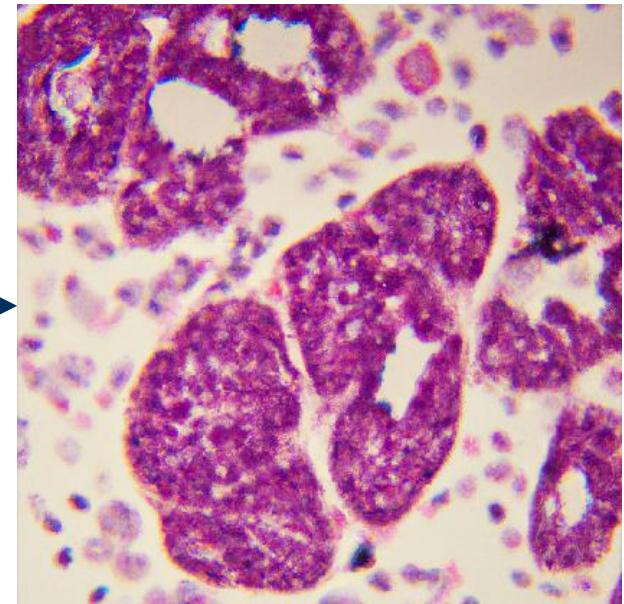
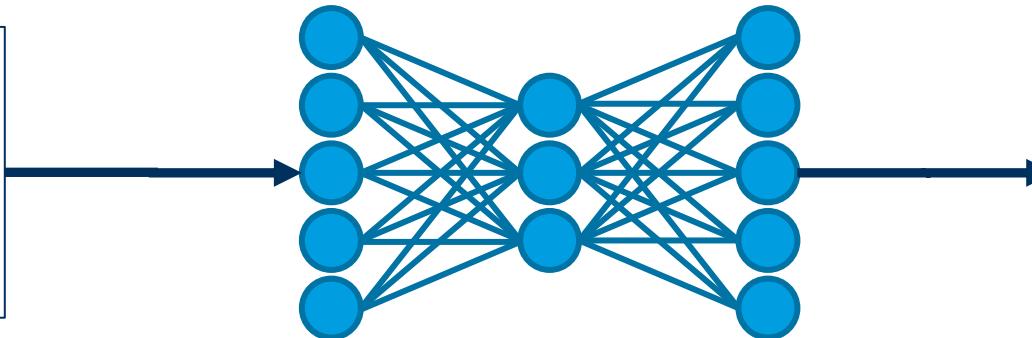
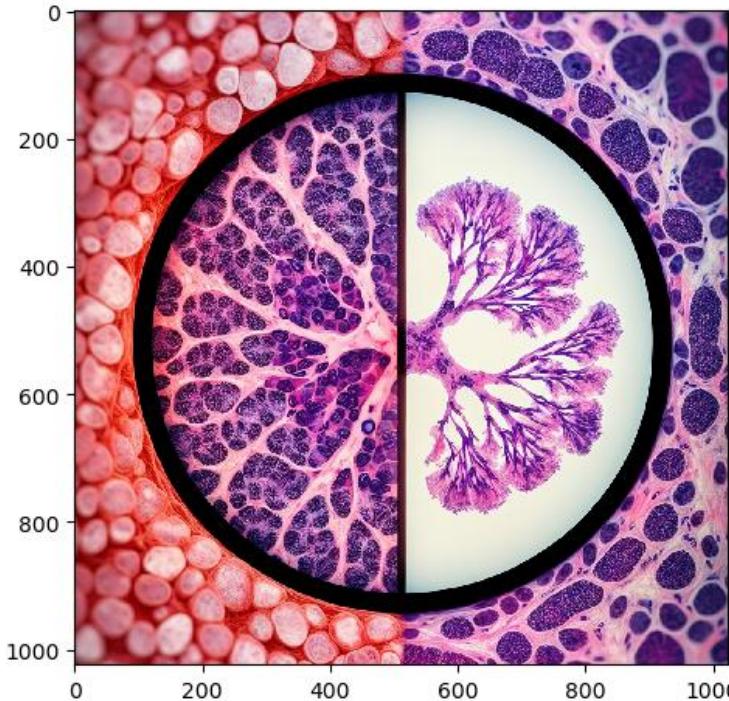


Image generation in Python: Dall-E

Is Dall-E 2 more capable of creating realistic microscopy images than Dall-E 3?

```
histology = prompt_image('a histology image of lung cancer cells and some healthy tissue')
imshow(histology)
```

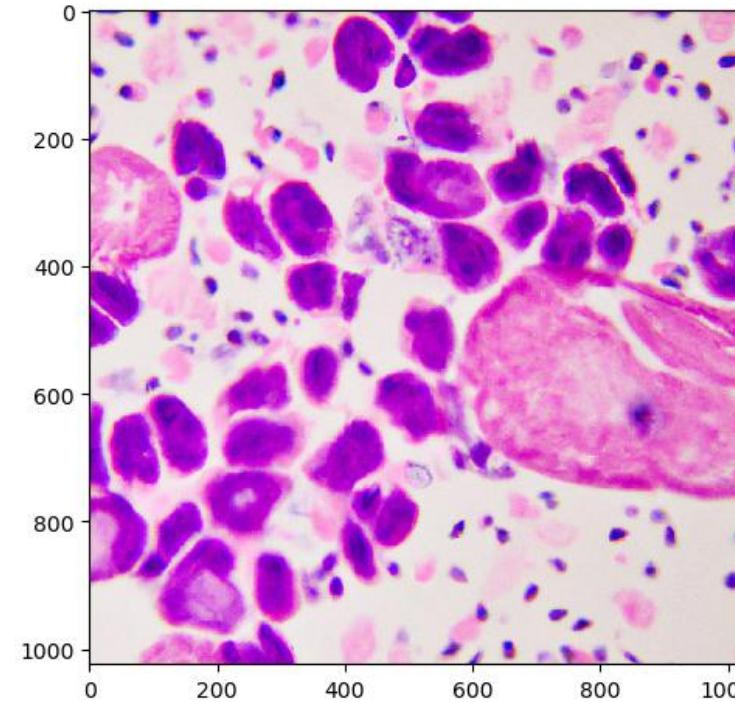
```
<matplotlib.image.AxesImage at 0x1d9eda5bd90>
```



Dall-E 3

```
histology = prompt_image('a histology image of lung cancer cells and some healthy tissue',
                         model='dall-e-2')
imshow(histology)
```

```
<matplotlib.image.AxesImage at 0x1d9edac6fd0>
```



Dall-E 2

Image Generation LLMs

Challenges: physical / physiological hallucinations

Prompt: "Draw a close-up picture of people's hands"

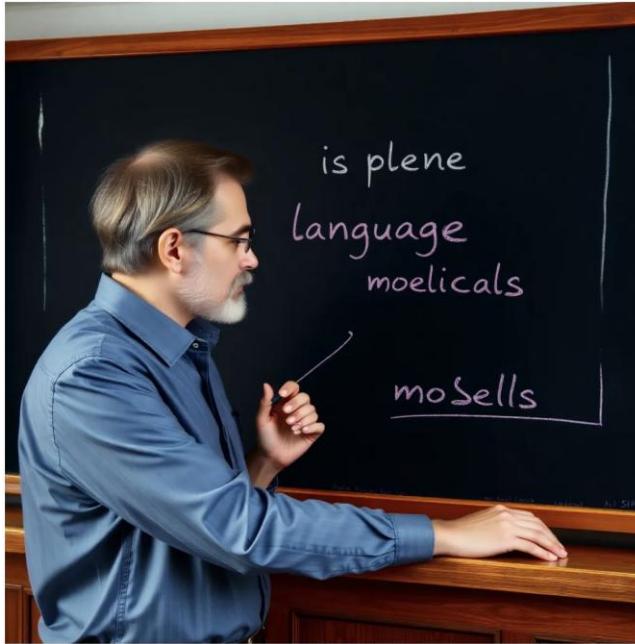


Source: <https://replicate.com/black-forest-labs/flux-schnell-lora>

Image Generation LLMs

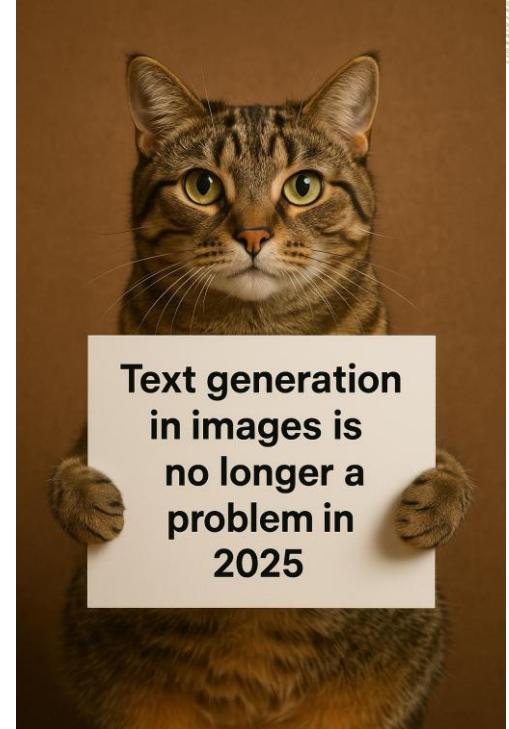
Challenges: Text in images

Prompt: 'please draw a picture of a professor writing "large language models" on a blackboard'



Source: <https://replicate.com/black-forest-labs/flux-schnell-lora>

Prompt: 'Please draw a picture of a cat holding a sign with the text "Text generation in images is no longer a problem in 2025".'



<https://chatgpt.com/share/685e2ea5-42f4-8001-9c56-6e3c6de2464e>

Image Generation LLMs

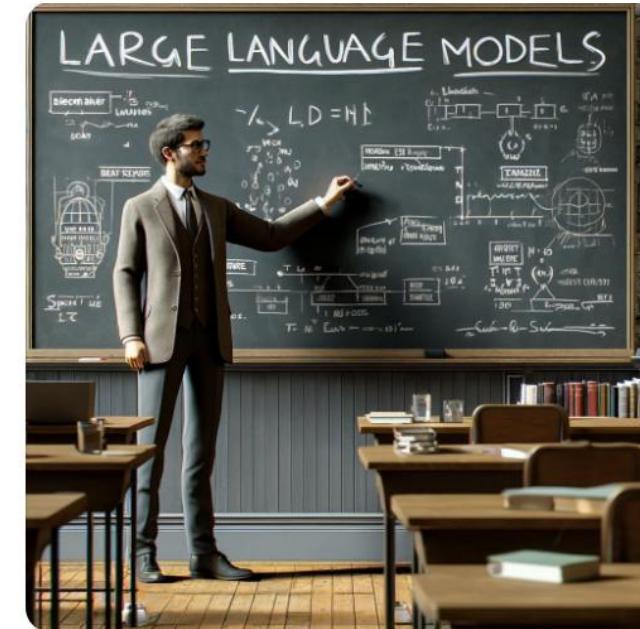
Challenges: Bias

Prompt: "a nice photo of a human"



Source: <https://replicate.com/stability-ai/stable-diffusion>

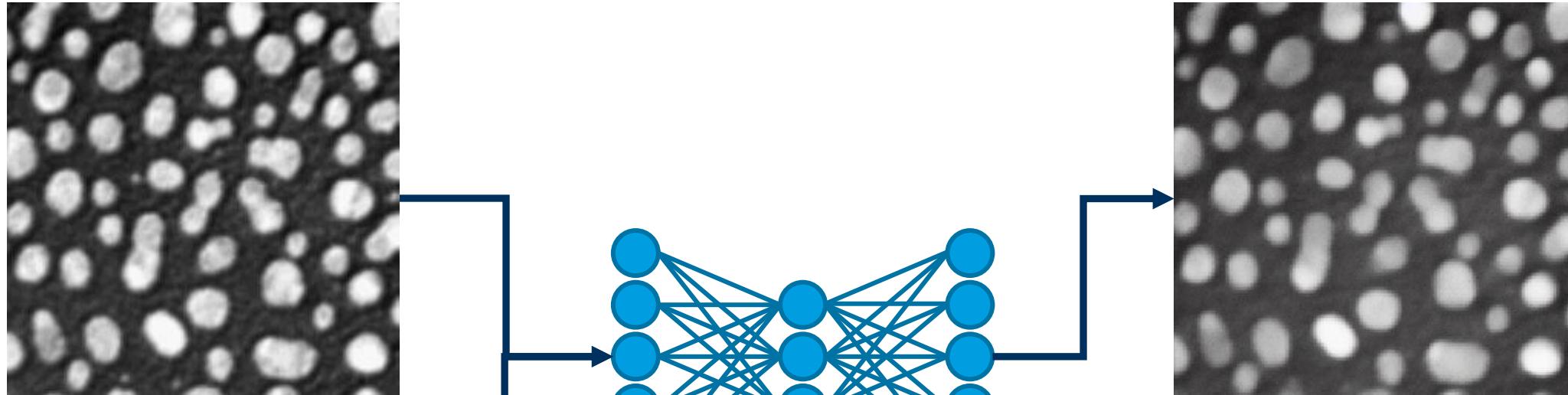
Prompt: 'please draw a picture of a professor writing "large language models" on a blackboard'



Source: <https://chatgpt.com/>

Image Generation LLMs

Image-to-image prompting, image variations



Blur the image

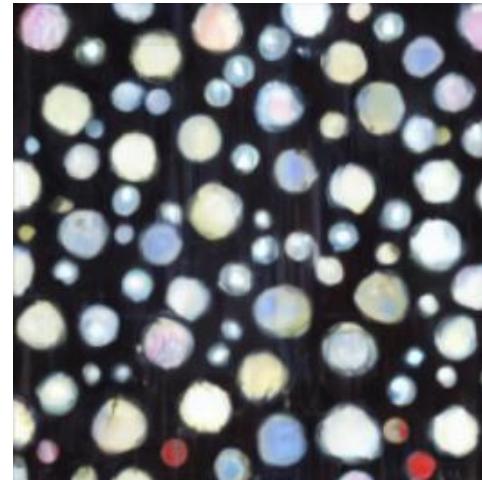
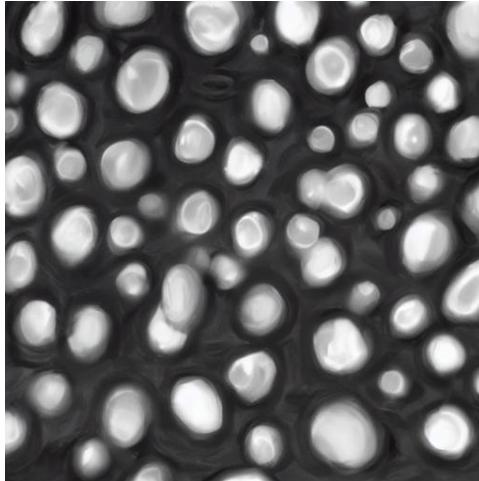
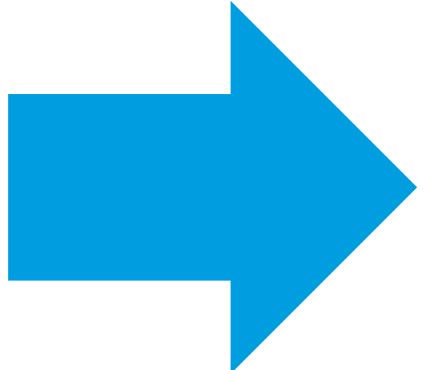
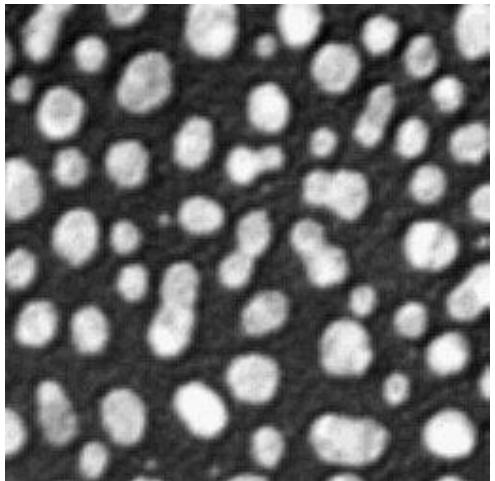
InstructPix2Pix: Learning to Follow Image Editing Instructions

Tim Brooks* Aleksander Holynski* Alexei A. Efros

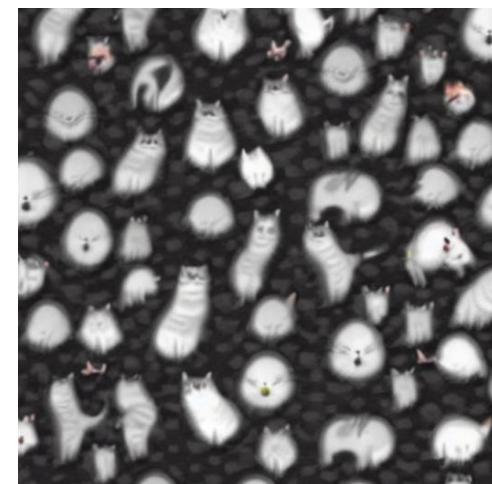
University of California, Berkeley

Image variation

Generate images, e.g. for augmenting data



Potentially useful to make algorithms more robust



Inpainting

Replacing regions in images
(also „Gap-filling“, „Replacing“)

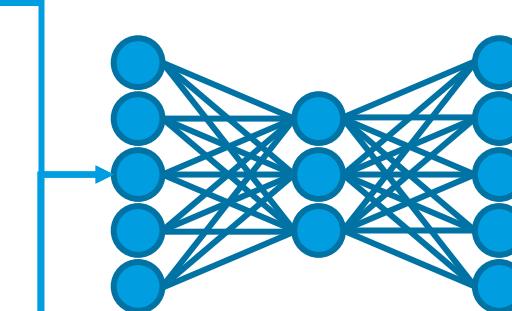
Raw image



Mask image



A black white
cat fur



Manipulated
image



Interesting
challenges for our
community ahead

Inpainting in Python: Huggingface

```
pipe = StableDiffusionInpaintPipeline.from_pretrained(  
    "stabilityai/stable-diffusion-2-inpainting",  
    torch_dtype=torch.float16  
)  
pipe = pipe.to("cuda")
```

Downloads
4.8 GB

Needs
Nvidia GPU



```
prompt = "A black white cat fur"  
image = pipe(prompt=prompt,  
            image=init_image,  
            mask_image=mask_image,  
            num_inference_steps=50,  
            width=512,  
            height=512,  
            num_images_per_prompt=1,  
            ).images[0]
```



Inpainting in Python: Huggingface

Check out the *model cards* online in the Huggingface hub.

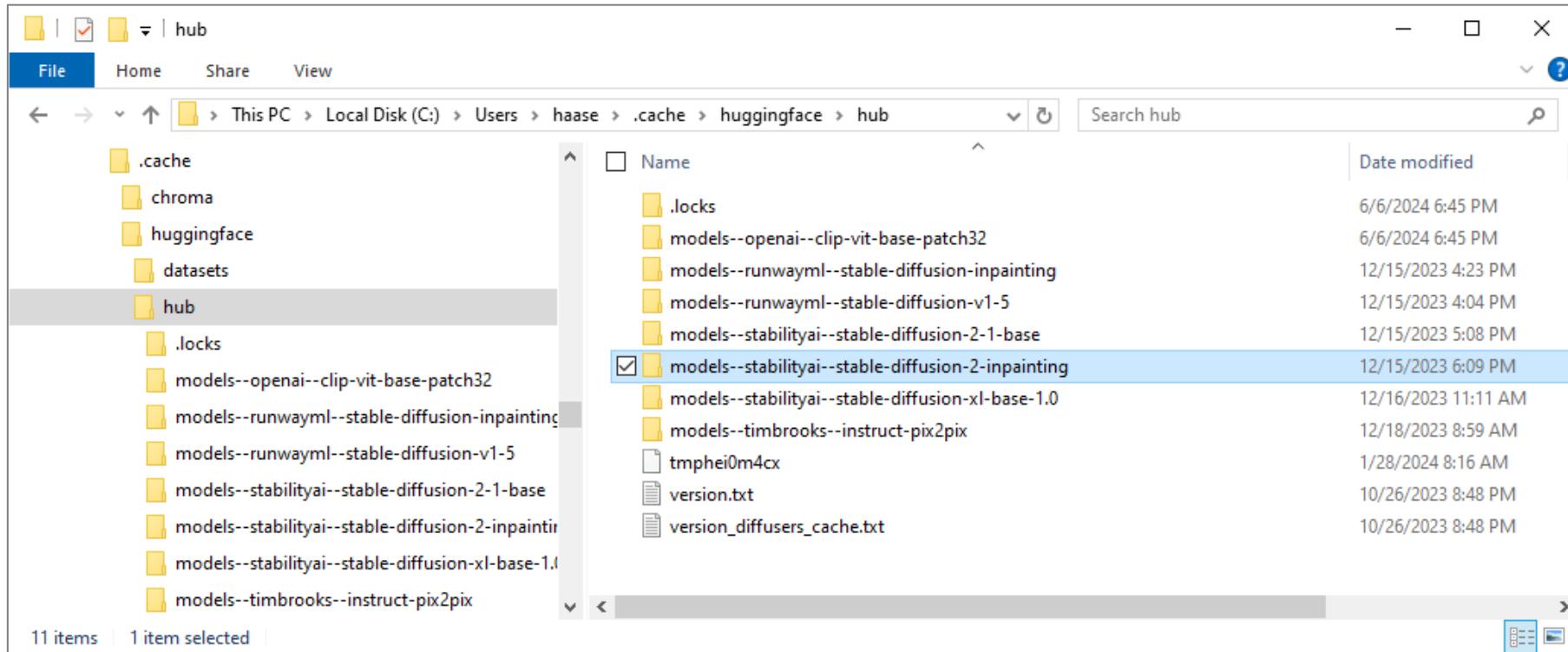
```
pipe = StableDiffusionInpaintPipeline.from_pretrained(  
    "stabilityai/stable-diffusion-2-inpainting",  
    torch_dtype=torch.float16  
)  
pipe = pipe.to("cuda")
```



The screenshots illustrate the Huggingface model card interface. The first screenshot shows the Python code for initializing the pipeline. The second screenshot shows the 'Model Details' section, which includes developer information (Robin Rombach, Patrick Esser), model type (Diffusion-based text-to-image generation model), language (English), license (CreativeML Open RAIL++-M License), and a detailed 'Model Description' block. The third screenshot shows the 'Examples' section, which provides a command-line example and a code snippet for running Stable Diffusion 2 inpainting.

Inpainting in Python: Huggingface

You find the downloaded models cached in your home directory
They are big! Clean up here from time to time.



Inpainting in Python: Dall-E

No need for a GPU, but costs



```
client = OpenAI()

response = client.images.edit(
    image=numpy_to_bytestream(resized_image_rgb),
    mask=numpy_to_bytestream(masked_rgba),
    prompt=prompt,
    n=1,
    size=f"{image_width}x{image_height}",
    model=model
)
```

2D RGB images
only

Supported: 256,
512, 1024 pixels

Size must match

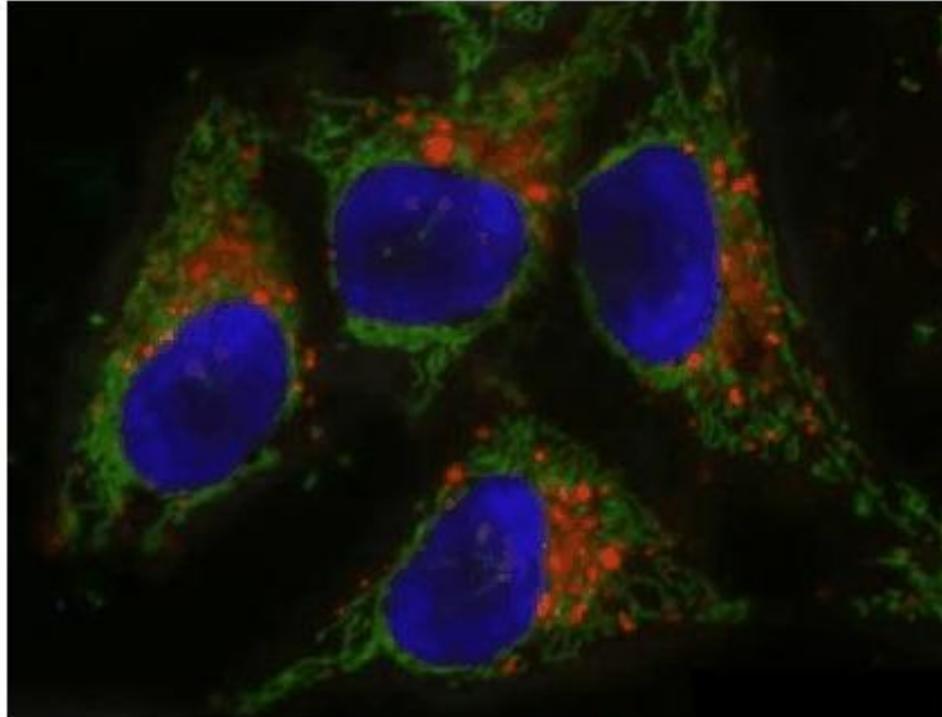


Result: List of URL(s)

New technologies bring new risks...

If you can generate images,
you can also generate parts of images....

[6]:



Curtain



672

Interesting
challenges for our
community ahead

Image manipulation detection

The noise pattern differs between raw and processed images...

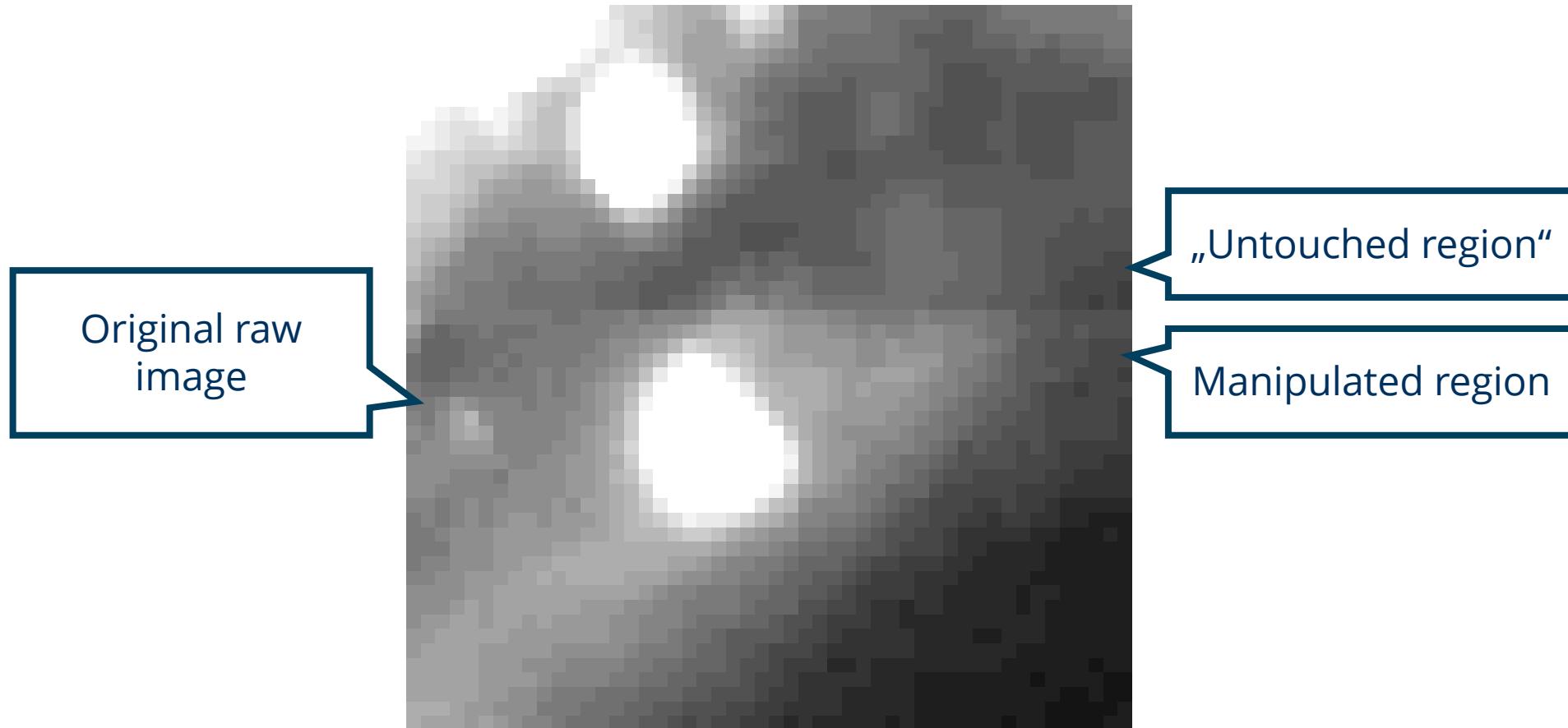
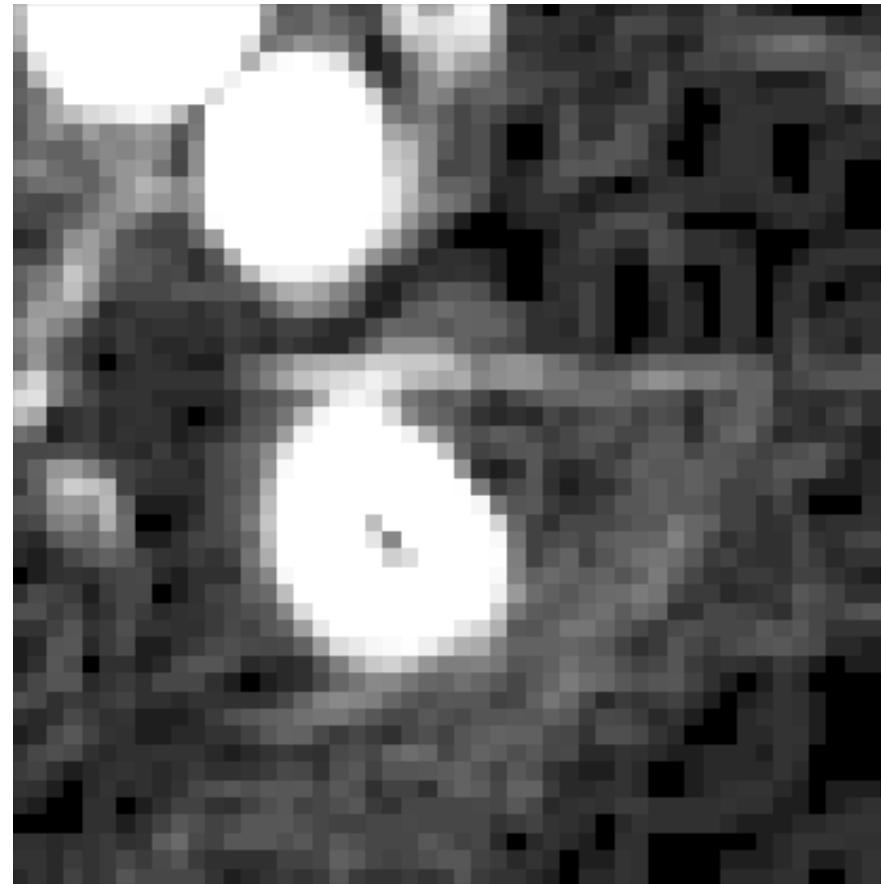


Image manipulation detection

e.g. by studying noise-patterns

Local standard deviation filter



„Untouched region“

Manipulated region

Image manipulation detection

e.g. by studying noise-patterns

Sobel filter



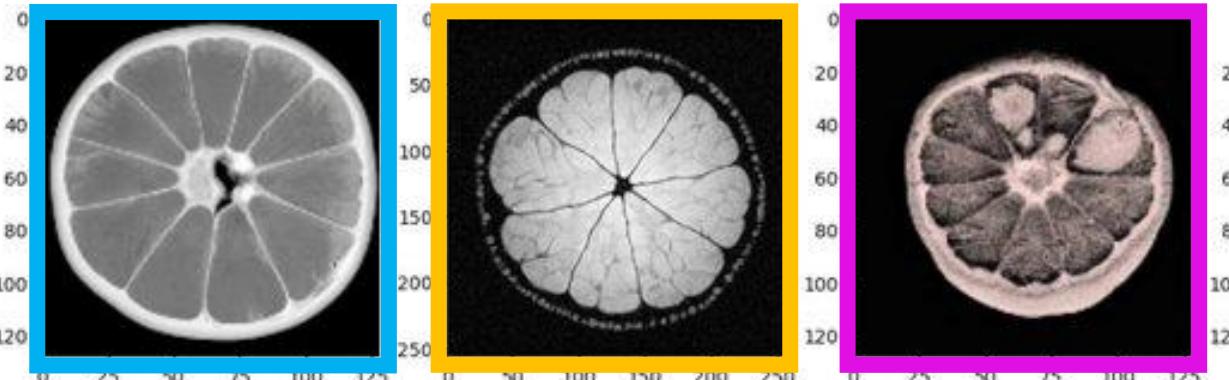
„Untouched region“

Manipulated region

Benchmarking image generation

When asking humans to evaluate results, make sure they are the right target audience

```
mri_prompt = """"  
A single, high resolution, black-white image of  
a realistically looking orange fruit slice  
imaged with T2-weighted magnetic resonance imaging (MRI).  
"""
```



Robert Haase @haesleinhuepf · 1h

Fun poll time! Which of these images shows a real MRI image of an orange? (Credits: licensed CC-BY 4.0 by Alexandr Khrapichev, University of Oxford; the other images were generated by @openai's DALL-E)

Please vote below, RT and if you can explain why, please comment! 😊

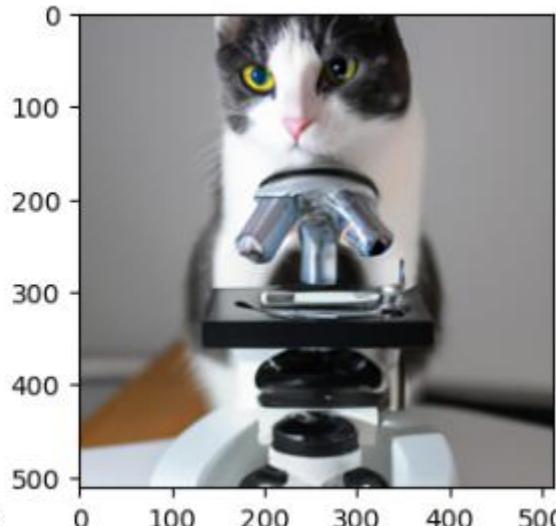
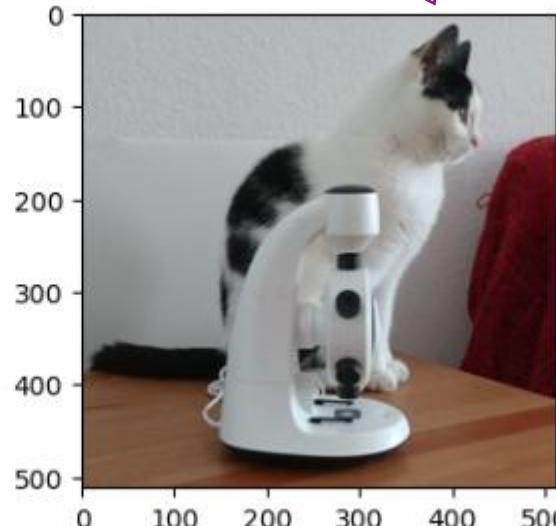
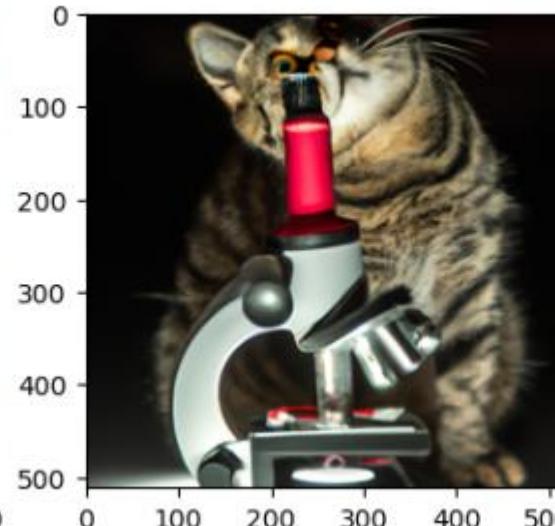
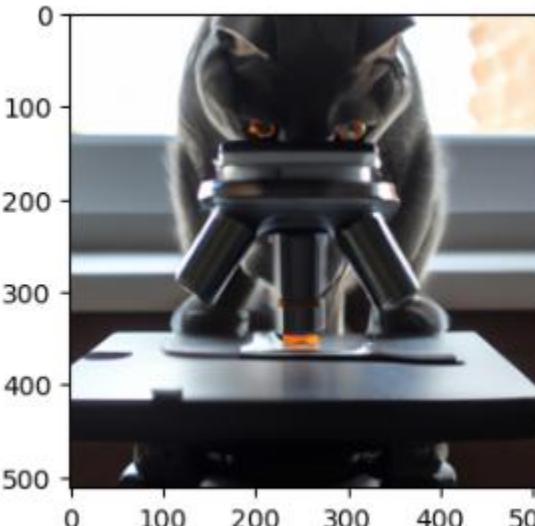
Option	Percentage
1	24.6%
2	56.1%
3	7%
4	12.3%

57 votes · 18 hours left

Benchmarking image generation

Prompt engineering to optimize images

```
cat_microscope_prompt = """  
Image of a cat sitting behind a microscope.  
"""
```



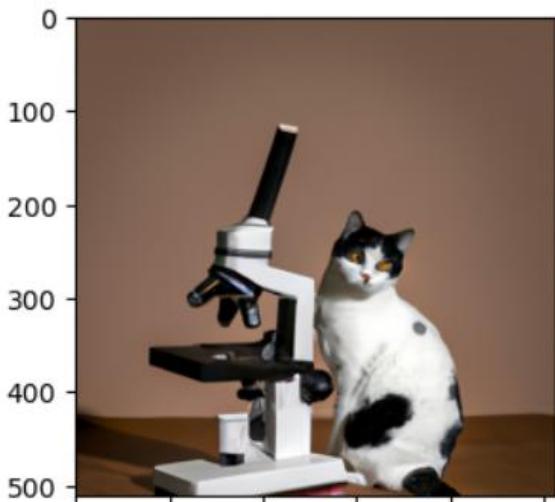
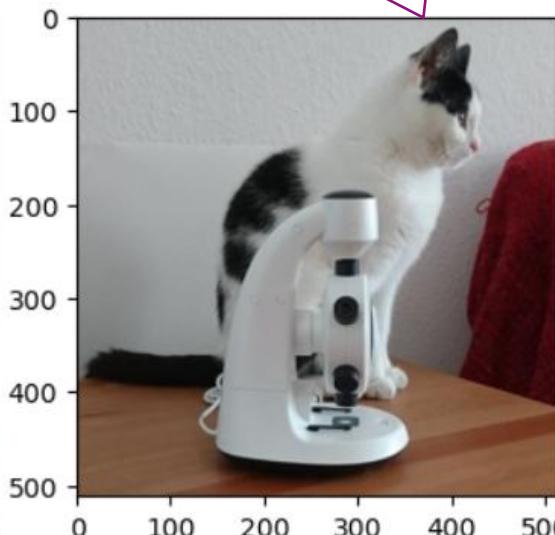
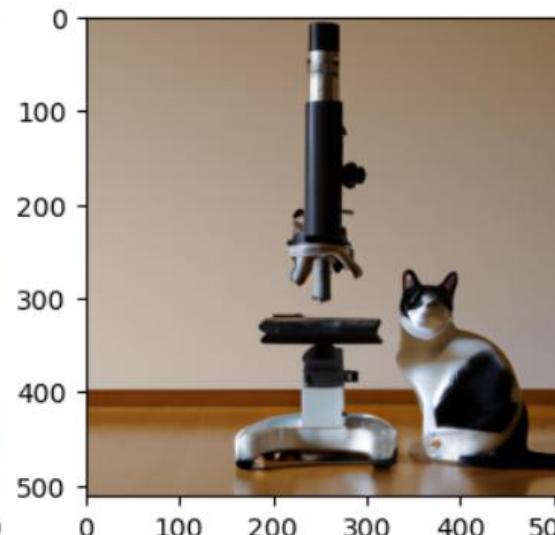
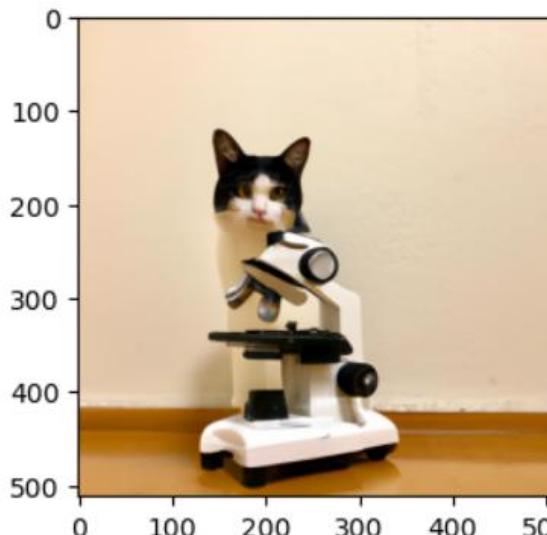
One cat
is real.

Benchmarking image generation

Prompt engineering to optimize images

```
[5]: cat_microscope_prompt = """  
Image of a cat sitting behind a microscope.  
Both are on a brown floor in front of a white wall.  
The cat is mostly white and has some black dots.  
The cat sits straight.  
The cat is a bit larger than the microscope.  
"""
```

One cat
is real.



CLIP scores

Contrastive Language-Image Pre-Training (CLIP)

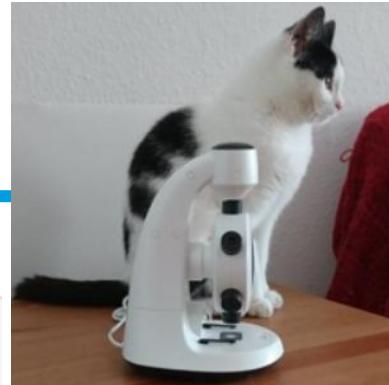
- For image describing

Here: Similarity between image and prompt

```
from torchmetrics.multimodal.clip_score import CLIPScore  
metric = CLIPScore(model_name_or_path="openai/clip-vit-base-patch16")
```

```
score = metric(torch.as_tensor(image), "cat")  
score.detach()
```

```
tensor(25.3473)
```



```
score = metric(torch.as_tensor(image), "microscope")  
float(score.detach())
```

```
30.786287307739258
```

CLIP scores

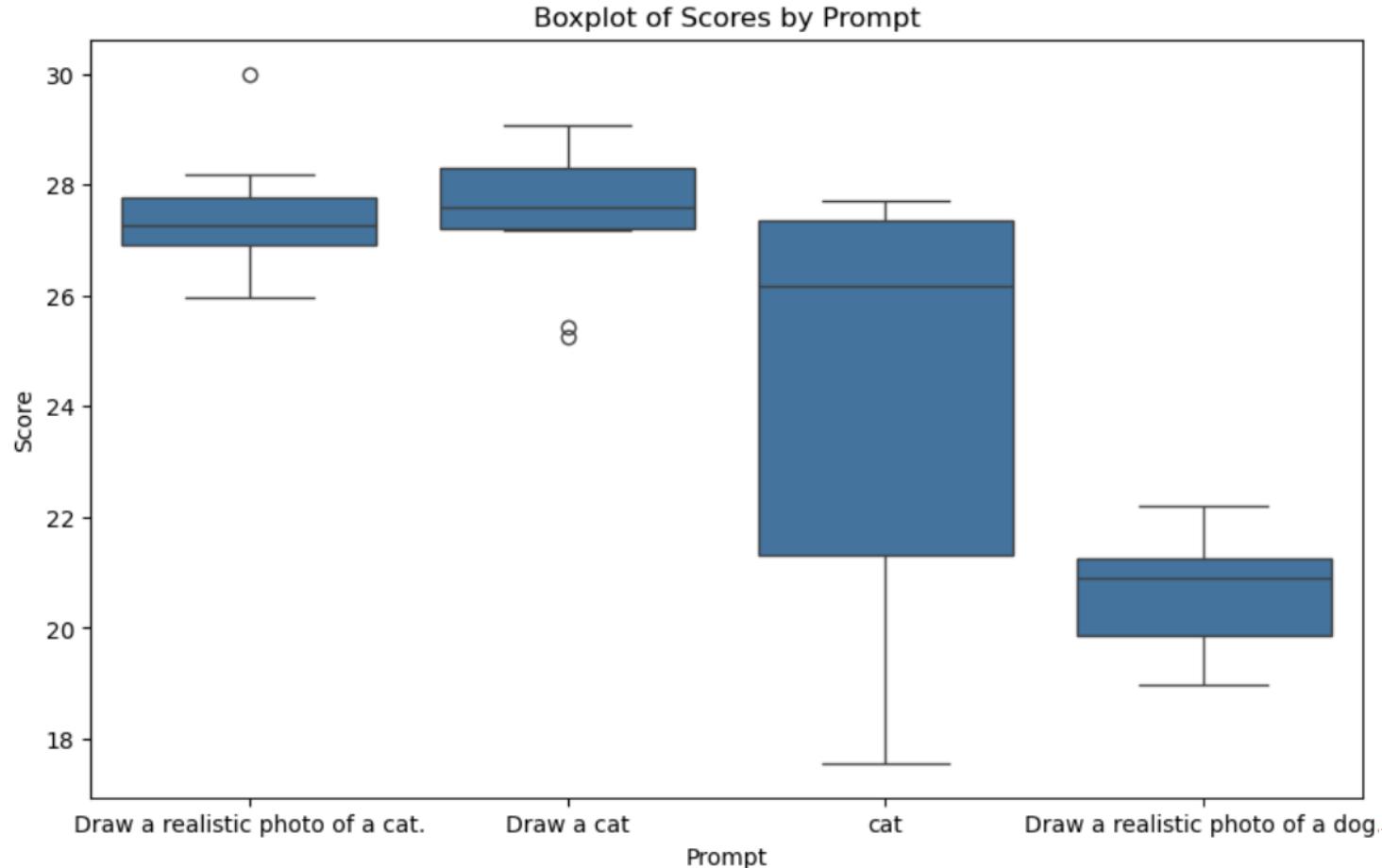
Example: Prompt optimization

```
num_attempts = 10
prompts = ["Draw a realistic photo of a cat.",
           "Draw a cat",
           "cat",
           "Draw a realistic photo of a dog."]

data = {"prompt": [],
        "score": []}
for prompt in prompts:
    for i in range(num_attempts):
        image = pipe(prompt, disable_tqdm=True).images[0] → Image generation
        score = metric(torch.as_tensor(np.array(image)), "cat") → Quality measurement
        data["score"].append(float(score.detach()))
        data["prompt"].append(prompt)
```

CLIP scores

Example: Prompt optimization



Always have a control experiment!

Summary

- Multi-modal models combine image, text and [...] data
- Combination of pre-existing model architectures
- APIs not standardized (yet)

Trade-off:

- LLM capacity / size limited
- Multi-modal LLMs presumably perform worse compared to more specialized models
 - Example: Ask a vision-language model to write code
- High potential for mixture-of-experts models and multi-agent systems using small specialized language models

Exercises

Robert Haase

Funded by



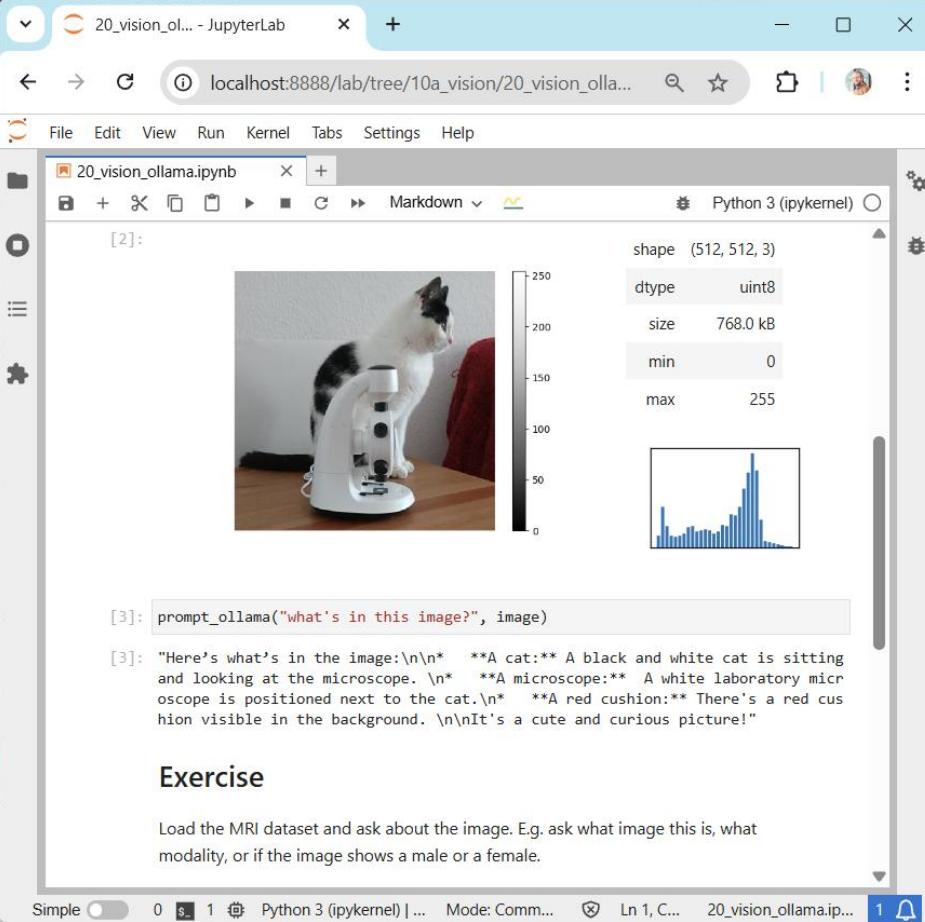
Bundesministerium
für Bildung
und Forschung



Diese Maßnahme wird gefördert durch die Bundesregierung
aufgrund eines Beschlusses des Deutschen Bundestages.
Diese Maßnahme wird mitfinanziert durch Steuermittel auf
der Grundlage des von den Abgeordneten des Sächsischen
Landtags beschlossenen Haushaltes.

Vision language models

Elaborate on the limits of VLMs



localhost:8888/lab/tree/10a_vision/20_vision_olla...

File Edit View Run Kernel Tabs Settings Help

20_vision_ollama.ipynb

[2]:



shape (512, 512, 3)
dtype uint8
size 768.0 kB
min 0
max 255

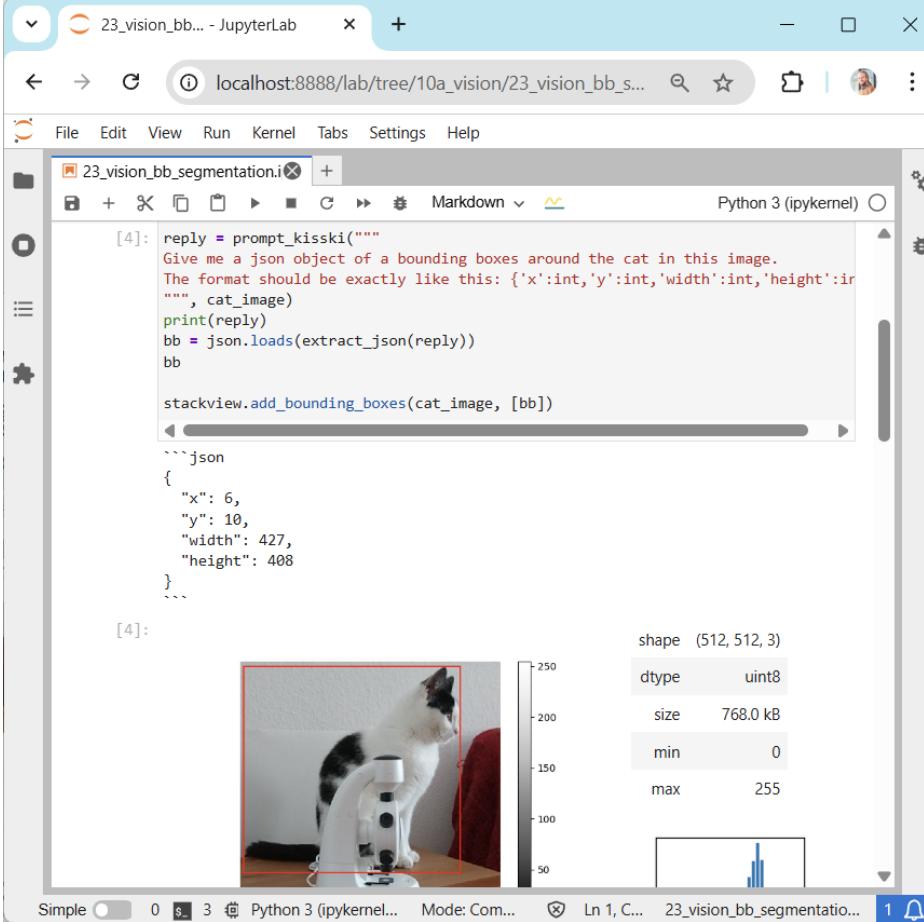
[3]: prompt_ollama("what's in this image?", image)

[3]: "Here's what's in the image:\n* **A cat:** A black and white cat is sitting and looking at the microscope.\n* **A microscope:** A white laboratory microscope is positioned next to the cat.\n* **A red cushion:** There's a red cushion visible in the background.\n\nIt's a cute and curious picture!"

Exercise

Load the MRI dataset and ask about the image. E.g. ask what image this is, what modality, or if the image shows a male or a female.

Simple 0 Python 3 (ipykernel) | ... Mode: Comm... Ln 1, C... 20_vision_ollama.ip... 1



localhost:8888/lab/tree/10a_vision/23_vision_bb_se...

File Edit View Run Kernel Tabs Settings Help

23_vision_bb_segmentation.ipynb

[4]: reply = prompt_kisski("")
Give me a json object of a bounding boxes around the cat in this image.
The format should be exactly like this: {'x':int,'y':int,'width':int,'height':int},
cat_image)
print(reply)
bb = json.loads(extract_json(reply))
bb

stackview.add_bounding_boxes(cat_image, [bb])

...json
{
"x": 6,
"y": 10,
"width": 427,
"height": 408
}
...

[4]:

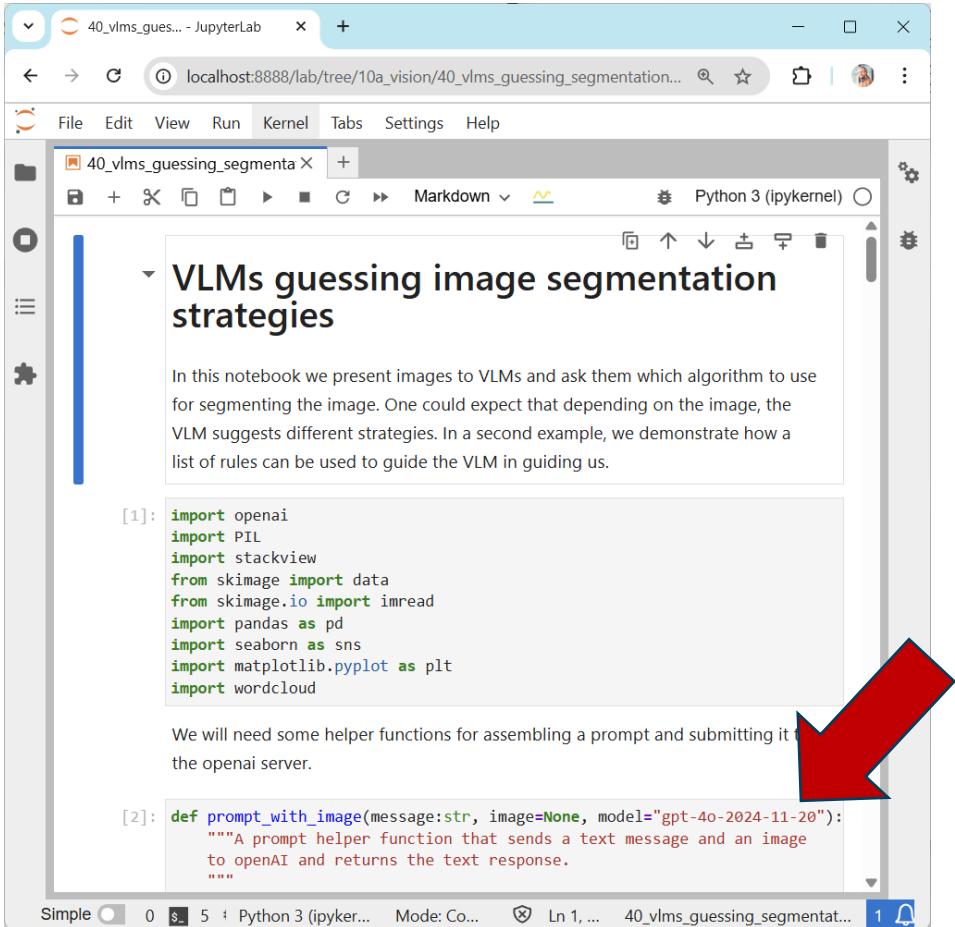


shape (512, 512, 3)
dtype uint8
size 768.0 kB
min 0
max 255

Simple 0 Python 3 (ipykernel) | ... Mode: Comm... Ln 1, C... 23_vision_bb_segmentatio... 1

Vision language models

Elaborate on the limits of VLMs



The screenshot shows a Jupyter Notebook interface with the title '40_vlms_guessing_segmentation'. The notebook contains the following content:

VLMs guessing image segmentation strategies

In this notebook we present images to VLMs and ask them which algorithm to use for segmenting the image. One could expect that depending on the image, the VLM suggests different strategies. In a second example, we demonstrate how a list of rules can be used to guide the VLM in guiding us.

```
[1]: import openai
import PIL
import stackview
from skimage import data
from skimage.io import imread
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import wordcloud
```

We will need some helper functions for assembling a prompt and submitting it to the openai server.

```
[2]: def prompt_with_image(message:str, image=None, model="gpt-4o-2024-11-20"):
    """A prompt helper function that sends a text message and an image to openAI and returns the text response.
    """

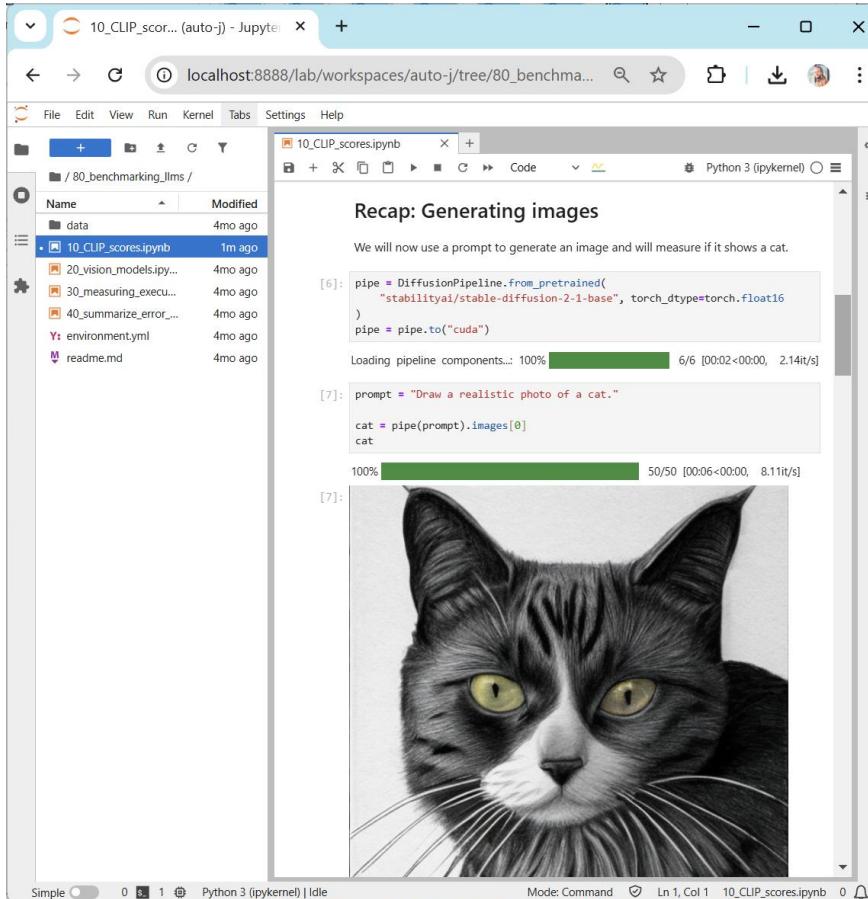
```

Simple 0 5 + Python 3 (ipyker... Mode: Co... Ln 1, ... 40_vlms_guessing_segmentat... 1

Try to do this with ollama or kisski instead.

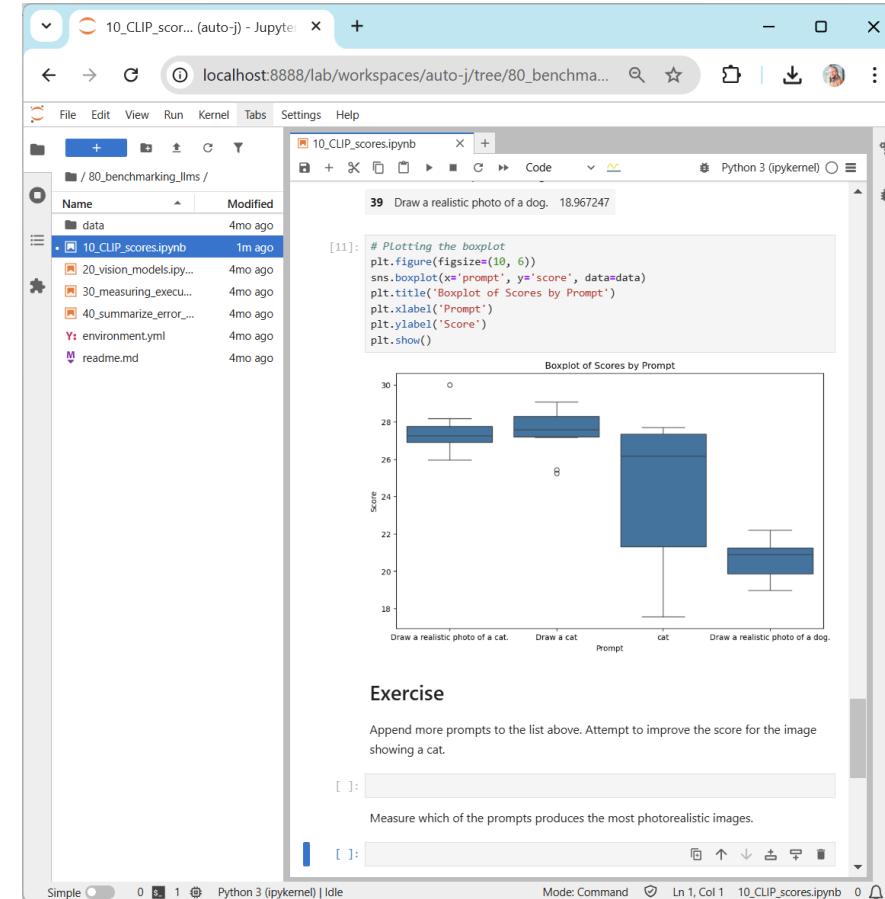
CLIP Scores

Improve given prompts and measure the improvement



The screenshot shows a Jupyter Notebook interface. On the left, a sidebar lists files: 10_CLIP_scores.ipynb (modified 1m ago), 20_vision_models.ipynb, 30_measuring_execu..., 40_summarize_error..., environment.yml, and readme.md. The main notebook cell [6] contains Python code to initialize a diffusion pipeline and generate an image from a prompt. The output shows a black and white image of a cat's face. Cell [7] shows the generated image.

```
pipe = DiffusionPipeline.from_pretrained("stabilityai/stable-diffusion-2-1-base", torch_dtype=torch.float16)
pipe = pipe.to("cuda")  
  
prompt = "Draw a realistic photo of a cat."  
  
cat = pipe(prompt).images[0]  
cat
```



Exam preparation

Use the opportunity to try out this Vision-RAG system written by Lea Gihlein, who attended this course last year.

- Generate exam questions
- See related Slides

Beware, not all slides might be related to the question!

```
# Setting up the model
endpoint = "http://localhost:11434/v1"
token = "whatever"
model = "gemma3:4b"
```

```
question_topic = "Linear filtering"
```

1. What is the core operation behind a linear filter, and how is it mathematically described?
2. Name three examples of linear filters presented in the slides.
3. Explain the difference between intrinsically explainable and black-box algorithms, using Linear Regression and Deep Neural Networks as examples.

The collage consists of four rectangular panels, each containing a different slide from the ScaDS.AI repository:

- Top Left:** A slide titled "Linear filters" showing two examples of convolution operations. It includes a "Terminology" section with definitions of "convolve" and "convolution operator", and a "Examples" section listing Mean, Gaussian blur, Sobel, and Laplace filters.
- Top Right:** A slide titled "Quiz: Noise removal" featuring a grayscale image of a cell with noise. It asks, "The median filter is a ...". Below the image are two options: "Linear filter" (represented by a pink square) and "Non-linear filter" (represented by a green square).
- Bottom Left:** Another slide titled "Linear filters" that explains what linear filters do, defines filter kernels as matrices, and describes convolution as multiplying surrounding pixels according to a matrix. It shows a "Mean Filter, 3x3 kernel" example.
- Bottom Right:** A slide titled "Explainability" comparing "Intrinsically explainable AI-algorithms" (like Linear Regression, shown with the equation $f(x_1, x_2) = w_1x_1 + w_2x_2$) and "Black-Box AI-algorithms" (like Deep Neural Networks, shown with a diagram of a neural network). It notes that if w_1 is much bigger than w_2 , the result depends much more on x_1 compared to x_2 .

Preliminary schedule

April 11th 2025 – Introduction Microscopy, Bioimage Analysis, Research Software Management

~~April 18th 2025 – Good Friday (Karfreitag)~~

April 25th 2025 - Microscopy Image Processing

May 2nd 2025 – Segmentation of cells and nuclei

May 9th 2025 – Feature Extraction, Data Viz, Quality Assurance

Handout exam pre-requisite
complex exercise

May 16th 2025 – Big data, parallel processing & tiled image computing

May 23th 2025 – Introduction to Machine Learning for bio-image analysis

May 30st 2025 – Deep Learning for image denoising + segmentation

June 6th 2025 – Large Language Models and Prompt Engineering

[June 13th 2025 – DevOps Hands-on tutorial @ Data Week, optional <https://www.mathcs.uni-leipzig.de/veranstaltungsdetail/termin/data-week-leipzig-2025>]

June 20th 2025 – Large Language Models for Code Generation

June 27th 2025 – Vision Language Models and Agents

July 4th 2025 – Research Data Management

Submission deadline
complex exercise

July 11nd 2024 – Summary, exam preparation

8 weeks