# Image Generation and Vision Language Models

## Robert Haase
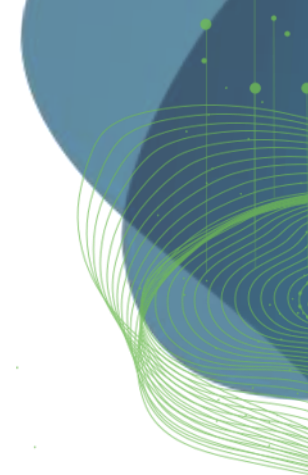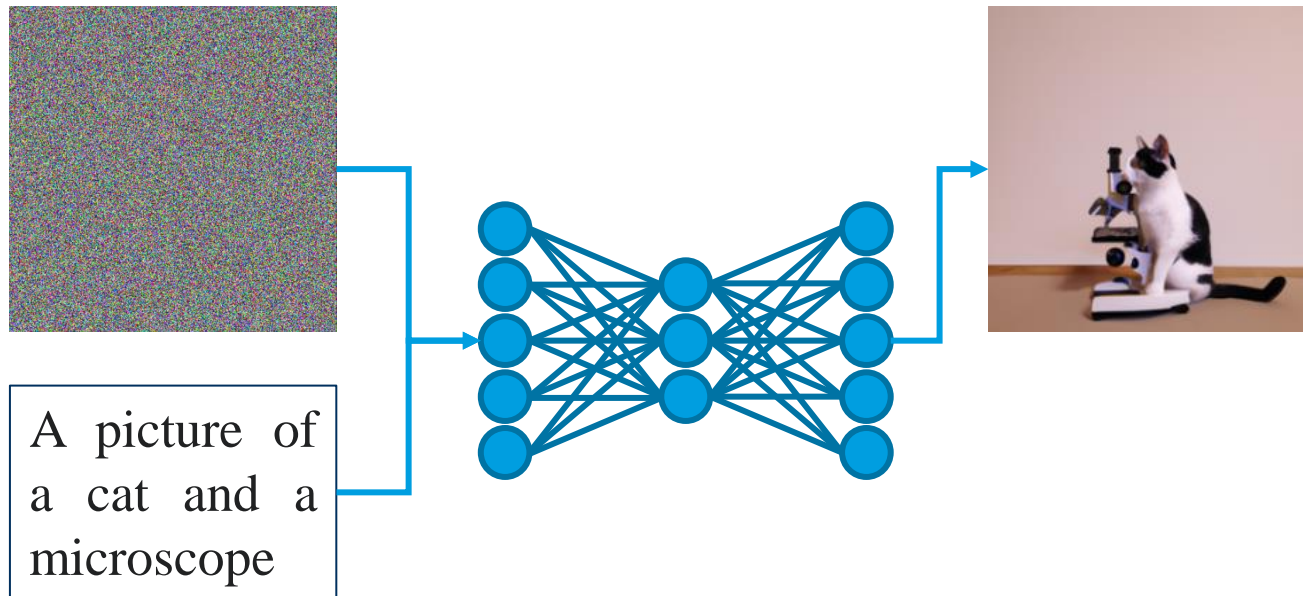
CENTER FOR SCALABLE DATA ANALYTICS AND ARTIFICIAL INTELLIGENCE

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Image Generation

„text-to-image"

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Variational Auto-Encoder

„image-to-image"

Cat picture source: Ramesh et al. 2021, licensed
CC-BY 4.0 https://arxiv.org/pdf/2102.12092.pdf

# Image Generation



Decoder

„Image of a cat" → Large Language Model → "Bottleneck", "Embedding" → Output
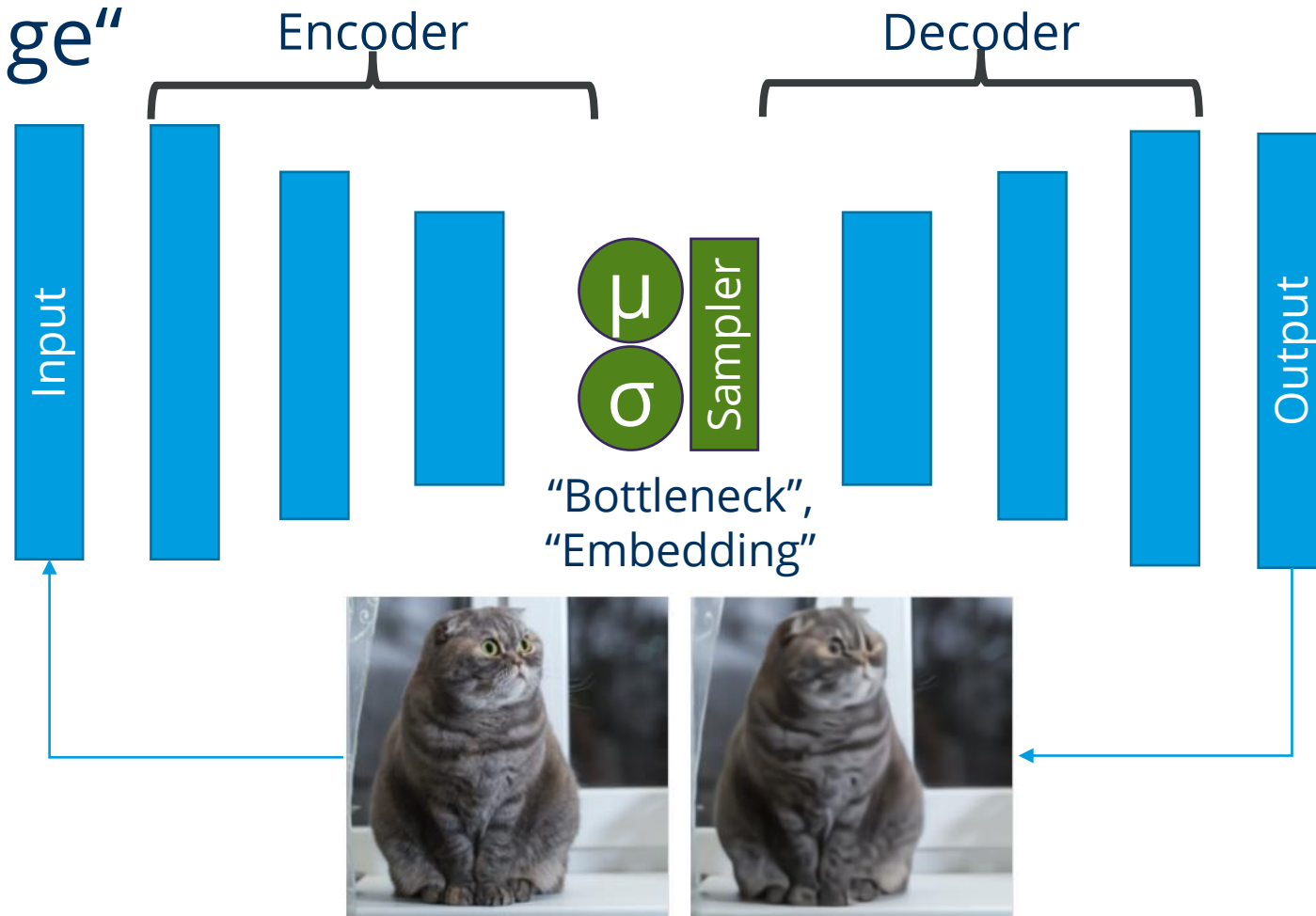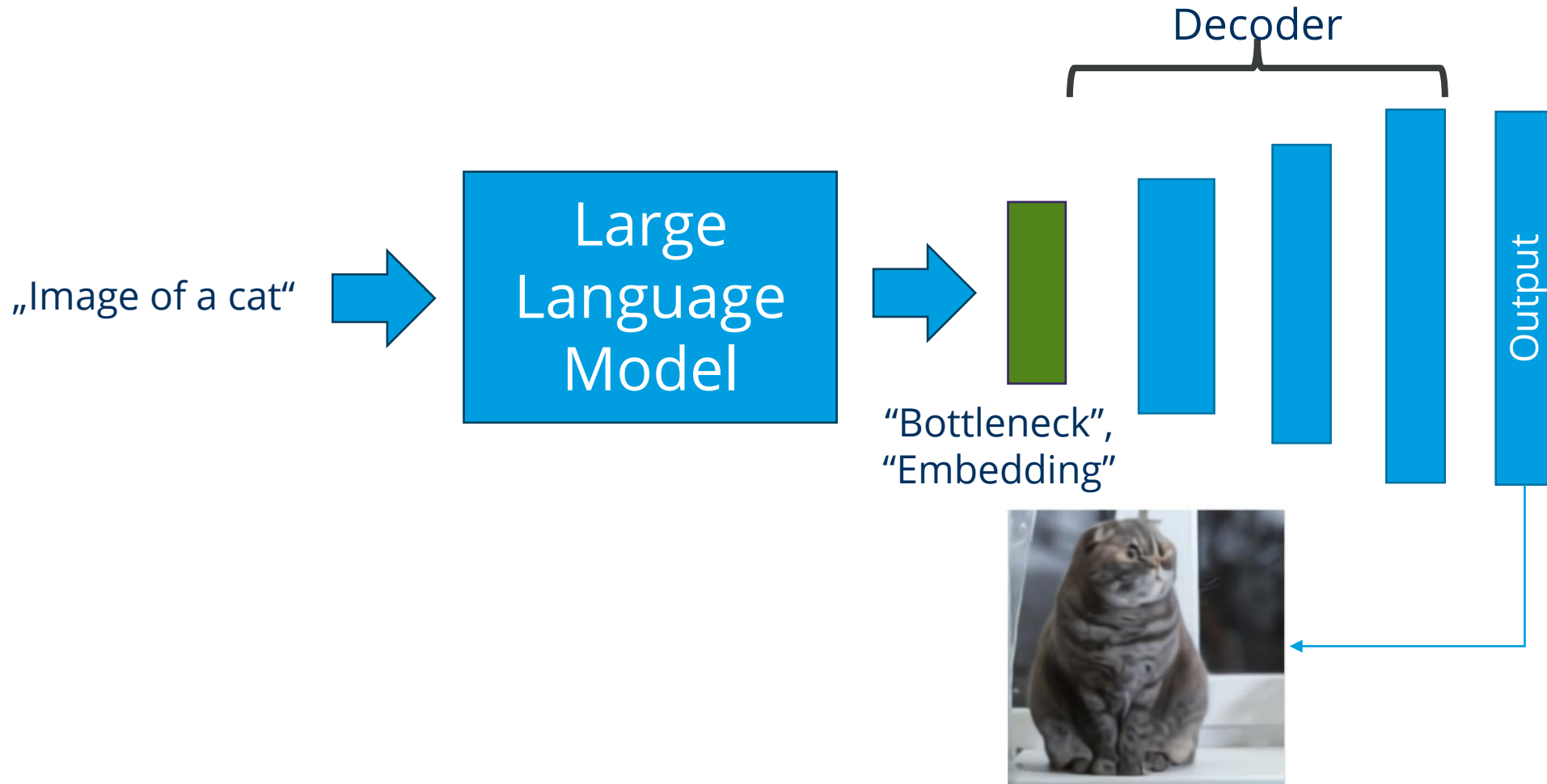
Slide 4

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

ScaDS.AI DRESDEN LEIPZIG

# Stable Diffusion

Diffusion: reverse denoising

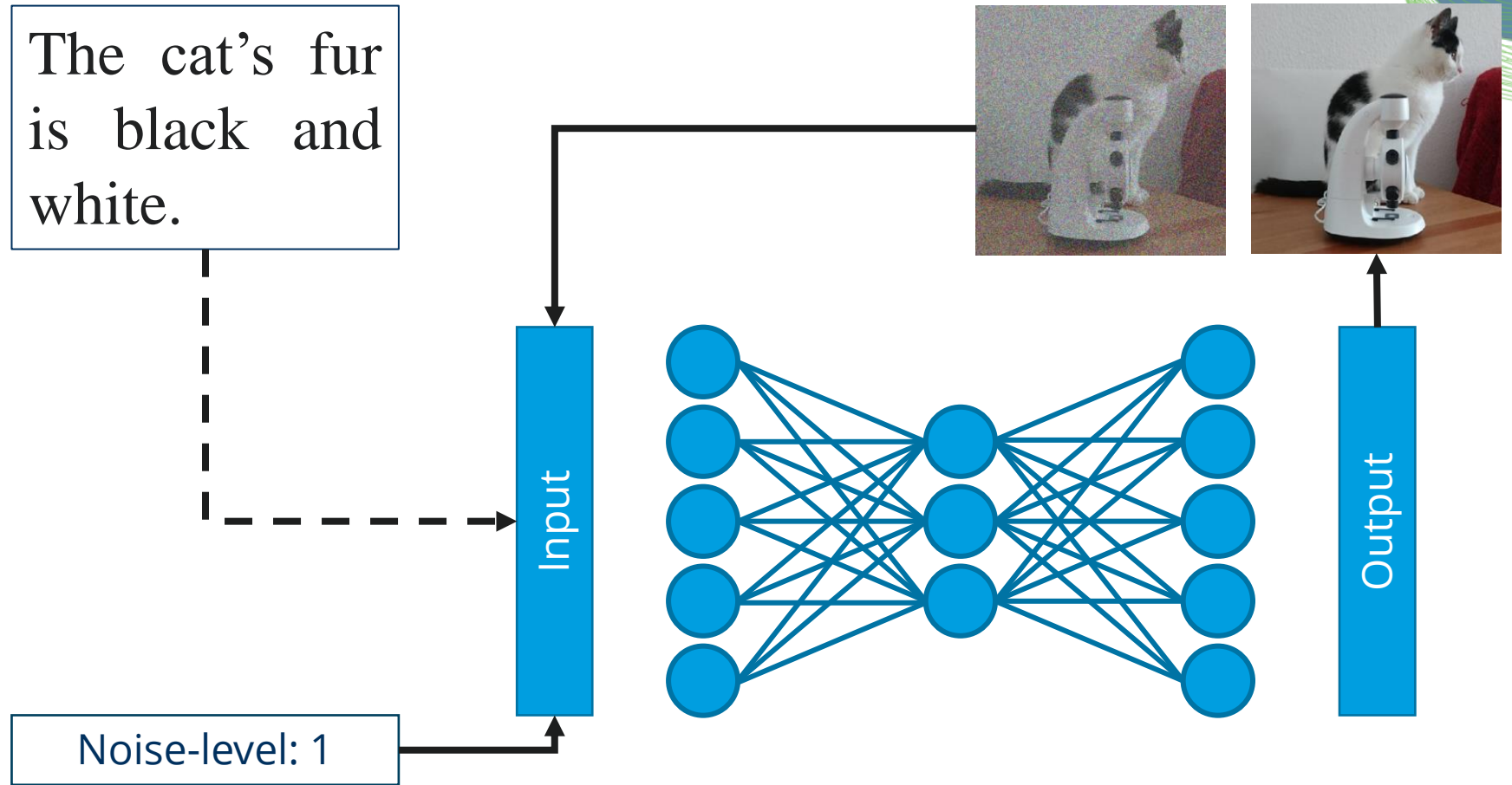Simplified from: Rombach et al. 2021
https://arxiv.org/abs/2112.10752

# Stable Diffusion

Diffusion: reverse denoising

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

Source: Rombach et al. 2021
https://arxiv.org/abs/2112.10752

# How does it work?

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

The cat's fur is black and white.



Input

Output

Noise-level: 1

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# How does it work?

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



The cat's fur is black and white.

Noise-level: 2

Input

Output

# How does it work?

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



The cat's fur is black and white.

Input

Output

Noise-level: 3

# How does it work?

Prediction is iterative denoising of:

Pure noise + text prompt

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# How does it work?

Reminder:

- Word embeddings

- Attention



The cat's fur is black and white.

Self-attention

0  1  4  3  2  1

The cat's fur is black and white.

Word Embedding

"microscope"

"cat"

"black"
"white"        "fur"

Input

Output

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Image generation in Python: Huggingface

Most Huggingface image-generation models require a GPU.

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

Downloads
4.8 GB

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

Needs
Nvidia GPU

# Image generation in Python: Huggingface

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

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Dall-E

OpenAI's model for image geneation based on diffusion models + CLIP transformer

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Image generation in Python: Dall-E

## No need for a GPU, but costs 💸 💸 💸

```python
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
```
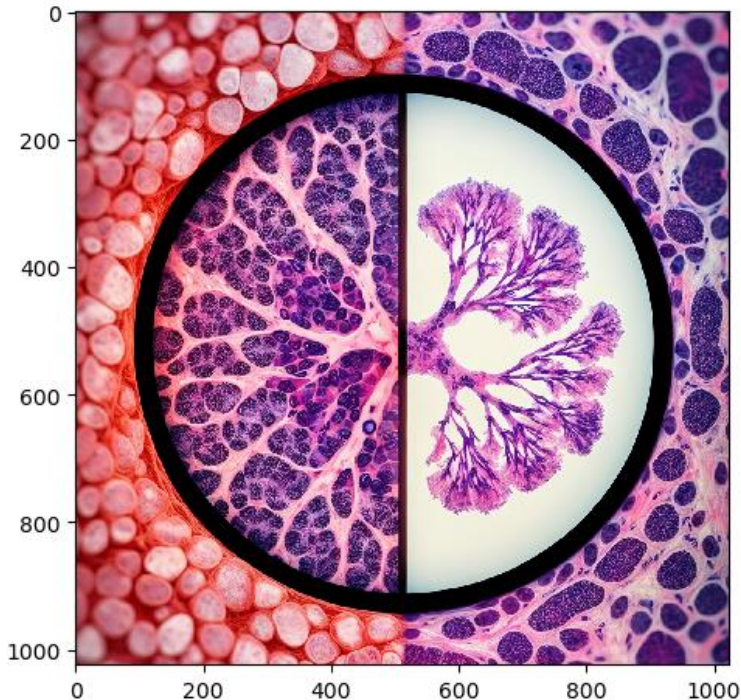
Works with
Dall-E 2 and 3

May soon also
work with gpt-4o

# Image generation in Python: Dall-E

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



Dall-E 3



Dall-E 2

# Inpainting

Replacing regions in images

(also „Gap-filling", „Replacing")

Raw image

Mask image

Manipulated image

A black white cat fur

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Inpainting in Python: Huggingface

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

Downloads
4.8 GB

Needs
Nvidia GPU

```python
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]
```
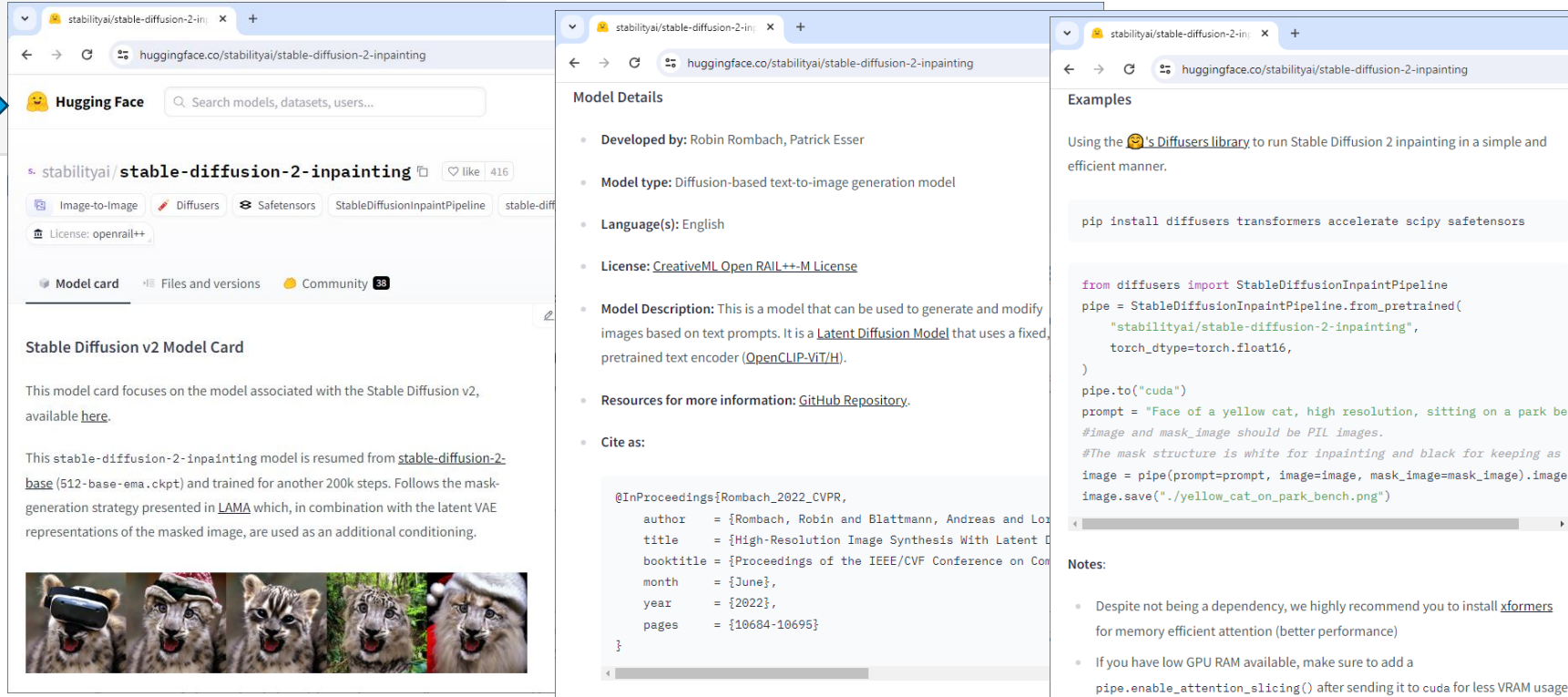
Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

Read more:
https://huggingface.co/docs/diffusers/api/
pipelines/stable_diffusion/inpaint

# Inpainting in Python: Huggingface

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

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

https://huggingface.co/stabilityai/
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.

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Inpainting in Python: Dall-E

No need for a GPU, but costs 💸 💸 💸

```python
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

Size must match

Supported: 256, 512, 1024 pixels

➡️ Result: List of URL(s)

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024
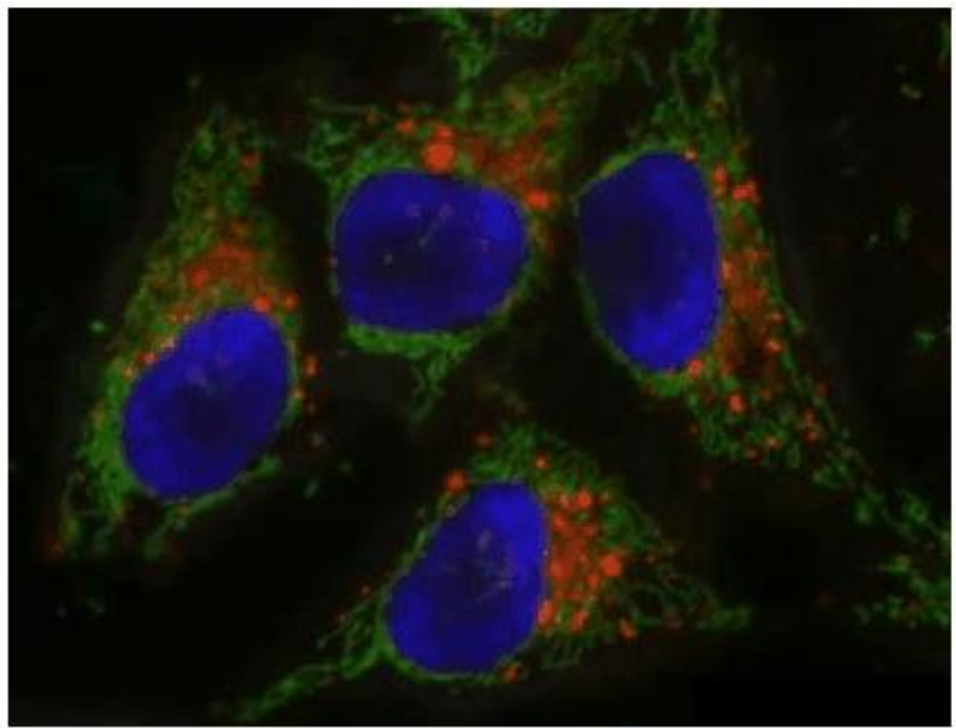
# New technologies bring new risks...

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

Interesting challenges for our community ahead

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

https://github.com/haesleinhuepf/darth-d/blob/main/demo/demo_replacing.ipynb

# Image manipulation detection

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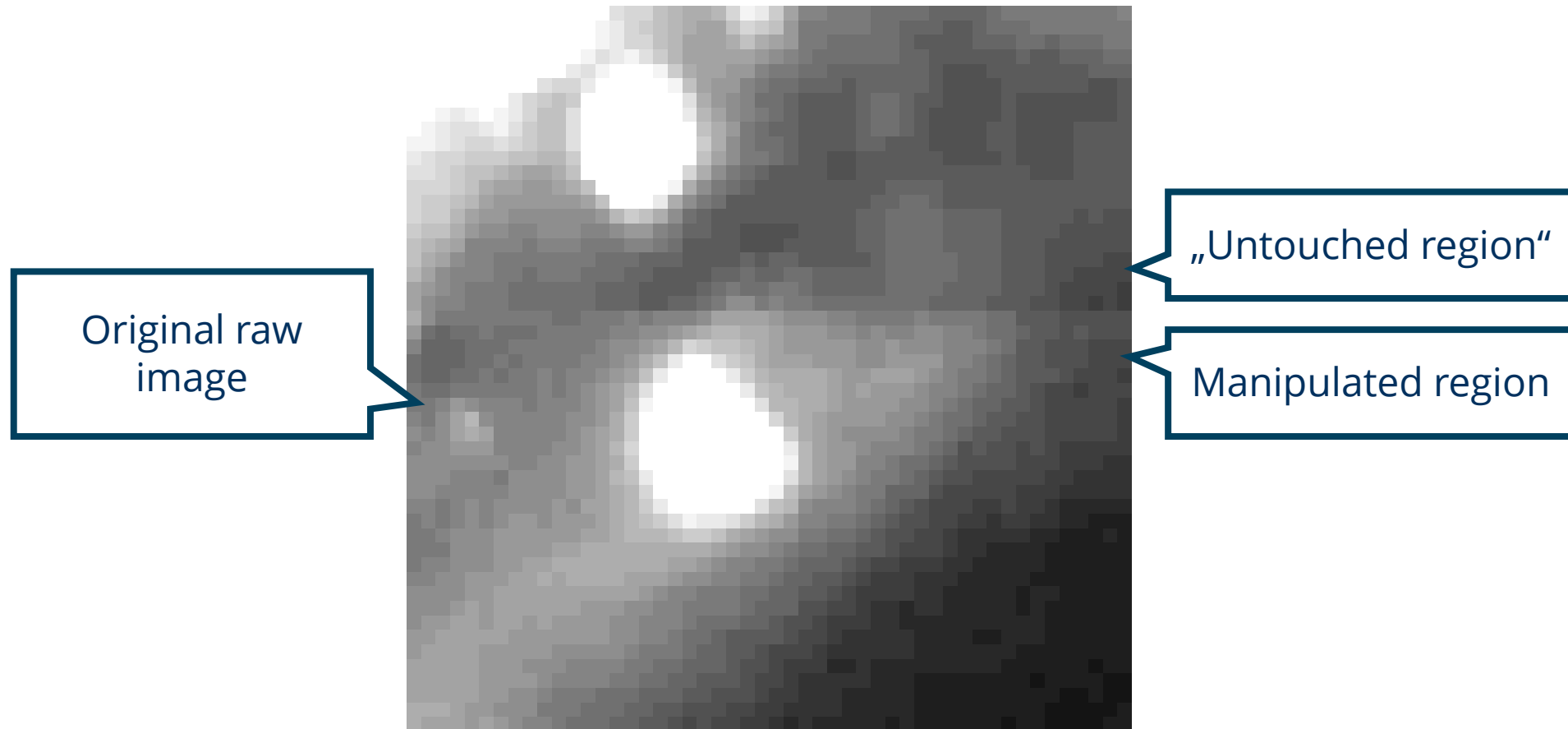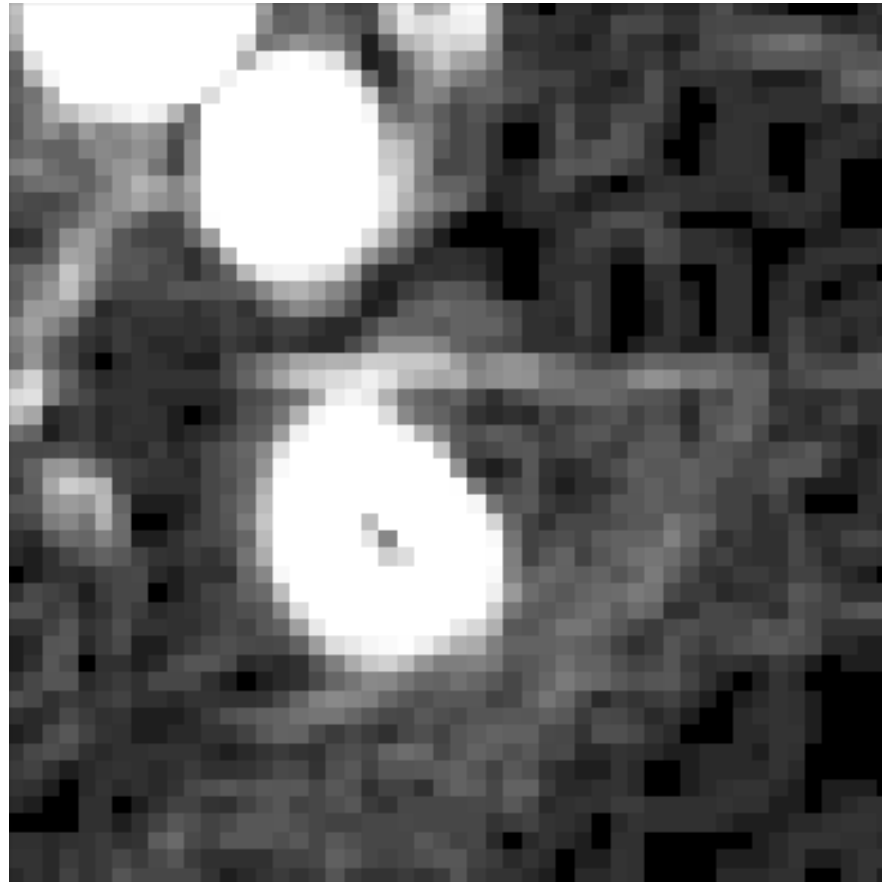


"Untouched region"

Original raw image

Manipulated region

# Image manipulation detection

e.g. by studying noise-patterns

Local standard deviation filter



„Untouched region"

Manipulated region

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Image manipulation detection

## e.g. by studying noise-patterns

Sobel filter



„Untouched region"

Manipulated region

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Vision Language Models

- Classifying images 🥱
- Describing images



A picture of a cat and a microscope

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Vision Language Models (VLMs)

Goal: Describe images



Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Variational Auto-Encoder



Encoder

Decoder

Input

Output

μ σ Sampler

"Bottleneck",
"Embedding"

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

Cat picture source: Ramesh et al. 2021, licensed CC-BY 4.0 https://arxiv.org/pdf/2102.12092.pdf

Slide 29

# Image classification

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Image classification -> image describing



Encoder

Input

"Bottleneck", "Embedding"

Large Language Model

„Image of a dog"

„Image of a cat"

„Image of a microscope"

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

ScaDS.AI DRESDEN LEIPZIG

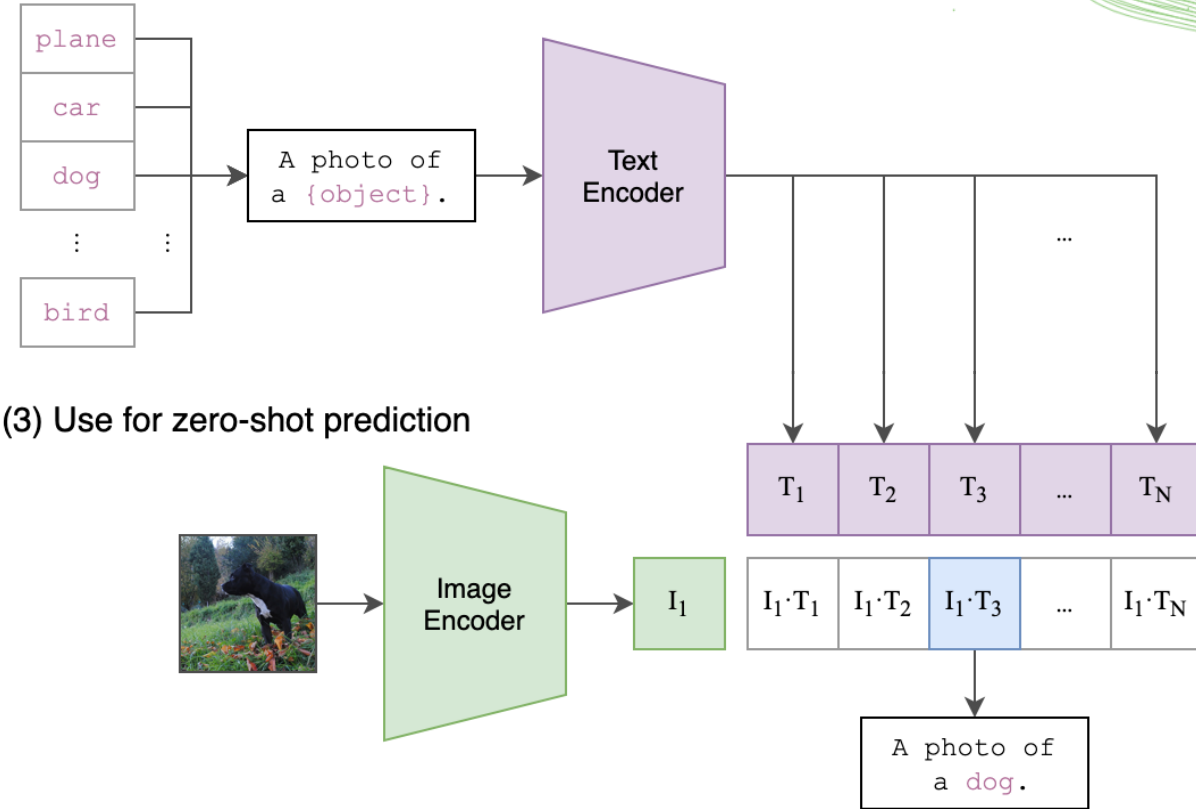TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Contrastive Language-Image Pre-Training

## „CLIP" Transformers



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

Source: Radford et al 2021, License MT
https://arxiv.org/abs/2103.00020
https://github.com/OpenAI/CLIP

# CLIP transformers in Python

## Using huggingface 🤗

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

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

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

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

...

```python
label_probabilities
```

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

```python
label_probabilities
```

```python
{'a photo of a cat': 0.135291740541458,
 'a photo of a dog': 0.001265904747436404,
 'a photo of a microscope': 0.8634429574012756}
```
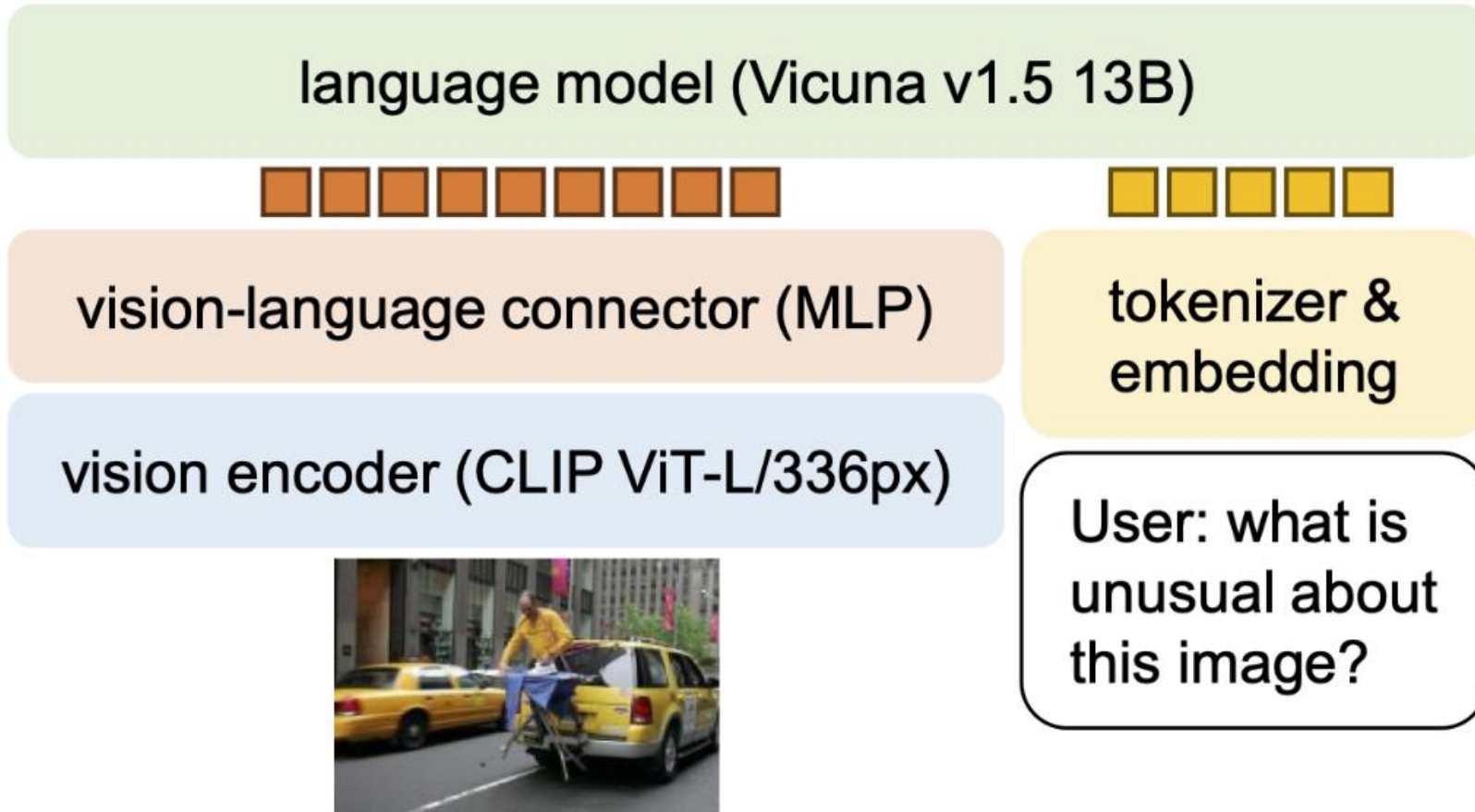
Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

Code example adapted from:
https://huggingface.co/docs/transformers/en/model_doc/clip

Slide 33

**TECHNISCHE UNIVERSITÄT DRESDEN**

**UNIVERSITÄT LEIPZIG**

# LLAVA

## Large Language and Vision Assistant

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# LLAVA 1.5

## Combining LLAVA with CLIP

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024
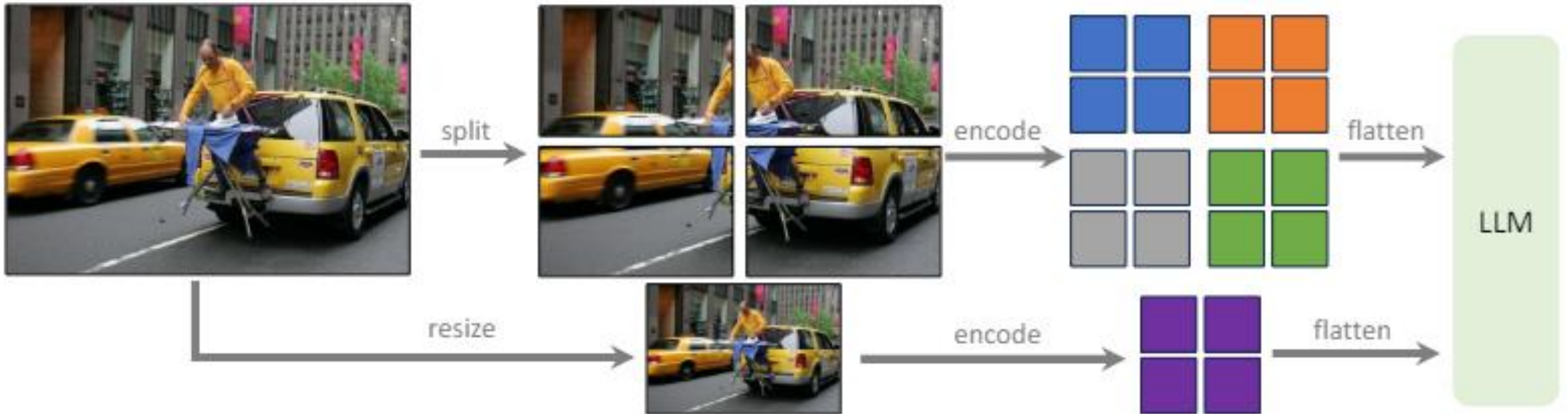
# LLAVA 1.5 HD

Giving the model multiple perspectives on the same scene

# Accessing VLMs using Python

API not standardized (yet)

```python
def prompt_chatGPT(prompt:str, image, model="gpt-4o"):
    """A prompt helper function that sends a message to openAI
    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/jpeg;base64,{base64_image}"
        }
    }]}]

    # setup connection to the LLM
    client = openai.OpenAI()

    # submit prompt
    response = client.chat.completions.create(
        model=model,
        messages=message
    )

    # extract answer
    return response.choices[0].message.content
```

```python
def prompt_ollama(prompt:str, image, model="llava"):
    """A prompt helper function that sends a message to ollama
    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': prompt,
        'images': [base64_image]
    }]

    # setup connection to the LLM
    client = openai.OpenAI(
        base_url = "http://localhost:11434/v1"
    )

    # submit prompt
    response = client.chat.completions.create(
        model=model,
        messages=message
    )

    # extract answer
    return response.choices[0].message.content
```

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

Slide 41

ScaDS.AI DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Exercises

## Robert Haase

Funded by

CENTER FOR SCALABLE DATA ANALYTICS AND
ARTIFICIAL INTELLIGENCE

# Exercise: Image generation

Try to identify and create realistically looking MRI images
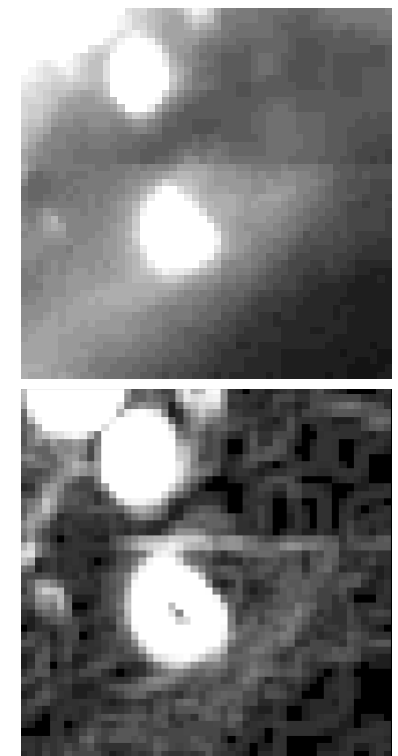
Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Exercise: Image manipulation

Inspect the image carefully, try to find the border of the manipulated region

Hint:

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024

# Exercise: Vision

Ask llava and gpt-4omni to describe an image *and* to produce Python code for analysing it.

Robert Haase
@haesleinhuepf
BIDS Lecture 12/14
June 18th 2024