

▼ Text-to-Image Generation with Stable Diffusion and OpenVINO™

Stable Diffusion is a text-to-image latent diffusion model created by the researchers and engineers from [CompVis](#), [Stability AI](#) and [LAION](#). It is trained on 512x512 images from a subset of the [LAION-5B](#) database. This model uses a frozen CLIP ViT-L/14 text encoder to condition the model on text prompts. With its 860M UNet and 123M text encoder. See the [model card](#) for more information.

General diffusion models are machine learning systems that are trained to denoise random gaussian noise step by step, to get to a sample of interest, such as an image. Diffusion models have shown to achieve state-of-the-art results for generating image data. But one downside of diffusion models is that the reverse denoising process is slow. In addition, these models consume a lot of memory because they operate in pixel space, which becomes unreasonably expensive when generating high-resolution images. Therefore, it is challenging to train these models and also use them for inference. OpenVINO brings capabilities to run model inference on Intel hardware and opens the door to the fantastic world of diffusion models for everyone!

Model capabilities are not limited text-to-image only, it also is able solve additional tasks, for example text-guided image-to-image generation and inpainting. This tutorial also considers how to run text-guided image-to-image generation using Stable Diffusion.

This notebook demonstrates how to convert and run stable diffusion model using OpenVINO.

Notebook contains the following steps:

1. Convert PyTorch models to ONNX format.
2. Convert ONNX models to OpenVINO IR format, using Model Optimizer tool.
3. Run Stable Diffusion pipeline with OpenVINO.

▼ Prerequisites

The following is needed only if you want to use the original model. If not, you do not have to do anything. Just run the notebook.

Note: The original model (for example, `stable-diffusion-v1-4`) requires you to accept the model license before downloading or using its weights. Visit the [stable-diffusion-v1-4 card](#) to read and accept the license before you proceed. To use this diffusion model, you must be a registered user in 🤗 Hugging Face Hub. You will need to use an access token for the code below to run. For more information on access tokens, refer to [this section of the documentation](#). You can login on Hugging Face Hub in notebook environment, using following code:

login to huggingfacehub to get access to pretrained model

```
from huggingface_hub import notebook_login, whoami
```

```
try: whoami() print('Authorization token already provided') except OSError: notebook_login()
```

This tutorial uses a Stable Diffusion model, fine-tuned using images from Midjourney v4 (another popular solution for text-to-image generation). You can find more details about this model on the [model card](https://huggingface.co/prompthero/openjourney). The same s

```
1 !pip install -r requirements.txt
```

```
ERROR: Could not open requirements file: [Errno 2] No such file or directory: 'requirements.txt'
```

▼ Create Pytorch Models pipeline

StableDiffusionPipeline is an end-to-end inference pipeline that you can use to generate images from text with just a few lines of code.

First, load the pre-trained weights of all components of the model.

```
1 !pip install diffusers
2 !pip install transformers
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: diffusers in /usr/local/lib/python3.9/dist-packages (0.14.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.9/dist-packages (from diffusers) (1.22.4)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.9/dist-packages (from diffusers)
Requirement already satisfied: requests in /usr/local/lib/python3.9/dist-packages (from diffusers) (2.27.1)
Requirement already satisfied: Pillow in /usr/local/lib/python3.9/dist-packages (from diffusers) (8.4.0)
Requirement already satisfied: importlib-metadata in /usr/local/lib/python3.9/dist-packages (from diffuser)
Requirement already satisfied: filelock in /usr/local/lib/python3.9/dist-packages (from diffusers) (3.10.7)
Requirement already satisfied: huggingface-hub<=0.10.0 in /usr/local/lib/python3.9/dist-packages (from dif
Requirement already satisfied: packaging>=20.9 in /usr/local/lib/python3.9/dist-packages (from huggingface
Requirement already satisfied: tqdm>=4.42.1 in /usr/local/lib/python3.9/dist-packages (from huggingface-hu
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.9/dist-packages (from huggingface-hub
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.9/dist-packages (from
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.9/dist-packages (from importlib-metadat
Requirement already satisfied: charset-normalizer==2.0.0 in /usr/local/lib/python3.9/dist-packages (from r
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.9/dist-packages (from requests
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.9/dist-packages (from requests->diff
Requirement already satisfied: urllib3<1.27,>=1.21.1 in /usr/local/lib/python3.9/dist-packages (from reque
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: transformers in /usr/local/lib/python3.9/dist-packages (4.27.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.9/dist-packages (from transformer
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.9/dist-packages (from transformers) (4
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.9/dist-packages (from transformers) (
Requirement already satisfied: huggingface-hub<1.0,>=0.11.0 in /usr/local/lib/python3.9/dist-packages (fro
Requirement already satisfied: filelock in /usr/local/lib/python3.9/dist-packages (from transformers) (3.1
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.9/dist-packages (from transformers) (
Requirement already satisfied: requests in /usr/local/lib/python3.9/dist-packages (from transformers) (2.2
Requirement already satisfied: tokenizers!=0.11.3,<0.14,>=0.11.1 in /usr/local/lib/python3.9/dist-packages
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.9/dist-packages (from transform
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.9/dist-packages (from
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.9/dist-packages (from requests->tran
Requirement already satisfied: urllib3<1.27,>=1.21.1 in /usr/local/lib/python3.9/dist-packages (from reque
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.9/dist-packages (from requests
Requirement already satisfied: charset-normalizer==2.0.0 in /usr/local/lib/python3.9/dist-packages (from r
```

```
1 from diffusers import StableDiffusionPipeline
2
3 pipe = StableDiffusionPipeline.from_pretrained("prompthero/openjourney").to("cpu")
4 text_encoder = pipe.text_encoder
5 text_encoder.eval()
6 unet = pipe.unet
7 unet.eval()
8 vae = pipe.vae
9 vae.eval()
10
11 del pipe
```

```
Fetching 15 files: 100% 15/15 [00:00<00:00, 916.19it/s]
```

```
Cannot initialize model with low cpu memory usage because `accelerate` was not found in the environment. D
...
pip install accelerate
...
/usr/local/lib/python3.9/dist-packages/transformers/models/clip/feature_extraction_clip.py:28: FutureWarni
warnings.warn(
`text_config_dict` is provided which will be used to initialize `CLIPTextConfig`. The value `text_config["
```

▼ Convert models to OpenVINO Intermediate representation (IR) format

OpenVINO supports PyTorch through export to the ONNX format. You will use `torch.onnx.export` function for obtaining ONNX model. You can learn more in the [PyTorch documentation](#). You need to provide a model object, input data for model tracing and a path for saving the model. Optionally, you can provide the target onnx opset for conversion and other parameters specified in documentation (for example, input and output names or dynamic shapes).

While ONNX models are directly supported by OpenVINO™ runtime, it can be useful to convert them to IR format to take advantage of advanced OpenVINO optimization tools and features. You will use OpenVINO Model Optimizer tool for conversion model to IR format and compression weights to **FP16** format.

The model consists of three important parts:

- Text Encoder for creation condition to generate image from text prompt.
- Unet for step by step denoising latent image representation.
- Autoencoder (VAE) for encoding input image to latent space (if required) and decoding latent space to image back after generation.

Let us convert each part.

▼ Text Encoder

The text-encoder is responsible for transforming the input prompt, for example, "a photo of an astronaut riding a horse" into an embedding space that can be understood by the U-Net. It is usually a simple transformer-based encoder that maps a sequence of input tokens to a sequence of latent text embeddings.

Input of the text encoder is the tensor `input_ids` which contains indexes of tokens from text processed by tokenizer and padded to maximum length accepted by model. Model outputs are two tensors: `last_hidden_state` - hidden state from the last MultiHeadAttention layer in the model and `pooler_out` - Pooled output for whole model hidden states. You will use `opset_version=14`, because model contains `triu` operation, supported in ONNX only starting from this opset.

```
1 from pathlib import Path
2 import torch
3 !pip install opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0
```

```
Requirement already satisfied: pandas~=1.3.5 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: opencv-dev==2022.3.0 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: scipy>=1.8 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: addict>=2.4.0 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: protobuf<4.0.0,>=3.18.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: fastjsonschema==2.15.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: onnx<=1.12,>=1.8.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: tensorflow<=2.9.3,>=2.5 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: torchvision<=0.14.0,>=0.9.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: torch<=1.13.0,>=1.8.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: yacs>=0.1.8 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: typing-extensions>=3.6.2.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.9/dist-packages (from pandas~=1.3.5 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.9/dist-packages (from pandas~=1.3.5 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: charset-normalizer==2.0.0 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.9/dist-packages (from requests>=2.25.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.9/dist-packages (from requests>=2.25.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: urllib3<1.27,>=1.21.1 in /usr/local/lib/python3.9/dist-packages (from requests>=2.25.1 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: flatbuffers<2,>=1.12 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: opt-einsum==2.3.2 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: tensorflow-estimator<2.10.0,>=2.9.0rc0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: setuptools in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: wrapt>=1.11.0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: astunparse>=1.6.0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: tensorboard<2.10,>=2.9 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: six>=1.12.0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: h5py>=2.9.0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: libclang>=13.0.0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: keras-preprocessing>=1.1.1 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: gast<=0.4.0,>=0.2.1 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: grpcio<2.0,>=1.24.3 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: tensorflow-io-gcs-filesystem>=0.23.1 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: packaging in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: google-pasta>=0.1.1 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: keras<2.10.0,>=2.9.0rc0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
Requirement already satisfied: absl-py>=1.0.0 in /usr/local/lib/python3.9/dist-packages (from tensorflow==2.9.3 in /usr/local/lib/python3.9/dist-packages (from opencv-dev[pytorch,ONNX,tensorflow2]==2022.3.0))
```

```
Requirement already satisfied: rsa<5,>=3.1.4 in /usr/local/lib/python3.9/dist-packages (from google-auth)
Requirement already satisfied: requests-oauthlib>=0.7.0 in /usr/local/lib/python3.9/dist-packages (from google-auth)
Requirement already satisfied: importlib-metadata>=4.4 in /usr/local/lib/python3.9/dist-packages (from werkzeug)
Requirement already satisfied: MarkupSafe>=2.1.1 in /usr/local/lib/python3.9/dist-packages (from werkzeug)
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.9/dist-packages (from importlib-metadata)
Requirement already satisfied: pyasn1<0.5.0,>=0.4.6 in /usr/local/lib/python3.9/dist-packages (from pyasn1)
Requirement already satisfied: oauthlib>=3.0.0 in /usr/local/lib/python3.9/dist-packages (from requests-oauthlib)
```

```
1 import gc
2 from pathlib import Path
3 import torch
4
5 TEXT_ENCODER_ONNX_PATH = Path('text_encoder.onnx')
6 TEXT_ENCODER_OV_PATH = TEXT_ENCODER_ONNX_PATH.with_suffix('.xml')
7
8
9 def convert_encoder_onnx(text_encoder: StableDiffusionPipeline, onnx_path:Path):
10     """
11     Convert Text Encoder model to ONNX.
12     Function accepts pipeline, prepares example inputs for ONNX conversion via torch.export,
13     Parameters:
14         pipe (StableDiffusionPipeline): Stable Diffusion pipeline
15         onnx_path (Path): File for storing onnx model
16     Returns:
17         None
18     """
19     if not onnx_path.exists():
20         input_ids = torch.ones((1, 77), dtype=torch.long)
21         # switch model to inference mode
22         text_encoder.eval()
23
24         # disable gradients calculation for reducing memory consumption
25         with torch.no_grad():
26             # infer model, just to make sure that it works
27             text_encoder(input_ids)
28             # export model to ONNX format
29             torch.onnx.export(
30                 text_encoder, # model instance
31                 input_ids, # inputs for model tracing
32                 onnx_path, # output file for saving result
33                 input_names=['tokens'], # model input name for onnx representation
34                 output_names=['last_hidden_state', 'pooler_out'], # model output names for onnx representation
35                 opset_version=14 # onnx opset version for export
36             )
37         print('Text Encoder successfully converted to ONNX')
38
39
40 if not TEXT_ENCODER_OV_PATH.exists():
41     convert_encoder_onnx(text_encoder, TEXT_ENCODER_ONNX_PATH)
42     !mo --input_model $TEXT_ENCODER_ONNX_PATH --compress_to_fp16
43     print('Text Encoder successfully converted to IR')
44 else:
45     print(f"Text encoder will be loaded from {TEXT_ENCODER_OV_PATH}")
46
47 del text_encoder
48 gc.collect()
```

```
Text encoder will be loaded from text_encoder.xml
81
```

▼ U-net

Unet model has three inputs:

- **sample** - latent image sample from previous step. Generation process has not been started yet, so you will use random noise.
- **timestep** - current scheduler step.
- **encoder_hidden_state** - hidden state of text encoder.

Model predicts the **sample** state for the next step.

```
1 !pip install onnx
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: onnx in /usr/local/lib/python3.9/dist-packages (1.12.0)
Requirement already satisfied: protobuf<=3.20.1,>=3.12.2 in /usr/local/lib/python3.9/dist-packages (from onnx) (3.20.1)
Requirement already satisfied: numpy>=1.16.6 in /usr/local/lib/python3.9/dist-packages (from onnx) (1.22.4)
Requirement already satisfied: typing-extensions>=3.6.2.1 in /usr/local/lib/python3.9/dist-packages (from onnx) (4.5.0)
```

```
1 # !git clone https://github.com/microsoft/vcpkg.git
2 # %cd vcpkg
3 # !./bootstrap-vcpkg.bat # For powershell
4 # !./bootstrap-vcpkg.sh # For bash
5 # !./vcpkg install onnx
```

```
1 import numpy as np
2
3 UNET_ONNX_PATH = Path('unet/unet.onnx')
4 UNET_OV_PATH = UNET_ONNX_PATH.parents[1] / 'unet.xml'
5
6
7 def convert_unet_onnx(unet:StableDiffusionPipeline, onnx_path:Path):
8     """
9     Convert Unet model to ONNX, then IR format.
10    Function accepts pipeline, prepares example inputs for ONNX conversion via torch.export,
11    Parameters:
12        pipe (StableDiffusionPipeline): Stable Diffusion pipeline
13        onnx_path (Path): File for storing onnx model
14    Returns:
15        None
16    """
17    if not onnx_path.exists():
18        # prepare inputs
19        encoder_hidden_state = torch.ones((2, 77, 768))
20        latents_shape = (2, 4, 512 // 8, 512 // 8)
21        latents = torch.randn(latents_shape)
22        t = torch.from_numpy(np.array(1, dtype=float))
23
24        # model size > 2Gb, it will be represented as onnx with external data files, you will store it in s
25        onnx_path.parent.mkdir(exist_ok=True, parents=True)
26        unet.eval()
27
28        with torch.no_grad():
29            torch.onnx.export(
30                unet,
31                (latents, t, encoder_hidden_state), str(onnx_path),
32                input_names=['latent_model_input', 't', 'encoder_hidden_states'],
33                output_names=['out_sample'],
34                opset_version = 11,
35                verbose=True
36            )
37        print('Unet successfully converted to ONNX')
38
39    if not UNET_OV_PATH.exists():
40        convert_unet_onnx(unet, UNET_ONNX_PATH)
41        del unet
42        gc.collect()
43        !mo --input_model $UNET_ONNX_PATH --compress_to_fp16
44        print('Unet successfully converted to IR')
45    else:
46        del unet
47        print(f"Unet will be loaded from {UNET_OV_PATH}")
48    gc.collect()
```

```
/usr/local/lib/python3.9/dist-packages/diffusers/models/unet_2d_condition.py:526: TracerWarning: Converting a tensor to a numpy array. This operation can lead to errors when using tracing.
  if any(s % default_overall_up_factor != 0 for s in sample.shape[-2:]):
/usr/local/lib/python3.9/dist-packages/diffusers/models/resnet.py:185: TracerWarning: Converting a tensor to a numpy array. This operation can lead to errors when using tracing.
  assert hidden_states.shape[1] == self.channels
/usr/local/lib/python3.9/dist-packages/diffusers/models/resnet.py:190: TracerWarning: Converting a tensor to a numpy array. This operation can lead to errors when using tracing.
  assert hidden_states.shape[1] == self.channels
/usr/local/lib/python3.9/dist-packages/diffusers/models/resnet.py:112: TracerWarning: Converting a tensor to a numpy array. This operation can lead to errors when using tracing.
  assert hidden_states.shape[1] == self.channels
```

```

/usr/local/lib/python3.9/dist-packages/diffusers/models/resnet.py:125: TracerWarning: Converting a tensor
if hidden_states.shape[0] >= 64:
/usr/local/lib/python3.9/dist-packages/diffusers/models/unet_2d_condition.py:651: TracerWarning: Converting a tensor to a
if not return_dict:
/usr/local/lib/python3.9/dist-packages/torch/onnx/symbolic_helper.py:817: UserWarning: You are trying to export a module
ONNX's Upsample/Resize operator did not match Pytorch's Interpolation until opset 11. Attributes to determine the mode
We recommend using opset 11 and above for models using this operator.
warnings.warn(
Unet successfully converted to ONNX
[ INFO ] The model was converted to IR v11, the latest model format that corresponds to the source DL framework.
Find more information about API v2.0 and IR v11 at https://docs.openvino.ai/latest/openvino\_2\_0\_transition
[ SUCCESS ] Generated IR version 11 model.
[ SUCCESS ] XML file: /content/unet.xml
[ SUCCESS ] BIN file: /content/unet.bin
Unet successfully converted to IR
0

```

▼ VAE

The VAE model has two parts, an encoder and a decoder. The encoder is used to convert the image into a low dimensional latent representation, which will serve as the input to the U-Net model. The decoder, conversely, transforms the latent representation back into an image.

During latent diffusion training, the encoder is used to get the latent representations (latents) of the images for the forward diffusion process, which applies more and more noise at each step. During inference, the denoised latents generated by the reverse diffusion process are converted back into images using the VAE decoder. When you run inference for text-to-image, there is no initial image as a starting point. You can skip this step and directly generate initial random noise.

As the encoder and the decoder are used independently in different parts of the pipeline, it will be better to convert them to separate models.

```

1 VAE_ENCODER_ONNX_PATH = Path('vae_encoder.onnx')
2 VAE_ENCODER_OV_PATH = VAE_ENCODER_ONNX_PATH.with_suffix('.xml')
3
4
5 def convert_vae_encoder_onnx(vae: StableDiffusionPipeline, onnx_path: Path):
6     """
7     Convert VAE model to ONNX, then IR format.
8     Function accepts pipeline, creates wrapper class for export only necessary for inference part,
9     prepares example inputs for ONNX conversion via torch.export,
10    Parameters:
11        pipe (StableDiffusionInstructPix2PixPipeline): InstructPix2Pix pipeline
12        onnx_path (Path): File for storing onnx model
13    Returns:
14        None
15    """
16    class VAEEncoderWrapper(torch.nn.Module):
17        def __init__(self, vae):
18            super().__init__()
19            self.vae = vae
20
21        def forward(self, image):
22            h = self.vae.encoder(image)
23            moments = self.vae.quant_conv(h)
24            return moments
25
26    if not onnx_path.exists():
27        vae_encoder = VAEEncoderWrapper(vae)
28        vae_encoder.eval()
29        image = torch.zeros((1, 3, 512, 512))
30        with torch.no_grad():
31            torch.onnx.export(vae_encoder, image, onnx_path, input_names=[
32                'init_image'], output_names=['image_latent'])
33        print('VAE encoder successfully converted to ONNX')
34
35
36 if not VAE_ENCODER_OV_PATH.exists():
37     convert_vae_encoder_onnx(vae, VAE_ENCODER_ONNX_PATH)
38     !mo --input_model $VAE_ENCODER_ONNX_PATH --compress_to_fp16
39     print('VAE encoder successfully converted to IR')
40 else:

```

```

41     print(f"VAE encoder will be loaded from {VAE_ENCODER_OV_PATH}")
42
43     VAE_DECODER_ONNX_PATH = Path('vae_decoder.onnx')
44     VAE_DECODER_OV_PATH = VAE_DECODER_ONNX_PATH.with_suffix('.xml')
45
46
47     def convert_vae_decoder_onnx(vae: StableDiffusionPipeline, onnx_path: Path):
48         """
49         Convert VAE model to ONNX, then IR format.
50         Function accepts pipeline, creates wrapper class for export only necessary for inference part,
51         prepares example inputs for ONNX conversion via torch.export,
52         Parameters:
53             pipe (StableDiffusionInstructPix2PixPipeline): InstructPix2Pix pipeline
54             onnx_path (Path): File for storing onnx model
55         Returns:
56             None
57         """
58         class VAEDecoderWrapper(torch.nn.Module):
59             def __init__(self, vae):
60                 super().__init__()
61                 self.vae = vae
62
63             def forward(self, latents):
64                 latents = 1 / 0.18215 * latents
65                 return self.vae.decode(latents)
66
67         if not onnx_path.exists():
68             vae_decoder = VAEDecoderWrapper(vae)
69             latents = torch.zeros((1, 4, 64, 64))
70
71             vae_decoder.eval()
72             with torch.no_grad():
73                 torch.onnx.export(vae_decoder, latents, onnx_path, input_names=[
74                     'latents'], output_names=['sample'])
75             print('VAE decoder successfully converted to ONNX')
76
77
78         if not VAE_DECODER_OV_PATH.exists():
79             convert_vae_decoder_onnx(vae, VAE_DECODER_ONNX_PATH)
80             !mo --input_model $VAE_DECODER_ONNX_PATH --compress_to_fp16
81             print('VAE decoder successfully converted to IR')
82         else:
83             print(f"VAE decoder will be loaded from {VAE_DECODER_OV_PATH}")
84
85     del vae

```

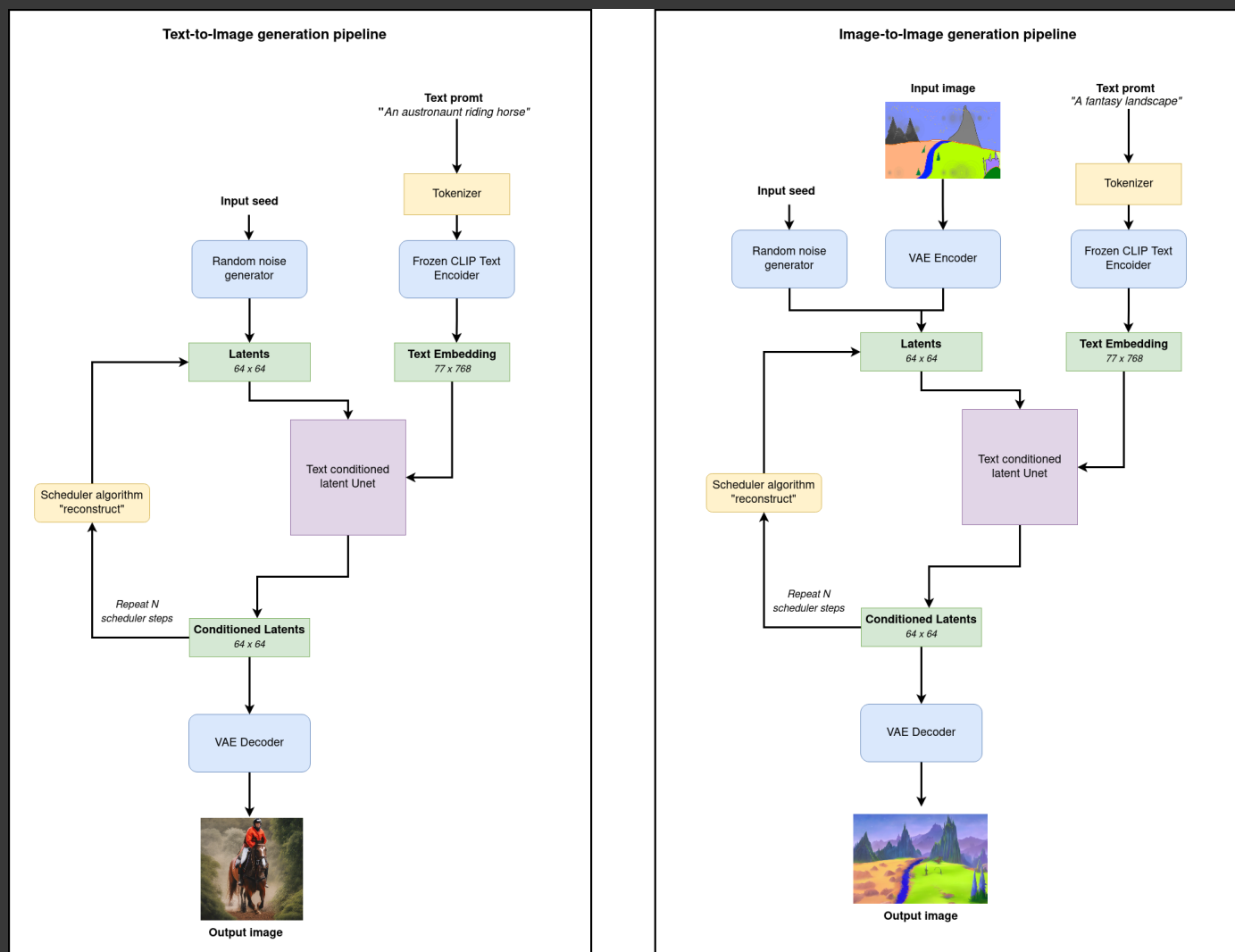
```

/usr/local/lib/python3.9/dist-packages/torch/onnx/_internal/jit_utils.py:258: UserWarning: Constant folding
_C._jit_pass_onnx_node_shape_type_inference(node, params_dict, opset_version)
/usr/local/lib/python3.9/dist-packages/torch/onnx/utils.py:687: UserWarning: Constant folding - Only steps
_C._jit_pass_onnx_graph_shape_type_inference(
/usr/local/lib/python3.9/dist-packages/torch/onnx/utils.py:1178: UserWarning: Constant folding - Only step
_C._jit_pass_onnx_graph_shape_type_inference(
VAE encoder successfully converted to ONNX
[ INFO ] The model was converted to IR v11, the latest model format that corresponds to the source DL fram
Find more information about API v2.0 and IR v11 at https://docs.openvino.ai/latest/openvino\_2\_0\_transition
[ SUCCESS ] Generated IR version 11 model.
[ SUCCESS ] XML file: /content/vae_encoder.xml
[ SUCCESS ] BIN file: /content/vae_encoder.bin
VAE encoder successfully converted to IR
/usr/local/lib/python3.9/dist-packages/torch/onnx/_internal/jit_utils.py:258: UserWarning: The shape infer
_C._jit_pass_onnx_node_shape_type_inference(node, params_dict, opset_version)
/usr/local/lib/python3.9/dist-packages/torch/onnx/utils.py:687: UserWarning: The shape inference of prim::
_C._jit_pass_onnx_graph_shape_type_inference(
/usr/local/lib/python3.9/dist-packages/torch/onnx/utils.py:1178: UserWarning: The shape inference of prim::
_C._jit_pass_onnx_graph_shape_type_inference(
VAE decoder successfully converted to ONNX
[ INFO ] The model was converted to IR v11, the latest model format that corresponds to the source DL fram
Find more information about API v2.0 and IR v11 at https://docs.openvino.ai/latest/openvino\_2\_0\_transition
[ SUCCESS ] Generated IR version 11 model.
[ SUCCESS ] XML file: /content/vae_decoder.xml
[ SUCCESS ] BIN file: /content/vae_decoder.bin
VAE decoder successfully converted to IR

```

Prepare Inference Pipeline

Putting it all together, let us now take a closer look at how the model works in inference by illustrating the logical flow.



As you can see from the diagram, the only difference between Text-to-Image and text-guided Image-to-Image generation in approach is how initial latent state is generated. In case of Image-to-Image generation, you additionally have an image encoded by VAE encoder mixed with the noise produced by using latent seed, while in Text-to-Image you use only noise as initial latent state. The stable diffusion model takes both a latent image representation of size 64×64 and a text prompt is transformed to text embeddings of size 77×768 via CLIP's text encoder as an input.

Next, the U-Net iteratively *denoises* the random latent image representations while being conditioned on the text embeddings. The output of the U-Net, being the noise residual, is used to compute a denoised latent image representation via a scheduler algorithm. Many different scheduler algorithms can be used for this computation, each having its pros and cons. For Stable Diffusion, it is recommended to use one of:

- [PNDM scheduler](#)
- [DDIM scheduler](#)
- [K-LMS scheduler](#)(you will use it in your pipeline)

Theory on how the scheduler algorithm function works is out of scope for this notebook. Nonetheless, in short, you should remember that you compute the predicted denoised image representation from the previous noise representation and the predicted noise residual. For more information, refer to the recommended [Elucidating the Design Space of Diffusion-Based Generative Models](#)

The *denoising* process is repeated given number of times (by default 50) to step-by-step retrieve better latent image representations. When complete, the latent image representation is decoded by the decoder part of the variational auto encoder.


```

1 import inspect
2 from typing import List, Optional, Union, Dict
3
4 import PIL
5 import cv2
6
7 from transformers import CLIPTokenizer
8 from diffusers.pipeline_utils import DiffusionPipeline
9 from diffusers.schedulers import DDIMScheduler, LMSDiscreteScheduler, PNDMScheduler
10 from openvino.runtime import Model
11
12
13 def scale_fit_to_window(dst_width:int, dst_height:int, image_width:int, image_height:int):
14     """
15     Preprocessing helper function for calculating image size for resize with preserving original aspect ratio
16     and fitting image to specific window size
17
18     Parameters:
19         dst_width (int): destination window width
20         dst_height (int): destination window height
21         image_width (int): source image width
22         image_height (int): source image height
23     Returns:
24         result_width (int): calculated width for resize
25         result_height (int): calculated height for resize
26     """
27     im_scale = min(dst_height / image_height, dst_width / image_width)
28     return int(im_scale * image_width), int(im_scale * image_height)
29
30
31 def preprocess(image: PIL.Image.Image):
32     """
33     Image preprocessing function. Takes image in PIL.Image format, resizes it to keep aspect ratio and fit
34     then converts it to np.ndarray and adds padding with zeros on right or bottom side of image (depends fr
35     converts data to float32 data type and change range of values from [0, 255] to [-1, 1], finally, conver
36     The function returns preprocessed input tensor and padding size, which can be used in postprocessing.
37
38     Parameters:
39         image (PIL.Image.Image): input image
40     Returns:
41         image (np.ndarray): preprocessed image tensor
42         meta (Dict): dictionary with preprocessing metadata info
43     """
44     src_width, src_height = image.size
45     dst_width, dst_height = scale_fit_to_window(
46         512, 512, src_width, src_height)
47     image = np.array(image.resize((dst_width, dst_height),
48                                 resample=PIL.Image.Resampling.LANCZOS))[None, :]
49     pad_width = 512 - dst_width
50     pad_height = 512 - dst_height
51     pad = ((0, 0), (0, pad_height), (0, pad_width), (0, 0))
52     image = np.pad(image, pad, mode="constant")
53     image = image.astype(np.float32) / 255.0
54     image = 2.0 * image - 1.0
55     image = image.transpose(0, 3, 1, 2)
56     return image, {"padding": pad, "src_width": src_width, "src_height": src_height}
57
58
59 class OVStableDiffusionPipeline(DiffusionPipeline):
60     def __init__(
61         self,
62         vae_decoder: Model,
63         text_encoder: Model,
64         tokenizer: CLIPTokenizer,
65         unet: Model,
66         scheduler: Union[DDIMScheduler, PNDMScheduler, LMSDiscreteScheduler],
67         vae_encoder: Model = None,
68     ):
69         """
70         Pipeline for text-to-image generation using Stable Diffusion.
71         Parameters:
72             vae (Model):
73                 Variational Auto-Encoder (VAE) Model to decode images to and from latent representations.

```

```

74     text_encoder (Model):
75         Frozen text-encoder. Stable Diffusion uses the text portion of
76         [CLIP](https://huggingface.co/docs/transformers/model_doc/clip#transformers.CLIPTextModel),
77         the clip-vit-large-patch14(https://huggingface.co/openai/clip-vit-large-patch14) variant.
78     tokenizer (CLIPTokenizer):
79         Tokenizer of class CLIPTokenizer(https://huggingface.co/docs/transformers/v4.21.0/en/model_
80     unet (Model): Conditional U-Net architecture to denoise the encoded image latents.
81     scheduler (SchedulerMixin):
82         A scheduler to be used in combination with unet to denoise the encoded image latents. Can b
83         DDIMScheduler, LMSDiscreteScheduler, or PNDMScheduler.
84     """
85     super().__init__()
86     self.scheduler = scheduler
87     self.vae_decoder = vae_decoder
88     self.vae_encoder = vae_encoder
89     self.text_encoder = text_encoder
90     self.unet = unet
91     self._text_encoder_output = text_encoder.output(0)
92     self._unet_output = unet.output(0)
93     self._vae_d_output = vae_decoder.output(0)
94     self._vae_e_output = vae_encoder.output(0) if vae_encoder is not None else None
95     self.height = self.unet.input(0).shape[2] * 8
96     self.width = self.unet.input(0).shape[3] * 8
97     self.tokenizer = tokenizer
98
99     def __call__(
100         self,
101         prompt: Union[str, List[str]],
102         image: PIL.Image.Image = None,
103         num_inference_steps: Optional[int] = 50,
104         guidance_scale: Optional[float] = 7.5,
105         eta: Optional[float] = 0.0,
106         output_type: Optional[str] = "pil",
107         seed: Optional[int] = None,
108         strength: float = 1.0,
109         gif: Optional[bool] = False,
110         **kwargs,
111     ):
112         """
113         Function invoked when calling the pipeline for generation.
114         Parameters:
115             prompt (str or List[str]):
116                 The prompt or prompts to guide the image generation.
117             image (PIL.Image.Image, *optional*, None):
118                 Intinal image for generation.
119             num_inference_steps (int, *optional*, defaults to 50):
120                 The number of denoising steps. More denoising steps usually lead to a higher quality image
121                 expense of slower inference.
122             guidance_scale (float, *optional*, defaults to 7.5):
123                 Guidance scale as defined in Classifier-Free Diffusion Guidance(https://arxiv.org/abs/2207.
124                 guidance_scale is defined as `w` of equation 2.
125                 Higher guidance scale encourages to generate images that are closely linked to the text pro
126                 usually at the expense of lower image quality.
127             eta (float, *optional*, defaults to 0.0):
128                 Corresponds to parameter eta ( $\eta$ ) in the DDIM paper: https://arxiv.org/abs/2010.02502. Only
129                 [DDIMScheduler], will be ignored for others.
130             output_type (`str`, *optional*, defaults to "pil"):
131                 The output format of the generate image. Choose between
132                 [PIL](https://pillow.readthedocs.io/en/stable/): PIL.Image.Image or np.array.
133             seed (int, *optional*, None):
134                 Seed for random generator state initialization.
135             gif (bool, *optional*, False):
136                 Flag for storing all steps results or not.
137         Returns:
138             Dictionary with keys:
139                 sample - the last generated image PIL.Image.Image or np.array
140                 iterations - *optional* (if gif=True) images for all diffusion steps, List of PIL.Image.Ima
141         """
142         if seed is not None:
143             np.random.seed(seed)
144
145         if isinstance(prompt, str):
146             batch_size = 1

```

```

147     elif isinstance(prompt, list):
148         batch_size = len(prompt)
149     else:
150         raise ValueError(f"`prompt` has to be of type `str` or `list` but is {type(prompt)}")
151
152     img_buffer = []
153     # get prompt text embeddings
154     text_input = self.tokenizer(
155         prompt,
156         padding="max_length",
157         max_length=self.tokenizer.model_max_length,
158         truncation=True,
159         return_tensors="np",
160     )
161     text_embeddings = self.text_encoder(text_input.input_ids)[self._text_encoder_output]
162     # here `guidance_scale` is defined analog to the guidance weight `w` of equation (2)
163     # of the Imagen paper: https://arxiv.org/pdf/2205.11487.pdf . `guidance_scale = 1`
164     # corresponds to doing no classifier free guidance.
165     do_classifier_free_guidance = guidance_scale > 1.0
166     # get unconditional embeddings for classifier free guidance
167     if do_classifier_free_guidance:
168         max_length = text_input.input_ids.shape[-1]
169         uncond_input = self.tokenizer(
170             [""] * batch_size, padding="max_length", max_length=max_length, return_tensors="np"
171         )
172         uncond_embeddings = self.text_encoder(uncond_input.input_ids)[self._text_encoder_output]
173
174         # For classifier free guidance, you need to do two forward passes.
175         # Here you concatenate the unconditional and text embeddings into a single batch
176         # to avoid doing two forward passes
177         text_embeddings = np.concatenate([uncond_embeddings, text_embeddings])
178
179     # set timesteps
180     accepts_offset = "offset" in set(inspect.signature(self.scheduler.set_timesteps).parameters.keys())
181     extra_set_kwargs = {}
182     if accepts_offset:
183         extra_set_kwargs["offset"] = 1
184
185     self.scheduler.set_timesteps(num_inference_steps, **extra_set_kwargs)
186     timesteps, num_inference_steps = self.get_timesteps(num_inference_steps, strength)
187     latent_timestep = timesteps[:1]
188
189     # get the initial random noise unless the user supplied it
190     latents, meta = self.prepare_latents(image, latent_timestep)
191
192     # prepare extra kwargs for the scheduler step, since not all schedulers have the same signature
193     # eta ( $\eta$ ) is only used with the DDIMScheduler, it will be ignored for other schedulers.
194     # eta corresponds to  $\eta$  in DDIM paper: https://arxiv.org/abs/2010.02502
195     # and should be between [0, 1]
196     accepts_eta = "eta" in set(inspect.signature(self.scheduler.step).parameters.keys())
197     extra_step_kwargs = {}
198     if accepts_eta:
199         extra_step_kwargs["eta"] = eta
200
201     for i, t in enumerate(self.progress_bar(timesteps)):
202         # expand the latents if you are doing classifier free guidance
203         latent_model_input = np.concatenate([latents] * 2) if do_classifier_free_guidance else latents
204         latent_model_input = self.scheduler.scale_model_input(latent_model_input, t)
205
206         # predict the noise residual
207         noise_pred = self.unet([latent_model_input, t, text_embeddings])[self._unet_output]
208         # perform guidance
209         if do_classifier_free_guidance:
210             noise_pred_uncond, noise_pred_text = noise_pred[0], noise_pred[1]
211             noise_pred = noise_pred_uncond + guidance_scale * (noise_pred_text - noise_pred_uncond)
212
213         # compute the previous noisy sample  $x_t \rightarrow x_{t-1}$ 
214         latents = self.scheduler.step(torch.from_numpy(noise_pred), t, torch.from_numpy(latents), **extra_step_kwargs)
215         if gif:
216             image = self.vae_decoder(latents)[self._vae_d_output]
217             image = self.postprocess_image(image, meta, output_type)
218             img_buffer.extend(image)
219

```

```

220     # scale and decode the image latents with vae
221     image = self.vae_decoder(latents)[self._vae_d_output]
222
223     image = self.postprocess_image(image, meta, output_type)
224     return {"sample": image, 'iterations': img_buffer}
225
226 def prepare_latents(self, image:PIL.Image.Image = None, latent_timestep:torch.Tensor = None):
227     """
228     Function for getting initial latents for starting generation
229
230     Parameters:
231         image (PIL.Image.Image, *optional*, None):
232             Input image for generation, if not provided random noise will be used as starting point
233         latent_timestep (torch.Tensor, *optional*, None):
234             Predicted by scheduler initial step for image generation, required for latent image mixing
235     Returns:
236         latents (np.ndarray):
237             Image encoded in latent space
238     """
239     latents_shape = (1, 4, self.height // 8, self.width // 8)
240     noise = np.random.randn(*latents_shape).astype(np.float32)
241     if image is None:
242         # if you use LMSDiscreteScheduler, let's make sure latents are multiplied by sigmas
243         if isinstance(self.scheduler, LMSDiscreteScheduler):
244             noise = noise * self.scheduler.sigmas[0].numpy()
245         return noise, {}
246     input_image, meta = preprocess(image)
247     moments = self.vae_encoder(input_image)[self._vae_e_output]
248     mean, logvar = np.split(moments, 2, axis=1)
249     std = np.exp(logvar * 0.5)
250     latents = (mean + std * np.random.randn(*mean.shape)) * 0.18215
251     latents = self.scheduler.add_noise(torch.from_numpy(latents), torch.from_numpy(noise), latent_times
252     return latents, meta
253
254 def postprocess_image(self, image:np.ndarray, meta:Dict, output_type:str = "pil"):
255     """
256     Postprocessing for decoded image. Takes generated image decoded by VAE decoder, unpad it to initila
257     normalize and convert to [0, 255] pixels range. Optionally, convertes it from np.ndarray to PIL.Ima
258
259     Parameters:
260         image (np.ndarray):
261             Generated image
262         meta (Dict):
263             Metadata obtained on latents preparing step, can be empty
264         output_type (str, *optional*, pil):
265             Output format for result, can be pil or numpy
266     Returns:
267         image (List of np.ndarray or PIL.Image.Image):
268             Postprocessed images
269     """
270     if "padding" in meta:
271         pad = meta["padding"]
272         (_, end_h), (_, end_w) = pad[1:3]
273         h, w = image.shape[2:]
274         unpad_h = h - end_h
275         unpad_w = w - end_w
276         image = image[:, :, :unpad_h, :unpad_w]
277     image = np.clip(image / 2 + 0.5, 0, 1)
278     image = np.transpose(image, (0, 2, 3, 1))
279     # 9. Convert to PIL
280     if output_type == "pil":
281         image = self.numpy_to_pil(image)
282         if "src_height" in meta:
283             orig_height, orig_width = meta["src_height"], meta["src_width"]
284             image = [img.resize((orig_width, orig_height),
285                               PIL.Image.Resampling.LANCZOS) for img in image]
286     else:
287         if "src_height" in meta:
288             orig_height, orig_width = meta["src_height"], meta["src_width"]
289             image = [cv2.resize(img, (orig_width, orig_height))
290                     for img in image]
291     return image
292

```

```

293 def get_timesteps(self, num_inference_steps:int, strength:float):
294     """
295     Helper function for getting scheduler timesteps for generation
296     In case of image-to-image generation, it updates number of steps according to strength
297
298     Parameters:
299         num_inference_steps (int):
300             number of inference steps for generation
301         strength (float):
302             value between 0.0 and 1.0, that controls the amount of noise that is added to the input image
303             Values that approach 1.0 enable lots of variations but will also produce images that are not
304     """
305     # get the original timestep using init_timestep
306     init_timestep = min(int(num_inference_steps * strength), num_inference_steps)
307
308     t_start = max(num_inference_steps - init_timestep, 0)
309     timesteps = self.scheduler.timesteps[t_start:]

```

▼ Configure Inference Pipeline

First, you should create instances of OpenVINO Model.

```

1 from openvino.runtime import Core
2 core = Core()
3 text_enc = core.compile_model(TEXT_ENCODER_OV_PATH, 'AUTO')

```

```

1 unet_model = core.compile_model(UNET_OV_PATH, 'AUTO')

```

```

1 vae_decoder = core.compile_model(VAE_DECODER_OV_PATH, 'AUTO')
2 vae_encoder = core.compile_model(VAE_ENCODER_OV_PATH, 'AUTO')

```

Model tokenizer and scheduler are also important parts of the pipeline. Let us define them and put all components together

```

1 from transformers import CLIPTokenizer
2 from diffusers.schedulers import LMSDiscreteScheduler
3
4 lms = LMSDiscreteScheduler(
5     beta_start=0.00085,
6     beta_end=0.012,
7     beta_schedule="scaled_linear"
8 )
9 tokenizer = CLIPTokenizer.from_pretrained('openai/clip-vit-large-patch14')
10
11 ov_pipe = OVStableDiffusionPipeline(
12     tokenizer=tokenizer,
13     text_encoder=text_enc,
14     unet=unet_model,
15     vae_encoder=vae_encoder,
16     vae_decoder=vae_decoder,
17     scheduler=lms
18 )

```

Downloading (...)olve/main/vocab.json: 100%  961k/961k [00:00<00:00, 9.61MB/s]

Downloading (...)olve/main/merges.txt: 100%  525k/525k [00:00<00:00, 6.13MB/s]

Downloading (...)cial_tokens_map.json: 100%  389/389 [00:00<00:00, 25.3kB/s]

Downloading (...)okenizer_config.json: 100%  905/905 [00:00<00:00, 62.1kB/s]

▼ Text-to-Image generation

Now, you can define a text prompt for image generation and run inference pipeline. Optionally, you can also change the random generator seed for latent state initialization and number of steps.

Note: Consider increasing `steps` to get more precise results. A suggested value is 50, but it will take longer time to process.

```
1 import ipywidgets as widgets
2
3 text_prompt = widgets.Text(value='A city of Africa with tall buildings. a golden sun shining', description=
4 num_steps = widgets.IntSlider(min=1, max=50, value=20, description='steps:')
5 seed = widgets.IntSlider(min=0, max=100000000, description='seed: ', value=42)
6 widgets.VBox([text_prompt, seed, num_steps])
```

your text

seed: 42

steps: 50

```
1 print('Pipeline settings')
2 print(f'Input text: {text_prompt.value}')
3 print(f'Seed: {seed.value}')
4 print(f'Number of steps: {num_steps.value}')
```

```
Pipeline settings
Input text: A city of Africa with tall buildings. a golden sun shining
Seed: 42
Number of steps: 50
```

```
1 result = ov_pipe(text_prompt.value, num_inference_steps=num_steps.value, seed=seed.value)
```

100% 50/50 [02:53<00:00, 3.44s/it]

Finally, let us save generation results. The pipeline returns several results: `sample` contains final generated image, `iterations` contains list of intermediate results for each step.

```
1 final_image = result['sample'][0]
2 if result['iterations']:
3     all_frames = result['iterations']
4     img = next(iter(all_frames))
5     img.save(fp='result.gif', format='GIF', append_images=iter(all_frames), save_all=True, duration=len(all
6 final_image.save('result.png')
```

Now is show time!

```
1 import ipywidgets as widgets
2
3 text = '\n\t'.join(text_prompt.value.split('.'))
4 print("Input text:")
5 print("\t" + text)
6 display(final_image)
```

Input text:

cyberpunk cityscape like Tokyo New York with tall buildings at dusk golden hour cinematic lighting
A golden daylight, hyper-realistic environment
Hyper and intricate detail, photo-realistic
Cinematic and volumetric light
Epic concept art
Octane render and Unreal Engine, trending on artstation



```
1 import ipywidgets as widgets
2
3 text = '\n\t'.join(text_prompt.value.split('.'))
4 print("Input text:")
5 print("\t" + text)
6 display(final_image)
```

Input text:

A city of Paris with tall buildings
a golden sun shining



```
1 import ipywidgets as widgets
2
3 text = '\n\t'.join(text_prompt.value.split('.'))
4 print("Input text:")
5 print("\t" + text)
6 display(final_image)
```

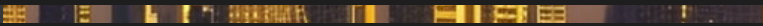


Input text:

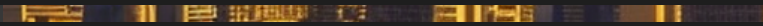
A city of Africa with tall buildings
a golden sun shining



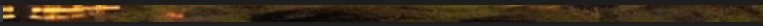
1



1



1



1

1

1

1

✓ 0s completed at 02:09

