

intel[®] ai
summit
英特爾 AI 科技論壇

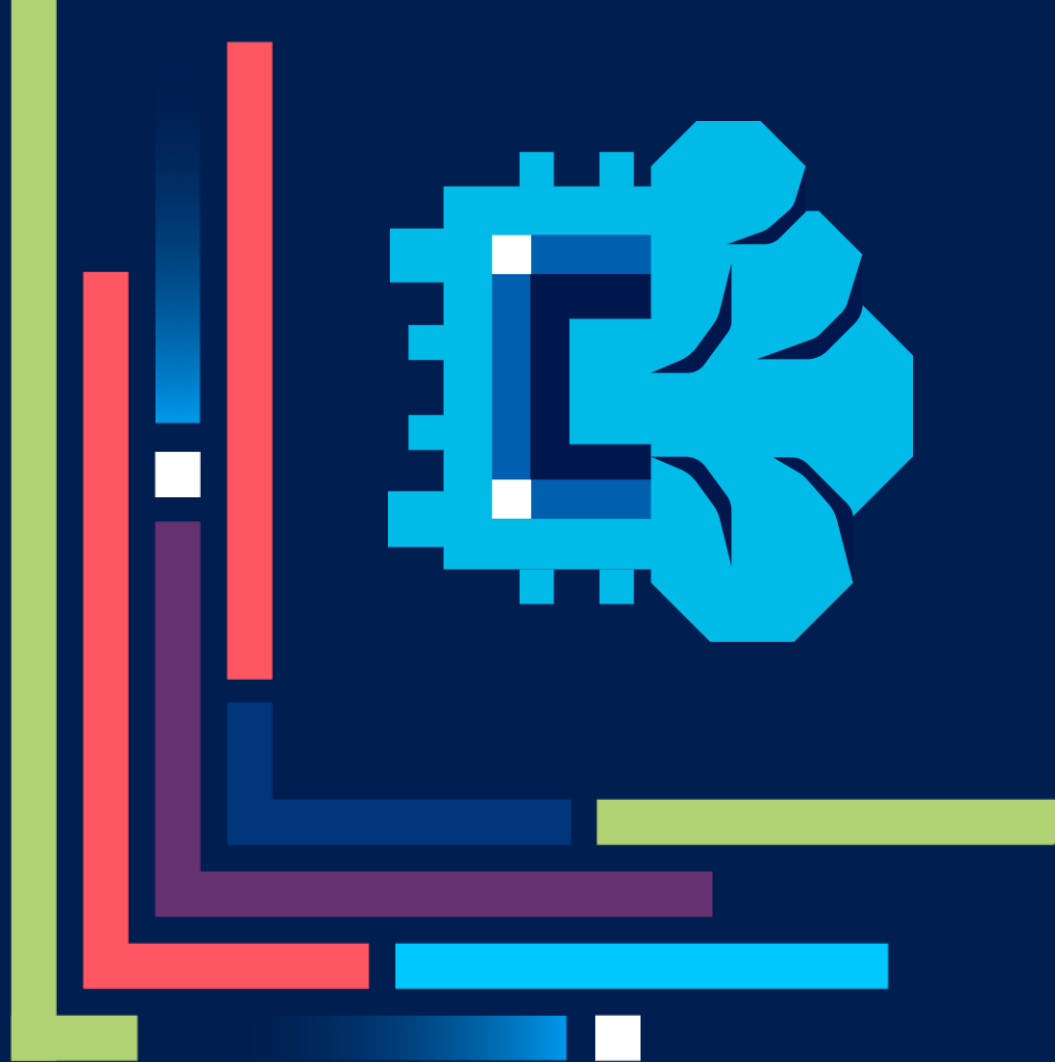
Bringing AI Everywhere

Generative AI Development and Application on Client

Nancy Cheng

Intel Sales Application Engineer

March 27th, 2024



Agenda

AI PC & What can OpenVINO™ do?

Riffusion

Converting/Enabling the models

Integration into Audacity® + Demo

Summary



Three AI Engines

Heterogenous execution of AI workloads embraces the best practices in AI software design.

GPU

High Throughput
Ideal for AI-accelerated digital content creation

NPU

Low Power
Ideal for sustained AI workloads and AI offload for battery life

CPU

Fast Response
Ideal for low-latency AI workloads

PyTorch

TensorFlow

Keras

TensorFlow Lite

ONNX

Caffe

PaddlePaddle

1 Model

OpenVINO™

2 Optimization

Optimized Performance

CPU

intel.
ATOM™

intel.
CORE™

intel.
XEON™

arm

NEW IN
2023

GPU

intel.
iRIS[®]xe

intel.
ARC™

intel.
DATA CENTER
GPU

GRAPHICS

GRAPHICS

FLEX SERIES

NPU

intel.
MOViDIUS™

FPGA

intel.
FPGA
AI Suite

3 Deployment

Windows

Linux

macOS

Software Optimizations

Additional
Configuration

Autobatching

Default
Inference
Precision

Threads
Scheduling

Shared
Memory

Model
Caching

Pre- Post-
Processing

Asynchronous
Mode

Model
Compression

OpenVINO™ Notebooks

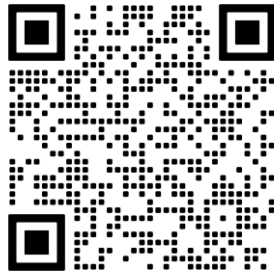


100+

Demos

LLMs, Stable Diffusion, Whisper, GPT,
YOLOv5/v8, CLIPS, Object Detection
and Segmentation, Image Classification,
Human Pose Estimation, and much more!

OpenVINO™ Notebooks Installation

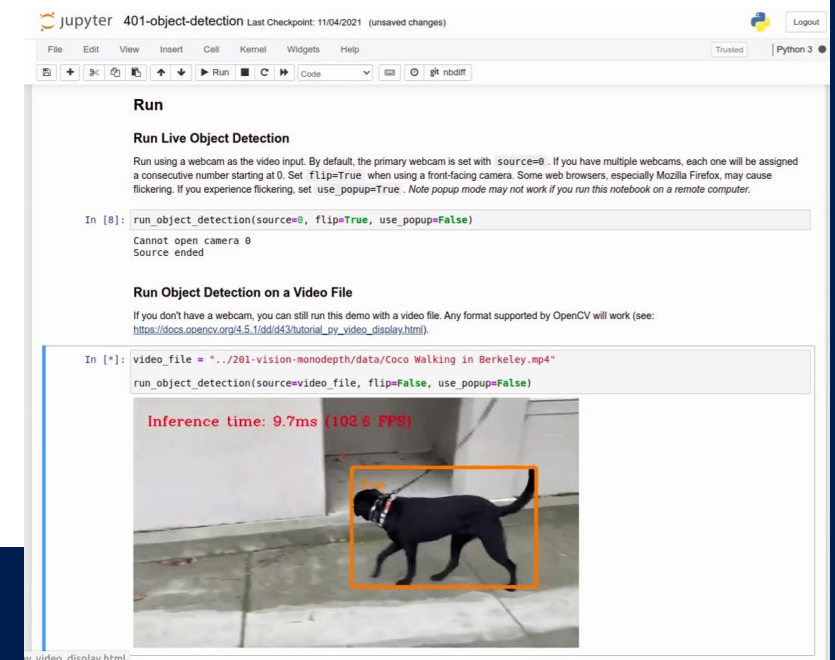


```
# Step 1: Clone the Repository
$ git clone
https://github.com/openvinotoolkit/openvino_notebooks.git
$ cd openvino_notebooks

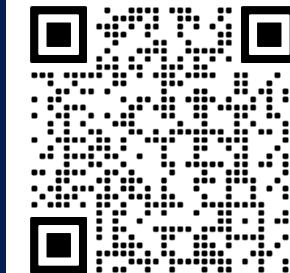
# Step 2: Create the Environment
$ python3 -m venv openvino_env
$ source openvino_env/bin/activate

# Step 3: Install Requirements
$ python -m pip install --upgrade pip
$ pip install -r requirements.txt

# Step 4: Launch the Notebooks
$ jupyter lab
```



Installation



Install OpenVINO™ 2024.0

Version	2024.0 Recommended	Nightly Build	2023.3 LTS	2022.3.1 LTS Includes NCS2/HDDL support		
Operating System	Windows		macOS		Linux	
Distribution	OpenVINO Archives Includes NPU plugin	PIP Includes NPU plugin Python API only	GitHub Source	Gitee Source	Docker	Conda
	vcpkg Source		Conan		npm JavaScript API only	
Install	<pre># Step 1: Create virtual environment python -m venv openvino_env</pre>					
	<pre># Step 2: Activate virtual environment openvino_env\Scripts\activate</pre>					
	<pre># Step 3: Upgrade pip to latest version python -m pip install --upgrade pip</pre>					
	<pre># Step 4: Download and install the package pip install openvino==2024.0.0</pre>					

[Installation Instructions](#) [Previous Releases](#)

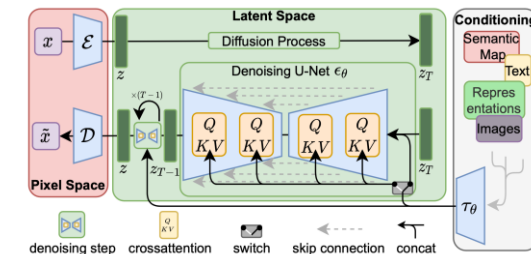
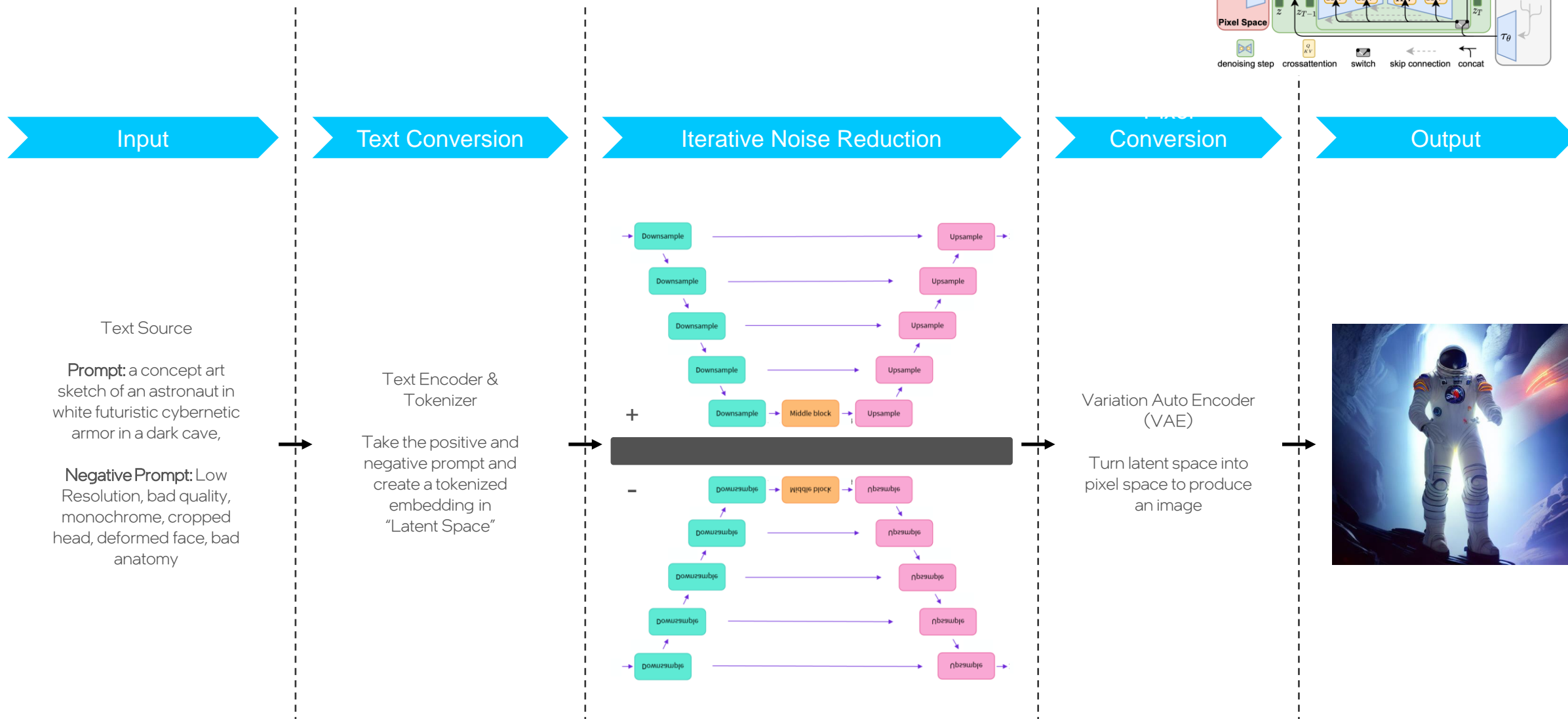
Advanced Optimization tool available separately: [Learn about NNCF](#)

OpenVINO™
NOTEBOOKS

Stable Diffusion



Stable Diffusion : Text to Image

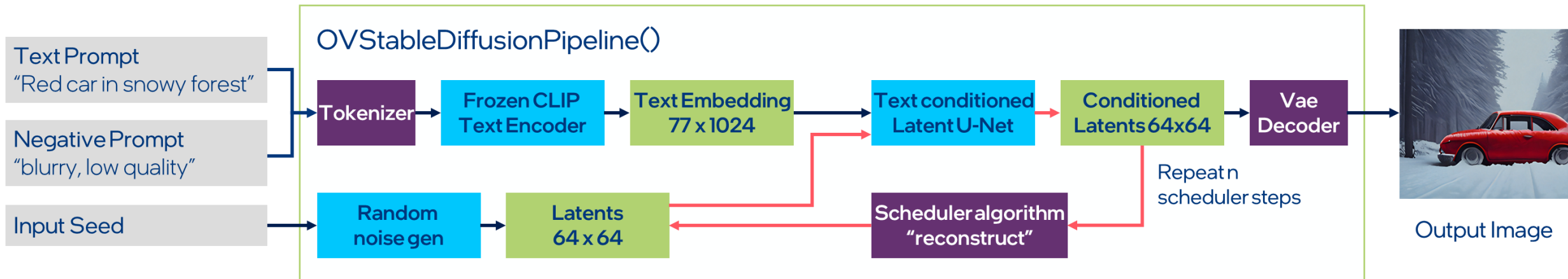


OpenVINO™ Stable Diffusion

High performance inference

SD Pipeline

Inference Pipeline

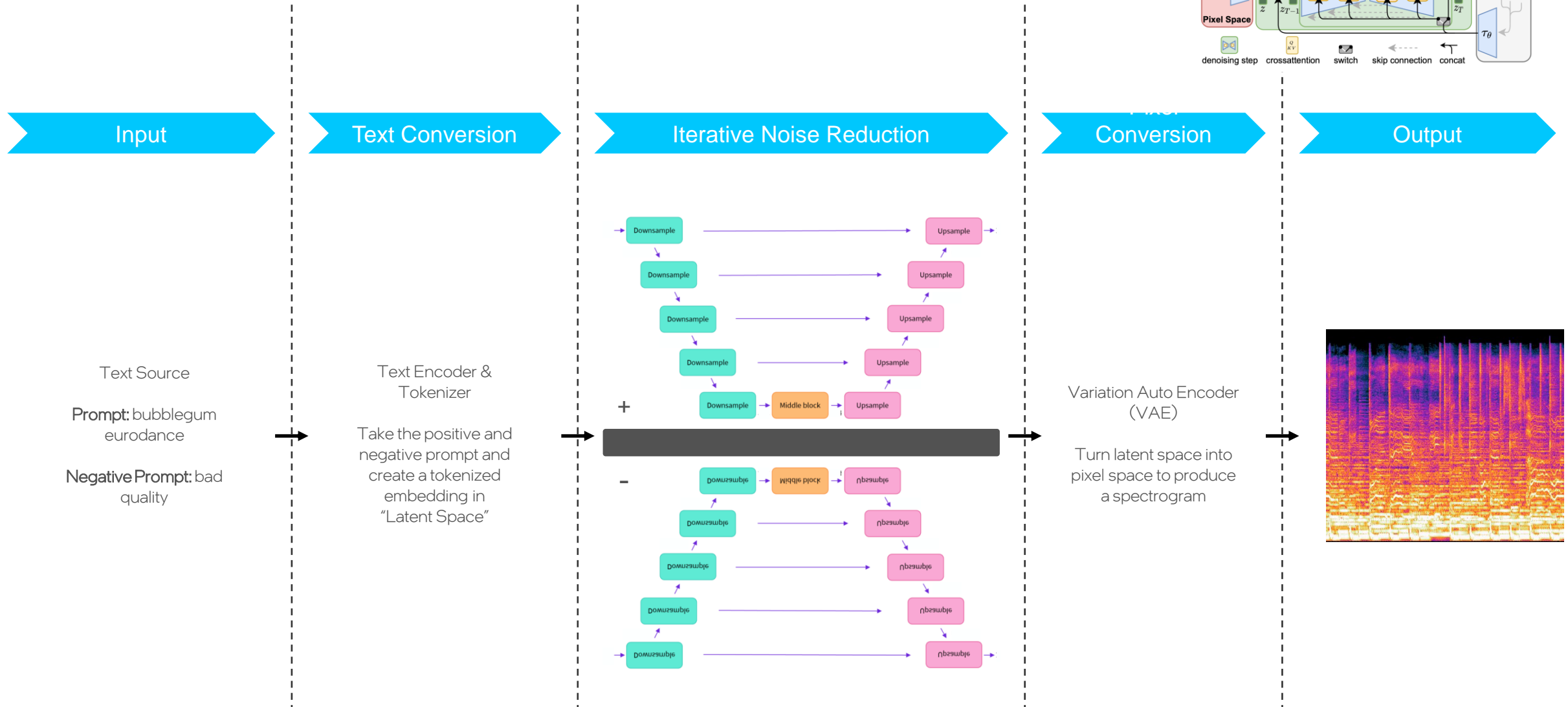


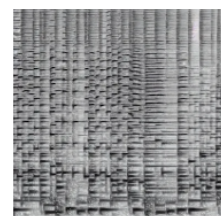
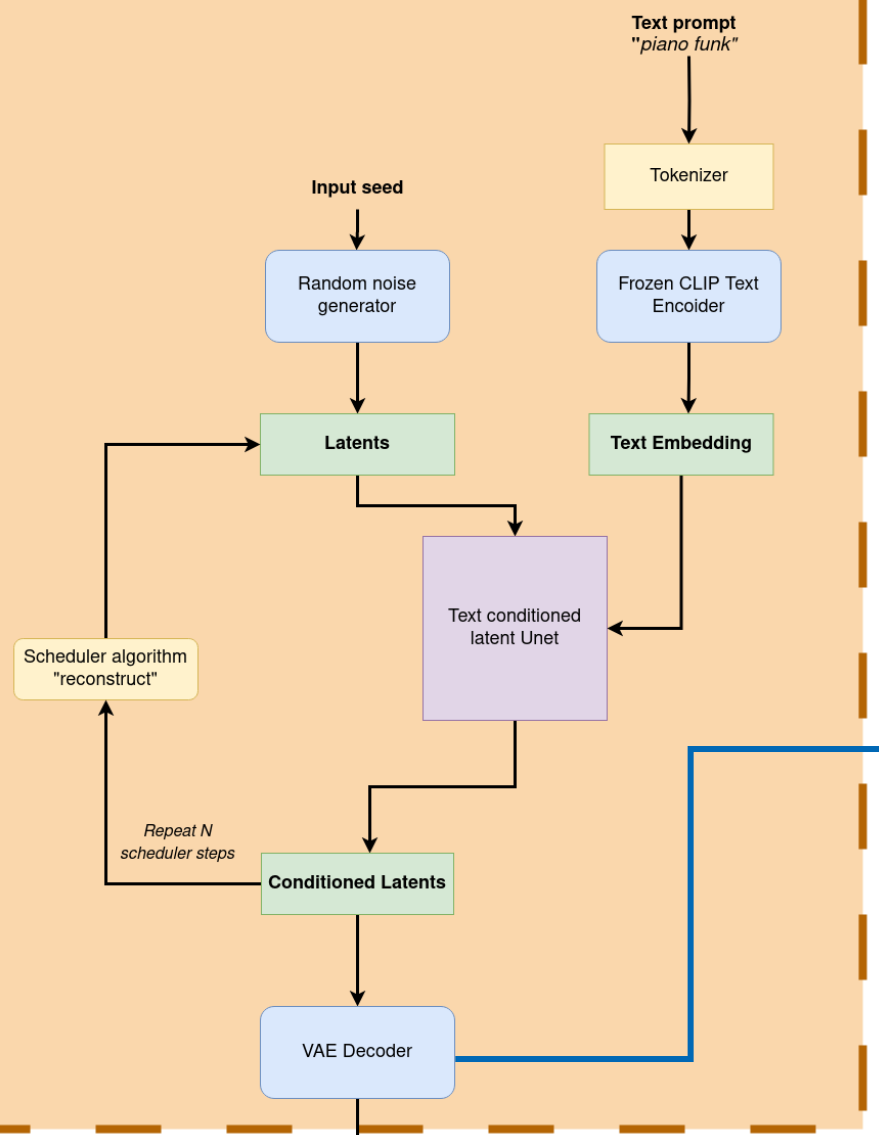


OpenVINO™ NOTEBOOKS

Riffusion Text-to-music

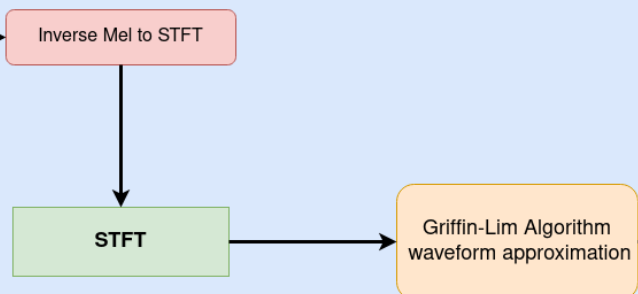
Riffusion : Text to Image (of audio!)





Output spectrogram

Post-processing



Reconstructed audio signal



Audacity®



Audacity®

Search ...



HOME

ABOUT ▾

DOWNLOAD ▾

HELP ▾

CONTACT ▾

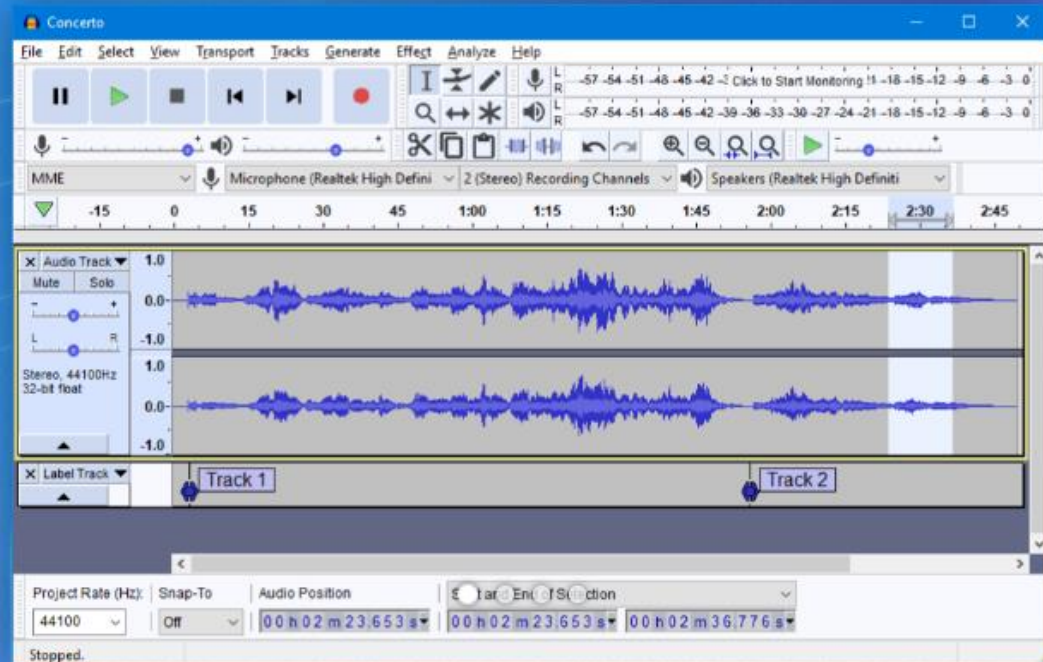
GET INVOLVED ▾

COPYRIGHT

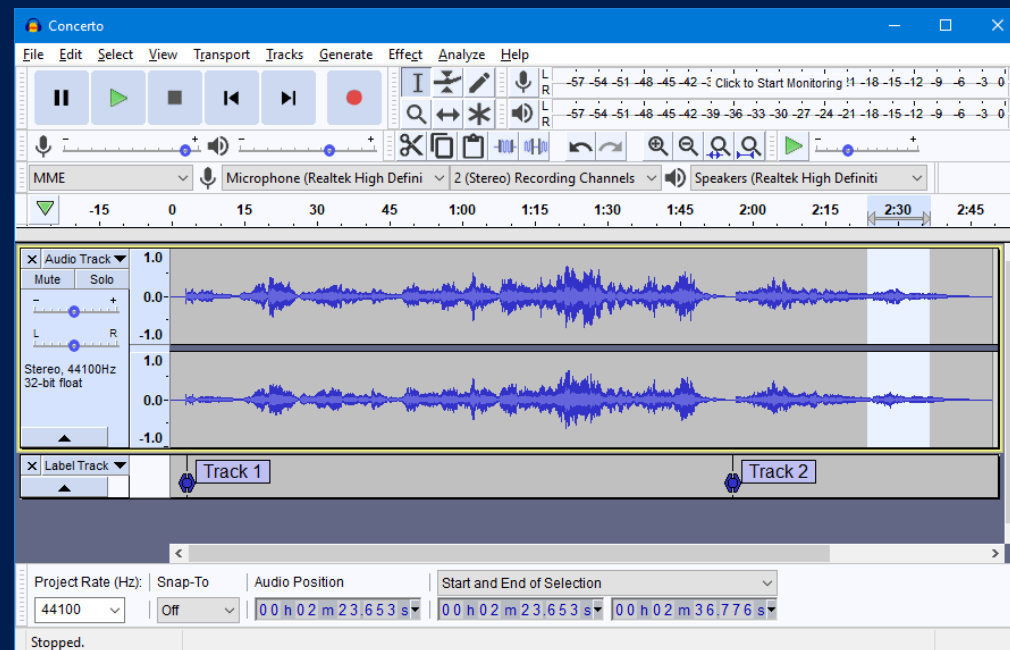
Free, open source,
cross-platform
audio software

Audacity is an easy-to-use, multi-track audio editor
and recorder for Windows, macOS, GNU/Linux and
other operating systems.

Audacity is free, open source software.



Riffusion in Audacity®

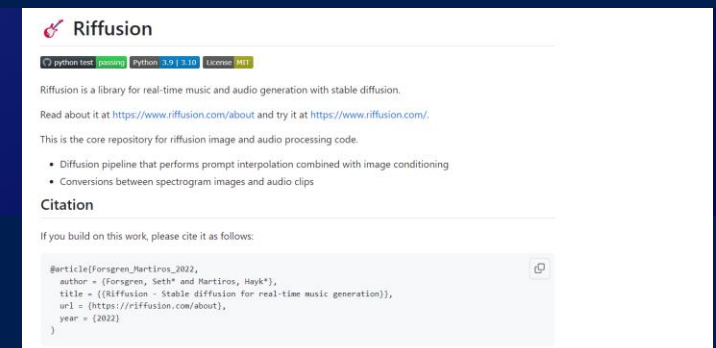
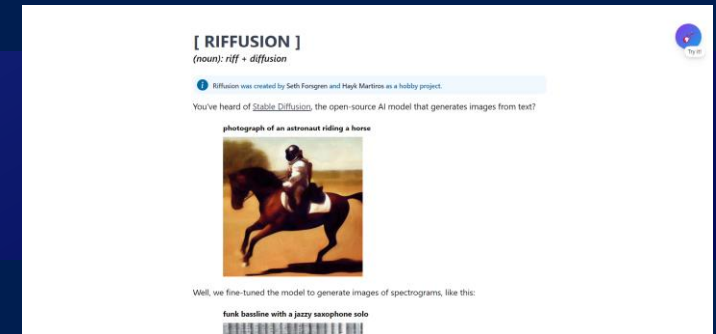
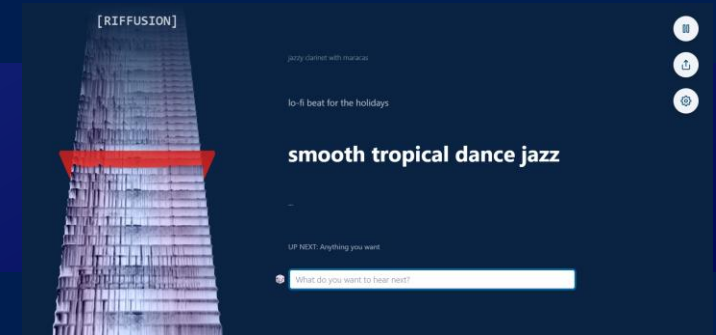


Acknowledgment to Riffusion Project

<https://www.riffusion.com>

<https://www.riffusion.com/about>

<https://www.github.com/riffusion/riffusion>



Integrating into Audacity® – Python to C++

Python APIs clearly map to C++ APIs, making the transition straightforward

Python APIs from Notebooks

```
from openvino.runtime import Core, Model

self.core = Core()

# Device to target such as CPU, GPU and now NPU
device = "NPU"

# Setup Caching to improve model load in future
self.core.set_property({'CACHE_DIR':
os.path.join(cache_dir, 'cache')})

# Produce a compiled-model object for a given device
self.compiled_model =
self.core.compile_model(model_path, device)

# Create the inference request
self.infer_request =
self.compiled_model.create_infer_request()
```

C++ APIs for integration into Application

```
#include <openvino/openvino.hpp>

ov::Core core;

// Device to target such as CPU, GPU and now NPU
std::string device = "NPU"

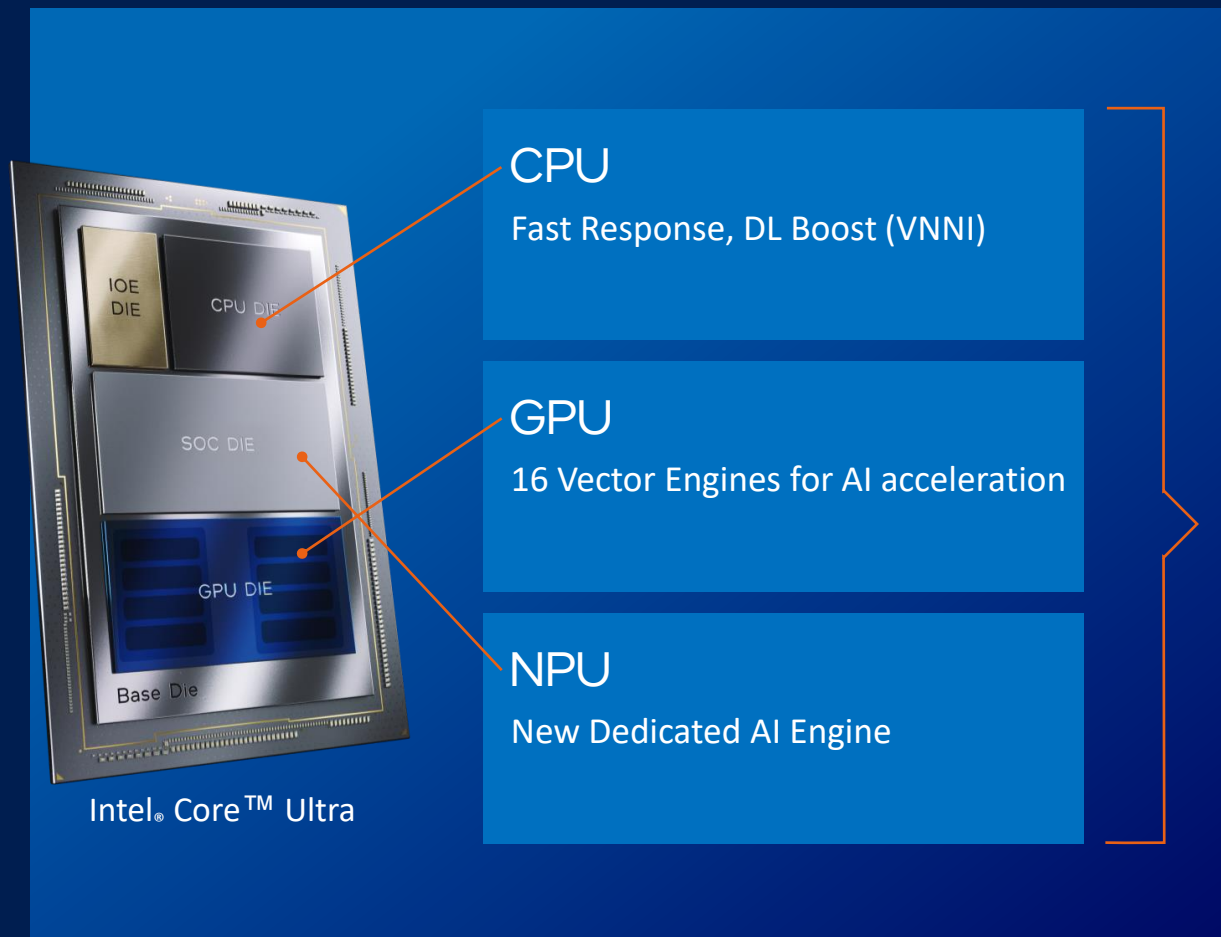
// Setup Caching to improve model load in future
core.set_property(ov::cache_dir("cache"));

// Produce a compiled-model object for a given device
ov::CompiledModel compiled_model =
core.compile_model(model_path, device);

// Create the inference request
ov::InferRequest infer_request =
compiled_model.create_infer_request();
```

Integrating into Audacity® – Using Different Device Types

OpenVINO™ supports multiple devices easily



C++ APIs for integration into Application

```
#include <openvino/openvino.hpp>

ov::Core core;

//Find all supported devices on this system
std::vector<std::string> devices =
core.get_available_devices();

// Device to target such as CPU, GPU and now NPU
std::string device = "NPU"

// Setup Caching to improve model load in future
core.set_property(ov::cache_dir("cache"));

// Produce a compiled-model object for a given device
ov::CompiledModel compiled_model =
core.compile_model(model_path, device);
```

Integrating into Audacity® – Using Different Model Types

OpenVINO™ supports many formats directly, avoiding the need for conversion

C++ APIs for integration into Application

```
#include <openvino/openvino.hpp>

ov::Core core;

// Device to target such as CPU, GPU and now NPU
std::string device = "NPU"

// We can choose between a bunch of different model formats
//std::string model_path = "path/to/text_encoder.xml";
//std::string model_path = "path/to/text_encoder.tflite";
std::string model_path = "path/to/text_encoder.onnx";

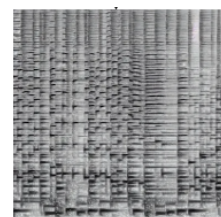
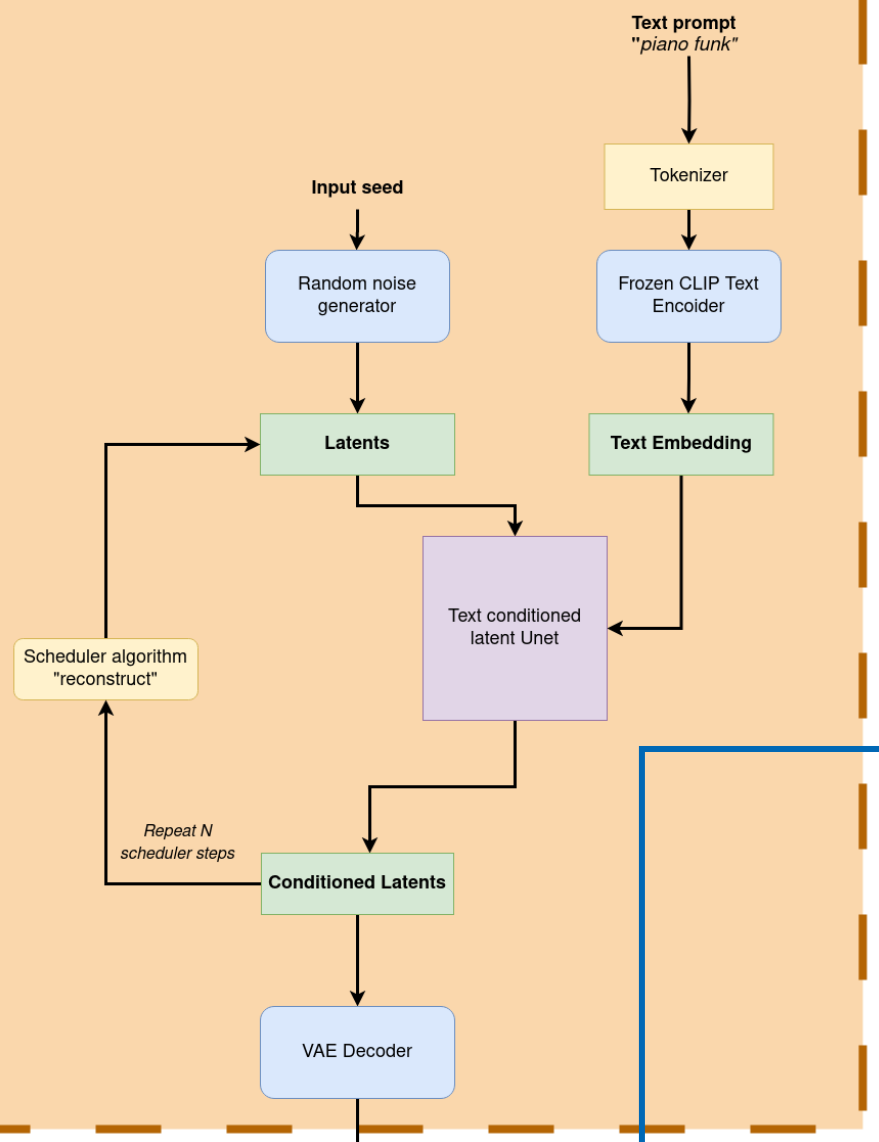
// Produce a compiled-model object for a given device
ov::CompiledModel compiled_model = core.compile_model(model_path, device);

// Create the inference request
ov::InferRequest infer_request = compiled_model.create_infer_request();
```



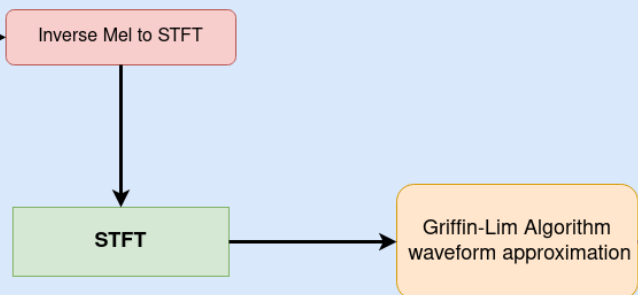
Caffe





Output spectrogram

Post-processing

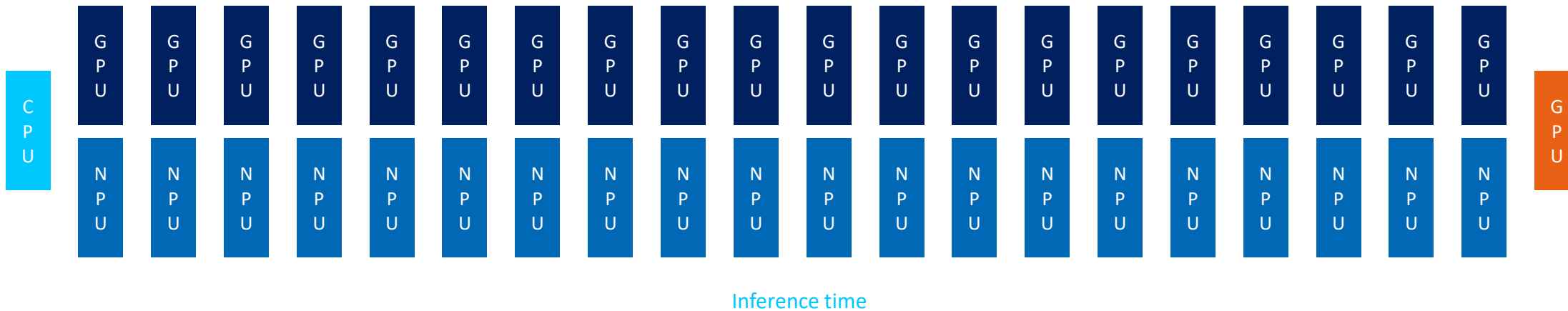


Reconstructed audio signal



Audacity® – Riffusion Model Device Assignments

Utilizing NPU in parallel with GPU maximizes throughput

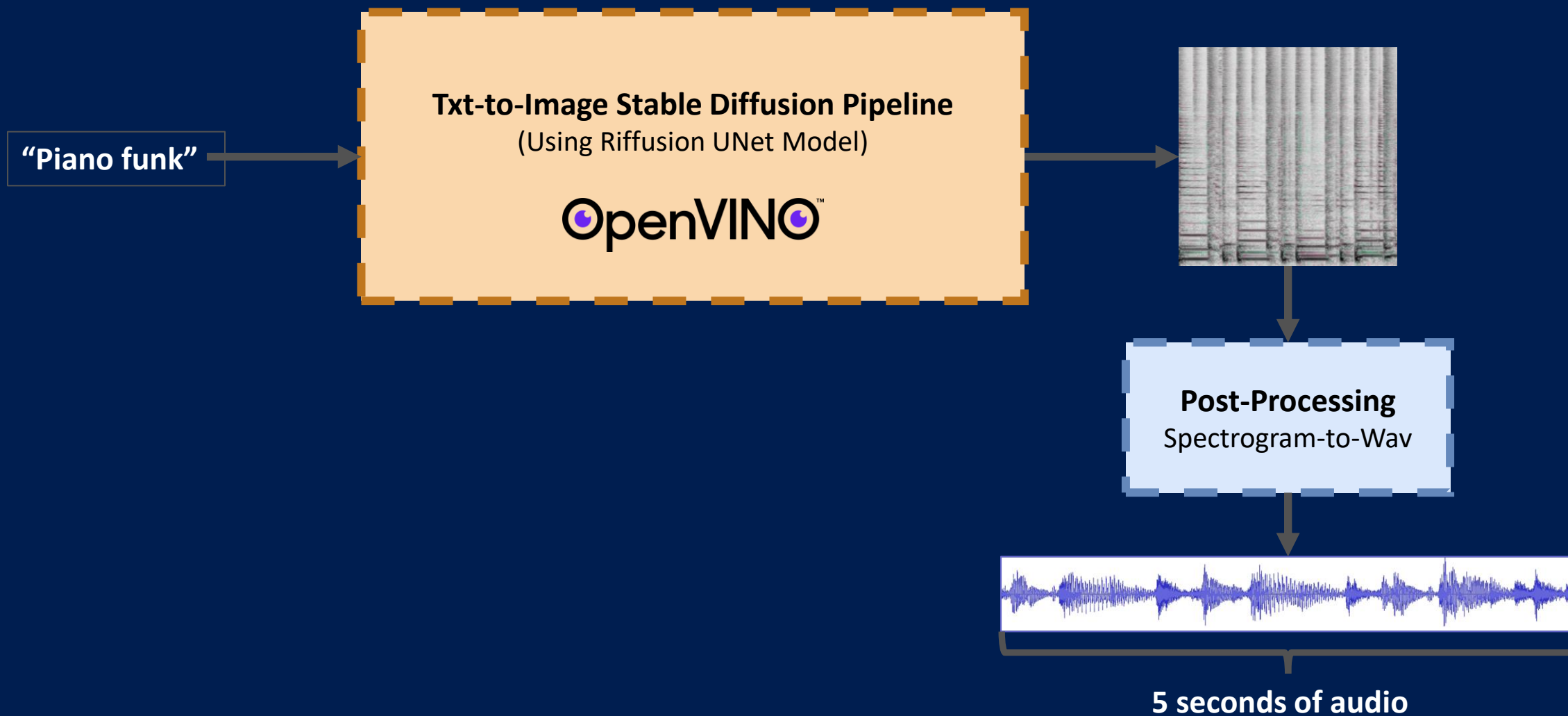


```
// Text Encoder on Core CPU
std::string device = "CPU"
std::string model_path = "path\to\text_encoder.onnx"
ov::CompiledModel text_enc_model = core.compile_model(model_path, device);
```

```
// VAE decoder on integrated GPU
device = "GPU"
model_path = "path\to\vae_decoder.onnx"
ov::CompiledModel vae_enc_model = core.compile_model(model_path, device);
```

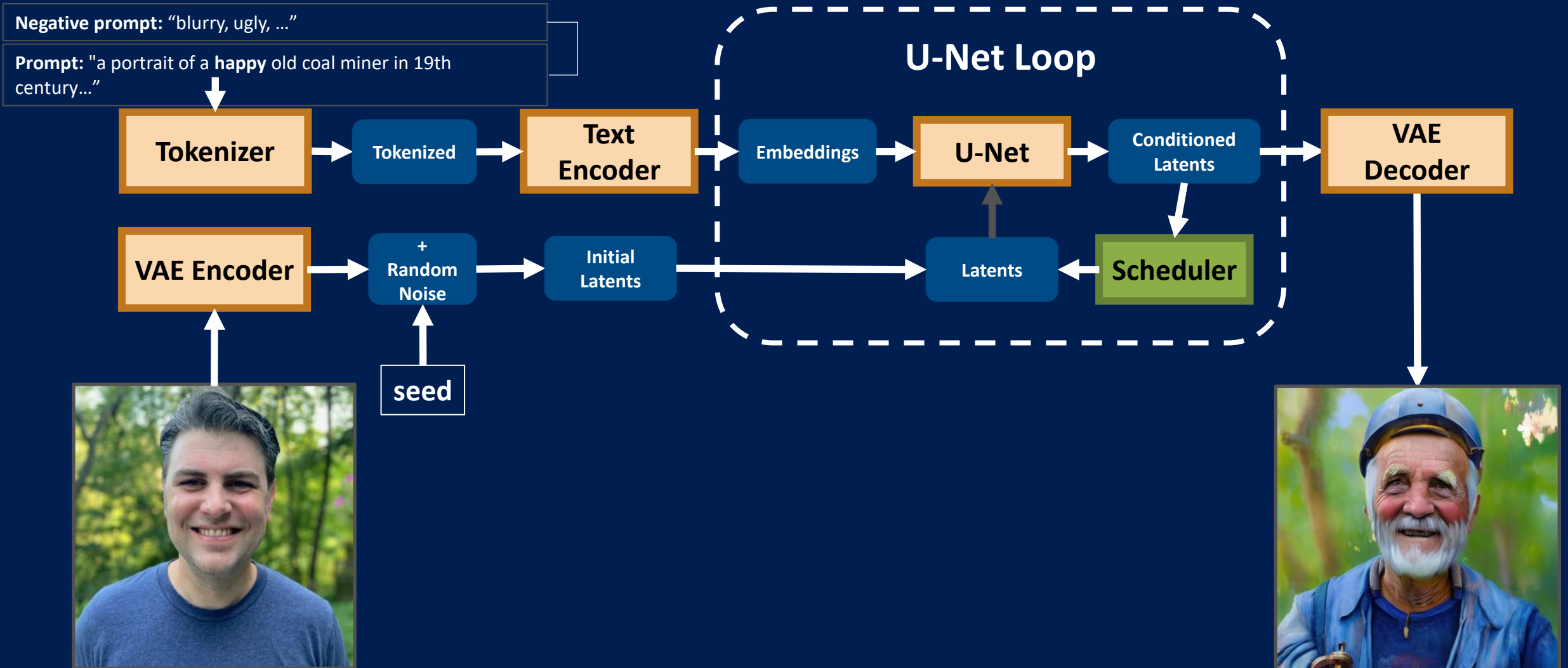
```
// Unet Positive Prompt on integrated GPU
device = "GPU"
model_path = "path\to\unet.xml"
ov::CompiledModel unet_pos_model = core.compile_model(model_path, device);
```

```
// Unet Negative Prompt on new integrated Neural Processor
device = "NPU"
model_path = "path\to\unet.xml"
ov::CompiledModel unet_neg_model = core.compile_model(model_path, device);
```

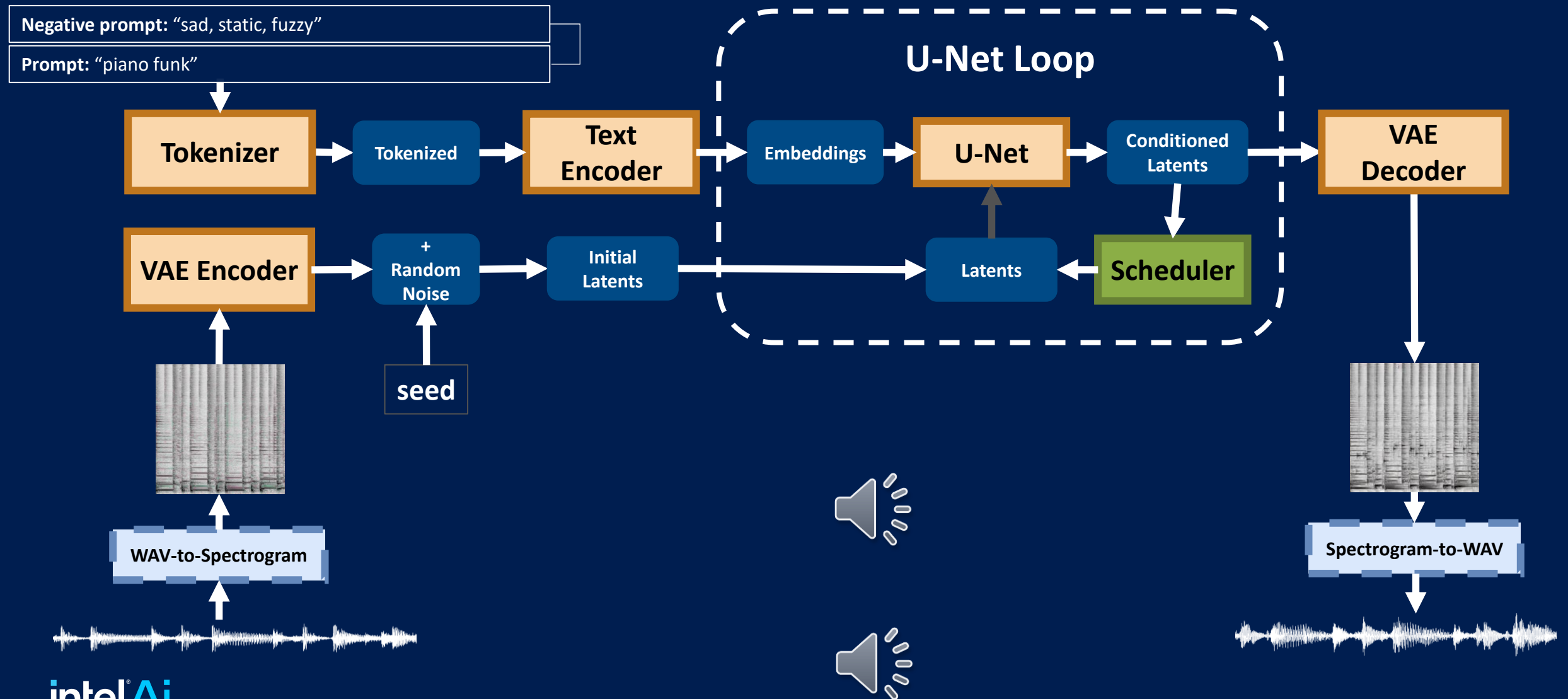



How can we generate audio segments longer than 5 seconds?

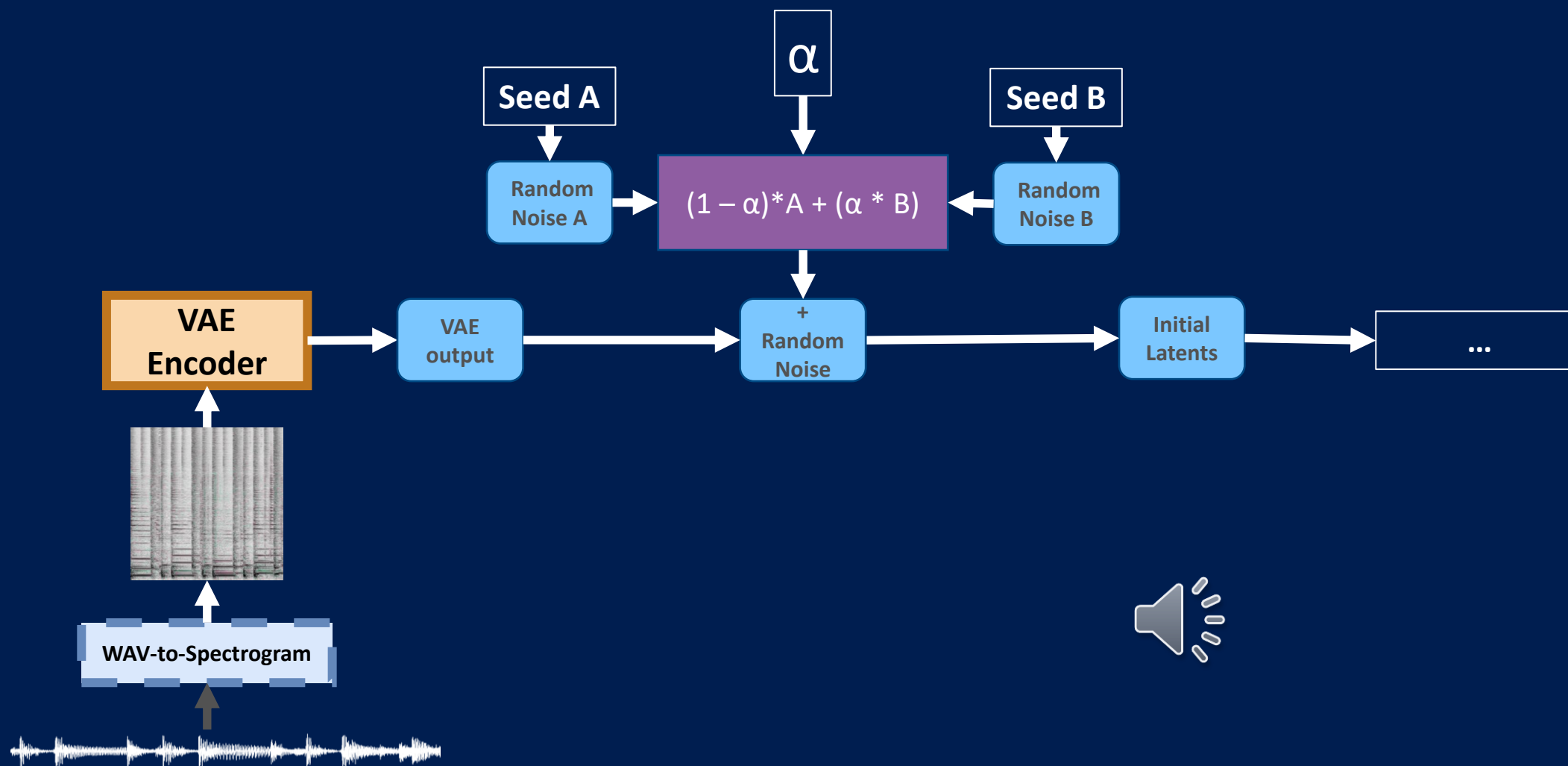
Stable Diffusion Image-to-Image Pipeline



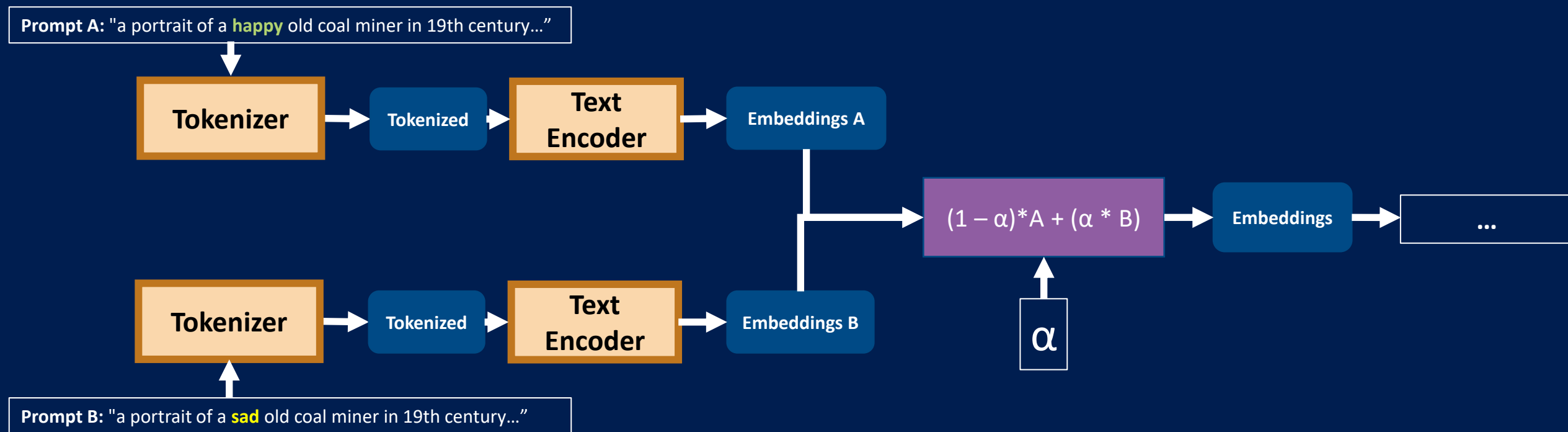
Stable Diffusion Image-to-Image Pipeline



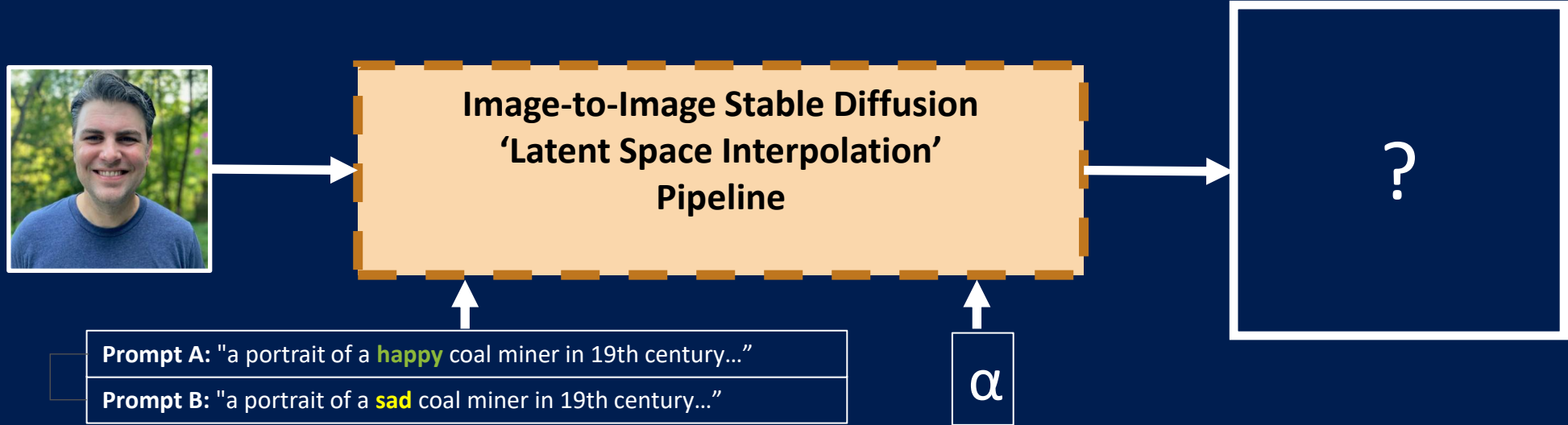
Creating Smooth Transitions – Seed Interpolation



Creating Smooth Transitions – Prompt Interpolation



Creating Smooth Transitions (In Images!)



$\alpha = 0.$



0.25



0.5



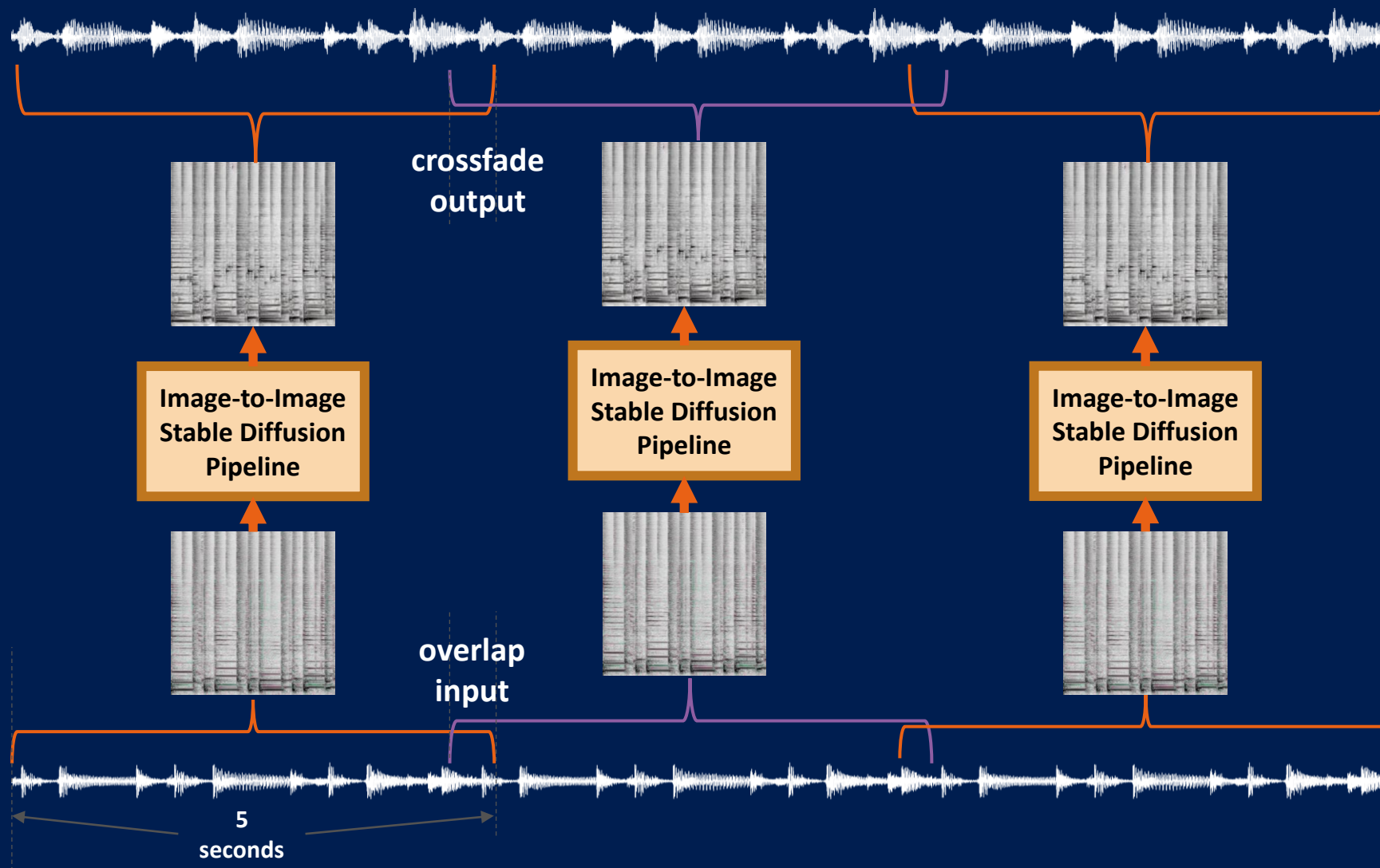
0.75



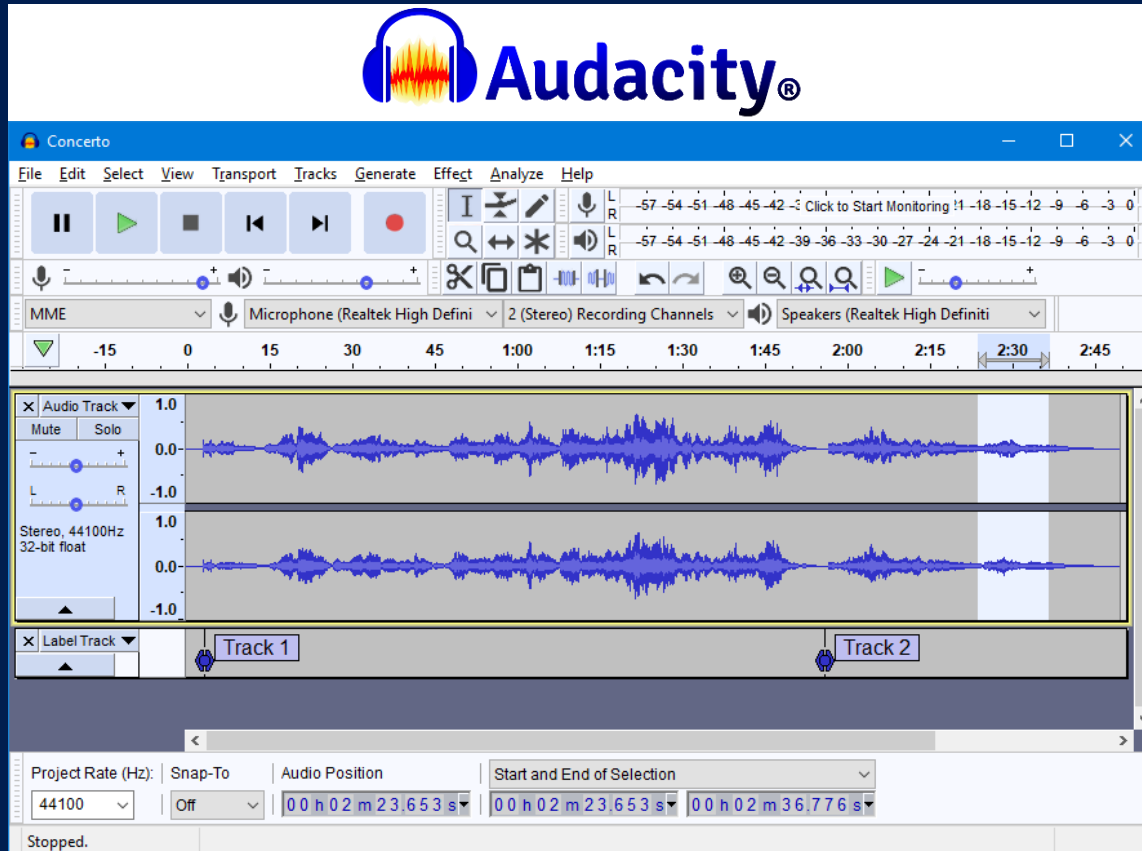
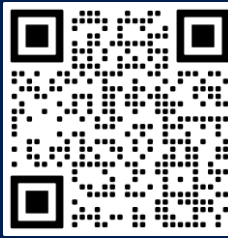
1.



Music Style Remix



OpenVINO™ AI Plugins for Audacity®



OpenVINO™ AI Plugins

Noise
Suppression

Music Separation

Transcription

wHISperC++

Music Generation
& Style Remix
[RIFUSION]

"We at Audacity are thrilled to be partnering with Intel to help bring powerful, open, and most importantly free AI tools to a mass audience. With time, we expect these kinds of initiatives to produce a new kind of creative environment for musicians, podcasters and audiophiles - a worthy successor to the traditional audio tools that have typified the last 20 years."

*—Martin Keary, Head of Product,
Audacity®*

OpenVINO™ AI Plugins for Audacity® Music Generation

The screenshot displays the Audacity software interface with the 'OpenVINO Music Generation' plugin. The 'Generate' menu item is selected, opening a dialog box with the following settings:

- Mode:** Simple (dropdown)
- Normalize:** -20.0 dB
- Duration:** 00 h 00 m 10.000 s
- What Kind of Music?:** piano funk
- Seed Image:** og_beat
- Strength:** 0.75
- Seed:** (empty field)
- Guidance Scale:** 7.5
- Num Inference Steps:** 20
- Text Encoder Device:** CPU
- UNet + Device:** GPU
- UNet - Device:** NPU
- VAE Decoder Device:** GPU
- VAE Encoder Device:** GPU
- Scheduler:** EulerDiscreteScheduler

The 'Generate' button is highlighted with a blue border. A secondary dialog box, 'OpenVINO Music Style Remix', is also visible, showing a progress bar and time information: Elapsed Time: 00:00:02, Remaining Time: 00:03:25. The 'Generate' button in this dialog is also highlighted with a blue border.

Numbered callouts 1 through 12 are present on the image, pointing to various UI elements:

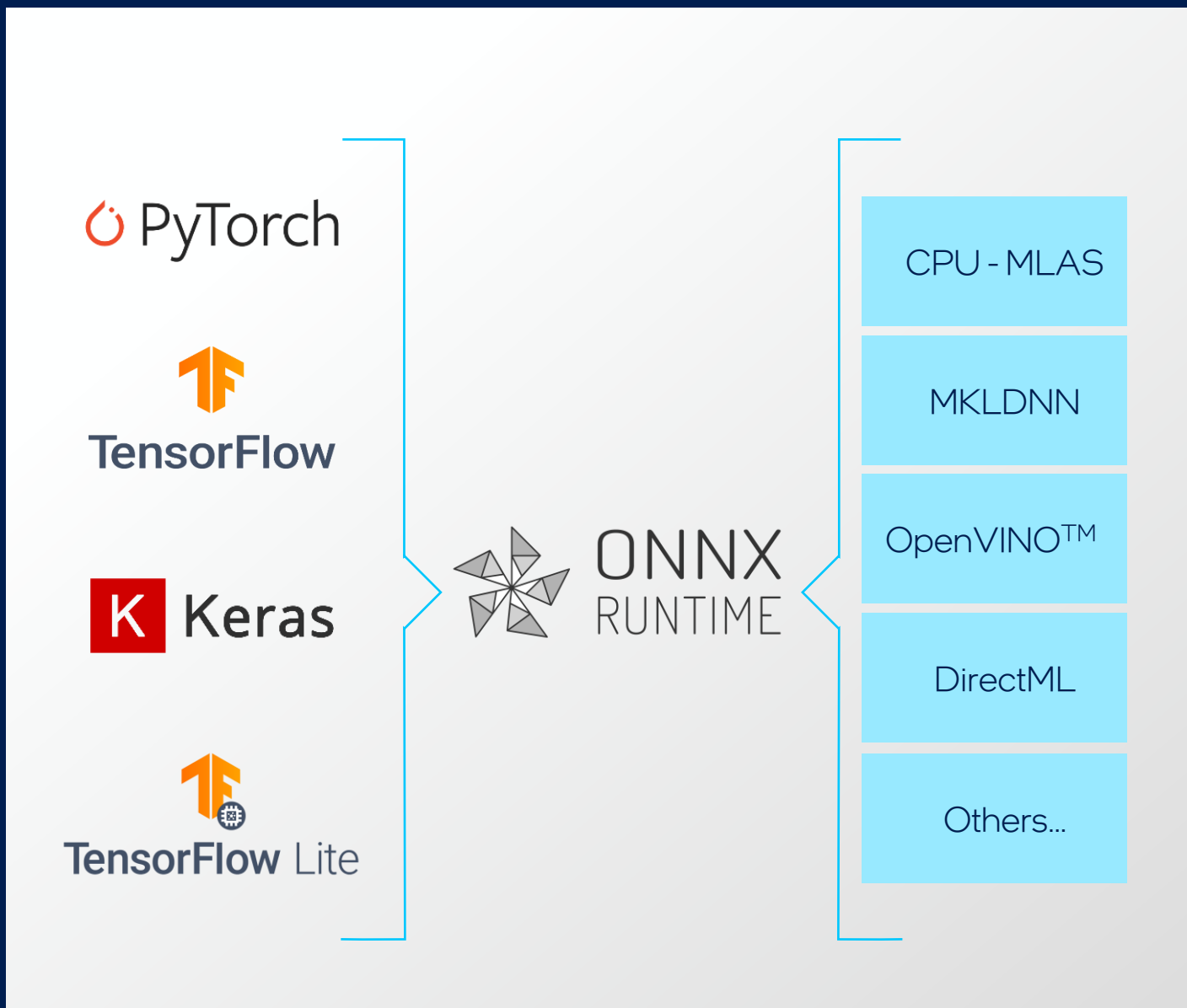
- 1: 'Generate' menu item
- 2: 'OpenVINO Music Generation...' menu item
- 3: 'Presets & settings' tab
- 4: 'What Kind of Music?' dropdown
- 5: 'Seed Image' dropdown
- 6: 'Strength' input field
- 7: 'Seed' input field
- 8: 'Guidance Scale' input field
- 9: 'Num Inference Steps' input field
- 10: 'Unload Models' button
- 11: 'Generate' button in the 'OpenVINO Music Style Remix' dialog
- 12: 'Generate' button in the 'OpenVINO Music Generation' dialog

Using ONNX* with OpenVINO™ EP

ONNX* enables cross platform development, utilizing execution providers for multiple hardware backends

OpenVINO™ Execution Provider for ONNX* provides optimized performance on Intel® Platforms

DirectML Execution provider enables cross platform scalability



Transforming the PC Experience through AI

Strong set of AI capabilities, including the new NPU for power efficient AI offload

Optimized Performance using OpenVINO™, ONNX*, and DirectML Frameworks

Transition from PoC to Productization easily using OpenVINO™ Notebooks



Intel® Core™ Ultra

CPU

Fast Response, DL Boost (VNNI)

Ideal for light-weight, single inference low-latency AI tasks

GPU

Performance Parallelism & Throughput, 16 Vector Engines

Ideal for AI infused in Media/3D/Render pipeline

NPU

Dedicated AI Engine for New & Improved Experiences

Ideal for power-efficient, sustained AI and AI offload

SW Tools



ONNX

OpenVINO™



DirectML

What's Next?

Download and try it for yourself!

https://github.com/openvinotoolkit/openvino_notebooks



Audacity® plugins to be published:

<https://github.com/intel/openvino-plugins-ai-audacity>



Explore the many Intel supported open-source projects on Github:

- Stable Diffusion for GIMP

<https://github.com/intel/openvino-ai-plugins-gimp>



- Plugins for OBS Studio

<https://github.com/intel/openvino-plugins-for-obs-studio>



- Many more:

<https://github.com/intel>



OpenVINO™ Documentation:

<https://docs.openvino.ai/>





The Art of Performance: Inference Optimization of Gen AI and LLMs with OpenVINO™ on AI PC

Zhuo Wu

AI Software Evangelist

intel.

Notices and Disclaimers

For notices, disclaimers, and details about performance claims, visit www.intel.com/PerformanceIndex or scan the QR code:



© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

intel[®] Ai
summit

Thank You!

