# Agenda

- Workshop pre-requisites
- Intel® AI Optimizations
- AI Handson Workshop

# Workshop Pre-requisites

- Register for accessing Intel® Developer Cloud (5 mins)
  - visit ➔ cloud.intel.com
  - Sign up --> Create Account
- Setup SSH access to Intel® Developer Cloud (5 mins)
  - Login to Intel® Developer Cloud
- Clone the workshop repository on Developer cloud
  - https://tinyurl.com/oneapi-ai-workshop
- Laptop with open internet access (preferred)

# Intel® Developer Cloud

# Intel® Developer Cloud

**a service platform for developing and running workloads in Intel®-optimized deployment environments with the latest Intel® processors**

- Landing page :
  - https://cloud.intel.com

- Instructions to get started:
  - **http://tinyurl.com/ReadmeIDC**

# Intel® AI Optimizations

# Diverse Compute Requirements

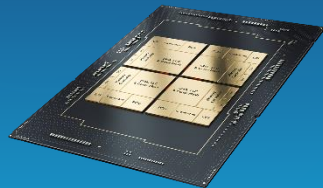**Diverse accelerators needed to meet today's performance requirements:**

48% of developers target heterogeneous systems
that use more than one kind of processor or core[1]

**Developer Challenges: Multiple Architectures, Vendors, and Programming Models**
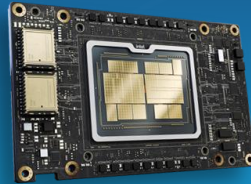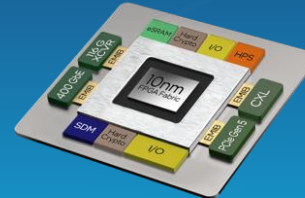
**Open, Standards-based, Multiarchitecture Programming**

oneAPI

**CPU**

**GPU**

**FPGA**

**Other Accelerators**

1 - https://evansdata.com/reports/viewRelease.php?reportID=40

# oneAPI Industry Initiative

## Break the Chains of Proprietary Lock-in

### Freedom to Make Your Best Choice

- C++ programming model for multiple architectures and vendors
- Cross-architecture code reuse for freedom from vendor lock-in

### Realize all the Hardware Value

- Performance across CPU, GPUs, FPGAs, and other accelerators
- Expose and exploit cutting-edge features of the latest hardware

### Develop & Deploy Software with Peace of Mind

- Open industry standards provide a safe, clear path to the future
- Interoperable with familiar languages and programming models including Fortran, Python, OpenMP, and MPI
- Powerful libraries for acceleration of domain-specific functions



**Application Workloads Need Diverse Hardware**

**Middleware & Frameworks**

TensorFlow   PyTorch   learn   NumPy   XBOOST   OpenVINO   …

**oneAPI Industry Specification**

Direct Programming

SYCL (C++)

API-Based Programming

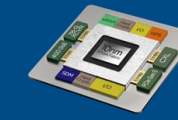| Math | Threading | Parallel STL | Ray Tracing |
| Analytics/ML | DNN | ML Comm | Volumetric Rendering |
| Video Processing | Signal Processing | Image Processing | Image Denoise |

Low-Level Hardware Interface (oneAPI Level Zero)

CPU   GPU   FPGA   Other Accelerators

The productive, smart path to freedom for accelerated computing from the economic and technical burdens of proprietary programming models

# oneAPI Industry Momentum

## End Users

@rchanan · Accrad Accelerated Radiology · FOUNDRY MODO · BioDataAnalysis GmbH — Next Generation Bio-Medical Image Analysis · GeoEast · BENTLEY · CINESITE · LAIKA · GE · kt · AUTODESK ARNOLD · 中国石油集团东方地球物理勘探有限责任公司 BGP INC.,CHINA NATIONAL PETROLEUM CORPORATION · iOmniscient Autonomous AI · EURECOM Sophia Antipolis · Brightskies · PHILIPS · 大势智慧 DASPATIAL · WeBank · allegro.ai · ILLUMINATION MACGUFF · Verizon · MediaKind · SAMSUNG MEDISON · TANGENT ANIMATION

## National Labs

Argonne NATIONAL LABORATORY · UT-BATTELLE Oak Ridge National Laboratory · CERN openlab · CINECA · Peraton Labs · Laboratório Nacional de Computação Científica · lrz Leibniz Supercomputing Centre of the Bavarian Academy of Sciences and Humanities · SANKHYA SUTRA Labs

## ISVs & OSVs

codeplay · AI SINGAPORE · CANONICAL · QITI CERTIFACE · ANACONDA · Ansys · MAX PLANCK COMPUTING & DATA FACILITY · MPCDF · SAS · ES · Red Hat OpenShift Data Science · KATANA GRAPH · CGG · AIBLE · CHAOSGROUP · MAXON · spirent · FOUNDRY · mercenaries engineering Guerilla · Hisense Medical 海信医疗 · SENAI CIMATEC FIEB SISTEMA · SAP · E4 COMPUTER ENGINEERING · YUAN · GIGASPACES innovate with confidence · MEGH COMPUTING · KFBIO technology for health. · Tech Mahindra · AsiaInfo 亚信科技 · SUSE We adapt. You succeed. · vmware · UNITED IMAGING

## OEMs & SIs

BittWare a molex company · Hewlett Packard Enterprise · DELL Technologies · Atos · Lenovo · HCL · rENIAC · MEGWARE SUPERCOMPUTING · TECHNOLOGY

## Universities & Research Institutes

LOBACHEVSKY UNIVERSITY · TECHNION Israel Institute of Technology · UNIVERSITY OF CAMBRIDGE · CTG · 中国科学院计算技术研究所 INSTITUTE OF COMPUTING TECHNOLOGY,CHINESE ACADEMY OF SCIENCES · Ben-Gurion University of the Negev · Berkeley UNIVERSITY OF CALIFORNIA · UCDAVIS UNIVERSITY OF CALIFORNIA · UNIVERSIDAD COMPLUTENSE MADRID · OLD DOMINION UNIVERSITY · 北京大学软件与微电子学院 School of Software & Microelectronics · ZIB · THE UNIVERSITY OF TENNESSEE KNOXVILLE · THE UNIVERSITY OF UTAH · ILLINOIS · University of Stuttgart Germany · UNIVERSITY OF OREGON · TÉCNICO LISBOA · inesc id lisboa · University College London · Durham University · FACULTY OF MATHEMATICS AND PHYSICS Charles University · PURDUE UNIVERSITY Elmore Family School of Electrical and Computer Engineering · Indian Institutes of Technology Delhi / Kharagpur / Roorkee · SDSC SAN DIEGO SUPERCOMPUTER CENTER · CDAC · UKRI Science and Technology Facilities Council Scientific Computing · TACC · Indian Institute of Science Bangalore · Indian Institute of Science Education & Research Pune · IISER PUNE · UNIVERSIDAD DE MÁLAGA · URZ HEIDELBERG UNIVERSITY COMPUTING CENTRE · Northern Illinois University NIU · University of BRISTOL · Stockholm University · KTH

## CSPs & Frameworks

Microsoft Azure · Alibaba Cloud · Google Cloud · TensorFlow · Taboola · Tencent 腾讯 · DataRobot · Baidu 百度 · PaddlePaddle 飞桨 · NAVER CLOVA · PyTorch

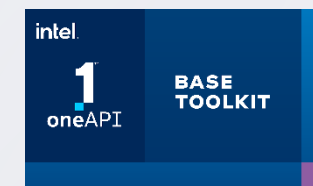**oneAPI DevSummit South-East Asia 2023**

# Intel® oneAPI Toolkits

## A complete set of proven developer tools expanded from CPU to Accelerators
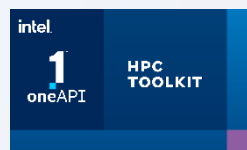
### Intel® oneAPI **Base** Toolkit

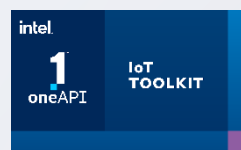A core set of high-performance libraries and tools for building C++, SYCL and Python applications

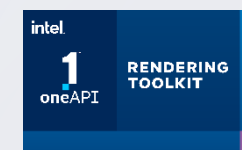## Add-on Domain-specific Toolkits

**Intel® oneAPI Tools for HPC**
Deliver fast Fortran, OpenMP & MPI applications that scale

**Intel® oneAPI Tools for IoT**
Build efficient, reliable solutions that run at network's edge

**Intel® oneAPI Rendering Toolkit**
Create performant, high-fidelity visualization applications

## Toolkits powered by oneAPI

**Intel® AI Analytics Toolkit**
Accelerate machine learning & data science pipelines end-to-end with optimized DL frameworks & high-performing Python libraries

**OpenVINO™**

**Intel® Distribution of OpenVINO™ Toolkit**
Deploy high performance inference & applications from edge to cloud

Latest version available 2023.1

# Intel® AI Analytics Toolkit

Accelerate end-to-end AI and data analytics pipelines with libraries optimized for Intel® architectures

## Who needs this product?

Data scientists, AI researchers, ML and DL developers,
AI application developers

## Top Features/Benefits

- Deep learning performance for training and inference with Intel ® Optimized DL frameworks and tools

- Drop-in acceleration for data analytics and machine learning workflows with compute-intensive Python packages

intel.
AI ANALYTICS TOOLKIT

## Intel® AI Analytics Toolkit

### Deep Learning

Intel® Optimization for TensorFlow

Intel® Optimization for PyTorch

Intel® Neural Compressor

Model Zoo for Intel® Architecture

### Machine Learning

Intel® Extension for Scikit-learn

Intel-optimized XGBoost

### Data Analytics

Intel® Distribution of Modin

OmniSci Backend

### Intel-optimized Python

| NumPy | SciPy | Numba | Pandas | Data Parallel Python |

CPU     GPU

Hardware support varies by individual tool. Architecture support will be expanded over time.

Get the Toolkit HERE or via these locations
https://software.intel.com/content/www/us/en/develop/tools/oneapi/download.html#aikit

| Intel ® Installer | Docker | Apt, Yum | Conda | Intel® DevCloud |

oneAPI

# Intel® Optimization for Tensorflow

## What's New for TensorFlow Optimization?

- oneDNN is in official TensorFlow release!

-  The platforms use the Intel® oneAPI Deep Neural Network Library (oneDNN), an open-source, cross-platform performance library for Deep-Learning applications

- Enable those Intel ® oneDNN CPU optimizations by setting the environment variable TF_ENABLE_ONEDNN_OPTS=1 for the official x86-64 TensorFlow after v2.5.

- Since TensorFlow **v2.9** and above, the oneAPI Deep Neural Network Library (oneDNN) optimizations are enabled by default

## Features

- Operator optimizations: Replace default (Eigen) kernels by highly-optimized kernels (using Intel® oneDNN)

- Graph optimizations: Fusion, Layout Propagation

- System optimizations: Threading model

# Intel® Optimization for Tensorflow

## Features

- Supports  FP32, FP16, Bfloat16, and int8.
- Leverages Intel® DL Boost, AVX512 instructions and processor capabilities
- Fused operations for optimized performance

## Support Matrix

- Compilers: Intel® oneAPI DPC++ / C++ Compilers
- OS: Linux, Windows, macOS
- CPU: Intel ® Atom, Intel® Core™, Intel® Xeon®, Intel® Xeon® Scalable processors
- GPU: Intel® Processor Graphics Gen9, Intel® Processor Graphics Gen 12

| Category | Functions |
|---|---|
| Compute intensive operations | • (De-)Convolution<br>• Inner Product<br>• RNN (Vanilla, LSTM, GRU)<br>• GEMM |
| Memory bandwidth limited operations | • Pooling<br>• Batch Normalization<br>• Local Response Normalization<br>• Layer Normalization<br>• Elementwise<br>• Binary elementwise<br>• Softmax<br>• Sum<br>• Concat<br>• Shuffle |
| Data manipulation | • Reorder |

# oneDNN Integration with TensorFlow

## Features

- Replaces compute-intensive standard TF ops with highly optimized custom oneDNN ops
- Aggressive op fusions to improve performance of Convolutions and Matrix Multiplications
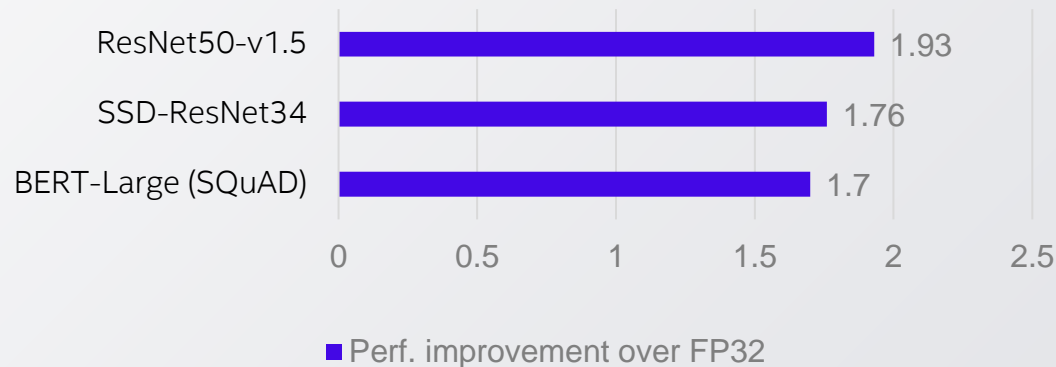- Primitive caching to reduce overhead of calling oneDNN Graphics Gen 12



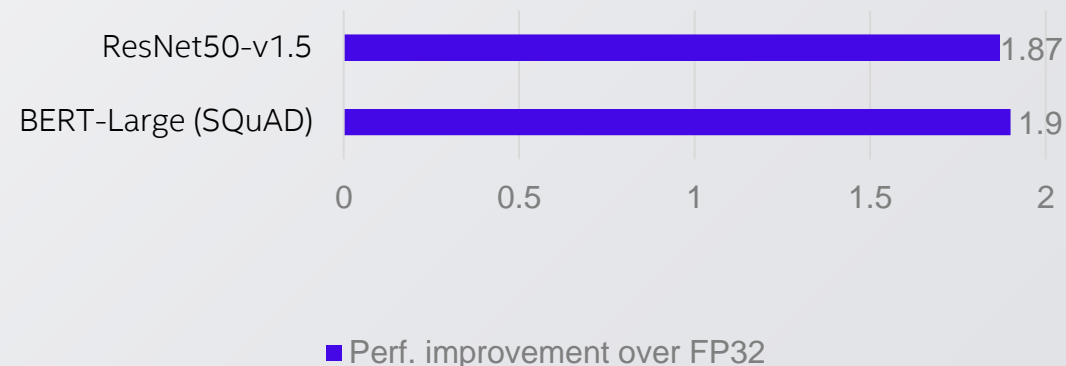Aggressive Fusions of BERT ops (MatMul + BiasAdd + GELU) into oneDNN op

https://www.intel.com/content/www/us/en/docs/onednn/developer-guide-reference/2023-1/gelu.html

# BFloat16 Data Type



FP32

| 8 bits | 23 bits |

S E E E E E E E E M M M M M M M M M M M M M M M M M M M M M M M

Same dynamic range and simple conversion

BF16

S E E E E E E E E M M M M M M M

7 bits

Delivers required level of accuracy

# Bfloat16 Optimization

- Bfloat16 - 16-bit data type with the same dynamic range as FP32

- Benefits
  - Reduced bandwidth
  - Improved performance with hardware support

- Easy to use
  - No special handling for loss scaling
  - No hyperparameter tuning for training, can reuse FP32 hyperparameters

- Up to 2x improvement on training and inference with negligible accuracy loss (< 0.20%)

- AMP (Automatic Mixed Precision) in tensorflow automatically converts model to use bfloat16 data type.

- Supports both Keras and arbitrary graph based models.

## Mixed precision training with Bfloat16

| Model | Value |
|---|---|
| ResNet50-v1.5 | 1.93 |
| SSD-ResNet34 | 1.76 |
| BERT-Large (SQuAD) | 1.7 |

0    0.5    1    1.5    2    2.5

■ Perf. improvement over FP32

## Inference with Bfloat16

| Model | Value |
|---|---|
| ResNet50-v1.5 | 1.87 |
| BERT-Large (SQuAD) | 1.9 |

0    0.5    1    1.5    2

■ Perf. improvement over FP32

# How to Install Intel ® optimization for Tensorflow

- Intel®optimization for Tensorflow is included in AI kit. If you have AI Analytics toolkit Intel-Tensorflow conda environment can be activated.

- Install via Pip: pip install intel-tensorflow==2.11.0

- For Stock-tensorflow: pip install tensorflow==2.11.0
  **Since TensorFlow v2.9, the oneAPI Deep Neural Network Library (oneDNN) optimizations are enabled by default.**

- With Conda: conda install tensorflow -c intel

# Handson Workshop

# Github Repo

https://tinyurl.com/oneapi-ai-workshop

# IDC Access Architecture

- SLURM based
- Access valid for 20 days
- Expires if unused for last 7 days
- 20 GB NFS storage



## Picture it this way

nodes in queues
pvc-shared
and
pvc
are identically configured
same CPUs, same four PC cards
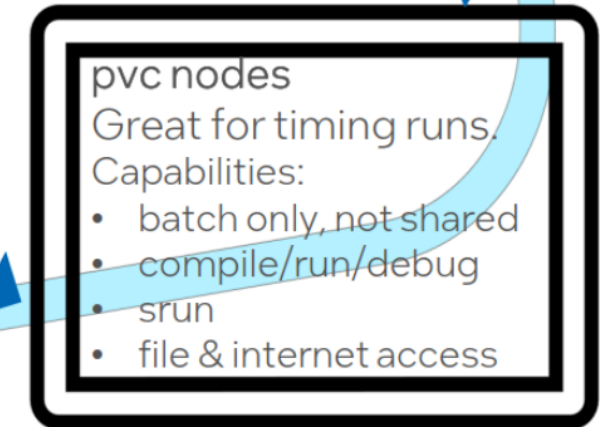(single tile PVCs – but four of them!)

Intel Developer Cloud (IDC)

ssh

**Head Node**
No development tools here.
Capabilities:
- srun
- file & internet access

srun

**pvc-shared nodes**
Great for development.
Capabilities:
- interactive & shared
- compile/run/debug
- srun
- file & internet access

srun

**pvc nodes**
Great for timing runs.
Capabilities:
- batch only, not shared
- compile/run/debug
- srun
- file & internet access

srun

# Intel Data Center GPU Max Series Products & Form Factor

| | Max 1550 GPU (600W OAM) | Max 1350 GPU (450W OAM) | Max 1100 GPU (300W PCIe) |
|---|---|---|---|
| Architecture | | $X^e$ HPC | |
| $X^e$ Cores | 128 | 112 | 56 |
| Memory | HBM2E 128 GB | HBM2E 96 GB | HBM2E 48 GB |
| Cache | L1 64 MB<br>L2 408 MB | L1 48 MB<br>L2 216 MB | L1 28 MB<br>L2 108 MB |
| Max TDP | 600W | 450W | 300W |
| Form Factor | | OAM | PCIe AIC |
| Host Interconnect | | PCIe Gen5 | |
| Physical Ports | | $X^e$ Link 53 GB/s<br>16 ports | $X^e$ Link 53 GB/s<br>6 ports |

# One Generation → 30x AI Performance Gain



**TensorFlow**

**one DNN** / oneAPI

**Intel® Neural Compressor**

**4th Gen Intel® Xeon® Scalable Processor**

| Baseline | 1.5x | 3.9x | 4.8x |
|---|---|---|---|
| 13.06 images/s | 20.54 images/s | 81.66 images/s | 394 images/s |
| (FP32) Official TensorFlow on 3rd Gen Intel® Xeon® Scalable Processor | (FP32) TensorFlow with oneDNN enabled | (INT8) Model quantization with Intel® Neural compressor | (INT8) Intel® AMX optimization on Sapphire Rapids |

**3rd Gen Intel® Xeon™ Scalable Processors**

**4th Gen Intel® Xeon® Scalable processors**

# Slurm Commands

- sinfo – Lists available partitions and node allocations
- squeue – lists the queued jobs
- srun – Sends a job to the queue for execution
- scancel – Deletes a queued job

# Intel® Extension for TensorFlow*

- Intel® Extension for TensorFlow* is a heterogeneous, high performance deep learning extension plugin based on TensorFlow PluggableDevice interface to bring Intel XPU(GPU, CPU, etc) devices into TensorFlow .

- Good performance using default ITEX setting with no code change

- More performance optimizations with minor code change using simple frontend Python API

- GitHub: https://github.com/intel/intel-extension-for-tensorflow



Intel® Extension for TensorFlow* PyPI packages and dependencies

# Intel® Extension for TensorFlow* - Features

**Features:**

Auto Mixed Precision (AMP)

support of AMP with BFloat16 and Float16 operations

Channels Last

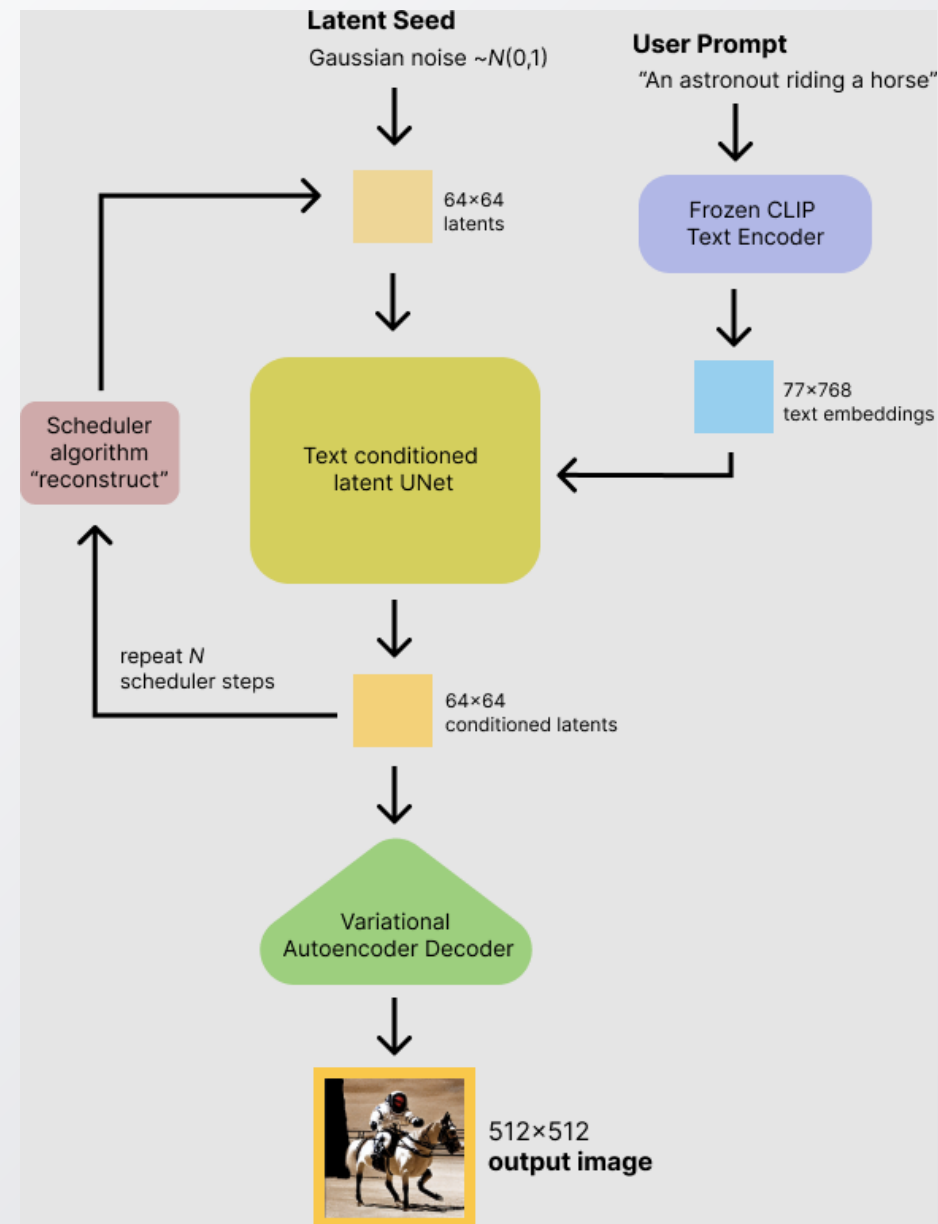support of channels_last (NHWC) memory format

DPC++ Extension

mechanism to create operators with custom DPC++ kernels running on the XPU device

Optimized Fusion

support of SGD/AdamW fusion for both FP32 and BF16 precision
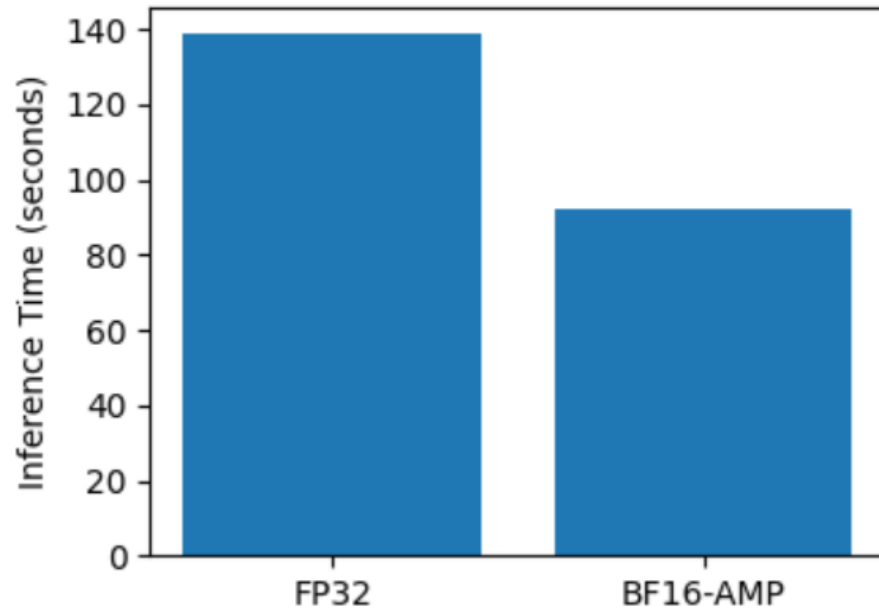
a set of fusion patterns for inference

# Stable Diffusion

- Latent(space) Diffusion Models
  - like DALL-E, Midjourney etc.

- Various tasks text2image, inpainiting, image2image etc.

- 3 main components:
  - Text Encoder – CLIPText
  - Diffusion Model – Unet
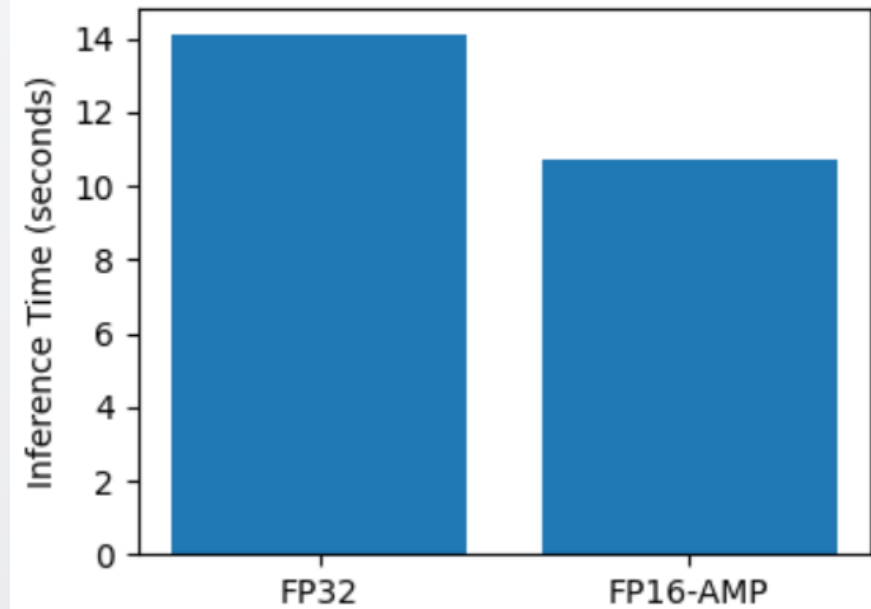  - Image Decoder – VAE

- ~ 1B parameters



https://huggingface.co/blog/stable_diffusion

# Performance Estimates SPR vs PVC

# Thank you!!

oneAPI DevSummit Southeast

FAQ!!

oneAPI

oneAPI DevSummit Southeast