# AI on Intel Architecture

Dr. Séverine Habert, AI Engineering Manager
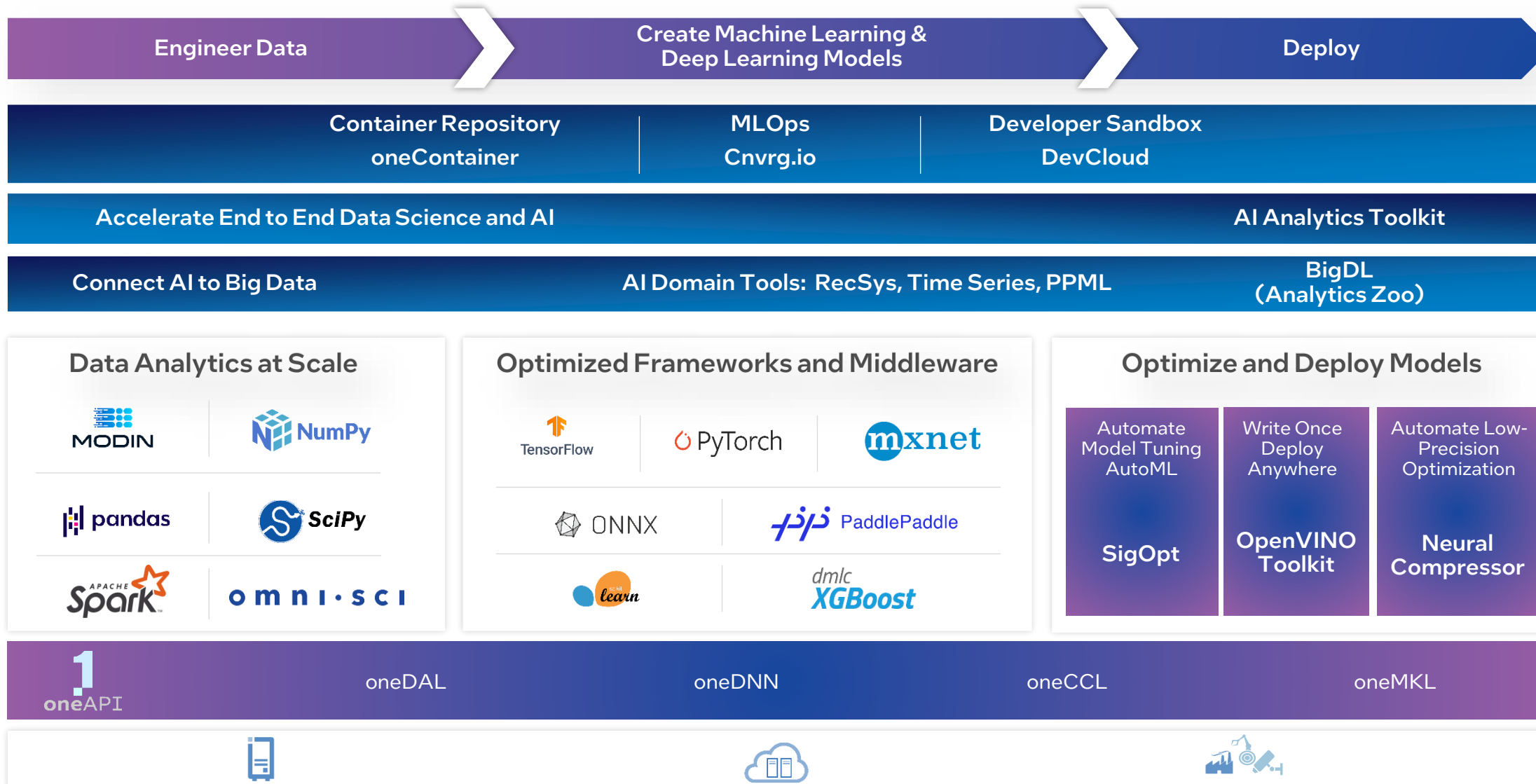
intel.

# Agenda

- oneAPI AI Analytics Toolkit
  - Intel Distribution for Python
  - Classic Machine Learning Libraries
  - Deep Learning Frameworks
  - Intel Neural Compressor
- BigDL
- OpenVINO
- SigOpt

# AI Software Ecosystem and Intel Tools

| Engineer Data | Create Machine Learning & Deep Learning Models | Deploy |
|---|---|---|

| Container Repository oneContainer | MLOps Cnvrg.io | Developer Sandbox DevCloud |
|---|---|---|

| Accelerate End to End Data Science and AI | AI Analytics Toolkit |
|---|---|

| Connect AI to Big Data | AI Domain Tools: RecSys, Time Series, PPML | BigDL (Analytics Zoo) |
|---|---|---|

## Data Analytics at Scale

MODIN  NumPy

pandas  SciPy

Apache Spark  omni·sci

## Optimized Frameworks and Middleware

TensorFlow  PyTorch  mxnet

ONNX  PaddlePaddle

learn  dmlc XGBoost

## Optimize and Deploy Models

| Automate Model Tuning AutoML | Write Once Deploy Anywhere | Automate Low-Precision Optimization |
|---|---|---|
| SigOpt | OpenVINO Toolkit | Neural Compressor |

**1 oneAPI**  oneDAL  oneDNN  oneCCL  oneMKL

Intel Corporation

intel

# oneAPI AI Analytics toolkit

intel.

# Intel's oneAPI Ecosystem

## Built on Intel's Rich Heritage of CPU Tools Expanded to XPUs
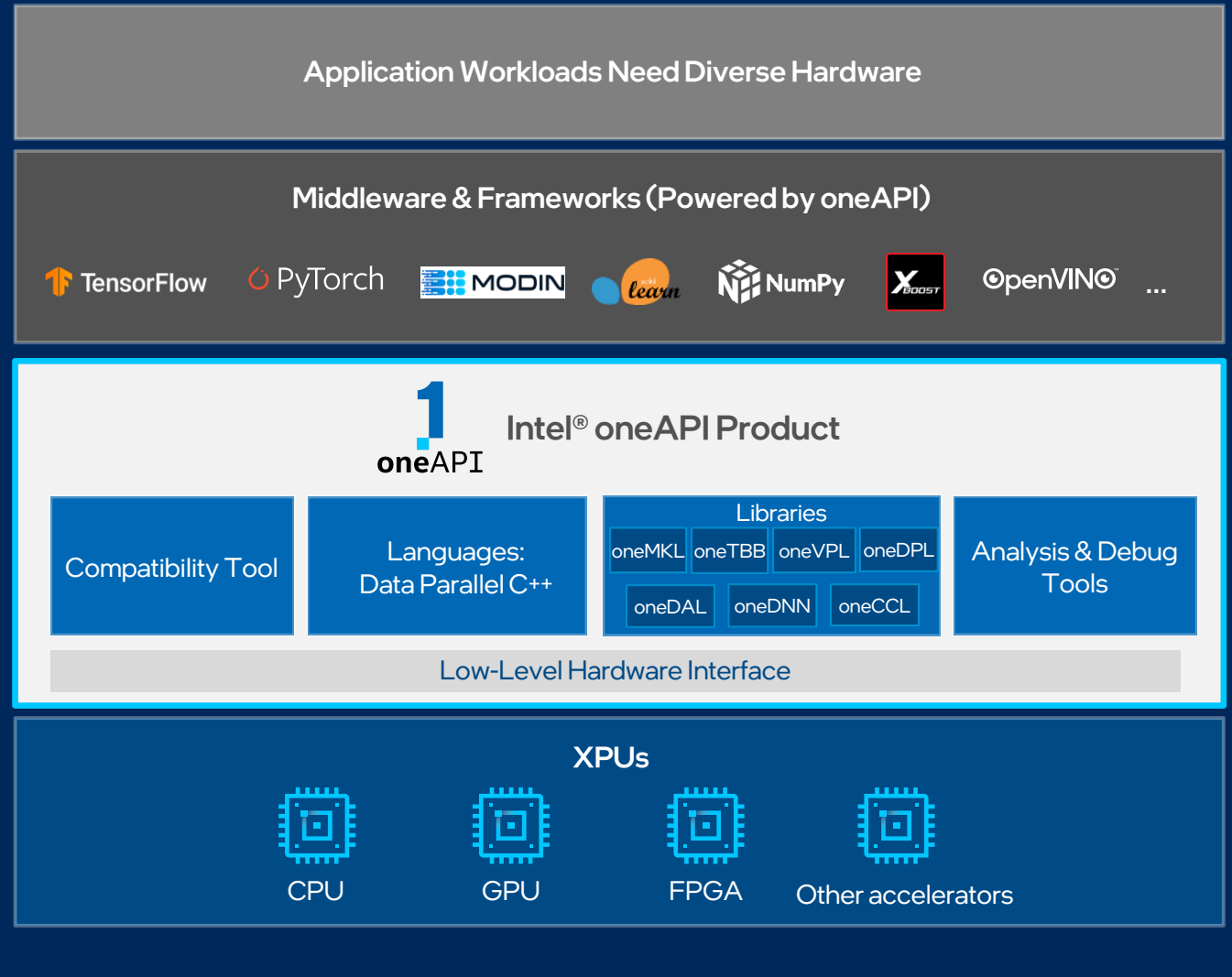
oneAPI

A cross-architecture language based on C++ and SYCL standards

Powerful libraries designed for acceleration of domain-specific functions

A complete set of advanced compilers, libraries, and porting, analysis and debugger tools
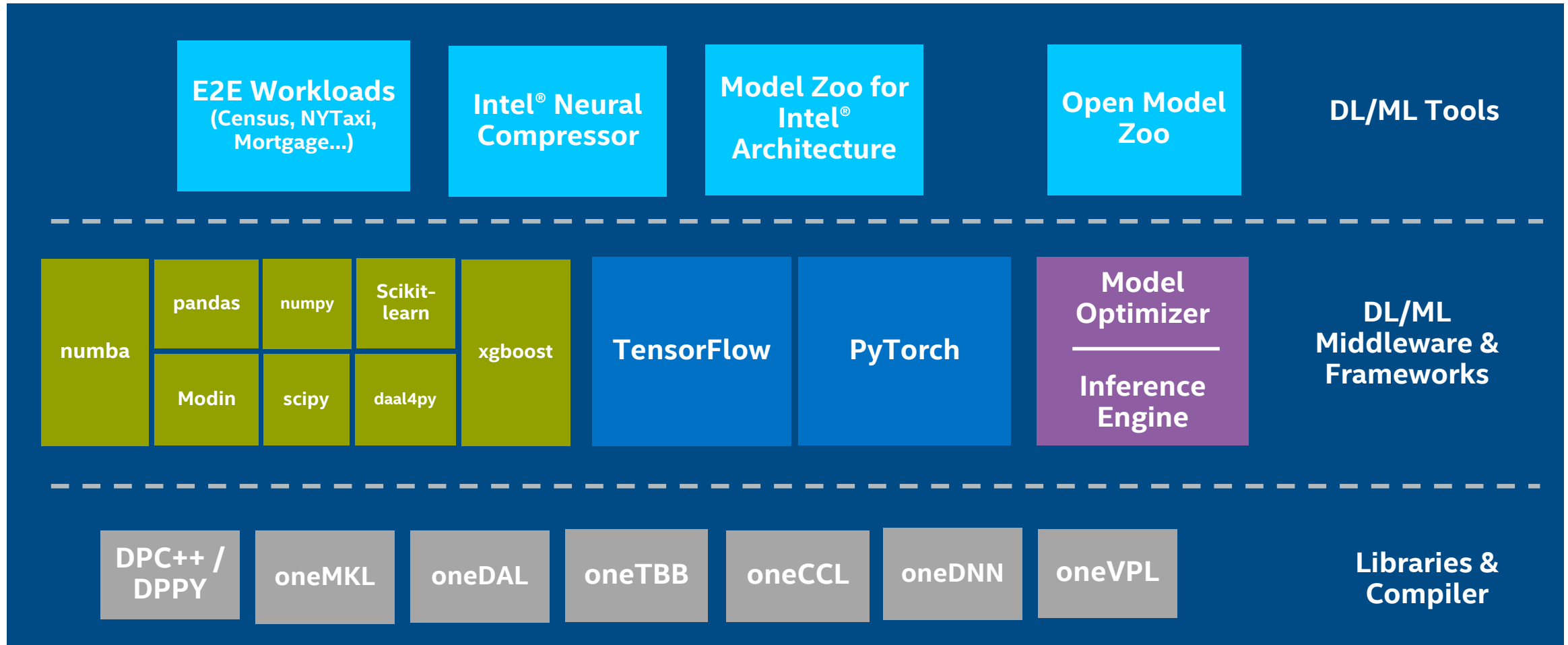
Powered by oneAPI

Frameworks and middleware that are built using one or more of the oneAPI industry specification elements, the DPC++ language, and libraries listed on oneapi.com.

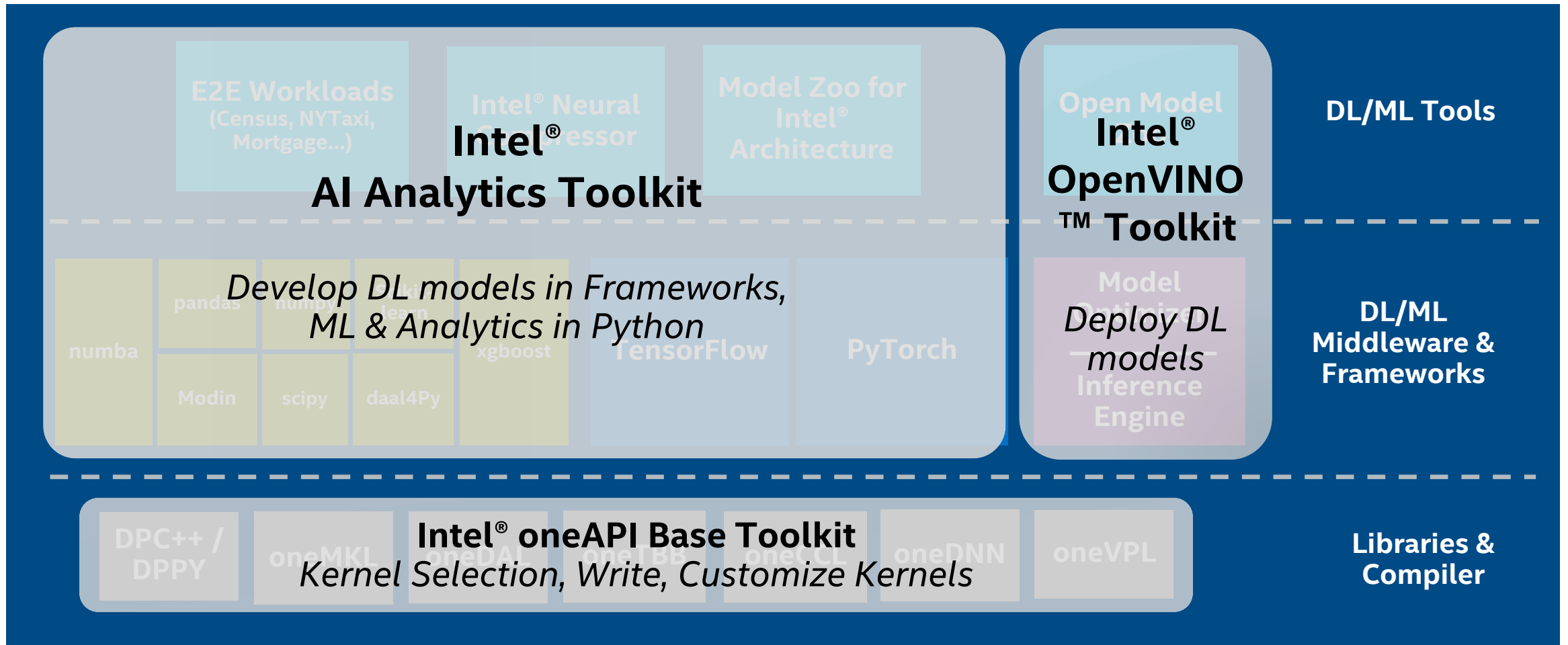**Application Workloads Need Diverse Hardware**

**Middleware & Frameworks (Powered by oneAPI)**

TensorFlow    PyTorch    MODIN    learn    NumPy    XBOOST    OpenVINO    ...

**Intel® oneAPI Product**

| Compatibility Tool | Languages: Data Parallel C++ | Libraries | Analysis & Debug Tools |
|---|---|---|---|
| | | oneMKL  oneTBB  oneVPL  oneDPL | |
| | | oneDAL  oneDNN  oneCCL | |

Low-Level Hardware Interface

**XPUs**

CPU    GPU    FPGA    Other accelerators

[Available Now](#)

Visit software.intel.com/oneapi for more details
Some capabilities may differ per architecture and custom-tuning will still be required. Other accelerators to be supported in the future.

intel    5

# AI Software Stack for Intel® XPUs

**Intel offers a robust software stack to maximize performance of diverse workloads**

| E2E Workloads (Census, NYTaxi, Mortgage...) | Intel® Neural Compressor | Model Zoo for Intel® Architecture | Open Model Zoo | **DL/ML Tools** |
|---|---|---|---|---|

| numba | pandas | numpy | Scikit-learn | xgboost | TensorFlow | PyTorch | Model Optimizer ——— Inference Engine | **DL/ML Middleware & Frameworks** |
| | Modin | scipy | daal4py | | | | | |

| DPC++ / DPPY | oneMKL | oneDAL | oneTBB | oneCCL | oneDNN | oneVPL | **Libraries & Compiler** |
|---|---|---|---|---|---|---|---|

intel
6

# AI Software Stack for Intel® XPUs

**Intel offers a robust software stack to maximize performance of diverse workloads**

E2E Workloads (Census, NYTaxi, Mortgage...)

Intel® Neural Compressor

Model Zoo for Intel® Architecture

**Intel® AI Analytics Toolkit**

Open Model

**Intel® OpenVINO™ Toolkit**

**DL/ML Tools**

numba

pandas

numpy

scikit-learn

xgboost

Modin

scipy

daal4Py

*Develop DL models in Frameworks, ML & Analytics in Python*

TensorFlow

PyTorch

Model Optimizer

*Deploy DL models*

Inference Engine

**DL/ML Middleware & Frameworks**

DPC++ / DPPY

oneMKL

oneDAL

oneTBB

oneCCL

oneDNN

oneVPL

**Intel® oneAPI Base Toolkit**
*Kernel Selection, Write, Customize Kernels*

**Libraries & Compiler**

**Full Set of AI ML and DL Software Solutions Delivered with Intel's oneAPI Ecosystem**

intel

# Intel® AI Analytics Toolkit

## Powered by oneAPI

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

## Who Uses It?

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

| Deep Learning | Data Analytics & Machine Learning | |
|---|---|---|
| Intel® Optimization for TensorFlow | **Accelerated Data Frames** | |
| | Intel® Distribution of Modin | OmniSci Backend |
| Intel® Optimization for PyTorch | **Intel® Distribution for Python** | |
| Intel® Neural Compressor | XGBoost · Scikit-learn · Daal-4Py | |
| Model Zoo for Intel® Architecture | NumPy · SciPy · Pandas | |

**Samples and End2End Workloads**

CPU    GPU

Supported Hardware Architechures[1]

Hardware support varies by individual tool. Architecture support will be expanded over time.
Other names and brands may be claimed as the property of others.

Get the Toolkit HERE or via these locations

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

# Classical Machine Learning

# Intel® Distribution for Python

# Intel® Distribution for Python

- Intel® Distribution for Python covers major usages in HPC and Data Science

- Achieve faster Python application performance — right out of the box — with minimal or no changes to a code

- Accelerate NumPy*, SciPy*, and scikit-learn* with integrated Intel® Performance Libraries such as Intel® oneMKL (Math Kernel Library) and Intel® oneDAL (Data Analytics Library)

- By default, already integrated in Anaconda

# Choose Your Download Option

| Python Solutions | Download Options |
|---|---|
| Tools and frameworks to accelerate end-to-end data science and analytics pipelines | Intel® AI Analytics Toolkit |
| Develop fast, performant Python code with essential computational packages | Intel® Distribution for Python |
| Optimized Python packages from package managers and containers | Conda \| YUM \| APT \| Docker |
| Develop in the Cloud | Intel® DevCloud |

Linux*  Windows*

OS X*

* Also available in the Intel® oneAPI Base Toolkit

# Intel Extension for Scikit-learn

# THE MOST POPULAR ML PACKAGE FOR PYTHON*

# Intel(R) Extension for Scikit-learn

## Common Scikit-learn

```python
from sklearn.svm import SVC

X, Y = get_dataset()

clf = SVC().fit(X, y)
res = clf.predict(X)
```

Scikit-learn mainline

## Scikit-learn with Intel CPU opts

```python
import daal4py as d4p
d4p.patch_sklearn()

from sklearn.svm import SVC

X, Y = get_dataset()



clf = SVC().fit(X, y)
res = clf.predict(X)
```

**Available** through Intel conda
(conda install daal4py –c intel)
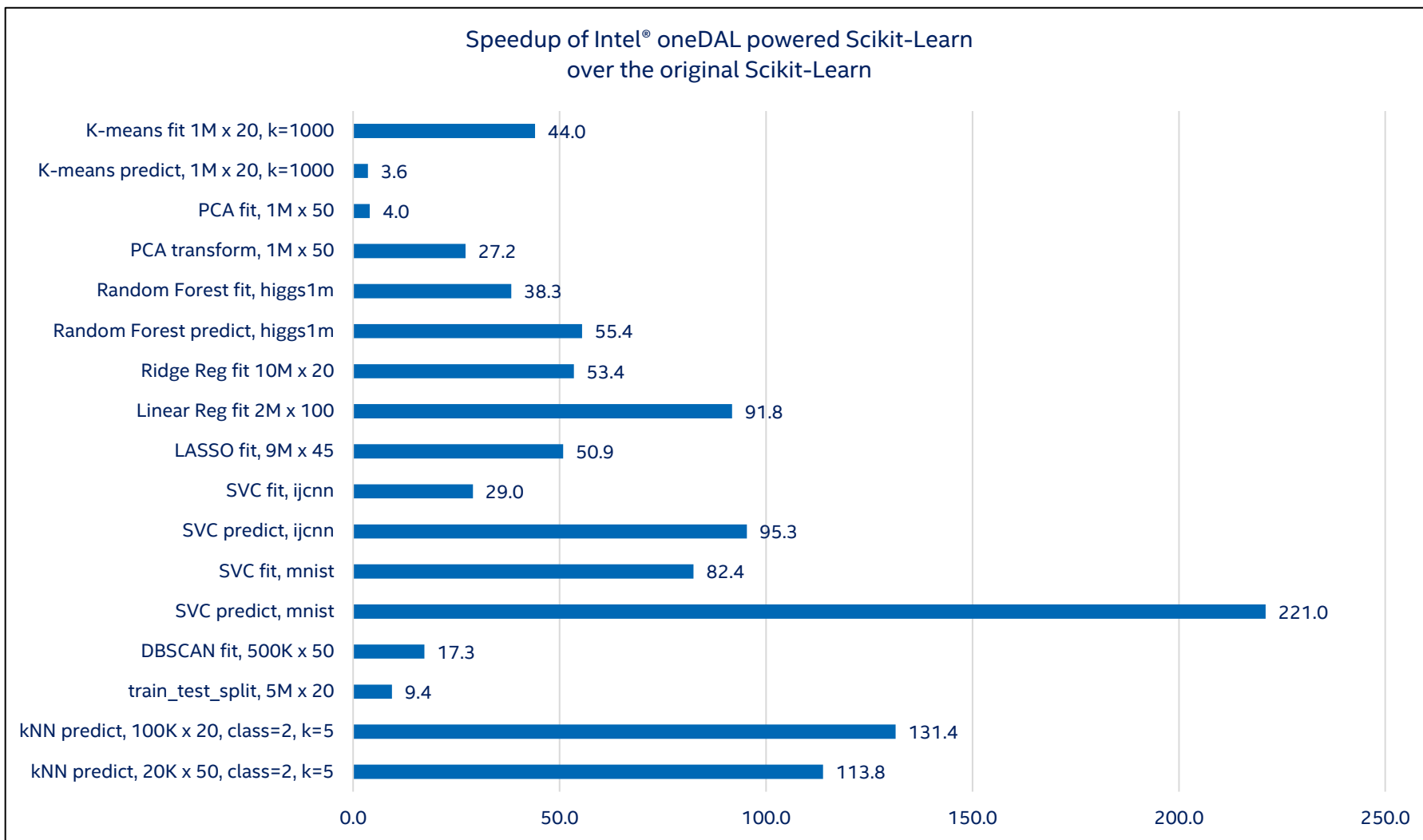
```
> python -m daal4py <your-scikit-learn-script>
```

## Same Code, Same Behavior

**PASSED**

- Scikit-learn, <u>not</u> scikit-learn-*like*
- Scikit-learn conformance (mathematical equivalence) defined by Scikit-learn Consortium, continuously vetted by public CI

Monkey-patch any scikit-learn*
on the command-line

intel.

# Intel optimized Scikit-Learn

**Speedup of Intel® oneDAL powered Scikit-Learn over the original Scikit-Learn**

| Benchmark | Speedup |
|-----------|---------|
| K-means fit 1M x 20, k=1000 | 44.0 |
| K-means predict, 1M x 20, k=1000 | 3.6 |
| PCA fit, 1M x 50 | 4.0 |
| PCA transform, 1M x 50 | 27.2 |
| Random Forest fit, higgs1m | 38.3 |
| Random Forest predict, higgs1m | 55.4 |
| Ridge Reg fit 10M x 20 | 53.4 |
| Linear Reg fit 2M x 100 | 91.8 |
| LASSO fit, 9M x 45 | 50.9 |
| SVC fit, ijcnn | 29.0 |
| SVC predict, ijcnn | 95.3 |
| SVC fit, mnist | 82.4 |
| SVC predict, mnist | 221.0 |
| DBSCAN fit, 500K x 50 | 17.3 |
| train_test_split, 5M x 20 | 9.4 |
| kNN predict, 100K x 20, class=2, k=5 | 131.4 |
| kNN predict, 20K x 50, class=2, k=5 | 113.8 |

## Same Code, Same Behavior

**✔ PASSED**

- Scikit-learn, <u>not</u> scikit-learn-*like*
- Scikit-learn conformance (mathematical equivalence) defined by Scikit-learn Consortium, continuously vetted by public CI

HW: Intel Xeon Platinum 8276L CPU @ 2.20GHz, 2 sockets, 28 cores per socket;
Details: https://medium.com/intel-analytics-software/accelerate-your-scikit-learn-applications-a06cacf44912
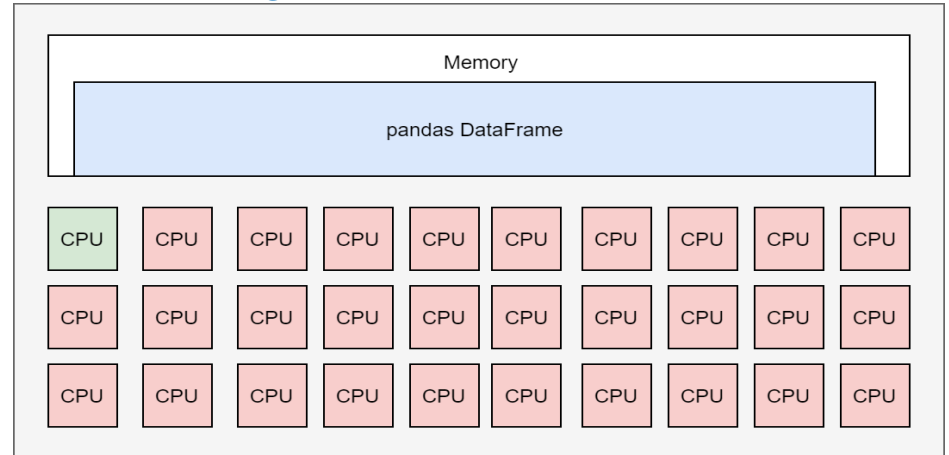
# Modin

# Intel distribution of Modin

- Pandas is a Python package for data manipulation and analysis that offers data structures and operations for manipulating numerical tables and time series

- **Modin =** Pandas + Scalability
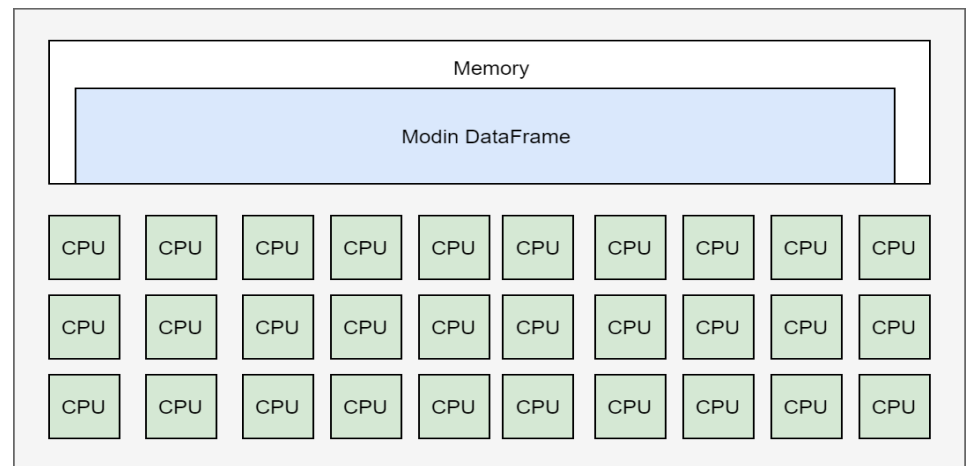
- As simple as **import modin.pandas as pd**

```python
import pandas as pd
```

- In opposition to Pandas, Modin will use all available cores on CPU

- No need to know how many cores your system has, and no need to specify how to distribute the data

- You can get speed-up even on a laptop

- As of 0.9 version, Modin supports 100% of Pandas API

Pandas* on Big Machine

| Memory |
| --- |
| pandas DataFrame |

| CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU |
| CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU |
| CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU |

Modin on Big Machine

| Memory |
| --- |
| Modin DataFrame |

| CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU |
| CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU |
| CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU | CPU |

intel.

# Modin

```python
import modin.pandas as pd
import numpy as np


def run_etl():

    def cat_converter(x):
        if x is '':
            return np.int32(0)
        else:
            return np.int32(int(x, 16))

    names = [f"column_{i}" for i in range(40)]
    converter= {names[i]: cat_converter for i in range(14, 40)}

    df = pd.read_csv('data.csv', delimiter='\t', names=names,
                    converters=converter)

    count_y = df.groupby("column_0")["0"].count()

    return df, count_y


df, count_y = run_etl()
```
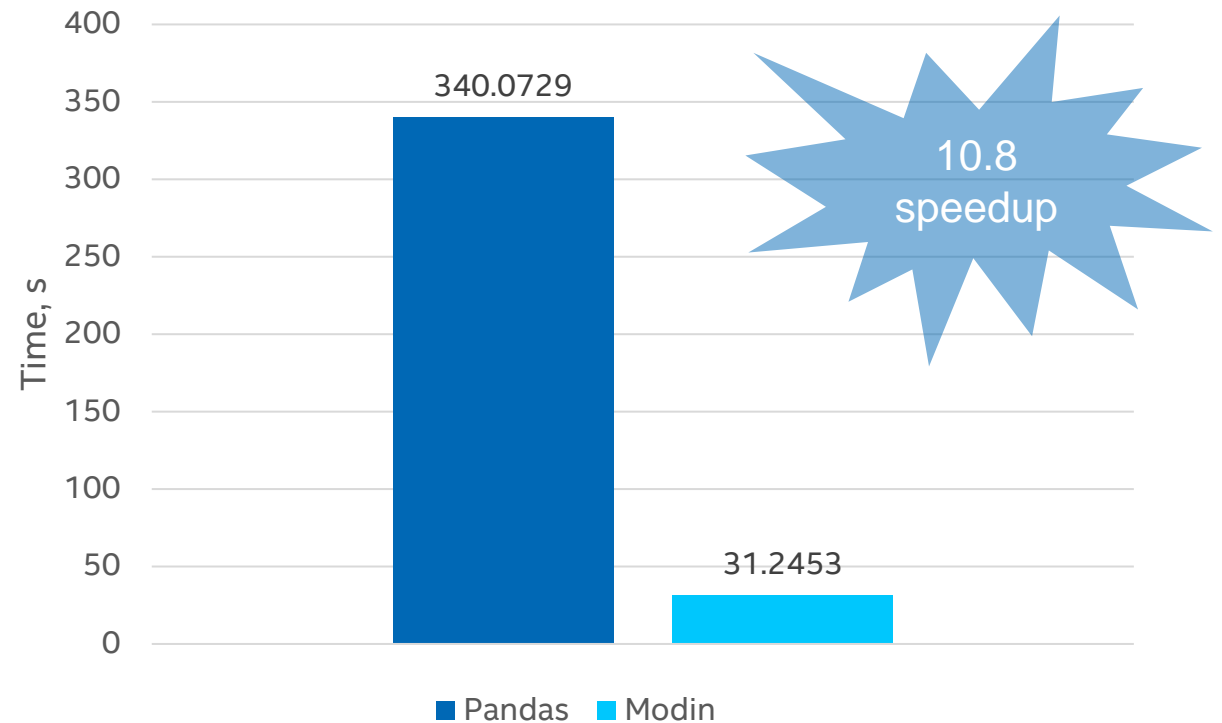
## Execution time Pandas vs. Modin[ray]



10.8 speedup

340.0729 (Pandas)

31.2453 (Modin)

Time, s

■ Pandas   ■ Modin

Intel® Xeon™ Gold 6248 CPU @ 2.50GHz, 2x20 cores

▪ Dataset size: 2.4GB

# XGBoost

# Gradient Boosting – Overview

Gradient Boosting:

- Boosting algorithm (Decision Trees - base learners)
- Solve many types of ML problems (classification, regression, learning to rank)
- Highly-accurate, widely used by Data Scientists
- Compute intensive workload
- Known implementations: XGBoost*, LightGBM*, CatBoost*, Intel® oneDAL, …

# Gradient Boosting Acceleration – gain sources

**Pseudocode for XGBoost* (0.81) implementation**

```
def ComputeHist(node):
  hist = []
  for i in samples:
    for f in features:
      bin = bin_matrix[i][f]
      hist[bin].g += g[i]
      hist[bin].h += h[i]
  return hist


def BuildLvl:
  for node in nodes:
    ComputeHist(node)

  for node in nodes:
    for f in features:
      FindBestSplit(node, f)

  for node in nodes:
    SamplePartition(node)
```

**Pseudocode for Intel® oneDAL implementation**

```
def ComputeHist(node):
  hist = []
  for i in samples:
    prefetch(bin_matrix[i + 10])
    for f in features:
      bin = bin_matrix[i][f]
      bin_value = load(hist[2*bin])
      bin_value = add(bin_value, gh[i])
      store(hist[2*bin], bin_value)
  return hist


def BuildLvl:
  parallel_for node in nodes:
    ComputeHist(node)

  parallel_for node in nodes:
    for f in features:
      FindBestSplit(node, f)

  parallel_for node in nodes:
    SamplePartition(node)
```

Memory prefetching to mitigate irregular memory access

Usage uint8 instead of uint32

SIMD instructions instead of scalar code

Nested parallelism

Advanced parallelism, reducing seq loops

Usage of AVX-512, vcompress instruction (from Skylake)
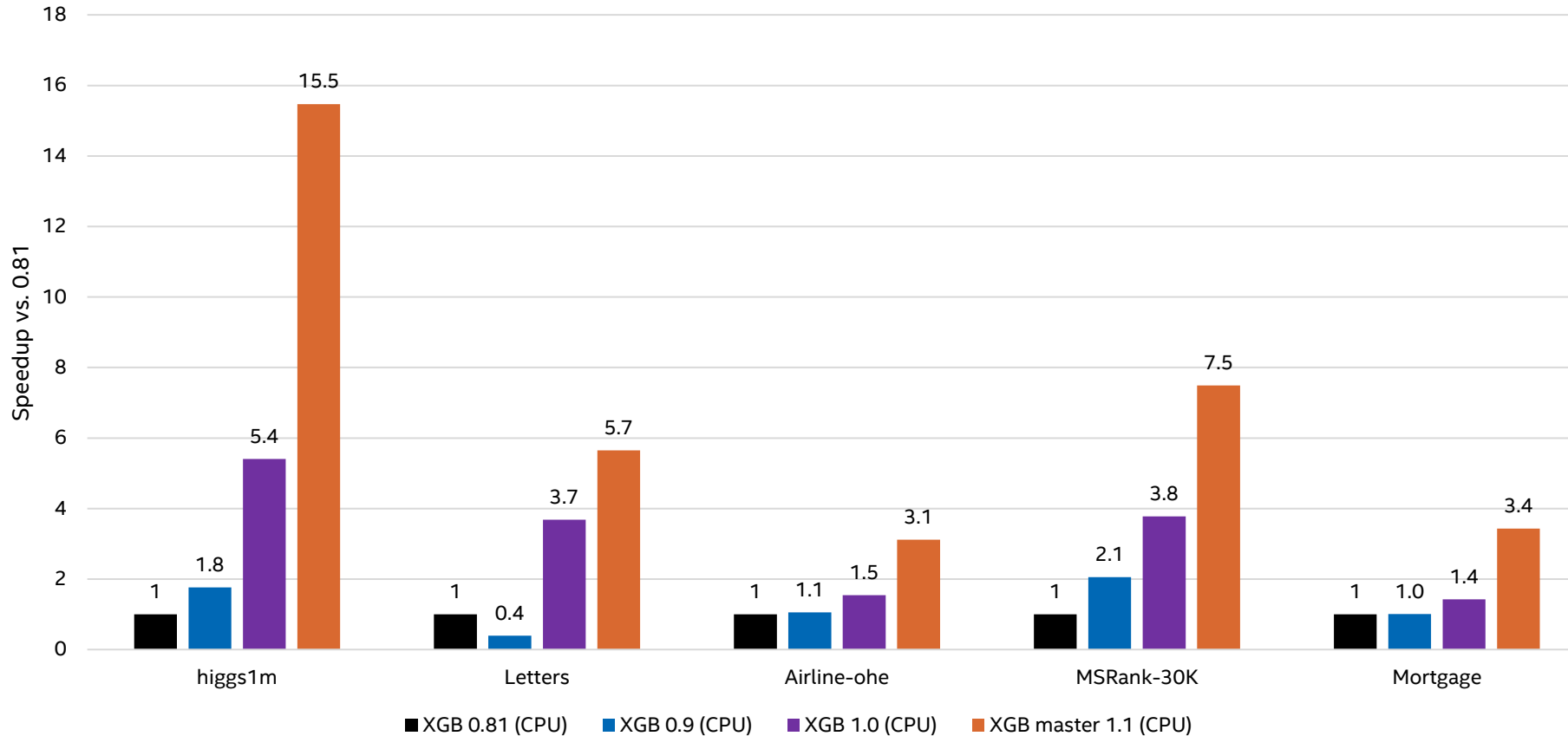
Training stage

Legend:
Moved from Intel® oneDAL to XGBoost (v1.3)

Already available in Intel® DAAL, potential optimizations for XGBoost*

# XGBoost* fit CPU acceleration ("hist" method)

## XGBoost fit - acceleration against baseline (v0.81) on Intel CPU



Speedup vs. 0.81

| | higgs1m | Letters | Airline-ohe | MSRank-30K | Mortgage |
|---|---|---|---|---|---|
| XGB 0.81 (CPU) | 1 | 1 | 1 | 1 | 1 |
| XGB 0.9 (CPU) | 1.8 | 0.4 | 1.1 | 2.1 | 1.0 |
| XGB 1.0 (CPU) | 5.4 | 3.7 | 1.5 | 3.8 | 1.4 |
| XGB master 1.1 (CPU) | 15.5 | 5.7 | 3.1 | 7.5 | 3.4 |

**+ Reducing memory consumption**

| memory, Kb | Airline | Higgs1m |
|---|---|---|
| Before | 28311860 | 1907812 |
| #5334 | 16218404 | 1155156 |
| reduced: | 1.75 | 1.65 |

**CPU configuration**: c5.24xlarge AWS Instance, CLX 8275 @ 3.0GHz, 2 sockets, 24 cores per socket, HT:on, DRAM (12 slots / 32GB / 2933 MHz)
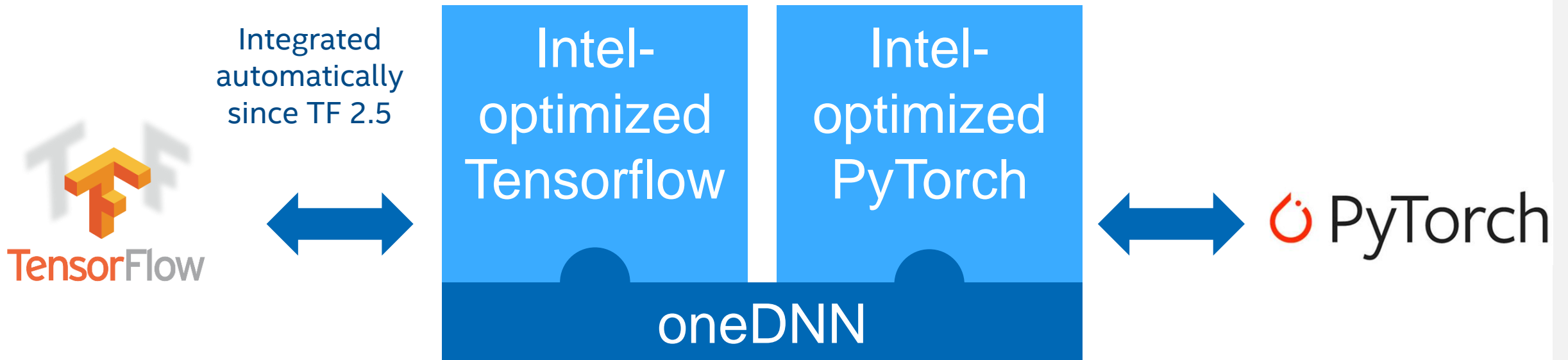
# Deep Learning

intel ®
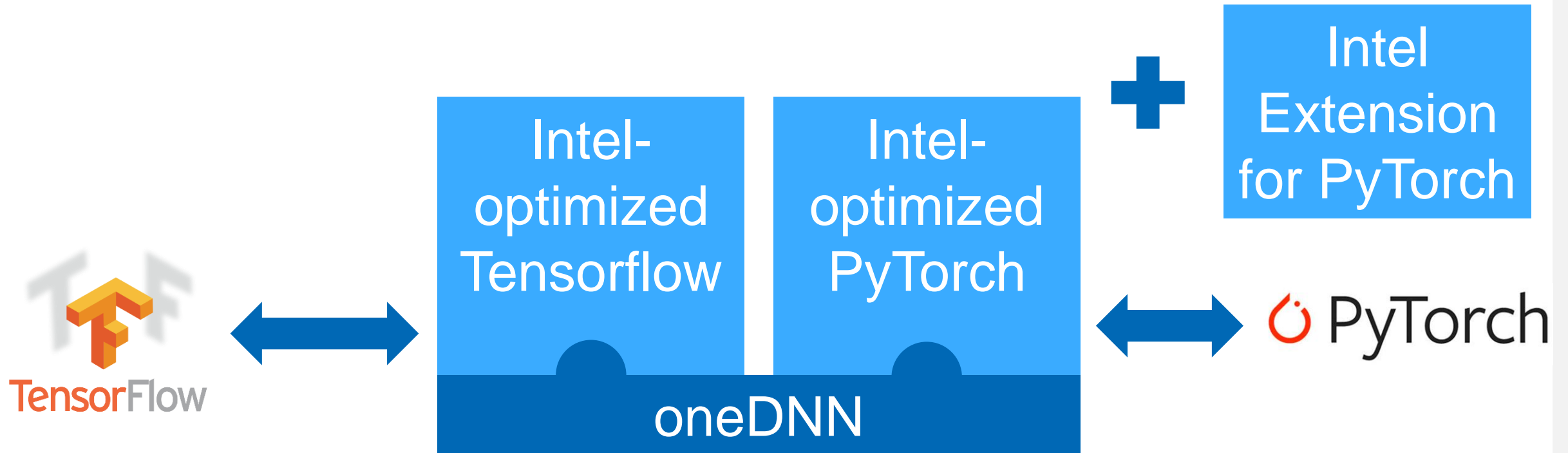
# Intel-optimized Deep Learning frameworks

# Intel-optimized Deep Learning Frameworks

- Intel-optimized DL frameworks are drop-in replacement,

  - No front code change for the user

- Optimizations are upstreamed automatically (TF) or on a regular basis (PyTorch) to stock frameworks

  - TF: Optimizations are integrated automatically since TF 2.5 and are activated after setting up TF_ENABLE_ONEDNN_OPTS=1

**TensorFlow**

Integrated automatically since TF 2.5

⟷

| Intel-optimized Tensorflow | Intel-optimized PyTorch |
|---|---|

**oneDNN**

⟷

**PyTorch**

# Intel-optimized Deep Learning Frameworks

- Intel Extension for PyTorch is an additional module for functions not supported in standard PyTorch (such as mixed precision and dGPU support)

- As they offer more aggressive optimizations, they offer bigger speed-up for training and inference

# Intel® oneAPI Deep Neural Network Library (oneDNN)

## Basic Information

- ▪ Features
- • Training: float32, bfloat16[1]
- • Inference: float32, bfloat16[1], float16[1], and int8[1]
- • Runs on Intel CPU and GPU

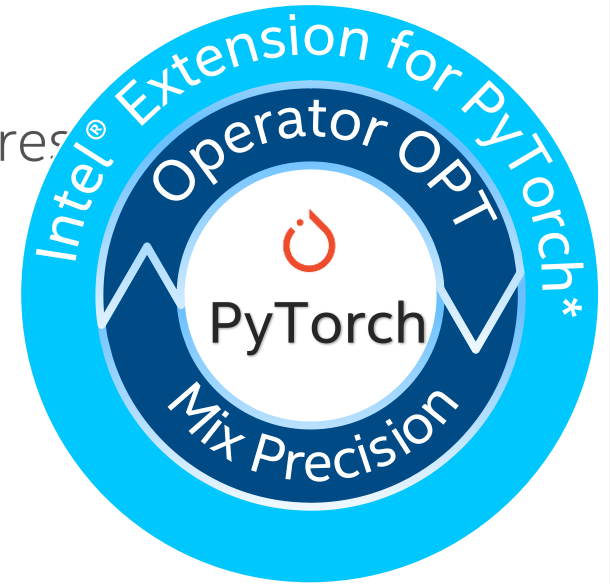| | Intel® oneDNN |
|---|---|
| Convolution | 2D/3D Direct Convolution/Deconvolution, Depthwise separable convolution 2D Winograd convolution |
| Inner Product | 2D/3D Inner Production |
| Pooling | 2D/3D Maximum 2D/3D Average (include/exclude padding) |
| Normalization | 2D/3D LRN across/within channel, 2D/3D Batch normalization |
| Eltwise (Loss/activation) | ReLU(bounded/soft), ELU, Tanh; Softmax, Logistic, linear; square, sqrt, abs, exp, gelu, swish |
| Data manipulation | Reorder, sum, concat, View |
| RNN cell | RNN cell, LSTM cell, GRU cell |
| Fused primitive | Conv+ReLU+sum, BatchNorm+ReLU |
| Data type | f32, bfloat16, s8, u8 |

[1] Low precision data types are supported only for platforms where hardware acceleration is available

# Optimizations

1. <u>Operator optimizations</u>: Replace default kernels by highly-optimized kernels (using Intel® oneDNN)

2. <u>Memory layout optimizations:</u> set optimal layout for each kernel, while minimizing memory changes in between kernels

3. <u>Graph optimizations</u>: Fusion, Layout Propagation

# Intel® Extension for PyTorch* (IPEX)

- Buffer the PRs for stock Pytorch

- Provide users with the up-to-date Intel software/hardware features

- Streamline the work to integrate oneDNN

- Unify user experiences on Intel CPU and GPU

### Operator Optimization
➢ Customized operators
➢ Auto graph optimization

### Mix Precision
➢ Accelerate PyTorch operator by LP
➢ Simplify the data type conversion

### Optimal Optimizer
➢ Split Optimizer (e.g., split-sgd)
➢ Fused Optimizer

# Ease-of-Use User-Facing API (v1.10.x~)

## For Float32

```python
import torch
import torchvision.models as models

model = models.resnet50(pretrained=True)
model.eval()
data = torch.rand(1, 3, 224, 224)

model = model.to(memory_format=torch.channels_last)
data = data.to(memory_format=torch.channels_last)

#################### code changes ####################
import intel_extension_for_pytorch as ipex
model = ipex.optimize(model)
######################################################

with torch.no_grad():
    model(data)
```

# Ease-of-Use User-Facing API (v1.10.x~)

## For BFloat16

```python
import torch
import torchvision.models as models

model = models.resnet50(pretrained=True)
model.eval()
data = torch.rand(1, 3, 224, 224)

model = model.to(memory_format=torch.channels_last)
data = data.to(memory_format=torch.channels_last)

#################### code changes ####################
import intel_extension_for_pytorch as ipex
model = ipex.optimize(model, dtype=torch.bfloat16)
#####################################################

with torch.no_grad():
    with torch.cpu.amp.autocast():
        model(data)
```

# Usage

```python
import torch
import intel_pytorch_extension

class Model(torch.nn.Module):
  def __init__(self):
    super(Model, self).__init__()
    self.conv2d = torch.nn.Conv2d(3, 5, 5)

  def forward(self, input):
    res = self.conv2d(input)
    return res

input = torch.randn(5, 3, 9, 9)
model = Model()
model = model.to('xpu')
input = input.to('xpu')
res = model(input)
```

```python
import torch
import intel_extension_for_pytorch as ipex

class Model(torch.nn.Module):
  def __init__(self):
    super(Model, self).__init__()
    self.conv2d = torch.nn.Conv2d(3, 5, 5)

  def forward(self, input):
    res = self.conv2d(input)
    return res

input = torch.randn(5, 3, 9, 9)
model = Model()
model = ipex.optimize(model, dtype=torch.float32, level="O1")
input = input.to(memory_format=torch.channels_last)
res = model(input)
```
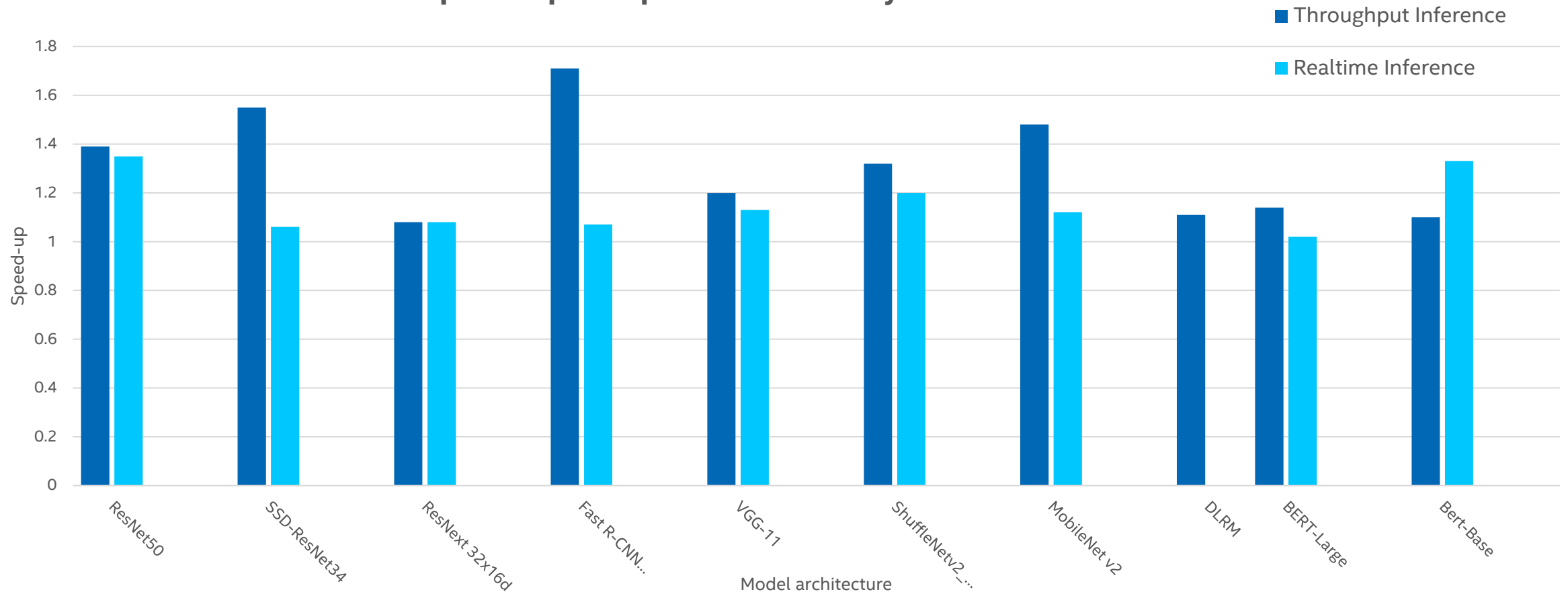
Prior to v1.10                          v1.10

# Intel Extension for PyTorch benchmark

**Speed-up compared to stock PyTorch for Float32**



Legend:
- Throughput Inference (dark blue)
- Realtime Inference (light blue)

Y-axis: Speed-up (0 to 1.8)
X-axis: Model architecture

Models: ResNet50, SSD-ResNet34, ResNext 32x16d, Fast R-CNN.., VGG-11, ShuffleNetv2_.., MobileNet v2, DLRM, BERT-Large, Bert-Base

# Intel Neural Compressor (INC)

# INC: Intel Neural Compressor

- Intel Neural Compressor is an open-source Python library to create low-precision inference solutions on popular deep-learning frameworks

- It supports quantization to BF16/INT8, pruning, knowledge distillation and graph optimizations
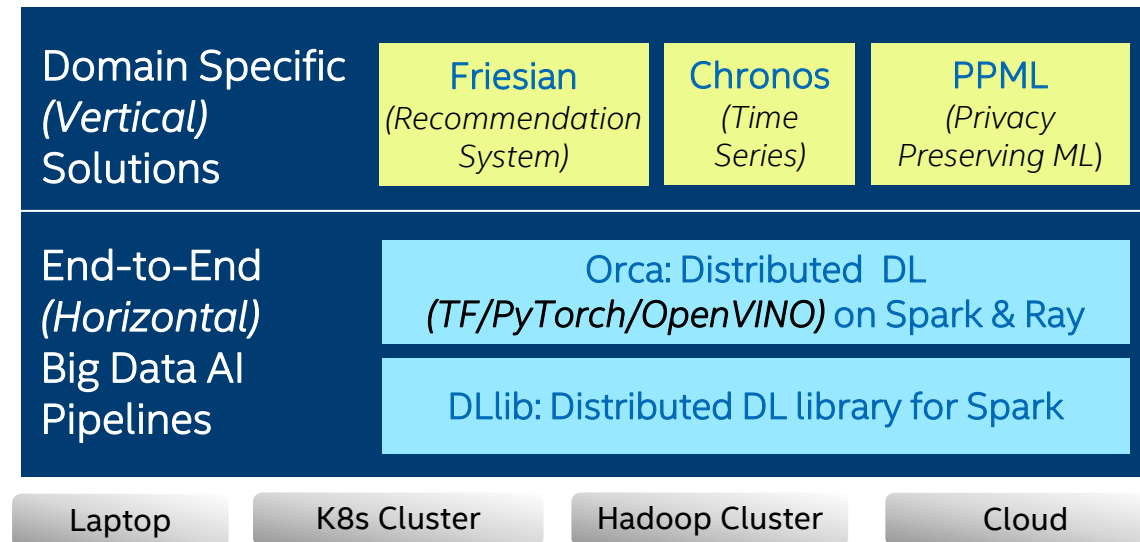
| Model | TensorFlow | PyTorch | ONNX | MXNet |
|---|---|---|---|---|

**Intel® Neural Compressor Architecture**

**User-Facing APIs**

Quantization, Pruning, Knowledge Distillation, Graph Optimization, …

**Compressions**

**Quantization**
- Post training static quantization
- Post training dynamic quantization
- Quantization-aware training

**Pruning**
- magnitude pruning
- Gradient sensitivity pruning

Knowledge Distillation

**Mix Precision**
- FP32 -> INT8/BF16
- FP32 -> BF16
- FP32 -> Opt Fp32

**Auto Tuning**

Tuning Strategies

**Backends**

TensorFlow, PyTorch, ONNX Runtime, MXNet, Engine

**Hardware platforms**

Intel CPU        Intel GPU

intel.

# BigDL

# BigDL for Big Data AI

## Technology Stack

- Bringing AI to Big Data software ecosystem
- Leading with "IA differentiated" domain-specific solutions

| Domain Specific (Vertical) Solutions | Friesian (Recommendation System) | Chronos (Time Series) | PPML (Privacy Preserving ML) |
|---|---|---|---|
| End-to-End (Horizontal) Big Data AI Pipelines | Orca: Distributed DL (TF/PyTorch/OpenVINO) on Spark & Ray | | |
| | DLlib: Distributed DL library for Spark | | |

Laptop    K8s Cluster    Hadoop Cluster    Cloud

## Value of Big Data AI Toolkit

| Value | Example Users |
|---|---|
| Rich *software ecosystem* for Big Data processing on IA | *Mastercard, BBVA, Alibaba Cloud, Inspur, etc.* |
| Better *E2E productivity and performance* for AI pipelines | *Burger King, SK Telecom, JD.com, Midea, etc.* |
| *Domain-specific AI solutions* for Big Data | *Ant Financial, Capgemini, Mavenir, UnionPay, etc.* |

# OpenVINO

Intel Confidential

intel

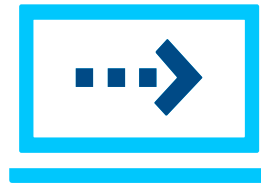# Intel® Distribution of OpenVINO™ Toolkit

- Tool Suite for High-Performance, Deep Learning Inference
- Fast, accurate real-world results using high-performance, AI and computer vision inference deployed into production across Intel® architecture from edge to cloud
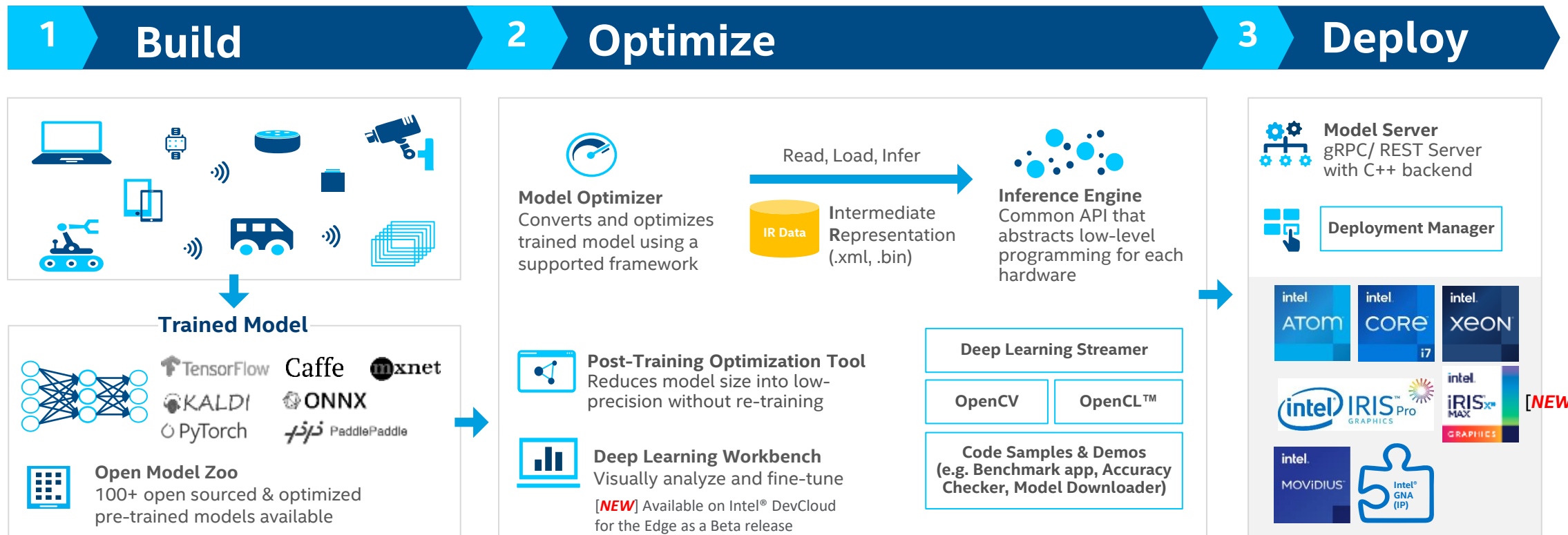
High-Performance,
Deep Learning Inference

Streamlined Development,
Ease of Use

Write Once,
Deploy Anywhere

- Enables deep learning inference from the edge to cloud.
- Supports heterogeneous execution across Intel accelerators, using a common API for the Intel® CPU, Intel® Integrated Graphics, Intel® Gaussian & Neural Accelerator, Intel® Neural Compute Stick 2, Intel® Vision Accelerator Design with Intel® Movidius™ VPUs.

- Speeds time-to-market through an easy-to-use library of CV functions and pre-optimized kernels.
- Includes optimized calls for CV standards, including OpenCV* and OpenCL™.

# Three steps for the Intel® Distribution of OpenVINO™ toolkit

**1  Build**  |  **2  Optimize**  |  **3  Deploy**

## Build



**Trained Model**

TensorFlow · Caffe · mxnet · KALDI · ONNX · PyTorch · PaddlePaddle

**Open Model Zoo**
100+ open sourced & optimized pre-trained models available

## Optimize

**Model Optimizer**
Converts and optimizes trained model using a supported framework

Read, Load, Infer →

**IR Data**

**I**ntermediate **R**epresentation (.xml, .bin)

**Inference Engine**
Common API that abstracts low-level programming for each hardware

**Post-Training Optimization Tool**
Reduces model size into low-precision without re-training

**Deep Learning Workbench**
Visually analyze and fine-tune

[*NEW*] Available on Intel® DevCloud for the Edge as a Beta release

**Deep Learning Streamer**

**OpenCV**  |  **OpenCL™**

**Code Samples & Demos**
(e.g. Benchmark app, Accuracy Checker, Model Downloader)

## Deploy

**Model Server**
gRPC/ REST Server with C++ backend

**Deployment Manager**

intel ATOM · intel CORE i7 · intel XEON

intel IRIS Pro GRAPHICS · intel IRIS Xᵉ MAX GRAPHICS   [*NEW*]

intel MOVIDIUS · Intel® GNA (IP)

# Supported Frameworks

Breadth of supported frameworks to enable developers with flexibility

011010110110
110101101011
001011010100
011010110110
110101101011
001011010100
011010110110
110101101011
001011010100
011010110110
110101101011
001011010100
011010110110
110101101011
001011010100
011010110110
110101101011
001011010100

**Trained Model**

(and other tools via ONNX* conversion)

TensorFlow

Caffe

mxnet

KALDI

ONNX

OpenVINO™

Caffe2    Chainer    *dmlc* **XGBoost**    K    LibSVM

MATLAB    Cognitive Toolkit    mxnet    SIEMENS

Neural Network Libraries    PaddlePaddle    PyTorch    SAS    SIEMENS

*learn*    TensorFlow

**Supported Frameworks and Formats** ▸ https://docs.openvinotoolkit.org/latest/_docs_IE_DG_Introduction.html#SupportedFW
**Configure the Model Optimizer for your Framework** ▸ https://docs.openvinotoolkit.org/latest/_docs_MO_DG_prepare_model_Config_Model_Optimizer.html

# Model Optimization

Breadth of supported frameworks to enable developers with flexibility

**Model Optimizer** loads a model into memory, reads it, builds the internal representation of the model, optimizes it, and produces the **Intermediate Representation**.

Optimization techniques available are:

- Node merging

- Horizontal fusion

- Batch normalization to scale shift

- Fold scale shift with convolution

- Drop unused layers (dropout)

*Note:* Except for ONNX (.onnx model formats), all models have to be converted to an IR format to use as input to the Inference Engine

**Trained Model**

**Model Optimizer**

Read, Load, Infer

**IR Data**

**I**ntermediate **R**epresentation (.xml, .bin)

.**xml** – describes the network topology
.**bin** – describes the weights and biases binary data

# Optimal Model Performance Using the Inference Engine

## Core Inference Engine Libraries

- Create Inference Engine Core object to work with devices
- Read the network
- Manipulate network information
- Execute and pass inputs and outputs

## Device-specific Plugin Libraries

- For each supported target device, Inference Engine provides a plugin — a DLL/shared library that contains complete implementation for inference on this device.
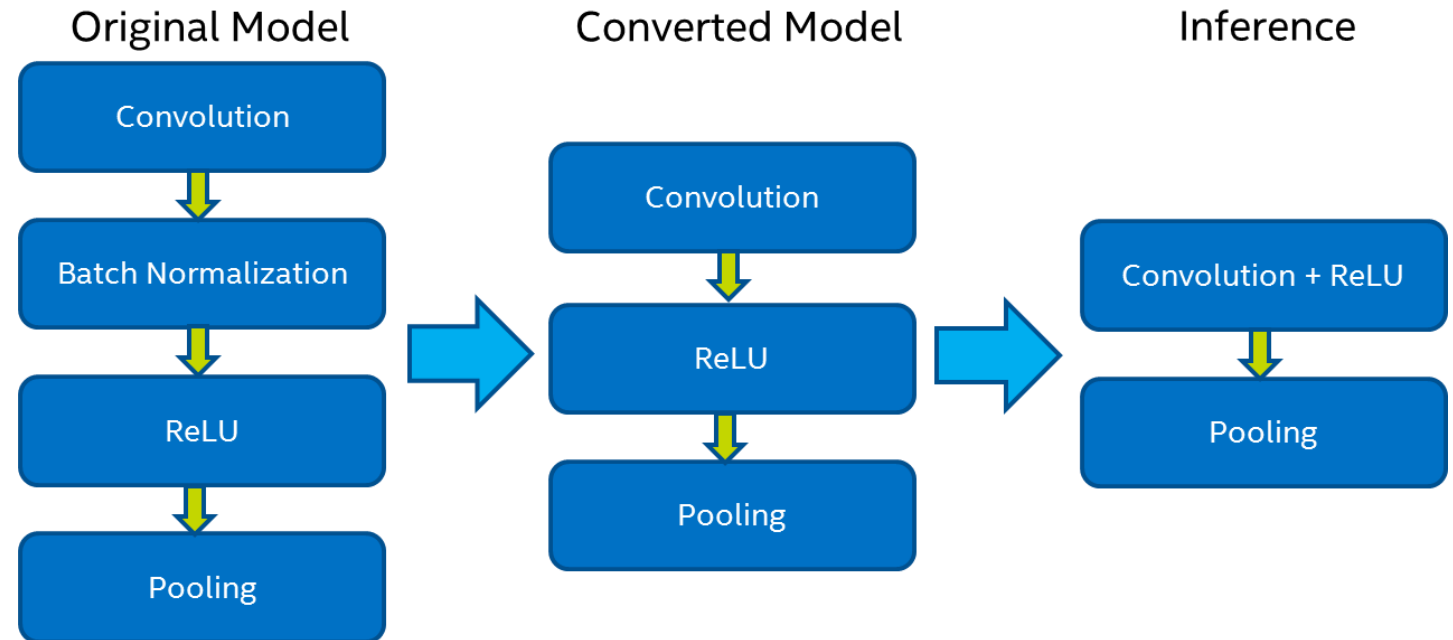
Applications

Inference Engine runtime

Inference Engine (Common API)

Multi-device plugin (optional but recommended - for full system utilization)

| OneDNN plugin | | GNA plugin | Myriad & HDDL plugins |
|---|---|---|---|
| Intrinsics | OpenCL™ | GNA API | Movidius API |

Plugin architecture

GPU = Intel CPU with integrated graphics/Intel® Processor Graphics/GEN

GNA = Gaussian mixture model and Neural Network Accelerator

# Model Optimizer: Linear Operation Fusing

- Example

1. Remove Batch normalization stage.

2. Recalculate the weights to 'include' the operation.

3. Merge Convolution and ReLU into one optimized kernel.

**Original Model**

```
Convolution
    ↓
Batch Normalization
    ↓
ReLU
    ↓
Pooling
```

**Converted Model**

```
Convolution
    ↓
ReLU
    ↓
Pooling
```

**Inference**

```
Convolution + ReLU
    ↓
Pooling
```

# Common Workflow for Using the Inference Engine API

**exec_net** = **ie.load_network**(network=net, device_name=device, num_requests=request_number)

```
Create Inference Engine Core object
```
→
```
Read the Intermediate Representation
```
→
```
Prepare inputs and outputs format
```
→
```
Load Network to device & Create infer request
```

**ie** = IECore()

**net** = **ie.read_network**(model=model_xml, weights=model_bin)

input_blob = next(iter(**net.inputs**))

output_blob = next(iter(**net.outputs**))

res = **exec_net.infer**(inputs={input_blob: in_frame})

```
Process the results
```
←
```
Run Inference
```
←
```
Prepare input frame
```

Inference loop

n, c, h, w = net.inputs[input_blob].**shape**

in_frame = cv2.**resize**(image, (w, h))

in_frame = in_frame.**transpose**((2, 0, 1))

in_frame = in_frame.**reshape**((n, c, h, w))

http://docs.openvinotoolkit.org/latest/_docs_IE_DG_Integrate_with_customer_application_new_API.html

# Pre-Trained Models and Public Models

Open-sourced repository of pre-trained models and support for public models

Use free **Pre-trained Models** to speed up development and deployment

Take advantage of the **Model Downloader** and other automation tools to quickly get started

Iterate with the **Accuracy Checker** to validate the accuracy of your models

## 100+ Pre-trained Models
*Common AI tasks*

Object Detection
Object Recognition
Reidentification
Semantic Segmentation
Instance Segmentation
Human Pose Estimation
Image Processing
Text Detection
Text Recognition
Text Spotting
Action Recognition
Image Retrieval
Compressed Models
Question Answering

## 100+ Public Models
*Pre-optimized external models*

Classification
Segmentation
Object Detection
Human Pose Estimation
Monocular Depth Estimation
Image Inpainting
Style Transfer
Action Recognition
Colorization

# OpenVINO as execution provider

- You can use OpenVINO Inference Engine as backend of other DL Inference Frameworks such as Tensorflow or ONNX Runtime



- Benefit: the advantages of OpenVINO (multiple HW support and acceleration) in your favorite framework

# OpenVINO™ Integration with TensorFlow*



Only 2 lines to be added to regular Tensorflow

```
1   # Installation steps
2   # more details : https://github.com/openvinotoolkit/openvino_tensorflow
3   #pip3 install -U pip==21.0.1
4   #pip3 install -U tensorflow==2.4.1
5   #pip3 install openvino-tensorflow
6
7   # Import package and set backend
8   import openvino_tensorflow
9   openvino_tensorflow.set_backend('GPU'))
10
11  # Load a TF Saved Model
12  model = tf.keras.models.load_model('resnet50_saved_model')
13
14  # Get the input size of the model
15  network_input_size = saved_model_loaded.input.shape()
16
17  # Resize the input image
18  resized_image = resize(input_image, network_input_size)
19
20  # Run inference
21  model.predict(resized_image)
```

CPU
GPU
MYRIAD
VAD-M

# SigOpt

# SigOpt

- SigOpt is the only experimentation platform that brings together:
  - Bayesian-based hyperparameter optimization tuning (including multi-metric optimization)



Objective Metric (AUC, Accuracy)

① Big Data
② Machine Learning Models
④ Better Models
③ SigOpt optimally suggests new hyperparameters

Model training

# SigOpt

- And experiment management



Example code

# Questions?

intel.