**Choose the Best Accelerated Technology**

# Intel AI Analytics Toolkit – Classical ML
## LRZ AI Workshop

Roy Allela – AI Software Solutions Engineer

21.04.2022

intel.

# Agenda

- Intel AI Analytics Toolkit

- Intel Distribution for Python

- Intel Distribution of Modin

- Intel(R) Extension for Scikit-learn

- XGBoost Optimizations

# Intel® AI Analytics Toolkit
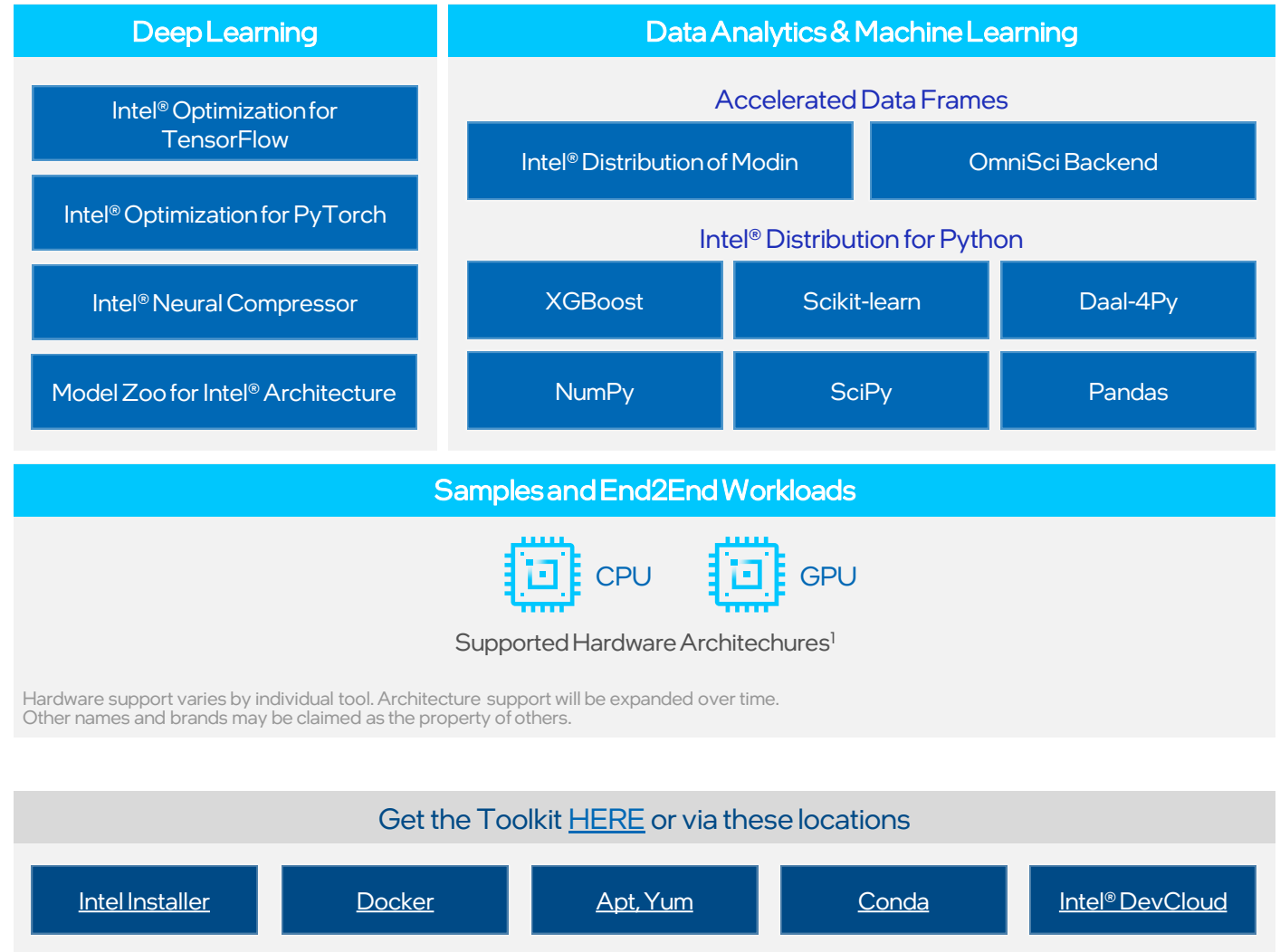
## Powered by oneAPI

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
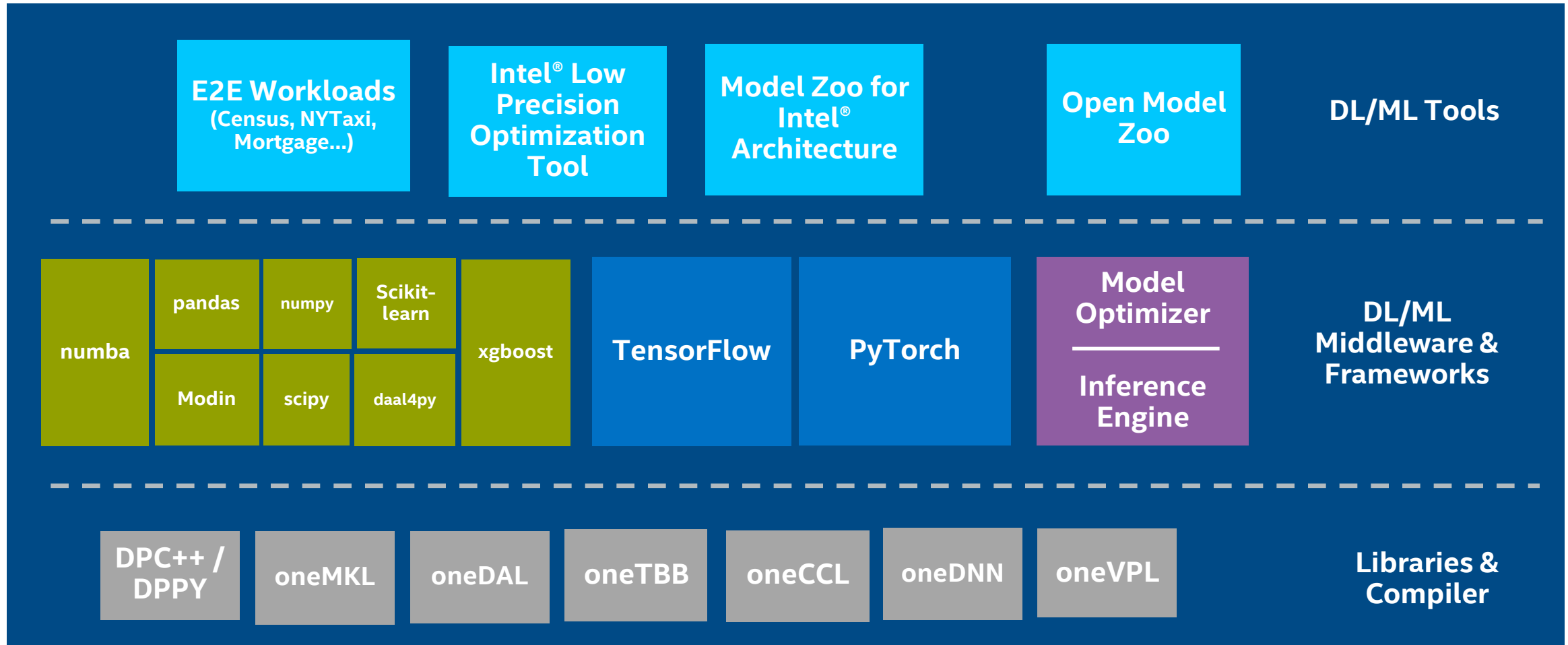
## Top Features/Benefits

- 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
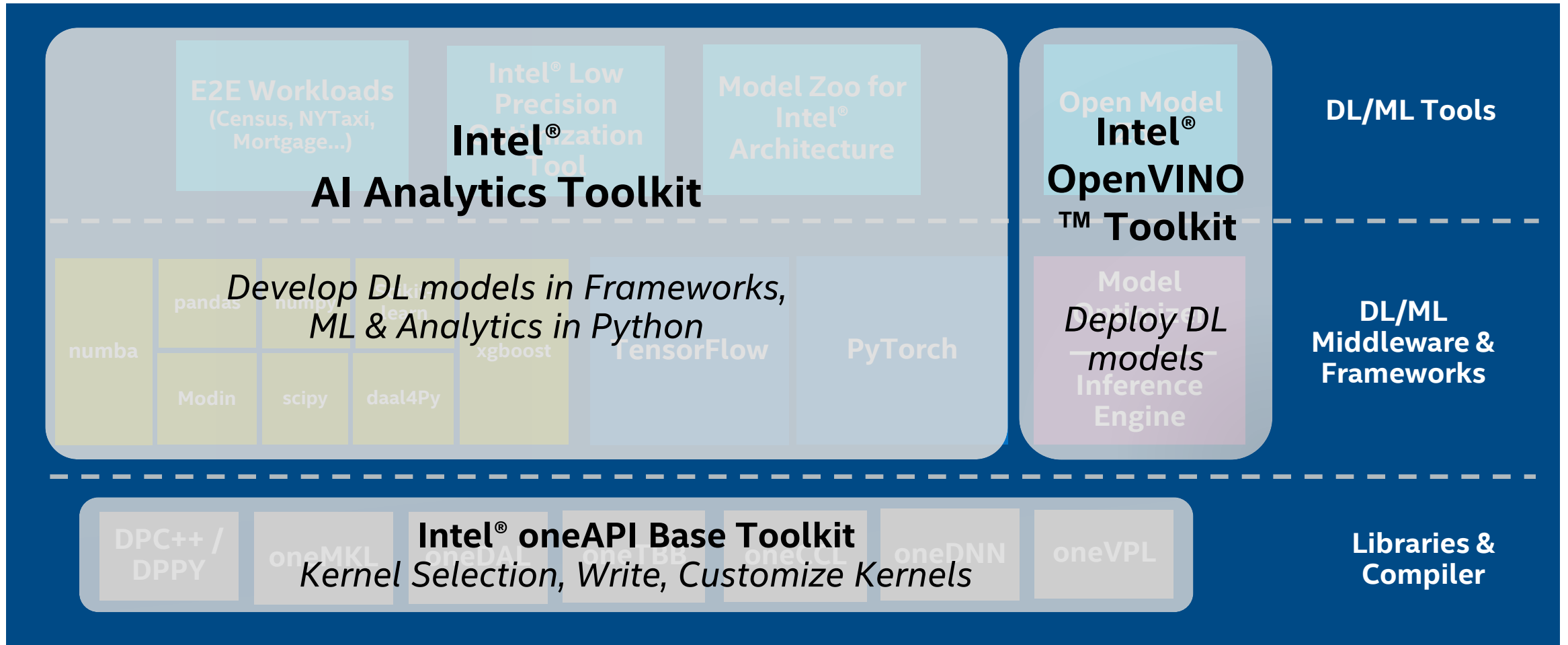
### Deep Learning

- Intel® Optimization for TensorFlow
- Intel® Optimization for PyTorch
- Intel® Neural Compressor
- Model Zoo for Intel® Architecture

### Data Analytics & Machine Learning

**Accelerated Data Frames**

| Intel® Distribution of Modin | OmniSci Backend |
|---|---|

**Intel® Distribution for Python**

| XGBoost | Scikit-learn | Daal-4Py |
|---|---|---|
| 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 |
|---|---|---|---|---|

# AI Software Stack for Intel® XPUs

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

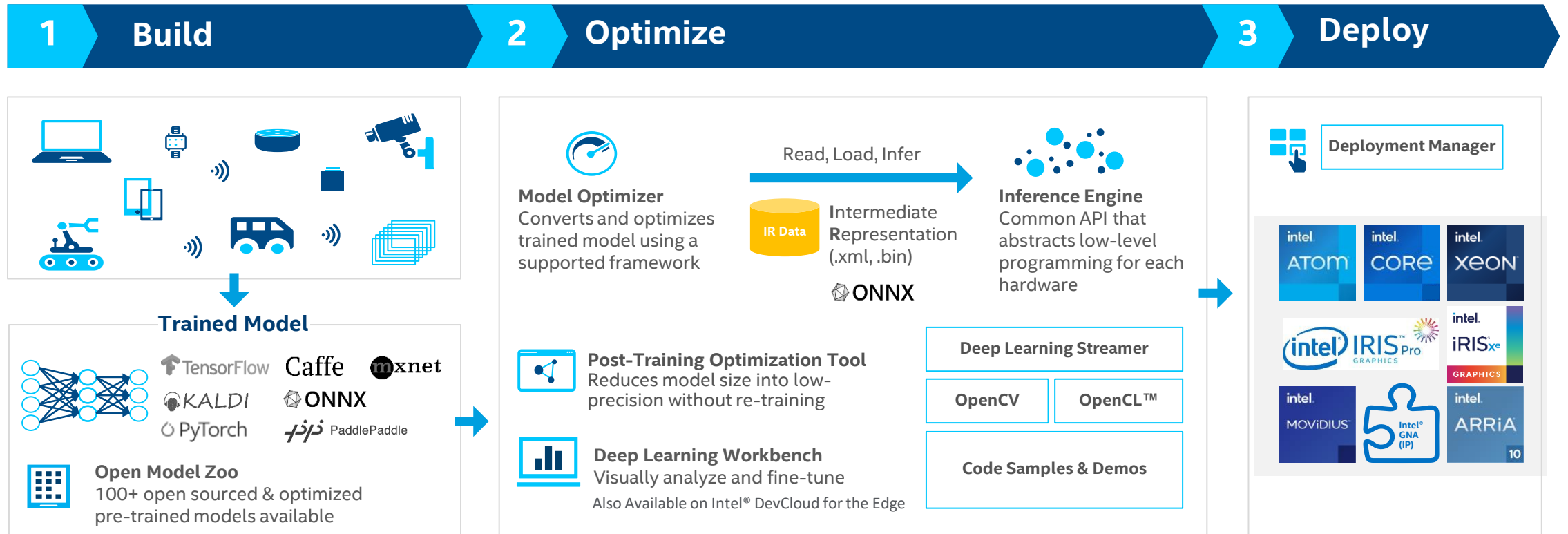| | | | | DL/ML Tools |
|---|---|---|---|---|
| **E2E Workloads** (Census, NYTaxi, Mortgage...) | **Intel® Low Precision Optimization Tool** | **Model Zoo for Intel® Architecture** | **Open Model Zoo** | |

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

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

# AI Software Stack for Intel® XPUs

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

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

Intel® Low Precision Optimization Tool

Model Zoo for Intel® Architecture

**Intel®
AI Analytics Toolkit**

Open Model

**Intel®
OpenVINO™ Toolkit**

**DL/ML Tools**

numba · pandas · numpy · 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**

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

| 1 Build | 2 Optimize | 3 Deploy |
| --- | --- | --- |

**Trained Model**

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

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

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

Read, Load, Infer

IR Data · **I**ntermediate **R**epresentation (.xml, .bin)

ONNX

**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
Also Available on Intel® DevCloud for the Edge

**Deep Learning Streamer**

**OpenCV** · **OpenCL™**

**Code Samples & Demos**

**Deployment Manager**

intel ATOM · intel CORE · intel XEON · intel IRIS Pro GRAPHICS · intel iRISxe GRAPHICS · intel MOViDIUS · Intel® GNA (IP) 5 · intel ARRiA 10

For workloads and configurations visit www.Intel.com/PerformanceIndex. Results may vary.

intel. 6

# Intel® Distribution for Python
## oneAPI Powered

Develop fast, performant Python code with this set of essential computational packages

## Who Uses It?

- Machine Learning Developers, Data Scientists, and Analysts can implement performance-packed, production-ready scikit-learn algorithms

- Numerical and Scientific Computing Developers can accelerate and scale the compute-intensive Python packages NumPy, SciPy, and mpi4py

- High-Performance Computing (HPC) Developers can unlock the power of modern hardware to speed up your Python applications

Initial GPU support enabled with Data Parallel Python

**~100 Packages Included**

### Intel® Distribution for Python

Numeric & Scientific Packages

| NumPy | SciPy | Numba |
|-------|-------|-------|

Machine Learning Packages

| Scikit-Learn | XGBoost* | daal4py |
|--------------|----------|---------|

Dataframe Packages

| Modin | Pandas | SDC |
|-------|--------|-----|

Runtimes

| OpenCL | DPC++ |
|--------|-------|

### Supported Hardware Architectures

CPU    GPU

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.

# Intel® Distribution for Python
## Developer Benefits

| Maximize Performance | Minimize Development Cost | Vast Ecosystem |
|---|---|---|
| **Performance Libraries, Parallelism, Multithreading, Language Extensions** | **Drop-in Python Replacement** | **Familiar usage and compatibility** |
| Near-native performance comes through acceleration of core Python numerical packages | Prebuilt optimized packages for numerical computing, machine/deep learning, HPC, & data analytics | Supports Python 3 |
| Accelerated NumPy/SciPy/scikit-learn with oneMKL & oneDAL | Data-Parallel Python provides cross-architecture XPU support | Supports conda & pip package managers |
| Data analytics, machine learning & deep learning with scikit-learn, XGBoost, Modin, daal4py | Conda build recipes included in packages | Packages available via conda, pip YUM/APT, Docker image on DockerHub |
| Scale with Numba*, Cython*, tbb4py, mpi4py, SDC | Free download & free for all uses including commercial deployment | Commercial support through the Intel® oneAPI Base Toolkit |

Optimized for latest Intel® architectures
Operating Systems: Windows*, Linux*, MacOS[1]*

**Intel® Architecture Platforms**     CPU     GPU     OTHER ACCEL.

# Intel® Modin Library

intel.

# Issue: Pandas Not Scaling to Larger Datasets

After a certain data size, need to change your API to handle more data

**100 MB+ of Data**

Easy to use,
difficult to scale

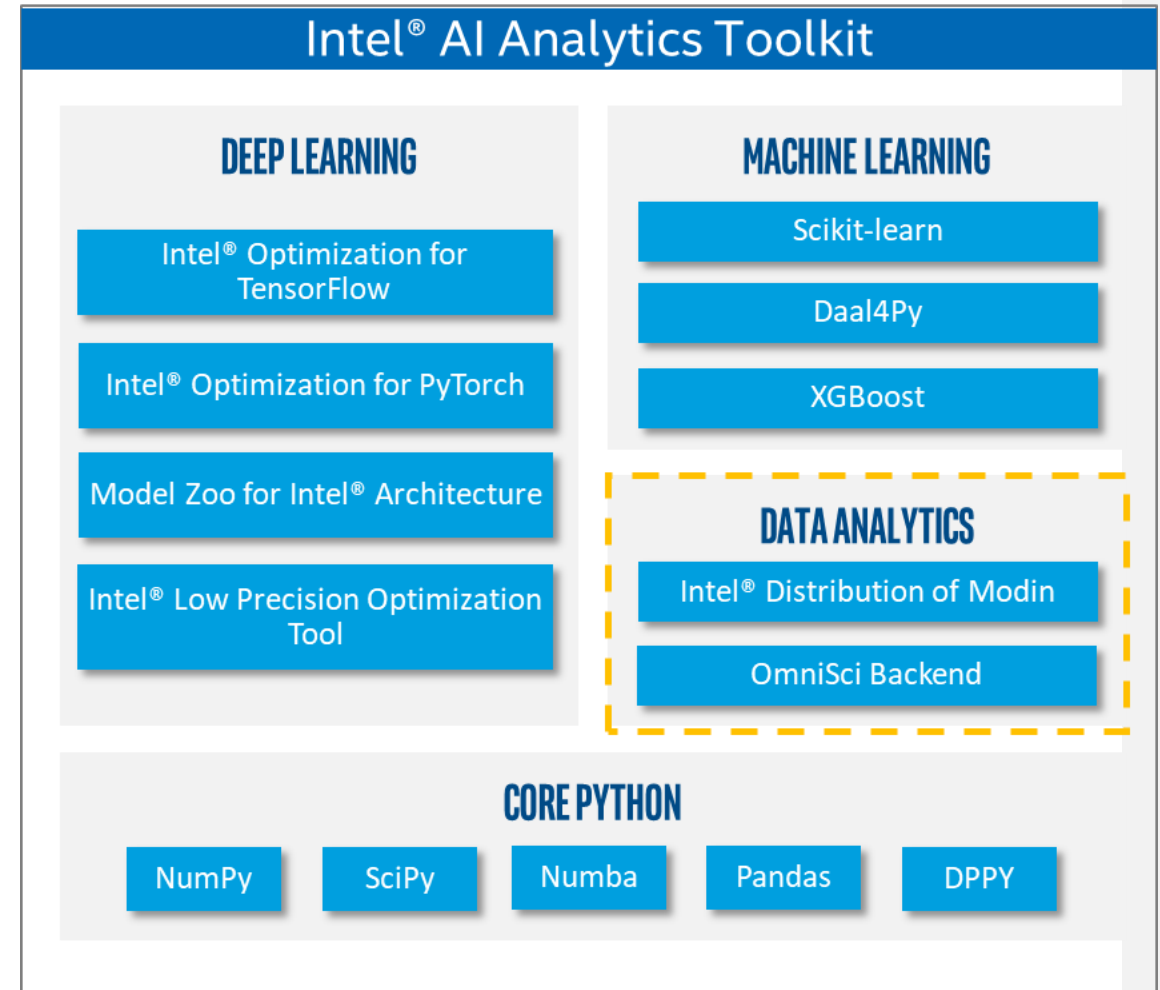**Increasing data size**

Easy to scale,
difficult to use

# Solution: Modin Pandas Scales to Big Datasets

Spend the time that would be used to change the workload's API, and use it to improve your workload and analysis
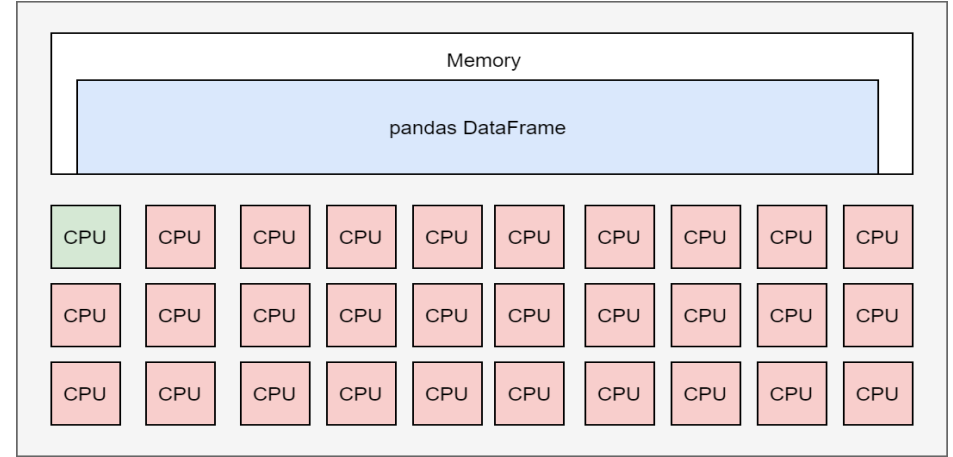


**Increasing data size**

Easy to use,
Easy to scale

**0-1TB+ of Data**

# Intel distribution of Modin

- Accelerate your Pandas* workloads across multiple cores and multiple nodes

- No upfront cost to learning a new API
  - import modin.pandas as pd

- In the backend, Intel Distribution of Modin is supported by Omnisci*, a performant framework for end-to-end analytics that has been optimized to harness the computing power of existing and emerging Intel® hardware

## Intel® AI Analytics Toolkit

### DEEP LEARNING

- Intel® Optimization for TensorFlow
- Intel® Optimization for PyTorch
- Model Zoo for Intel® Architecture
- Intel® Low Precision Optimization Tool

### MACHINE LEARNING

- Scikit-learn
- Daal4Py
- XGBoost

### DATA ANALYTICS

- Intel® Distribution of Modin
- OmniSci Backend

### CORE PYTHON

- NumPy
- SciPy
- Numba
- Pandas
- DPPY

# Intel distribution of Modin

- Recall: No upfront cost to learning a new API

  - import modin.pandas as pd

- Integration with the Python* ecosystem

- Integration with Ray*/Dask *clusters (Run on what you have, even on laptop!)

- To use Modin, you do not need to know how many cores your system has, and you do not need to specify how to distribute the data

**Pandas* on Big Machine**



**Modin on Big Machine**

# 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]



340,0729

**10.8 speedup**

31,2453

■ Pandas   ■ Modin

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

▪ Dataset size: 2.4GB

# Demo

intel.

# Intel® Extension for Scikit-Learn

intel.

# THE MOST POPULAR ML PACKAGE FOR PYTHON*

scikit learn

**Home**    **Installation**    **Documentation** ▼    **Examples**

Google Custom Search    **Search** ✕

## scikit-learn
*Machine Learning in Python*

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

## Classification

Identifying to which category an object belongs to.

**Applications**: Spam detection, Image recognition.
**Algorithms**: SVM, nearest neighbors, random forest, … — Examples

## Regression

Predicting a continuous-valued attribute associated with an object.

**Applications**: Drug response, Stock prices.
**Algorithms**: SVR, ridge regression, Lasso, … — Examples

## Clustering

Automatic grouping of similar objects into sets.

**Applications**: Customer segmentation, Grouping experiment outcomes
**Algorithms**: k-Means, spectral clustering, mean-shift, … — Examples

# Intel(R) Extension for Scikit-learn

### Common Scikit-learn

```
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

```
from sklearnex import patch_sklearn
patch_sklearn()

from sklearn.svm import SVC

X, Y = get_dataset()



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

**Available through:**
- conda install scikit-learn-intelex
- conda install –c intel scikit-learn-intelex
- conda install –c conda-forge scikit-learn-intelex
- pip install scikit-learn-intelex

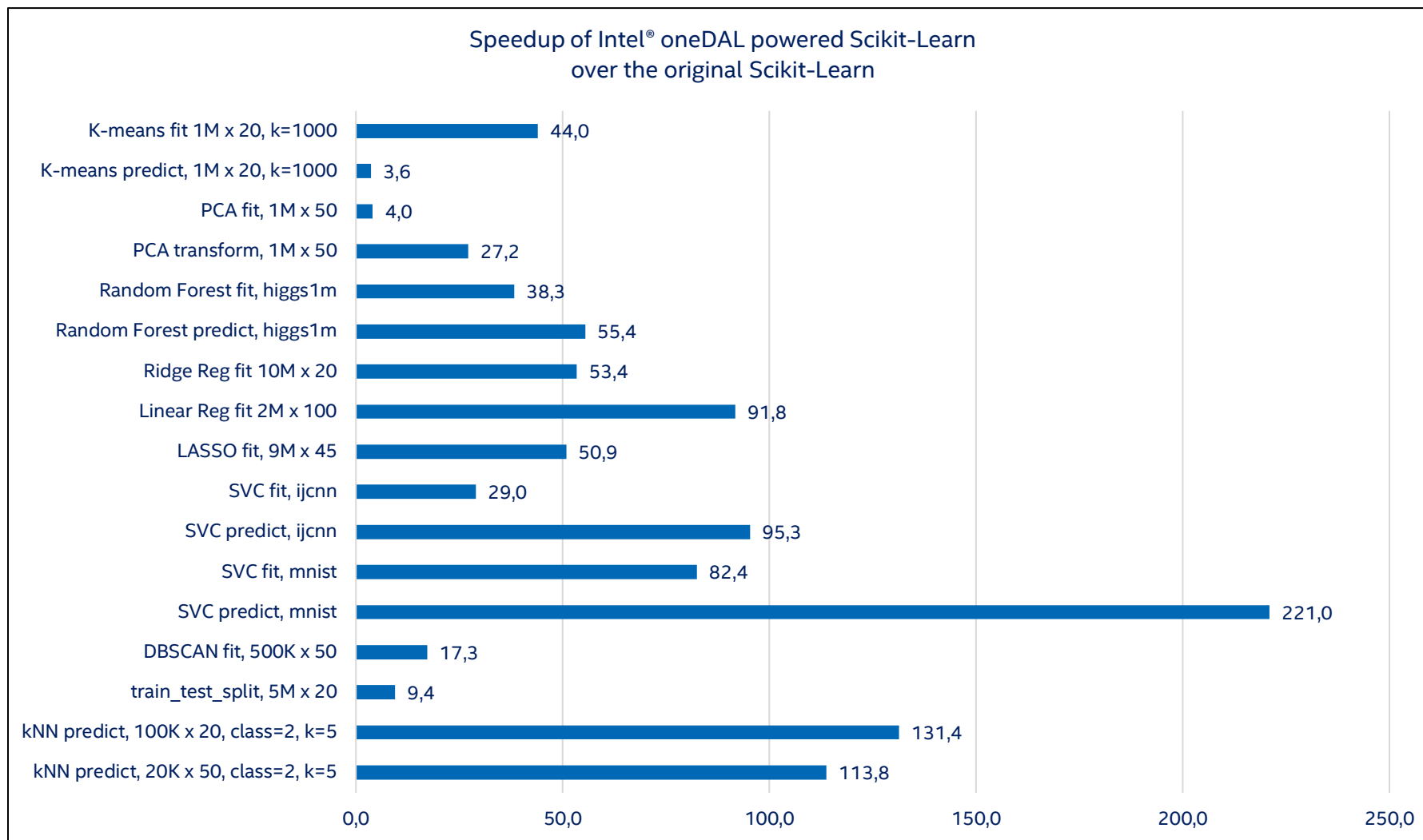## Same Code, Same Behavior

🔄 PASSED

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

# Available algorithms

- Accelerated IDP Scikit-learn algorithms:

- Linear/Ridge Regression

- Logistic Regression

- ElasticNet/LASSO

- PCA

- K-means

- DBSCAN

- SVC

- train_test_split(), assume_all_finite()

- Random Forest Regression/Classification – DAAL 2020.3

- kNN (kd-tree and brute force) – DAAL 2020.3

# 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 |

x-axis: 0,0 · 50,0 · 100,0 · 150,0 · 200,0 · 250,0

## Same Code, Same Behavior

PASSED

- Scikit-learn, not 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
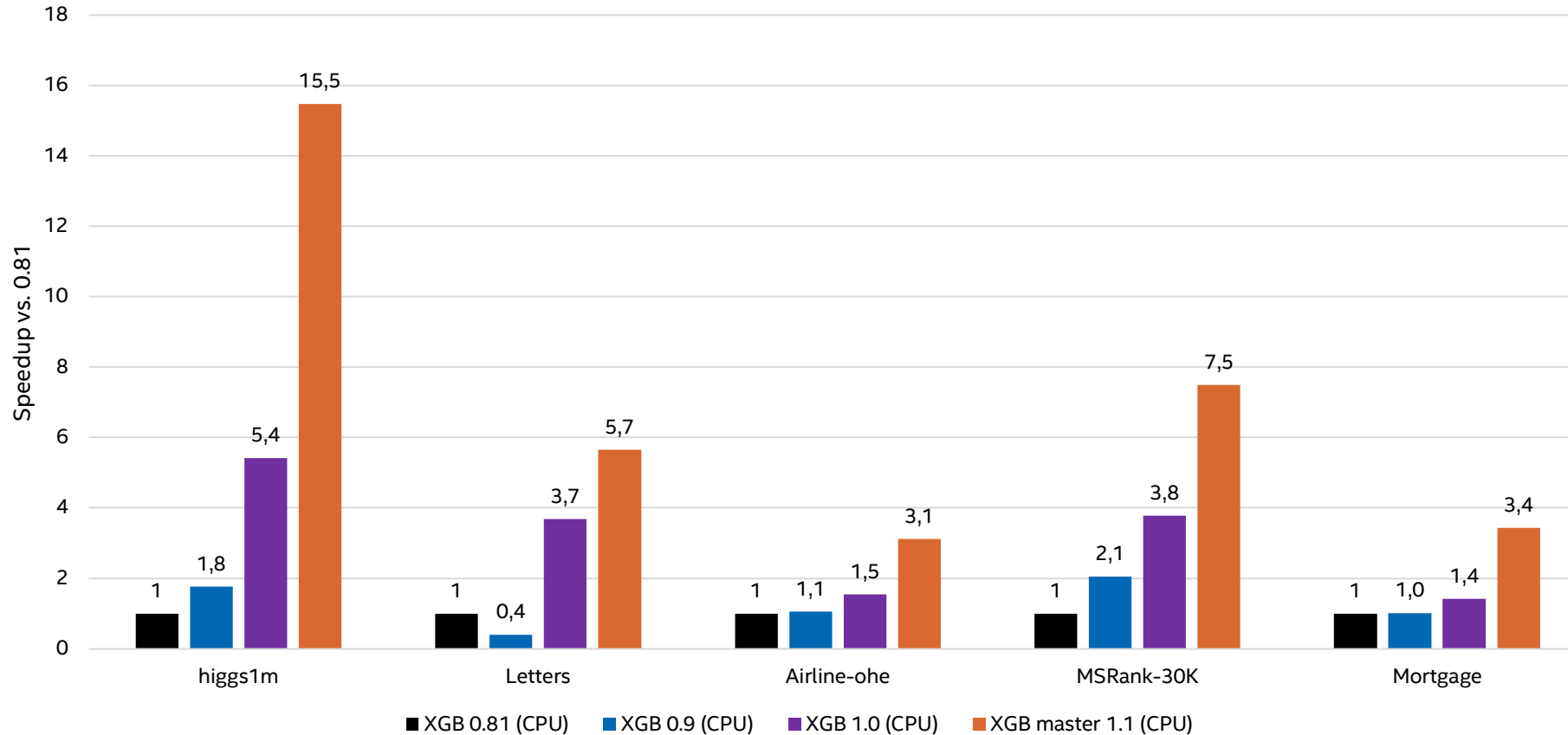
# Demo

intel

# XGBoost Library

intel.

# 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, …

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

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



**+ Reducing memory consumption**

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

Legend: ■ XGB 0.81 (CPU) ■ XGB 0.9 (CPU) ■ XGB 1.0 (CPU) ■ XGB master 1.1 (CPU)

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

# 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)

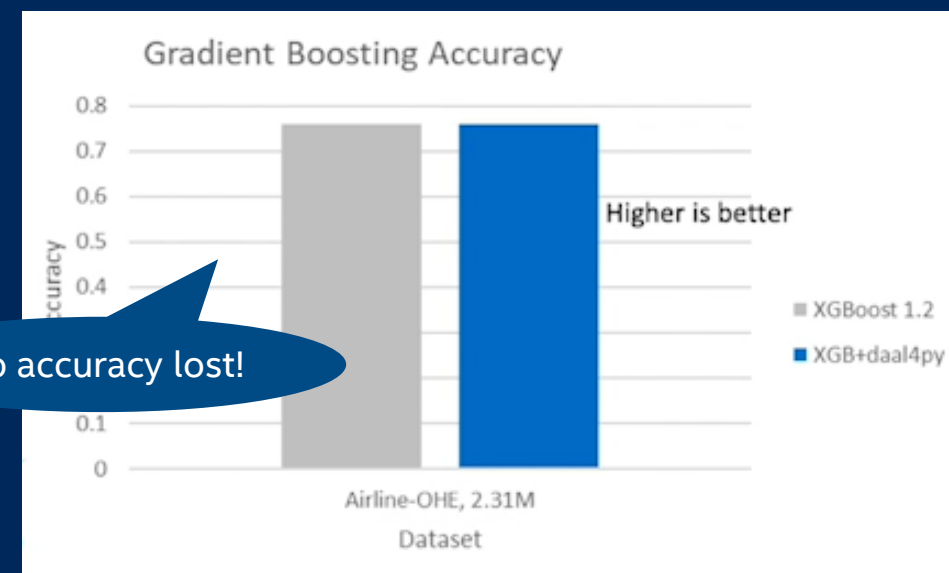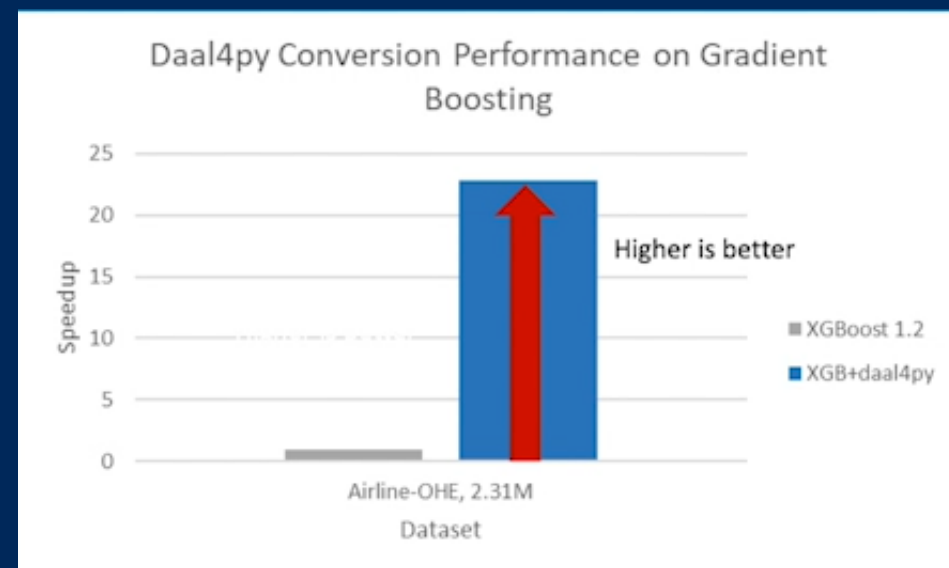Already available in Intel® oneDAL, potential optimizations for XGBoost*

# XGBoost* and LightGBM* Prediction Acceleration with Daal4Py

- Custom-trained XGBoost* and LightGBM* Models utilize Gradient Boosting Tree (GBT) from Daal4Py library for performance on CPUs

- No accuracy loss; 23x performance boost by simple model conversion into daal4py GBT:

```
# Train common XGBoost model as usual
xgb_model = xgb.train(params, X_train)

import daal4py as d4p

# XGBoost model to DAAL model
daal_model = d4p.get_gbt_model_from_xgboost(xgb_model)

# make fast prediction with DAAL
daal_prediction = d4p.gbt_classification_prediction(...).compute(X_test, daal_model)
```

- Advantages of daal4py GBT model:
  - More efficient model representation in memory
  - Avx512 instruction set usage
  - Better L1/L2 caches locality

For more complete information about performance and benchmark results, visit www.intel.com/benchmarks. See backup for configuration details.
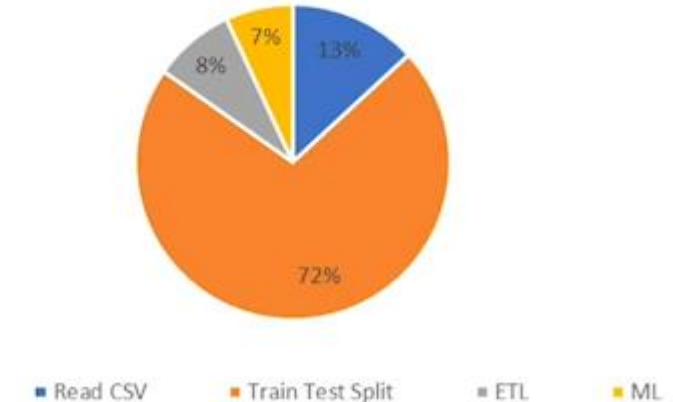


Daal4py Conversion Performance on Gradient Boosting

Higher is better

XGBoost 1.2
XGB+daal4py

Airline-OHE, 2.31M
Dataset



Gradient Boosting Accuracy

Higher is better

No accuracy lost!

XGBoost 1.2
XGB+daal4py
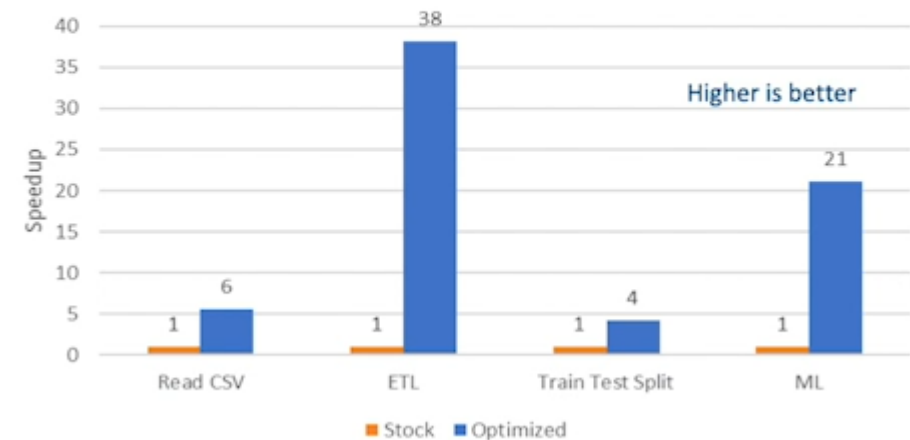
Airline-OHE, 2.31M
Dataset

# Demo

intel.

# End-to-End Data Pipeline Acceleration

- **Workload:** Train a model using 50yrs of Census dataset from IPUMS.org to predict income based on education

- **Solution:** Intel Modin for data ingestion and ETL, Daal4Py and Intel scikit-learn for model training and prediction

- **Perf Gains:**

  - Read_CSV (Read from disk and store as a dataframe) : **6x**

  - ETL operations : **38x**

  - Train Test Split : **4x**

  - ML training (fit & predict) with Ridge Regression : **21x**



End-to-End Time Breakdown : Census Education to Income



End-to-End Census: Speedup with optimized libraries

For more complete information about performance and benchmark results, visit www.intel.com/benchmarks. See backup for configuration details.

intel. 28

QnA

intel