**Choose the Best Accelerated Technology**

# Accelerate Machine Learning Workloads with the Intel AI Analytics Toolkit
## EUROCC AI Workshop

Roy Allela – AI Software Engineer

October 25th 2022

intel.

# Notices and Disclaimers

Performance varies by use, configuration and other factors. Learn more at www.Intel.com/PerformanceIndex
Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates.  See backup for configuration details.  No product or component can be absolutely secure.
You may not use or facilitate the use of this document in connection with any infringement or other legal analysis concerning Intel products described herein. You agree to grant Intel a non-exclusive, royalty-free license to any patent claim thereafter drafted which includes subject matter disclosed herein.
The products described may contain design defects or errors known as errata which may cause the product to deviate from published specifications.  Current characterized errata are available on request.
No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document, with the sole exception that a) you may publish an unmodified copy and b) code included in this document is licensed subject to the Zero-Clause BSD open source license (0BSD), https://opensource.org/licenses/0BSD. You may create software implementations based on this document and in compliance with the foregoing that are intended to execute on the Intel product(s) referenced in this document. No rights are granted to create modifications or derivatives of this document.
No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document, with the sole exception that code included in this document is licensed subject to the Zero-Clause BSD open source license (0BSD), http://opensource.org/licenses/0BSD.
No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.
Code names are used by Intel to identify products, technologies, or services that are in development and not publicly available. These are not "commercial" names and not intended to function as trademarks.
© 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.

# Agenda

- Intel® AI Analytics Toolkit

- Intel® Distribution for Python

- Intel® Distribution of Modin

- Intel® Extension for Scikit-learn

- Intel® Optimization for XGBoost

intel. 3

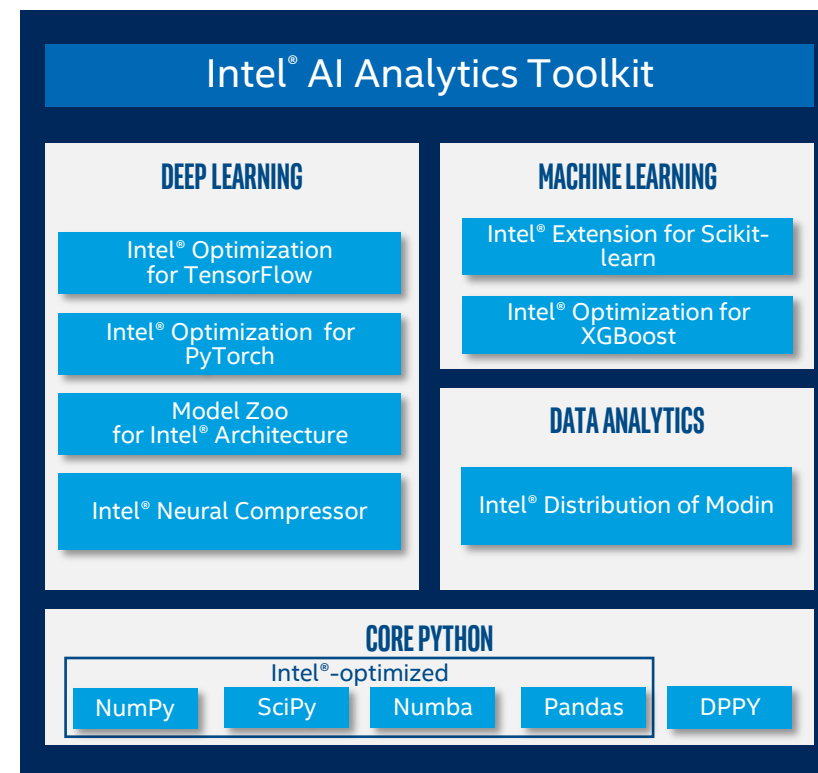# Intel® AI Analytics Toolkit
Powered by oneAPI

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

## Who Uses It?

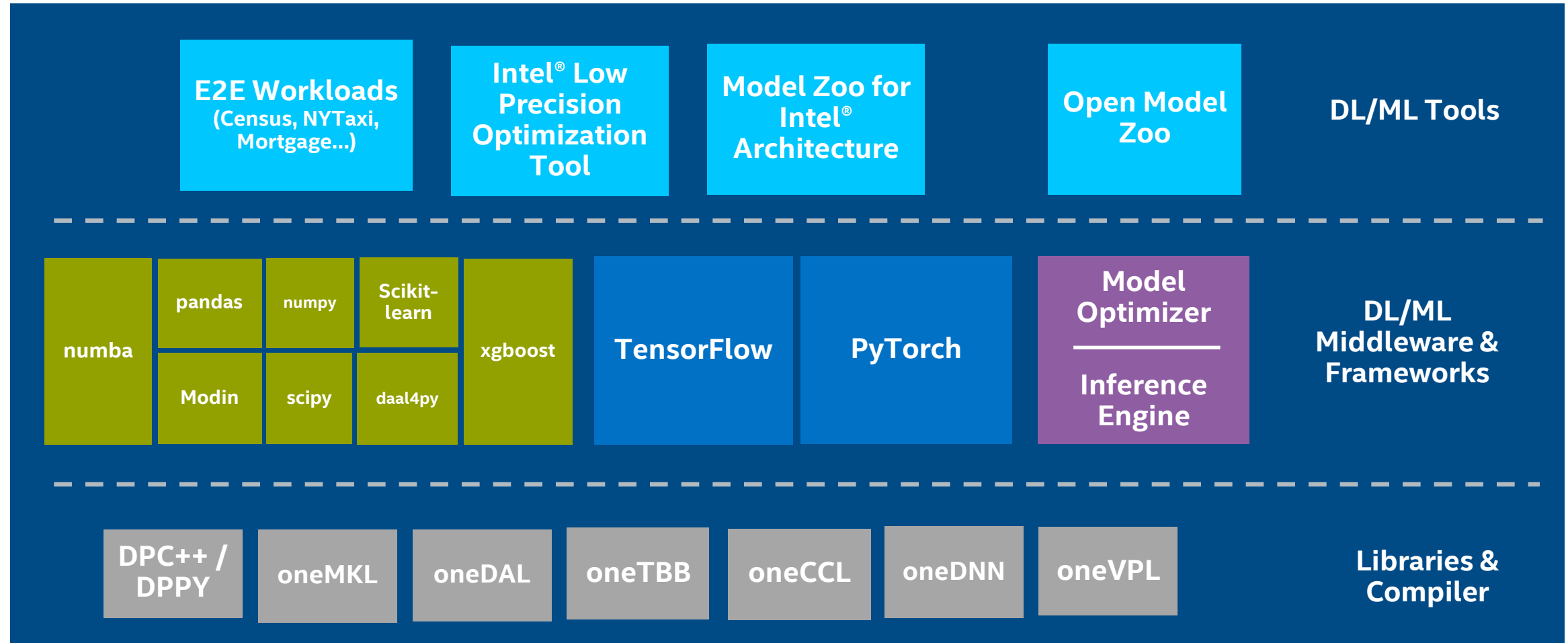Data scientists, AI researchers, ML and DL developers, AI application developers

## 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

### Intel® AI Analytics Toolkit

**DEEP LEARNING**

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

**MACHINE LEARNING**

- Intel® Extension for Scikit-learn
- Intel® Optimization for XGBoost

**DATA ANALYTICS**

- Intel® Distribution of Modin

**CORE PYTHON**

Intel®-optimized

| NumPy | SciPy | Numba | Pandas | DPPY |

# 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 |

# Why Use the Intel® AI Analytics Toolkit ?

## Accelerate Performance

Maximize machine learning performance for multiple architectures (Intel® CPUs/GPUs) with tools & components built using oneAPI libraries

## Streamline End2End Workflows

Get the latest AI Analytics optimizations in one place that work seamlessly together; scale end-to-end workflows fast

No need to download and integrate multiple external packages together

## Speed Development

Reduce the learning curve with drop-in replacement for Python packages with minimal to no code changes

Get started quickly with samples, pre-trained models, and end-to-end workloads

# 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

**~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 Performance Optimization with NumPy and SciPy

## The Layers of Quantitative Python

- The Python language is interpreted and has many type checks to make it flexible

- Each level has various tradeoffs; *NumPy\** value proposition is immediately seen

- For best performance, escaping the Python* layer early is best method

## Optimizations

- BLAS/LAPACK using oneMKL

- oneMKL-based FFT functionality

- Vectorized, threaded universal functions

- Use of Intel® C Compiler, and Intel® Fortran Compiler

- Aligned memory allocation

- Threaded memory copying

**Python**

Enforces *Global Interpreter Lock* (GIL) and is single-threaded, abstraction overhead. No

**NumPy**

Gets around the GIL (multi-thread and multi-core) *BLAS API* can be the bottleneck

**Intel® oneAPI Math Kernel Library (oneMKL)**

Gets around BLAS API bottleneck Much stricter typing Fastest performance level *Dispatches* to hardware vectorization

**oneMKL is included with Anaconda standard bundle; is Free for Python**

# Intel® Distribution of Modin

# 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**

**Increasing data size**

Easy to use,
difficult to scale

Easy to scale,
difficult to use

# Solution: Modin Pandas Scales to Big Datasets



**Increasing data size**

Easy to use,
Easy to scale

**0-1TB+ of Data**

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

# Intel Distribution of Modin

- Accelerate your Pandas workloads across multiple cores and multiple nodes

```python
import pandas as pd
|
```

- 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 a laptop!)

# Intel Distribution of Modin

- Modin transparently **distributes the data and computation across available cores**, unlike Pandas which only uses one core at a time

- 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()
```

- Dataset size: 2.4GB

### Execution time Pandas vs. Modin[ray]



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

# Intel® Extension for Scikit-Learn

# The Most Popular Package for Python



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

# Intel Extension for Scikit-Learn

## Performance on CLX compared to original Scikit-Learn : <u>Training</u>



Speedups of Intel® Extension for Scikit-learn training over the original Scikit-learn (higher is better)

# Intel Extension for Scikit-Learn
## Performance on CLX compared to original Scikit-Learn : <u>Inference</u>



Speedups of Intel® Extension for Scikit-learn inference over the original Scikit-learn (higher is better)

**Testing Date**: Performance results are based on testing by Intel as of June 8, 2021 and may not reflect all publicly available security updates.
**Configuration Details and Workload Setup:** c5.24xlarge AWS EC2 (3.0 GHz Intel Xeon Platinum 8275CL, two sockets, 24 cores per socket) Python 3.8, scikit-learn 0.24.2, scikit-learn-intelex 2021.2.3.
Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See configuration disclosure for details. Not product or component can be absolutely secure.
Performance varies by use, configuration, and other factors. Learn more at www.intel.com/PerformanceIndex. Your costs and results may vary

# XGBoost Optimizations

# 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® oneAPI
  Data Analytics Library (oneDAL), ...

intel

# Gradient Boosting Acceleration – Gain Sources

Pseudocode for XGBoost (0.81) implementation

Pseudocode for oneDAL implementation

```python
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)
```

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)

```python
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)
```

Training stage

Legend:

Moved from oneDAL to XGBoost (v1.3)

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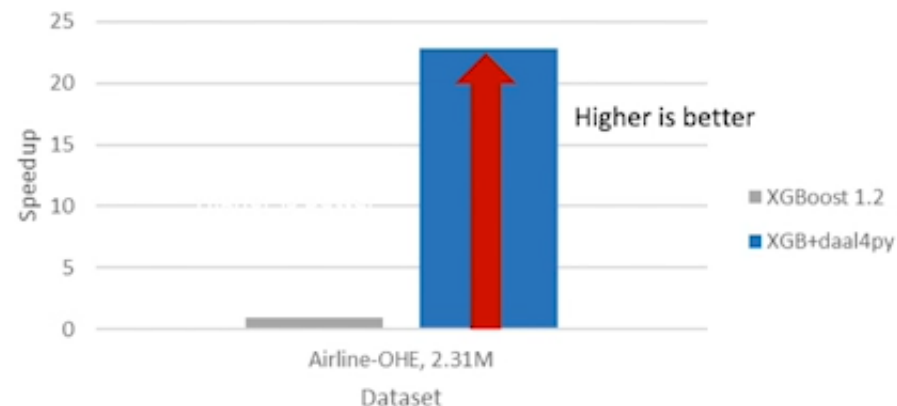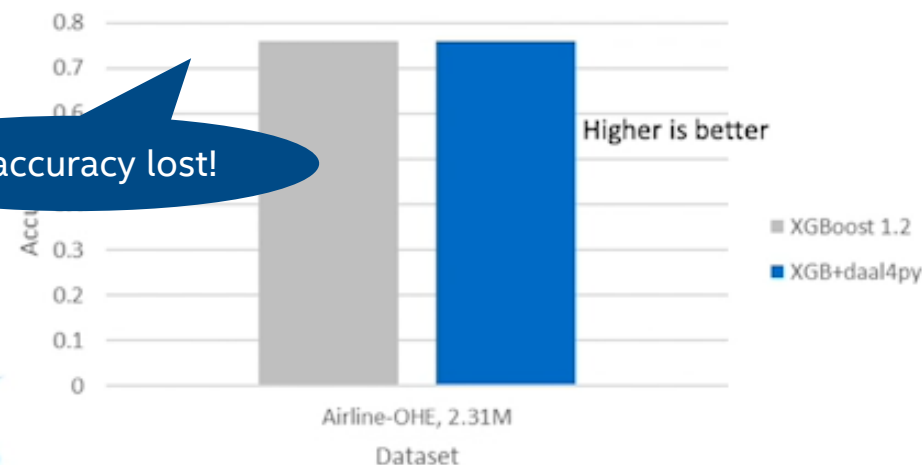
- Advantages of daal4py GBT model:
  - More efficient model representation in memory
  - Avx512 instruction set usage
  - Better L1/L2 caches locality

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



Gradient Boosting Accuracy

No accuracy lost!

Higher is better

# Demos

github.com/oneapi-src/oneAPI-samples/tree/master/AI-and-Analytics

# Takeaways

- Intel AI Analytics Toolkit offers tools to accelerate each stage of your pipeline: IDP, Modin, Intel Extension for Scikit-learn and XGBoost.

- Drop-in replacements with minimal code changes to speed up your development.

- Packages easily available through Conda, pip, Docker.

# QnA