

Accelerating the Open AI SW Ecosystem for AI Everywhere

Wei Li – VP/GM AI & Analytics



Imagine a world with Al everywhere...



AI Applications Today



3





5

... we are building the right Compute



...and models are growing in sizes & variety







Open Al Software Ecosystem Key Enabler



* Other names and brands may be claimed as the property of others

oneAPI



Open Specification for Accelerated Computing

Standards-Based Data Parallel Language

Standard Interfaces for Common Accelerator Libraries

Open-source implementations on diverse non-Intel CPU, GPU, FPGA, and AI solutions

An Open Project & Intel's Product



Intel's Implementation of the oneAPI Specification

First Customer Shipment – Dec 2020

Productive, Performant, Cross-Platform

Supports Intel CPU, GPU (integrated & discrete), and FPGA today

Realizing the vision of productive programming for accelerators, free from proprietary lock-in

11

inte

ΑΙΑ

Intel[®] oneAPI Software **Tools for AI & Analytics**

Popular AI frameworks and middleware are extended and optimized using one or more of the oneAPI industry specification elements

Can target CPUs, GPUs, and other accelerators



12

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 AI Software Strategy

to deliver

Simplicity, Productivity, and Performance

To go from Data By Optimizing End-to-End Workflows to Solutions Productivity Simplicity AI Performance For Every AI Workload **AIA**

intel

13

Performance

for Every AI workload







SSD-ResNet-34 Inference Throughput (Batch Size =1) For workloads and configurations visit www.intel.com/InnovationEventClaims. Results may vary.

15



Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

Leadership AI Performance

ResNet50 - Ponte Vecchio B step @1.4GHz vs A100(80G)





Intel data center GPU codenamed

Ponte Vecchio



Based on pre-production measurements. See backup for workloads and configurations. Results may vary. Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

AIA

AI







TensorFlow 2.9 now accelerated with Intel oneDNN by Default

Increased Productivity Faster Inference & Training Efficient Compute Utilization A erformance

Simplici

AI

Accelerating workflows for

laikA





•••

intel. 18

ÁİÁ

* Other names and brands may be claimed as the property of others

Simplicity

End-to-End Solutions For Every Industry



Apply Analytics & Machine Learning

on existing Intel environment



Build & Scale Quickly

With Optimized, Ready-to deploy solutions



Tangible Results

Without unnecessary complexity and specialized hardware





intel. + accenture

Industry Use Case Al Reference Kits

Open-Source and Prebuilt AI with Meaningful Enterprise Contexts



Productivity

AI

Simplicity

Jean-Luc Chatelain, CTO, Accenture Applied Intelligence

At Accenture, Jean-Luc Chatelain leads innovation and technology strategy focused on artificial intelligence, advanced analytics and cognitive automation. His team leads the development and roadmap of their flagship system of intelligence solution, the Applied Intelligence (AIP+).



intel. ²¹

Accenture AI VISION



All enterprises today must become digital



Al is the core engine that support this data driven transformation



Our AI vision is to enable our clients to deploy pragmatic AI to securely, responsibly augment the human business decision process to achieve best outcomes

intel.²²

AIA

oneAPI & AI





Speed is the #1 KPI for business



Enables you to deliver that speed at scale



Helps you

Hide plumbing complexity Simplify access speed Deploy everywhere



Intel + Accenture deliver acceleration with Industry Reference Kits

intel.²³

AIA

First Set of 4 Toolkits Releasing Today!

https://github.com/oneapi-src



24

Inside Our Al Reference Kits

intel. + accenture



Customer Care Agent Intent Enablement



To enable virtual agents to understand user intents in automated conversations using Natural Language Understanding (NLU). Customer care organizations need to reduce operational costs and yet offer a more natural and engaging conversational experience.

Experiment: build the NLU ML pipeline using airline travel dataset and PyTorch/BERT and use Intel Extension for PyTorch and Intel Neural Compressor boost inference times and reduce training cycles.



Productivity

Build Faster and Smarter



End-to-End Software Accelerators



Familiar Data Science Tools Optimized libraries and frameworks

5



One line of code

Unlocks End-to-End Performance gains



Source: https://techdecoded.intel.io/resources/one-line-code-changes-to-boost-pandas-scikit-learn-and-tensorflow-performance/#gs.bzkn2n

* Other names and brands may be claimed as the property of others

new



Hugging Face + intel.

Democratizing Accelerated Transformers on Intel Platforms

Inference Optimization Process from Days to Hours with Up to 4x Performance Speedup

intel

Xeon

Distributed Training through efficient compute scaling

:habana

Optimum Open-Source Library





* Other names and brands may be claimed as the property of others

intel.²⁹

AI & Big Data: Mastercard Recommender AI Service with Intel's BigDL for Accelerated End-to-End Development



Improving the customer experience, new product campaign performance, and better personalized recommendation accuracy by using the deep learning-based neural recommendation models



inte

ÂÂ

AI & High Performance Computing: CERN Large Hadron Collider (LHC) with Intel Neural Compressor for 10X Productivity





Sources: Conference paper: 10th International Conference on Pattern Recognition Applications and Methods Blog: https://www.nextplatform.com/2021/02/01/cern-uses-dlboost-oneapi-to-juice-inference-without-accuracy-loss

inte

ΑΙΑ

AI & Security: UPenn Federated Tumor Segmentation Initiative featuring Intel Federated Learning (OpenFL)





2.6%

BETTER

How much better does each institution do when training on the full data vs. just their own data?

17% better on the hold-out BraTS data on their own validation data



Secure Distributed Machine Learning

 Machine Learning collaboration without sharing sensitive data

(such as patient records, financial data, IP)

• Intel Software Guard Extensions protect data in use with HW-hardened enclaves

https://pypi.org/project/openfl/

TensorFlow OPyTorch K Keras

inte

AIA

intel. Al Builders

AI solutions accelerated by Intel



* Other names and brands may be claimed as the property of others

intel. ³³



* Other names and brands may be claimed as the property of others

Let's work together to bring Al Everywhere

Visit **developer.intel.com/ai** for more info



* Other names and brands may be claimed as the property of others

Thank You

intel