Intel Cloud Optimization Modules (ICOMs)



Date: August 22, 2023 Presenter: Benjamin Consolvo Al Solutions Engineer Manager, Intel

intel.

Agenda

- Executive Summary
- Technical Stack
- Content Package
- ICOM for AWS: XGBoost on SageMaker
- ICOM for AWS: XGBoost on Kubernetes
- ICOM for Microsoft Azure: XGBoost Pipeline on Kubernetes
- ICOM for Microsoft Azure: XGBoost Kubeflow Pipeline
- ICOM for AWS: GPT2-Small Distributed Training
- Next Steps

Executive Summary



What are Intel Cloud Optimization Modules?

The Intel Cloud Optimization Modules (ICOMs) are open-source codebases with codified Intel AI software optimizations and instructions built specifically for each Cloud Service Provider (CSP). The ICOMs are built with production AI developers in mind, leveraging popular AI frameworks within the context of cloud services.

Target Audience:

Enterprise Cloud AI Developer

Landing Page:

https://www.intel.com/content/www/us/en/developer/topic-technology/cloud-optimization.html

Technical Stack



Content Package

1. GitHub Repository

scaling_laws.ipynb	included distributed training NLP code	3 weeks ago
🗅 train.py	refactored the codebase using black	last week
🗅 trainer.py	updated readme	last week
Transformer_sizing.jpynb	included distributed training NLP code	3 weeks ago
README.md		0 ≔

intel.

Intel® Cloud Optimization Modules for AWS: GPT2-Small Distributed Training

The Intel Cloud Optimization Modules (ICOMs) are open-source codebases with codified Intel AI software optimizations and instructions built specifically for each Cloud Service Provider (CSP). The ICOMs are built with production AI developers in mind, leveraging popular AI frameworks within the context of cloud services.

Introduction

LLMs (large Language Models) are becoming ubiquitous, but in many cases, you don't need the full capability of the latest GPT models Additionally, when you have a specific task at hand, the perimanne of the biggest CPT model might not be optimal. Often, fine-tuning a small LLM on your dataset is sufficient. In this guide, you will learn how to fine-tune a GPT2-small (124M parameter) model on a cluster OFCVs on AWS. The objective here is not ta mire at a clartSPT-like A model, but rather to understand how to set up distributed training so that you can fine-tune to your specific objective. The end result of training here will result in a base LLM that can generate words (or tokens), but it will only be suitable for your secase when you uteria in a no your specific task and dataset.

The GPT2-Small model will be trained on the OpenWebText dataset in a distributed setting, using 3rd or 4th Gen. Intel® Xeon® Scalable Processors. The project builds upon the initial codebase of nanoGPT, by Andrej Karpathy.

2. Whitepaper

C GitHub

Solution Brief Cross-Industry Architecture Thete " Cloud Optimization Modules for AWS*: GPT2-Small Distributed Training

Authors: Ankur Singh, Benjamin Consolvo

Fine-tune GPT2-Small on the OpenWebText dataset in a distributed fashion on AWS on 3rd or 4th Generation Intel® Xeon® Scalable Processors.

Date: August 9, 2023

3. Cheatsheet

AWS Cluster Setup

For distributed training, our example deploys a cluster consisting of three 3rd Gen. Xeon CPUs (Amazon Elastic Compute Cloud* (EC2) <u>mói.4xlarge</u> instances) with an Ubuntu 22.04 AMI and 250 GB of storage. For maximum performance, we recommend using **bfloatl6** precision on the 4th Gen. Xeon® CPUs (<u>R7iz</u>) on AWS with the deep learning acceleration engine called <u>Intel® Advanced Matrix</u> Extensions (AMX).

If you are using a 4th Gen. Xeon CPU, you can verify AMX is present by running: lscpu | grep amx

You should see the following flags:

amx_bf16 amx_tile amx_int8

A construction A construction</l

4. Video Series (Coming Soon)

5. Office Hours – Registration



*Dependency on enterprise account managers

ICOM for AWS: XGBoost on SageMaker

<u>Overview</u>: SageMaker is a fully managed machine learning service on the AWS cloud. The motivation behind this platform is to make it easy to build robust machine learning pipelines on top of managed AWS cloud services. You can learn how to inject your custom training and inference code into a prebuilt SageMaker pipeline. This module enables you to use Intel® AI Analytics Toolkit accelerated software in SageMaker pipelines.

Open-Source Implementation: <u>https://github.com/intel/intel-cloud-optimizations-aws</u>



ICOM for AWS: XGBoost on Kubernetes

<u>Overview</u>: Build and deploy ML applications with XGBoost on AWS with Kubernetes with built-in Intel AI optimizations. We introduce the AWS services that we will use in the process, including Amazon Elastic Kubernetes Service (EKS), Amazon Elastic Container Registry (ECR), Amazon Elastic Compute Cloud (EC2), and Elastic Load Balancer (ELB).

Open-Source Implementation: <u>https://github.com/intel/intel-cloud-optimizations-aws</u>



ICOM for Microsoft Azure: XGBoost Pipeline on Kubernetes

Overview: Build and deploy highly available and scalable AI applications on Microsoft Azure with Kubernetes. The machine learning component of the module focuses on predicting the probability of a loan default using Intel® Optimization for XGBoost* and Intel® oneAPI Data Analytics Library (oneDAL) to accelerate model training and inference. We also demonstrate how to use incremental training of the XGBoost model as new data becomes available.

Open-Source Implementation: <u>https://github.com/intel/intel-cloud-optimizations-azure</u>



ICOM for Microsoft Azure: XGBoost Kubeflow Pipeline

<u>Overview</u>: Build and deploy accelerated AI applications on Kubeflow. The module is designed to maximize the performance and productivity of XGBoost with a loan default prediction problem. This set of reference architectures for Microsoft Azure also takes advantage of secure and confidential computing virtual machines leveraging Intel[®] Software Guard Extensions (Intel[®] SGX) on the Azure cloud.

Open-Source Implementation: https://github.com/intel/intel-cloud-optimizations-azure



ICOM for AWS: GPT2-Small Distributed Training

Overview: Learn how to fine-tune a GPT2-small (124M parameter) model on a cluster of CPUs on AWS. The objective here is not to arrive at a chatGPT-like AI model, but rather to understand how to set up distributed training so that you can fine-tune to your specific objective. The result will be a base LLM that can generate words (or tokens), but it will only be suitable for your use-case when you train it on your specific task and dataset. The GPT2-Small model is trained on the OpenWebText dataset in a distributed setting, using 3rd or 4th Gen. Intel® Xeon® Scalable Processors. The project builds upon the initial codebase of nanoGPT, by Andrej Karpathy.

Open-Source Implementation: <u>https://github.com/intel/intel-cloud-optimizations-aws</u>



And more!

Intel[®] Cloud Optimization Modules

This set of cloud-native open source reference architectures helps developers build and deploy Intel*optimized AI solutions on leading cloud providers, including Amazon Web Services (AWS)*, Microsoft Azure*, and Google Cloud Platform* service.

Each module or reference architecture includes a complete instruction set, all source code published on GitHub*, and a video walk-through.



Third-Party Cloud Optimization Modules

Microsoft Azure*

Amazon Web Services (AWS*)

Google Cloud Platform* Service

Amazon SageMaker* Amazon Elastic Kubernetes* Services Amazon Elastic Container Service (ECS) This module creates an ECS cluster using the latest Intel® Learn how to build and register images to Amazon Elastic Learn how to build accelerated machine learning pipelines Container Register (Amazon ECR) for the xgboost-daal4py app architecture available for the AWS ECS Service. using Amazon Elastic Kubernetes Services* (Amazon EKS) and the Lambda inference handler. with Intel® Optimization for XGBoost*. Amazon Elastic Compute Cloud (EC2) Amazon EKS Module with Terraform Amazon ElastiCache* for Redis This repository provides an example to create an Amazon EKS Configuration in this directory creates an AWS VM (instance). This module creates an Amazon ElastiCache* for Redis cluster cluster optimized on 3rd generation Intel® Xeon® Scalable based on Intel architecture and creates a new or existing VPC. The instance is created on a 3rd generation Intel® Xeon® processors (formerly code named Ice Lake). The example will This module uses the cache.r5.large by default, which is the Scalable processor by default. be creating an Amazon EKS cluster with an Amazon EKS latest Intel Xeon processor available at the time of this managed node group. module publication.

Next Steps

- Learn more about all of our Intel Cloud Optimization Modules here:
 - https://www.intel.com/content/www/us/en/developer/topic-technology/cloud-optimization.html
- Register for Office Hours here for help on your ICOM implementation:
 - https://software.seek.intel.com/SupportFromIntelExperts-Reg
- Come chat with us on our Intel DevHub Discord server to keep interacting with fellow developers:
 - https://discord.gg/rv2Gp55UJQ
- Stay connected with me on social media:
 - Benjamin Consolvo | AI Solutions Engineer Manager, Intel
 - LinkedIn: <u>https://linkedin.com/in/bconsolvo</u>
 - Twitter: <u>https://twitter.com/bpconsolvo</u>

