





COST EFFECTIVE MODEL DEPLOYMENT USING INTEL DEEP Learning deployment toolkit How to deploy alon your existing CPU



Prashant Shah, Intel Dave Chevalier, GE Healthcare - CT imaging 23-May-2018

(intel)



DEEP LEARNING @ GE HEALTHCARE CT



This Photo by Unknown Author is licensed under CC BY-SA



Tomo in the CAT scan machine Credit: Grahm S. Jones/Columbus Zoo and Aquarium

COMPUTERIZED AXIAL Tomography



COMPUTED TOMOGRAPHY (CT)

80+ million CT exams performed in United States each year

source: IMV 2016 CT Market Outlook

14% of Emergency Department (ED) visits result in a CT scan

source: "National Trends in Use of Computed Tomography in the Emergency Department." Annals of Emergency Medicine. Nov 2011.

Common CT procedures 1. Abdomen/Pelvis 2. Brain 3. Neck or Spine 4. Chest 5. CT Angiography source: IMV 2016 CT Market Outlook



PRINCIPLES OF COMPUTED TOMOGRAPHY





OurWorldInData.org/the-link-between-life-expectancy-and-health-spending-us-focus • CC BY-SA

AI AND HEALTHCARE ECONOMICS

CT EXAM COST BREAKDOWN

■ Radiologist ■ Other Personnel ■ Materials ■ Equipment

80% of the cost of a CT exam is human labor

Source: Dissecting Costs of CT Study: Application of TDABC (Time-driven Activity-based Costing) in a Tertiary Academic Center. Anzai Y, Heilbrun ME2, Haas D, Boi L, Moshre K, Minoshima S, Kaplan R, Lee VS. 2017 Feb; 24(2):200-208. doi: 10.1016/j.acra.2016.11.001. Epub 2016 Dec 14.



Productivity = cost savings

Automation of routine tasks - AI is the enabler

Inference deployment – keep the costs down

- Deploy AI on existing CPU
- CPU-based inference allows embedded AI across the full range of our applications and products

Reduce OpEx without increasing CapEx



AI: PUTTING THE "SEE" IN "CT"

AXIAL SLICE CLASSIFICATION



POTENTIAL USE CASES

- IMAGE QUALITY OPTIMIZATION ANATOMY SPECIFIC TECHNIQUES
- QUALITY ASSURANCE DID I SCAN WHAT I WANTED TO SCAN?
- POST-PROCESSING APPLICATION SELECTION
- DATA TAGGING FOR LATER RETRIEVAL, RESEARCH, AI TRAINING
- CONTENT-BASED RETRIEVAL COMPARE TO NORMALS



CT AXIAL IMAGE CLASSIFIER

ANATOMY DATASET:

bit

- 223 EXAMS, ~30,000 IMAGES •
- EACH AXIAL CT IMAGE LABELED INTO 6 ANATOMIC • REGIONS



INFERENCE PERFORMANCE GOAL

Clinical Goal: Keep pace with the imaging pipeline



Performance target: 100 images/sec or ~10msec latency

- up to 4 cores available without impacting the imaging pipeline



BASELINE BENCHMARKING WITH TENSORFLOW



Images/Sec (on 4 Intel Xeon[®] E5-2650 cores)

Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system. Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit. https://www.intel.com/performance

Optimization Notice: Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations, Intel does not guarantee availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessors, Presser effect to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by Unis not sets.

Configuration Details: * Xeon* processor E5-2650 v4 at 2.20GHz, and configured with 264 GB of memory, Intel* Solid State Drive Data Center 480 GB, and CentOS Linux* 7.4. 708; Tensorflow version 1.4 Compiled with MKL-DNN. Model and Dataset: GE proprietary





BASELINE BENCHMARKING WITH INTEL® DEEP LEARNING DEPLOYMENT TOOLKIT® (PART OF INTEL OPEN VINO™)



Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system. Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit: http://www.intel.com/performance

Optimization Notice: Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarante availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessors dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations on sets and other optimizations instruction sets covered by this induce.

Configuration Details: * Xeon* processor E5-2650 v4 at 2.20GHz, and configured with 264 GB of memory, Intel* Solid State Drive Data Center 480 GB, and CentOS Linux* 7.4.1708; Tensorflow version 1.4 Compiled with MKL-DNN. Model and Dataset: GE proprietary





OPTIMIZATIONS DONE ON THE GE AXIAL CT MODEL:



Images/Sec (4 Intel Xeon® E5-2650v4 Cores)



Batch size = 64



Exceeds GE's requirement with just 1 Intel® Xeon® core

Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system. Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit. https://www.intel.com/performance

Optimization Notice: Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations, Intel does not guarantee availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessors, Please refer to the applicable product User and Reference Guides for more information regarding the society is covered by this notice.





DEEP LEARNING DEPLOYMENT TOOLKIT (PART OF INTEL[®] OPEN VINO[™]



DEEP LEARNING DEPLOYMENT TOOLKIT INFERENCE WORKFLOW



Initialization

Load model and weights Set batch size (if needed) Load Inference Plugin (CPU, GPU, FPGA) Load network to plugin Allocate input, output buffers

Main loop

Fill input buffer with data Run inference Interpret output results







^aAlpha available †Beta available

[‡] Future

*Other names and brands may be claimed as the property of others. All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.

SUMMARY

- Using Intel Deep Learning Deployment Toolkit[®], we were able to achieve 6X the target performance, and 10X the performance over native Tensorflow without the need of add-in cards.
- Inference on CPUs provides GE the flexibility of deployment architecture. CPUs are ubiquitous in embedded devices, workstations and datacenter/cloud.
- DL DT can be used to efficiently deploy models on CPUs, Intel Integrated GPUs, FPGAs and Movidius inference accelerators.















(intel) Al

