# embedded VISIMA Summit

(intel)

Getting Efficient DNN Inference Performance: Is it Really All About the TOPS?

Gary Brown, Intel IOT Group September 22, 2020



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 www.intel.com/benchmarks.

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.

Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

© 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.

Intel is committed to respecting human rights and avoiding complicity in human rights abuses. See Intel's Global Human Rights Principles <u>https://www.intel.com/content/www/us/en/policy/policy-human-rights.html</u> Intel's products and software are intended only to be used in applications that do not cause or contribute to a violation of an internationally recognized human right.







The AI Inference Landscape and Growth of the Edge

Agenda



What about Benchmarking



Comparing AI Inference Silicon: Is it Really All About the TOPS



A New Era of Efficient Architectures









## **Edge Al Silicon Opportunity Is Growing Fast**





25%-35% CAGR In Al For Silicon Industry

Growth larger at the Edge, Compared To Data Center



### Al Inference at the Edge has Many Challenges





embedded





# How We Benchmark Silicon Platforms – What's Missing?



MIPS or IPC

### **TOPS or TFLOPS**







MLPerf





embedded

# **Benchmarks Making Progress, and Scope is Getting Larger**

#### embedded VISICN summit

### **Dimensions of Comparison:**

- Flexibility
- Compute Efficiency
- IOT Fortified
- Security
- Software Portability
- ...and many more





(intel)











### **An Example: Which is Higher Performance?**



Work accomplished per second

Child: 10 scoops/min

Commercial excavator: 2 scoop/min

Excavator: 10X slower?!



### Which Metric You Use Makes all the Difference







# A New Era of Efficient Architectures



## **Edge AI with Intel: Examining Compute Efficiency**





## Al Acceleration with 11<sup>th</sup> Gen Intel<sup>®</sup> Core<sup>®</sup> Processors

- The latest Intel® Core<sup>™</sup> processor for small form-factor servers, such as for digital signage with Al
- Optimized for industrial, public sector, transportation, retail, video analytics NVRs
- An efficient CPU with AI acceleration built-in using Intel® Deep Learning Boost
- Additional AI inference acceleration with Xe architecture





ember



## Intel<sup>®</sup> Xeon<sup>®</sup> Processors with Built-in AI Acceleration



- For larger servers with more video channel density
- DL inference built-into the CPU without the need for external AI accelerator
- Enough DL performance to run real-time inference across many channels of video
- No added Al accelerator card required.

#### **Optimized Software and Frameworks:**



### Caffe mxnet OpenVINO 🚱 ONNX 🙀 PaddlePaddle ÖPyTorch া 🕆 TensorFlow

Performance results are based on testing as of 12/25/2019 and may not reflect all publicly available updates. See configuration disclosure for details. No product or component can be absolutely secure. For more complete information about performance and benchmark results, visit <u>www.intel.com/benchmarks</u>. Refer to <u>http://software.intel.com/en-us/articles/optimization-notice</u> for more information regarding performance and optimization choices in Intel software products.



Configuration Details: 1-node, 2x Intel® Xeon® Gold 6258R on Intel Reference platform with 384 GB (12 slots / 32 GB / 2933) total memory, ucode 0x500002c, HT on, Turbo on, with Ubuntu19.10, 5.3.0-24-generic, AIXPRT Image Classification AIXPRT v1.01 Intel® distribution of OpenVINO version 2019 R3 ResNet50 v1, using INT8 with Intel® DL Boost, BS=4, no.of instances =56 as of 01/29/2020, test by Intel on 12/25/2019. 1-node, 2x Intel® Xeon® Gold 6152 on Intel Reference platform with 384 GB (12 slots / 32 GB / 2933) total memory, ucode 0x2000065, HT on, Turbo on, with Ubuntu19.10, 5.3.0-24-generic, Image Classification AIXPRT v1.01 Intel® distribution of OpenVINO version 2019 R3 ResNet50 v1, using INT8 with Intel® DL Boost, BS=4, no.of instances =56 as of 01/29/2020, test by Intel on 12/25/2019. 1-node, 2x Intel® Xeon® Gold 6152 on Intel Reference platform with 384 GB (12 slots / 32 GB / 2933) total memory, ucode 0x2000065, HT on, Turbo on, with Ubuntu19.10, 5.3.0-24-generic, Image Classification AIXPRT v1.01 Intel® distribution of OpenVINO version 2019 R3 ResNet-50 v1, Using INT8, BS=128, no.of instances =44 as of 01/29/2020, test by Intel on 12/25/2019.

## Gen 3 Intel<sup>®</sup> Movidius<sup>™</sup> VPU for Power-Efficient Edge AI

embedded VISICN Summit



### **How Developers Can Get Productive Quickly**



## INTEL® DEVCLOUD FOR THE EDGE

https://devcloud.intel.com/edge







- Simple benchmarks such as TOPS aren't effective at comparing architectures
- Better comparison is workload or application performance
- And there are many more dimensions of value in addition to performance, such as compute efficiency
- Intel's CPUs enable built-in AI inference, and VPUs offer efficiency
- Getting started is easy with Intel<sup>®</sup> DevCloud for the Edge









Intel® vision and edge AI products

Getting started with Intel® DevCloud

OpenVINO<sup>™</sup> toolkit documentation

https://www.intel.com/visionproducts

https://devcloud.intel.com/edge

https://docs.openvinotoolkit.org

