**Strata Data & AI Superstream Series**

# AI Inferencing with NLP at Scale with OpenVINO

Zoe Cayetano and Raymond Lo
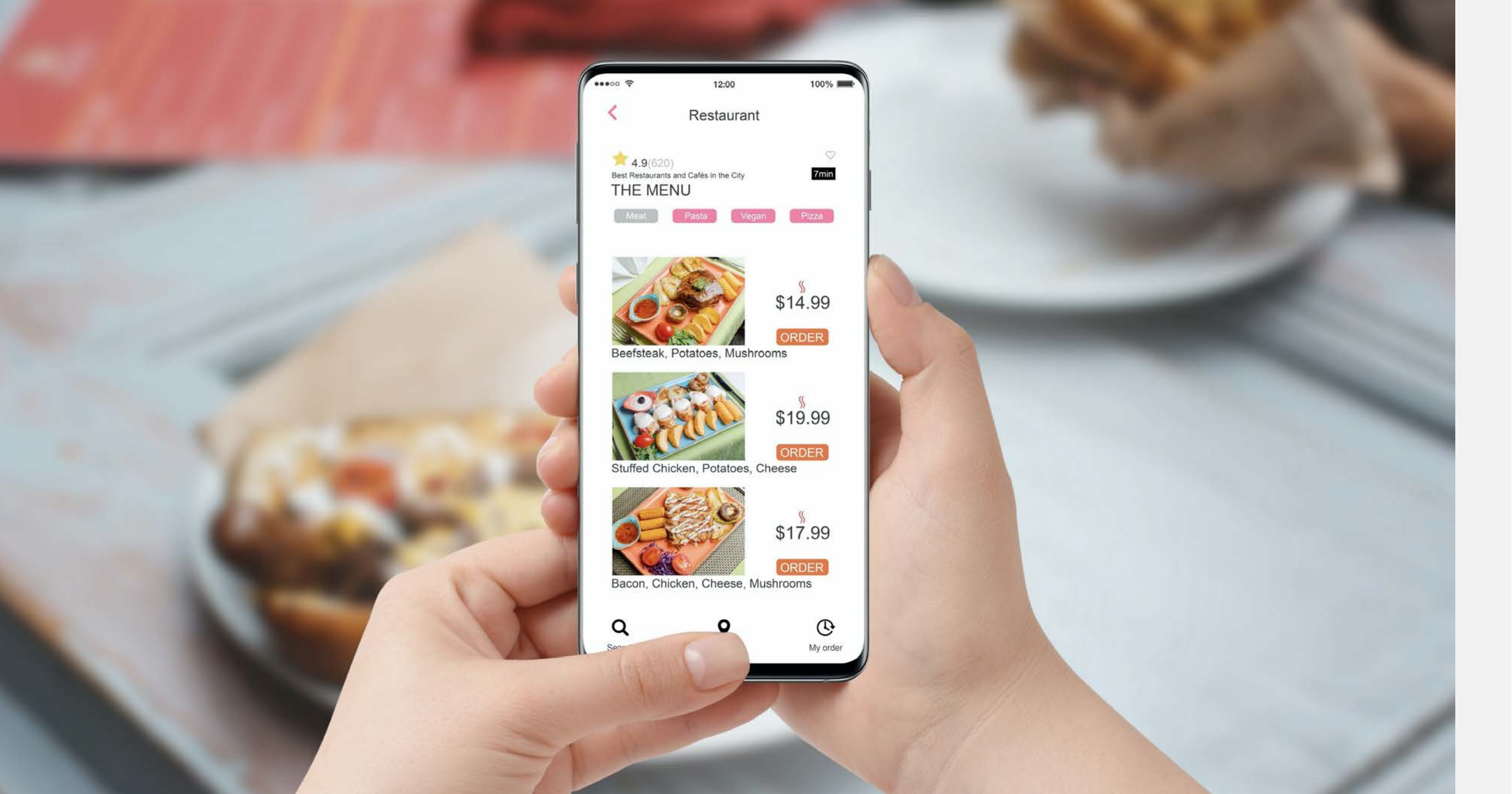
intel.

# Let's Make "Pasta" with NLP

intel.

intel.

3

# Voice-based Interactive Recipes App for Cooking
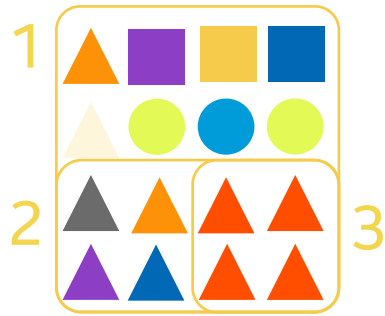
intel.

# Where we do get our data?

intel.

# AI Compute Considerations

## How do you determine the right computing for your AI needs?
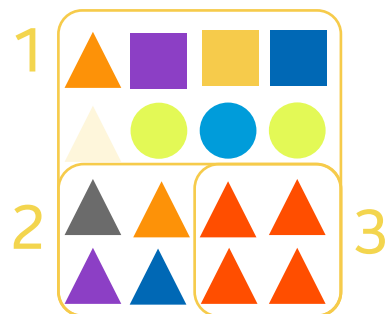


WORKLOADS

# AI Compute Considerations

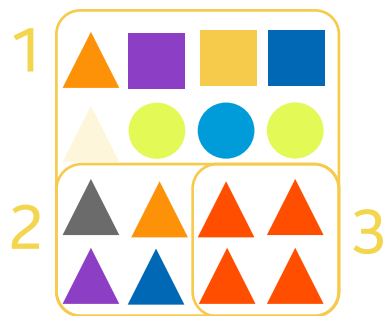## How do you determine the right computing for your AI needs?



WORKLOADS

REQUIREMENTS

# AI Compute Considerations

## How do you determine the right computing for your AI needs?
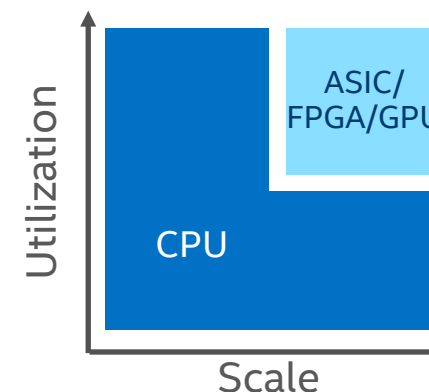


**WORKLOADS**

**REQUIREMENTS**

**DEMAND**

# Deep Learning Development Cycle

# Training vs Inference



Human

Bicycle

Strawberry

Lots of Labeled Data!

# Training vs Inference



*Human*

*Bicycle*

*Strawberry*

Lots of Labeled Data!

# Training vs Inference



*Human*

*Bicycle*

*Strawberry*

Lots of Labeled Data!

Forward

Backward

intel.

# Training vs Inference



*Human*

*Bicycle*

*Strawberry*

Lots of Labeled Data!

Forward

Backward

"Bicycle"

intel.

# Training vs Inference

*Human*

*Bicycle*

*Strawberry*

Lots of Labeled Data!

Forward

Backward

"Bicycle"

Model Weights

intel.

# Inference

??????

Inference

??????

Forward

# Inference



*??????*

Forward

90% = "Bicycle" ?

# Inference

??????

intel.

# Inference

??????

Forward

Inference

?????? 

Forward

99% = "Canon in D" ?

intel.

# Bert Fine-Tuning (SQuAD)



START / END SPAN

BERT

| C | $T_1$ | $T_N$ | $T_{[SEP]}$ | $T_1'$ | $T_M'$ |

| $E_{[CLS]}$ | $E_1$ | $E_N$ | $E_{[SEP]}$ | $E_1'$ | $E_M'$ |

| [CLS] | Tok 1 | Tok N | [SEP] | Tok 1 | Tok M |

QUESTION

PARAGRAPH

*Intel was founded in Mountain View, California, in 1968 by Gordon E. Moore (known for "Moore's law"), a chemist, and Robert Noyce, a physicist and co-inventor of t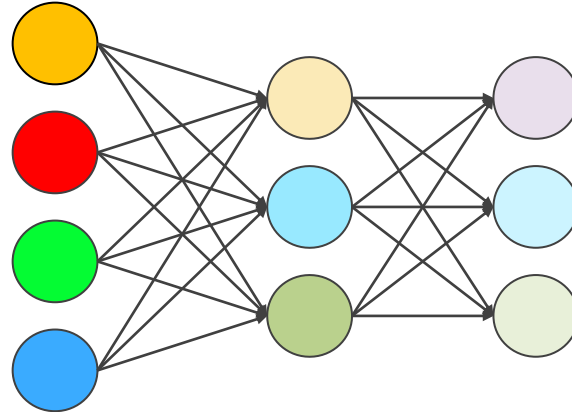he integrated circuit. Arthur Rock (investor and venture capitalist) helped them find investors, while Max Palevsky was on the board from an early stage.[23] Moore and Noyce had left Fairchild Semiconductor to found Intel. Rock was not an employee, but he was an investor and was chairman of the board.[24][25] The total initial investment in Intel was $2.5 million in convertible debentures (equivalent to $18.4 million in 2019) and $10,000 from Rock. Just 2 years later, Intel became a public company via an initial public offering (IPO), raising $6.8 million ($23.50 per share).[24] Intel's third employee was Andy Grove,[26] a chemical engineer, who later ran the company through much of the 1980s and the high-growth 1990s. ...*

START / END SPAN

BERT

QUESTION

PARAGRAPH

When was Intel founded?

Intel was founded in Mountain View, California, in 1968 by Gordon E. Moore (known for "Moore's law"), a chemist, and Robert Noyce, a physicist and co-inventor of the integrated circuit. Arthur Rock (investor and venture capitalist) helped them find investors, while Max Palevsky was on the board from an early stage.[23] Moore and Noyce had left Fairchild Semiconductor to found Intel. Rock was not an employee, but he was an investor and was chairman of the board.[24][25] The total initial investment in Intel was $2.5 million in convertible debentures (equivalent to $18.4 million in 2019) and $10,000 from Rock. Just 2 years later, Intel became a public company via an initial public offering (IPO), raising $6.8 million ($23.50 per share).[24] Intel's third employee was Andy Grove,[26] a chemical engineer, who later ran the company through much of the 1980s and the high-growth 1990s. ...

START / END SPAN

BERT

QUESTION          PARAGRAPH

Intel was founded in Mountain View, California, in **1968** by Gordon E. Moore (known for "Moore's law"), …

When was Intel founded?
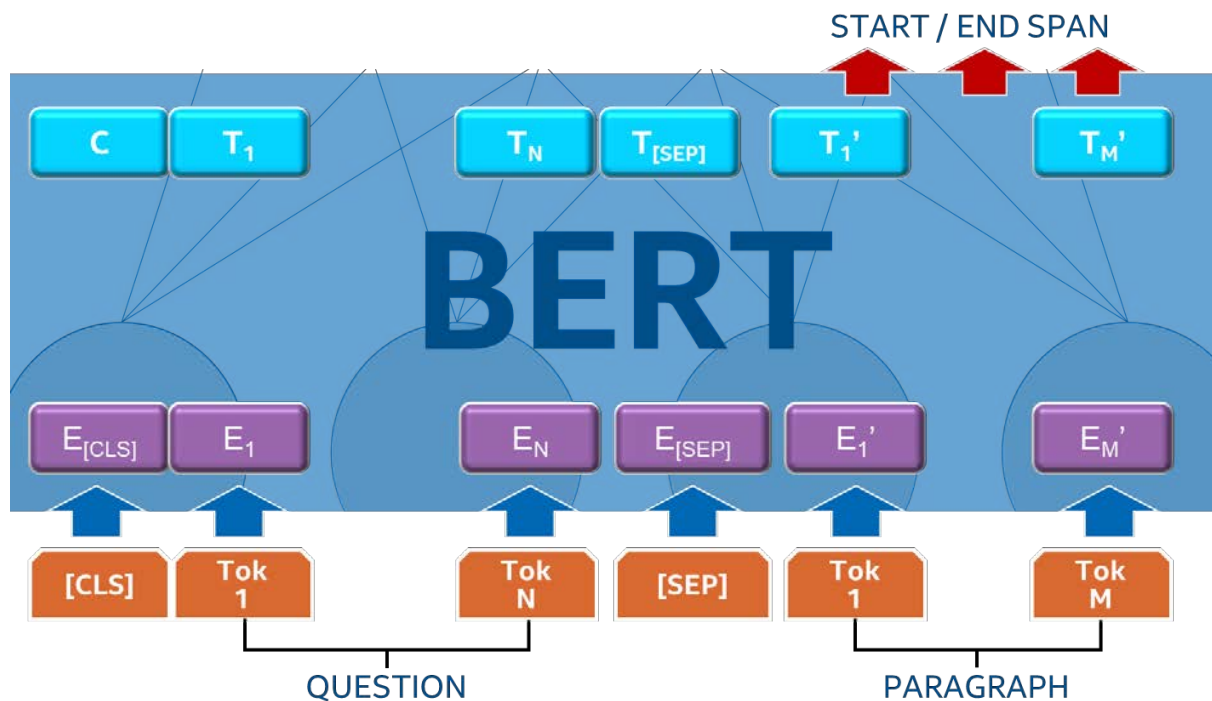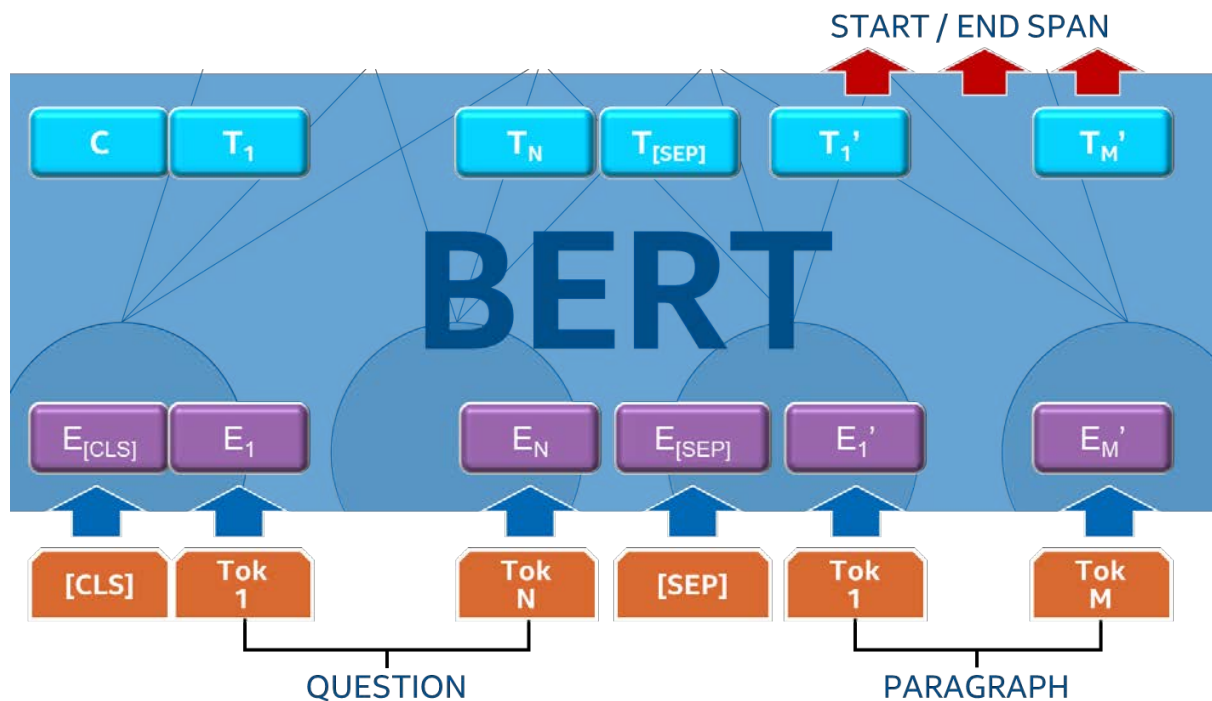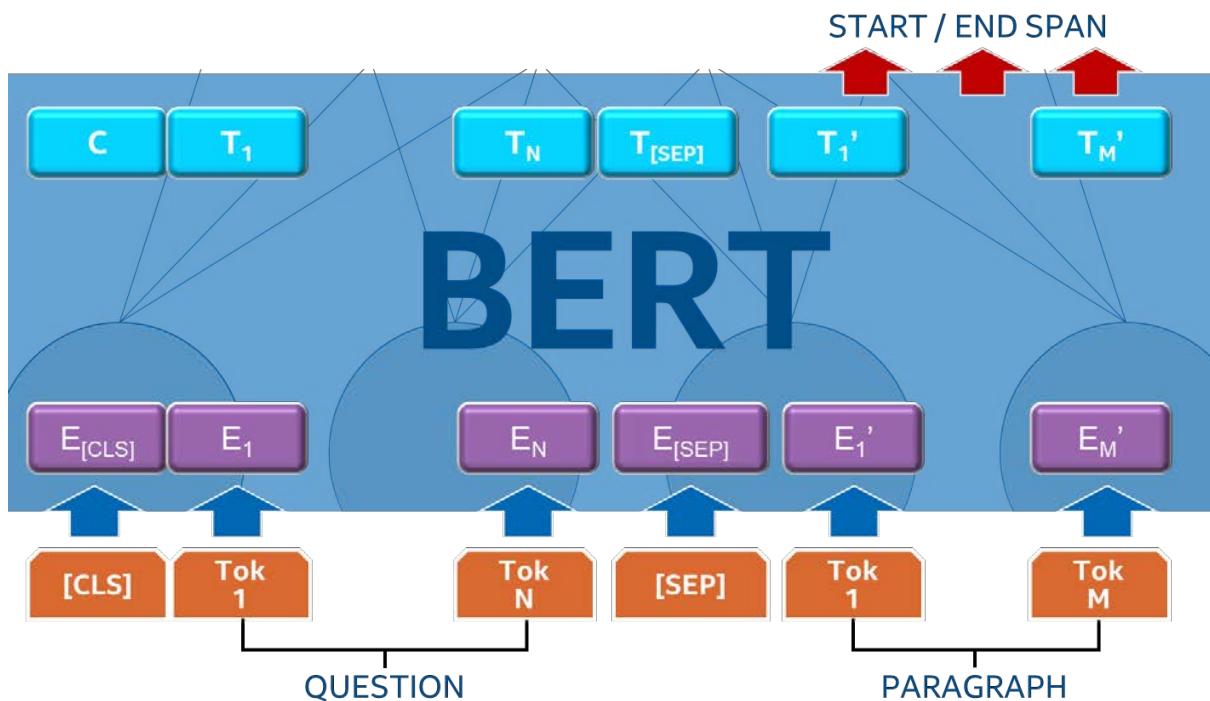
*Intel was founded in Mountain View, California, in 1968 by Gordon E. Moore (known for "Moore's law"), a chemist, and Robert Noyce, a physicist and co-inventor of the integrated circuit. Arthur Rock (investor and venture capitalist) helped them find investors, while Max Palevsky was on the board from an early stage.[23] Moore and Noyce had left Fairchild Semiconductor to found Intel. Rock was not an employee, but he was an investor and was chairman of the board.[24][25] The total initial investment in Intel was $2.5 million in convertible debentures (equivalent to $18.4 million in 2019) and $10,000 from Rock. Just 2 years later, Intel became a public company via an initial public offering (IPO), raising $6.8 million ($23.50 per share).[24] Intel's third employee was Andy Grove,[26] a chemical engineer, who later ran the company through much of the 1980s and the high-growth 1990s. …*

# 1. Build

TensorFlow  Caffe

KALDI  mxnet

ONNX

## Open Model Zoo

**100+** open sourced and optimized pre-trained models;

**80+** supported public models

# 2. Optimize



## Model Optimizer

Converts and optimizes trained model using a supported framework

# 3. Deploy



## Inference Engine

Common API that abstracts low-level programming for each hardware

# OpenVINO™ toolkit + BERT

- **Model calibration** is available via the Post-training Optimization Tool for TensorFlow

# OpenVINO™ toolkit + BERT

- Model calibration is available via the Post-training Optimization Tool for TensorFlow

- Model fine-tuning or re-training via PyTorch and HuggingFace recipe in **Neural Network Compression Framework**

# OpenVINO™ toolkit + BERT

- Model calibration is available via the Post-training Optimization Tool for TensorFlow

- Model fine-tuning or re-training via PyTorch and HuggingFace recipe in Neural Network Compression Framework

- Open-sourced full precision (FP32) and low precision (INT8) models and a demo

intel.

... FakeQuantize → MatMul → (Bias)Add → GELu → FakeQuantize ...

# VIDEO

Pad Thai Demo with Audio.mov

intel.

Can we do better?

# VIDEO

BERT Large (Conversational).mov

intel.

# Question Answering on SQuAD2.0

| RANK | MODEL | EM ↑ | F1 | PAPER | | CODE | RESULT | YEAR |
|------|-------|------|-----|-------|--|------|--------|------|
| 1 | **SA-Net on Albert** (ensemble) | 90.724 | 93.011 | | | | | 2020 |

# bert-large-uncased-whole-word-masking-squad-int8-0001



**Throughput** (higher is better) — Frames per Second (FPS)

**Latency** (lower is better) — Milliseconds

CPU INFERENCE ENGINES

| CPU | Throughput (FPS) | | Latency (ms) |
|---|---|---|---|
| Intel® Atom™ x5-E3940 | 0.12 | 0.21 | 5263.16 |
| Intel® Core™ i3-8100 | 1.26 | 1.9 | 537.63 |
| Intel® Core™ i5-8500 | 1.79 | 2.78 | 371.75 |
| Intel® Core™ i7-8700T | 2.03 | 3.17 | 346.02 |
| Intel® Core™ i7-10920X | 5.15 | 9.26 | 111.48 |
| Intel® Core™ i5-1145G7E CPU-only | 1.49 | 4.19 | 268.1 |
| Intel® Core™ i5-1145G7E GPU-only | 2.71 | 3.28 | 318.47 |
| Intel® Core™ i5-1145G7E GPU+CPU | 2.85 | 5.32 | 318.47 |
| Intel® Xeon® E-2124G | 1.49 | 2.22 | 460.83 |
| Intel® Xeon® Silver 4216R | 5.88 | 14.88 | 126.58 |
| Intel® Xeon® Gold 5218T | 6.09 | 15.43 | 120.63 |
| Intel® Xeon® Platinum 8270 | 14.8 | 29.7 | 106.04 |

# Get Started

Typical workflow from development to deployment

```
┌──────────────────┐      ┌──────────────────┐   ┌──────────────────┐   ┌──────────────────┐
│  Train a model   │      │   Run the Model  │   │   Intermediate   │   │  Deploy using the│
├──────────────────┤  ▶   │    Optimizer     │ ▶ │  Representation  │ ▶ │  Inference Engine│  ▶
│ Find a trained   │      │                  │   │                  │   │                  │
│     model        │      └──────────────────┘   │    .bin, .xml    │   └──────────────────┘
└──────────────────┘                             └──────────────────┘
```

intel XEON

intel CORE i7

intel ATOM

intel IRIS Pro GRAPHICS

intel MOViDIUS

intel ARRiA 10

# What's New in the 2021.1 Release

Typical workflow from development to deployment

Support for TensorFlow 2.x

Support for Tiger Lake (10th Generation Intel® Core® Processors)

Integration of DL Workbench and Intel® DevCloud for the Edge

Train a model

Find a trained model

Run the Model Optimizer

Intermediate Representation .bin, .xml

Deploy using the Inference Engine

New capabilities in OpenVINO™ Model Server

Introducing non-computer vision workloads

Support for GNA 2.0

# 35+
Open Source Deep Learning Demos

# VIDEO

Body Pose.mov

intel

# VIDEO

SF Inpainting.mov

# Resources and Community Support

Vibrant community of developers, enterprises and skills builders

QUALIFY > INSTALLATION > PREPARE > HANDS ON > SUPPORT

# Resources and Community Support

## Vibrant community of developers, enterprises and skills builders

| QUALIFY | INSTALLATION | PREPARE | HANDS ON | SUPPORT |
|---|---|---|---|---|
| ▪ Use a trained model and check if framework is supported<br><br>    *– or –*<br><br>▪ Take advantage of a pre-trained model from the Open Model Zoo | ▪ Download the Intel® OpenVINO™ toolkit package from Intel® Developer Zone, or by YUM or APT repositories<br><br>▪ Utilize the Getting Started Guide | ▪ Understand sample demos and tools included<br><br>▪ Understand performance<br><br>▪ Choose hardware option with Performance Benchmarks<br><br>▪ Build, test and remotely run workloads on the Intel® DevCloud for the Edge before buying hardware | ▪ Visualize metrics with the Deep Learning Workbench<br><br>▪ Utilize prebuilt, Reference Implementations to become familiar with capabilities<br><br>▪ Optimize workloads with these performance best practices<br><br>▪ Use the Deployment Manager to minimize deployment package | ▪ Ask questions and share information with others through the Community Forum<br><br>▪ Engage using #OpenVINO on Stack Overflow<br><br>▪ Visit documentation site for guides, how to's, and resources<br><br>▪ Attend training and get certified<br><br>▪ Ready to go to market? Tell us how we can help |

# Ready to get started?

Download directly from Intel for free:

https://software.intel.com/content/www/us/en/develop/tools/openvino-toolkit.html

# Ready to get started?

*Also available from* Intel's Edge Software Hub | Intel® DevCloud for the Edge | PIP | DockerHub | Dockerfile | Anaconda Cloud | YUM | APT

*Build from source:*

https://github.com/openvinotoolkit/openvino
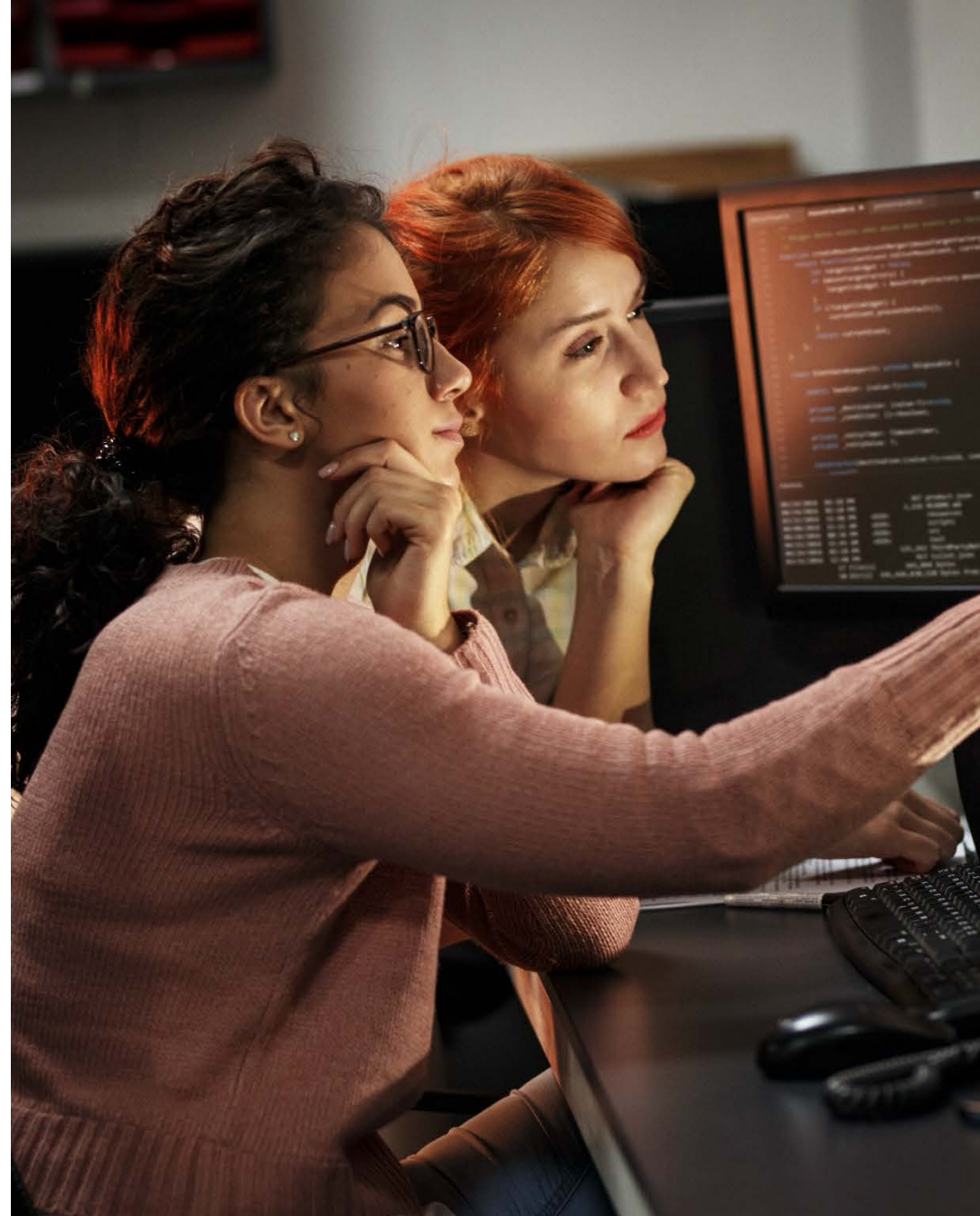
https://gitee.com/OpenVINO-Toolkit

intel.

# Notices and Disclaimers

- 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.

**Optimization Notice**
[1] 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 the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice. Notice revision #20110804.

[2] Software and workloads used in performance tests may have been optimized for performance only on microprocessors from Intel. 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. Consult other information and performance tests while evaluating potential purchases, including performance when combined with other products. For more information, see Performance Benchmark Test Disclosure. Source: Intel measurements, as of June 2017.