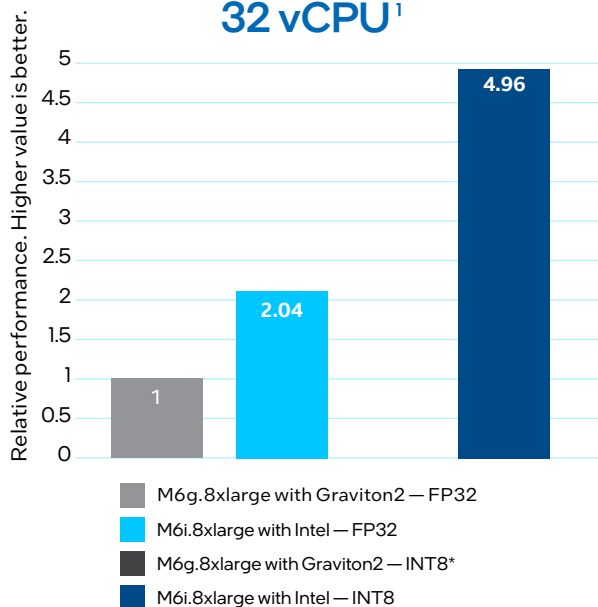# Faster Natural Language Processing with Intel

Natural language processing (NLP) delivers seamless integration with customers and propels your business forward. Take it to the next level on Intel.
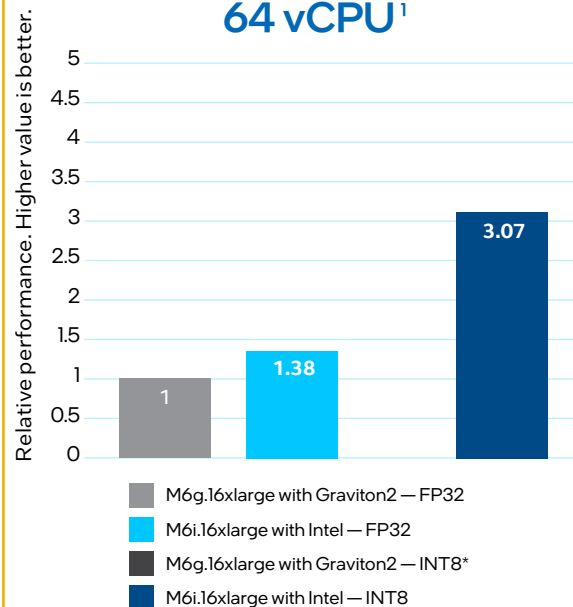
**intel** ®

## M6i with Intel Outperfoms M6g with Graviton2

### BERT-Large Inference Sentences Per Second 32 vCPU [1]

*Relative performance. Higher value is better.*



- 1 — M6g.8xlarge with Graviton2 — FP32
- 2.04 — M6i.8xlarge with Intel — FP32
- 4.96 — M6i.8xlarge with Intel — INT8

Legend:
- M6g.8xlarge with Graviton2 — FP32
- M6i.8xlarge with Intel — FP32
- M6g.8xlarge with Graviton2 — INT8*
- M6i.8xlarge with Intel — INT8

\* M6g INT8 support was not available when tested for Tensorflow in the BERT-large model used.

### BERT-Large Inference Sentences Per Second 64 vCPU [1]

*Relative performance. Higher value is better.*



- 1 — M6g.16xlarge with Graviton2 — FP32
- 1.38 — M6i.16xlarge with Intel — FP32
- 3.07 — M6i.16xlarge with Intel — INT8

Legend:
- M6g.16xlarge with Graviton2 — FP32
- M6i.16xlarge with Intel — FP32
- M6g.16xlarge with Graviton2 — INT8*
- M6i.16xlarge with Intel — INT8

\* M6g INT8 support was not available when tested for Tensorflow in the BERT-large model used.

## Solve Common Problems

- Analyze user-typed text more quickly
- Process results faster for predictive text
- Deliver meaningful automated responses to human users
- Remove communication barriers

## Ideal Uses

- Smart assistants
- Search
- Language translation
- Sentiment analysis

## Demanding Workloads

Processing input from customers and other users puts a heavy demand on compute resources. Deliver results fast by landing NLP applications on cloud instances with Intel processors.

# Cloud Performance Advantages

**intel**®

---

Up to **3.97x**

better price performance on NLP in M6i with INT8 precision than in M6g with FP32 precision (32vCPU)[2]

Up to **2.46x**

better price performance on NLP in M6i with INT8 precision than in M6g with FP32 precision (64vCPU)[3]

---

## Intentional Workload Placement

Choose from **AWS EC2 M6i, C6i instances,** and **DL1 instances** for your NLP workloads. Discover the benefits **Intel leadership in AI technology** brings to your business.

## Development Tools & Resources

Prepare, build, deploy, and scale AI solutions for your workloads. Visit **developer.intel.com/ai** for help getting the most out of your cloud.

**intel XEON**®

**Let's work together. Contact us.**

---