# GDDR Memory Enabling AI and High Performance Compute

## Wolfgang Spirkl, Fellow at Micron Technology

GTC, S9968, 20-March-2019

# Agenda

- The Demand for faster Memory and storage

- Competing Compute/Memory Solutions

- GDDR6 for AI applications and more

- Micron GDDR6 AI demonstration

Micron®

# Accelerated Data Cycle

## Driven by Increasing Data Value

- Creates continuous need to capture, process, move & store data

- Generates ever-increasing demand for memory & fast storage
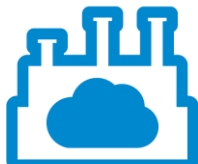
- AI is amplifying the Accelerated Cycle

Micron

# AI Landscape for Memory & Storage

## Data Center

### System/Infrastructure:

Cloud

On Prem.

### Workloads:

Model Development

Batch Training
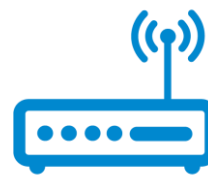
Batch Inference

### Memory/Storage:

| | |
|---|---|
| HBM | P-Mem |
| GDDR6 | 3D TLC |
| DDR4 | 3D QLC |

PERFORMANCE

## Smart Edge/Intelligent Endpoint

### System/Infrastructure:

Edge Compute

Smart Access Point

Autonomous Vehicle/Robot

Mobile Device

Smart IoT/ Sensor

### Workloads:

Online Training

On-demand Inference

### Memory/Storage:

| | |
|---|---|
| HBM | P-Mem |
| GDDR6 | 3D TLC |
| DDR4 | 3D QLC |
| LP4X | |

### Workloads:

Local Inference

01000110
01100001
01110
01110

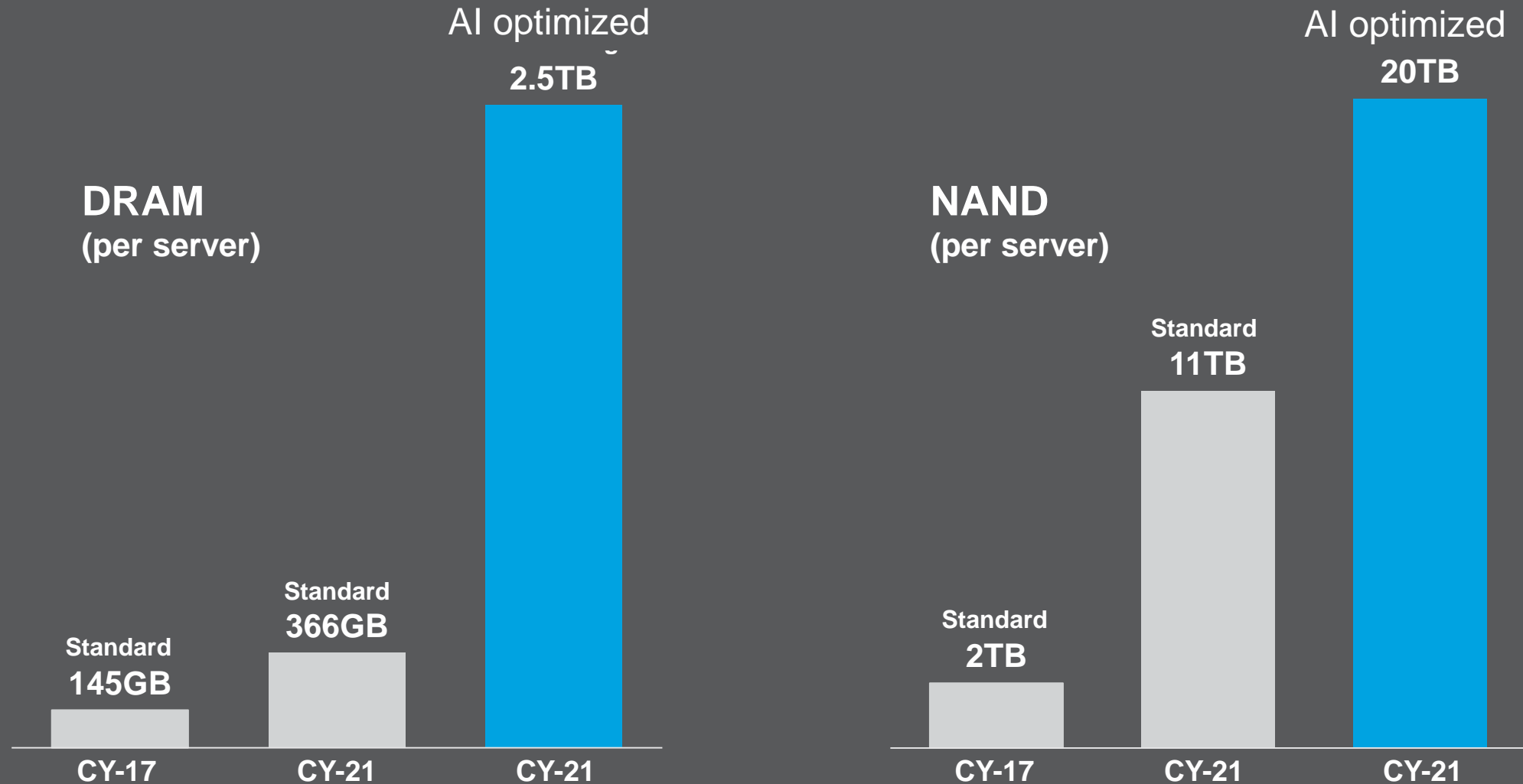Data Collection

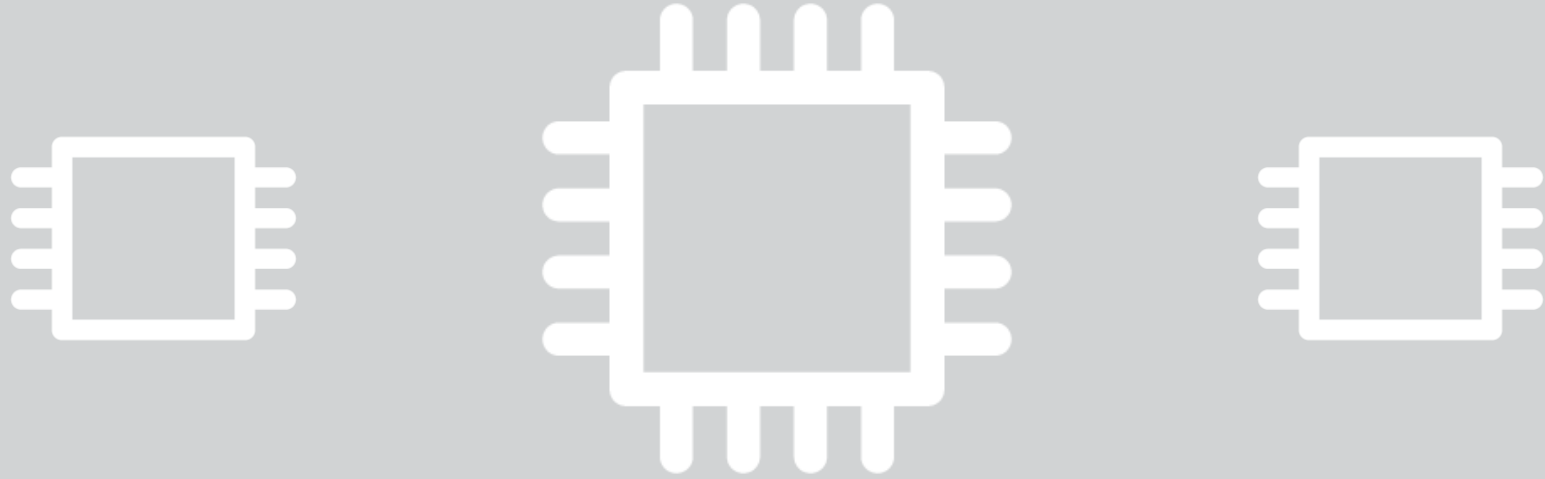### Memory/Storage:

| | |
|---|---|
| | 3D TLC |
| DDR4 | |
| LP4X | |

COST/POWER

Micron

# AI Workloads Unleash the Need For More Memory & Storage

Significant Growth Across Private, Public & Hybrid Cloud

**AI optimized**
**2.5TB**

**DRAM**
**(per server)**

**AI optimized**
**20TB**

**NAND**
**(per server)**

**Standard**
**11TB**

**Standard**
**366GB**

**Standard**
**145GB**

**Standard**
**2TB**

| CY-17 | CY-21 | CY-21 |
|-------|-------|-------|

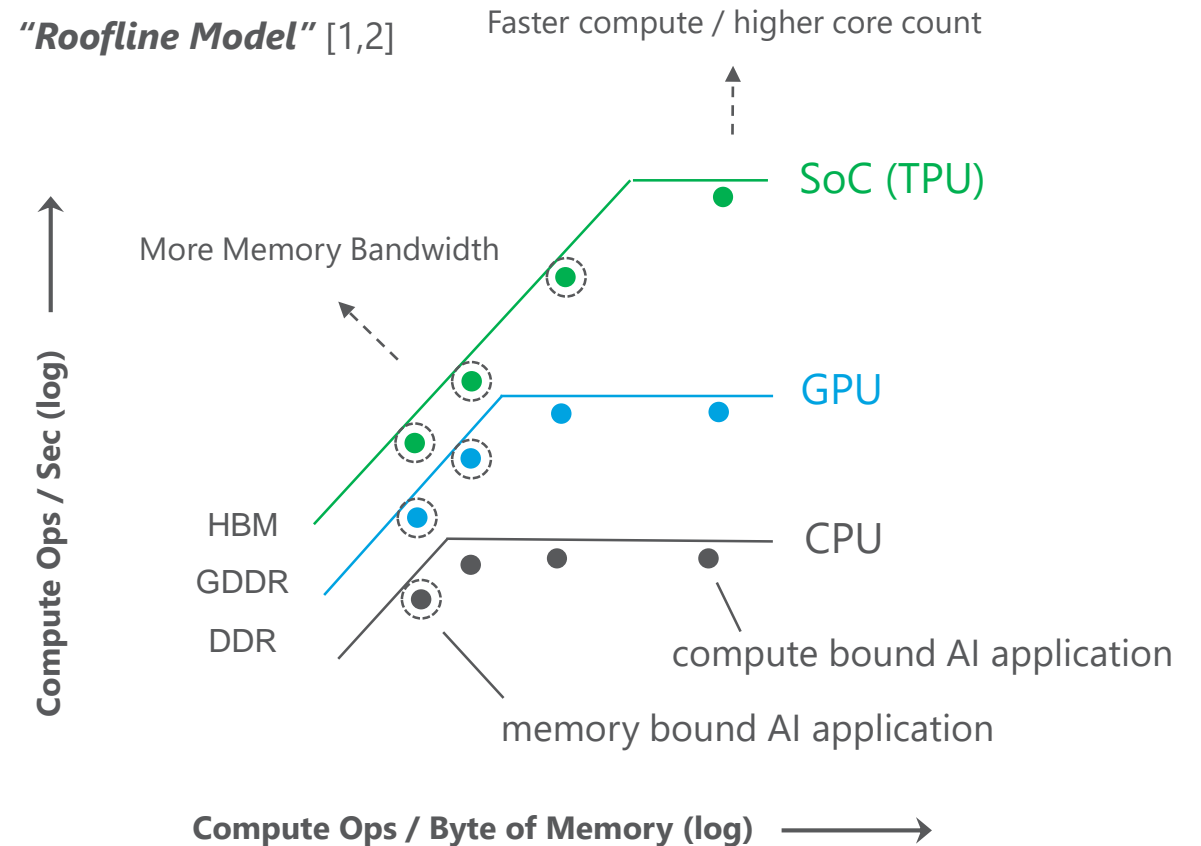| CY-17 | CY-21 | CY-21 |
|-------|-------|-------|

Source: Micron

Micron

- The Demand for faster Memory and storage

- **Competing Compute/Memory Solutions**

- GDDR6 for AI applications and more

- Micron GDDR6 AI demonstration

Micron

# AI Acceleration is Driving Demand for Memory Bandwidth

- AI accelerators increase compute performance
  - GPU, TPU, etc..

- Accelerated applications are more likely to be memory bound [3]

- Micron supports next gen technologies
  - GDDR6, HBM2E

*"Roofline Model"* [1,2]

Faster compute / higher core count

More Memory Bandwidth

SoC (TPU)

GPU

CPU

HBM
GDDR
DDR

Compute Ops / Sec (log)

compute bound AI application

memory bound AI application
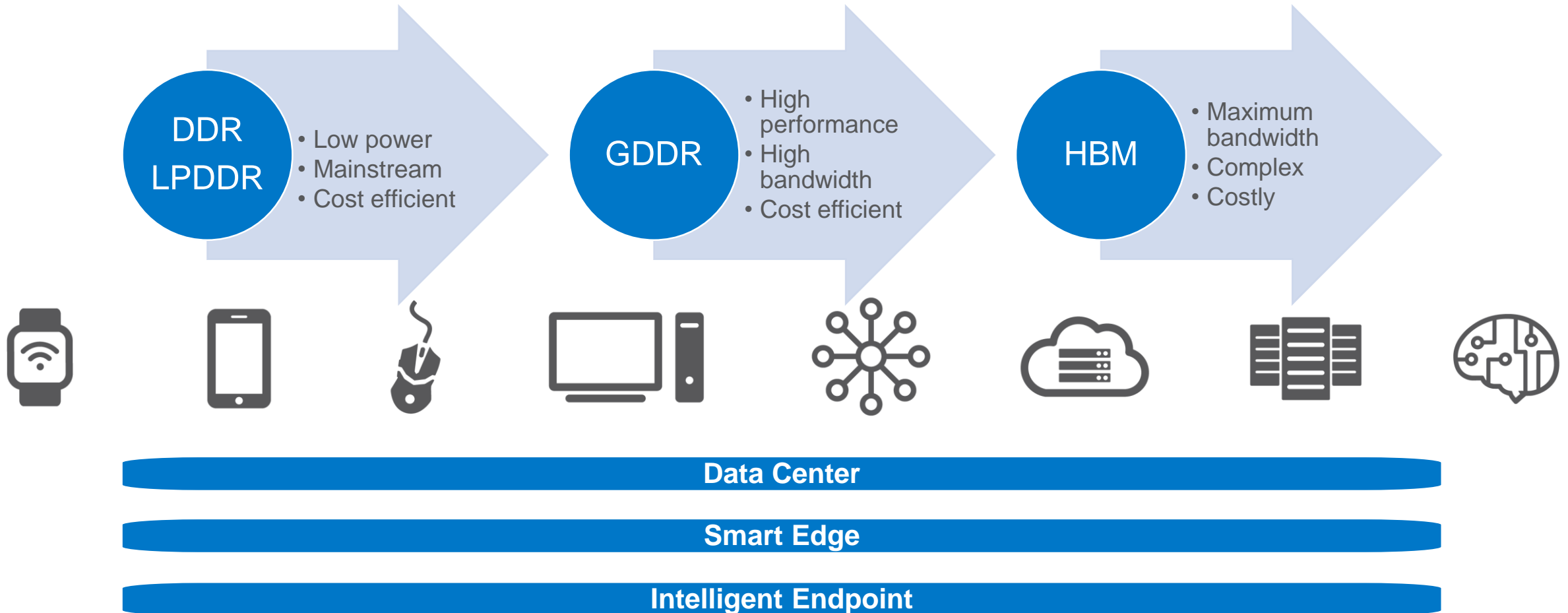
Compute Ops / Byte of Memory (log)

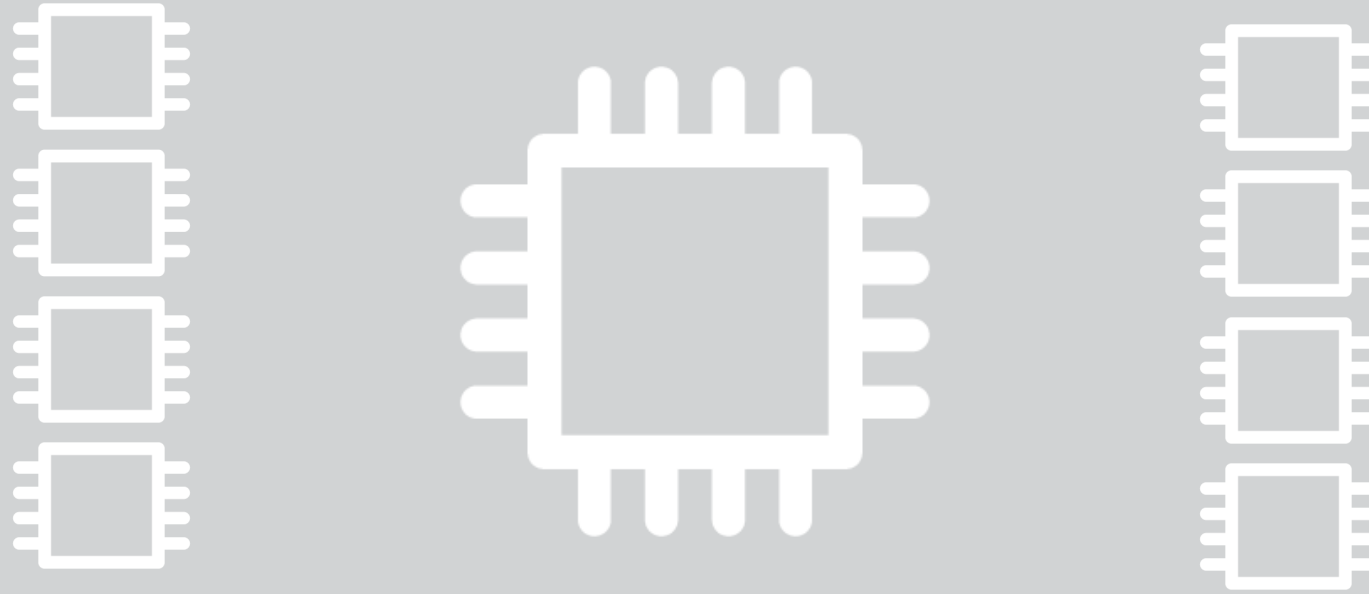[1] Jouppi, Norman, et al, 2017. In-Datacenter Performance Analysis of a TPU, ISCA
[2] Williams, S., Waterman, A. and Patterson, D., 2009. Roofline: an insightful visual performance model for multicore architectures. Communications of the ACM.
(3) Forrester report on memory and storage impact on AI

Micron

# Memory Options

**DDR LPDDR**
- Low power
- Mainstream
- Cost efficient

**GDDR**
- High performance
- High bandwidth
- Cost efficient

**HBM**
- Maximum bandwidth
- Complex
- Costly
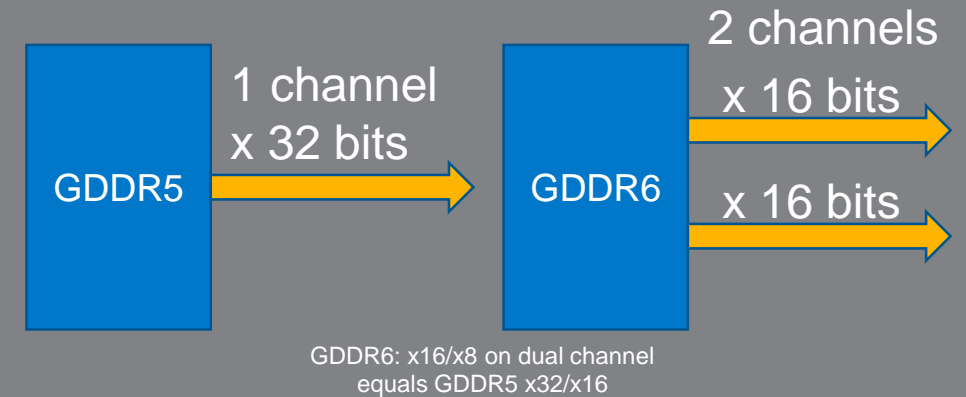
Data Center

Smart Edge

Intelligent Endpoint

Micron

The Demand for faster Memory and storage

Competing Compute/Memory Solutions

**GDDR6 for AI applications and more**

Micron GDDR6 AI demonstration

GTC 2019, Micron GDDR6

Micron®

# GDDR5/GDDR6 Features

| Feature | GDDR5 | GDDR6 |
|---|---|---|
| Density | 512Mb – 8Gb | 8Gb – 32Gb |
| VDD, VDDQ | 1.5V + 1.35V | 1.35V |
| VPP | N/A | 1.8V |
| Package | BGA-170 14mm x 12mm 0.8mm ball pitch | BGA-180 14mm x 12mm 0.75mm ball pitch |
| Signaling | POD15 / POD135 | POD135 |
| Data rate | ≤8 Gbps | ≤16 Gbps |
| I/O Width | x32/x16 | 2-ch x16/x8 |
| Access Granularity | 32B | 2-ch 32B each   or 1-ch 64B w/ PC mode |
| I/O Count | 61 | 62 / 74 |
| ABI, DBI | ✓ | ✓ |
| CRC | CRC-8 (BL8) | 2x CRC-8 (BL16); compressed 2x CRC-8 (BL8) |
| RDQS Mode | ✓ (BL8) | ✓ (BL16) |
| ODT | ✓ | ✓ |
| $V_{REFC}$ | external | external / internal |
| $V_{REFD}$ | ext. / int. | internal |
| Temp Sensor | ✓ | ✓ |

## Package Configuration



GDDR5  →  1 channel x 32 bits  →  GDDR6  →  2 channels x 16 bits / x 16 bits
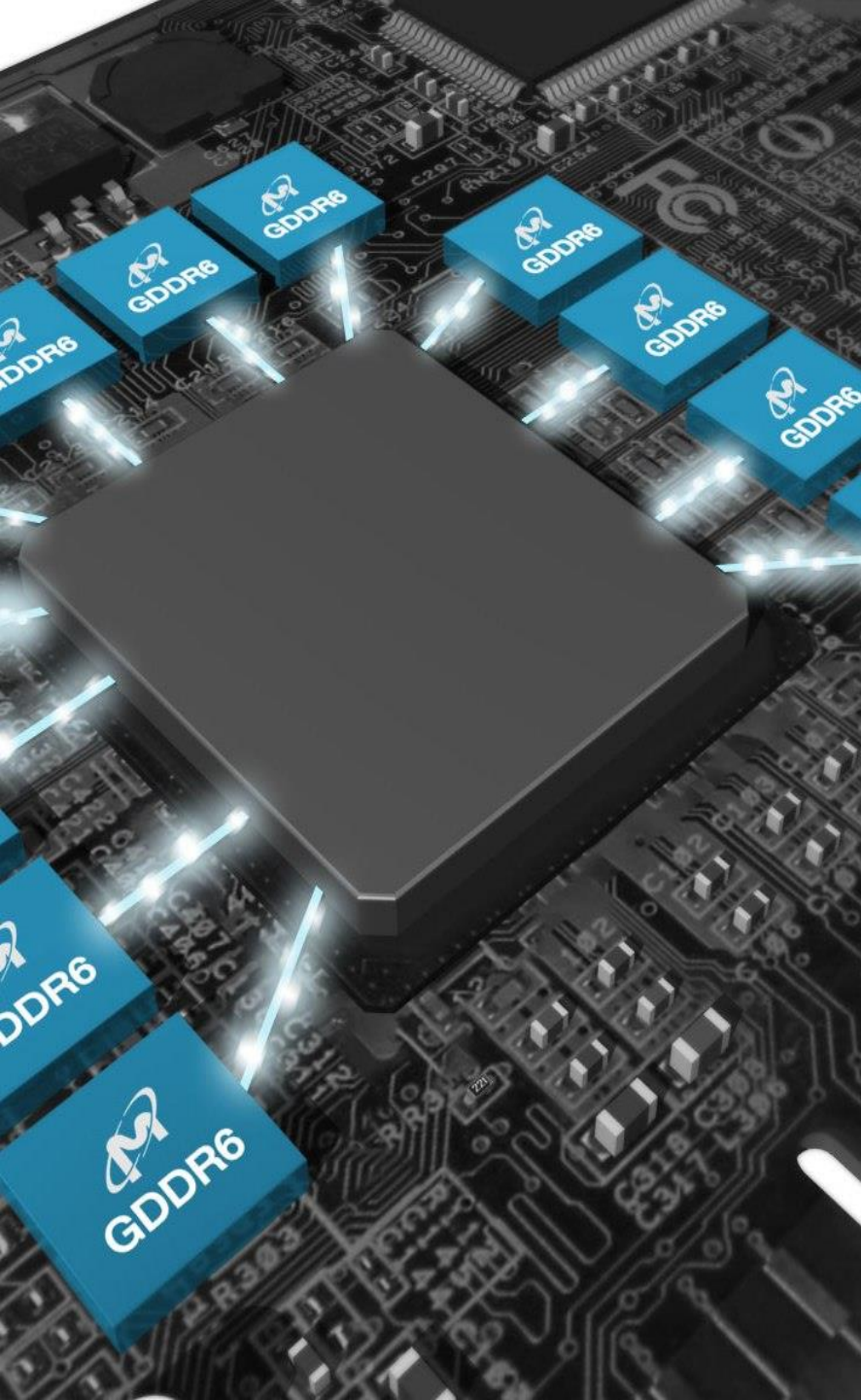
GDDR6: x16/x8 on dual channel equals GDDR5 x32/x16

Dual channel organization
- Maintains fine granularity (32 bytes per colum access)
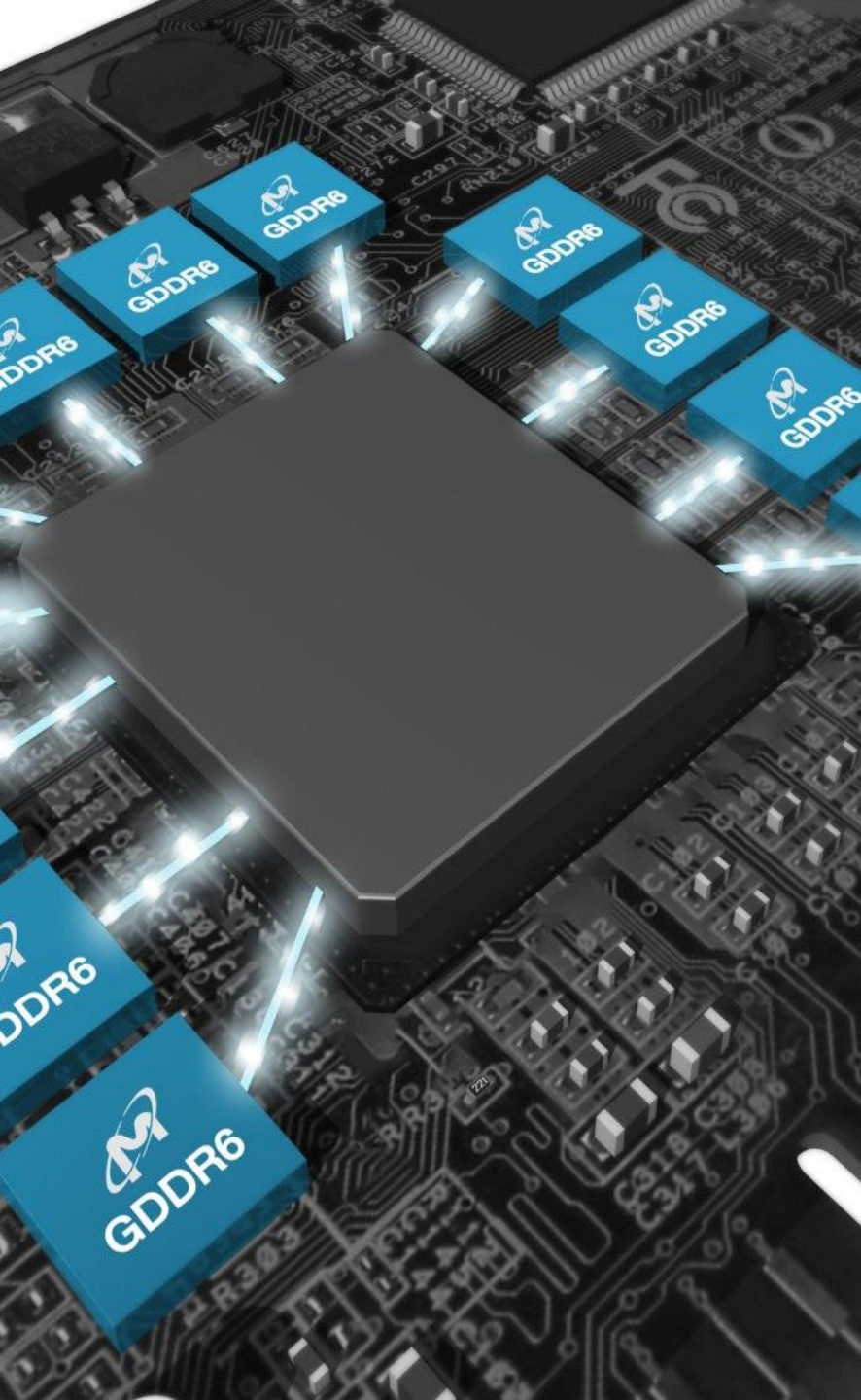- In spite of doubled prefetch size

New features for high data rates:
- Optimized signal ball-out for low-effort PCB design
- Per lane DFE and VREFD
- Transmitter equalization
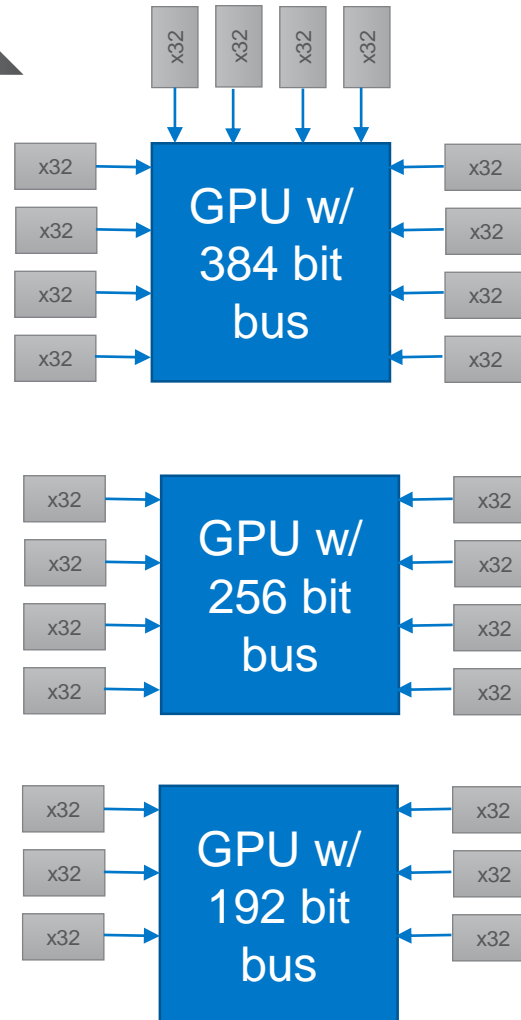
GTC 2019, Micron GDDR6

Micron

# Calculating GDDR Bandwidth

- Bandwidth =
  number of bits/s between GPU and memory

- Memory bus is like traffic lanes
  - More lanes, the greater the flow
  - Higher lane speed, the greater the flow

| Memory Bandwidth is | GDDR6 Example |
|---|---|
| number of memory components<br>X  number of lanes per component<br>X  Data rate per lane (Gbps) | 8<br>32<br>16 |
| Memory Bandwidth (GB/s) | 512 |

# GDDR Bandwidth / Memory Bus

GPU w/ 384 bit bus

GPU w/ 256 bit bus

GPU w/ 192 bit bus

| Technology | Speed (Gbps) | # of comp. | # of lanes | Memory bus (bit) | Bandwidth (GB/s) |
|---|---|---|---|---|---|
| HBM2 | 2 | 4 | 1024 | 4096 | 1024 |
| GDDR6 | 16 | 12 | 32 | 384 | 768 |
| GDDR6 | 14 | 12 | 32 | 384 | 672 |
| GDDR5X | 11 | 12 | 32 | 384 | 528 |
| GDDR5 | 7 | 12 | 32 | 384 | 336 |
| GDDR6 | 14 | 8 | 32 | 256 | 448 |
| GDDR5X | 11 | 8 | 32 | 256 | 352 |
| GDDR5 | 7 | 8 | 32 | 256 | 224 |

GTC 2019, Micron GDDR6

Micron®

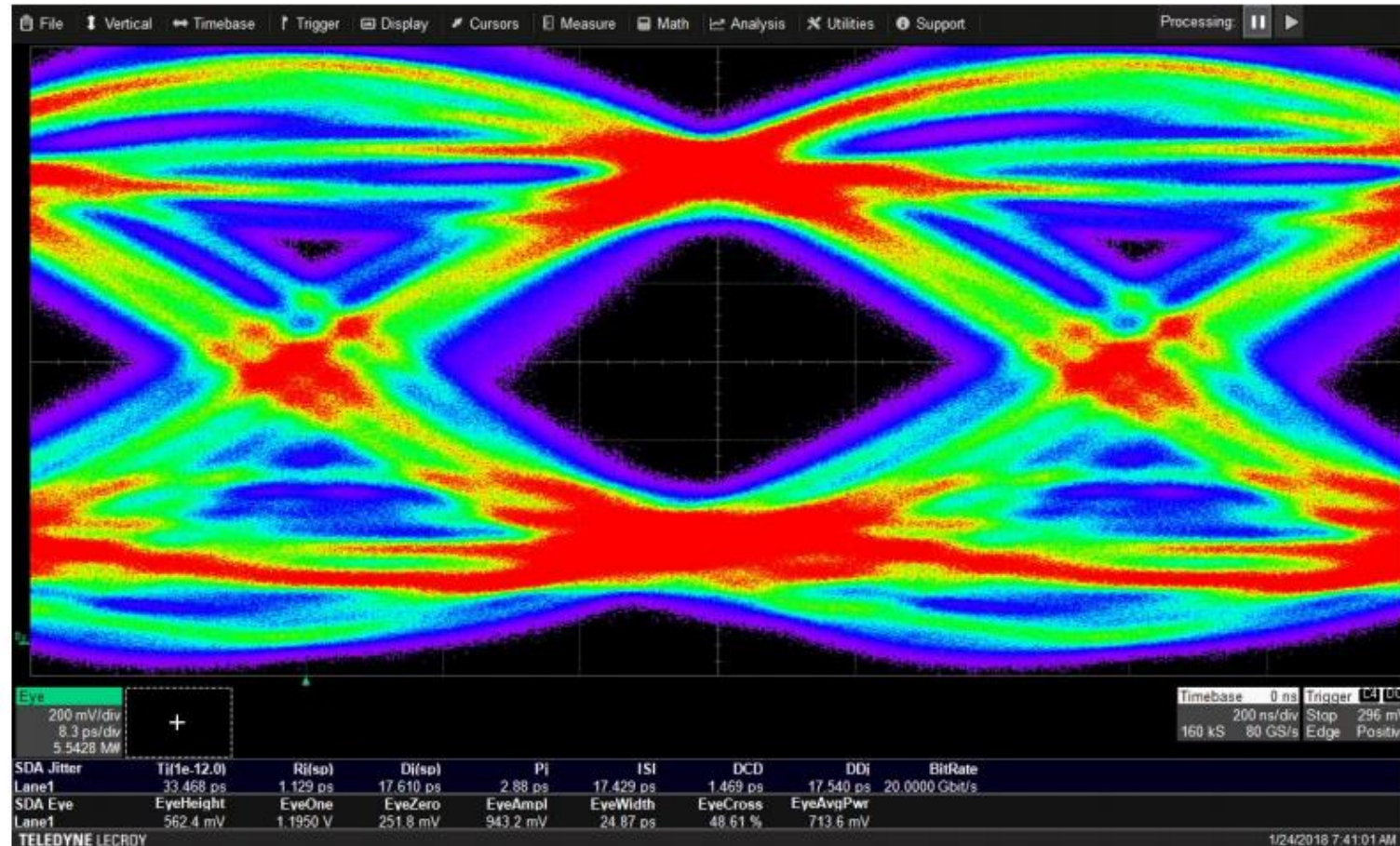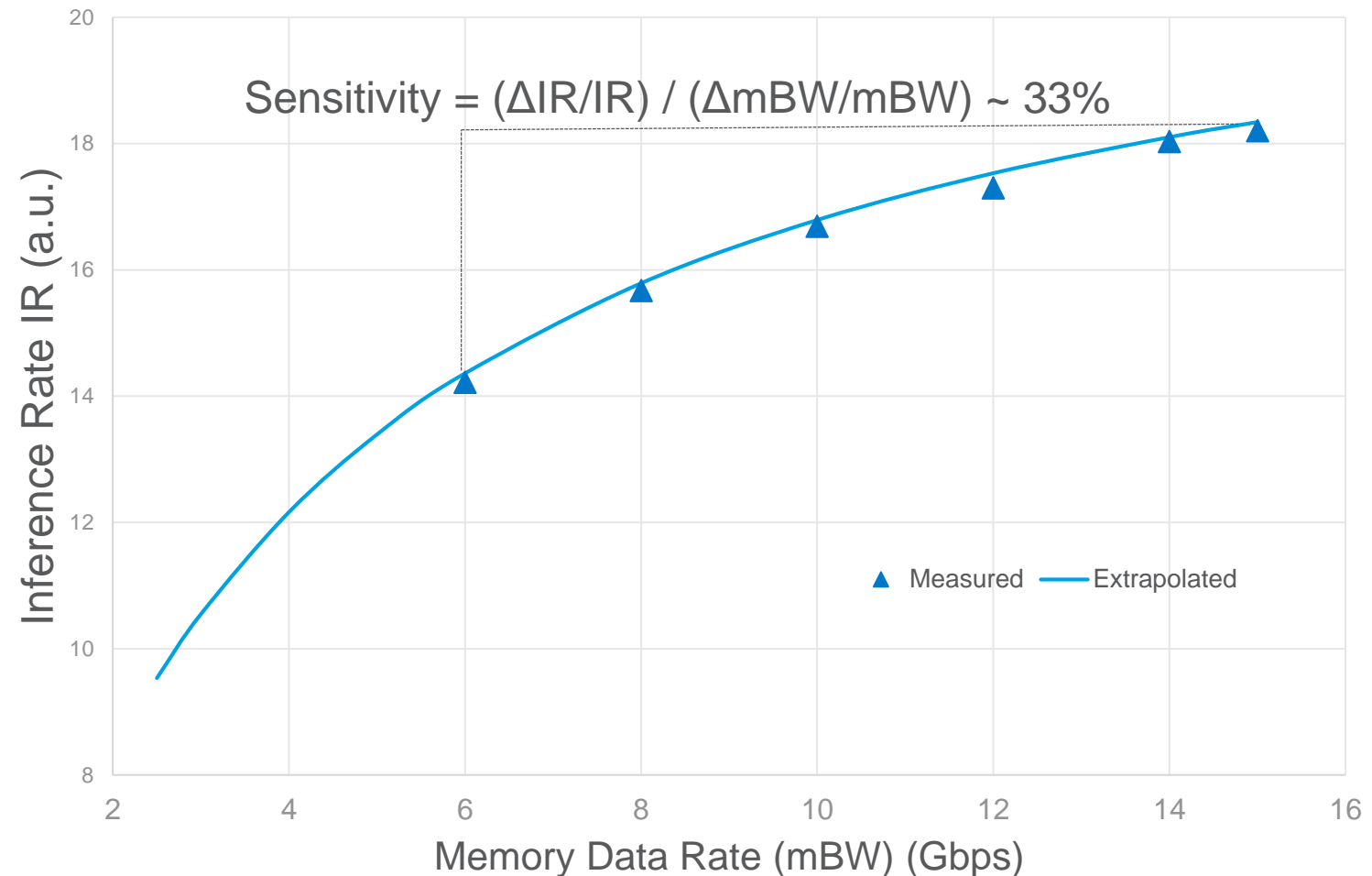# GDDR6 20 Gbps Data Eye
Measured performance beyond the specification



Figure 15: Measured 20Gb/s data eye based on a PRBS6 pattern

https://www.micron.com/-/media/client/global/documents/products/white-paper/16gb_s_and_beyond_w_single_endedio_in_graphics_memory.pdf?la=en

GTC 2019, Micron GDDR6

**Micron**

# Speech Recognition Craves Memory Bandwidth

- "Deep Speech" recognition application
  - Baidu Research's AI algorithm (https://arxiv.org/pdf/1412.5567.pdf)
  - Mozilla's tensorflow implementation
  - Speech-to-text benchmark for AI hardware (https://github.com/mozilla/DeepSpeech)
- Hardware
  - NVIDIA RTX 2080 Ti
  - 11GB GDDR6
    - 384 bit bus @14Gb/s/pin, 672GB/s
- Experiment setup
  - Adjust GDDR6 clock rate
  - Measure speech recognition inference rate:
    - $Inference\ rate = \dfrac{Audio\ file\ duration}{Inference\ time}$

Sensitivity = (ΔIR/IR) / (ΔmBW/mBW) ~ 33%

*X-axis:* Memory Data Rate (mBW) (Gbps)
*Y-axis:* Inference Rate IR (a.u.)

▲ Measured   — Extrapolated

GTC 2019, Micron GDDR6

Micron®

# AI Demonstrates the Need for Memory speed

SPEECH RECOGNITION DEMO – Micron Booth # 1713

# Conclusions

- AI Landscape demands higher performance memory to feed the compute needs

- Micron delivers a broad range of memory solutions for AI applications from data center to cloud to edge to endpoint devices

- GDDR6 high performance memory optimized for applications beyond graphics

Experience Micron speech recognition AI with GDDR6 in our booth 1713!

Micron