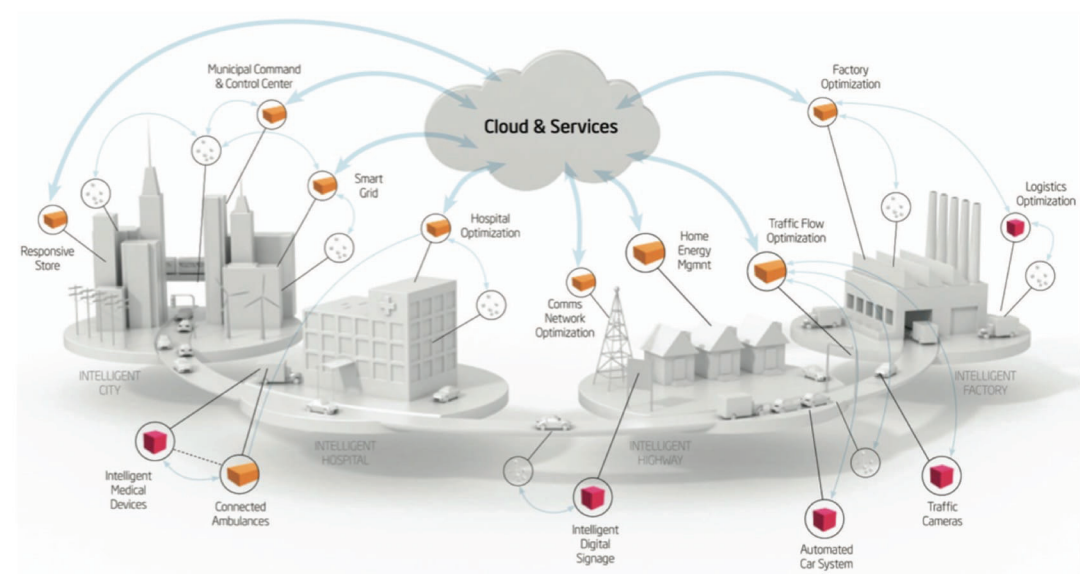


SESSION 39.3

Memory Systems for AI and Leading-edge Applications

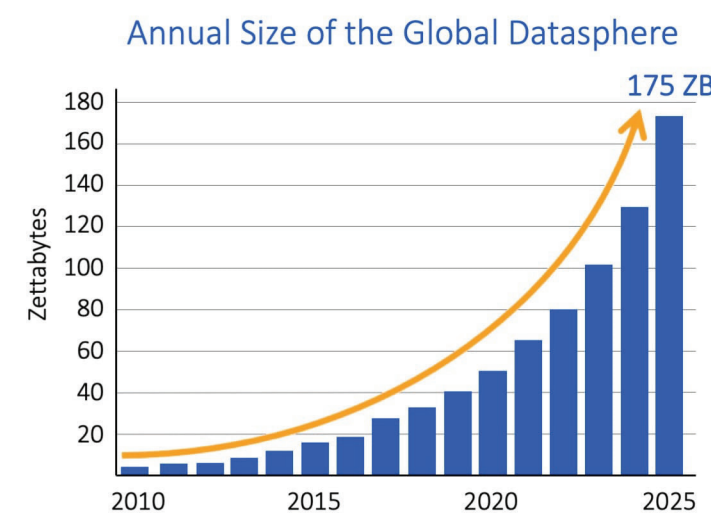
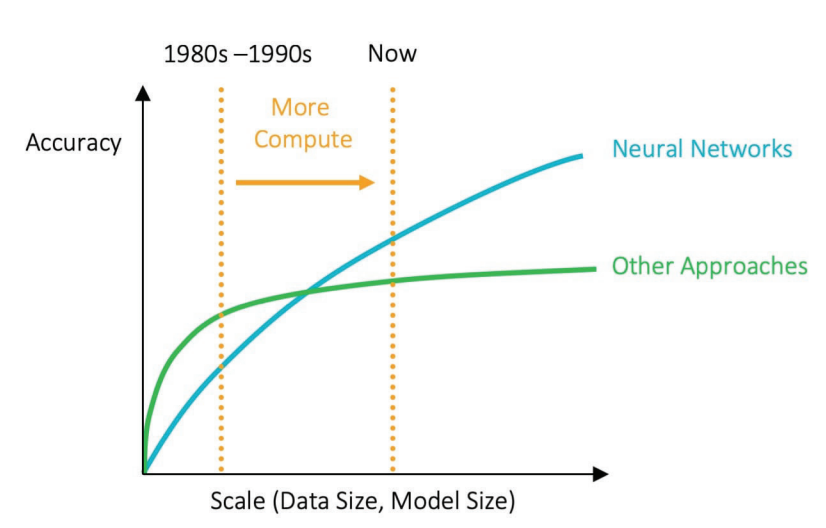
Pervasive Connectivity and the Internet of Things

- Memory, link, and storage performance must continue to increase
- Data and insights are increasingly more valuable, security a growing concern
- Compute and I/O power efficiency must also continue to improve



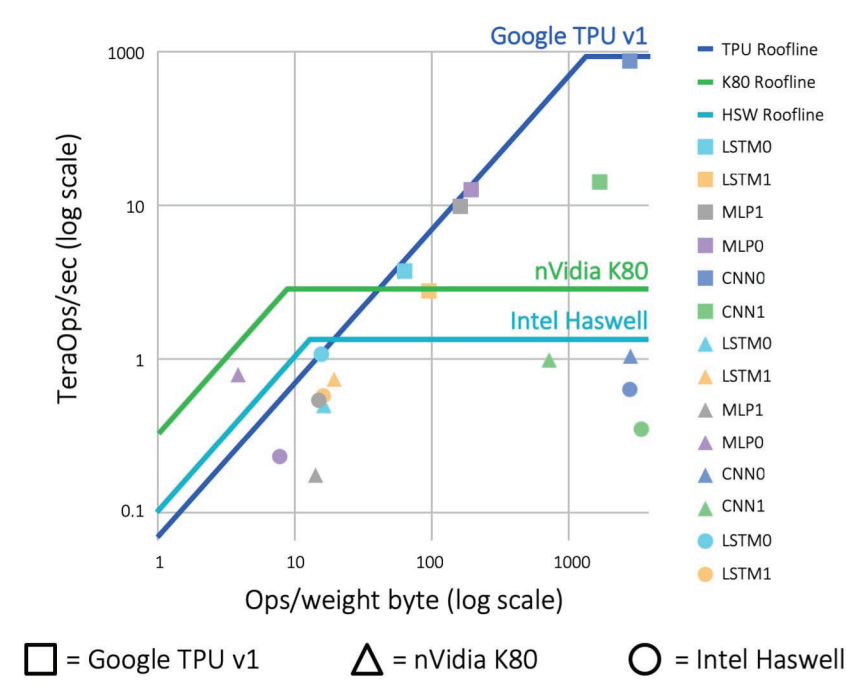
In our increasingly connected world, architectures are evolving to more efficiently capture, secure, move, and process the growing volume of digital data

Faster Compute + Big Data Enabling Explosive Growth in AI



- Faster compute and memory + large training sets have enabled modern AI
- Much more left to do, more performance needed
- Key challenges: Moore's Law ending, energy efficiency growing in importance

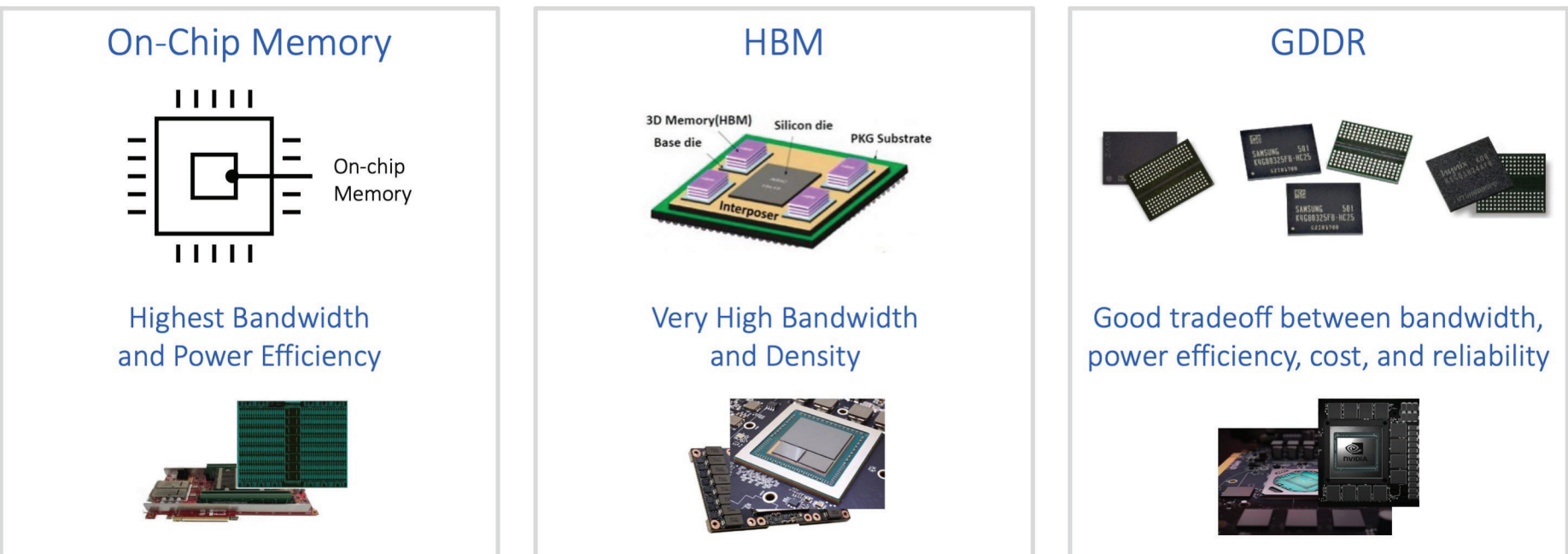
AI Accelerators Need Memory Bandwidth



- Inference on older, general purpose hardware (Haswell, K80) performs well, applications can benefit from compute and memory optimizations
- Inference on AI-specific silicon (Google TPU v1) largely limited by memory bandwidth

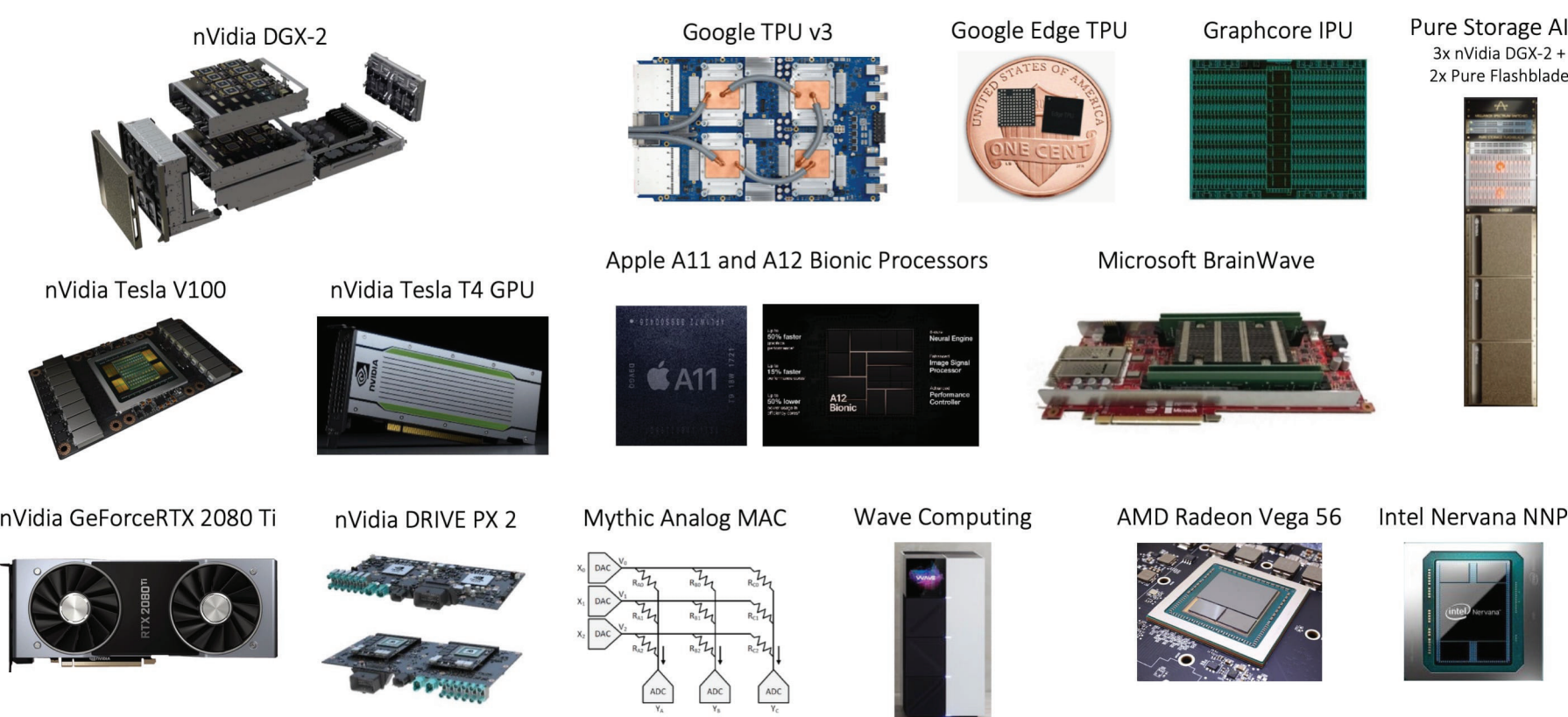
Memory bandwidth is a critical resource for AI applications

AI Needs Memory Bandwidth: Common AI Memory Systems



Multiple options suited to different needs

Example AI Hardware on the Market



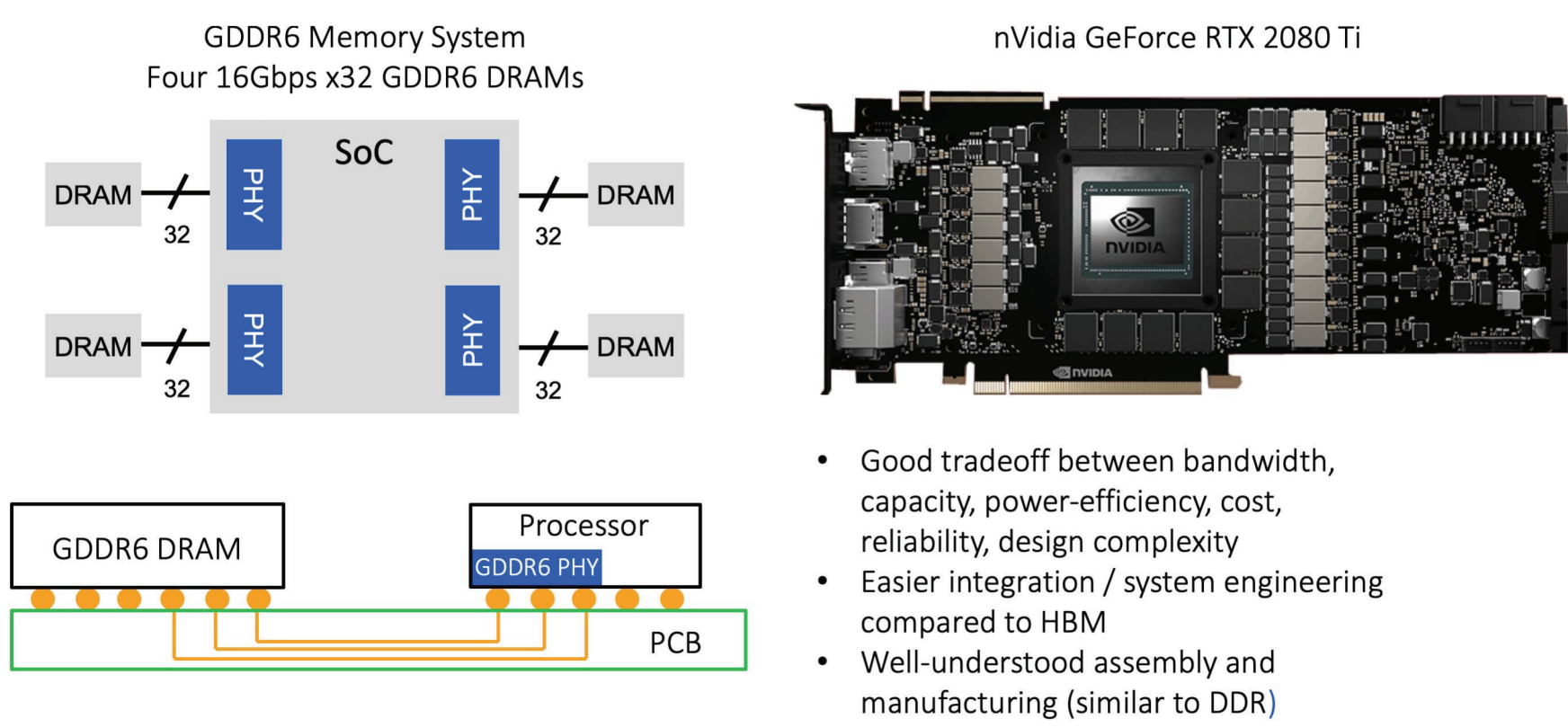
Memory System Comparison: 256GB/s GDDR6 vs. HBM2

	GDDR6 Memory System Four 16Gbps x32 GDDR6 DRAMs	HBM2 Memory System Single 2Gbps HBM2 Device
Total Bandwidth	256 GB/s	256 GB/s
Per-pin data rate	16 Gbps	2 Gbps
Relative Controller PHY Area ^[1]	1.5-1.75	1.0
Relative Controller PHY Power ^[1]	3.5-4.5	1.0
Interposer	None	Added cost ^[2]
Memory	Similar to GDDR5, DDR4	Stacked, adds cost ^[2]
		Cost advantage for GDDR6

[1] Source: Rambus Inc.
[2] Source: The Cost of HBM2 vs. GDDR6 & Why AMD Had to Use It, <https://www.gamernexus.net/guides/2022-vega-56-cost-of-hbm2-and-necessity-to-use-it>

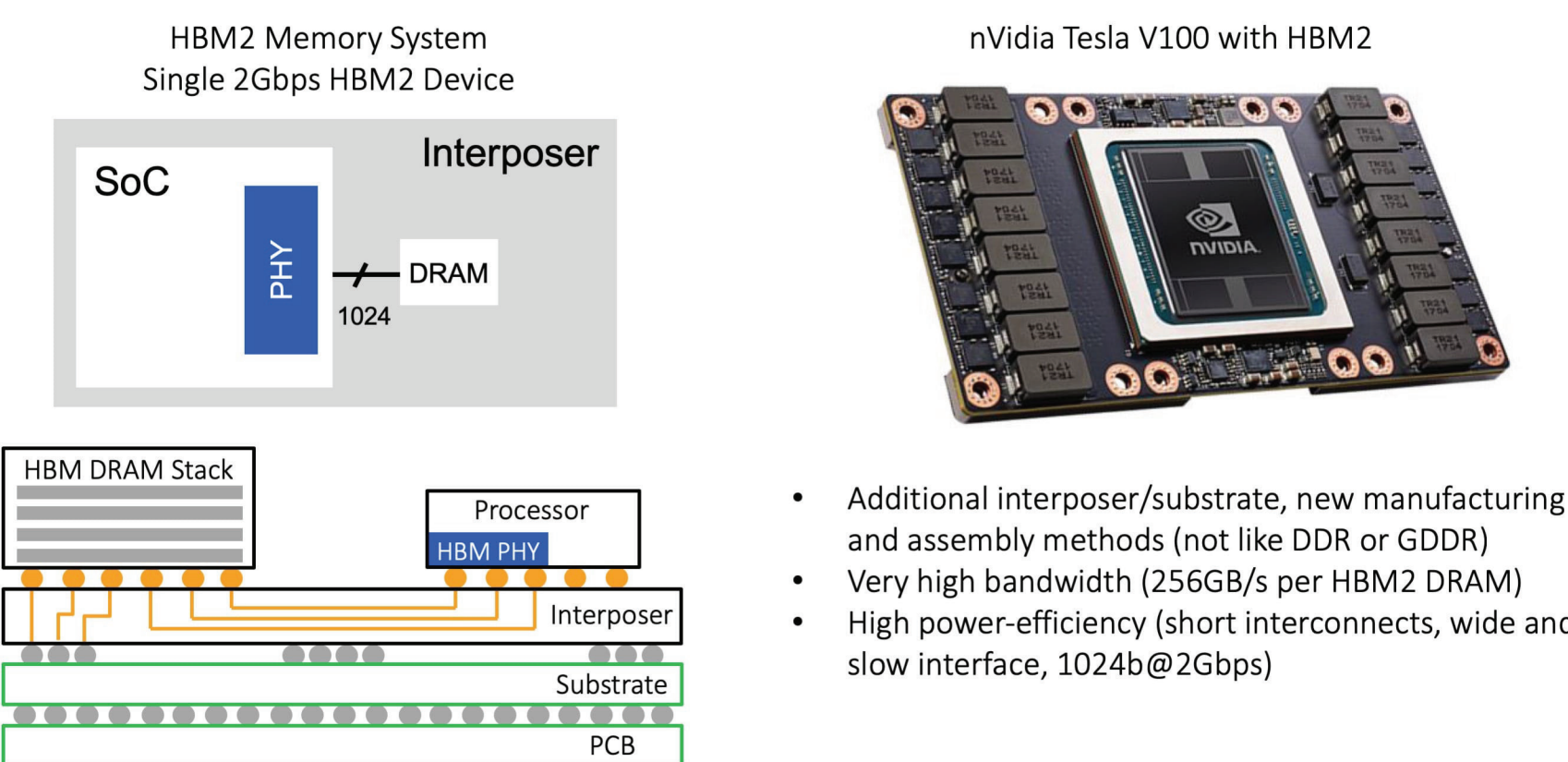
GDDR6 and HBM2 offer different system design tradeoffs

256GB/s GDDR6 Memory System



- Good tradeoff between bandwidth, capacity, power-efficiency, cost, reliability, design complexity
- Easier integration / system engineering compared to HBM
- Well-understood assembly and manufacturing (similar to DDR)

256GB/s HBM2 Memory System



- Additional interposer/substrate, new manufacturing and assembly methods (not like DDR or GDDR)
- Very high bandwidth (256GB/s per HBM2 DRAM)
- High power-efficiency (short interconnects, wide and slow interface, 1024b@2Gbps)

Summary/Conclusion

- AI driving the development of new silicon and new system architectures
- Memory bandwidth a critical resource for AI applications, memory systems are once again a hot topic in the semiconductor industry
- Multiple memory options to suit different AI application needs
 - On-chip memory: Highest bandwidth and power efficiency, lowest latency, but storage capacity limited
 - HBM: Extremely high bandwidth and power efficiency, but higher cost and more challenging integration and design complexity
 - GDDR: Good tradeoff between bandwidth, capacity, power efficiency, cost, reliability, design complexity

Frank Ferro

Sr. Director, Product Management, Rambus, Inc.

