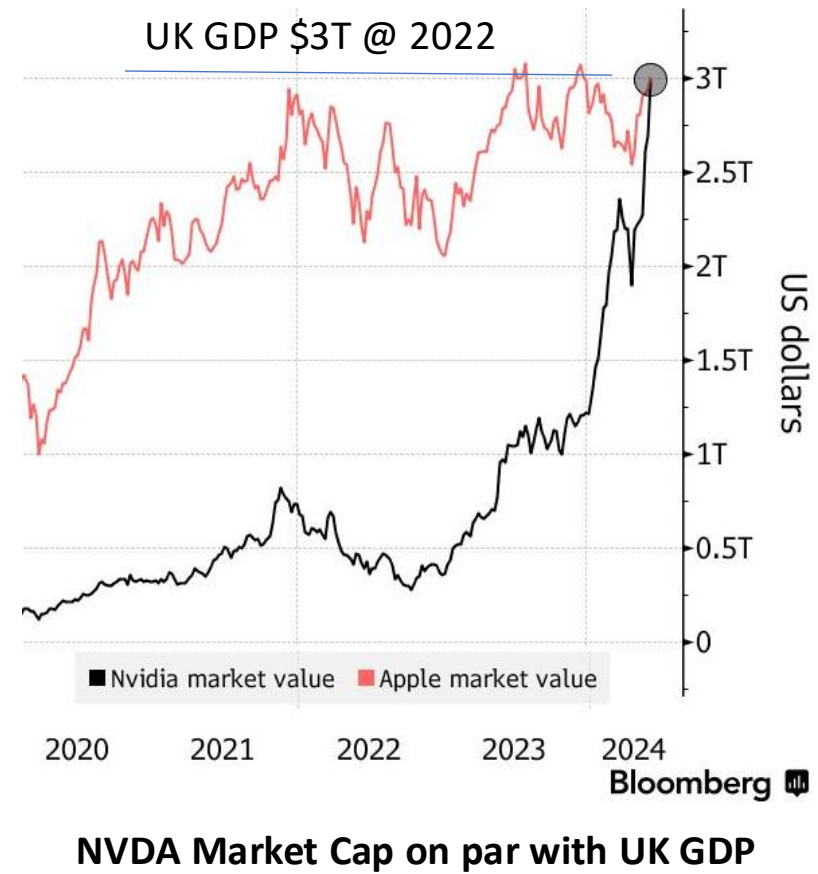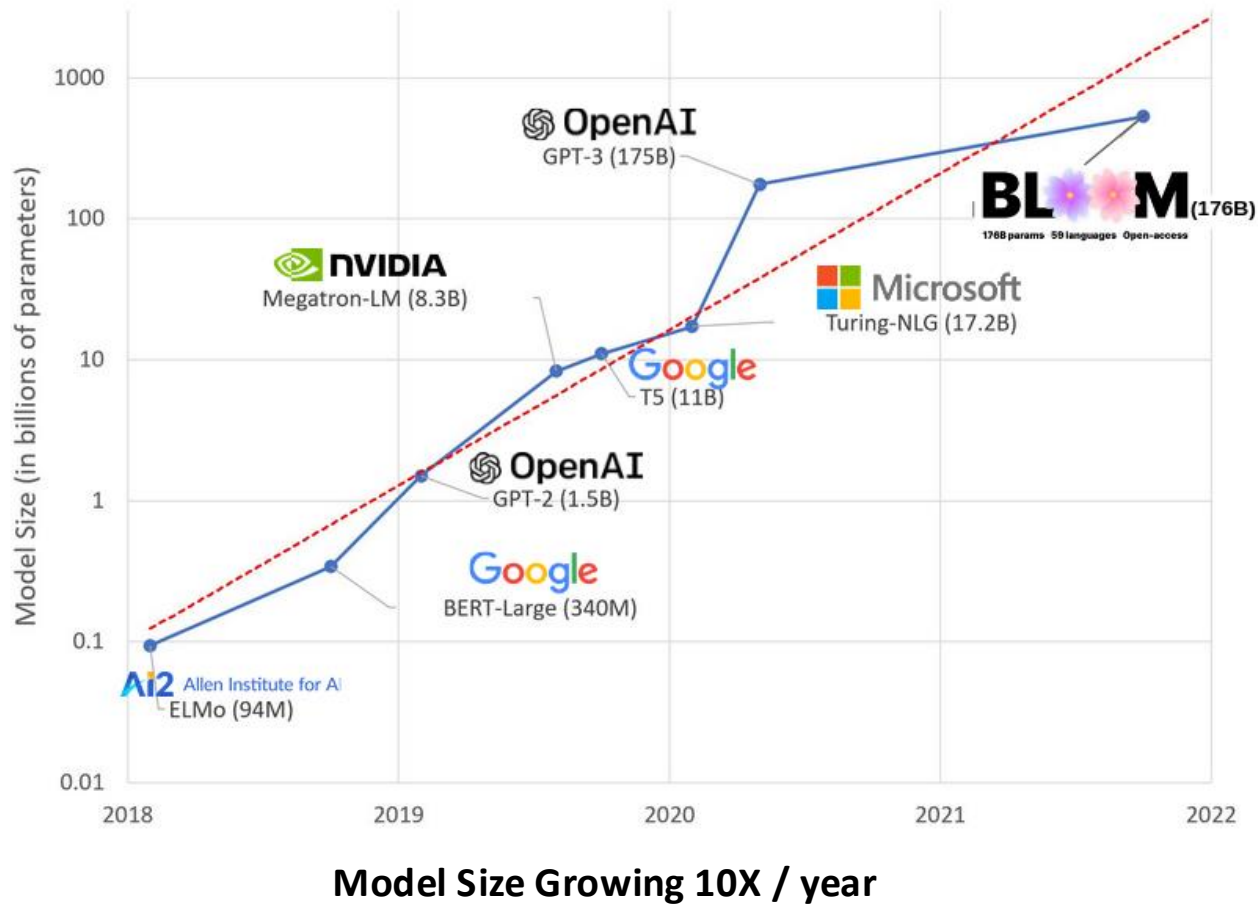# 高性能计算(HPC)的挑战与机会
# Challenges and Opportunities in HPC

爱普存储技术(杭州)有限公司　钟雷

apmemory

# Outline

- "Two Dark Clouds" over HPC Development
  高性能计算发展的 "两朵乌云"

  - Memory Bandwidth Limit ( or Information Bandwidth )
    内存带宽 / 信息带宽

  - Energy Bandwidth Limit: Power Delivery / Cooling
    能量带宽: 供电 / 冷却

- Inevitable 3DIC Trend & CapEx-Efficient Opportunity

  3DIC的必然趋势和与之对应的高效投资机会

apmemory

# AI's Incredible Growth



**Model Size Growing 10X / year**



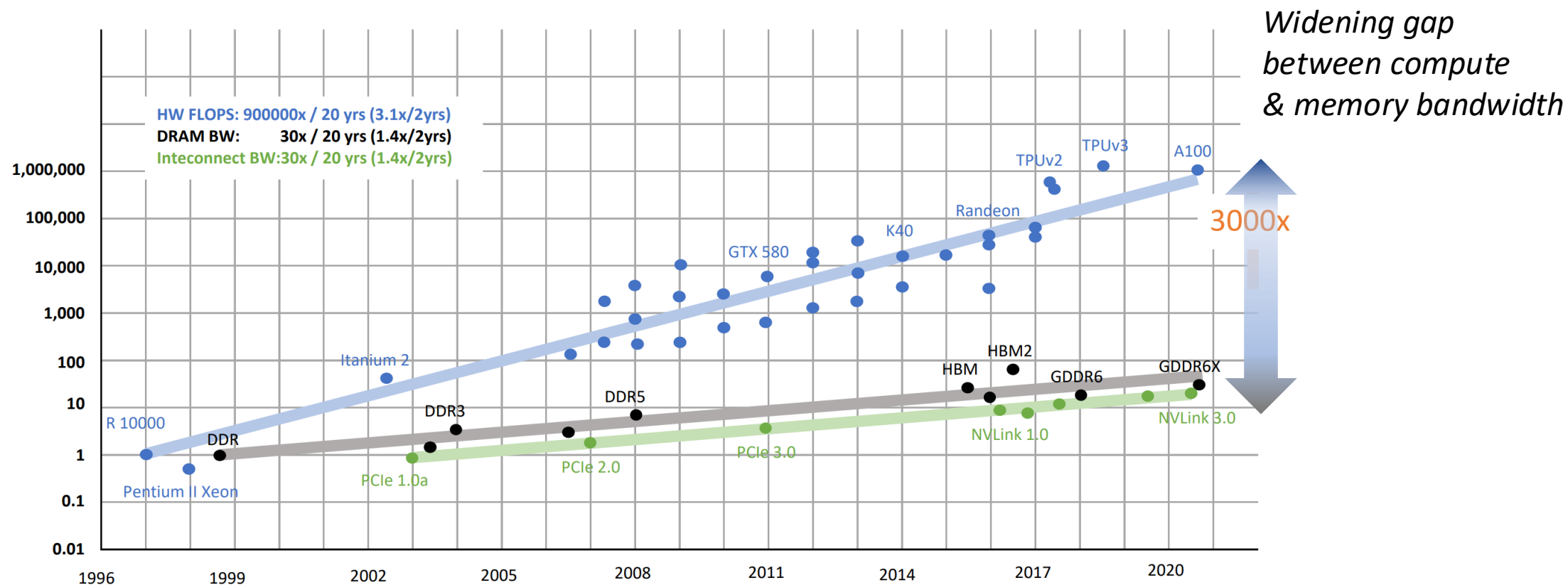**NVDA Market Cap on par with UK GDP**

# AI HPC Performance: "Two Dark Clouds"

- In 1900, Lord Kelvin called out "two dark clouds" of classical physics, which lead to the discovery of quantum physics and theory of relativity

- TCA's 50th Anniversary Special Event – the Taiwan Semiconductor Day Forum Debuts in Tokyo on April 2

- Dr. Chen Wen stated in his keynote speech
  - ① Memory Bandwidth Limit
  - ② Energy Bandwidth Limit

Source: https://show.computex.biz/NewsReleaseDetail.aspx?index=42488&category=68

# AI Computing Performance Bottleneck

- Hardware performance far outstrips memory bandwidth
- This is a problem, particularly for AI/ML computing



*Widening gap between compute & memory bandwidth*

**HW FLOPS:** 900000x / 20 yrs (3.1x/2yrs)
**DRAM BW:** 30x / 20 yrs (1.4x/2yrs)
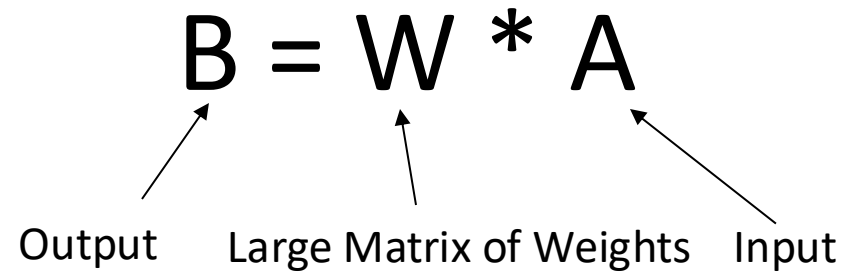**Inteconnect BW:** 30x / 20 yrs (1.4x/2yrs)

3000x

Source: Morgan Stanley Research, including (e) estimates.

# Performance Limit of Neuro-Network AI/ML

- Fundamental Computation of NN AI/ML
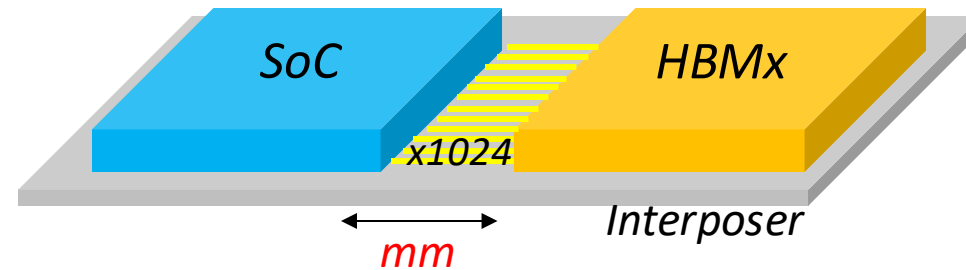
$$B = W * A$$

Output    Large Matrix of Weights    Input

- Measure of AI/ML performance

$$FLOPS = \frac{\# \text{ of FLOP}}{Sec} = \frac{\# \text{ of FLOP}}{Byte} \times \frac{Byte}{Sec}$$

*This is Memory bandwidth*

apmemory

# Scaling Limitations of HBM in 2.5D
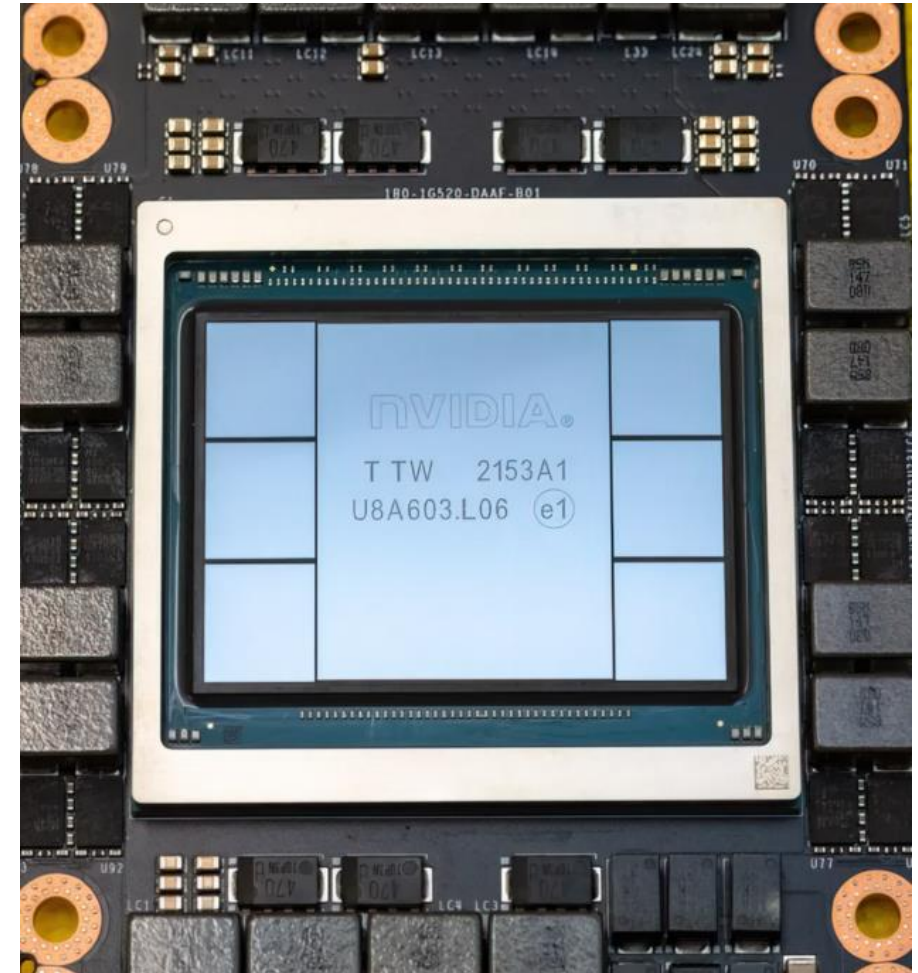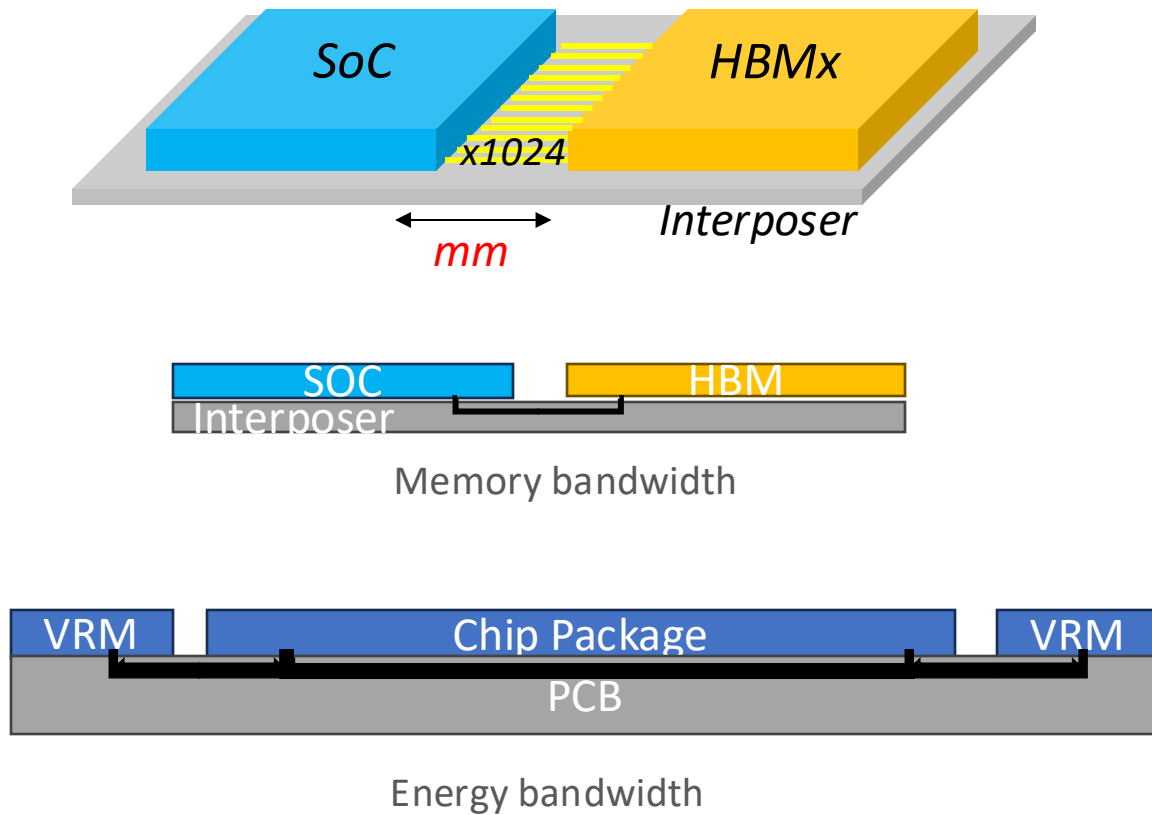


SoC

HBMx

x1024

mm

Interposer

# of connections ➔ Limits Bandwidth

"beach front" of SoC ➔ Limits # of HBM placement

Length of connections ➔ Limits Power consumption
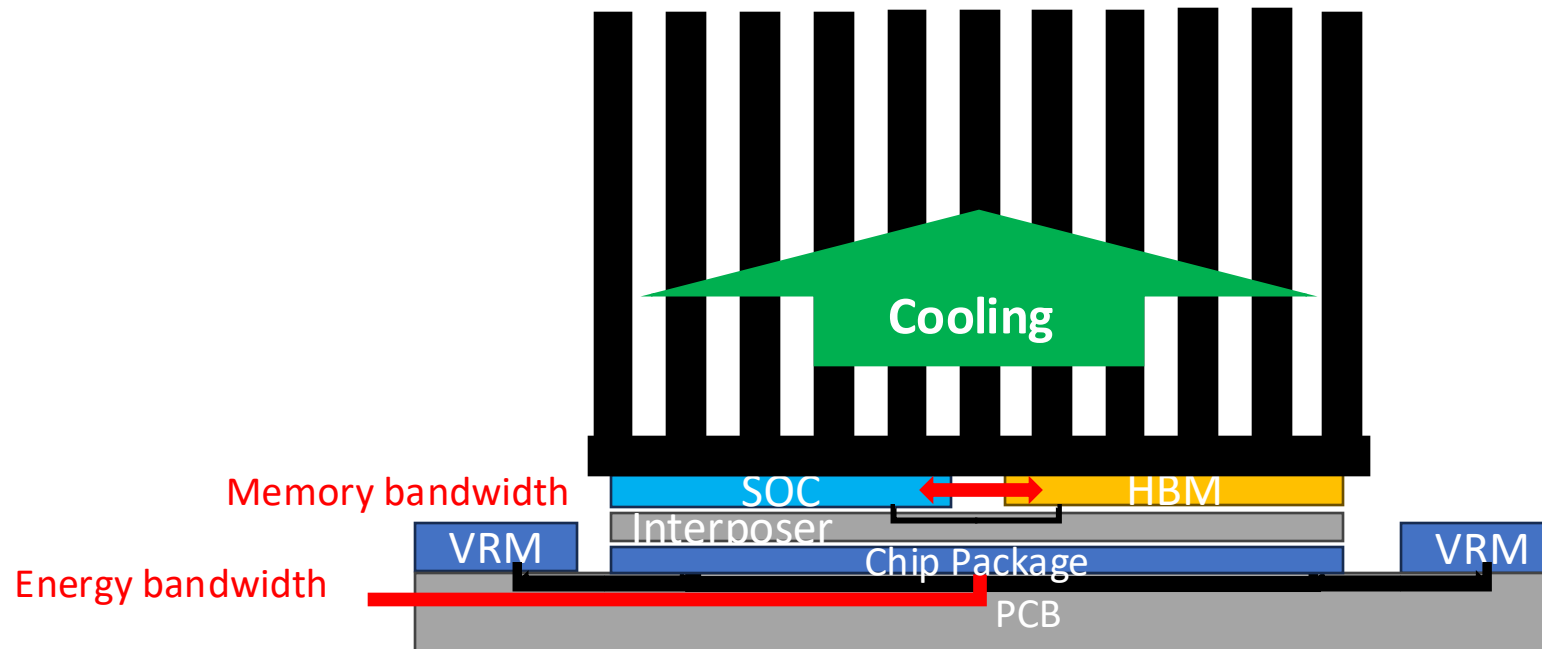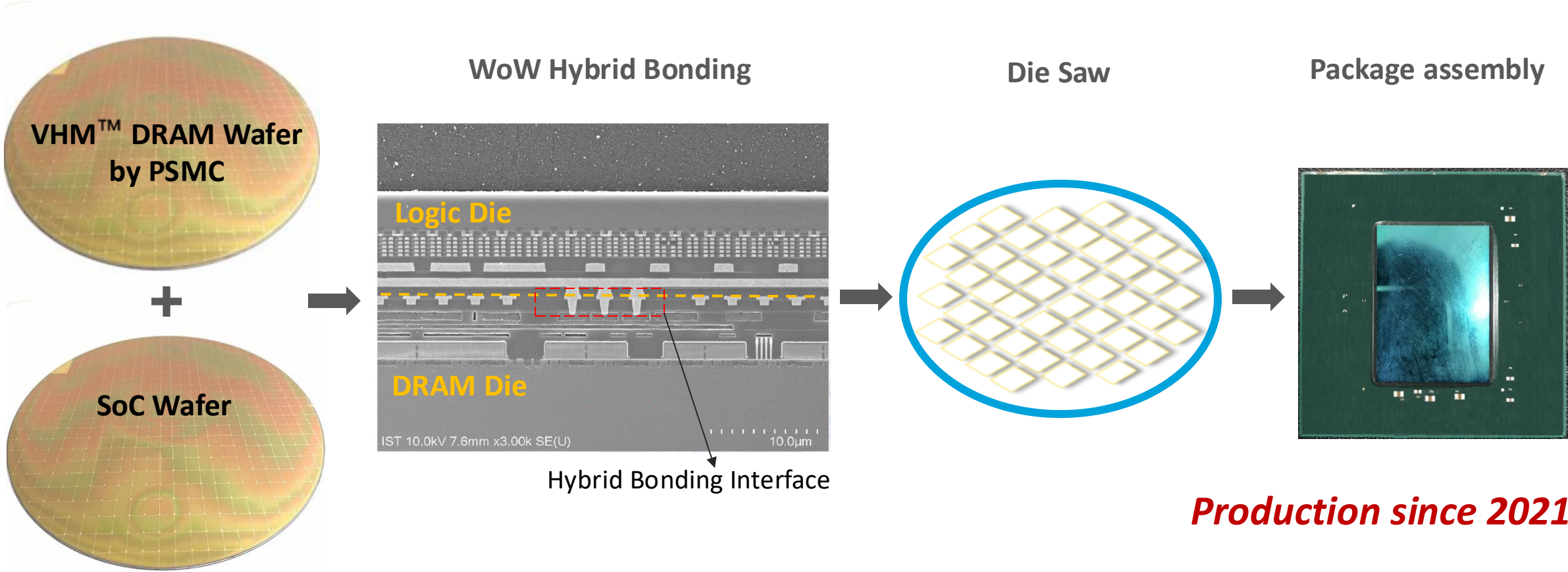
# State-of-the-art GPGPU



**Both memory and energy are delivered horizontally over a 2D plane**
**"Flux density" (bandwidth/area) limited by physics**

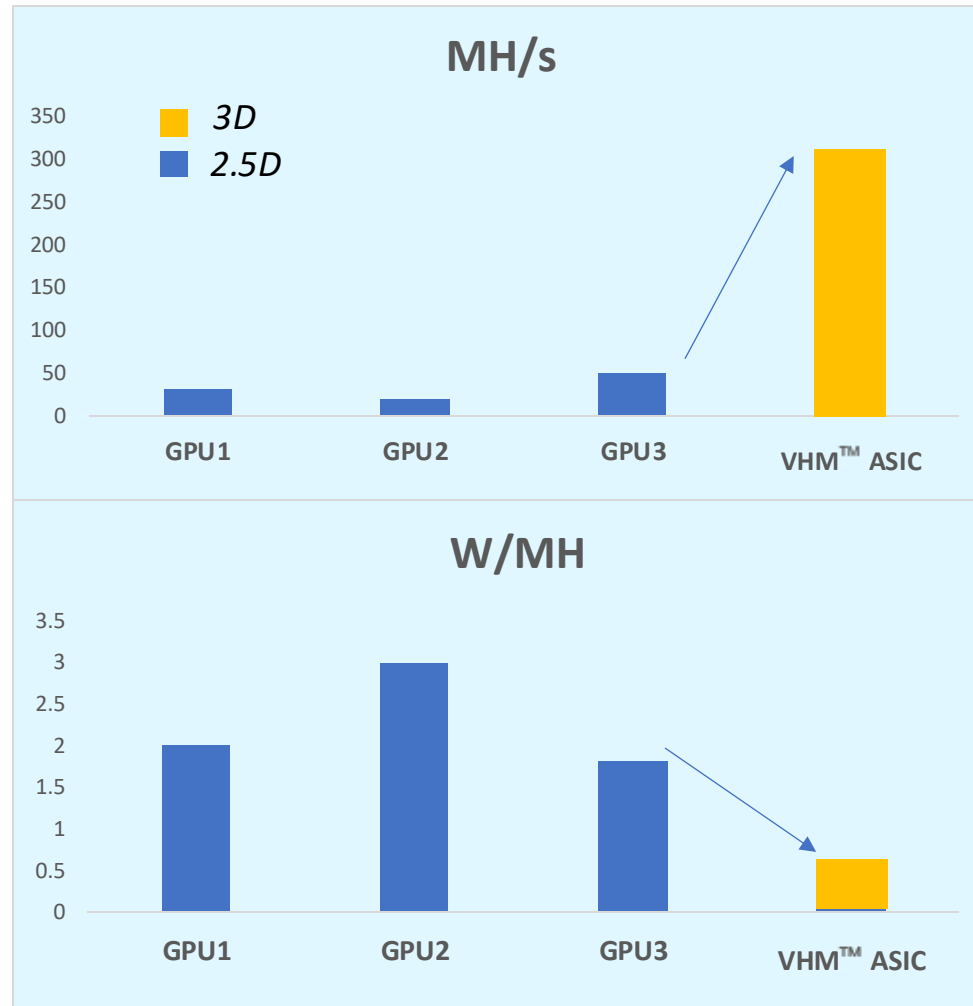NVIDIA H100

apmemory

# The 3<sup>rd</sup> Dimension

- Cooling is already done vertically in 3D

- Memory and Power must also be delivered vertically in 3D to breakthrough bandwidth limits

# Memory in 3D | VHM™ in Production Since 2021



**VHM™ DRAM Wafer by PSMC**

**+**

**SoC Wafer**

**WoW Hybrid Bonding**

Logic Die

DRAM Die

IST 10.0kV 7.6mm x3.00k SE(U)    10.0μm

Hybrid Bonding Interface

**Die Saw**

**Package assembly**

*Production since 2021*

# Actual Performance of 3DIC-enabled Systems



**MH/s**

Legend: 3D (yellow), 2.5D (blue)

350, 300, 250, 200, 150, 100, 50, 0

GPU1, GPU2, GPU3, VHM™ ASIC

**W/MH**

3.5, 3, 2.5, 2, 1.5, 1, 0.5, 0

GPU1, GPU2, GPU3, VHM™ ASIC

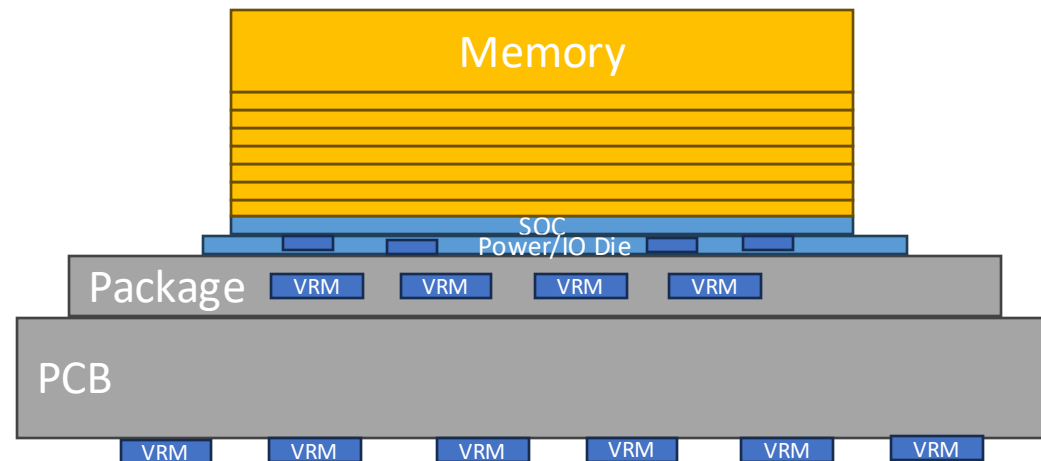*First Gen 3DIC-enabled Systems vs Leading GPUs*

## 10X Performance

## 1/3X Power

*\* Based on ETH mining algorithm*

apmemory

# Envisioned Future 3DIC HPC System

- Memory bandwidth in 3D: 10X proven, 1000X possible

- Energy bandwidth in 3D: TBD, >>2X
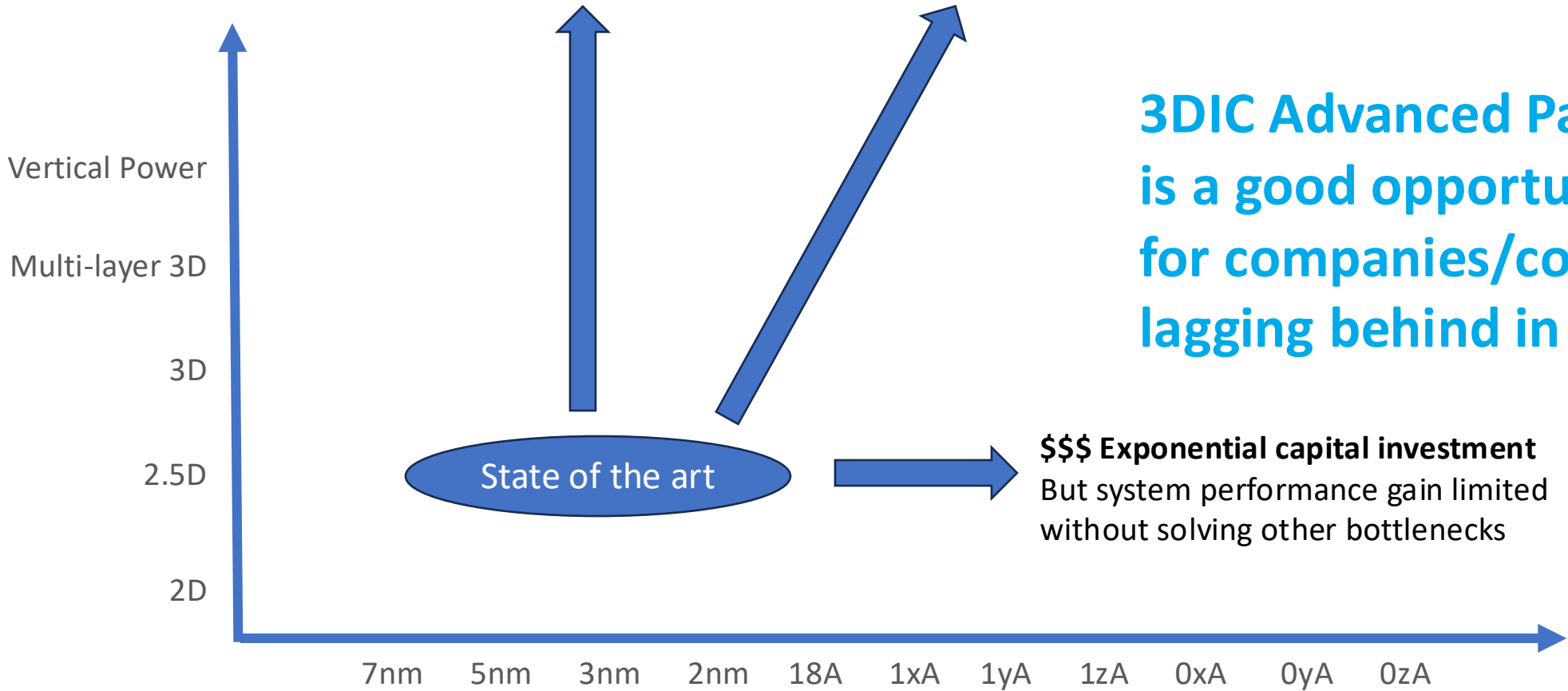  - VRM placed in various places vertically

# Future Roadmap: Scaling vs 3DIC Packaging

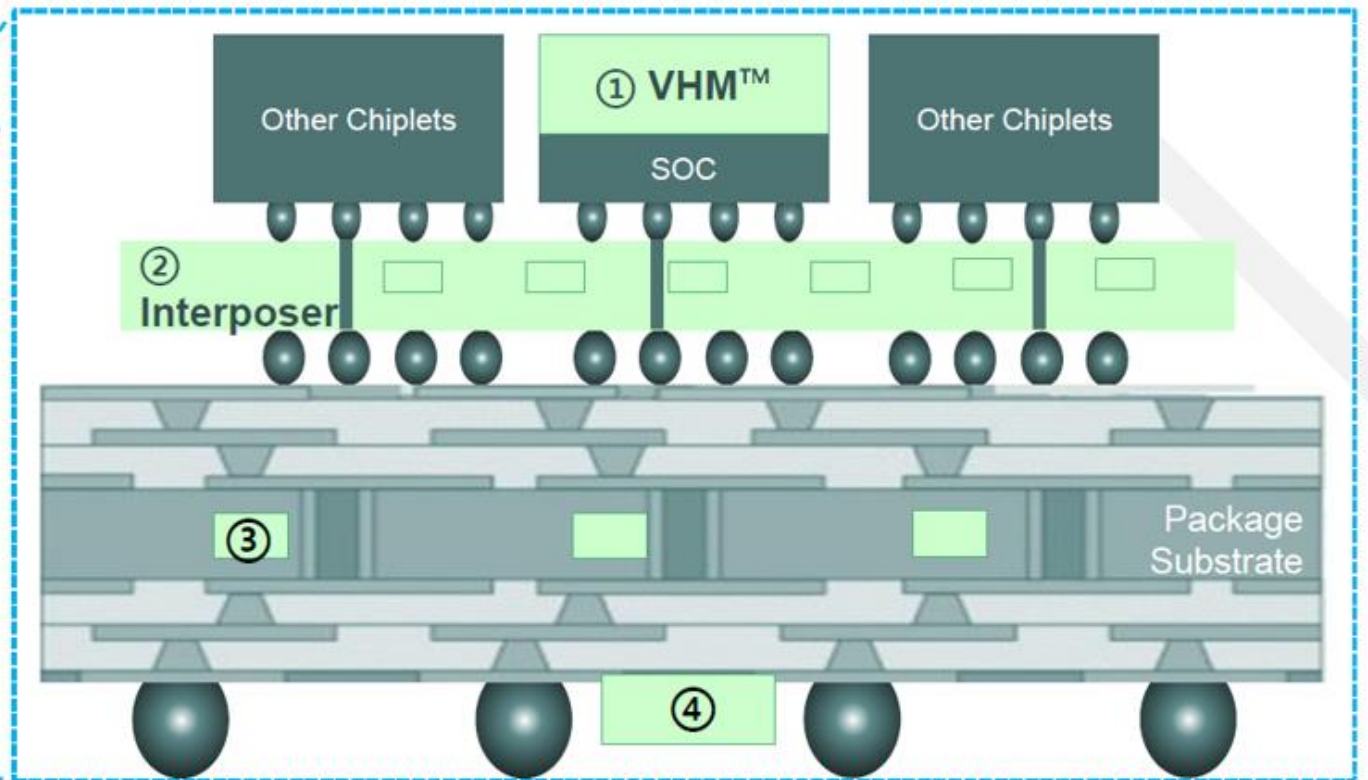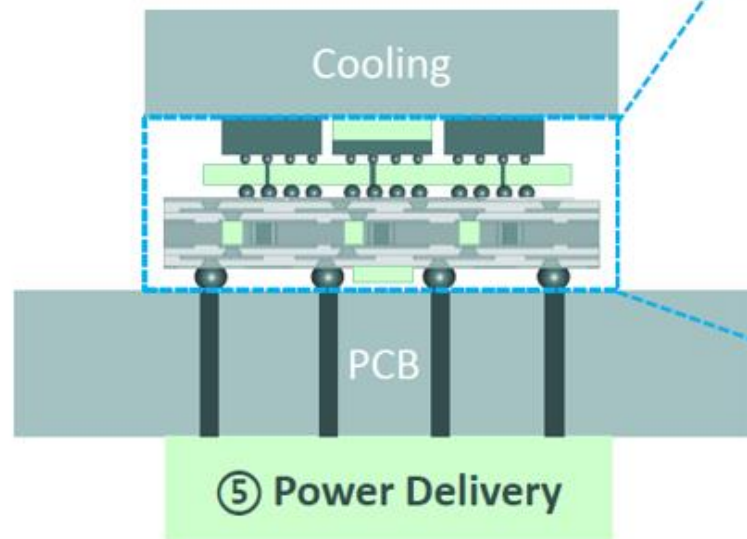**$ Relatively light Capex**
Bigger bang for the Buck!

**Likely scenario**

**3DIC Advanced Packaging
is a good opportunity
for companies/countries
lagging behind in scaling**

Vertical Power

Multi-layer 3D

3D

2.5D   State of the art

**$$$ Exponential capital investment**
But system performance gain limited
without solving other bottlenecks

2D

7nm   5nm   3nm   2nm   18A   1xA   1yA   1zA   0xA   0yA   0zA

# Positioning of **apmemory** solutions in AI/HPC

① VHM™ stacked with Compute

② Interposer with *S-SiCap™

③ S-SiCap™ embedded in Package Substrate

④ S-SiCap™ on landside

⑤ Power delivery solutions (future)



*S-SiCap™: Stack Silicon Capacitor, AP Memory's SiCap technology which uses a stack capacitor

# About apmemory



**IoTRAM™ – The Ideal Memory for IoT**

Leading in product innovation and market share

**VHM™ – The AI/HPC Memory Solution**

World's first 3D integration of DRAM and SoC chips

**S-SiCap™ – The Performance Enabler**

Highest capacitance density for high performance applications