

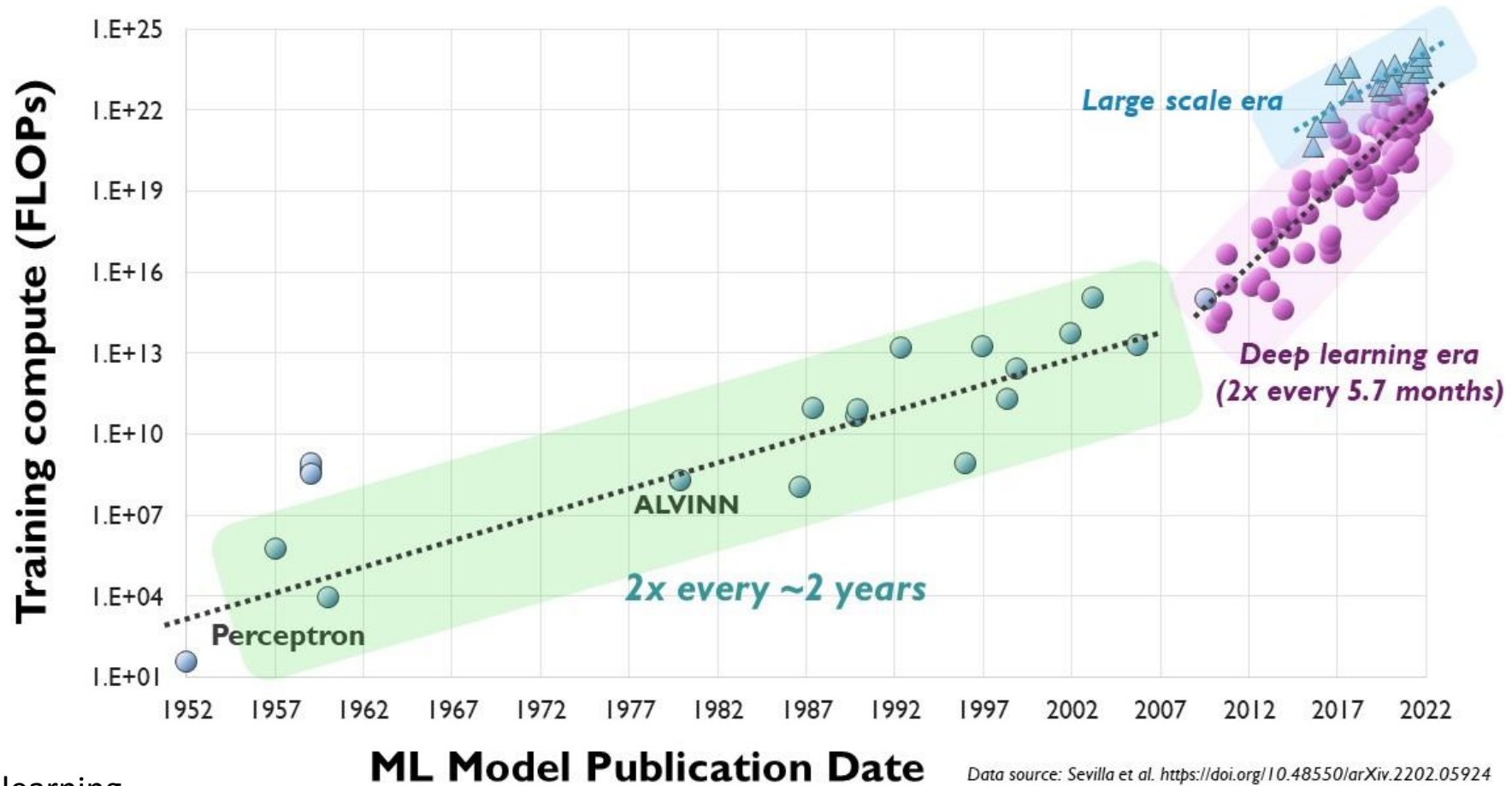
Challenges in Advanced Computing and Functionalities International Cooperation on Semiconductors

Future technologies in Advanced Computation

Nadine Collaert (imec) and Olivier Faynot (CEA-LETI)

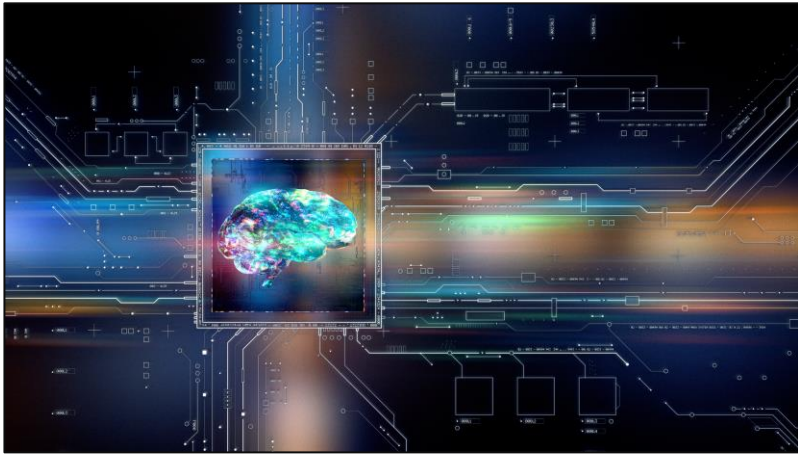
- Introduction: trends and challenges
- Computing roadmap:
 - CMOS device architecture
 - 2D materials for FEOL
 - New materials for BEOL
 - Lithography
 - CMOS 2.0
- Memory technologies
- Beyond Von Neumann disruptive approaches:
 - Near or in-memory computing
 - Quantum computing
- Heterogeneous integration: from chiplets to functional backside
- Analysis: EU and non-EU actors
- Conclusions

Compute needs for ML continue to grow



ML=machine learning

GPUs for Training



High throughput parallel compute
Very high memory bandwidth
Very high GPU-GPU bandwidth

AR/VR



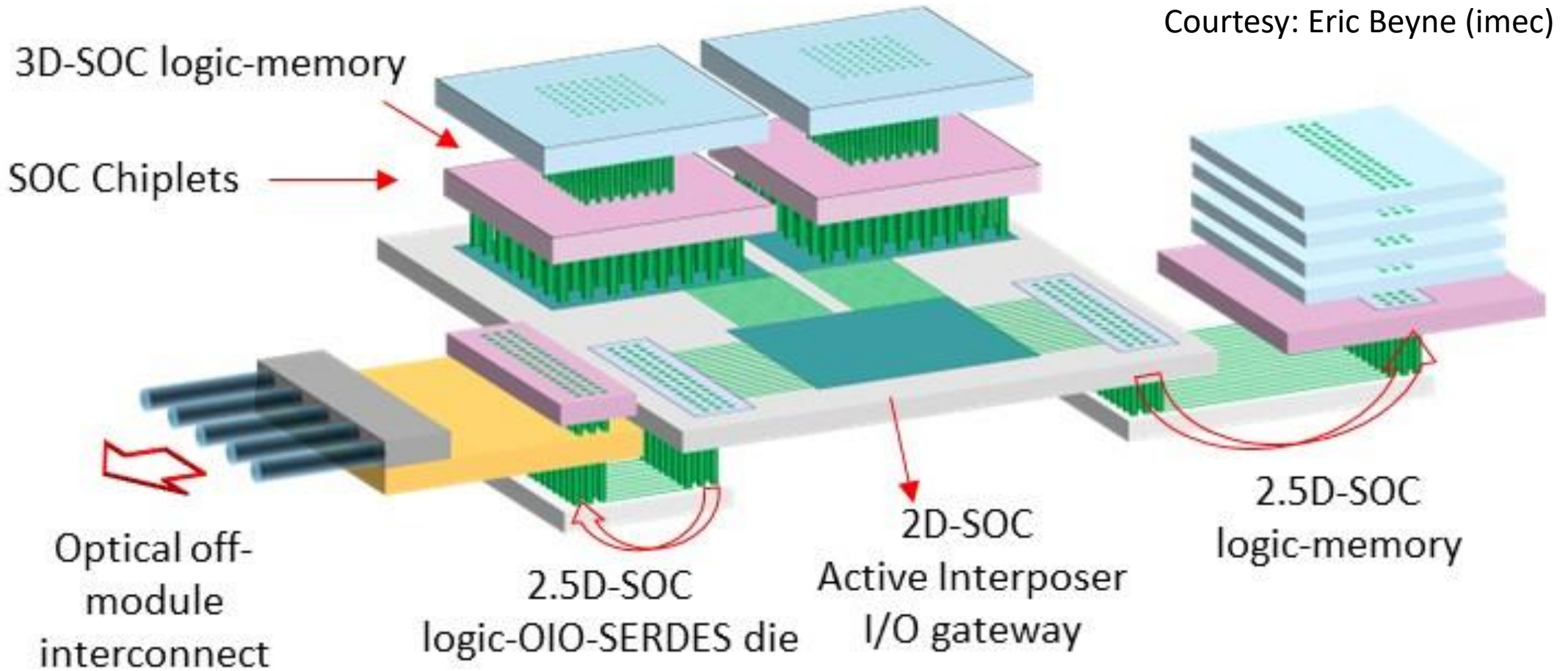
Low power
Ultra low latency
High memory bandwidth
Small form factor

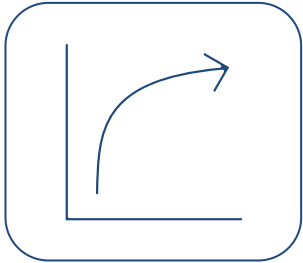
Autonomous driving



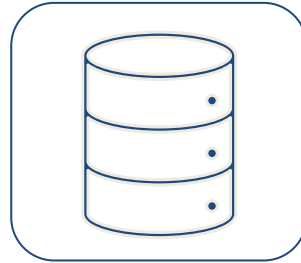
Multi-sensor fusion
Distributed real-time computation
Reliable and explainable AI

Future systems are heterogeneous





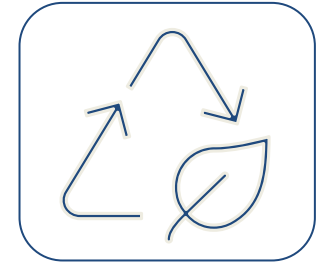
CMOS and
DRAM PPAC



Memory
Wall



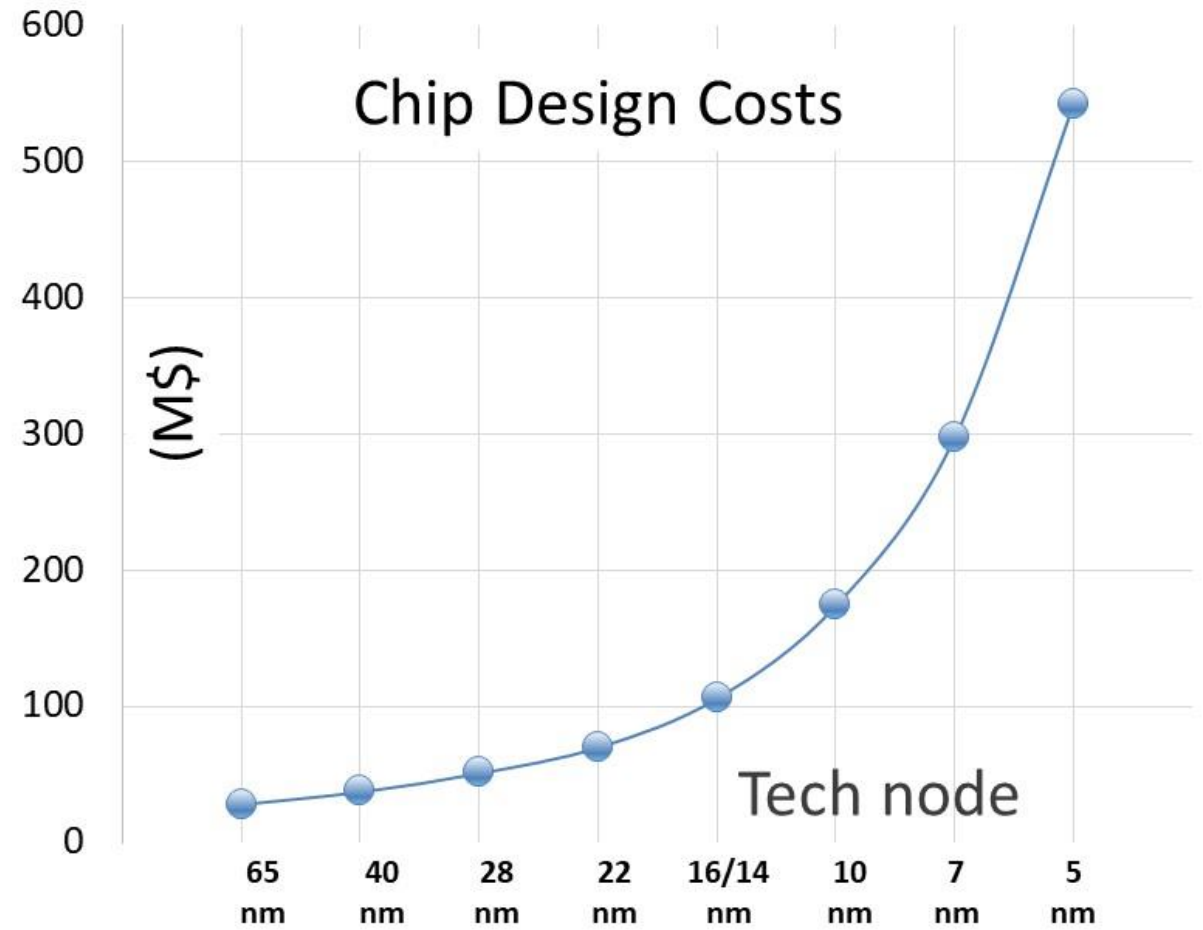
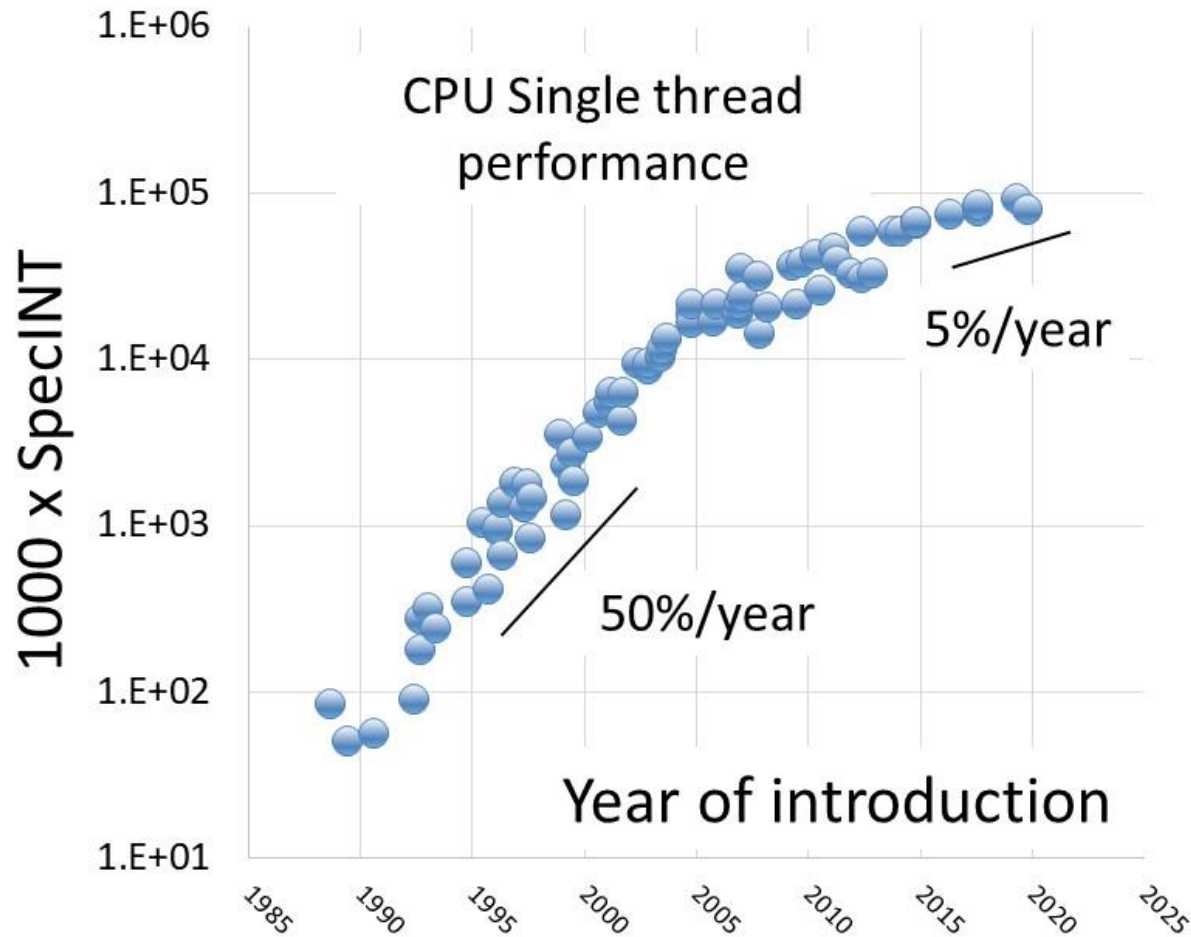
Power
Wall



Sustainable
Manufacturing

PPAC=Power-Performance-Area-Cost

Slowdown in System Performance and Increasing Costs

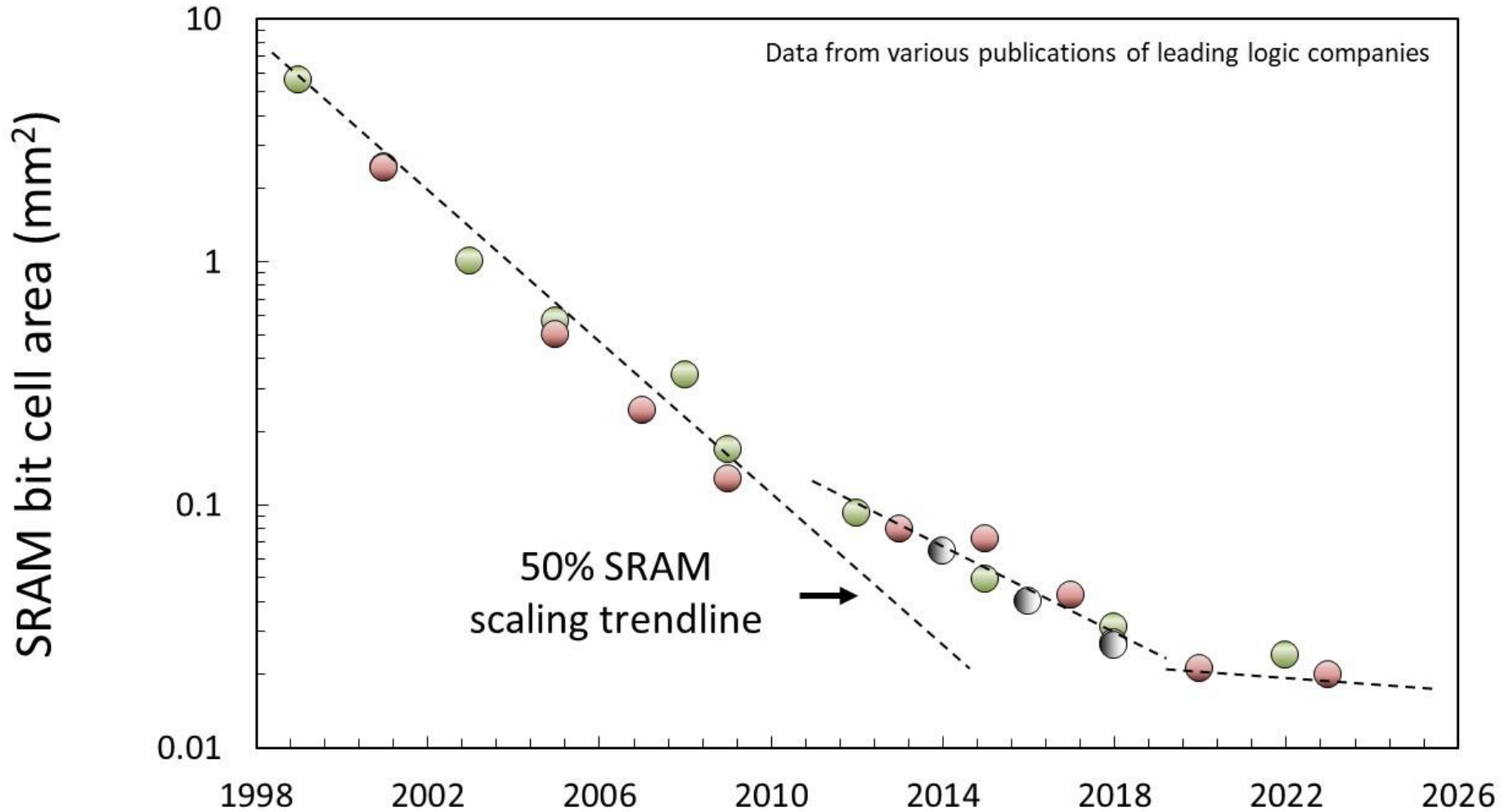


Based on original data plotted by M. Horowitz, F. Labonte, O. Shachan, K. Olukotun, L. Hammond, C. Batten. Additional data compiled by K. Rupp

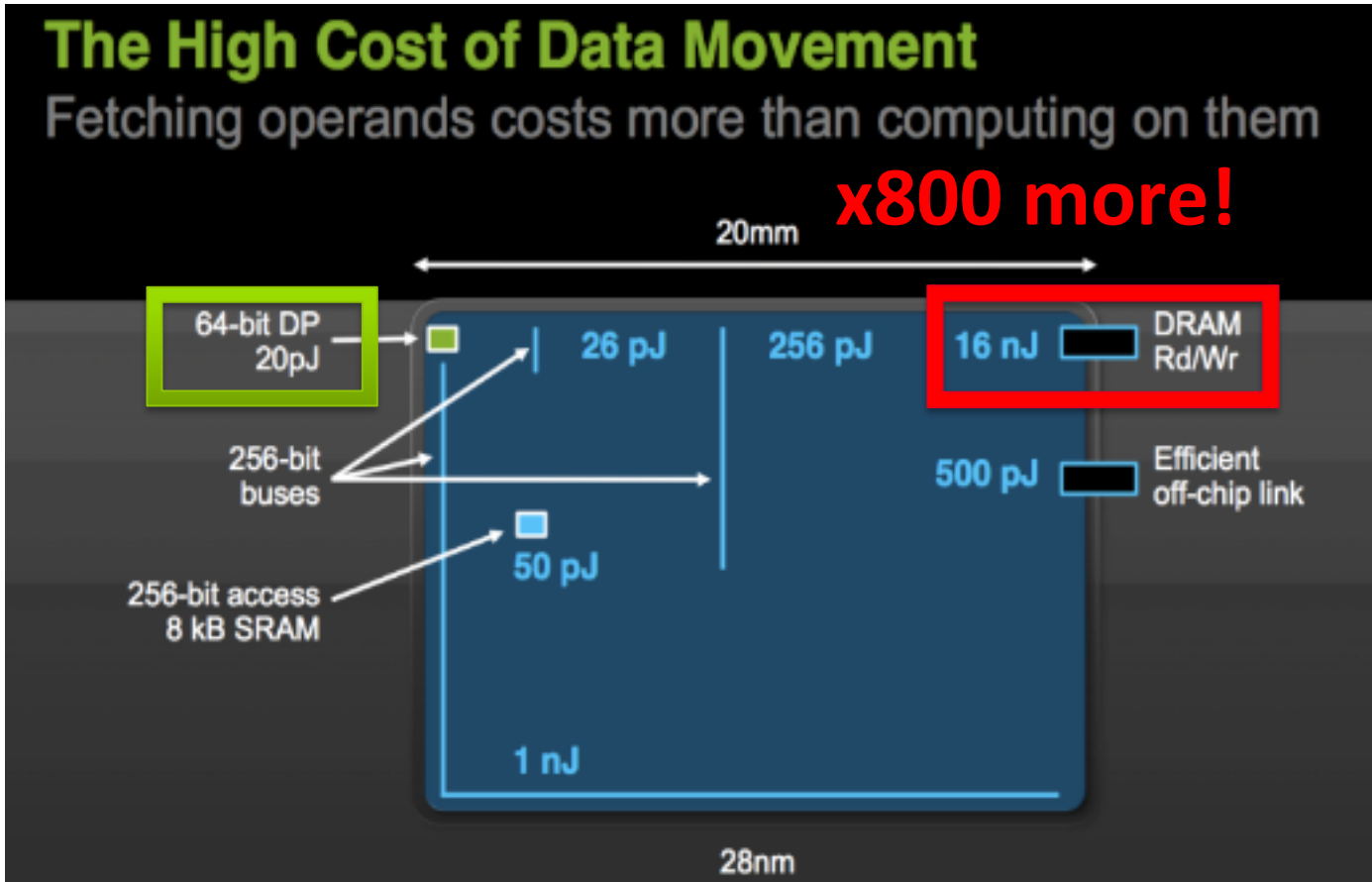
Source: "AI Chips and why they matter", S. Khan and A. Mann, 2020

[iference](#)

SRAM scaling has slowed down

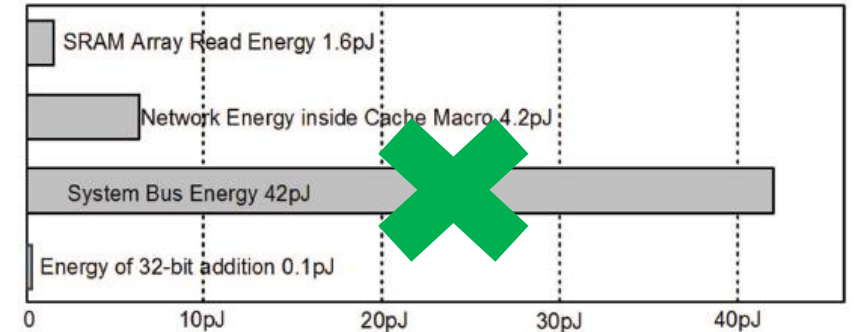


The cost of moving data



Bill Dally, "To ExaScale and Beyond", 2010

[J. Wang – ISSCC'19]



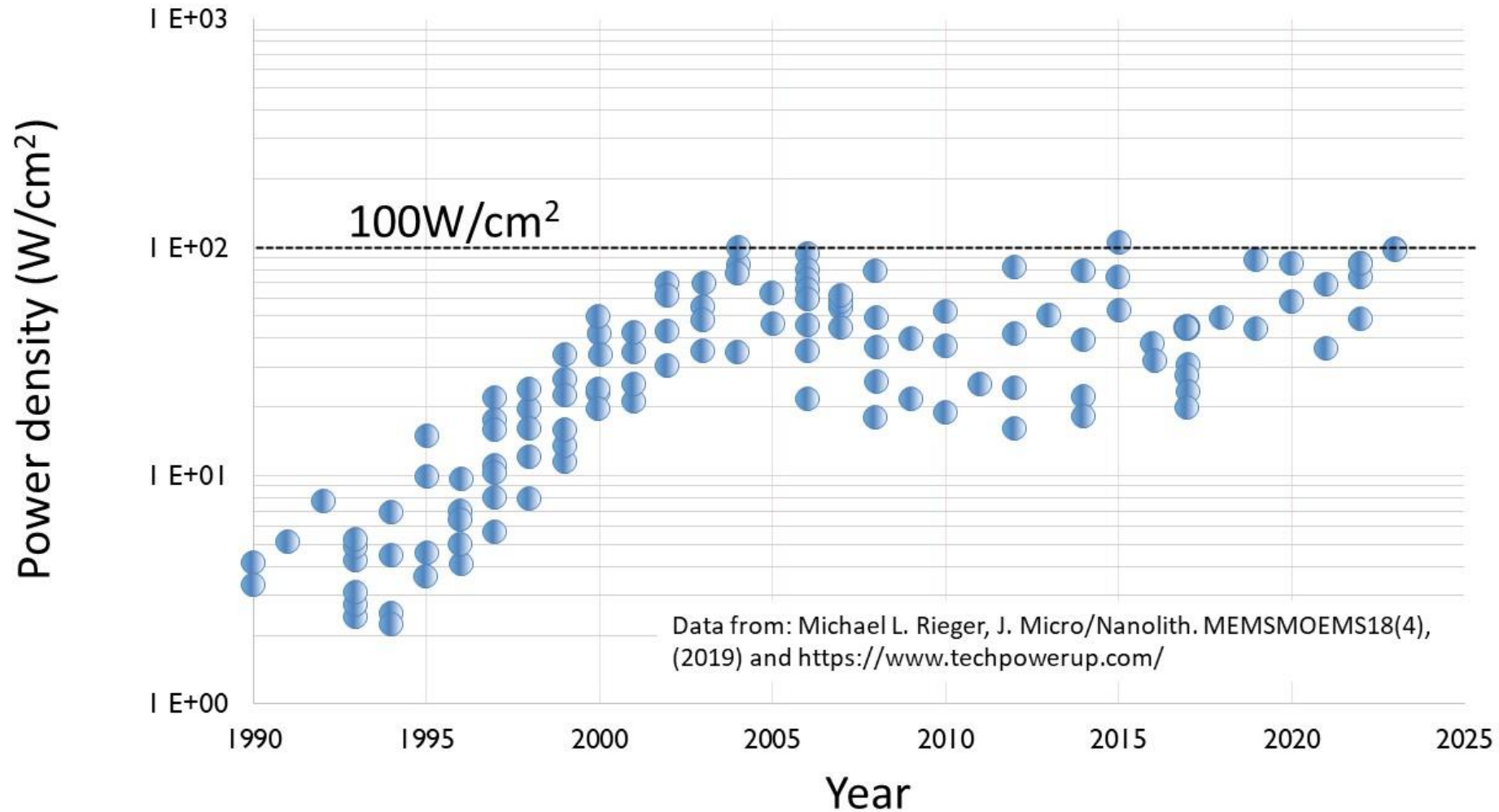
~90% of energy is in **data transfer**
 → IMC could lead **8x reduction**

Operation energy is negligible

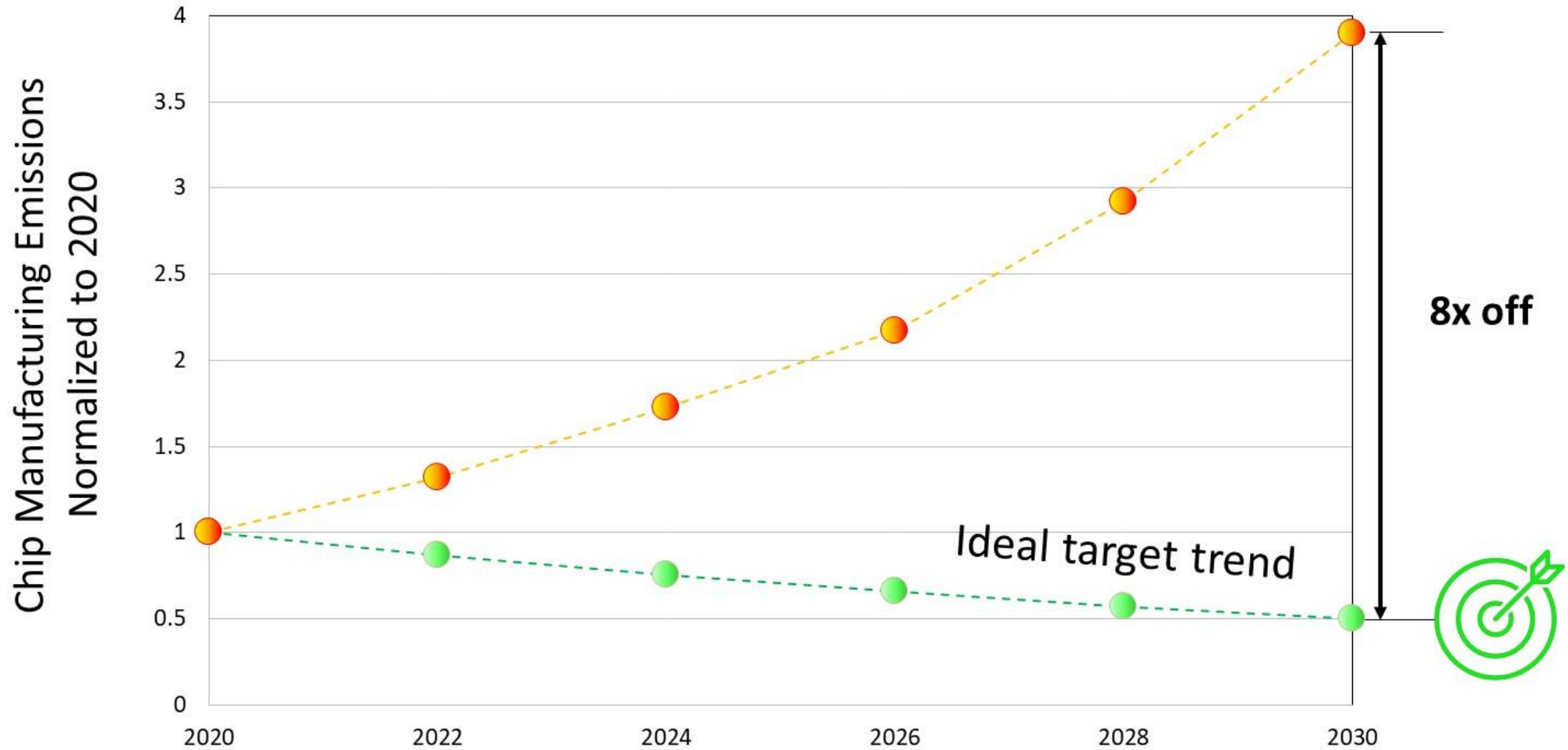


Memory access and control energies dominate

Chip Cooling Technologies Limit Power Density of Chips

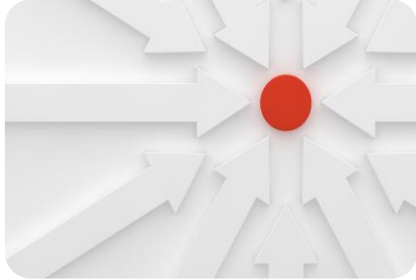


Carbon Emissions: “Do Nothing” Scenario



Constant electricity mix (0.49 kCO₂eq/kWh), Abatement and GHG global warming potential according to IPCC assumed for the years 2020-2030. Volume technology mix from IBS “Foundry Market Trends and Strategic Implications” Vol 30, N 12, Dec 2021. Logic nodes only. *imec.netzero emissions estimate of imec process nodes representative of foundry nodes.

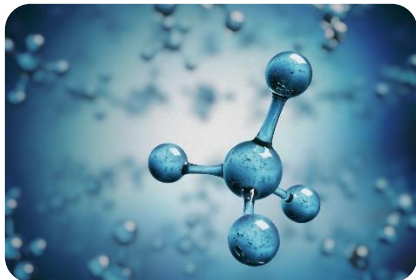
How to Enable Sustainable Manufacturing



Adopt 50% emissions reduction this decade as target



Must-do for fabs: green power, best abatement



Research new materials: F-gas and PFAS alternatives

The required gain in energy efficiency

**>1000x
by
2030**

CMOS scaling

Memory technologies

Disruptive Computing

Chiplet & 3D System

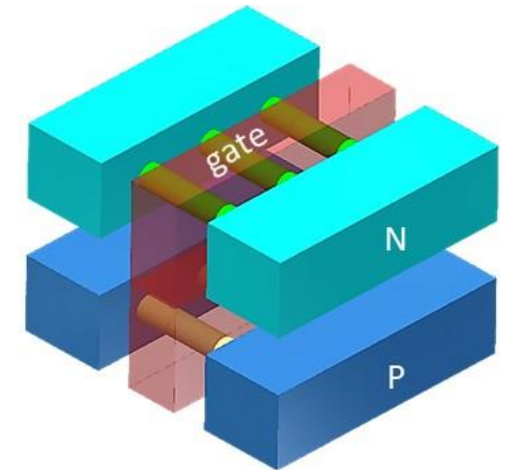
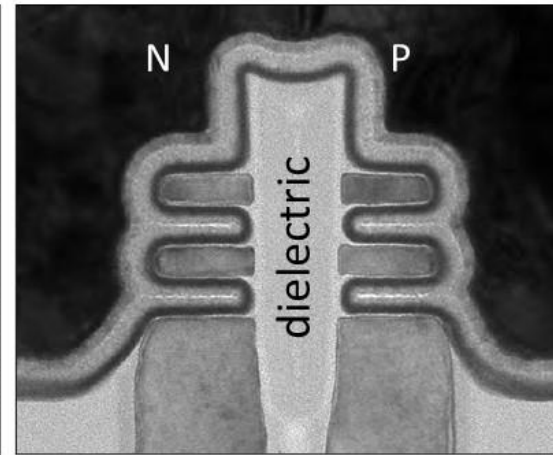
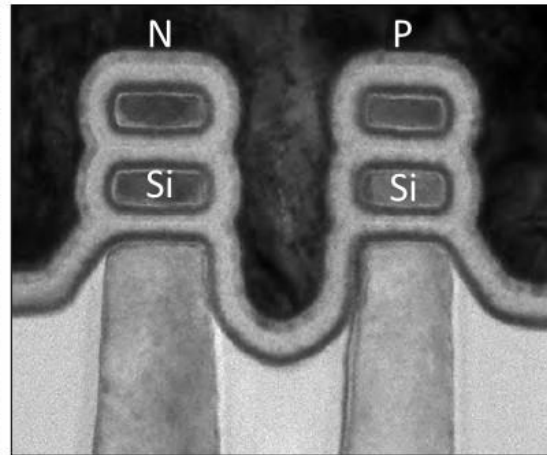
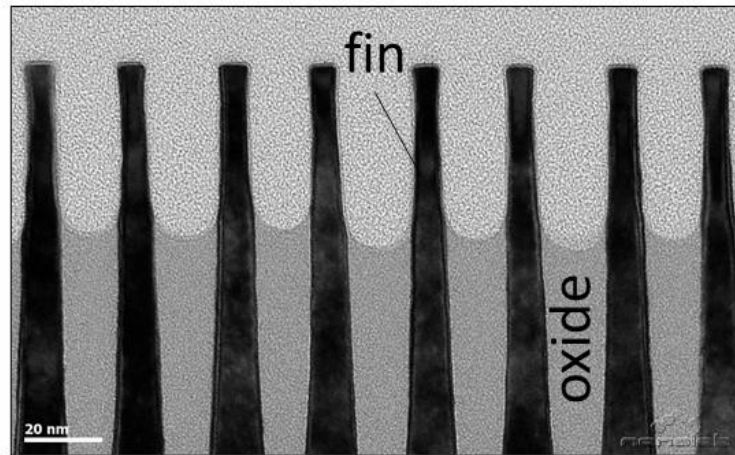
New device architectures to extend scaling

FINFET

NANOSHEETS

FORKSHEET

CFET



N14-N3

N2-A14

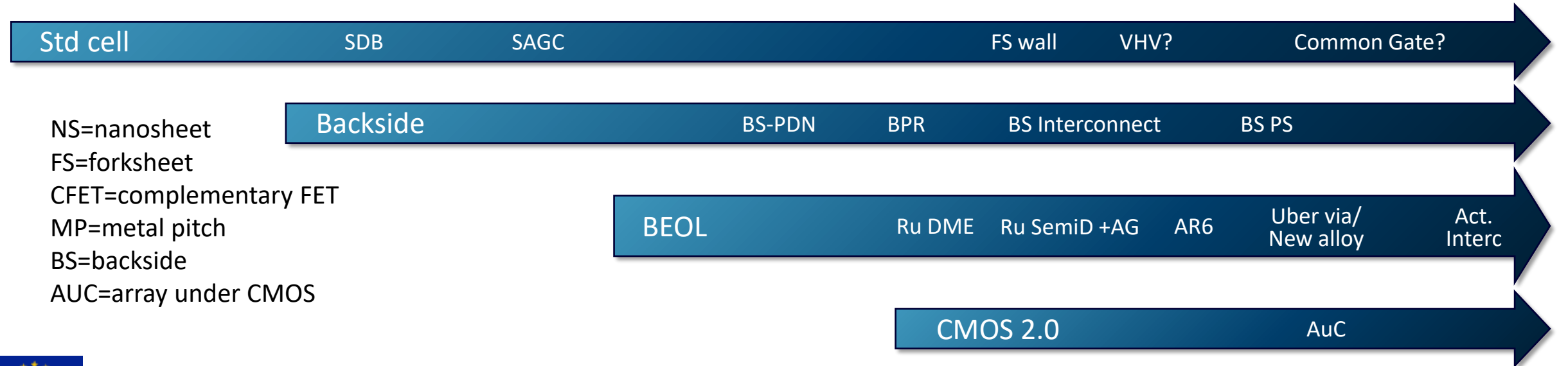
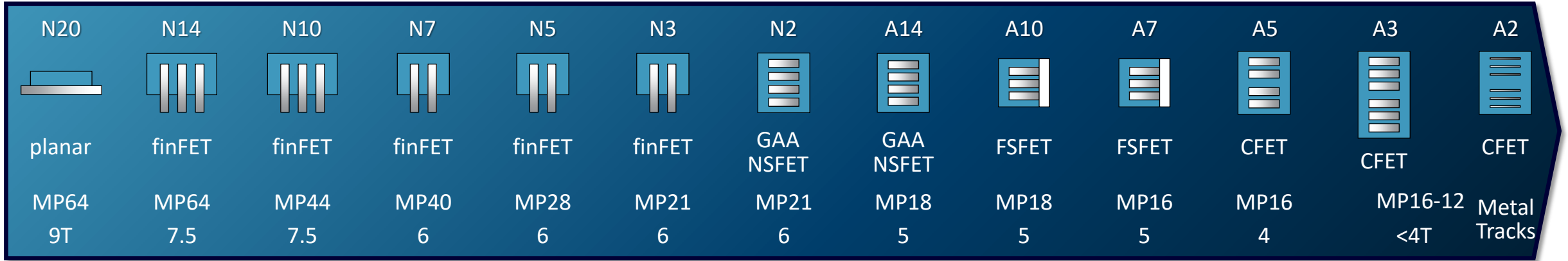
A10-A7

A5-A2

N. Collaert, "Advancements in IC Technologies: A look toward the future," in IEEE Solid-State Circuits Magazine, vol. 15, no. 3, pp. 80-86, 2023, doi: 10.1109/MSSC.2023.3280433.

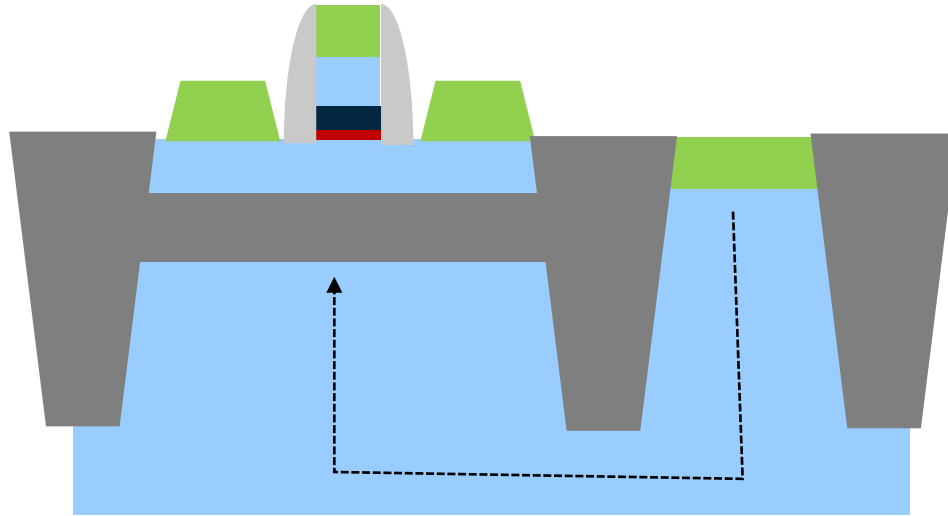
Logic scaling roadmap

A. Veloso, ICOS workshop, April 2023.

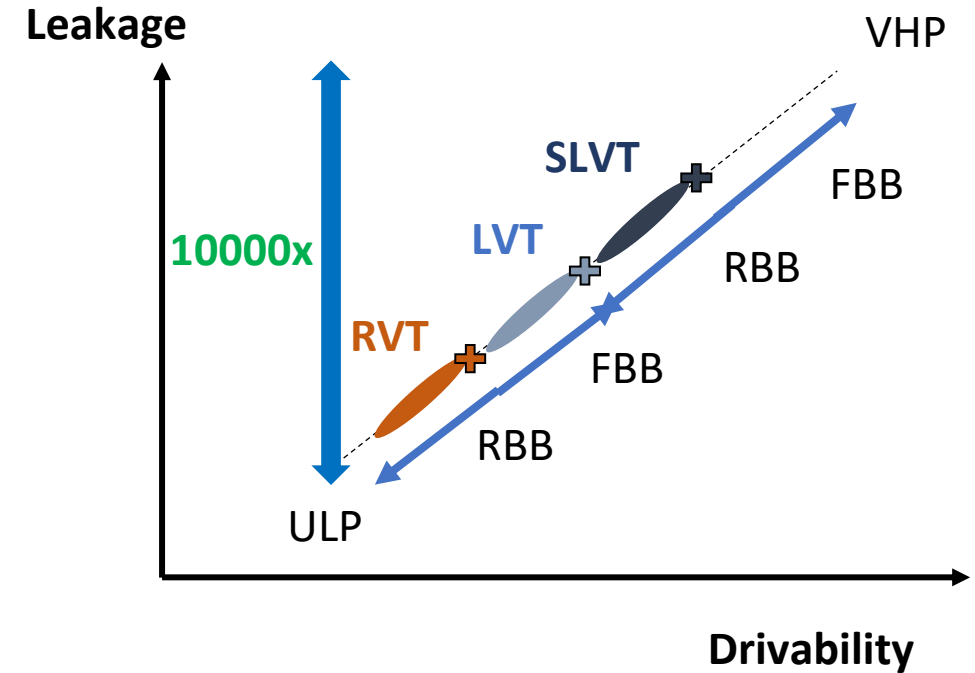


NS=nanosheet
FS=forksheets
CFET=complementary FET
MP=metal pitch
BS=backside
AUC=array under CMOS

FD-SOI Technology

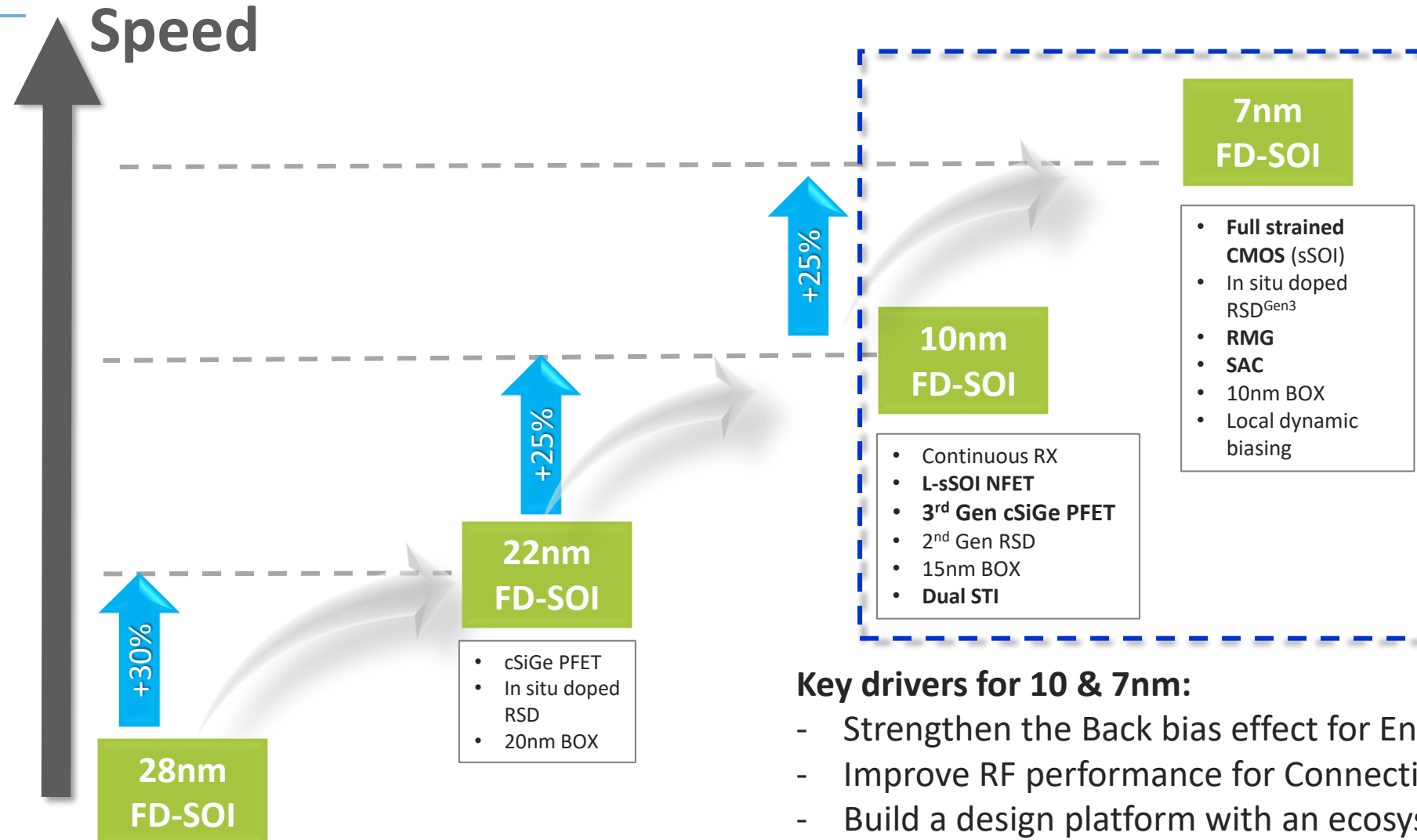


	22FDX	14nm FinFET	28nm Bulk	45nm PDSOI
f_T n-FET [GHz]	347	314	310	296
f_{max} n-FET [GHz]	371	180	161	342
f_T p-FET [GHz]	242 275 (mmWave)	285	185	-
f_{max} p-FET [GHz]	288 299 (mmWave)	140	104	-



RBB: Reverse Body Bias
FBB: Forward Body Bias
ULP: Ultra-Low Power
VHP: Very High Performance

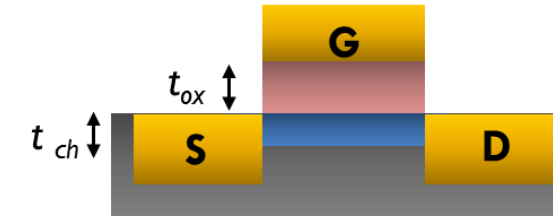
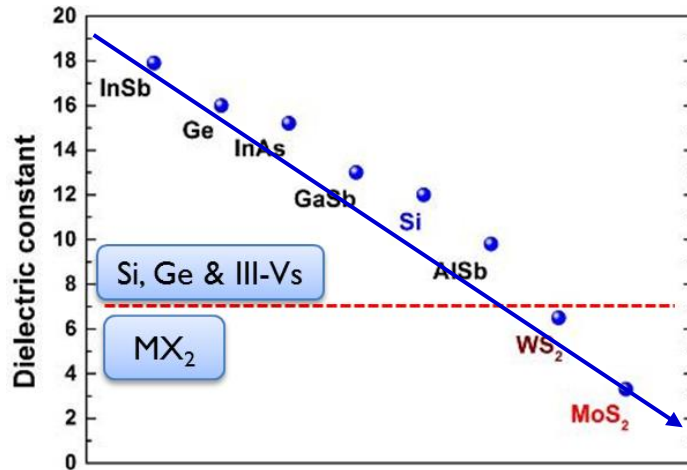
FD-SOI Technology Roadmap



Key drivers for 10 & 7nm:

- Strengthen the Back bias effect for Energy saving
- Improve RF performance for Connectivity
- Build a design platform with an ecosystem

2D Atomic Channels: Next generation logic devices

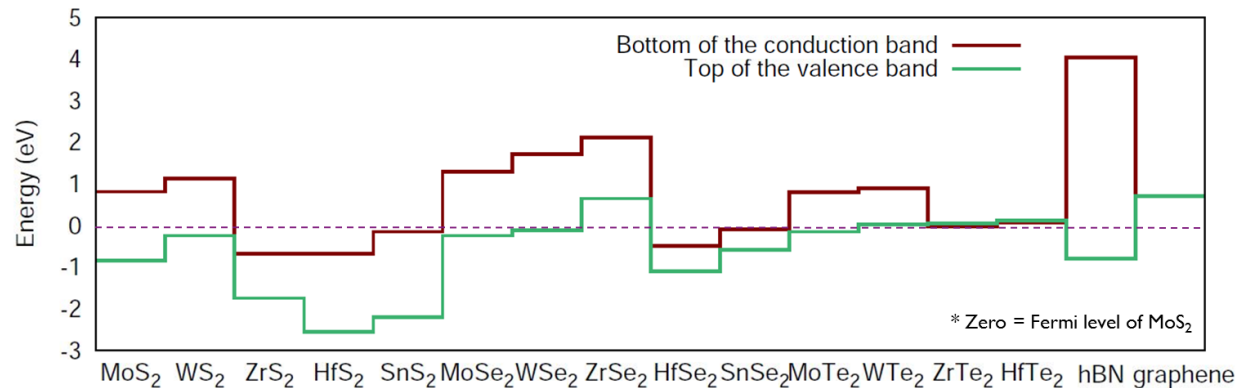


Characteristic length of short channel FETs:

$$\lambda = \sqrt{\frac{\epsilon_{ch}}{\epsilon_{ox}} t_{ch} \cdot t_{ox}}$$

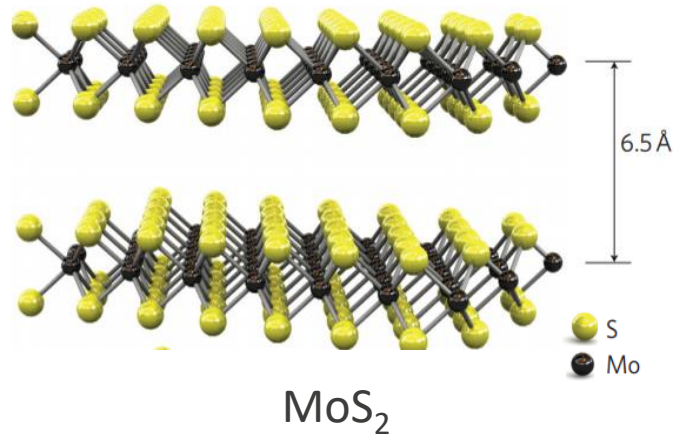
Expect reduced short channel effects in planar devices

Ultra-thin materials

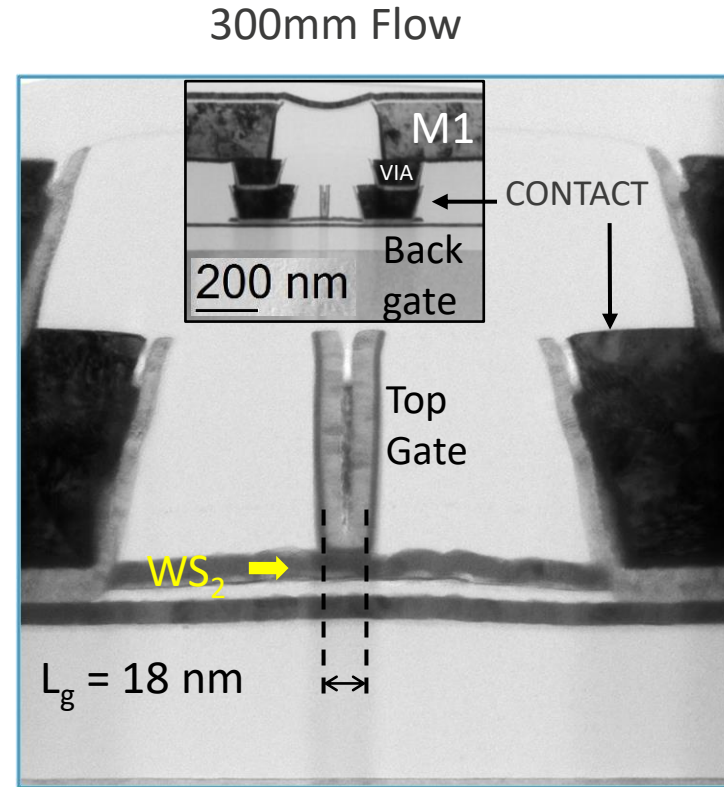


Choice of bandgaps and band alignment

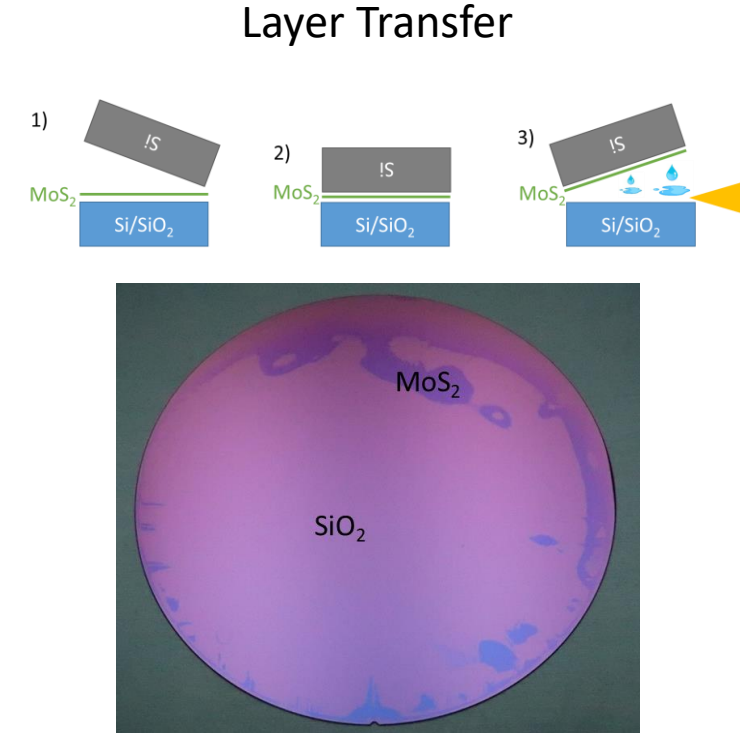
No/few dangling bonds at interfaces



Monolayer channel thickness enables gate length scaling while keeping high performance



I. Asselberghs et al, IEDM 2020

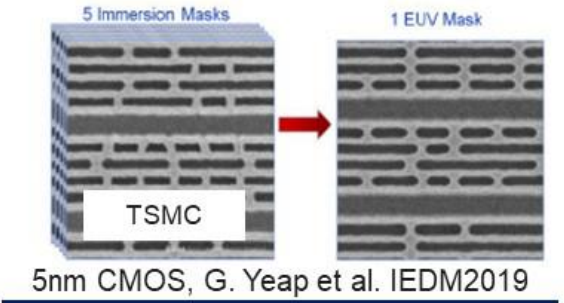
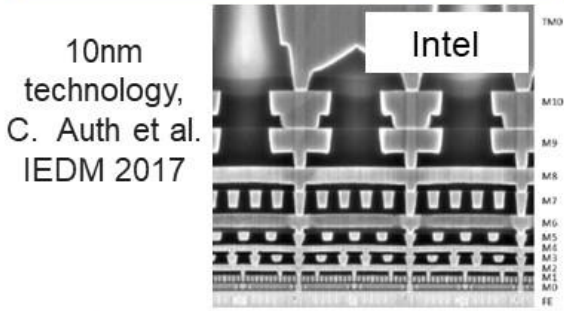


CEA-Leti, unpublished

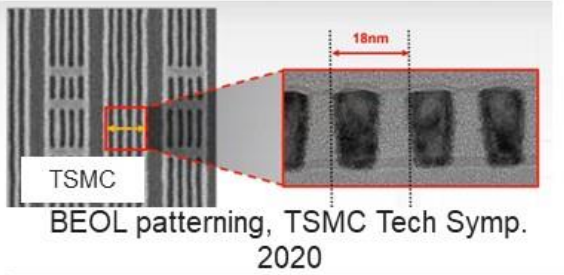
Industry BEOL trends

Courtesy: Zsolt Tokei (imec)

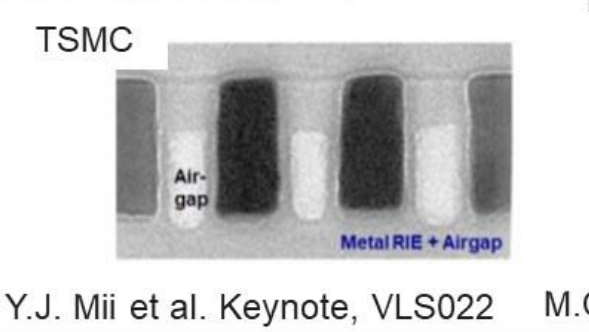
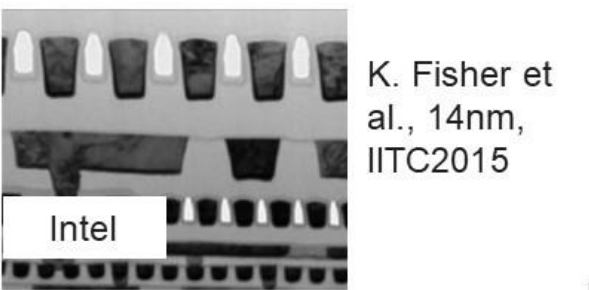
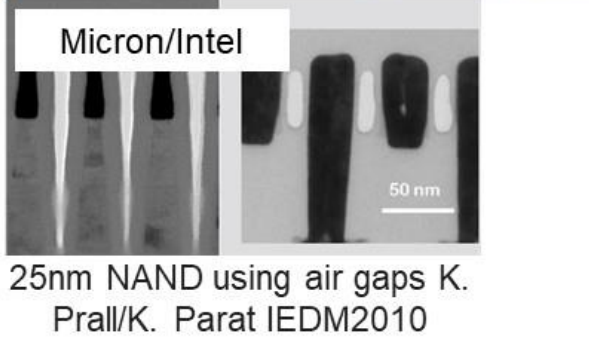
~36-40nm BEOL pitch



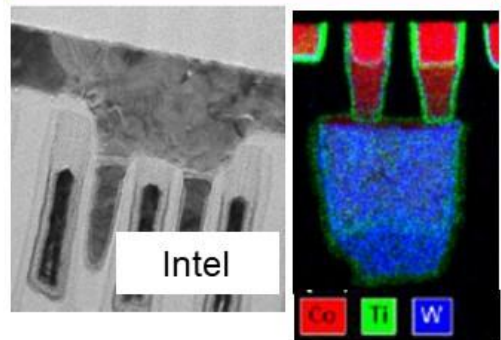
<20nm pitch



AGs in memory & logic



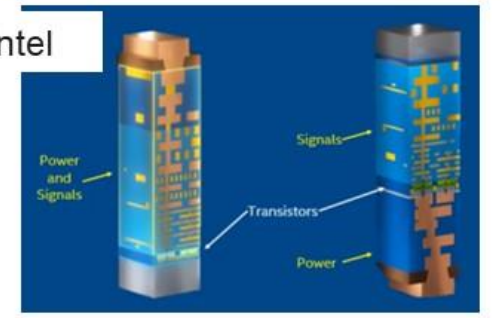
Co interconnects



Self-aligned contact, C. Auth et al. IEDM 2017

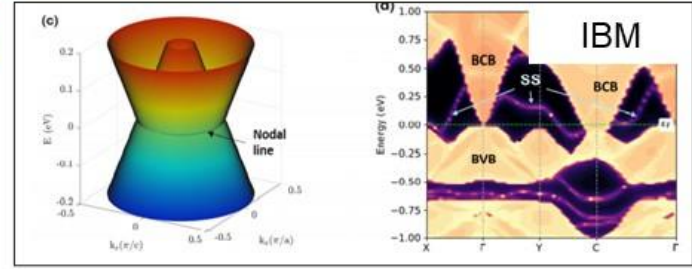
Reliability of Co, F. Griggo et al. IRPS 2018

Backside Power Delivery

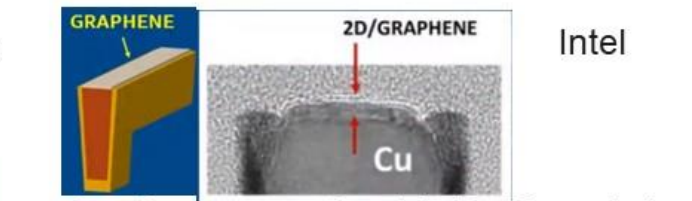


M.C. Mayberry, Keynote IITC 2020

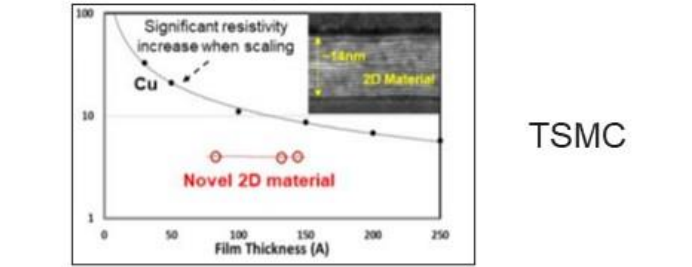
New materials on the horizon



Topological semi-metal, C.T. Chen et al. IEDM2020

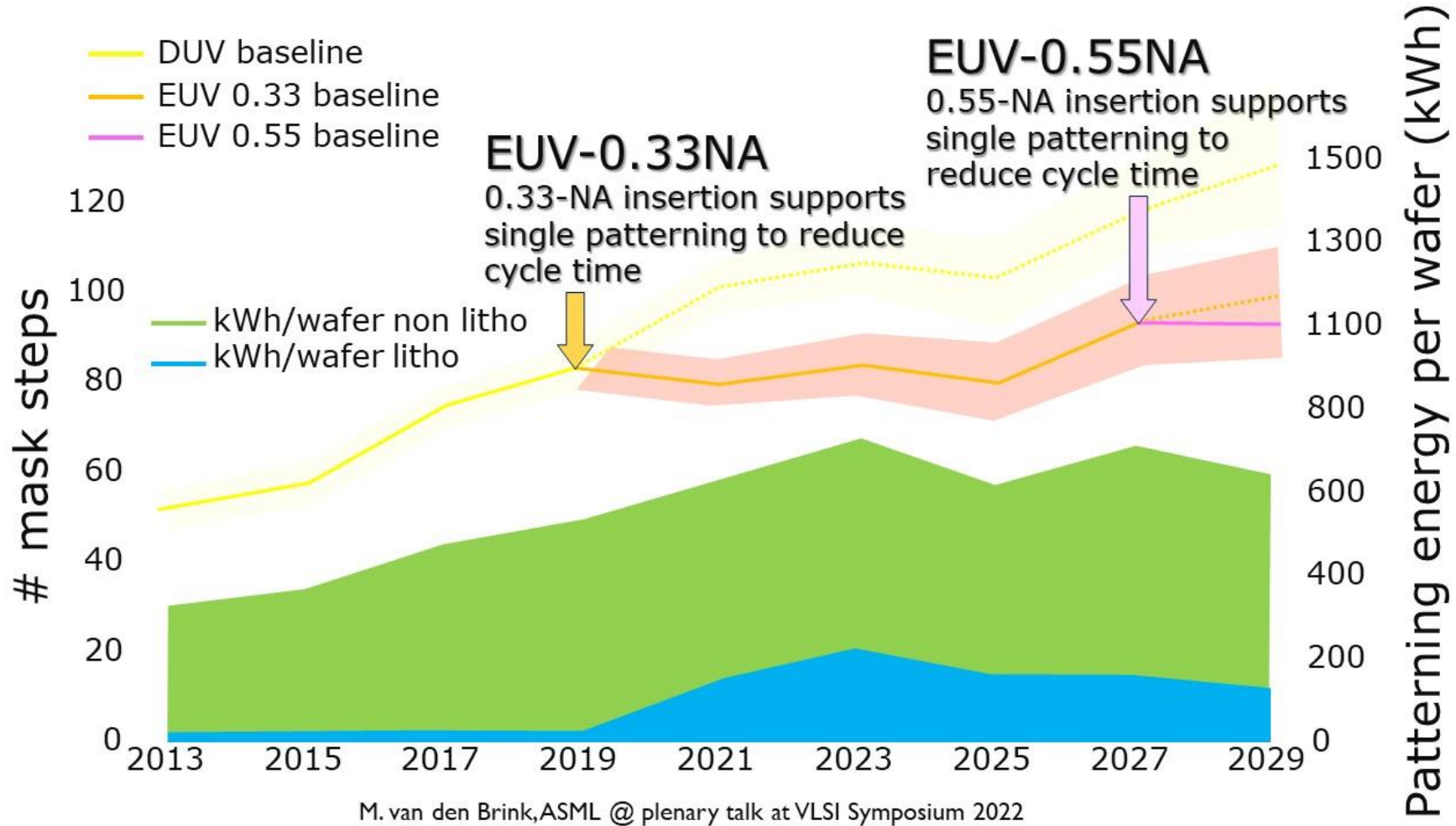


Graphene capped metal, R. Chau et al. Keynote IEDM2019

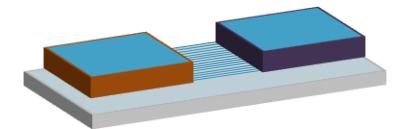
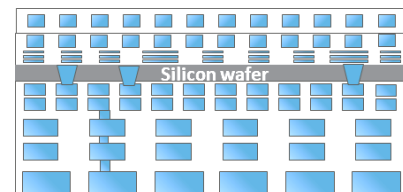
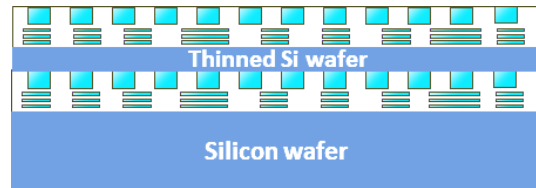
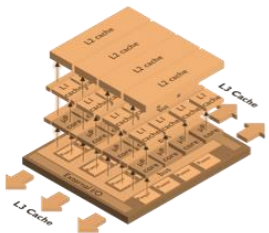
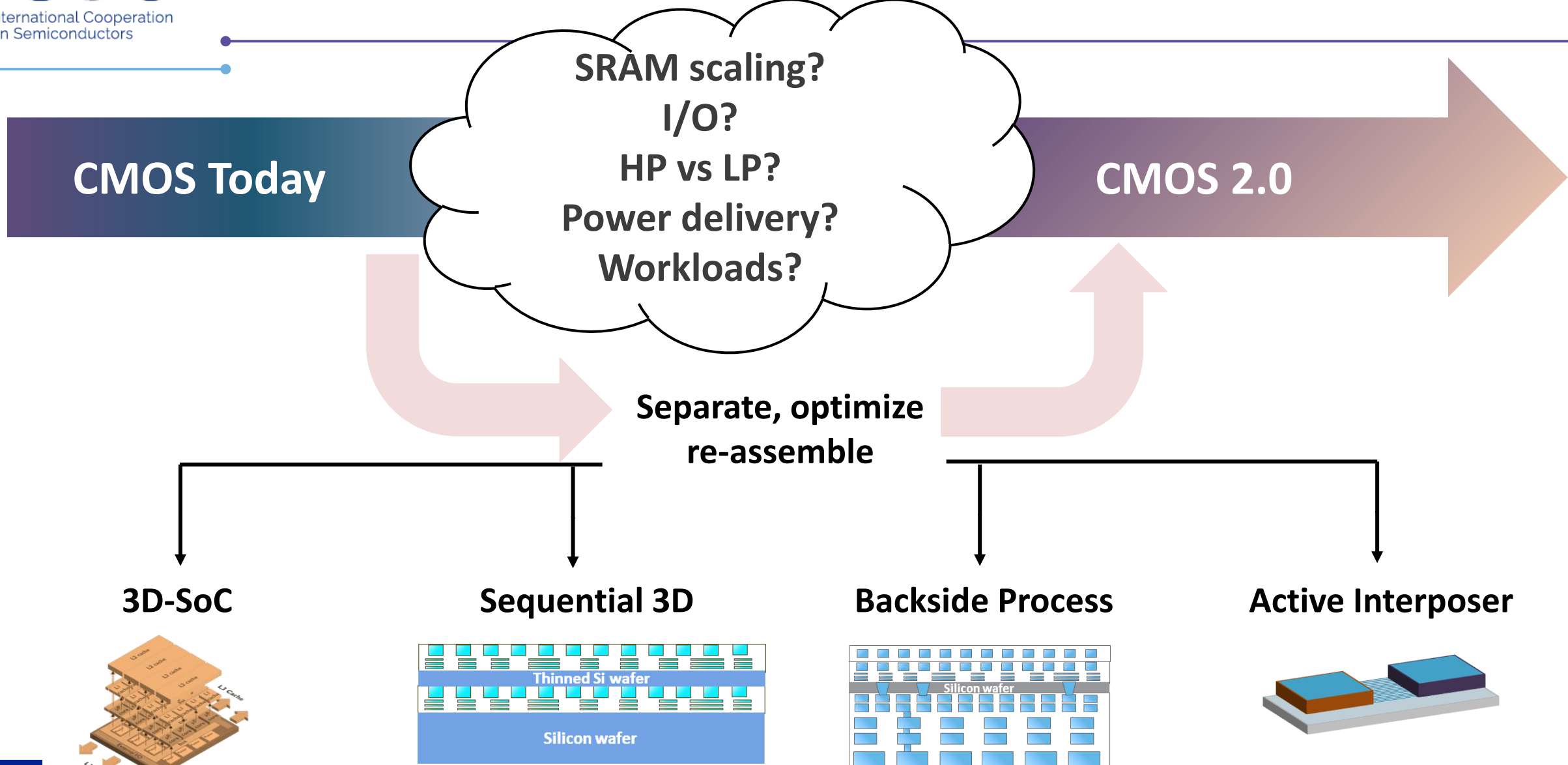


New conductor, Y.J Mii Keynote VLSI2022

EUV lithography key enabler for scaling



Today's Scaling Challenges Drive the Need for CMOS 2.0



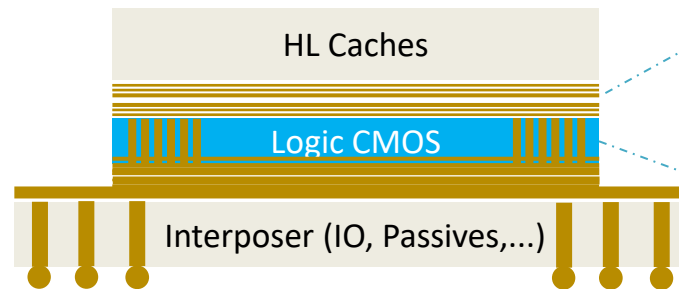
Past

CMOS 1.0
Homogeneous Platform



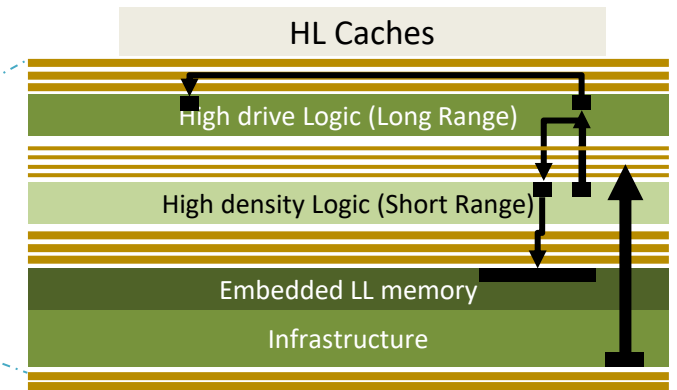
Present

Heterogeneous Systems



Future

CMOS 2.0
Heterogeneous Logic Platform



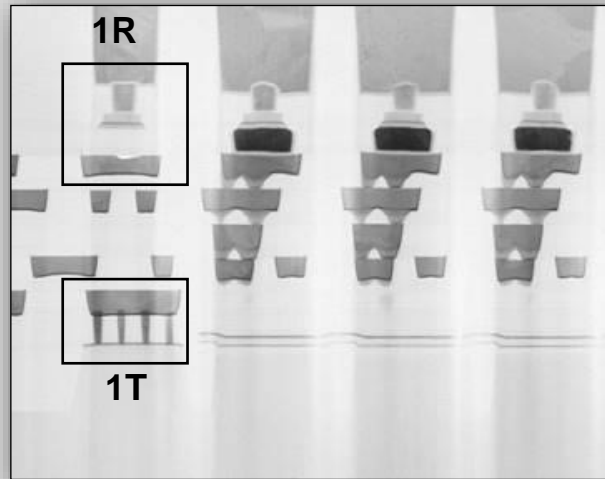
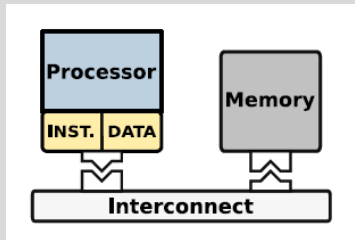
POWER WALL	Multi-core HW accelerators	Backside PDN	Functional Backside
BANDWIDTH WALL	Complex Cache /Memory hierarchy	Active Interposer	Active Interconnect
MEMORY WALL		Hybrid Bonding	Array-Under-CMOS Large embedded memory

Why Emerging Resistive Memories?

High dense on-chip memory

DRAM access is at least **1500x** more costly than a MAC operation in NN accelerators

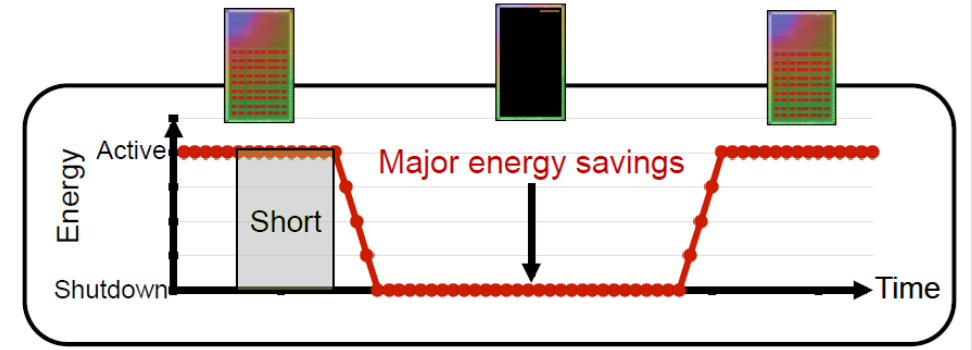
[F. Tu, et al., 2018 ACM/IEEE]



L. Grenouillet et al., 2021



Zero stand-by power thanks to non-volatility



10x better energy efficiency than embedded flash thanks to resistive memories



T. Wu et al., 2019

Emerging Non-Volatile Memories

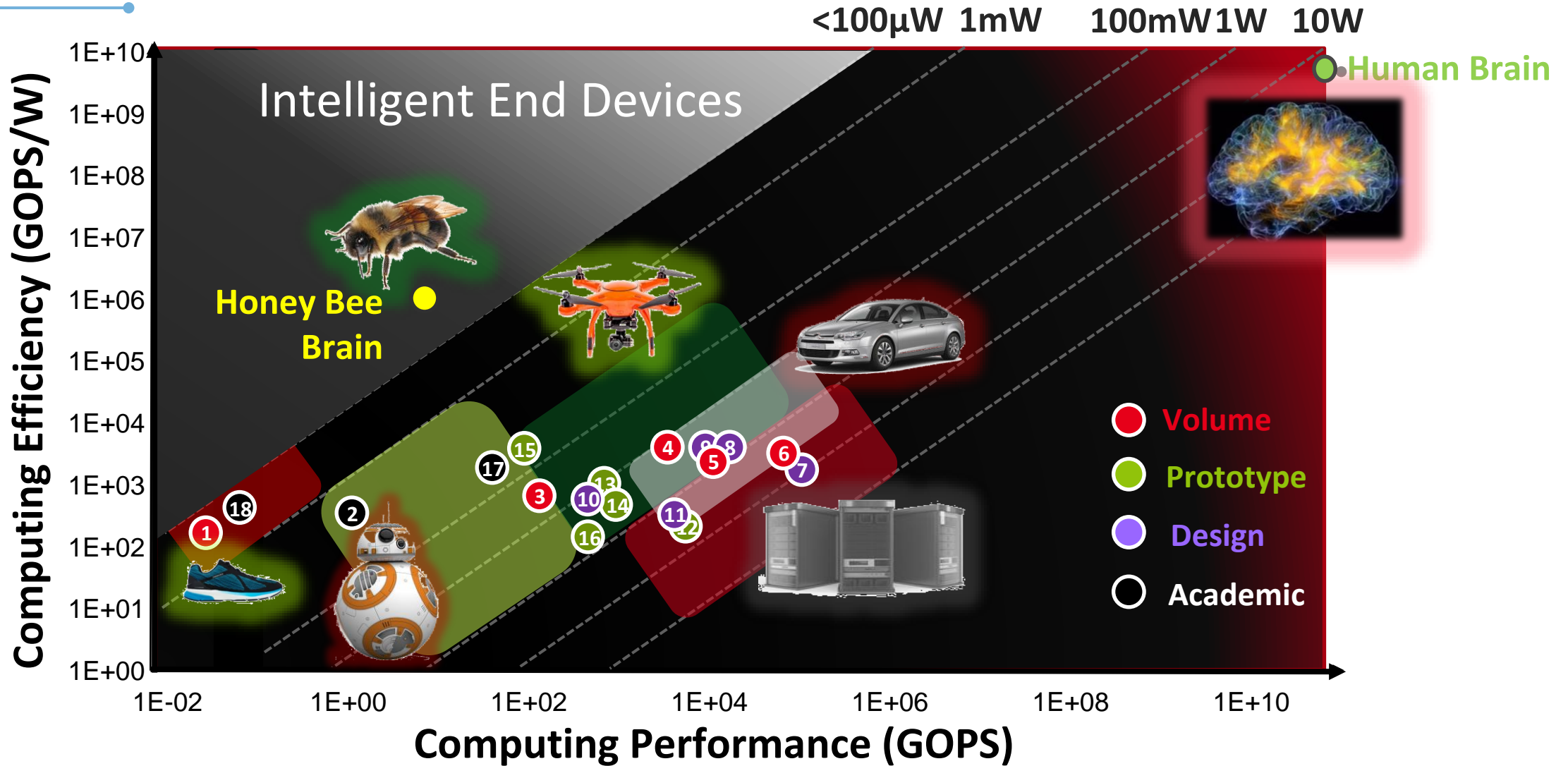
	NOR FLASH	MRAM	PCRAM	OxRAM	FeRAM (PZT)	FeRAM (HfO ₂)
Programming power	~200pJ/bit	~20pJ/bit	~300pJ/bit	~100pJ/bit	~10fJ/bit	~10fJ/bit
Write speed	20 μs	20 ns	10-100 ns	10-100 ns	<100ns	14ns @ 2.5V (SONY) 4ns @ 4.8V (LETI)
Endurance	10 ⁵ - 10 ⁶	10⁶-10¹⁵	10 ⁸	10 ⁵ – 10 ⁶ on 16 kbit	> 10¹⁵	> 10¹¹ single device 10⁶ – 10⁷ on 16 kbit
Retention	> 125°C	85°C - 165 °C	165°C	> 150°C	125°C	125°C
Extra masks	Very high (>10)	Limited (3-5)	Limited (3-5)	Low (2)	Low (2)	Low (2)
Process flow	Complex	Medium	Medium	Simple	Simple	Simple
Multi-Level Cell	Yes	No	Yes	Yes	No	No
Scalability	Bad	Medium	High	High	Medium	Poor (2D) High (3D)

Power Reduction by 10000!



Memory activity focus on embedded NVM for NOR flash replacement

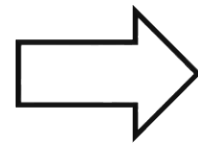
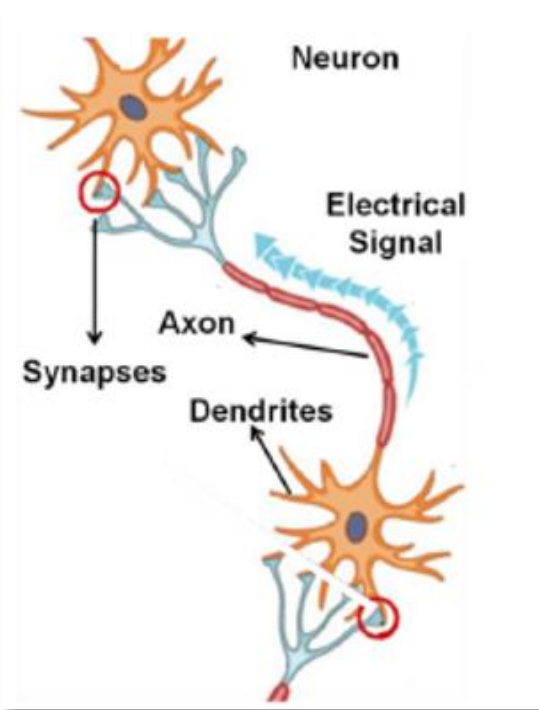
Energy Efficiency is far from biological Systems



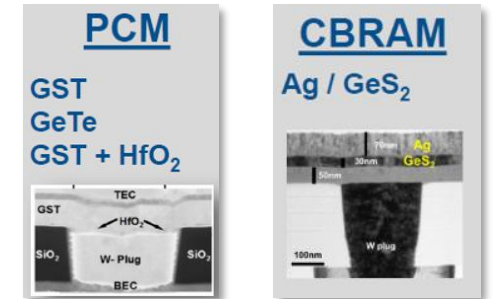
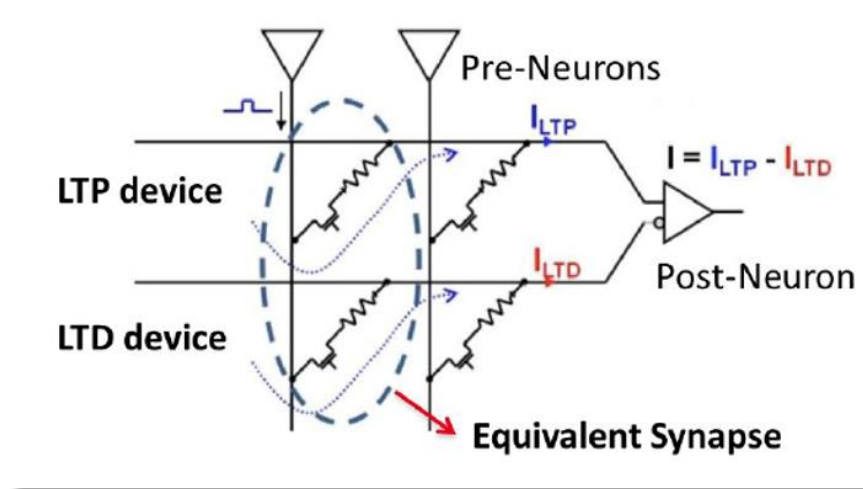
Computing Performance (GOPS)

Neuromorphic based RRAM circuit

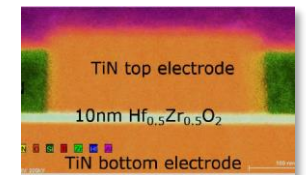
M. Suri et al, IEDM 2011.



2 PCRAM Example:



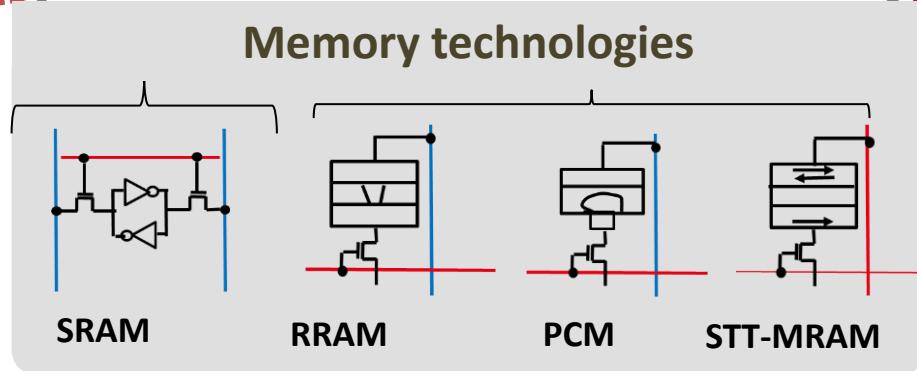
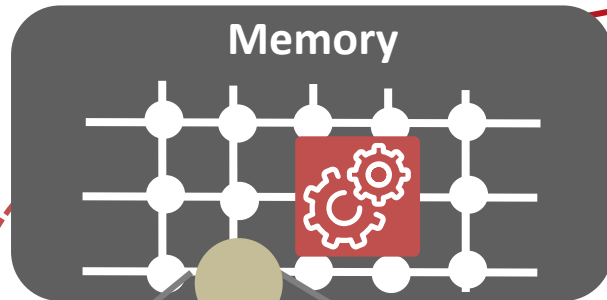
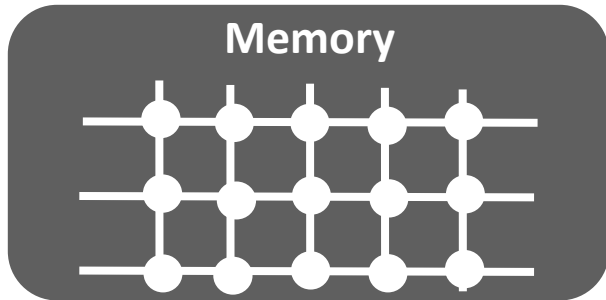
FeRAM



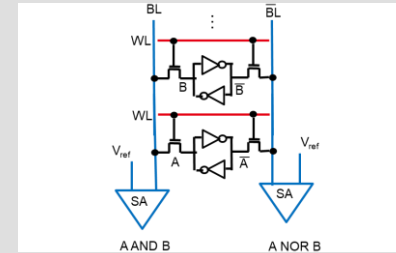
In memory computing

Von Neumann computing

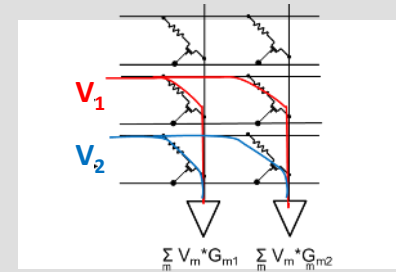
In-memory computing



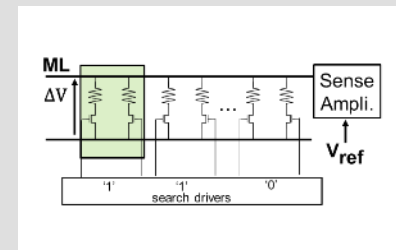
Logic and arithmetic operations



Digital operations using SRAM

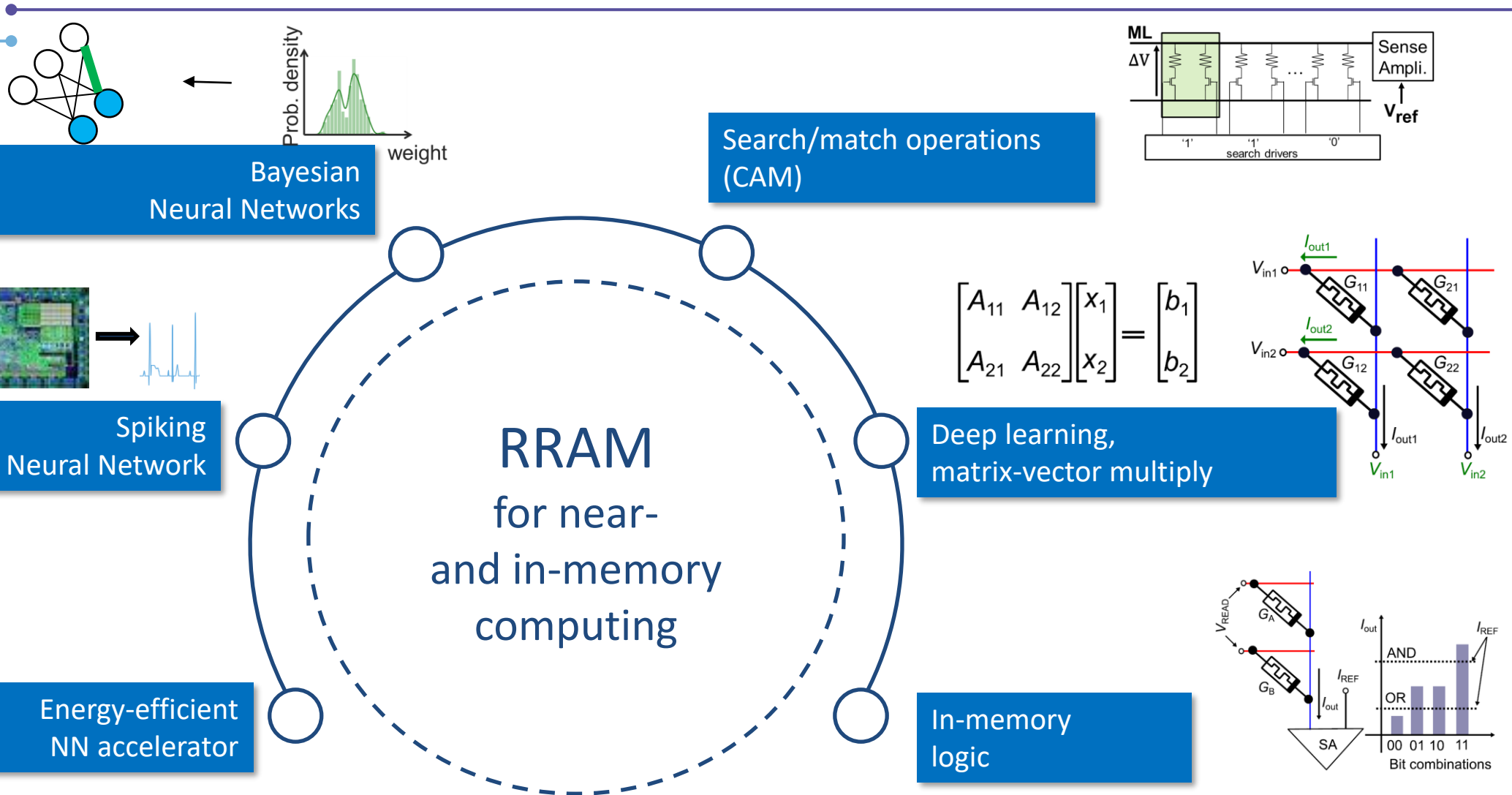


Analog: multiply and accumulate

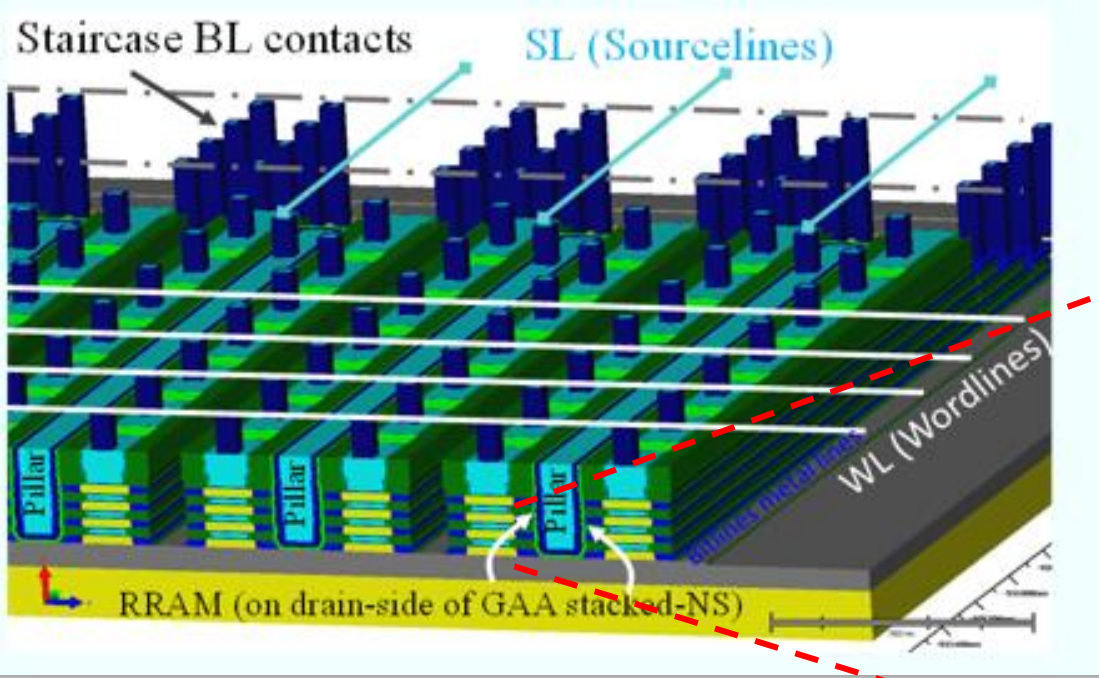


Searching / matching (CAM)

Near- & in-memory computing

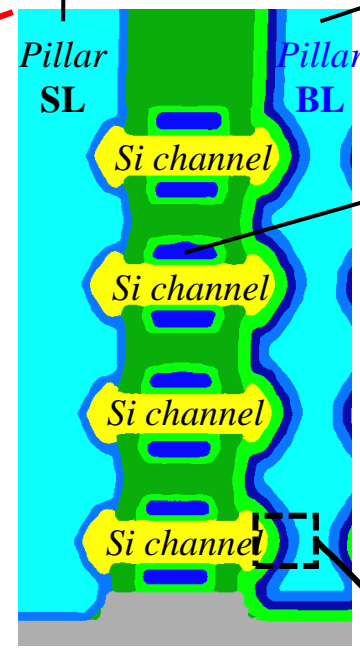


Towards In Memory Computing



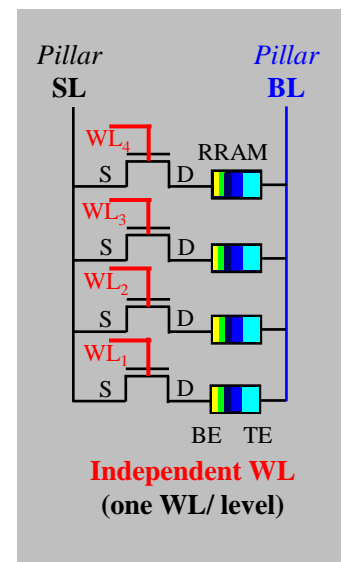
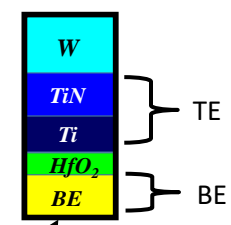
Common SL for all GAA stacked-NS

Common BL for all GAA stacked-NS



Courtesy of S. Barraud, Leti

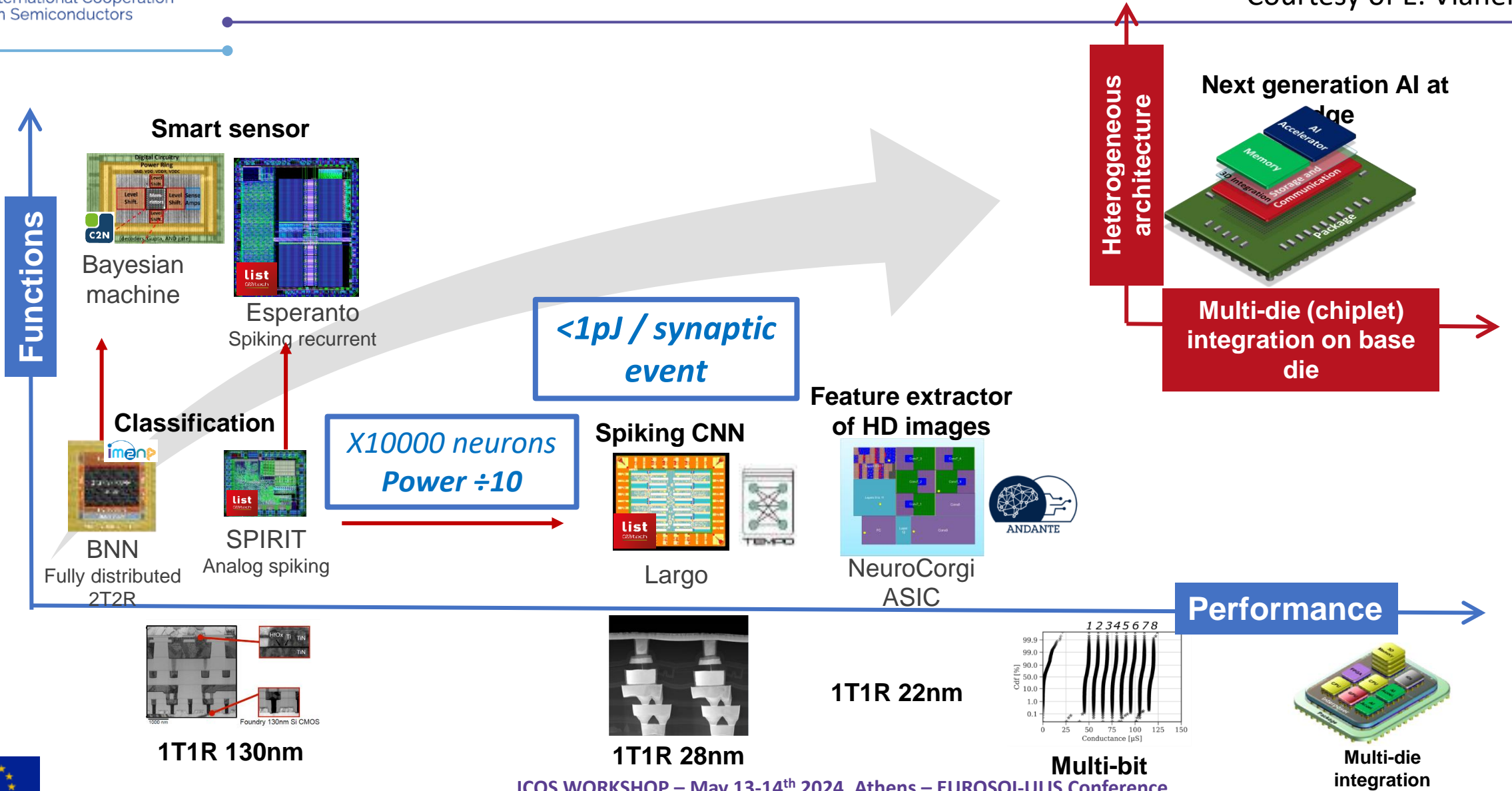
Independent WL
One GAA by level



Enabler for Hyperdimensional computing

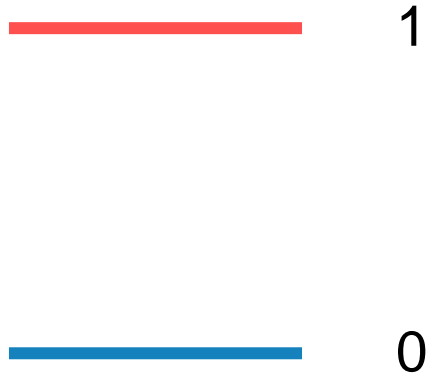
A possible Edge IA Roadmap

Courtesy of E. Vianello, Leti

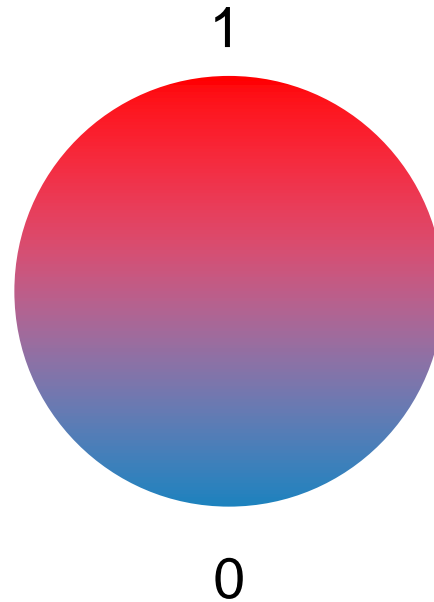


From bits to qubits

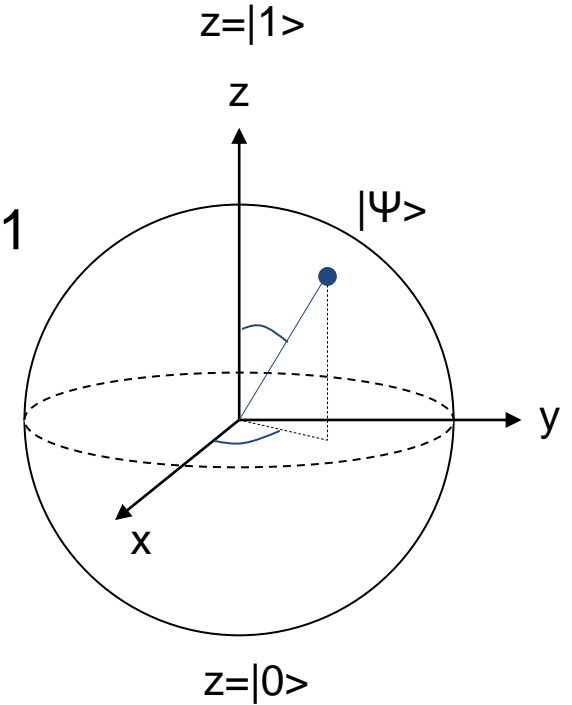
Bit



Qubit



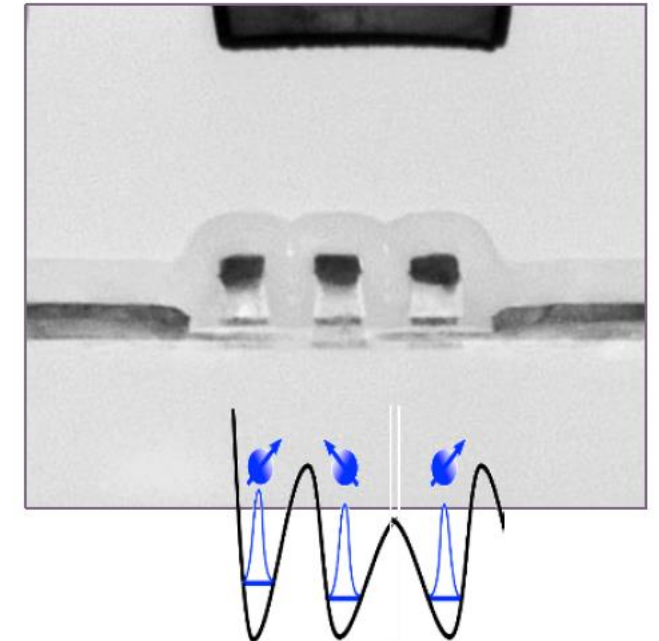
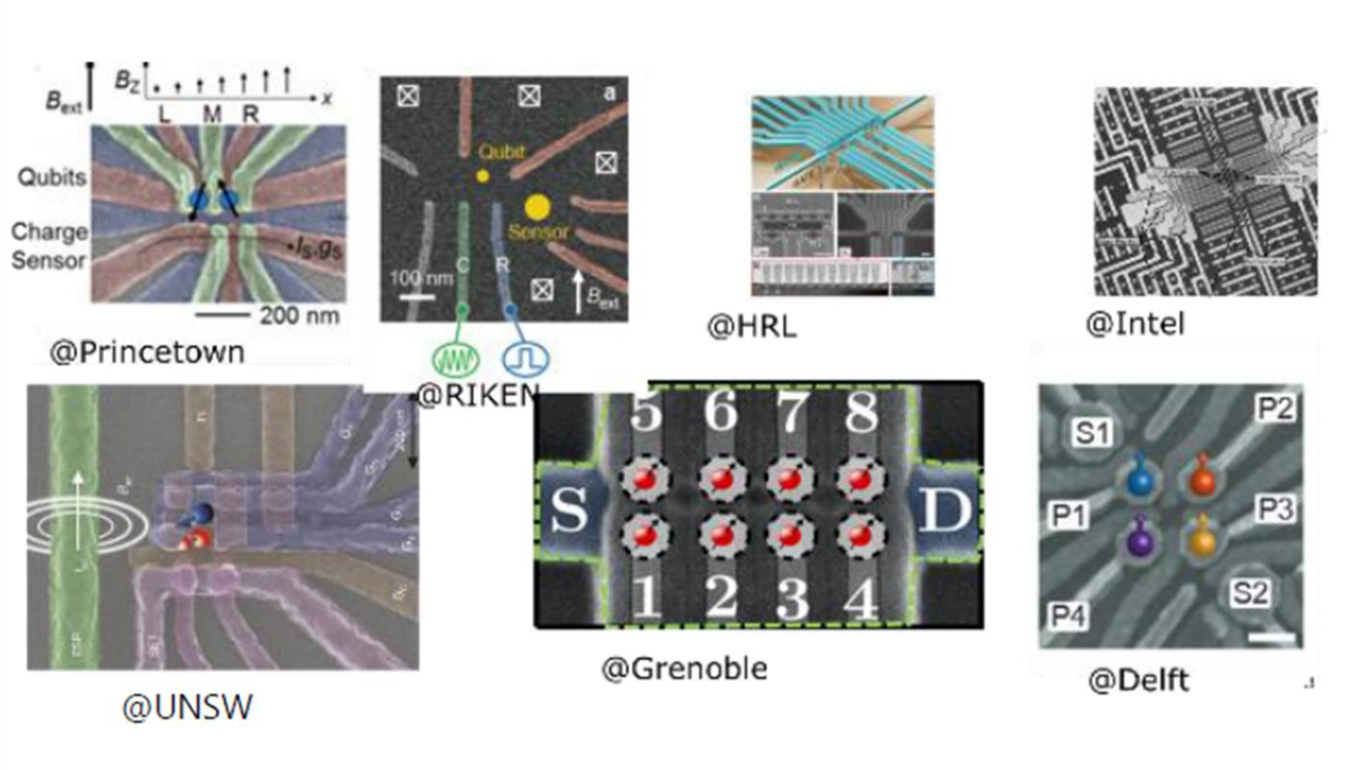
0 AND 1



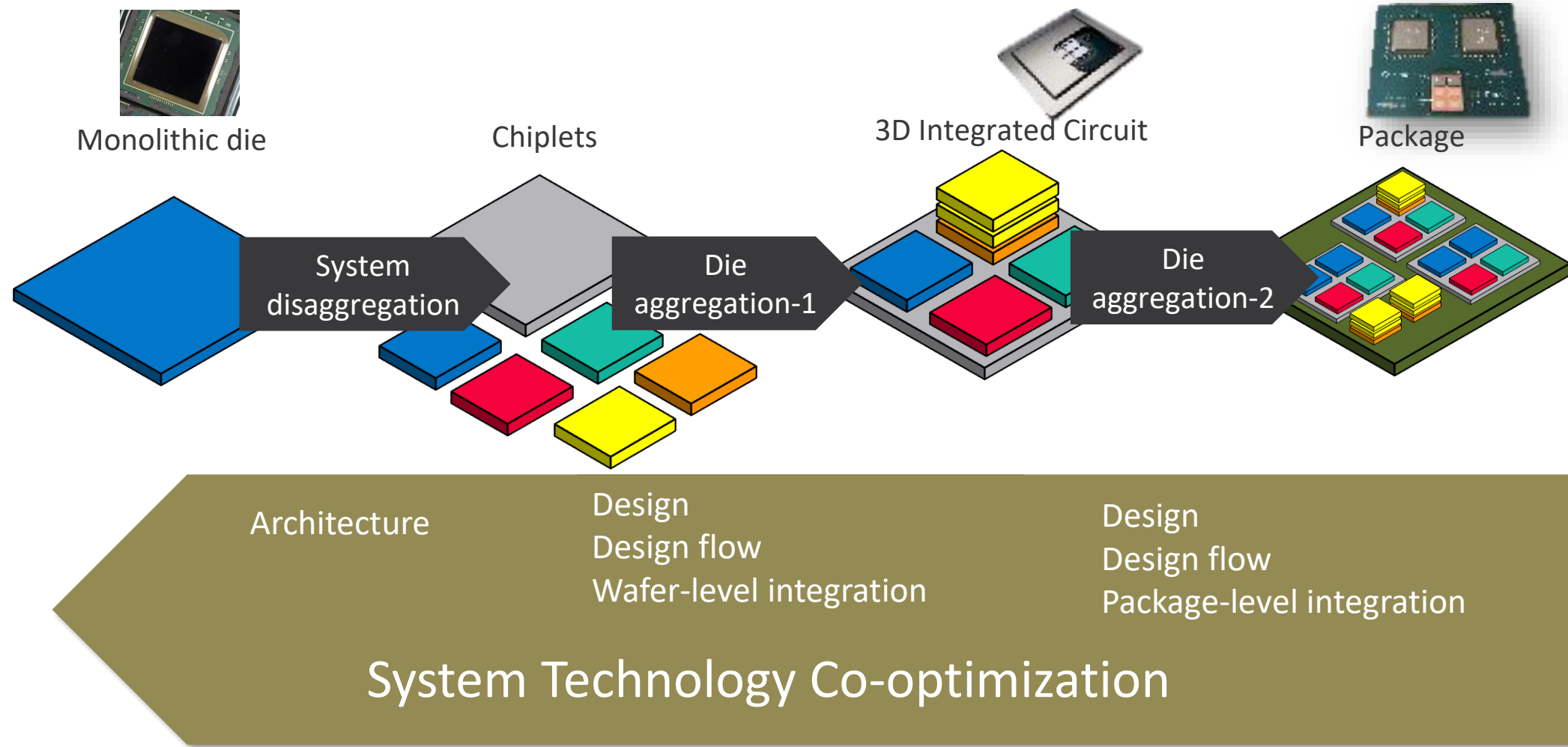
Many flavors of qubits

	Superconductor	Si spin qubit	Trapped ion	Photon
Size*	(100 μ m) ²	(100nm) ²	(1mm) ²	~(100 μ m) ²
1qubit fidelity	99.96%	99.93%	99.98%	
2qubit fidelity	~99.3%	>99%	99.9%	50% (measurement) 98% (gates)
Speed**	12-400 ns	~1 μ s	100 μ s	1 ms
Variability	3%	0.1%-0.5%	0.01%	0.5%
T° of operation	15mK	1K	10K	4K/10K
Entangled qubits	433 (IBM)	3 (TU) (6 - QuTech)	32 (IonQ)	70 (Pan-China)

- Quantum Physics to compute

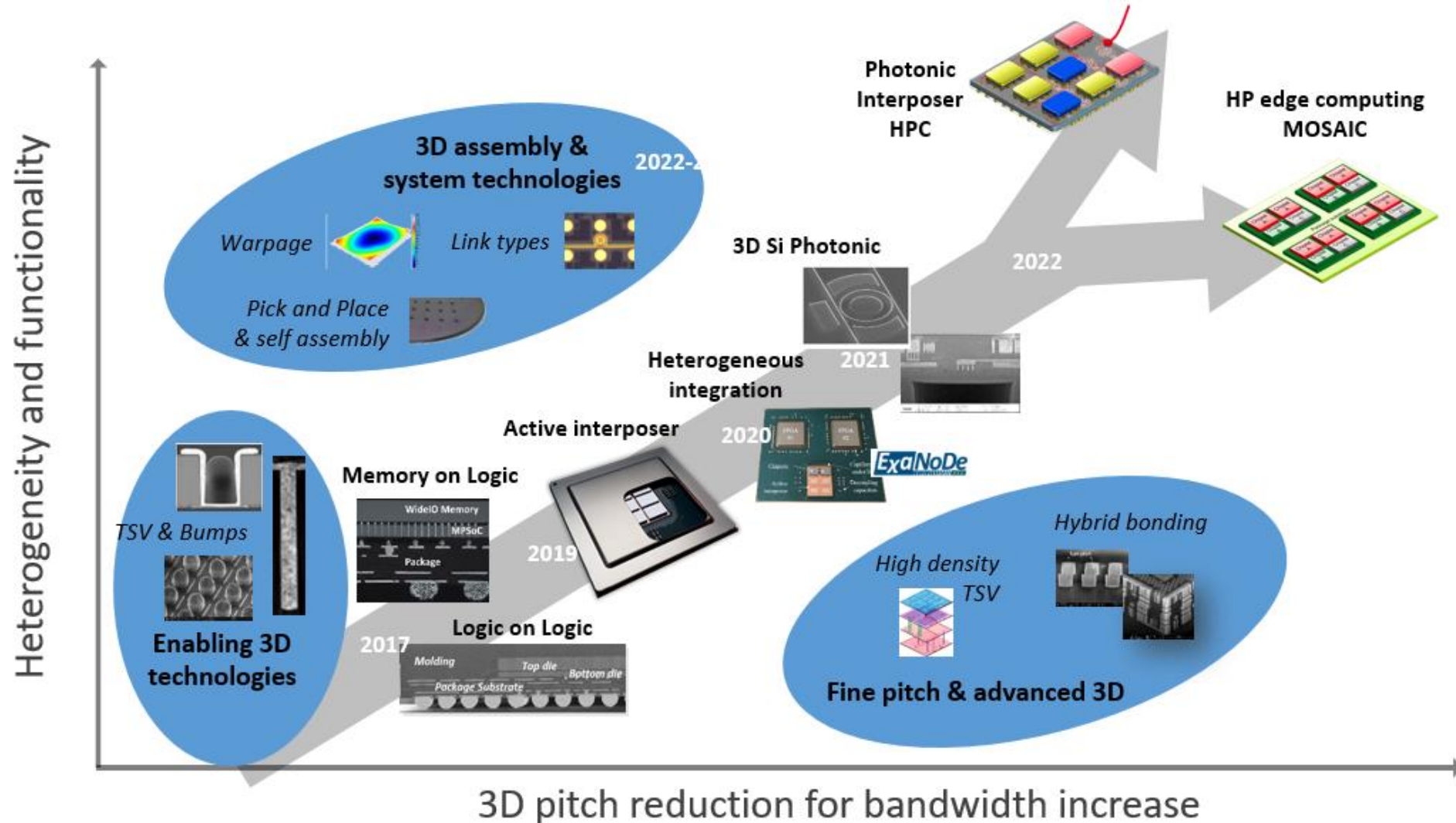


Chiptlets: the new IC design paradigm

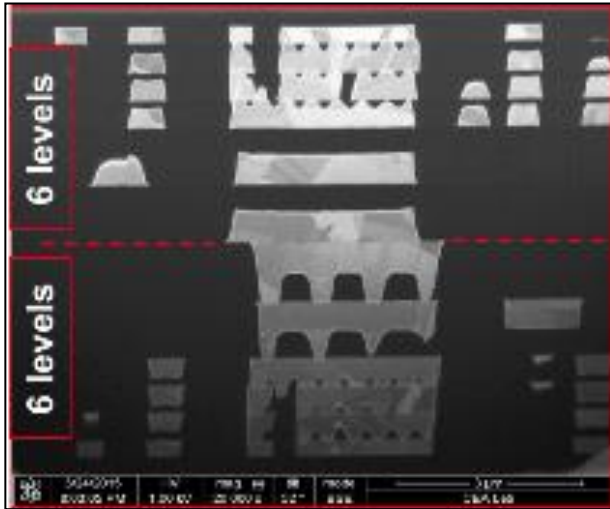


Up to 100x gain on Power Efficiency with 3D

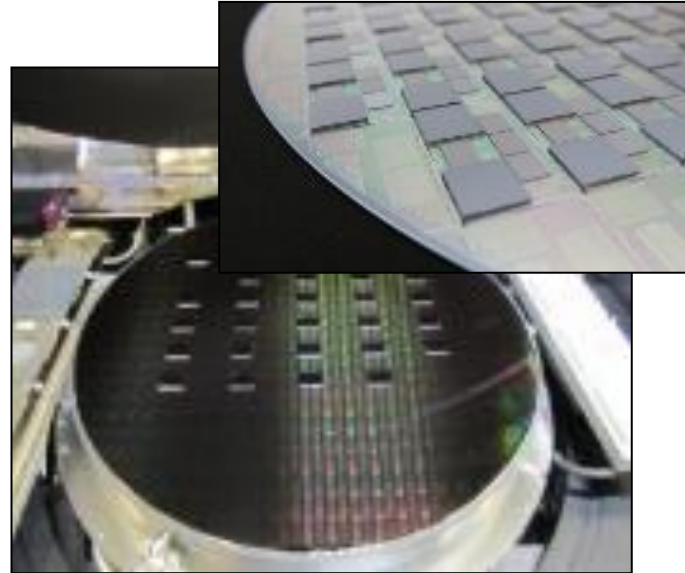
3D Tool Box for chiplet integration



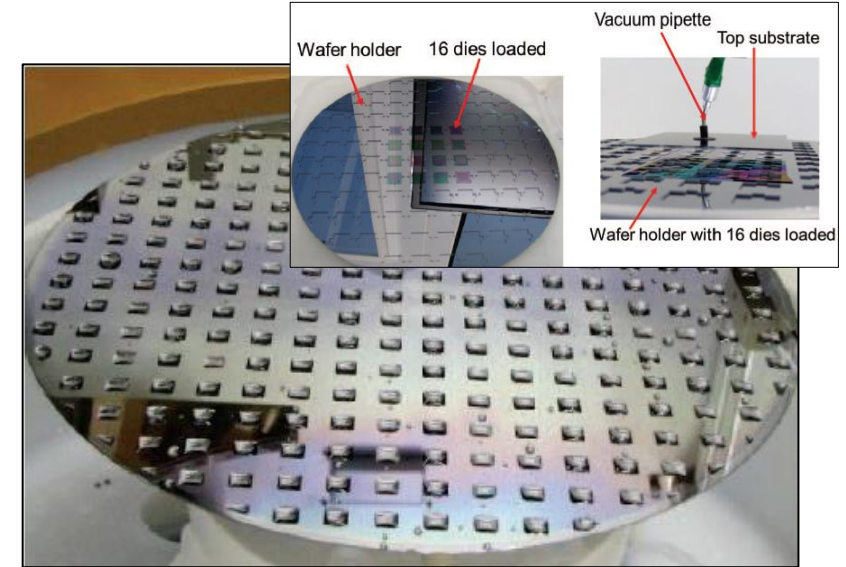
Hybrid Bonding Solutions



- › Direct bonding of metal and dielectric
- › Down to 1 micron pitch interconnects

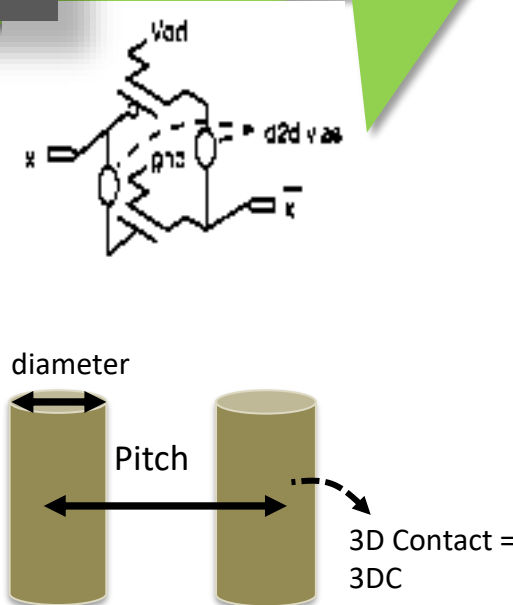
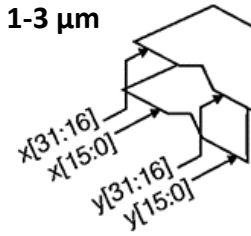
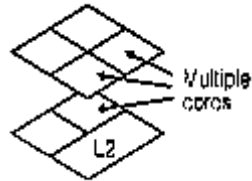
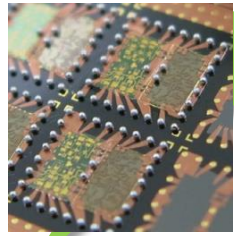
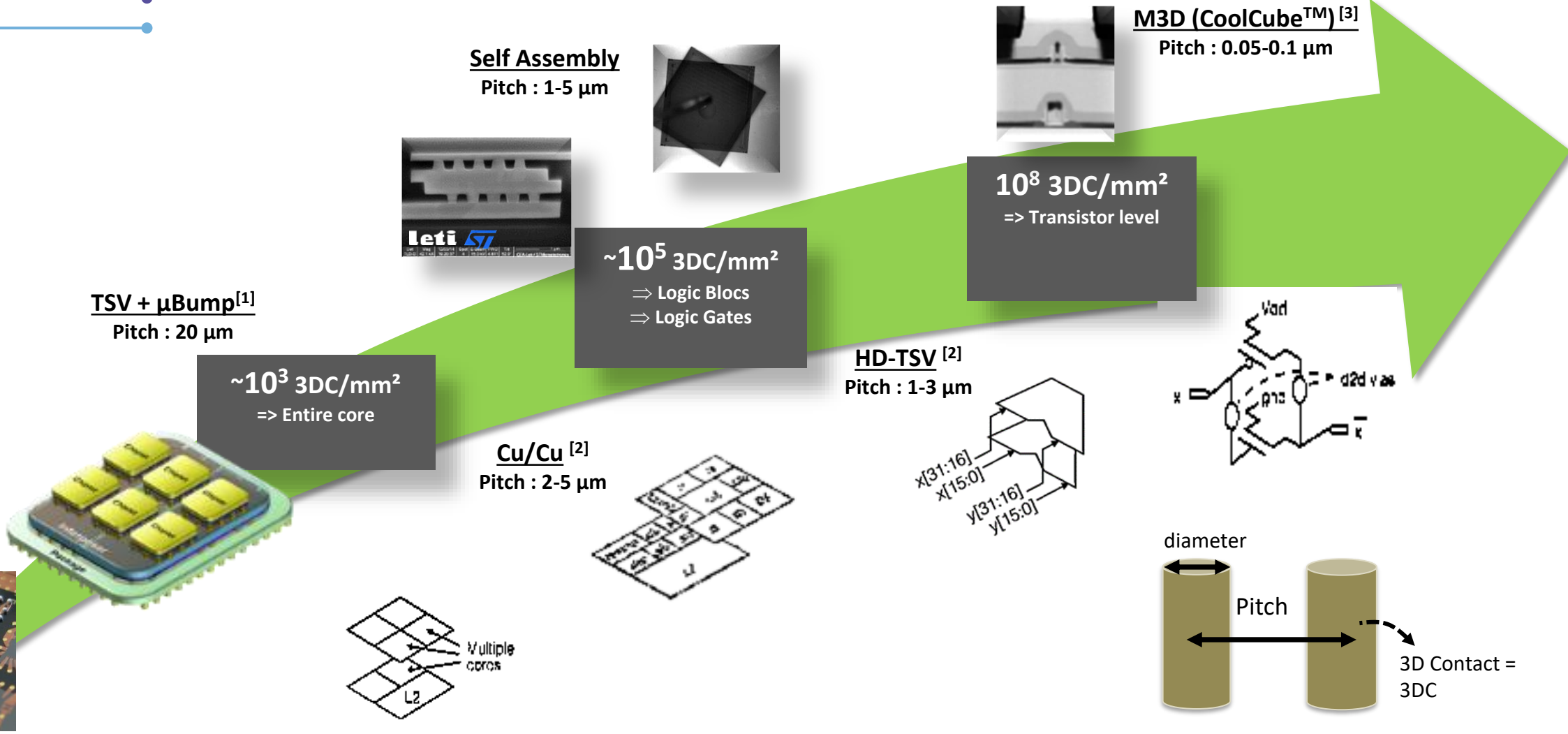


- › Wafer-to-wafer (W2W) or Die-to-wafer (D2W) technologies
- › High heterogeneity allowed by D2W

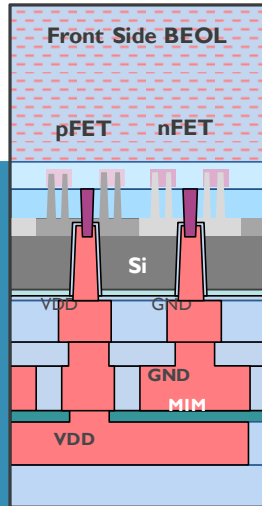


- › Collective D2W approaches
- › Self-assembly for high precision & high throughput

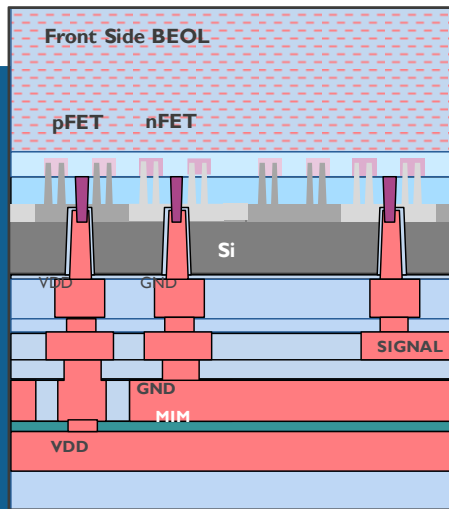
3D: from packaging to monolithic



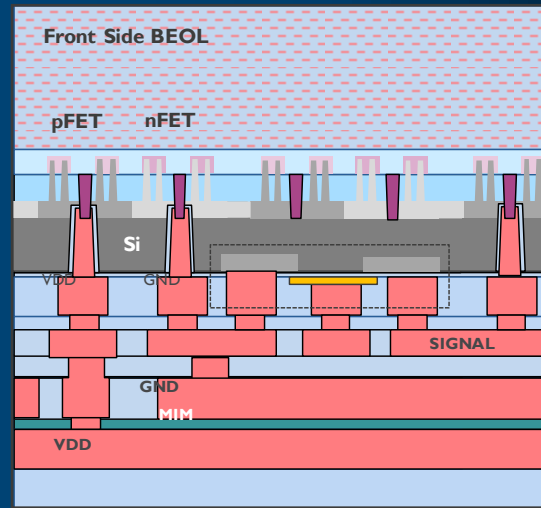
Backside Power Delivery



Backside Global Interconnect

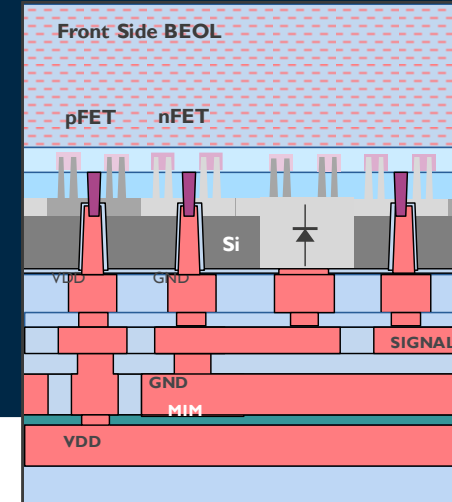


Backside Devices



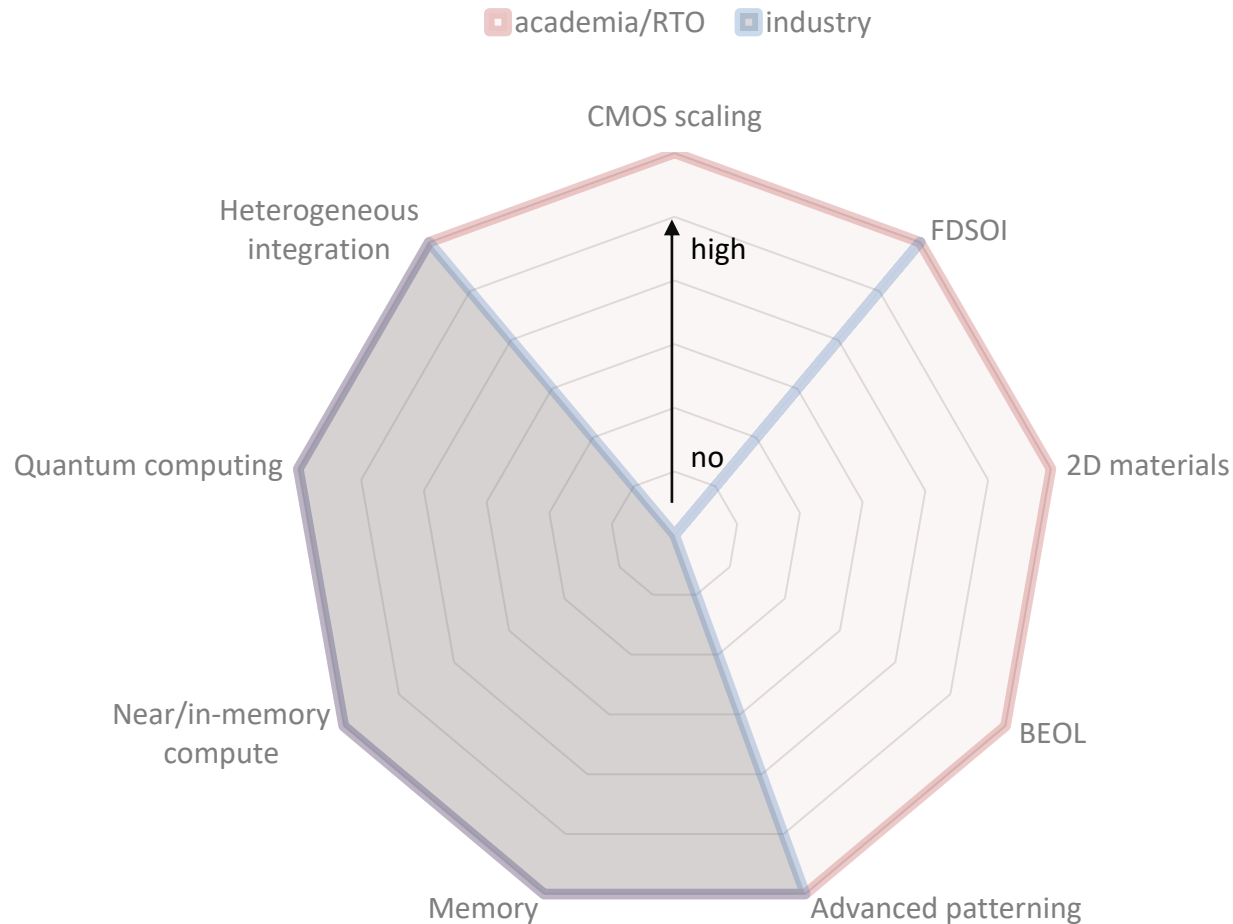
J. Ryckaert, ITF Japan 2022

Device Backside Extension



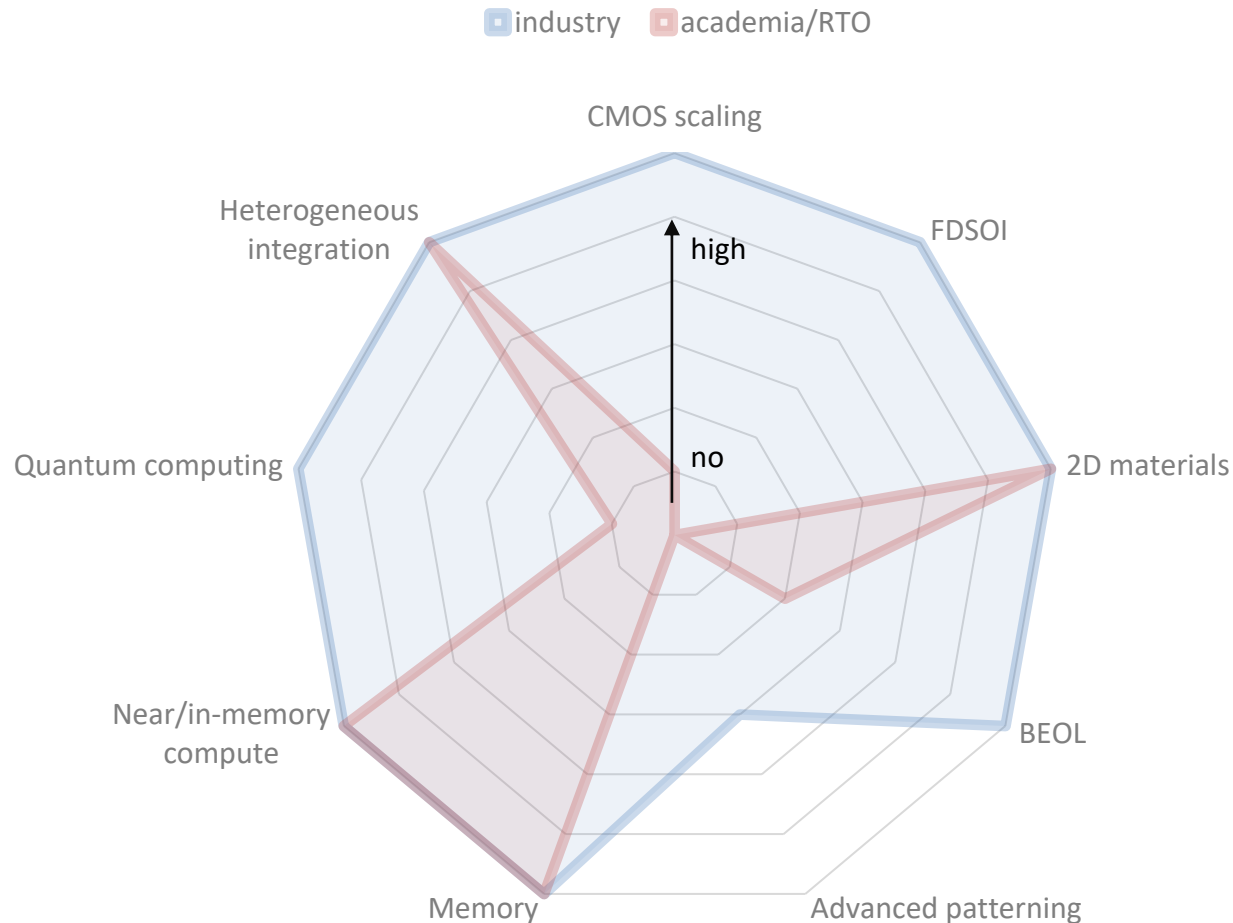
Enhancing system performance by migrating system functions to the backside

EU and non-EU actors - EU



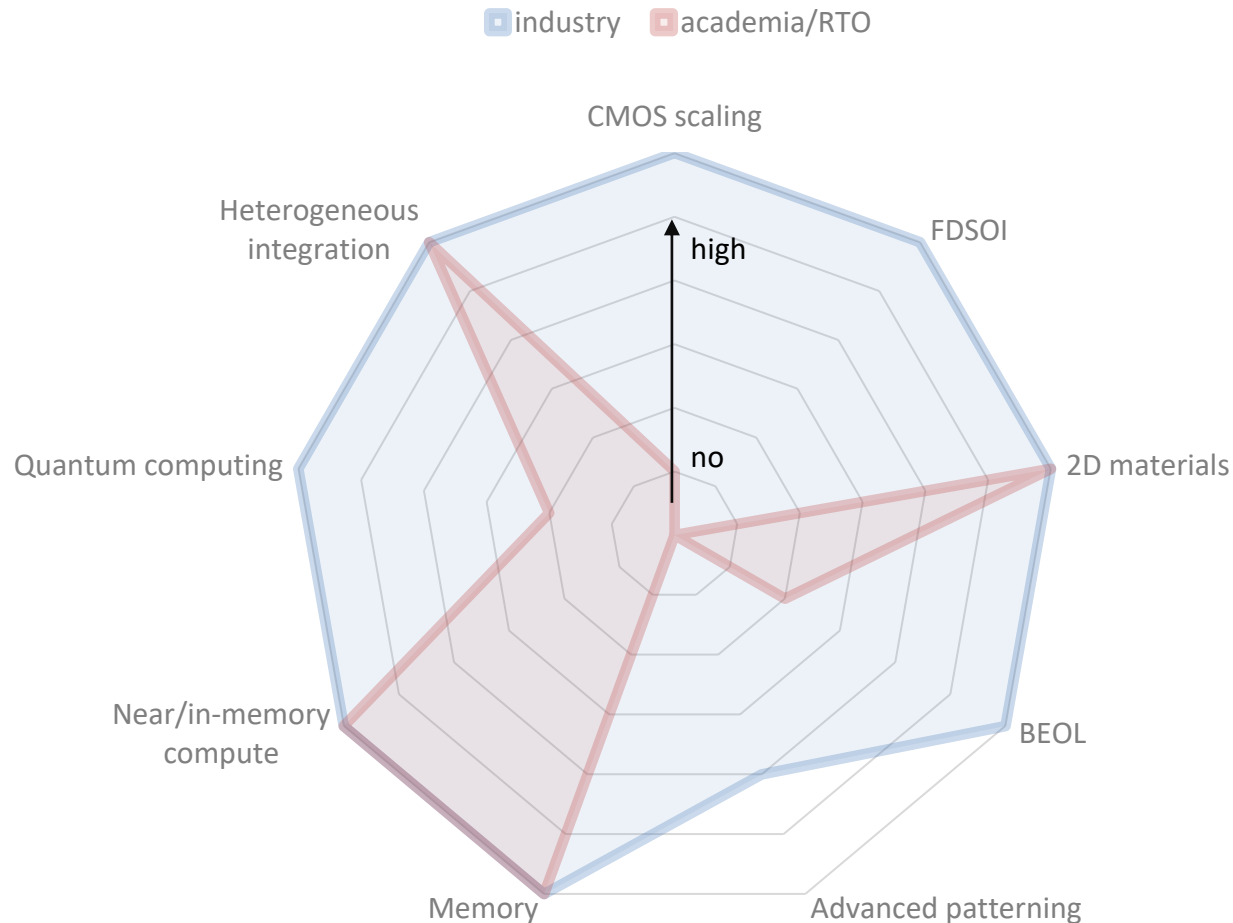
- R&D very strong in all areas of compute
- Unique strong position in EUV lithography
- In general, industrial EU players lacking to take up R&D

EU and non-EU actors - US



- Strong industrial activity in most areas of compute
- Weaker academic activity on traditional logic scaling
- Strong R&D in new materials, heterogeneous integration and memory

EU and non-EU actors - Asia



- Very similar to US
- Strong industrial activity in most areas of compute
- Weaker academic activity on traditional logic scaling
- Strong R&D in new materials, heterogeneous integration and memory

Summary

- Compute needs are growing at an unprecedented speed
- Innovations across different levels needed to enable chip & system performance enhancement
- Sustainability becoming an increasingly important metric for evaluating technology choices
- Europe is very strong in R&D in all advanced compute areas
- Manufacturing: several industry initiatives & pilot line programs
- Lack of Fabless companies that drive the needs in terms of advanced compute
- Strong initiatives in EU in start-up company creation



THANK YOU!

This project has received funding from the European Union's Horizon Europe research and innovation programme under GA N° 101092562

ICOS WORKSHOP – May 13-14th 2024, Athens

icos-semiconductors.eu