

Energy Efficient Computing in Nanoscale CMOS

Vivek De
Intel Fellow
Director of Circuit Technology Research
Intel Labs

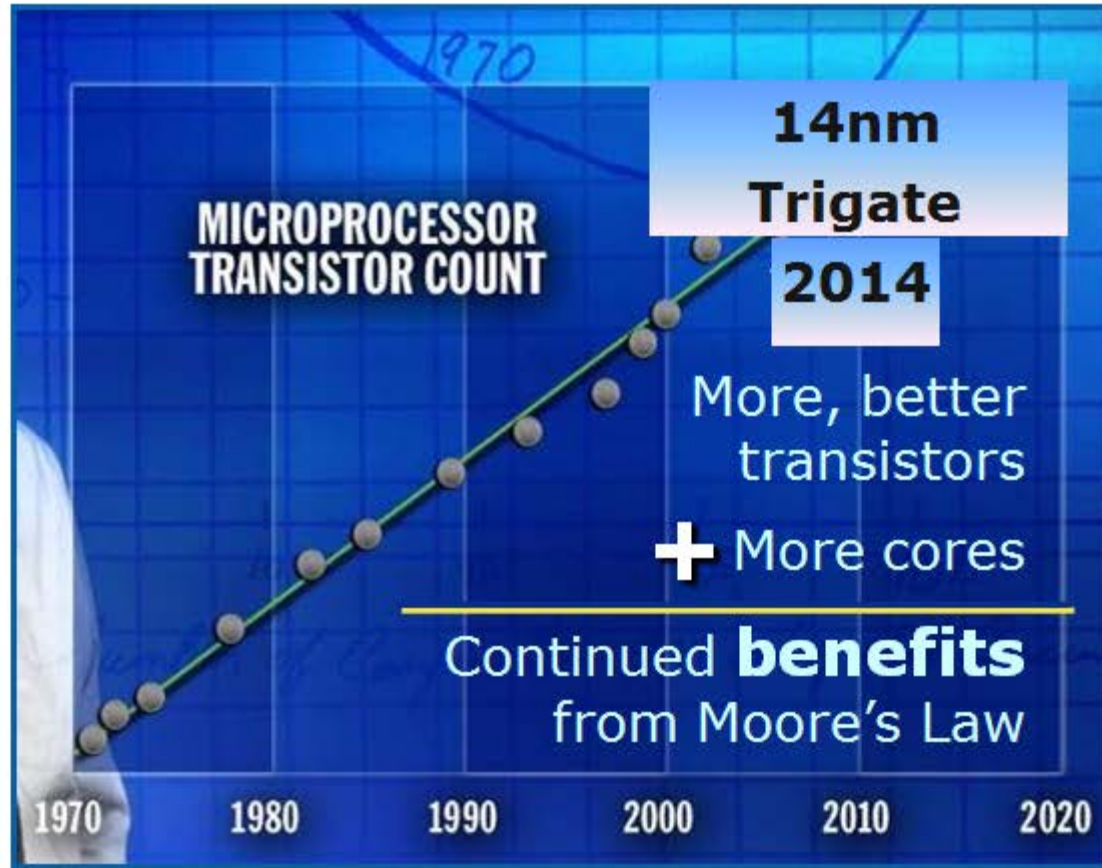


Internet of Everything (IoE)

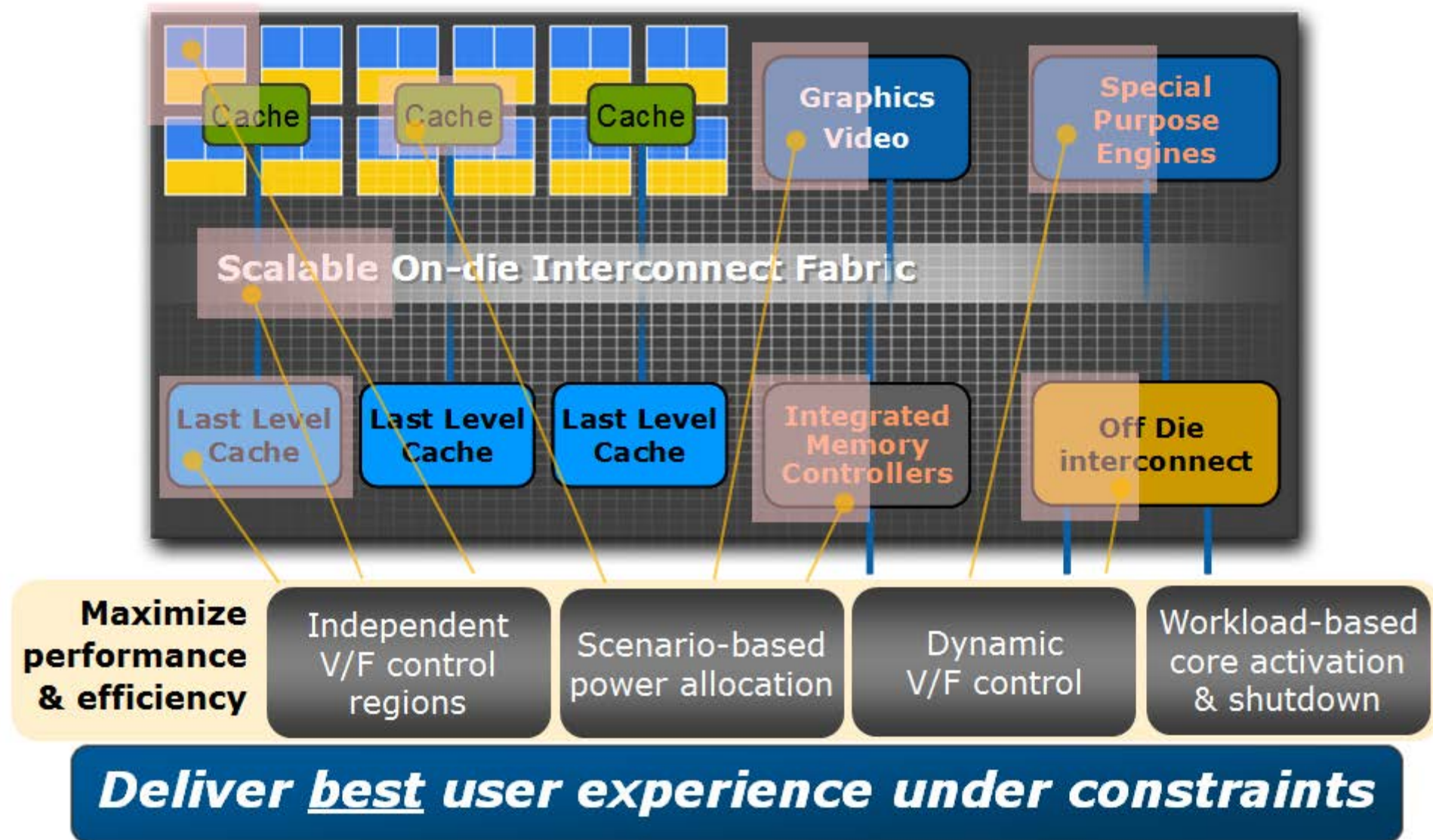


Paradigm shifts across hardware and software

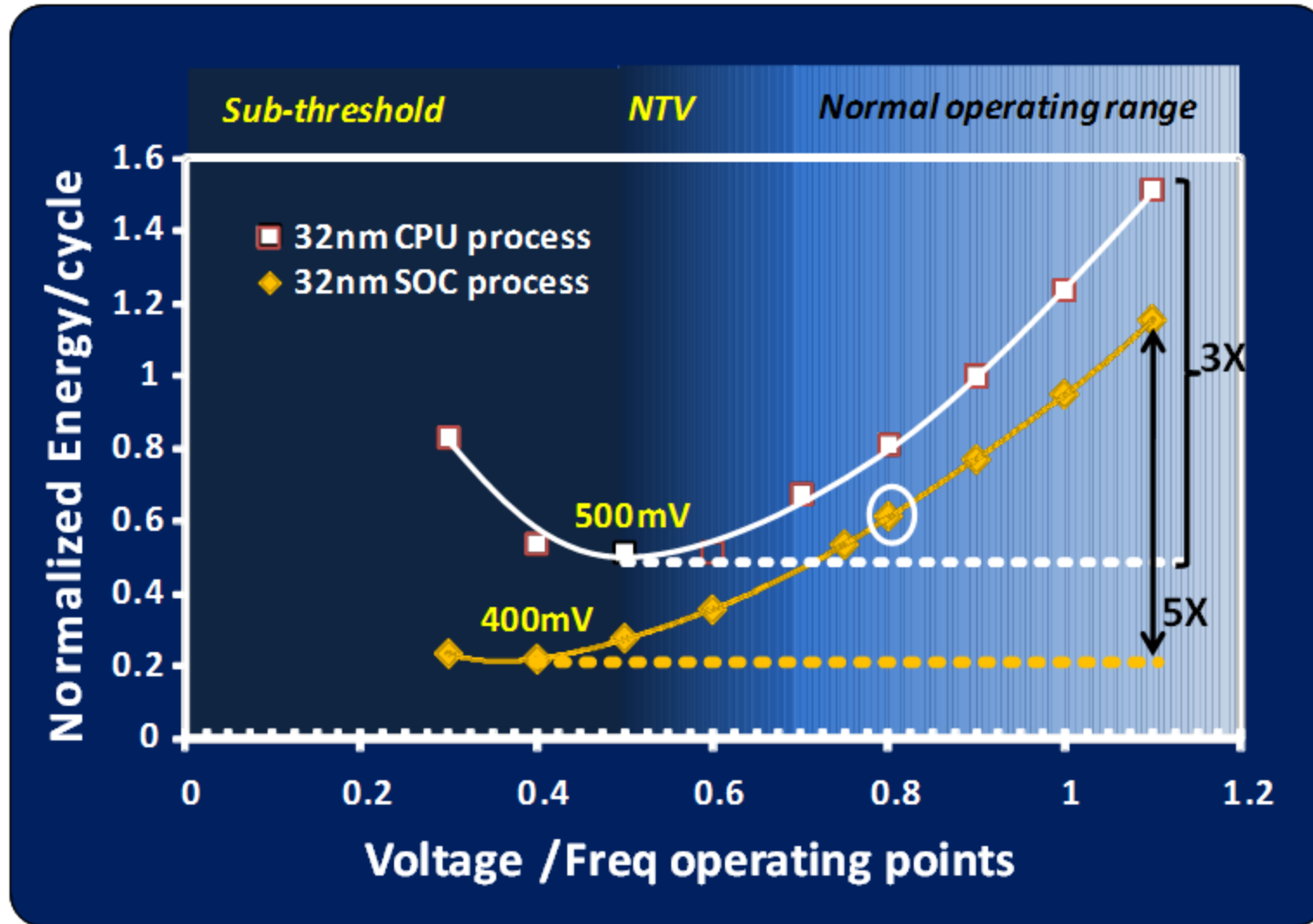
Moore's Law scaling



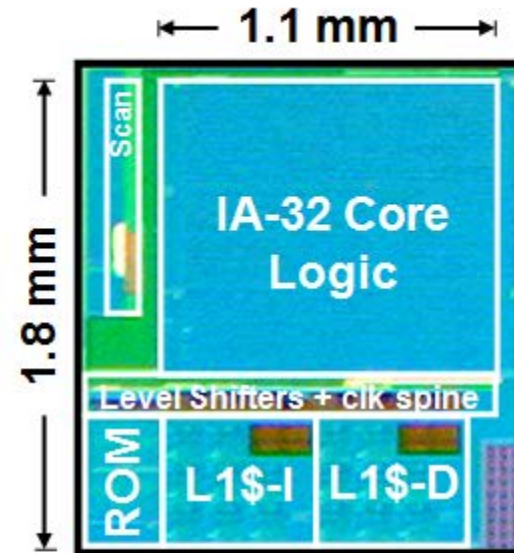
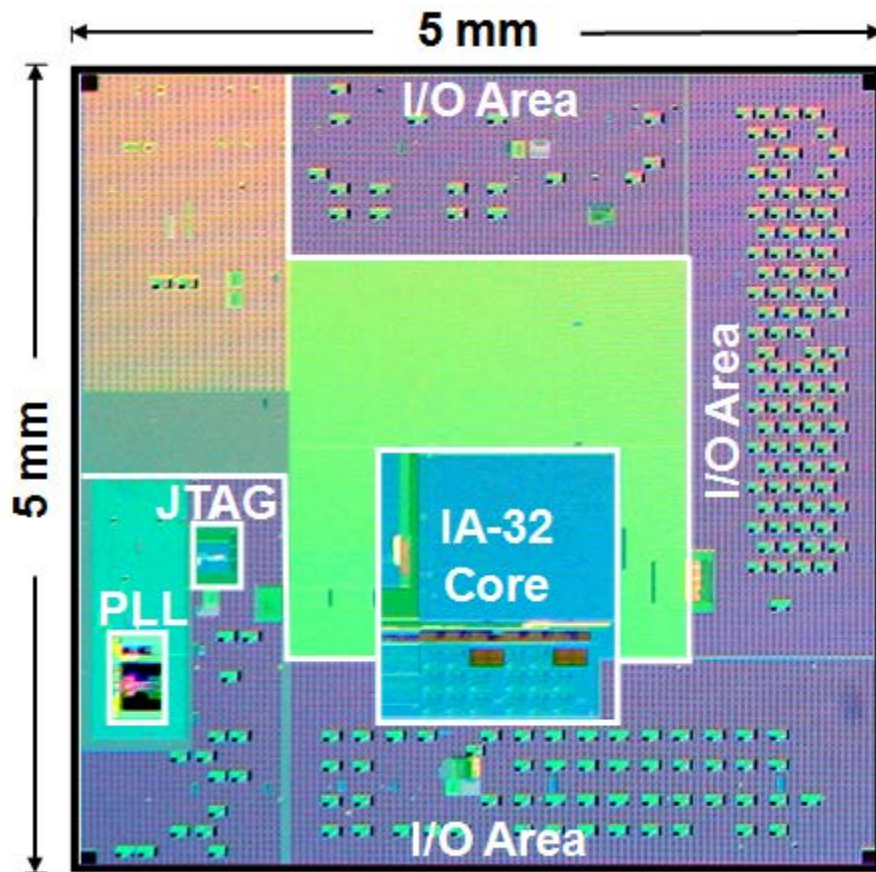
Dynamic platform control



Near Threshold Voltage (NTV) computing



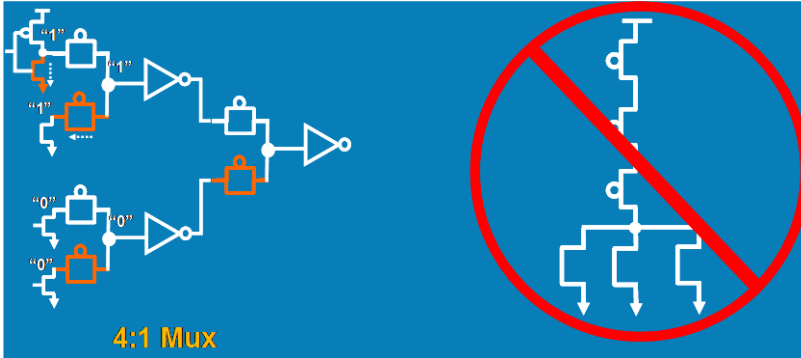
NTV IA processor



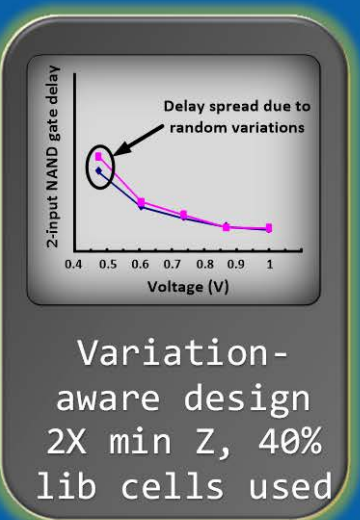
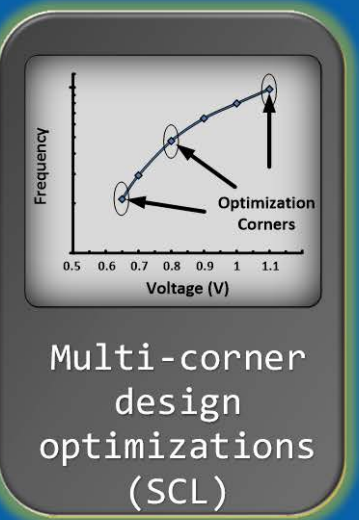
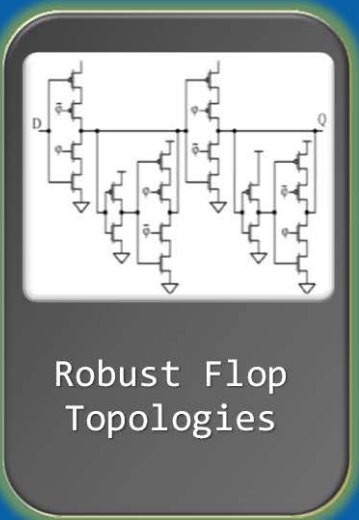
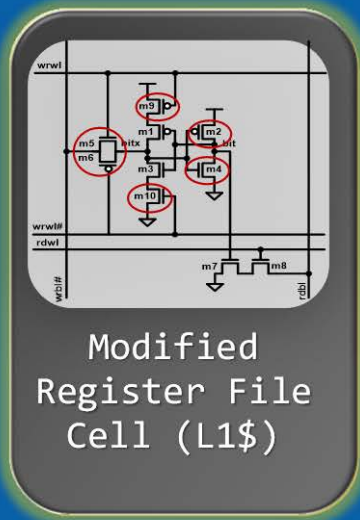
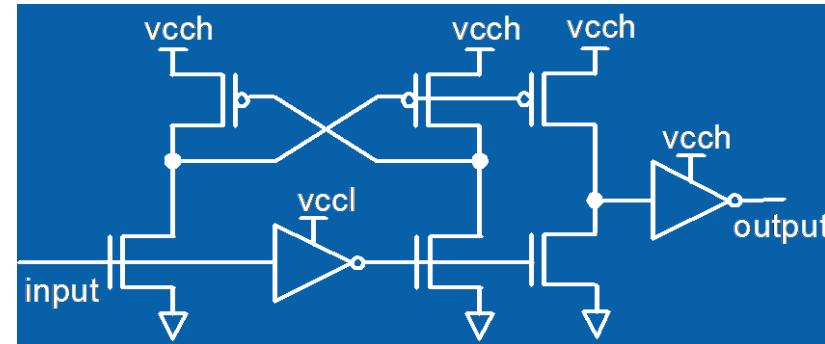
Technology	32nm High-K Metal Gate
Interconnect	1 Poly, 9 Metal (Cu)
Transistors	6 Million (Core)
Core Area	2mm ²

NTV design techniques

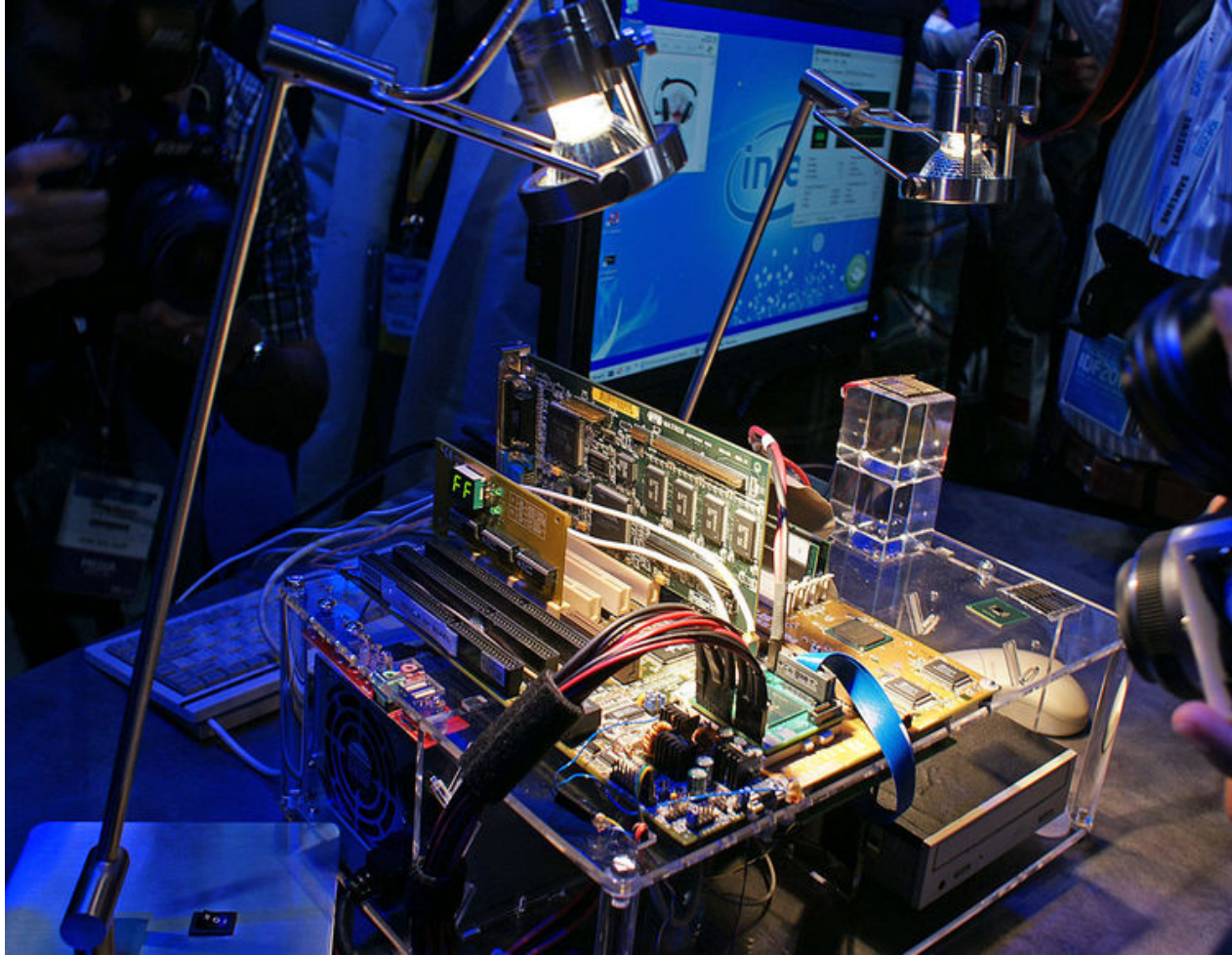
Narrow muxes No stack height > 2



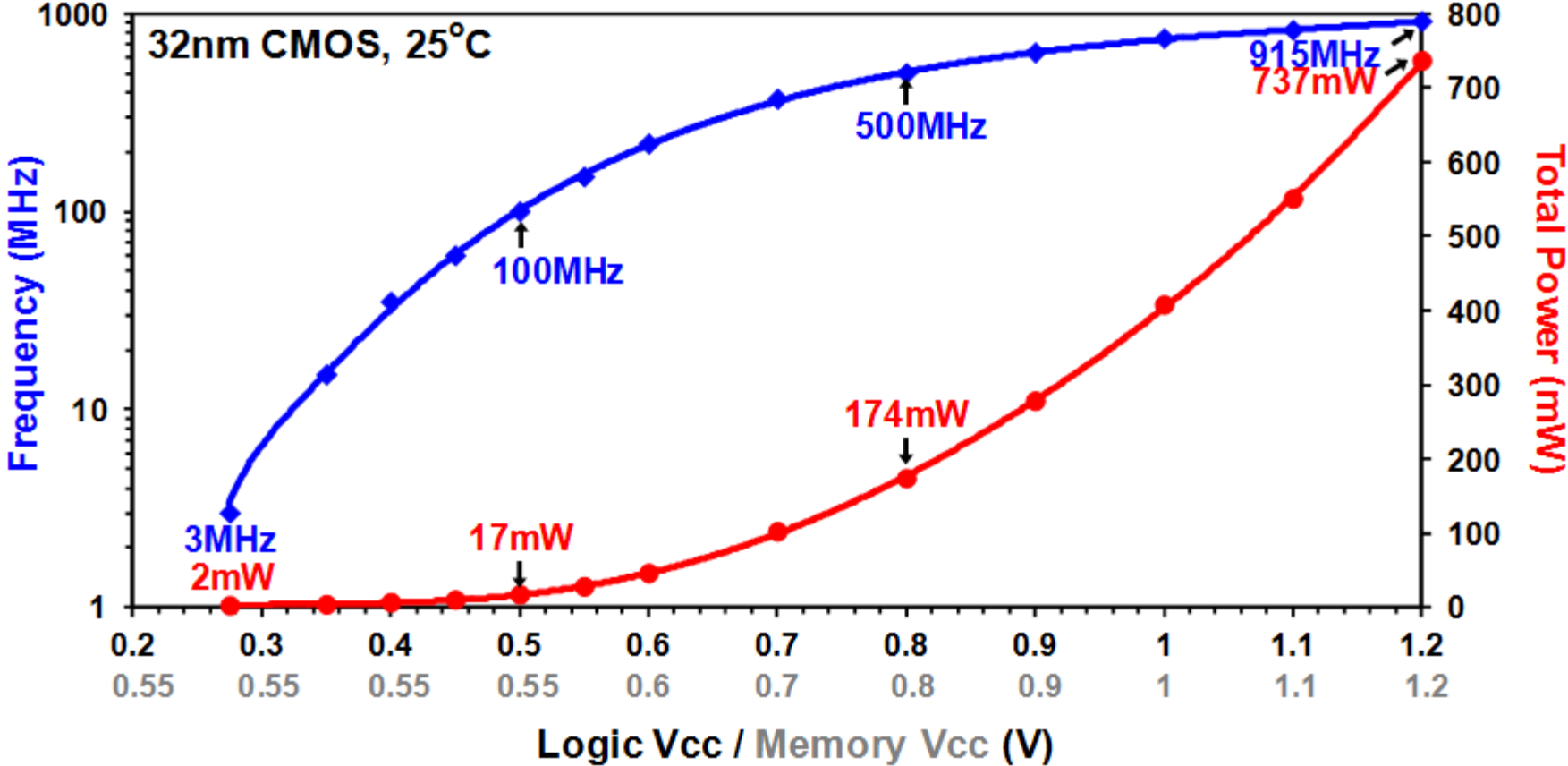
Robust level converters



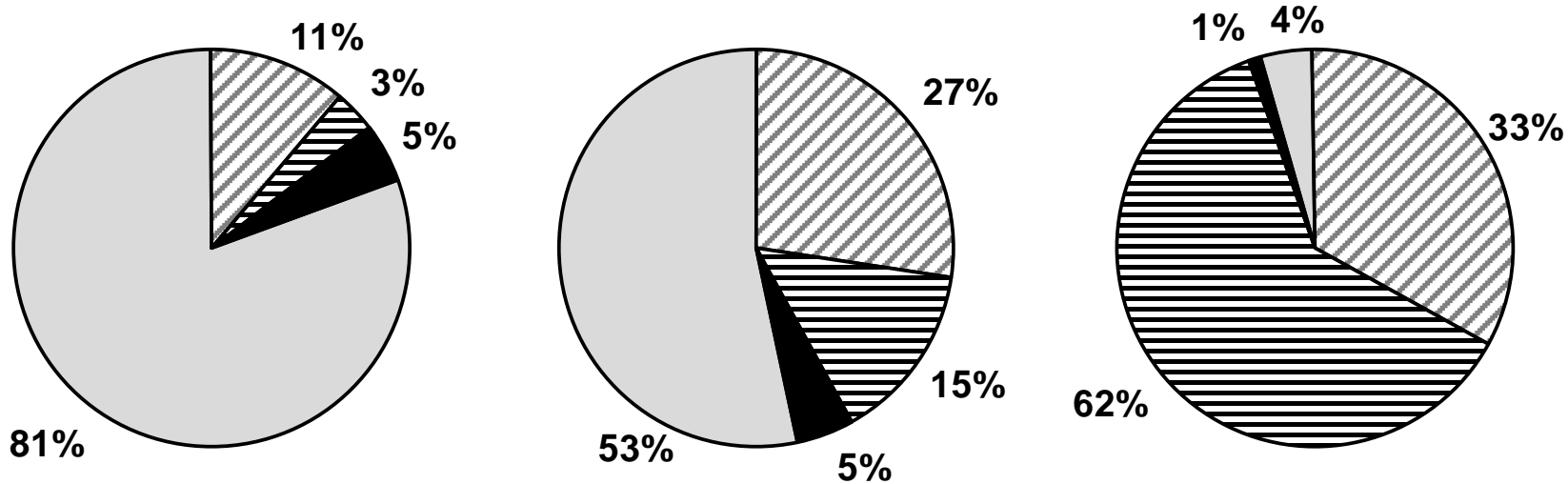
NTV IA – powered by solar cell!



Power performance measurements



Power components



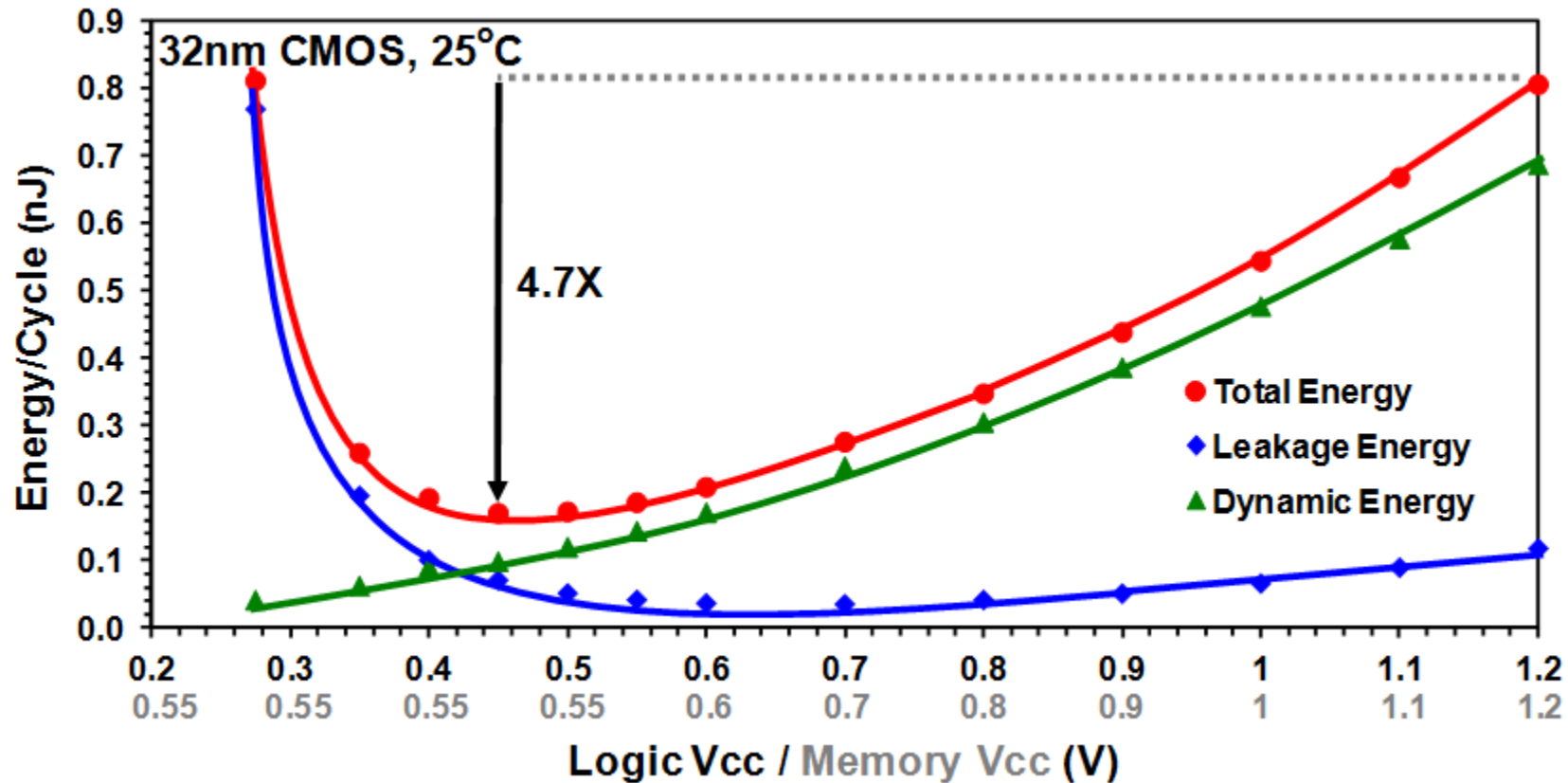
Vcc-max (Super-Threshold) Vcc-opt (Near-Threshold) Vcc-min (Sub-Threshold)

Logic Vcc: 1.2V
Memory Vcc: 1.2V

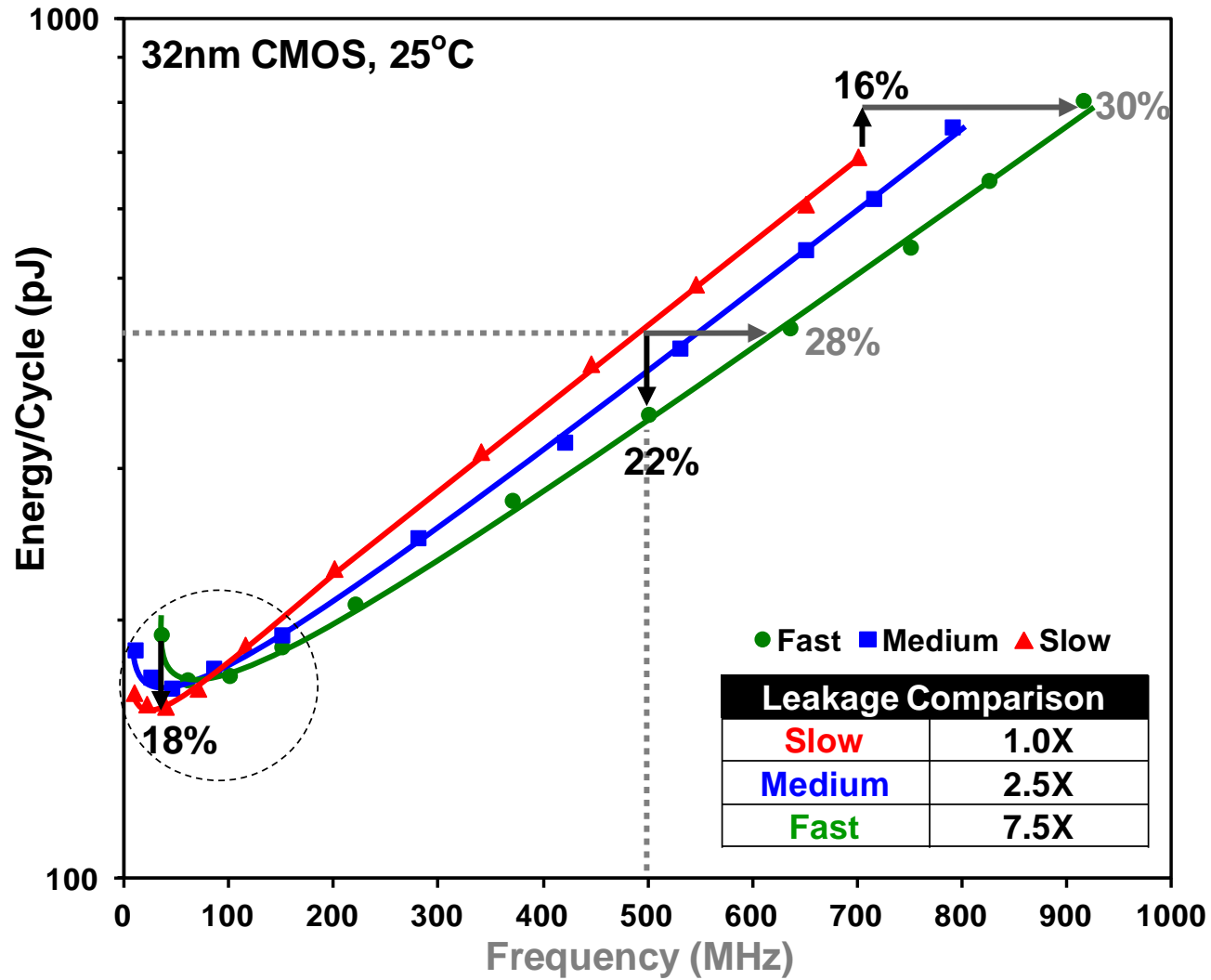
Logic Vcc: 0.45V
Memory Vcc: 0.55V

Logic Vcc: 0.28V
Memory Vcc: 0.55V

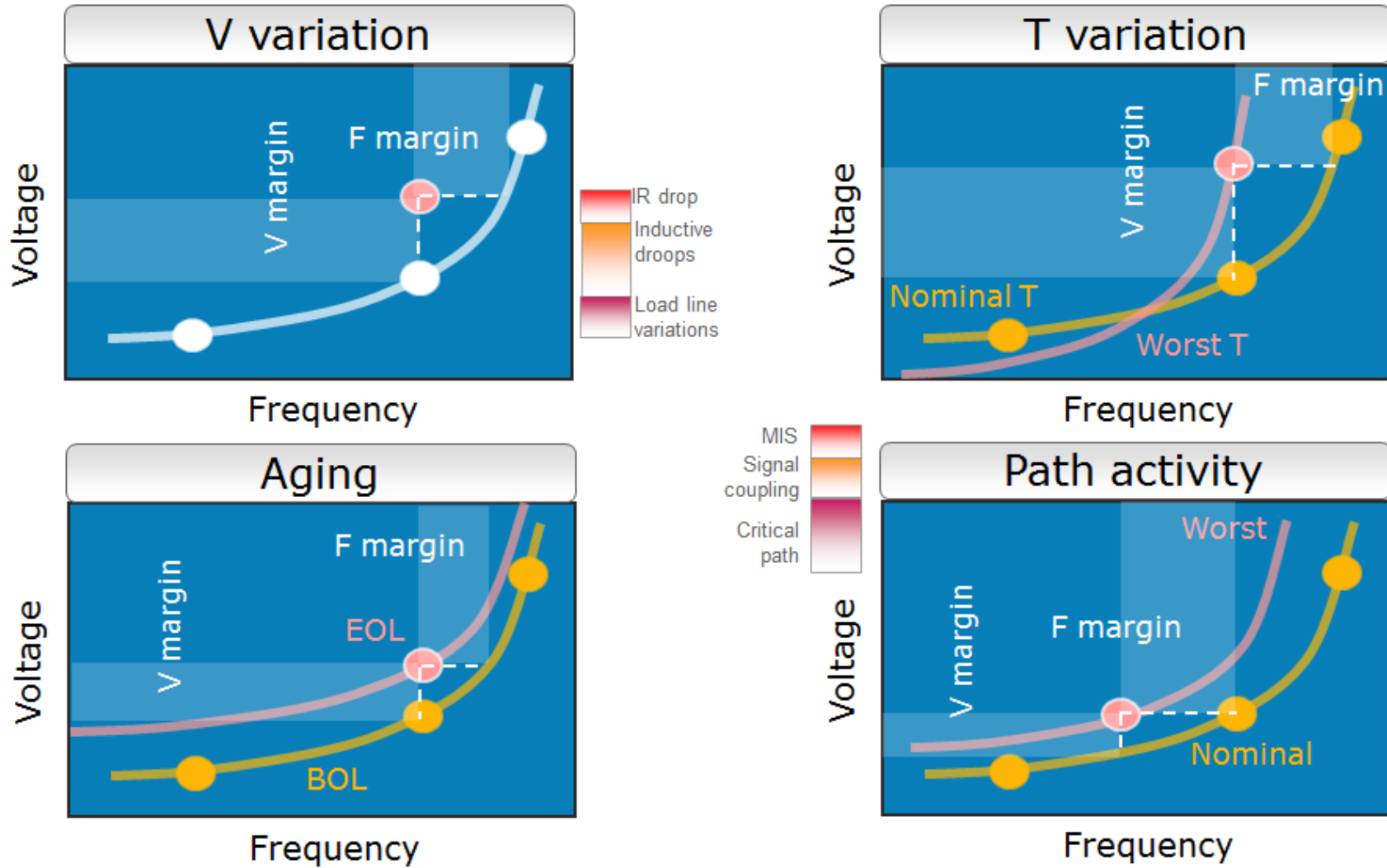
Minimum energy operation



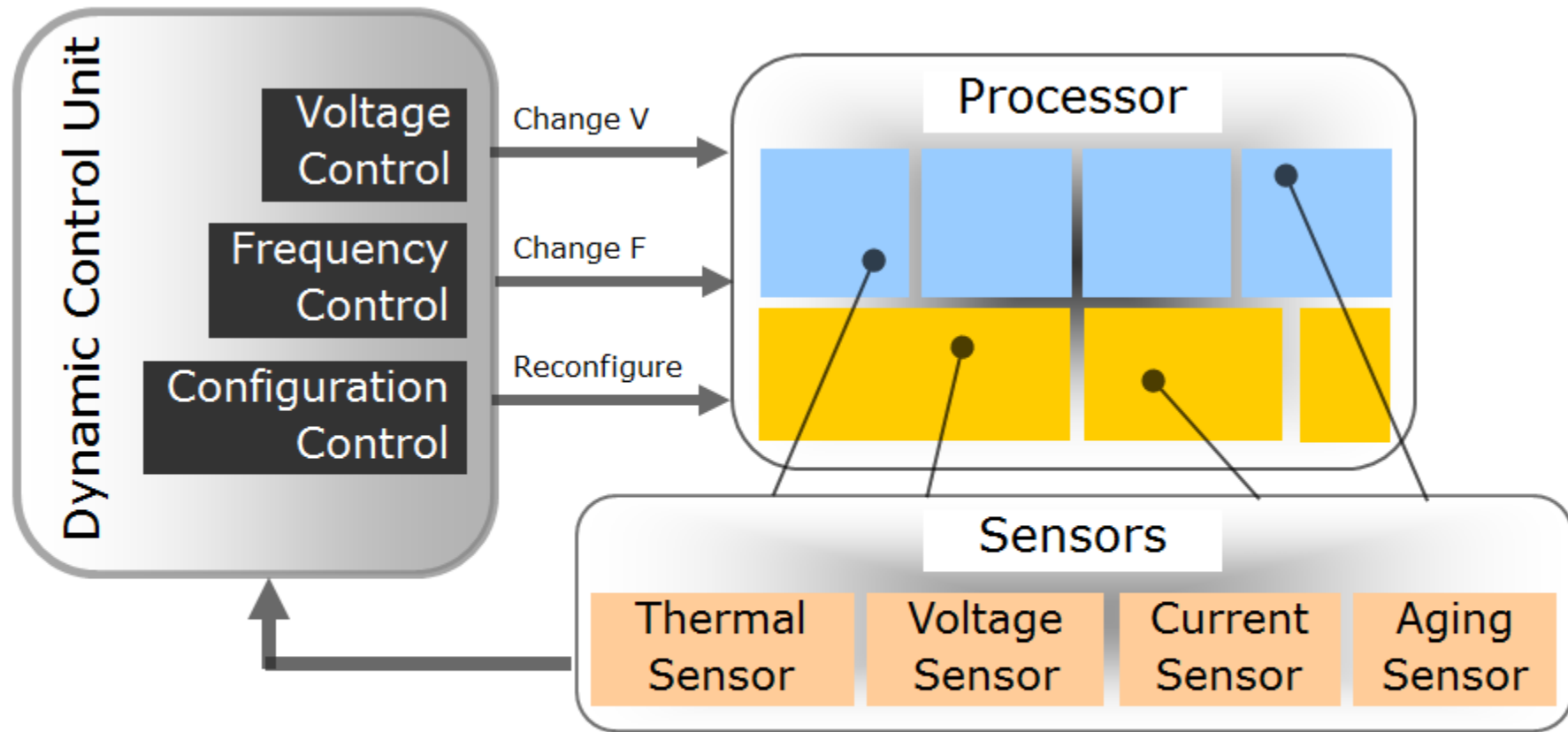
NTV and variability



Voltage-frequency margins

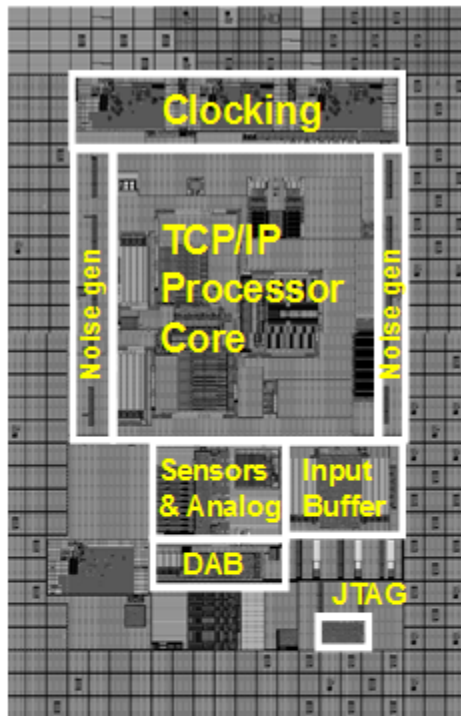


Dynamic adaptation & reconfiguration

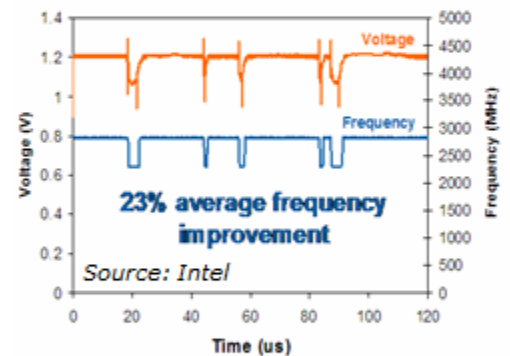
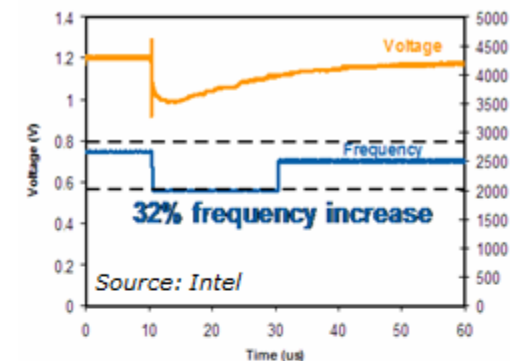
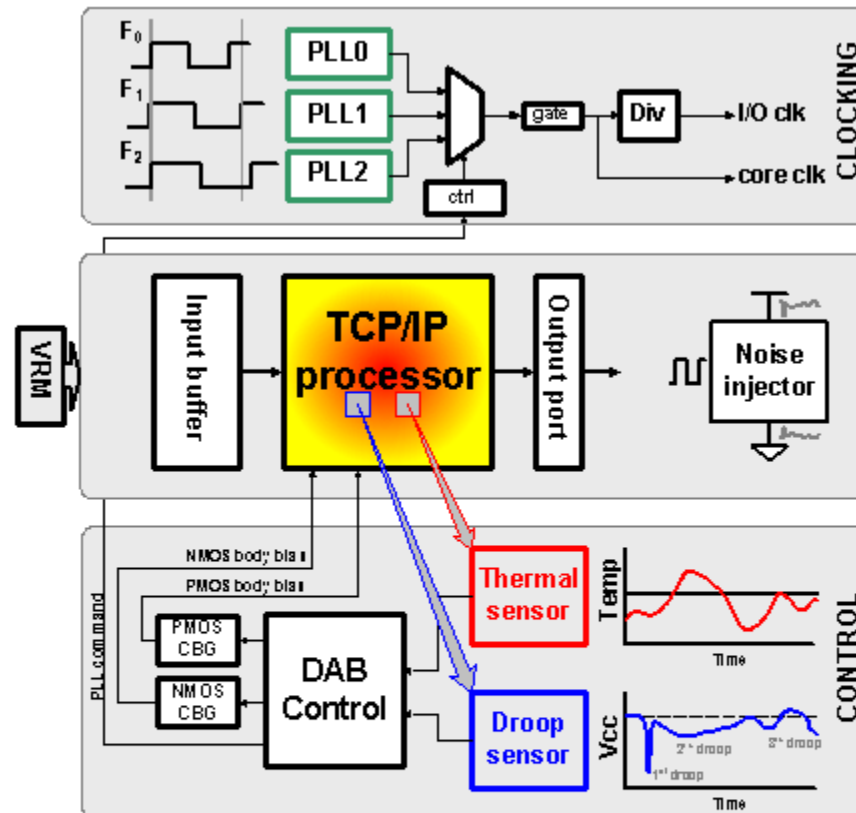


Adapt & reconfigure for best power-performance

Dynamic V & F adaptation



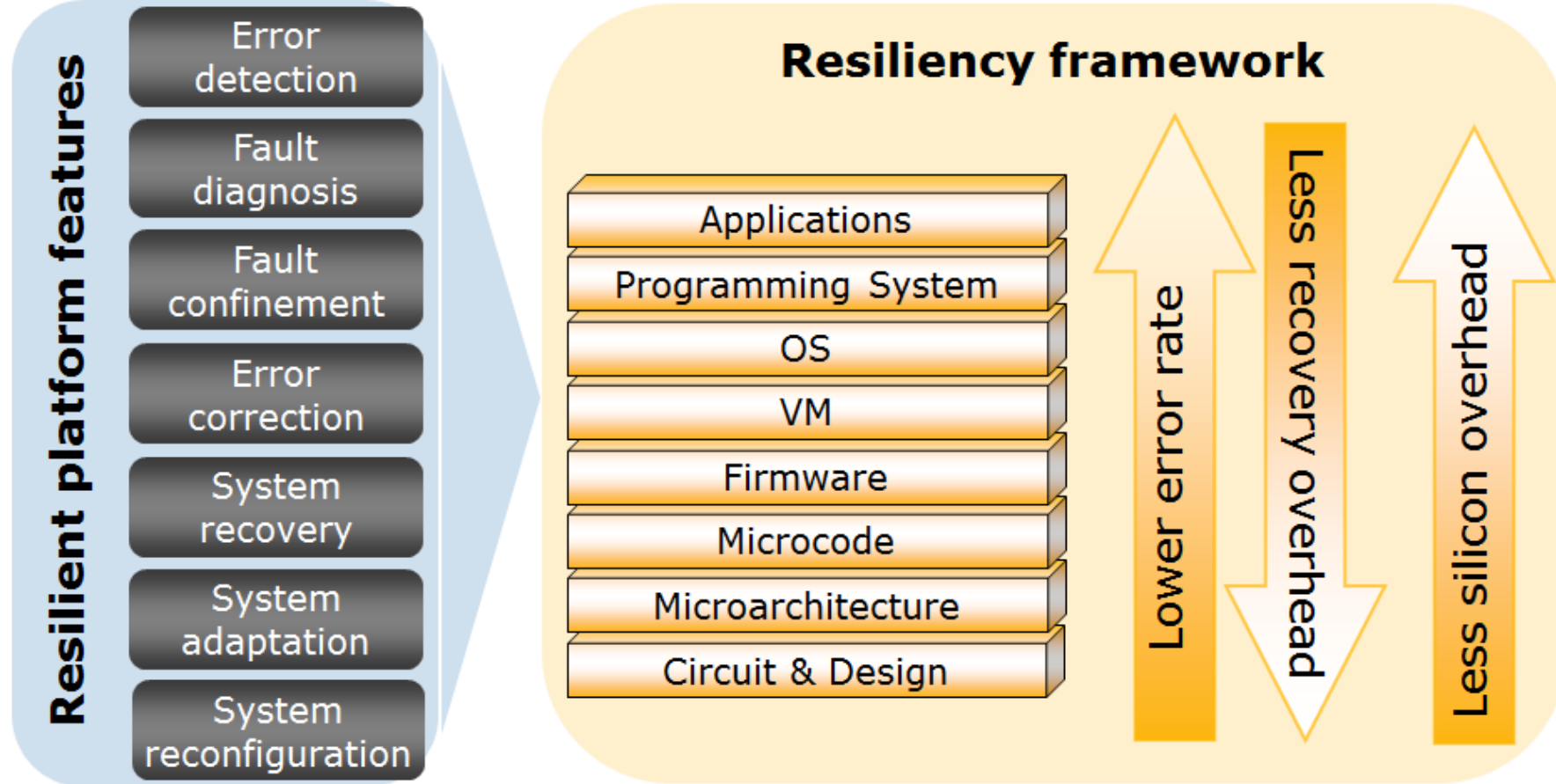
Prototype chip in 90nm



Environment-aware dynamic adaptation

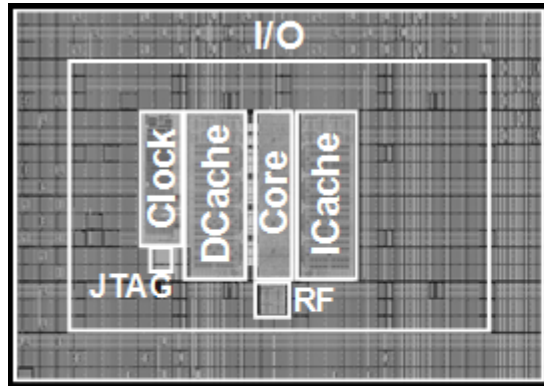
- Adapt F/V to V/T change → reduce V/T margin
- Adapt F/V to aging → reduce aging margin

Resilient platforms

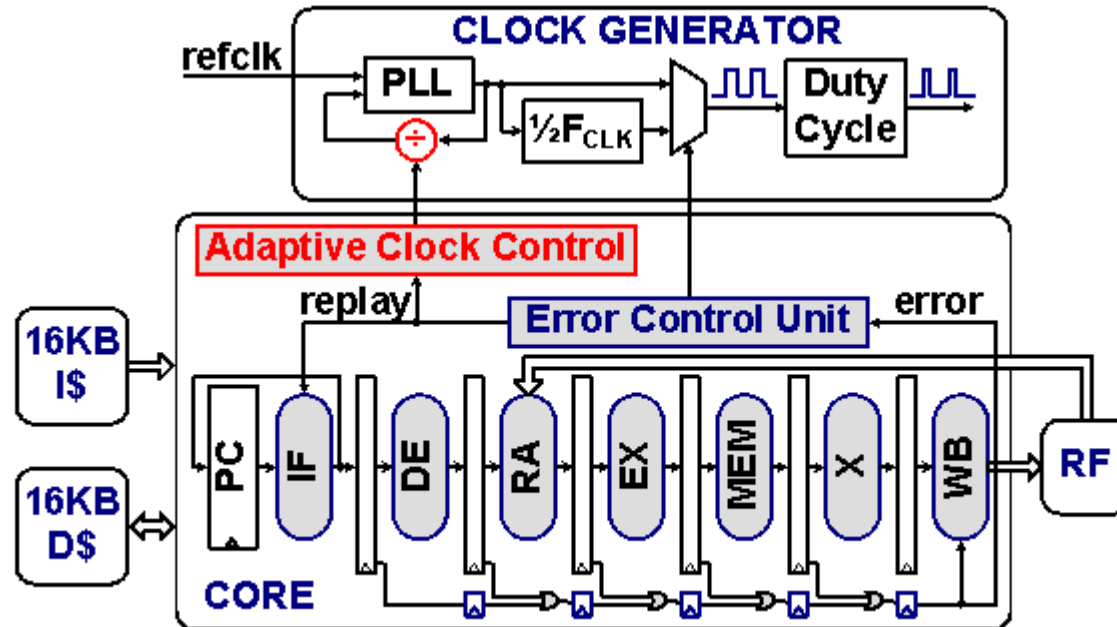
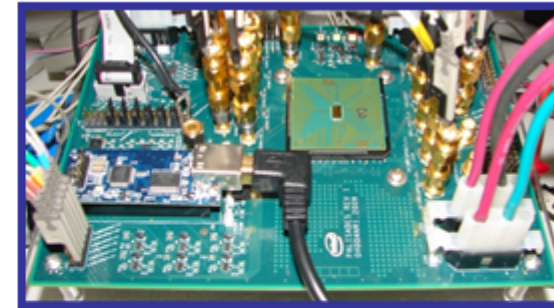


Resiliency for performance, efficiency & reliability

Resilient & adaptive core

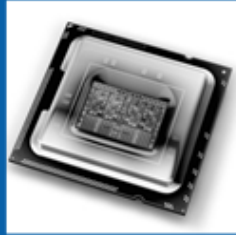


Technology	45nm CMOS
Die Area	13.64 mm ²
Core Area	0.39 mm ²
Core F _{MAX}	1.45GHz at 1.0V
Core Power	135mW at 1.0V

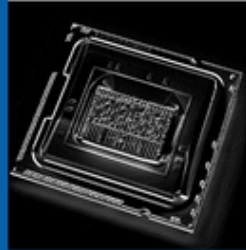


Performance & efficiency gains

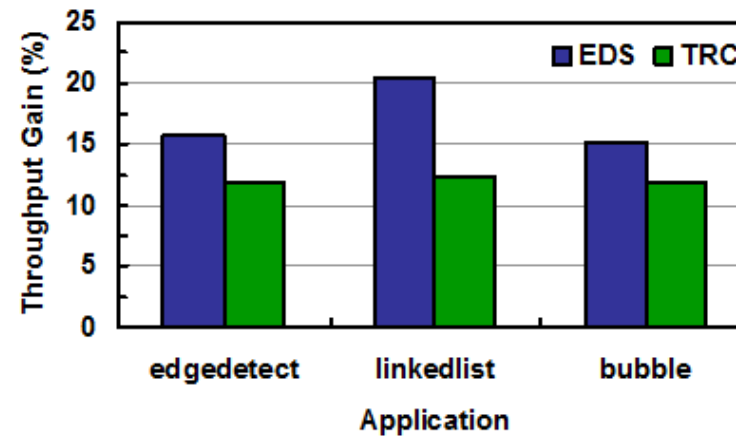
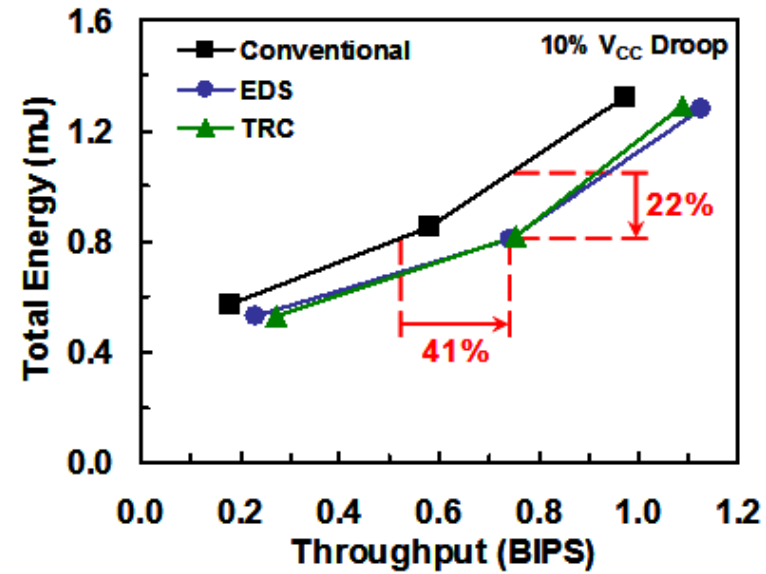
Input Image



Output Image:
Resiliency ON



Output Image:
Resiliency OFF



Integrated voltage regulators

Conversion

- Area efficient
- Scalable
- Persistent rail

Efficient

Distribution

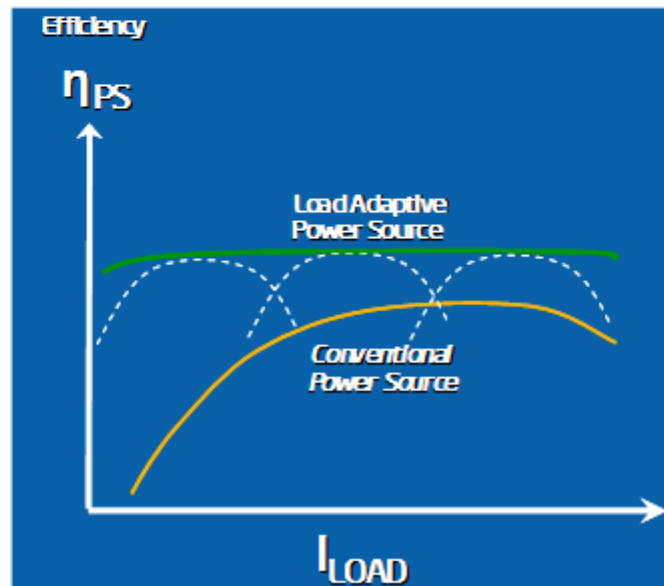
- Lower loss
- Higher fidelity
- Simpler

Low loss

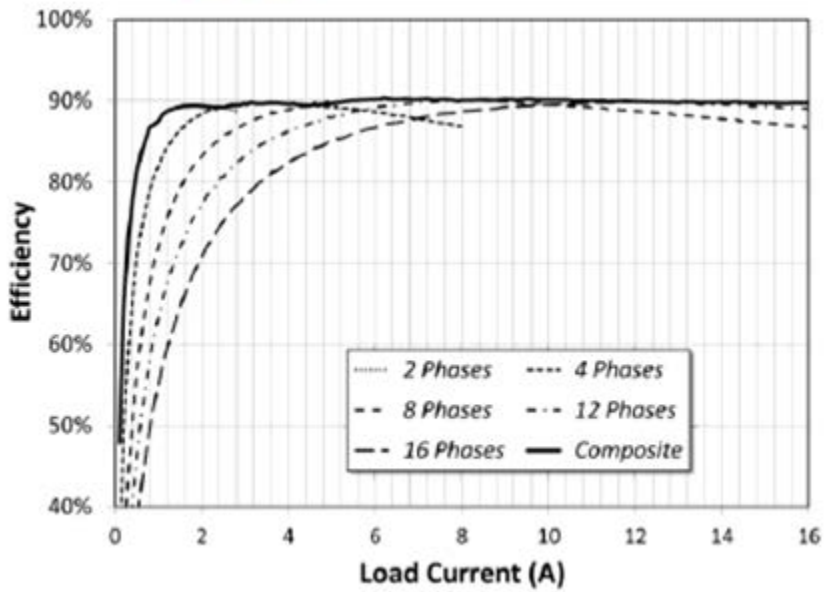
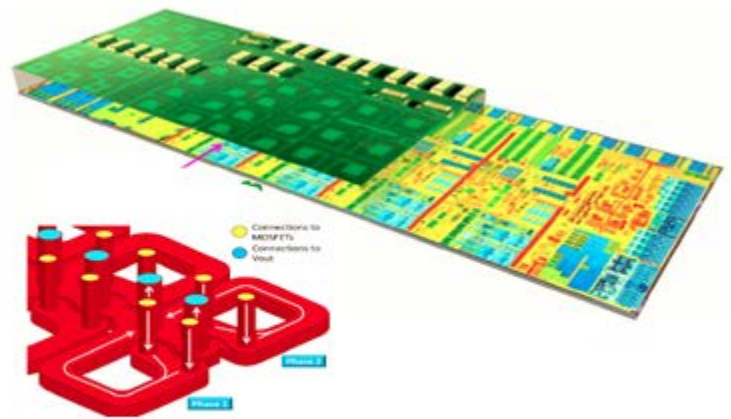
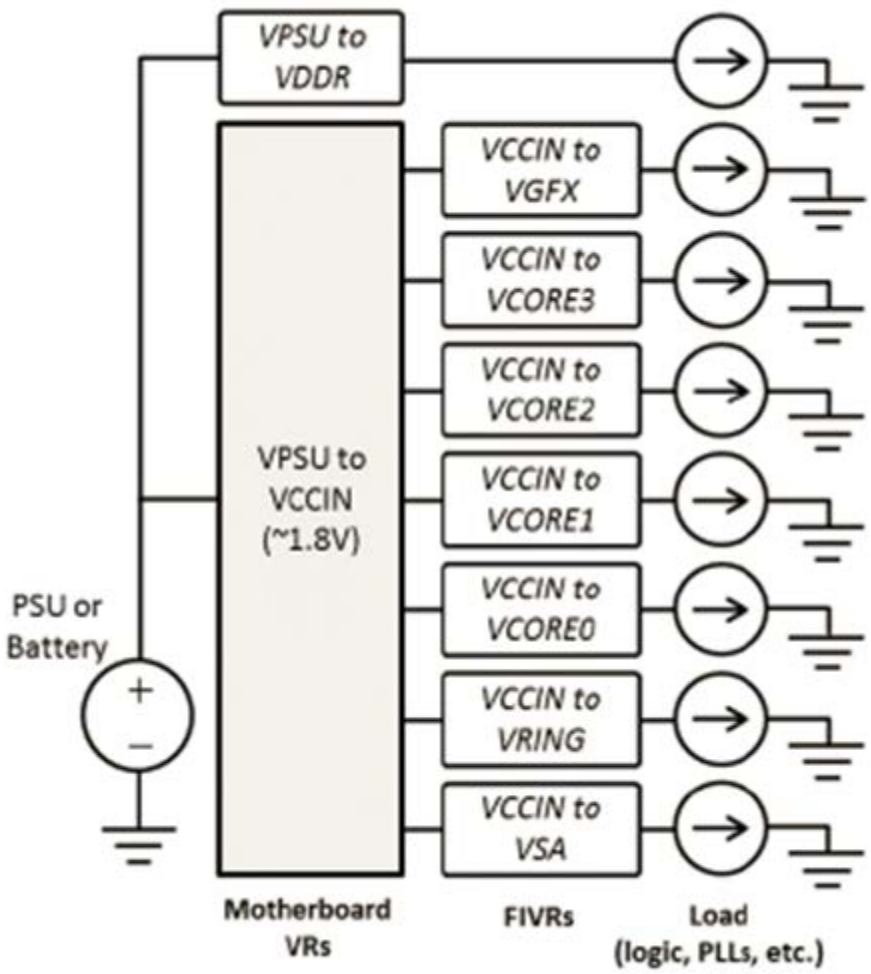
Control

- Fast & efficient
- Load adaptive
- Independent rails

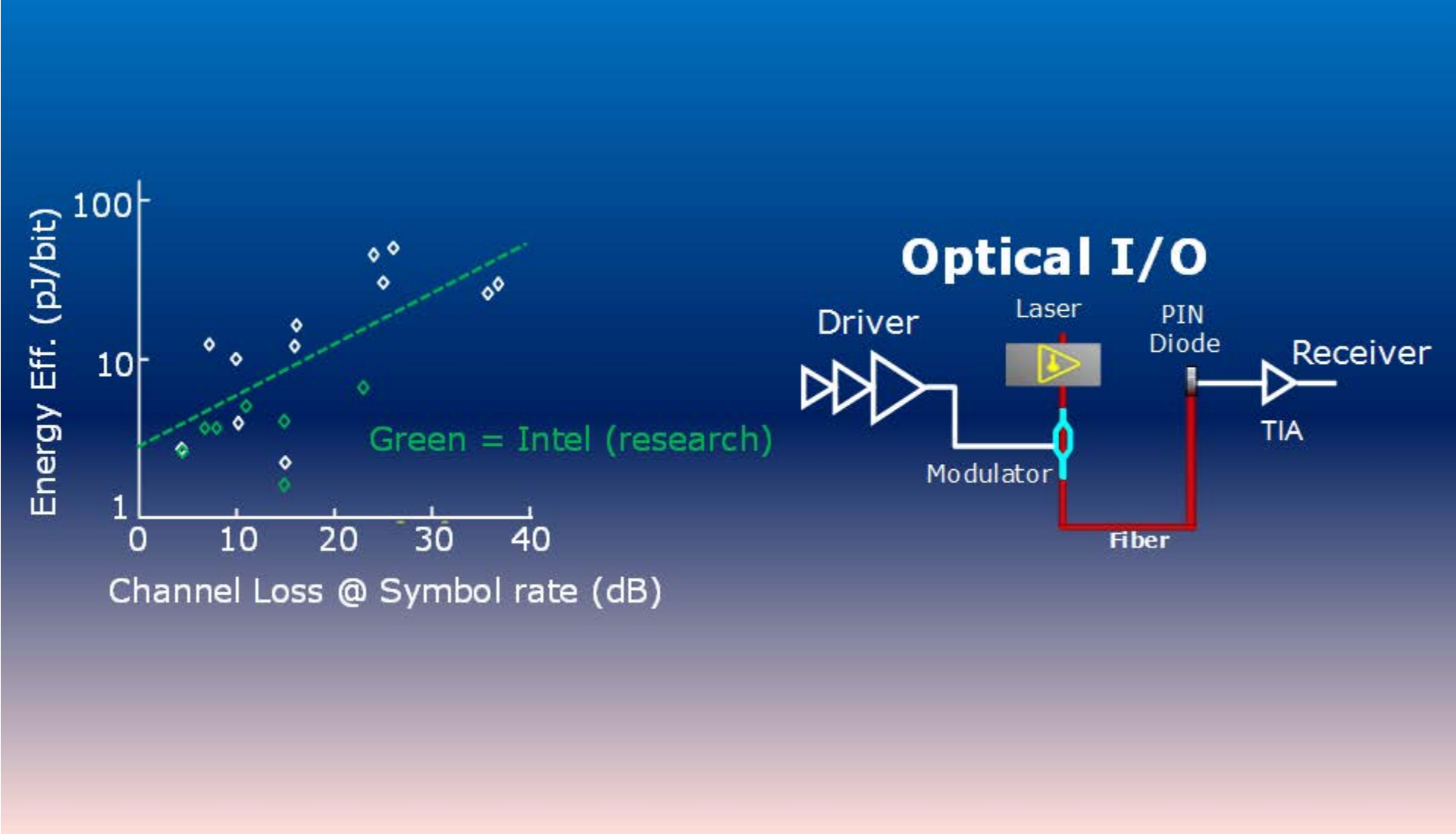
Fine-grain



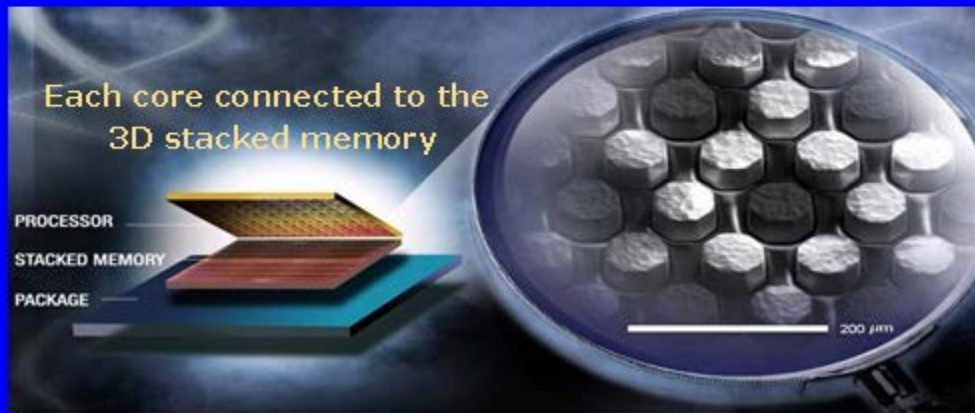
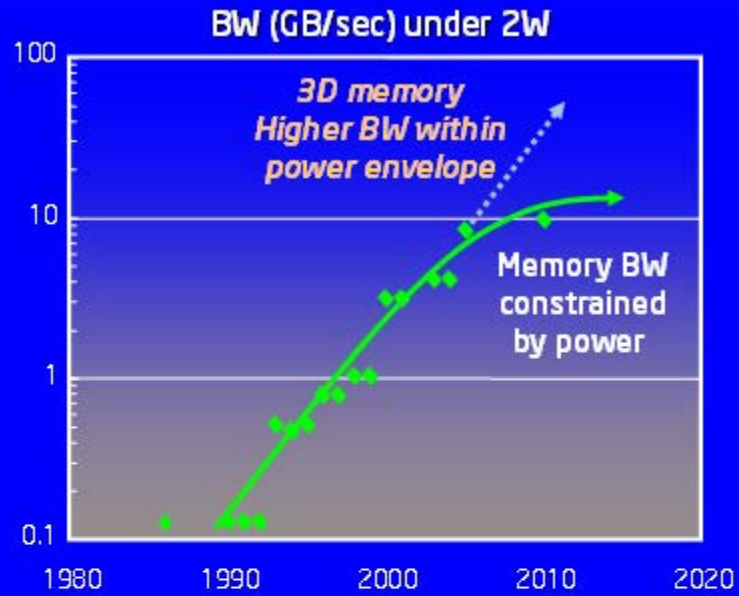
Fully integrated VR



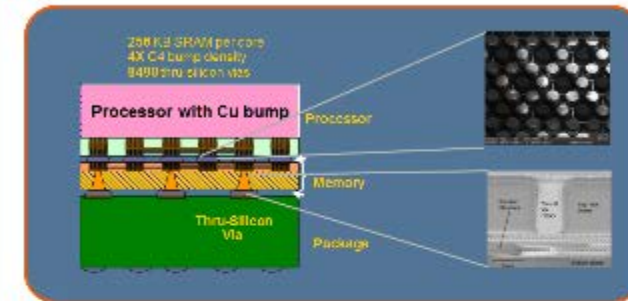
Energy efficient interconnects



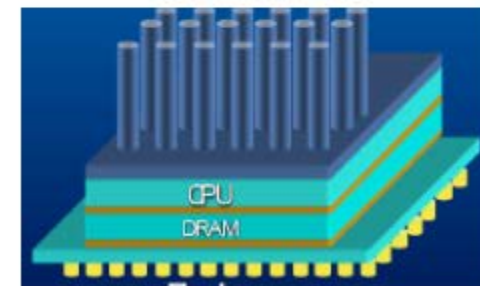
Memory capacity & bandwidth



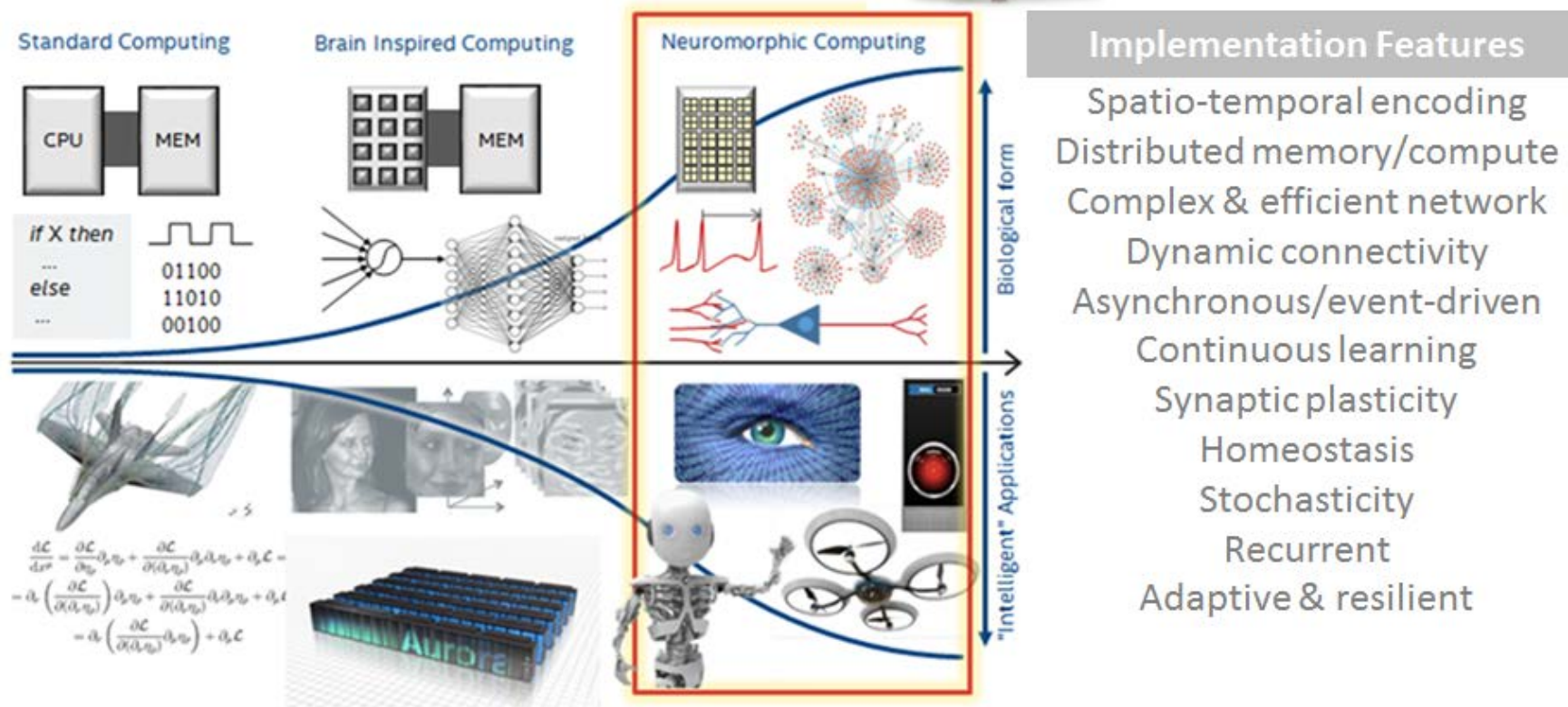
3D Integration: SRAM



3D Integration: DRAM

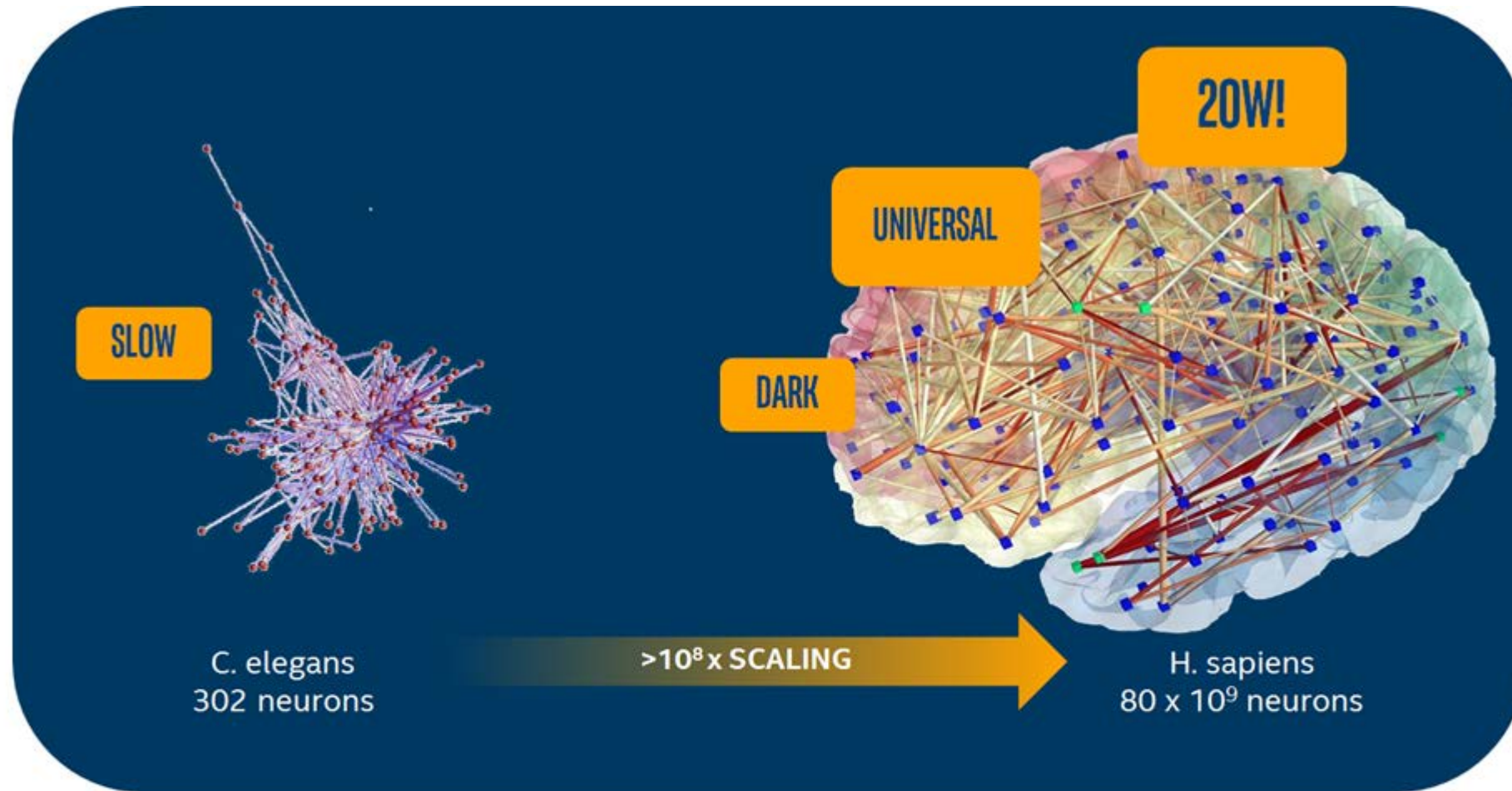


Efficient & scalable neuromorphic systems



Exploit integration, spike timing, sparsity, plasticity & resiliency

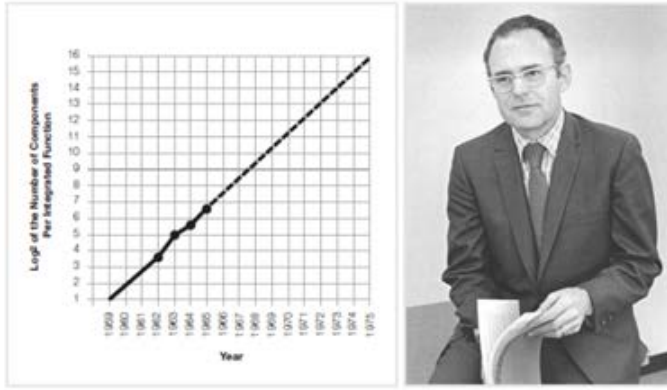
Efficient & scalable neuromorphic architecture



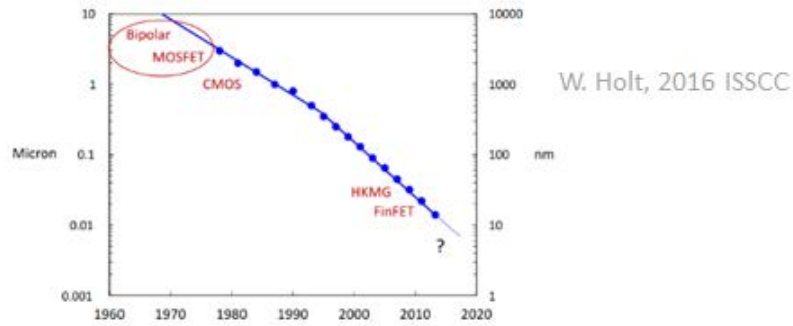
Phenomenal architecture efficiency & scalability in nature!

The next big leap...

Moore's Law: Economics AND Power

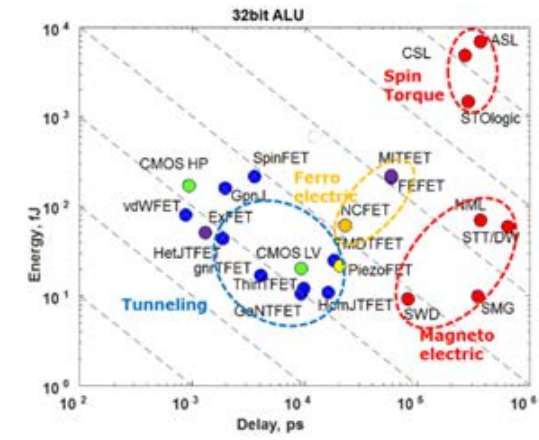
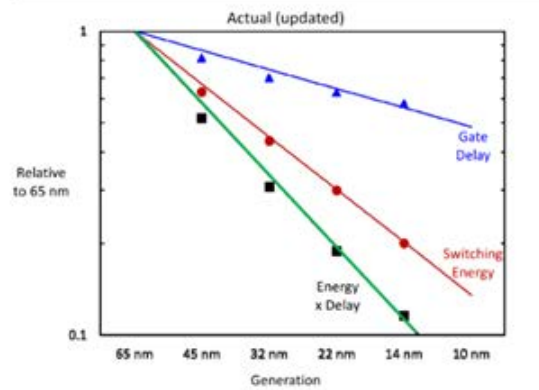


Heat problem
 "Will it be possible to remove the heat generated by tens of thousands of components in a single silicon chip?"
 "Cramming more components onto integrated circuits", Electronics, Volume 38, Number 8, April 19, 1965



W. Holt, 2016 ISSCC

CMOS & Beyond



D. Nikonov & I. Young, 2015 JXDC

Big leaps in energy efficiency trigger technology transitions!

LIMITLESS POSSIBILITIES

WE ARE HERE

