

Energy Efficient Computing in Nanoscale CMOS



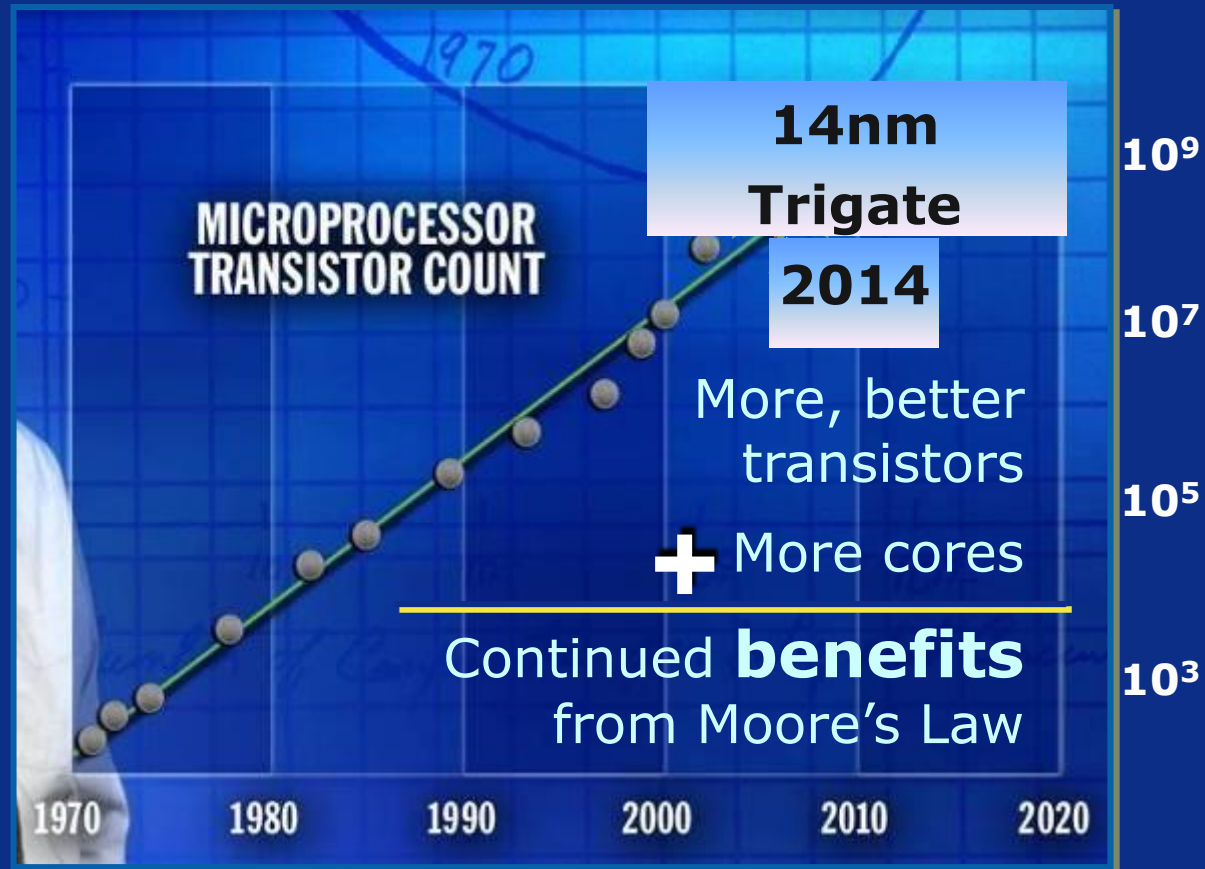
Vivek De
Intel Fellow
Director of Circuit Technology Research
Intel Labs

Internet of Everything (IoE)

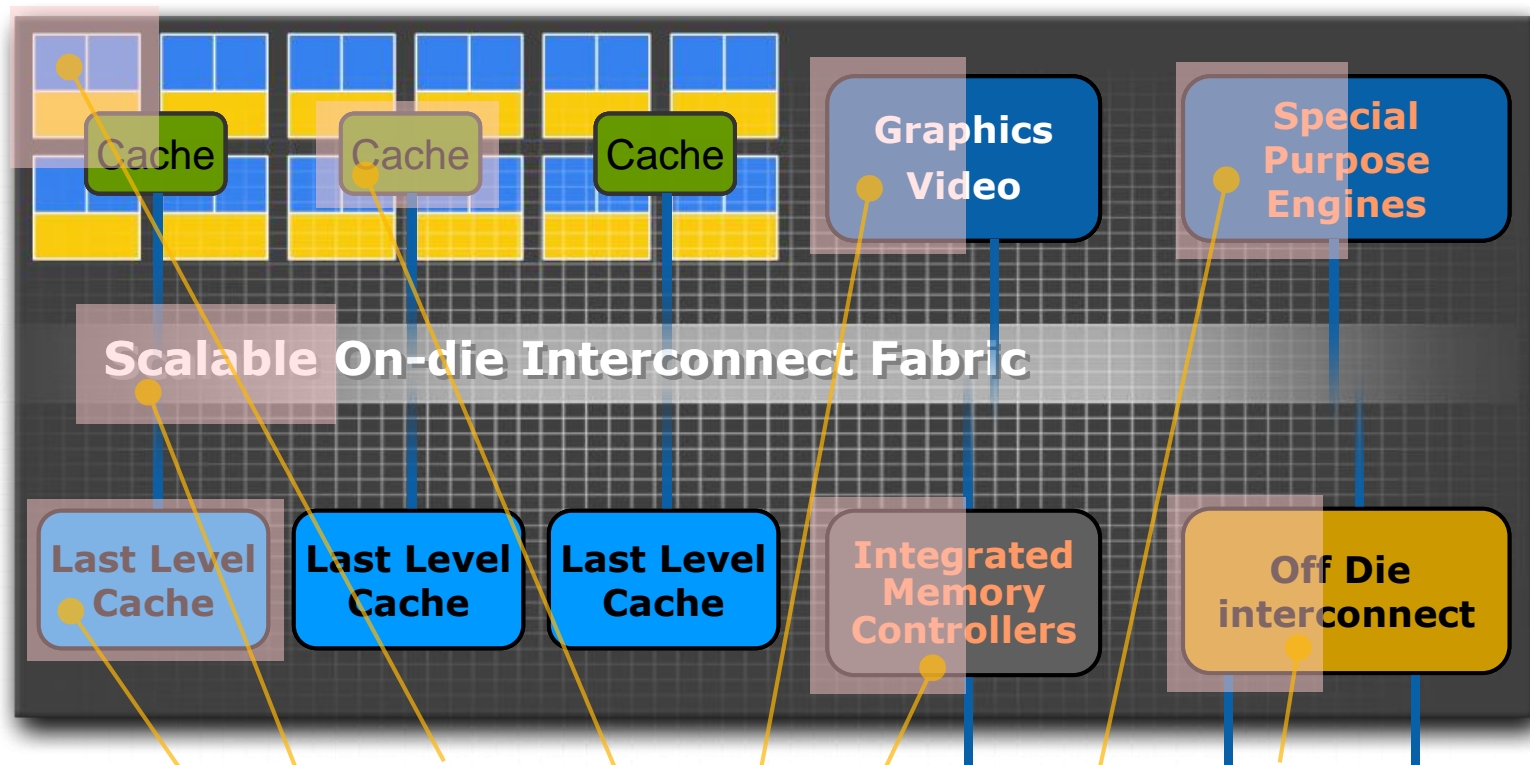


Need end-to-end energy efficiency

Moore's Law scaling



Dynamic platform control



Maximize performance & efficiency

Independent V/F control regions

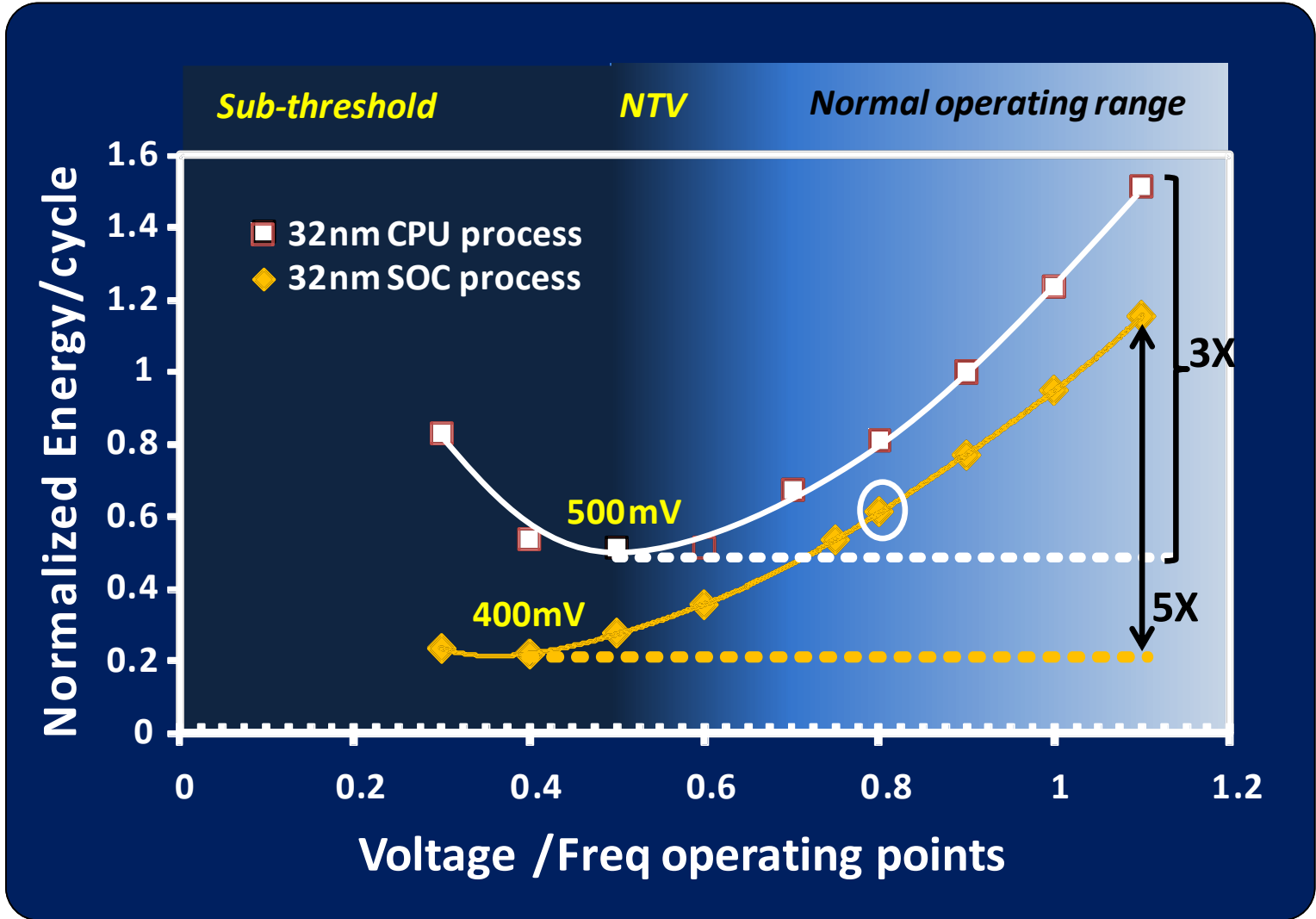
Scenario-based power allocation

Dynamic V/F control

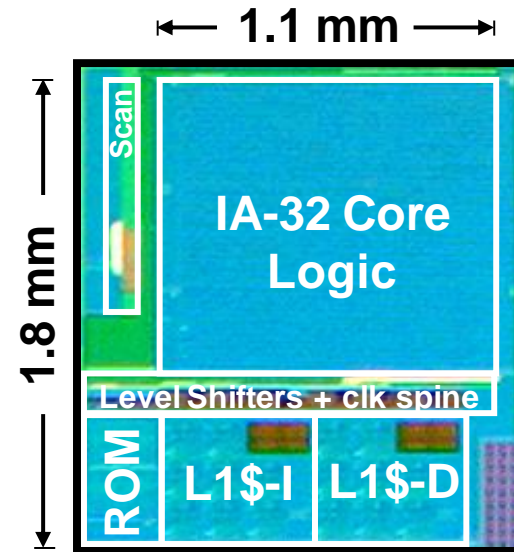
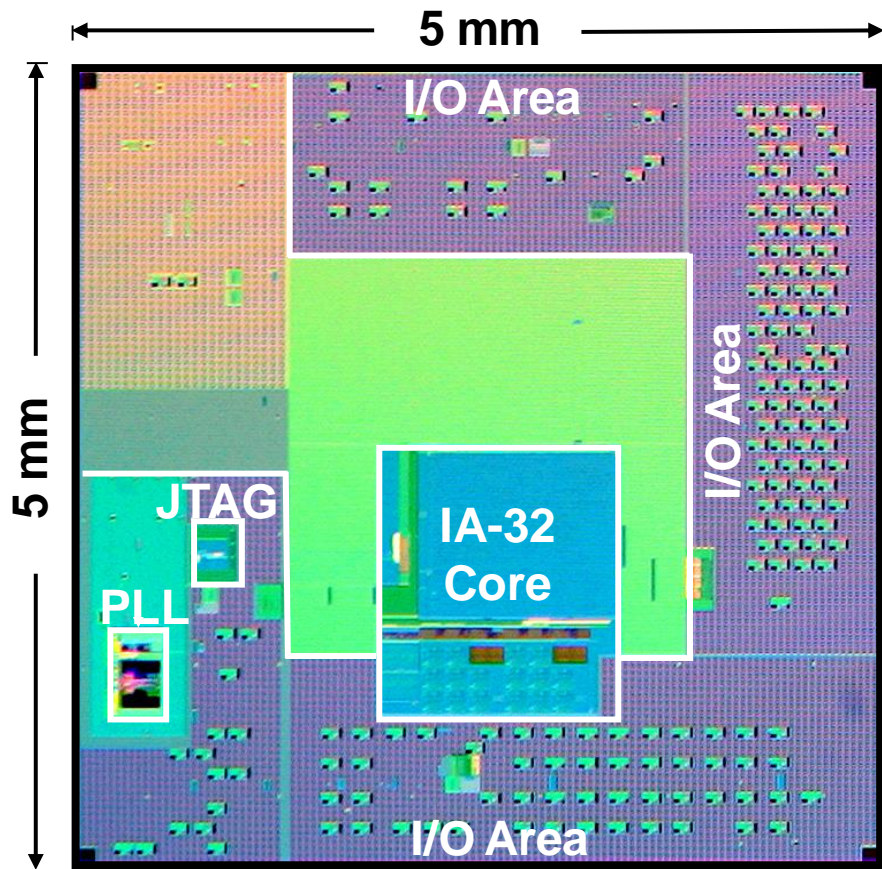
Workload-based core activation & shutdown

Deliver best user experience under constraints

Near Threshold Voltage (NTV) computing



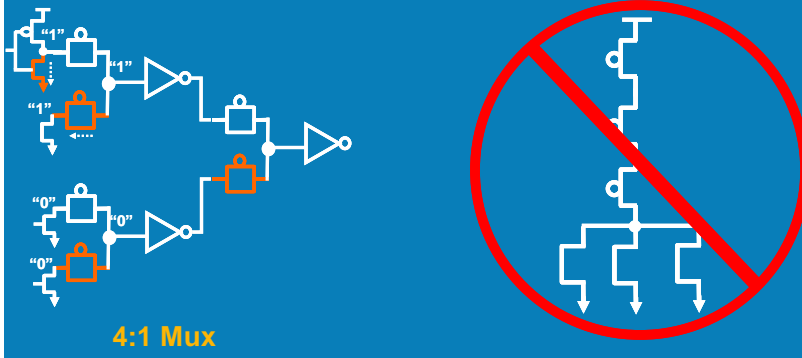
NTV IA processor



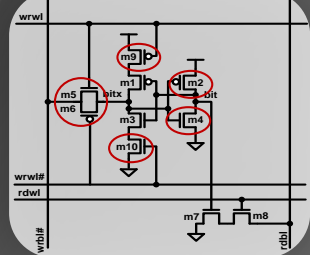
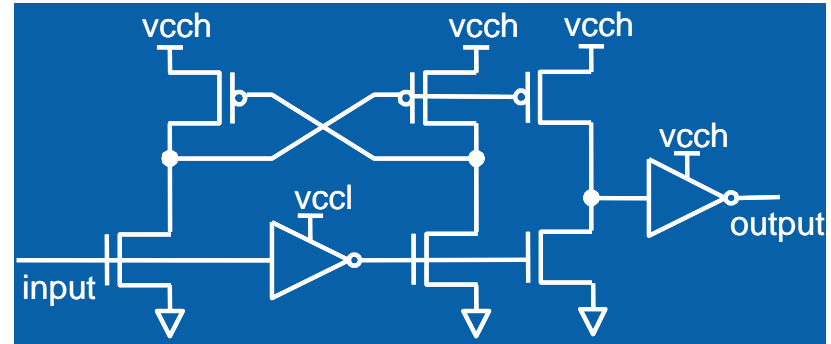
Technology	32nm High-K Metal Gate
Interconnect	1 Poly, 9 Metal (Cu)
Transistors	6 Million (Core)
Core Area	2mm ²

NTV design techniques

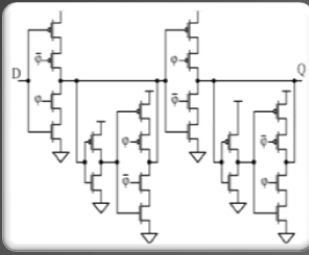
Narrow muxes No stack height > 2



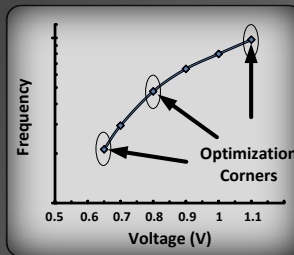
Robust level converters



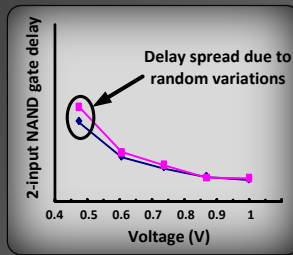
Modified Register File Cell (L1\$)



Robust Flop Topologies

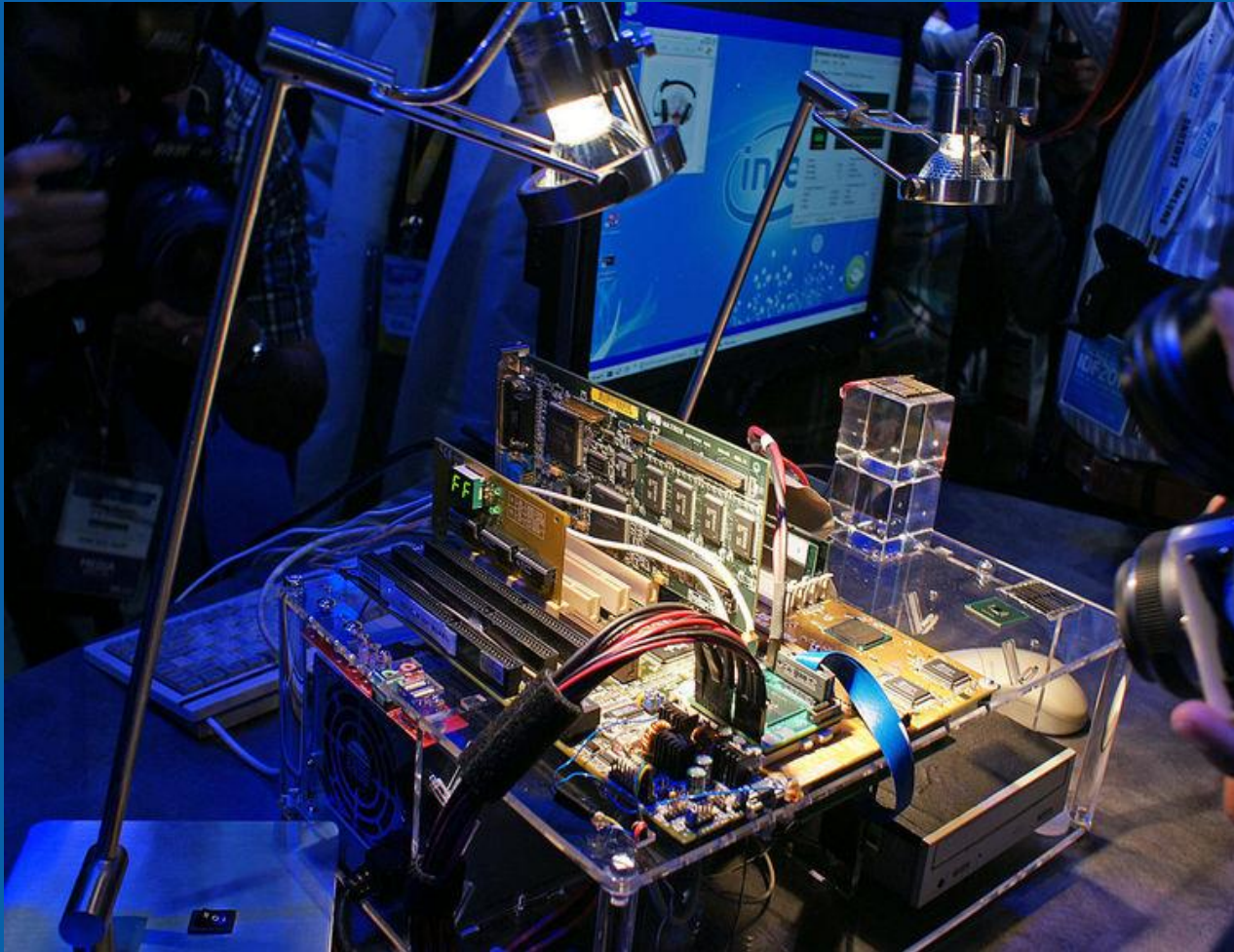


Multi-corner design optimizations (SCL)

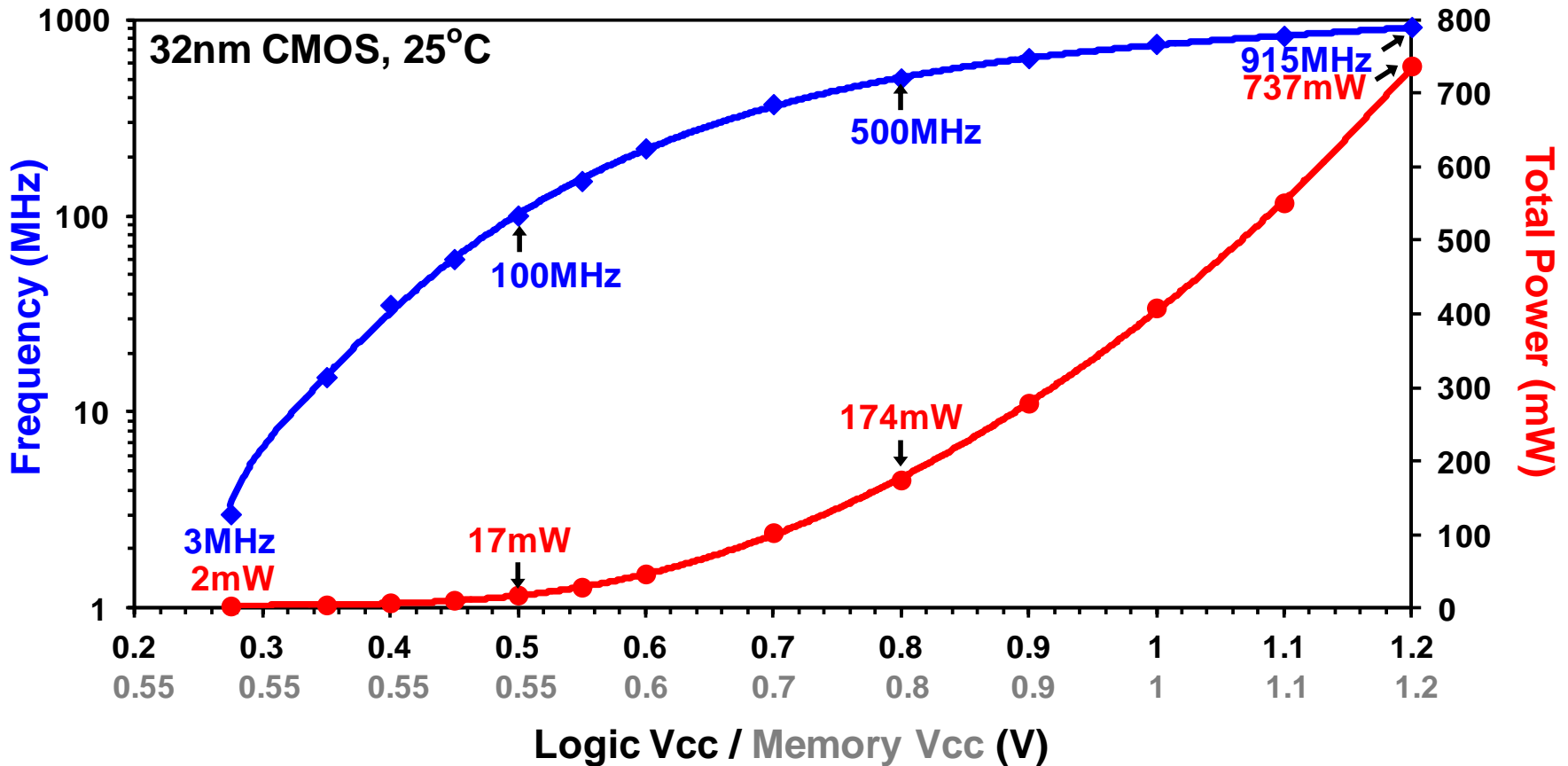


Variation-aware design
2X min Z, 40% lib cells used

NTV IA – powered by solar cell!



Power performance measurements



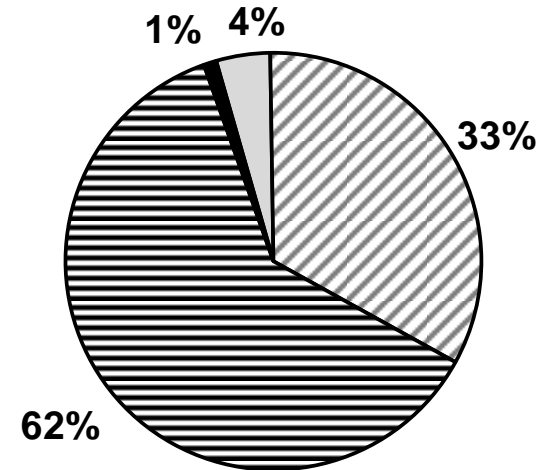
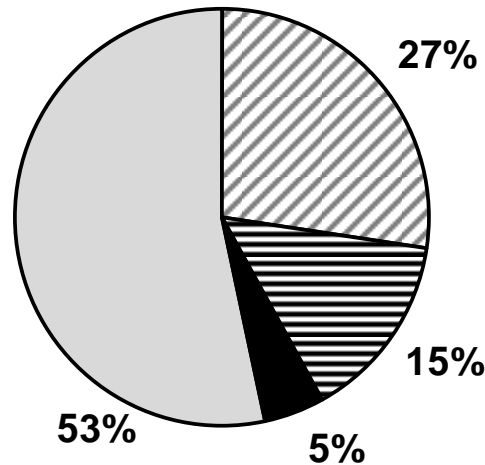
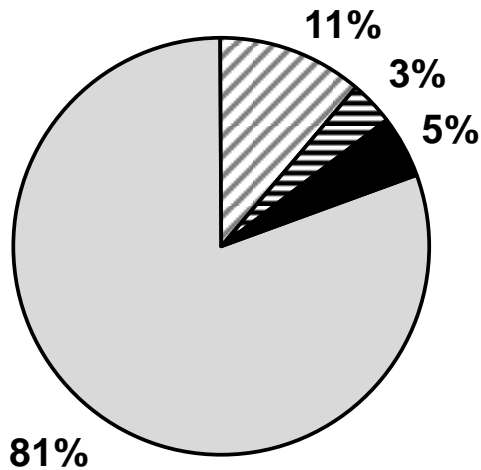
Power components

Logic Dynamic Power

Memory Dynamic Power

Logic Leakage Power

Memory Leakage Power



Vcc-max (Super-Threshold)

Vcc-opt (Near-Threshold)

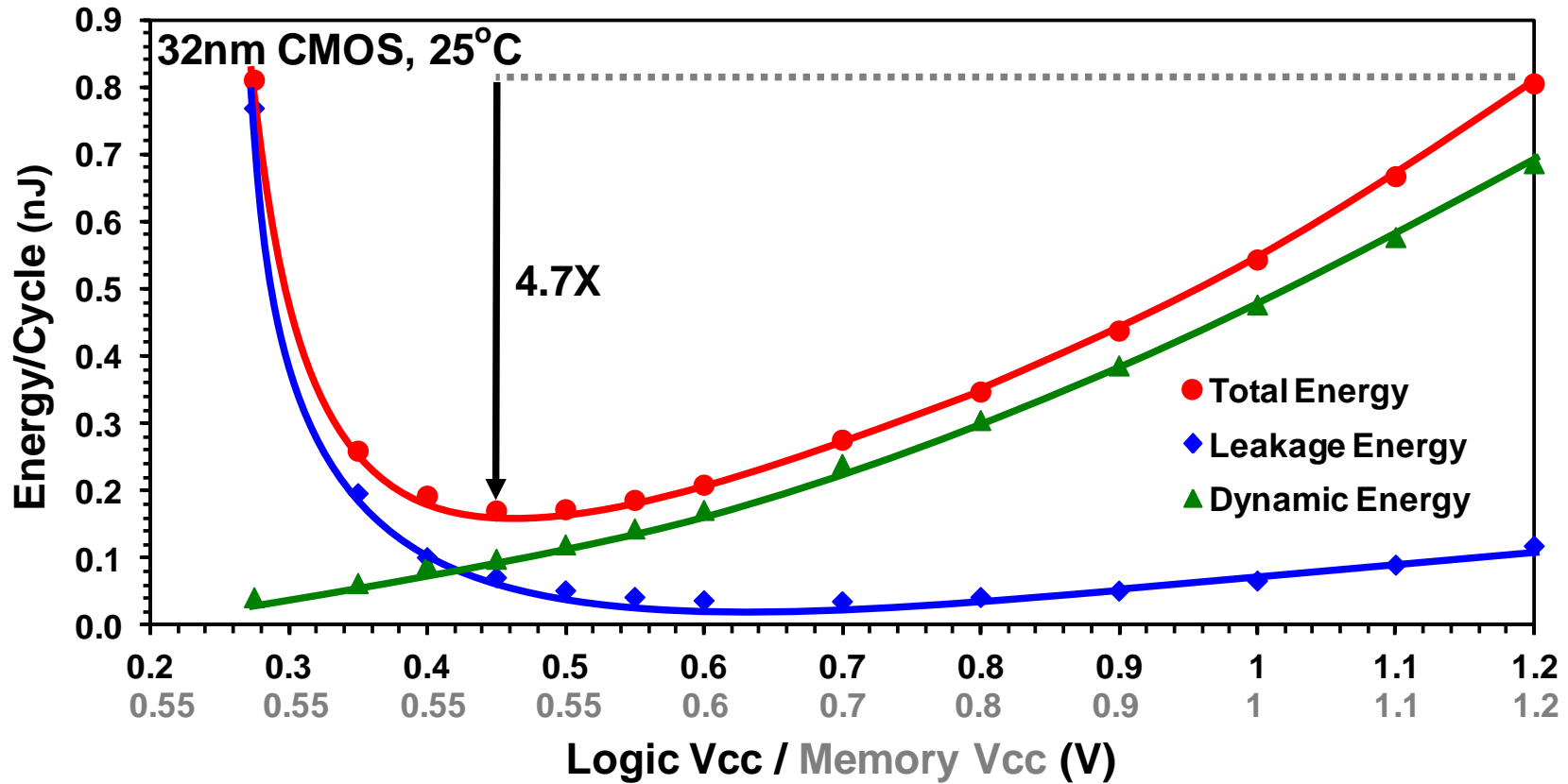
Vcc-min (Sub-Threshold)

Logic Vcc: 1.2V
Memory Vcc: 1.2V

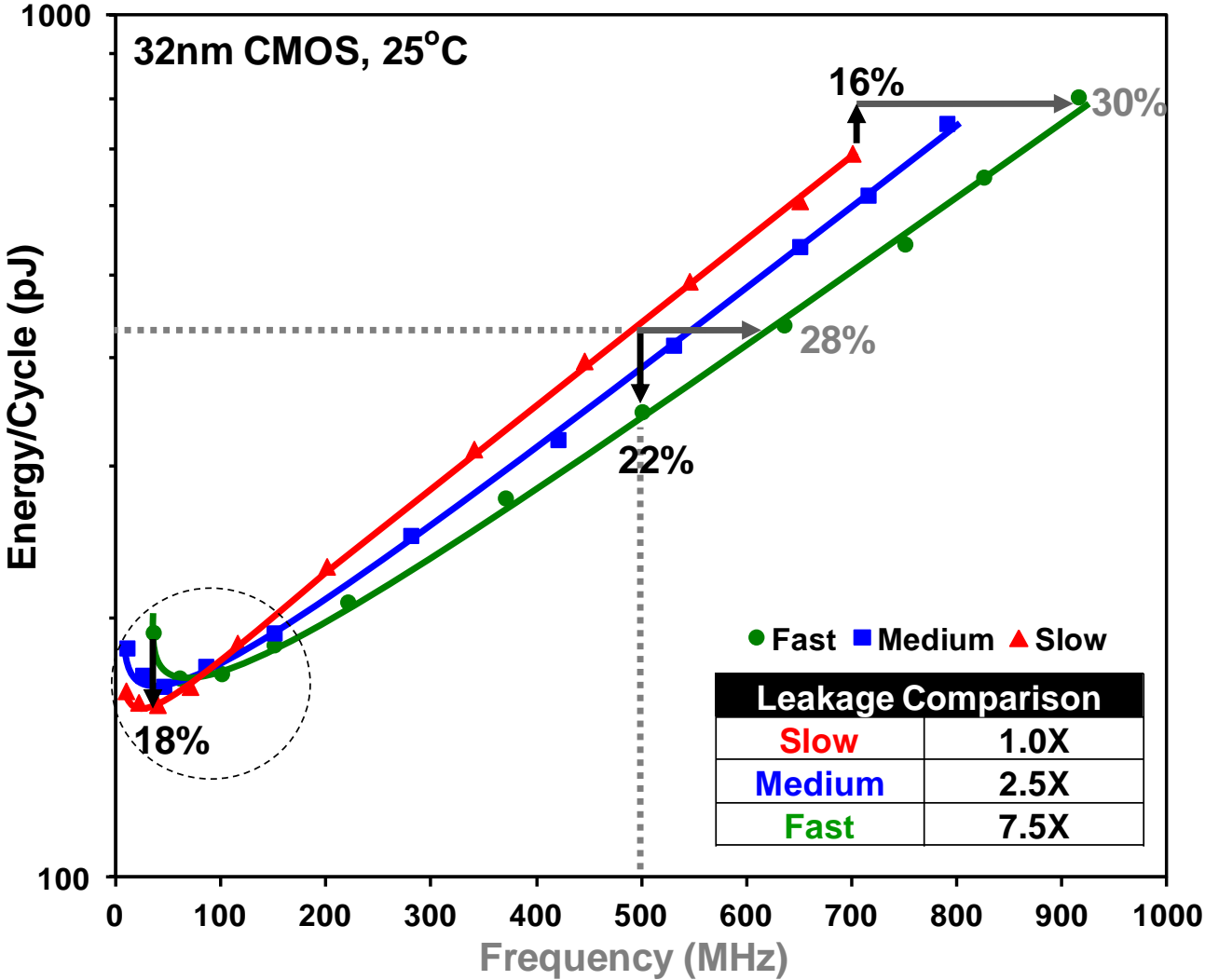
Logic Vcc: 0.45V
Memory Vcc: 0.55V

Logic Vcc: 0.28V
Memory Vcc: 0.55V

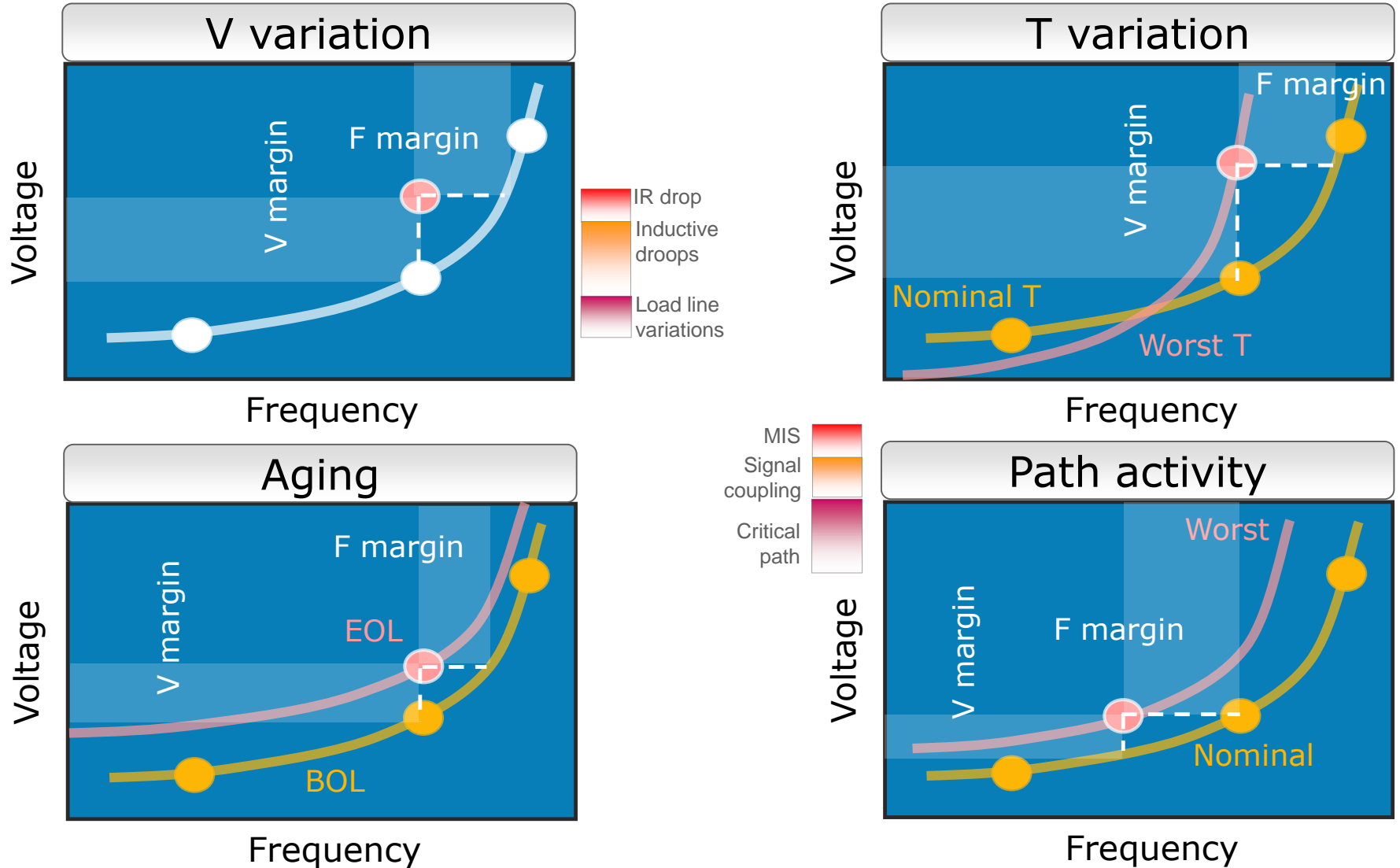
Minimum energy operation



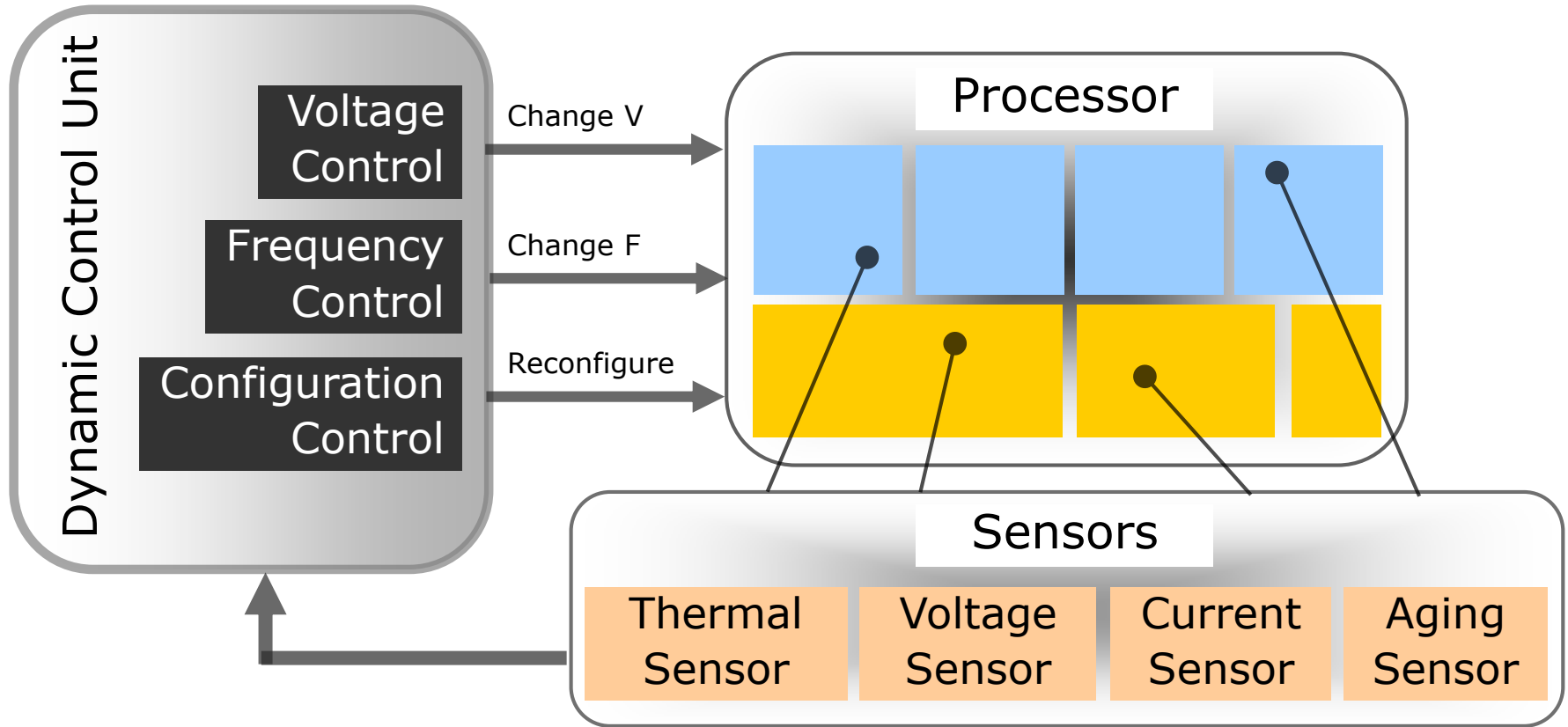
NTV and variability



Voltage-frequency margins

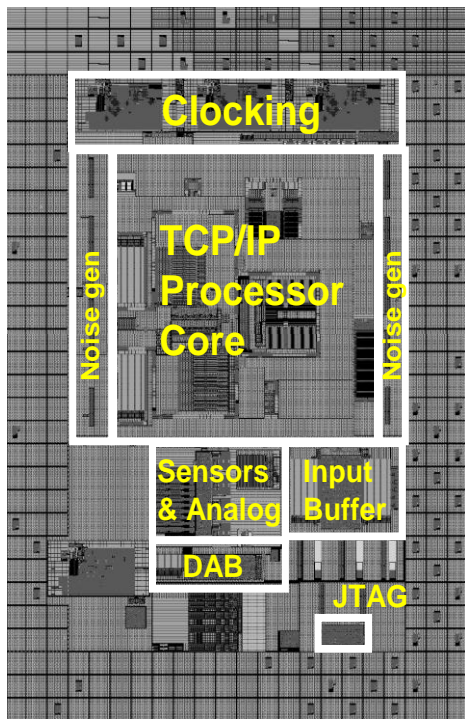


Dynamic adaptation & reconfiguration

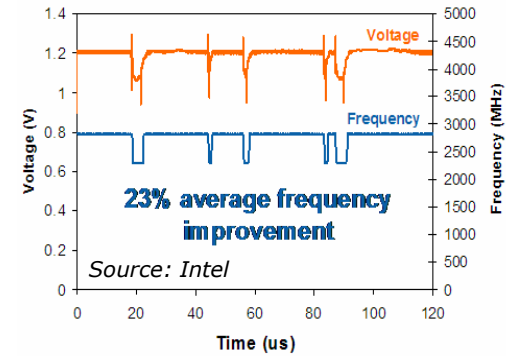
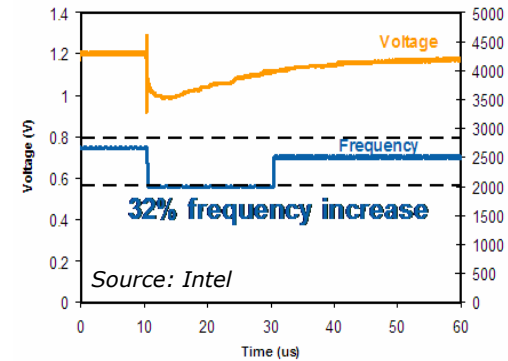
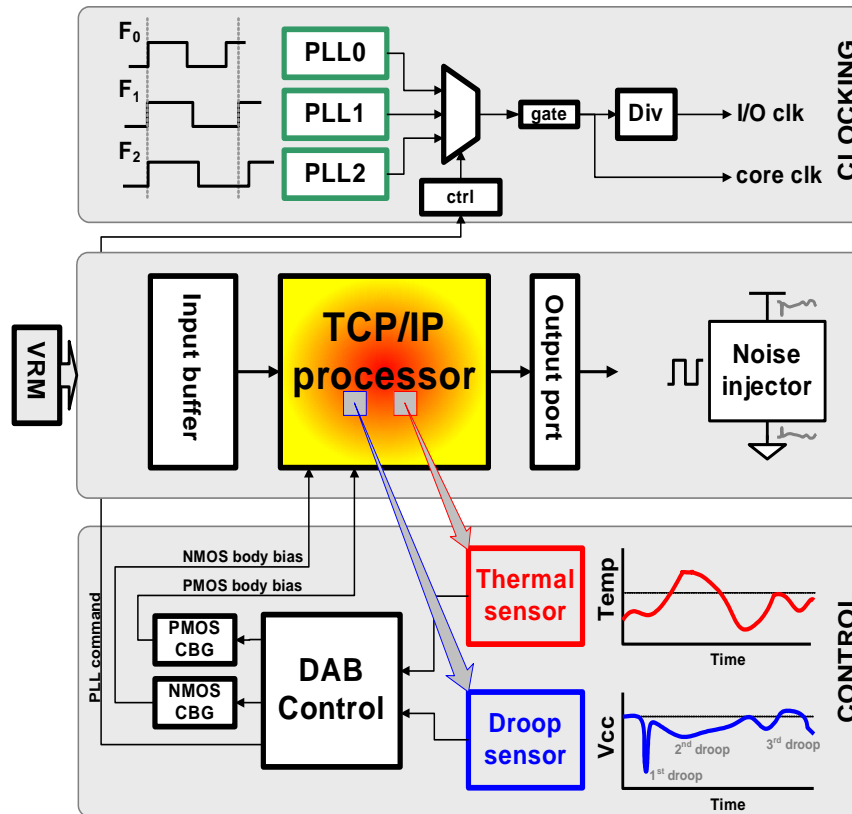


Adapt & reconfigure for best power-performance

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

Resilient platform features

Error detection

Fault diagnosis

Fault confinement

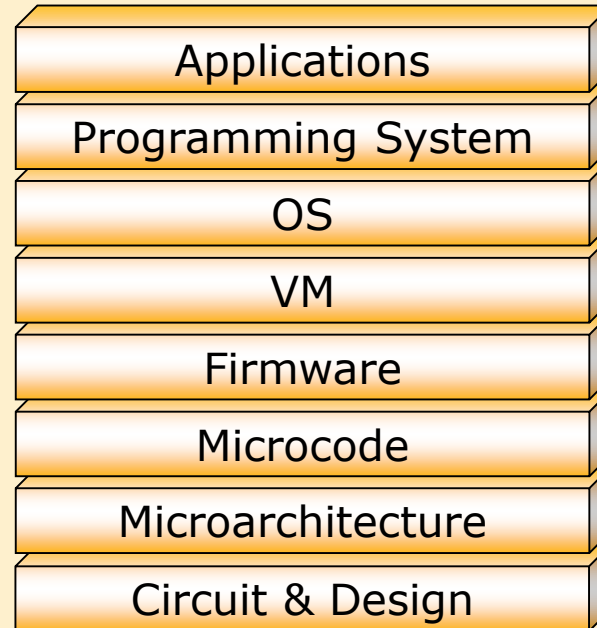
Error correction

System recovery

System adaptation

System reconfiguration

Resiliency framework



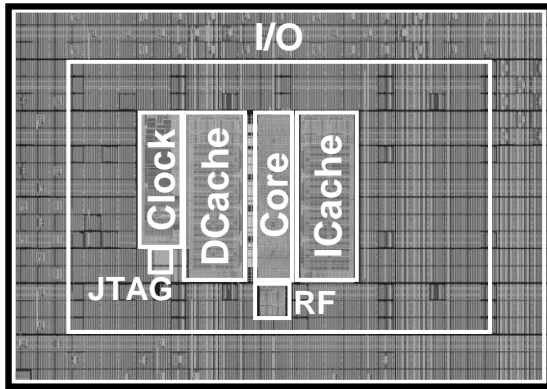
Lower error rate

Less recovery overhead

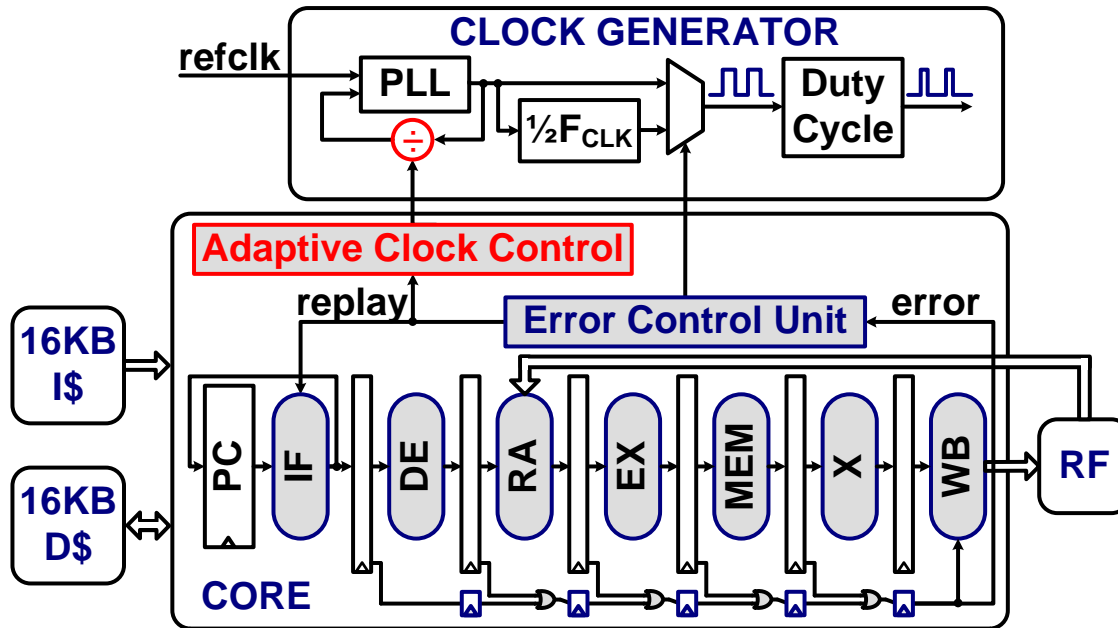
Less silicon overhead

Resiliency for performance, efficiency & reliability

Resilient & adaptive core

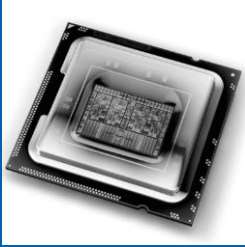


Technology	45nm CMOS
Die Area	13.64 mm ²
Core Area	0.39 mm ²
Core F _{MAX}	1.45GHz at 1.0V
Core Power	135mW at 1.0V

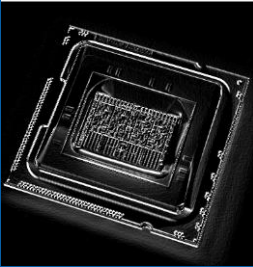


Performance & efficiency gains

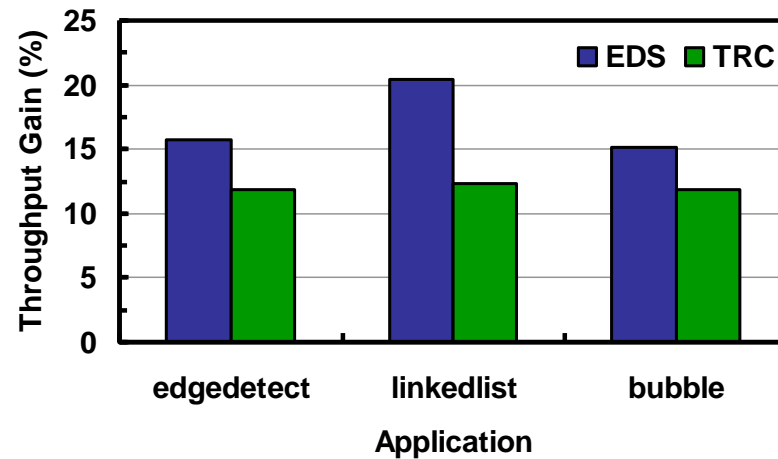
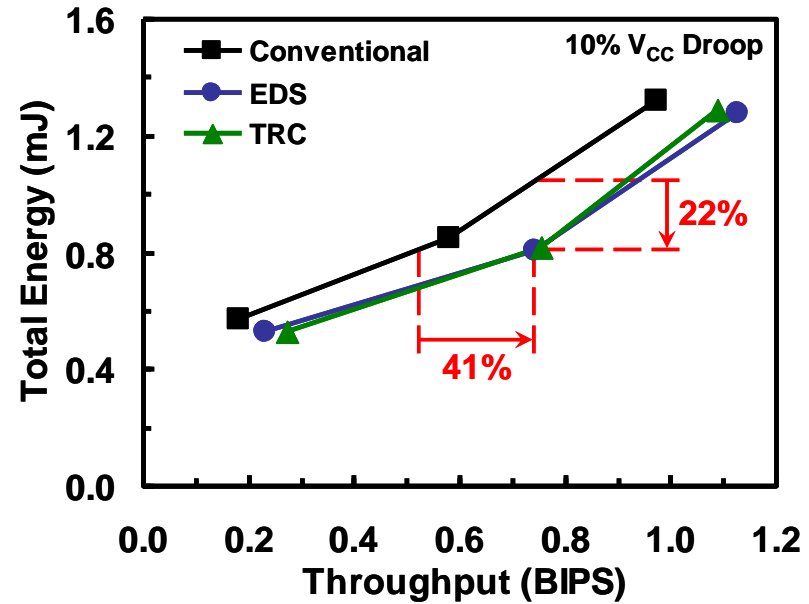
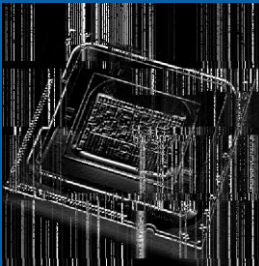
Input Image



Output Image:
Resiliency ON



Output Image:
Resiliency OFF



Integrated voltage regulators

Conversion

- Area efficient
- Scalable
- Persistent rail

Efficient

Distribution

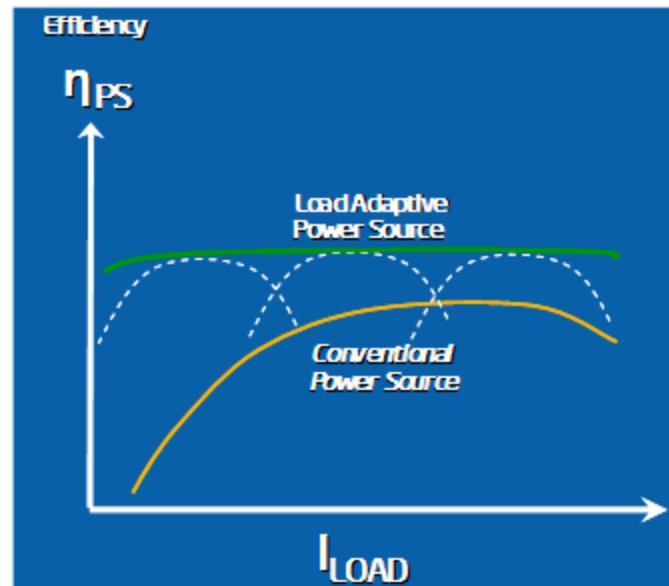
- Lower loss
- Higher fidelity
- Simpler

Low loss

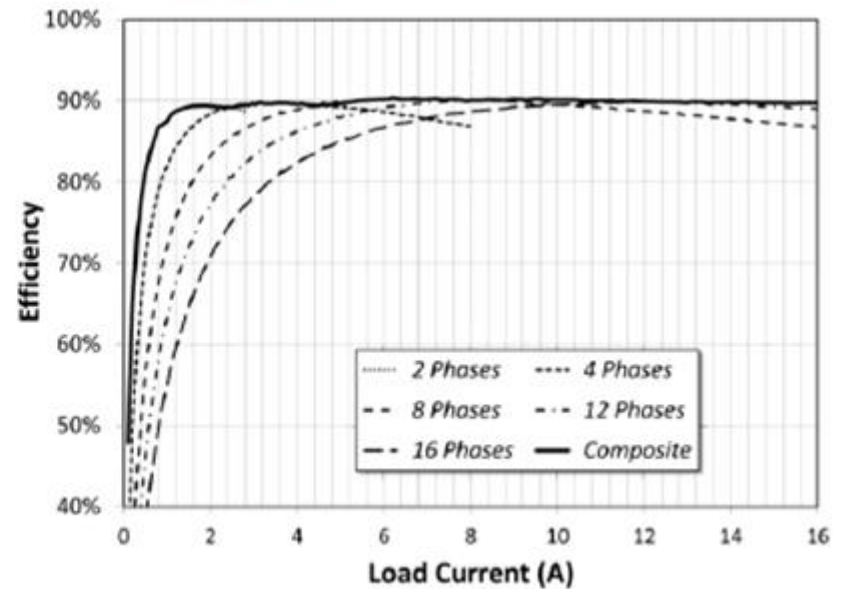
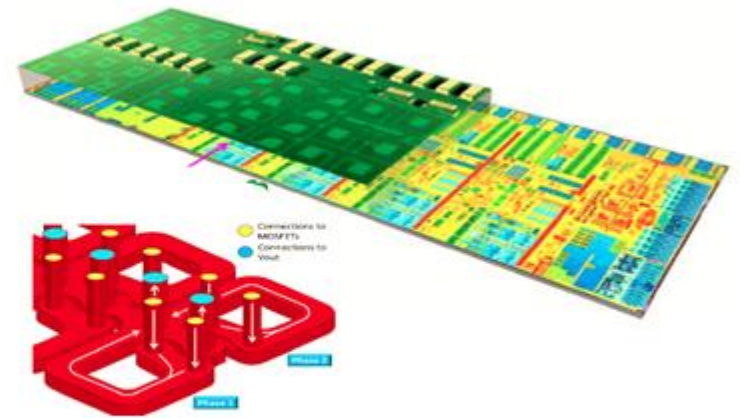
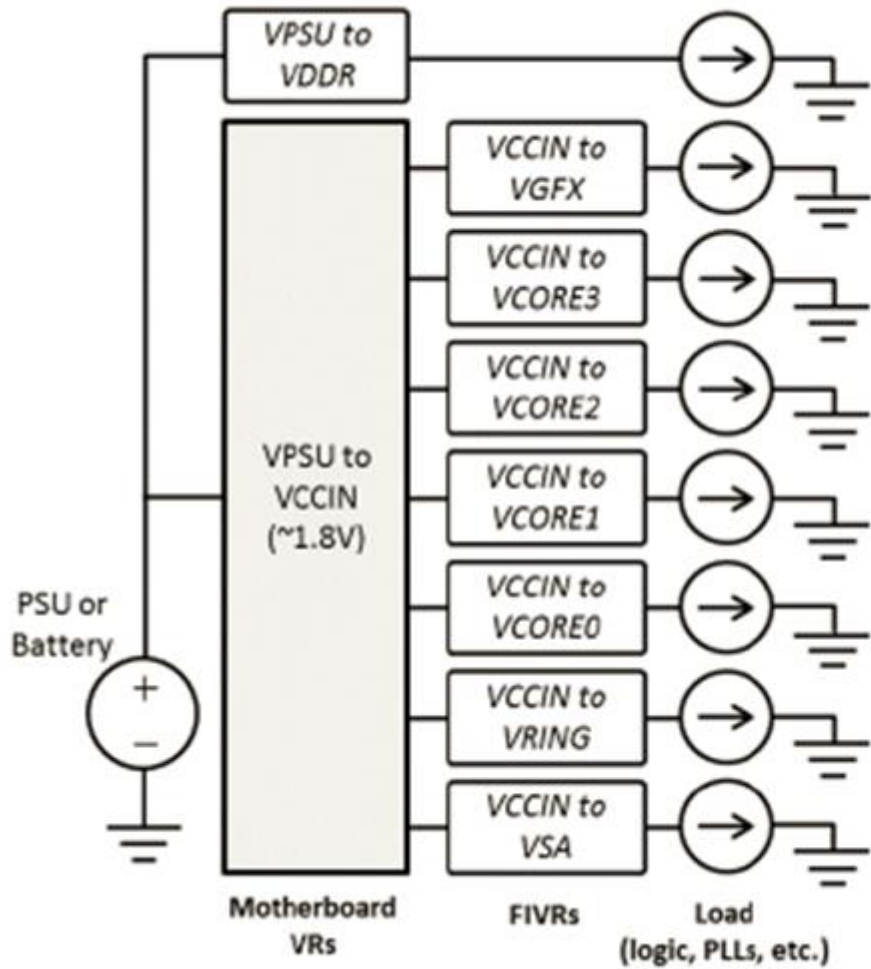
Control

- Fast & efficient
- Load adaptive
- Independent rails

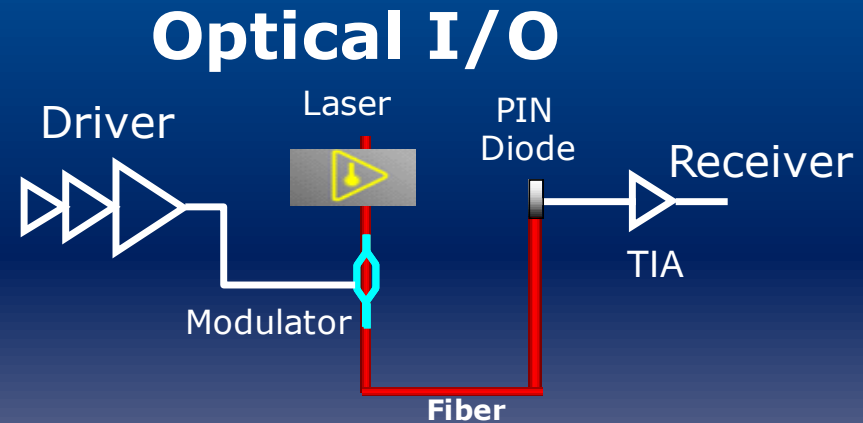
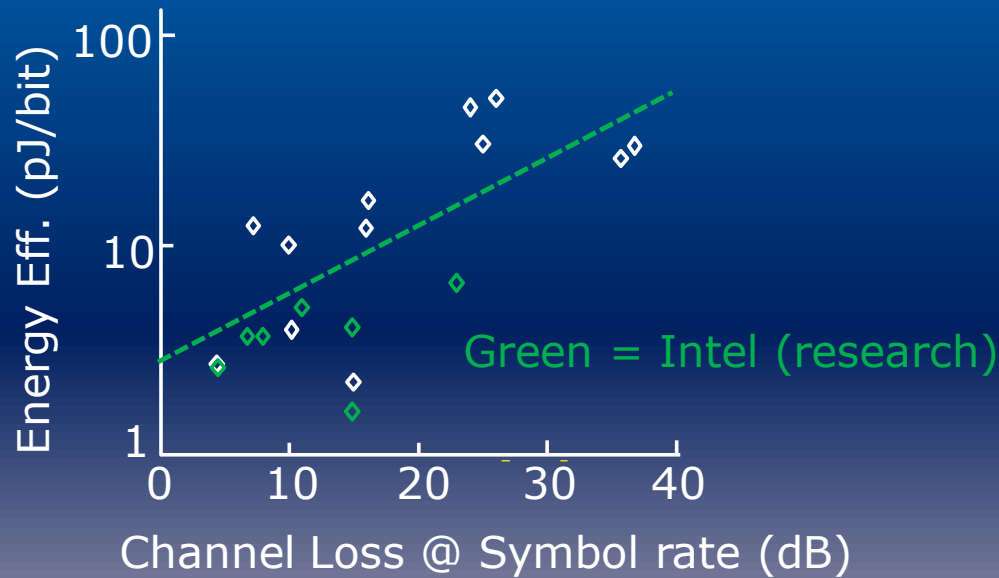
Fine-grain



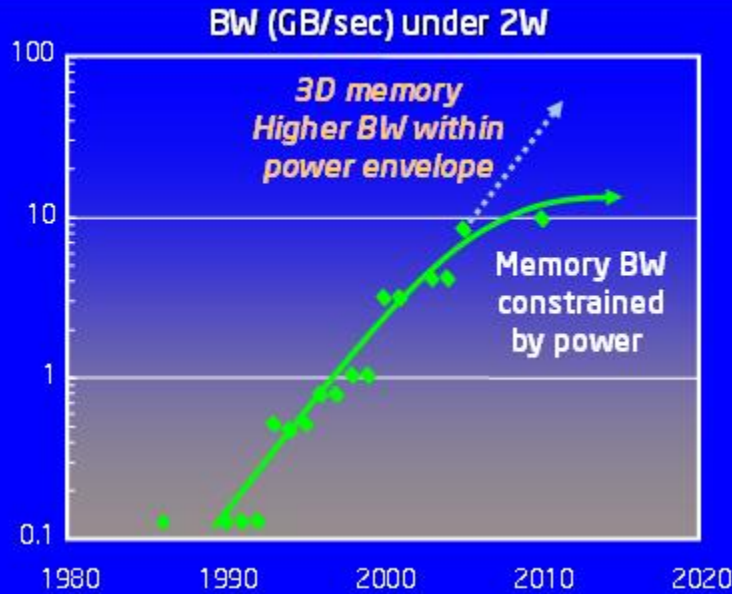
Fully integrated VR



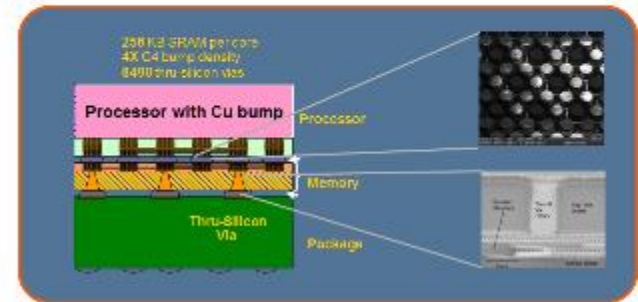
Energy efficient interconnects



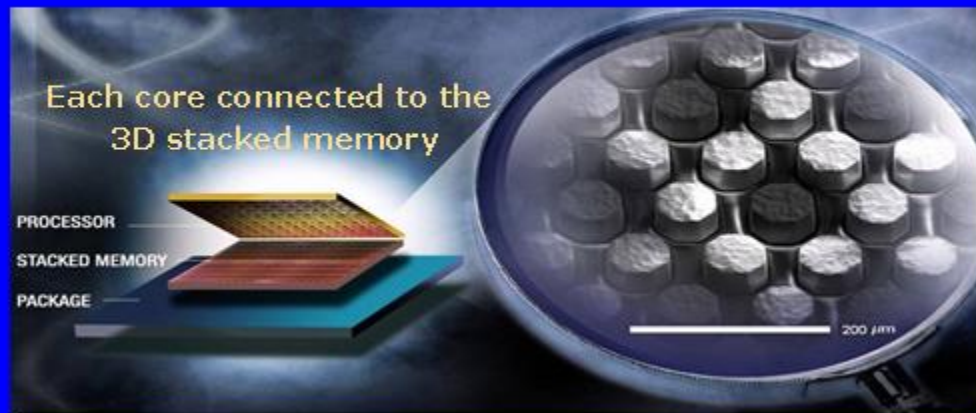
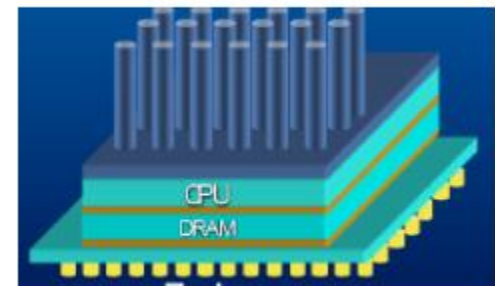
Memory capacity & bandwidth



3D Integration: SRAM

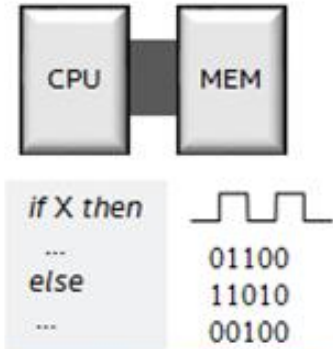


3D Integration: DRAM

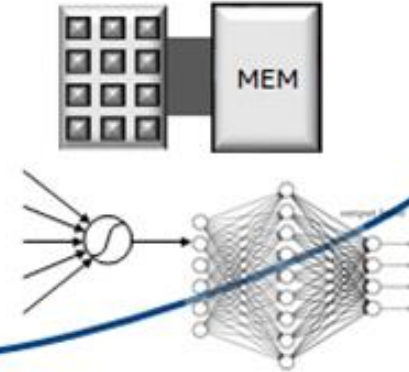


Neuromorphic computing

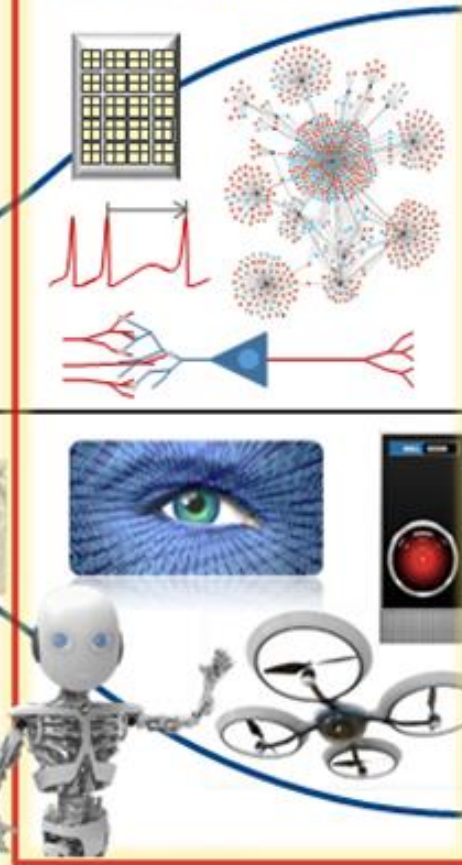
Standard Computing



Brain Inspired Computing



Neuromorphic Computing



Biological form

"Intelligent" Applications



$$\frac{d\mathcal{L}}{dx^p} = \frac{\partial \mathcal{L}}{\partial x^p} + \frac{\partial \mathcal{L}}{\partial (a_p \eta_p)} a_p \eta_p + \partial_p \mathcal{L} =$$

$$= \partial_p \left(\frac{\partial \mathcal{L}}{\partial (a_p \eta_p)} \right) a_p \eta_p + \frac{\partial \mathcal{L}}{\partial (a_p \eta_p)} a_p \eta_p + \partial_p \mathcal{L}$$

$$= \partial_p \left(\frac{\partial \mathcal{L}}{\partial (a_p \eta_p)} \right) a_p \eta_p + \partial_p \mathcal{L}$$

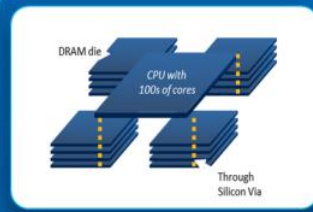
End-to-end efficiency for IoE



Extreme
Energy
Efficiency



Fine-Grain
Power
Management



Efficient
Memory
Subsystem



Self-Aware
Computing
Operation



Programming
for Extreme
Parallelism

System-Wide Breakthroughs Needed Across the Board