

The Era of Hyper Scaling in Electronics

Suman Datta

Contributors: Salahuddin, Bokor, Ramesh (Berkeley), Schlom, Jena (Cornell), Kummel (UCSD), George (Colorado), Raychowdhury, Swaminathan, Khan, Naemi, Bakir, Yu (Georgia Tech), Wong, Pop, Goodson, Wang (Stanford), Ye, Datta (Purdue), JP Wang (Minnesota), Mishra (UCSB), Iyer (UCLA), Cho (UTDallas), Winter (Wayne State), Hock (Illinois Tech), Fay, Niemier (Notre Dame)

Former Colleagues at Intel: Doyle, Kavalieros, Doczy, Metz, Dewey, Chau







- Past
- Present
- Future

Three Eras of Scaling





and hyper-scaling (future) S. Salahuddin, K. Ni, and S. Datta, Nature Electronics, Aug 2018

Geometric (Classical) Scaling



Table 1 Scaling Results for Circuit Performance		
Device or Circuit Parameter	Scaling Factor	
Device dimension t_{ox} , L , W Doping concentration N_a Voltage V Current I Capacitance $\epsilon A/t$ Delay time/circuit VC/I Power dissipation/circuit VI Power density VI/A	$ \frac{1/\kappa}{\kappa} \\ \frac{1/\kappa}{1/\kappa} \\ \frac{1/\kappa}{1/\kappa} \\ \frac{1/\kappa^2}{1} $	

10nm Transistor (circa 2002)



Physical Gate Lengths and Beyond" Intel Technology Journal, May 2002

Figure 7: Id-Vg characteristics at Vd=50mV and 0.75V for the 10nm experimental device.

Gate Voltage (V)

NIVERSITYOF

NOTRE DAME

105 (100) 085 (100)

Classical Scaling Ends





Classical scaling will not work. Need new solutions involving new materials, new device architecture and new switching mechanism

Equivalent (Effective) Scaling





3 key transistor innovations that shaped this era of effective scaling Lessons in serendipity embedded in each

Elevated Source Drain





Series resistance becomes bigger fraction of on-resistance

Partially Embedded Source Drain



Hole mobility responds in a remarkable fashion Local strain in PMOS and does not affect NMOS

NOTRE DAME Local Stressor – an act of serendipity



UNIVERSITY OF

SiGe S/D was intended for lower external resistance eSiGe S/D ends up as effective local stressor for the channel

Gate Insulator





"Transistors behave like resistors at small L unless the insulator thickness is reduced with L"

High-k metal-gate





Gate leakage increases with SiO₂ scaling -> running out of atoms High-k (HfO₂) Gate Stack enables T_{elec} scaling with low gate leakage

Gate Stack Challenges





S. Saito, et al., IEEE IEDM, Washington, DC, Dec., 2003.

Bulk & interface traps: Poor reliability New scattering modes: Poor mobility Technology Integration: Complexity and cost

Strained Si + High-k / MG stack



NOTRE DAME

Need strain + high-k + metal gate to recover mobility

S. Datta, J. Kavalieros, R. Chau, IEDM 2003

Gate-Last Enhances Channel Stress





Longitudinal compressive strain ε_{xx} is enhanced **Replacement metal gate acts as effective stress enhancers – another act of serendipity**

Multi-GateTransistor





"Even thin insulators cannot control leakage paths that are away from the gate"

From Single-Gate to Tri-Gate





Ultra thin body, double-gate or Tri-Gate architecture allows gate length scaling

Tri-gate / FinFET





Part of Tri-gate appeal is improved electrostatics Other part is **folded width**

Improved Variation – an act of serend Variation



S. Natarajan (Intel), IEDM 2014

Greg Yeric (ARM), IEDM 2015

Tri-gate was intended primarily for electrostatics **3D factor enabled effective width scaling with improved variation**

FinFET Evolution





Fin height and fin pitch evolution over last 3 generations Conductance density is increasing monotonically

Outline



• Past

- Present
- Future

Present





Present scaling relies on design technology co-design (DTCO)

S. Salahuddin and S. Datta, Nature Electronics, 2018

10nm (IDM) / 7nm (Foundry)





Feature	14 nm	10 nm
CPP	1	0.77
FP	1	0.84
MMP	1	0.51
Cell height	1	0.68
Cell width	1	0.55
Cell area	1	0.37
Transistor density (10 ⁶ mm ⁻²)	~40	~100

Ten nanometre CMOS logic technology

Through some unconventional approaches to improving transistor density and performance, the latest logic technology from Intel delivers 100 million transistors per square millimetre — and in the process, reaffirms Moore's law.

S. Datta Nature Electronics, Sept 2018

Extending FinFETs to GAA FETs





GAA NW FET is being explored extensively Hexagonal FinFETs as alternatives to FinFETs and GAA NWFETs

Beyond FinFETs – steep slope FETs



Ferroelectrics, insulator-metal phase transition materials, interband tunnel junctions are being explored for transistors with steep slope and low voltage operation

Negative Capacitance FETs





Negative capacitance (NC) FETs show short channel effect improvement over traditional FETs

Phase Transition FETs





Phase FETs show on-off ratio benefit over traditional FETs (albeit with hysteresis)

Tunnel FETs





Tunnel FETs show benefit over sub-threshold CMOS

[1] H. Zhao et al, IEEE EDL, Dec. 2010
[4] D. Sarkar. et al., Nature, Oct. 2015
[6] A. Villalon et al., VLSI 2012

[2] M. Noguchi et al., IEDM 2013

[5] L. Knoll et al., IEEE EDL, June 2013

[7] R. Pandey et al., VLSI 2015

[3] B. Ganjipour et al., ACS Nano, Apr. 2012

Steep Slope FETs





Outline



- Past
- Present
- Future

The Fourth Wave of Computing





Data intensive computing will drive the next era of hardware scaling

Source: https://www.economist.com/briefing/2015/02/26/the-truly-personal-computer https://www.statista.com/chart/12798/global-smartphone-shipments/





- Architecture: Ability to scale appropriately to meet increased workload demand
- Technology: Ability to scale beyond x-y dimension to meet system specifications of power, performance, area, function, form-factor and cost
 - Monolithic integration
 - Pseudolithic integration
 - Merged logic-memory fabrics
 - Neuro-inspired computing fabrics



Applications and Systems Driven Center for Energy-Efficient Integrated Nanotechnologies

ascent.nd.edu





ASCENT is sponsored jointly by the Semiconductor Research Corporation (SRC) and DARPA

Monolithic 3D





Levels of granularity:

- Die level
- Block level
- Gate level
- Transistor level

Grand Challenges:

- Protect bottom layer transistors
- □ Fabricate upper layer transistors
- □ Align top layer with bottom layer
- Fill high AR inter-layer vias with low resistivity metals
- Embed thermal management structures
- □ Assess system performance (ROI)

Monolithic 3D

Stanford)





Pop, Goodson (Stanford)

Heterogeneous Integration & Advanced Packaging





Why Heterogeneous Integration

- PCB ultimately limits size, weight, area, performance (SWAP) of microsystems
- Enables silicon IP reuse, reduces design cost and improves system performance

Grand Challenges

- Target 1um interconnect spacing and <20 um die-to-die spacing</p>
- Target aggregate data transfer rate of 1,000 Gb/s/mm at < 0.1pJ/bit over 50-500 um distance

Heterogeneous Integration Fabric





Beyond CMOS (spin)





Grand Challenges Efficient Write

- □ Spin-orbit torque (SOT)
- Magneto-electric (voltage controlled magnetism)
- Voltage controlled magnetic anisotropy (VCMA)

Efficient Read

- MTJ free read-out of magnetic information
 - (e.g. Inverse Rashba effect)

Beyond CMOS (Merged logic-memory)



NIVERSITY OF

NOTRE DAME

100 000 000 000

JP Wang, Salahuddin (Berkeley)

CoFeB

Ramesh (Berkeley)

Hardware for Al





Al hardware addresses different applications with appropriate energy efficiency and performance, but far from that of brain

Neuro-inspired Computing Fabric





Von Neumann



Brain

- Massively parallel architecture supported by dense connectivity
- Suitable for life long learning and decision making in changing environment
- □ Co-located logic and memory
- Based on Markov decision processes and resilient to instantaneous errors

□ Fully programmable

- Suitable for high precision computing
- Separate logic and memory
- Limited by memory bandwidth where read/writes are performed at course granularities

Neuro-inspired Computing Fabric



Artificial Neural Networks (ANNs)

Hidden Input Output

Digital Memory-centric **Synchronous**





Analog **Primitives Time-encoding** Asynchronous

Non-filamentary **RRAM Synapse**





FerroFET Analog Synpase



IMT Neuron



FeMFFT Synapse



Ferro Neuron



Raychowdhury, Yu, Khan (Georgia Tech)

Datta, Niemier (Notre Dame)

Salahuddin (Berkeley)

Key Takeaway





Lessons from the Past: Serendipity and not luck. It's being in the right place at the right time with the right training and mindset

Present state of Moore's Law: Not a Law of Physics. It's hard work by scientists and engineers, our ingenuity and audacity to compete and thrive, driven by free market economics

Future of Electronics: Materials scientists, device physicists, circuit designers, chip architects will collaborate to provide semiconductor innovations that provide system level benefit

"Such an exponential will continue for a very long time"

Thank You



