Dinesh Pamunuwa*, Matt Grange**, Axel Jantsch+, Sunil Rana*, Tyson Tian Qin*

*University of Bristol, Bristol UK

**Mentor Graphics, USA

+KTH Royal Institute of Technology, Stockholm, Sweden

# System Performance Analysis for Heterogeneous 3-D ICs and Emerging Technologies
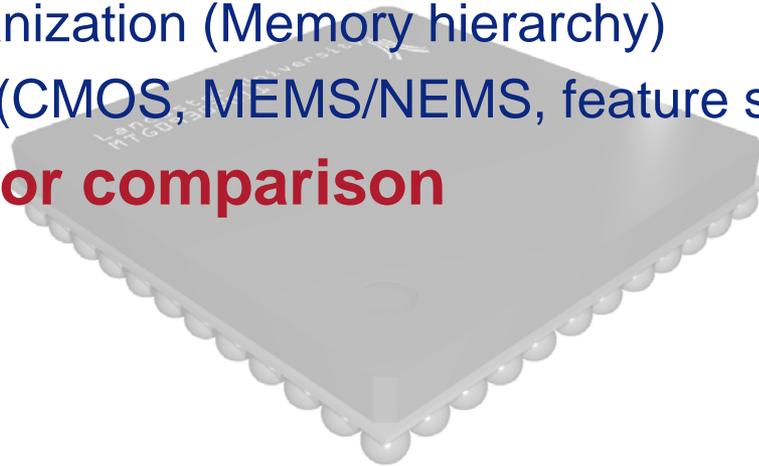
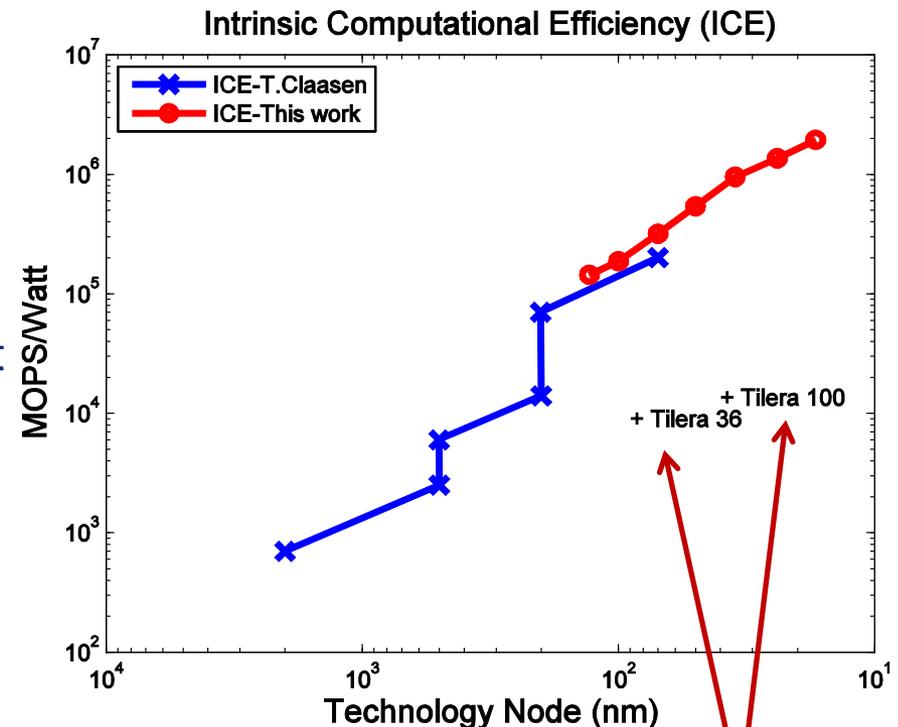University of BRISTOL

microelectronics research group µ

# Future Direction of 3-D ICs

❑ **3-D ICs promise performance / cost advantages for high performance digital applications as well as enhanced functionality**

❑ **A staggering amount of degrees of freedom exist**

   – Packages (2-D, Multi-chip Modules, 3-D die stacks)

   – Digital architecture (single to multi to many-core)

   – Routing architectures (buses, Networks-on-Chip, hybrids)

   – System organization (Memory hierarchy)

   – Technology (CMOS, MEMS/NEMS, feature size reduction)

❑ **Framework for comparison**

**microelectronics research group**

University of BRISTOL

# Intrinsic Computational Efficiency

❑ **The Intrinsic Computational Efficiency (ICE), proposed by T. Claasen, creates the maximum upper bound for computational capability of a silicon-based processor.**

– The entire silicon area of a processor is filled with the most fundamental computational unit, in this case we have used 32-bit adders.

– A real system could never achieve the same performance per Watt because this metric <u>ignores the overhead of control circuitry, interconnect, and memory</u>.



Intrinsic Computational Efficiency (ICE)

Two recent multi-core processors

**microelectronics research group**

University of BRISTOL

# Expanding the *ICE* to the *ECE*

❑ **The *ICE* gives the *maximum upper bound* on efficiency, but cannot account for realistic systems because it only considers the computational unit.**

❑ **We build upon the *ICE* by modelling the three fundamental operations of any processing unit:**
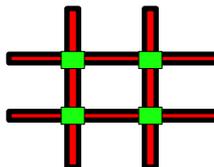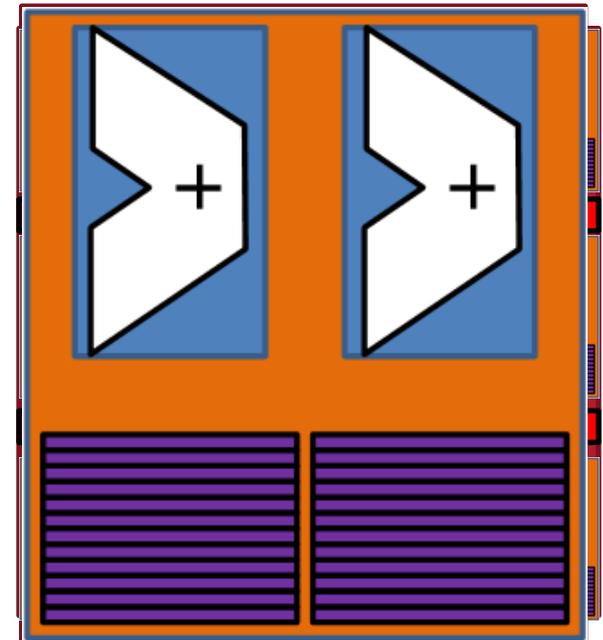
**Processing Core**

− **The computational operation**

− **The memory**

− **The interconnect**

**microelectronics research group**
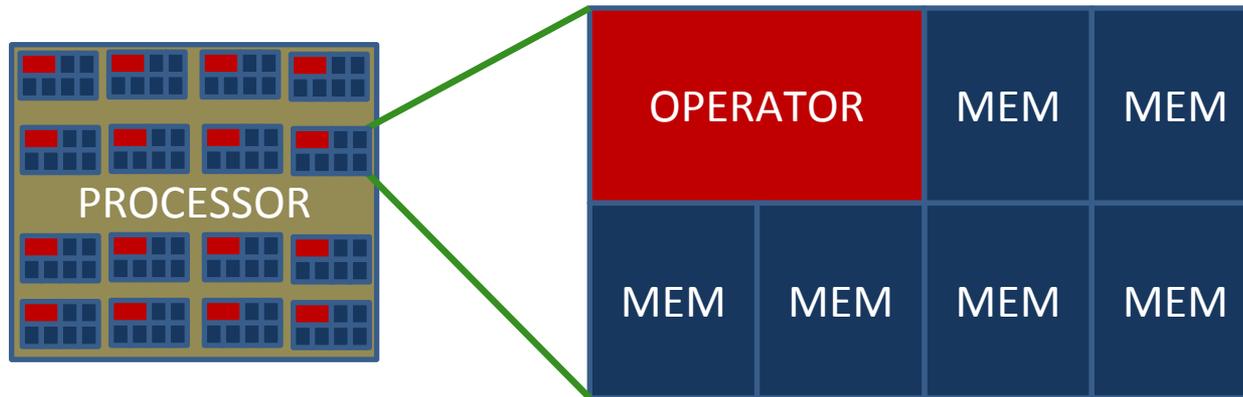
University of BRISTOL

# Temporal and Spatial Organization of Memory

□ **$\mu_s$**

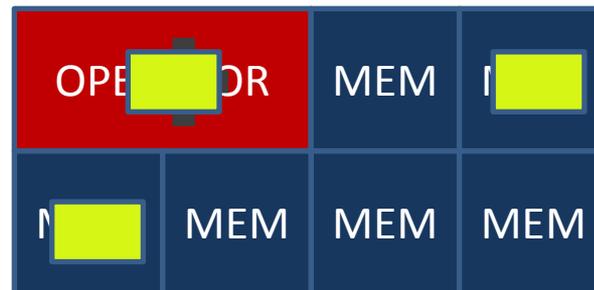   – **Gives us the amount of memory per operator. Think of it as the amount of on-chip cache available**



□ **$\mu_T$**

   – **Gives us the number of memory reads/writes per operation.**

2 reads, 1 write

$\mu_T = 3$

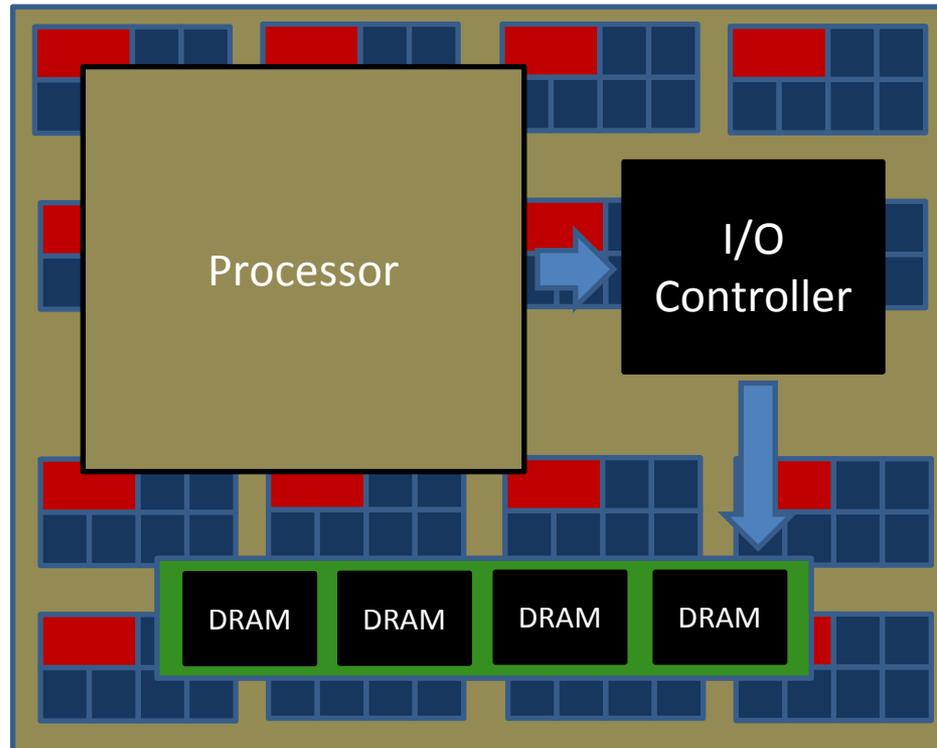**microelectronics research group**

University of BRISTOL

# On- or Off-chip Memory?

- ❑ **ω (0-1)**
  - – **Gives the ratio of on- to off-chip memory in the system. Off-chip memory requires exiting the die with I/O drivers to external chips.**
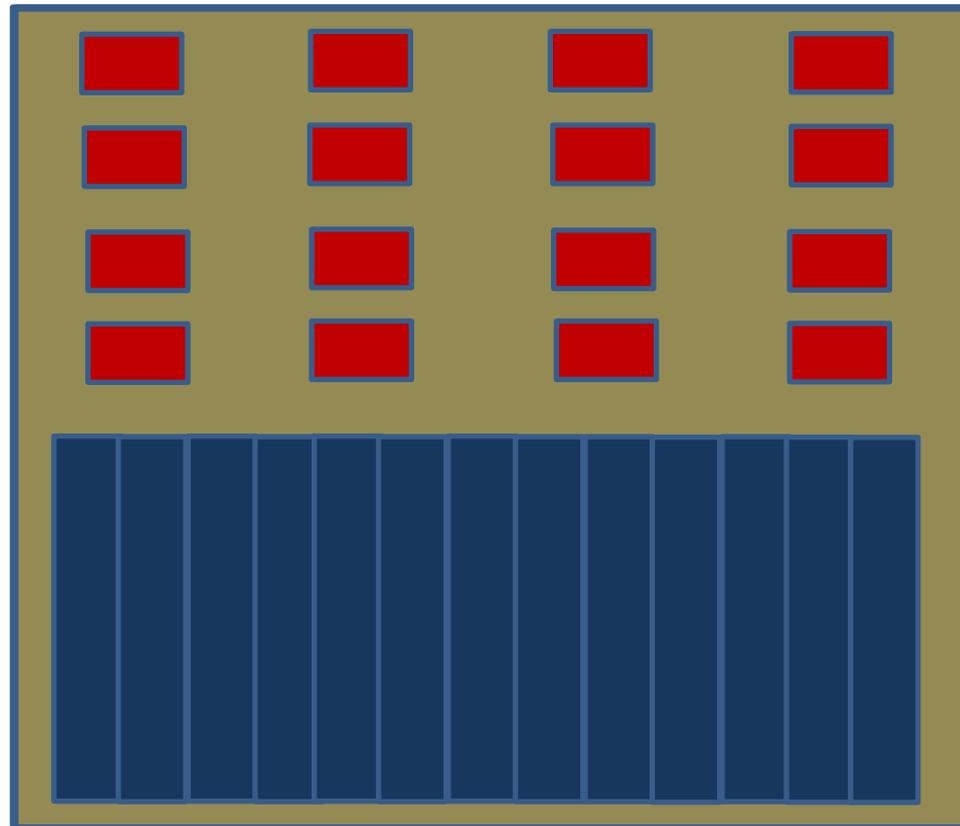
**0: all off-chip**

**microelectronics research group**

University of BRISTOL

# Memory Distribution Factor

❑ **Δ (0-1)**

– **Gives the distribution factor of the memory, or how close (local) it is the operator.**

**0:5abemi-local**

**microelectronics research group**

University of BRISTOL

# Effective Computational Efficiency

❑ **For each computation we must consider the expense of energy for the operation, interconnect (on and off-chip) and memory reads/writes.**

**Number of Memory Accesses/op**

**Energy for a 32-bit addition**

**On-chip memory energy**

**Off-chip memory energy**

$$EE^{tn}_{arch} = E^{tn}_{32} + \mu_T \left( \omega \left( e_1 + \Delta \times E\_int^{tn}_{arch} \right) + (1-\omega)\left( e_1 + E\_int^{tn}_{arch} + E_{offchip} \right) \right)$$

**Ratio of on-chip : off-chip memory**

$$ECE^{tn}_{arch} = \frac{1}{EE^{tn}_{arch}}$$

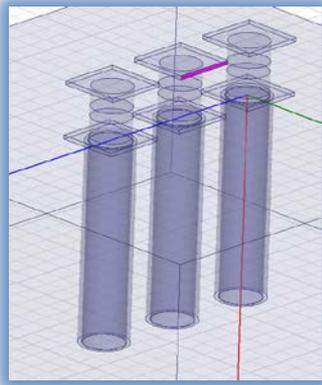← **Amount of computation possible within the envelope of 1 Joule**

# Modelling Hierarchy



System

Application

Architecture

Circuit

Physical

Cost

Perf

**Computation**

Traffic | Activity | Access

NoC, Bus | Arith unit | Mem org

Signalling | Gates | Cells

Parasitics | Devices | Devices

Thermal

**Interconnect** | **Logic** | **Memory**

**microelectronics research group**

University of BRISTOL

# Model Development

**Parasitic Extraction (Q3D)**

*Resistance (**R**), Inductance (**L**)*
*Conductance (**G**), Capacitance (**C**)*

**Geometrical Sweeps (length, radius, pitch, etc.)**
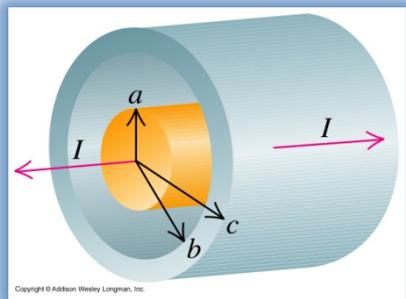
**Physical Sweeps (Topology, Substrate, etc.)**

$R \quad L$
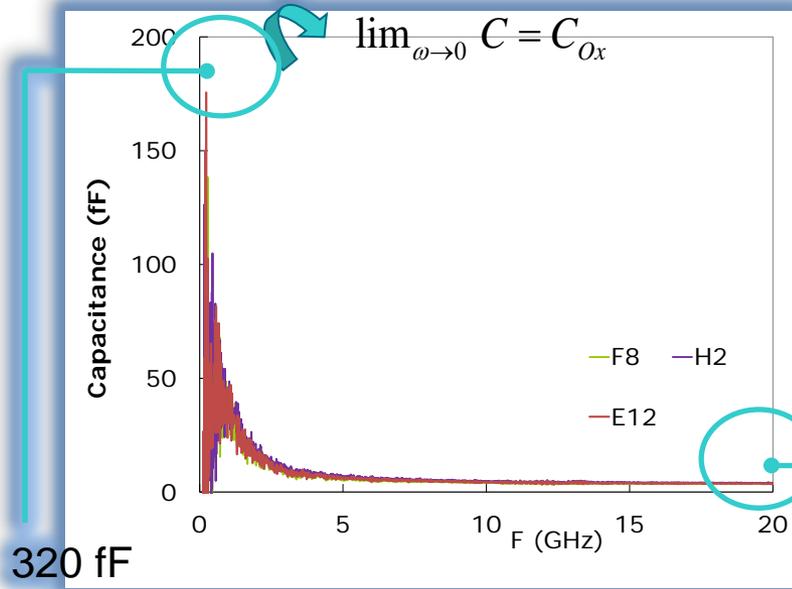
$G \quad C$

**Data Manipulation (MATLAB)**

**Modelling**

**Physical Principles**

**Verification**

$a$

$I$

$I$

$c$

$b$

Copyright © Addison Wesley Longman, Inc.

**microelectronics research group**

University of BRISTOL

# Model Verification

$$\lim_{\omega \to 0} C = C_{Ox}$$

Capacitance (fF)

200

150

100

50

0

0    5    10    15    20

F (GHz)

F8 — H2

E12

320 fF

(306.6)

4 fF

(4.5)

C — G  →  $C_{Si}$  $C_{ox}$  $G_{Si}$

DP, D43D Jun 2013        **microelectronics research group**

University of BRISTOL

# Stand-alone Parasitic Estimation Tool

# Modelling Hierarchy

**microelectronics research group**

University of BRISTOL

# Thermal Behaviour

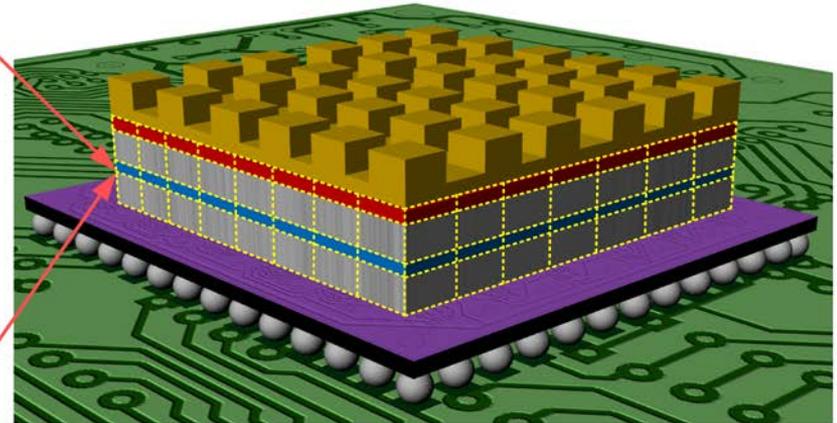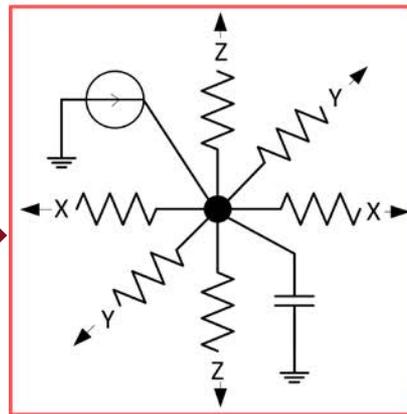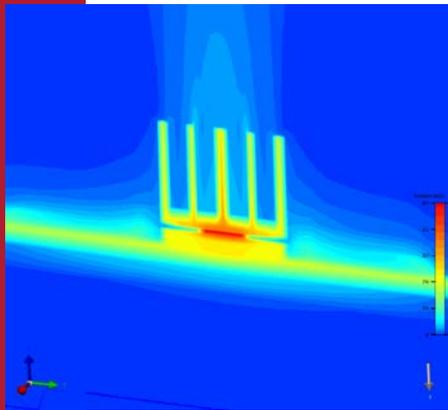- ❑ **Compact thermal models have been developed as part of the toolset to quickly predict the thermal behaviour.**
  - – Verification based on comparison with results from a Computational Fluid Dynamic (CFD) solver (FloTHERM).
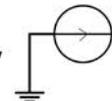  - – The thermal behaviour and limitations of 2-D and 3-D packages has been extracted from simulations.



Voltage at junction equivalent to temperature

Material Thermal Resistance

Material Thermal Absorption Capability (specific heat)
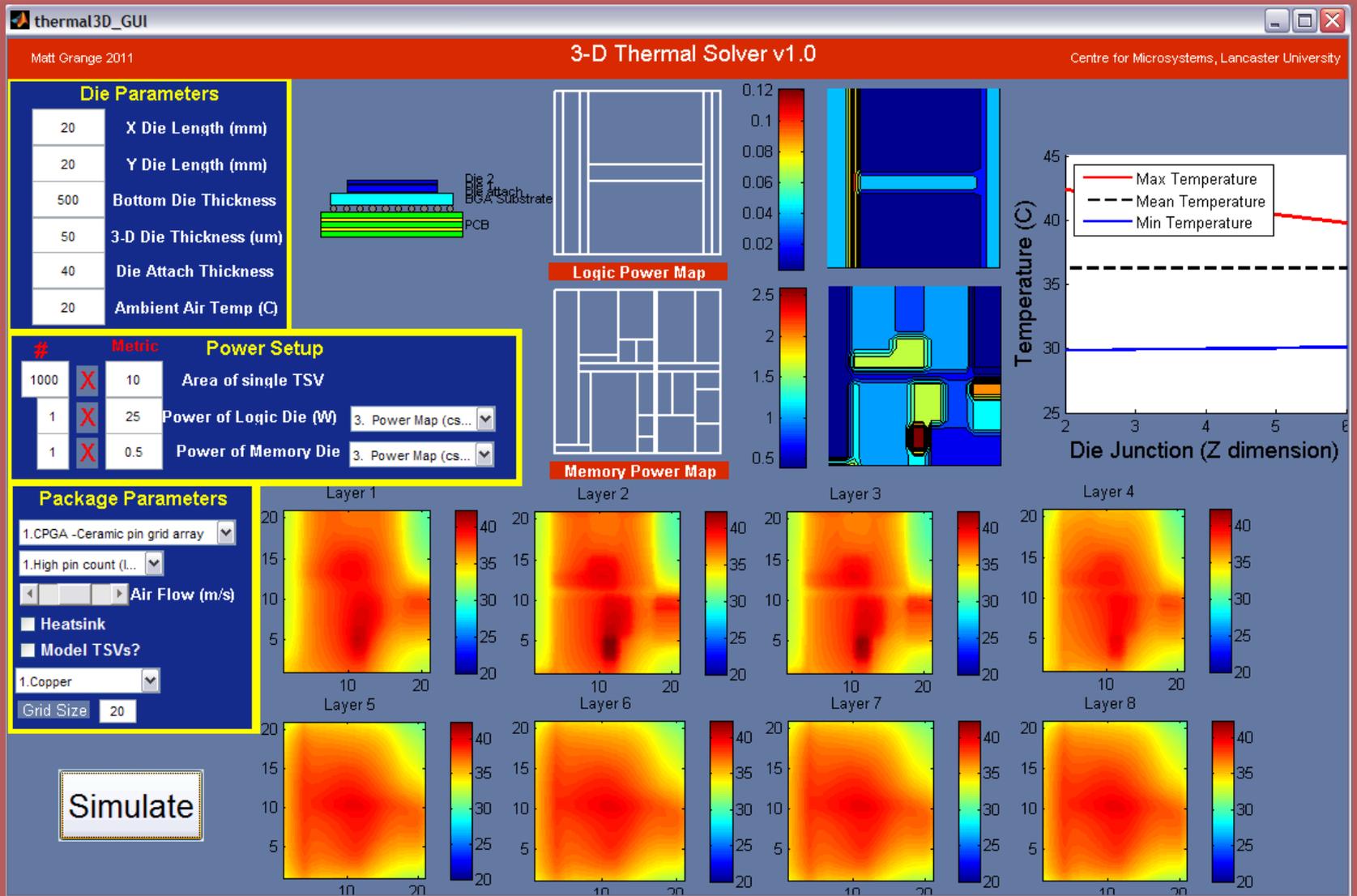
Power Injected into junction

$$R_{th} = \frac{L}{kA} \qquad T_j = qR_{th} + T_a \qquad \theta_{j-a} = \frac{T_j - T_a}{q}$$

University of BRISTOL

# Stand-alone 3-D Thermal Tool

University of BRISTOL

# Modelling Hierarchy



System — Cost, Perf, **Computation**

Application — Traffic, Activity, Access

Architecture — NoC, Bus; Arith unit; Mem org

Circuit — Signalling, Gates, Cells

Physical — Parasitics, Devices, Devices, Thermal

**Interconnect** **Logic** **Memory**

**microelectronics research group**

University of **BRISTOL**

# Scaling 2-D versus 3-D DRAM

## 2-D with off-chip DRAM

## 2-layer 3-D with in-stack DRAM

**microelectronics research group**

University of BRISTOL

# Fixed System: Intel 80 Core

| Param. | Type | Intel 80 Core | Equivalent |
|--------|------|---------------|------------|
| N | # of Layers | 1 (2-D) | 1-16 (3-D) |
| A | Die Area | 12.64×21.72 mm | 275 mm2/N |
| tn | Tech. node | 65 nm | 180-17 nm |
| b | Data width | 32-bit | 32 |
| µs | Memory/Operator | 2K SRAM/2 FPU | 1KB/Op |
| µt | Memory/Operation | App. Specific | 01-Mar |
| σ | Bus Sharing Ratio | 8x10 mesh/160 FPU | 18/160=0.11 |
| Δ | Memory Distribution | NoC Mesh | 0.01-0.1 |
| ω | On/off-chip mem | All on-chip | 1 |
| P | Power (W) | 20-230 | App. Specific |

**Implemented processors can be modelled by varying the parameters**

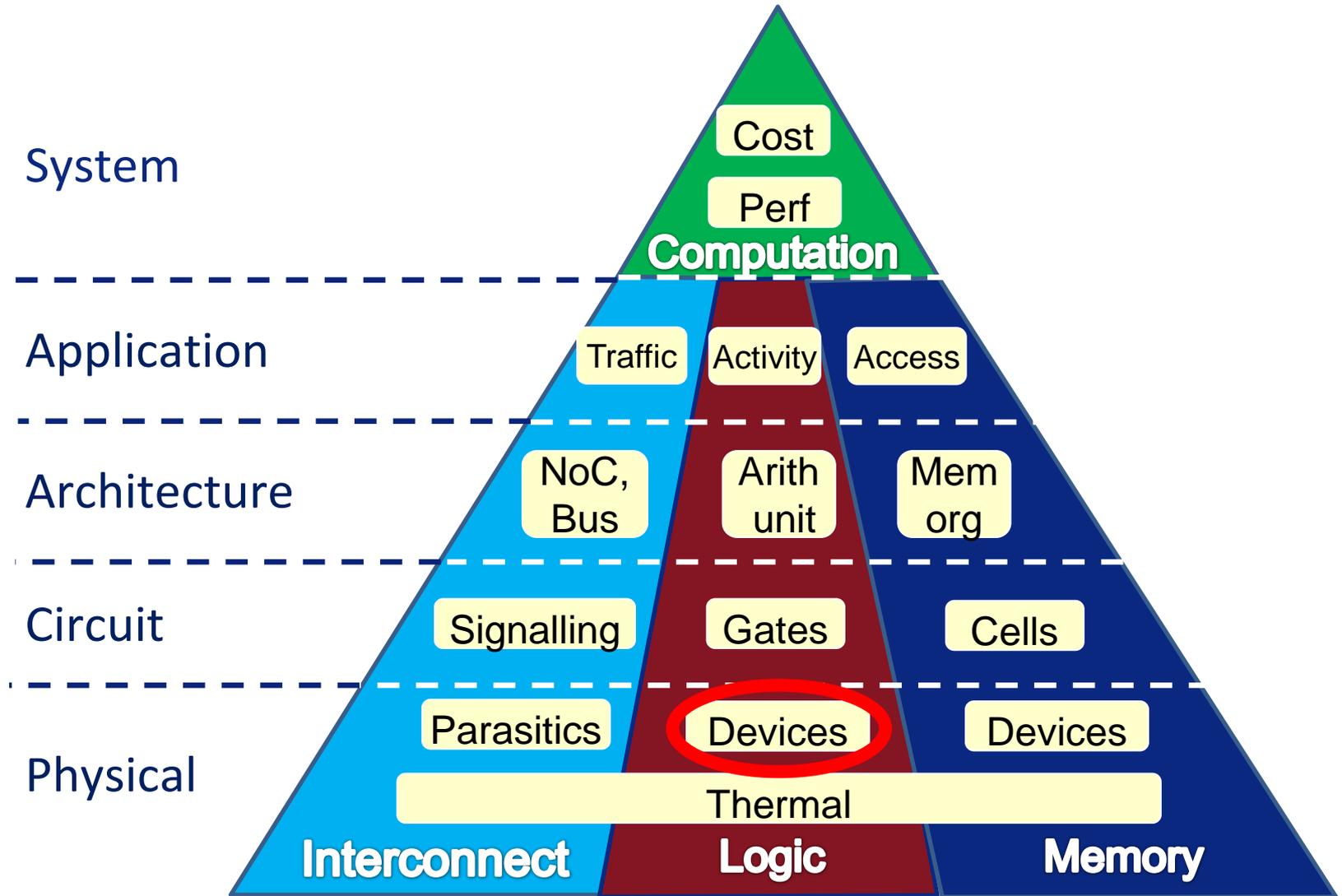**The 4-layer 80-core 3-D system at 90 nm is still better than a 2-D system at 65 nm**

**For every doubling of the stack height the computational efficiency increases by 20-30%**



Effective Computational Efficiency
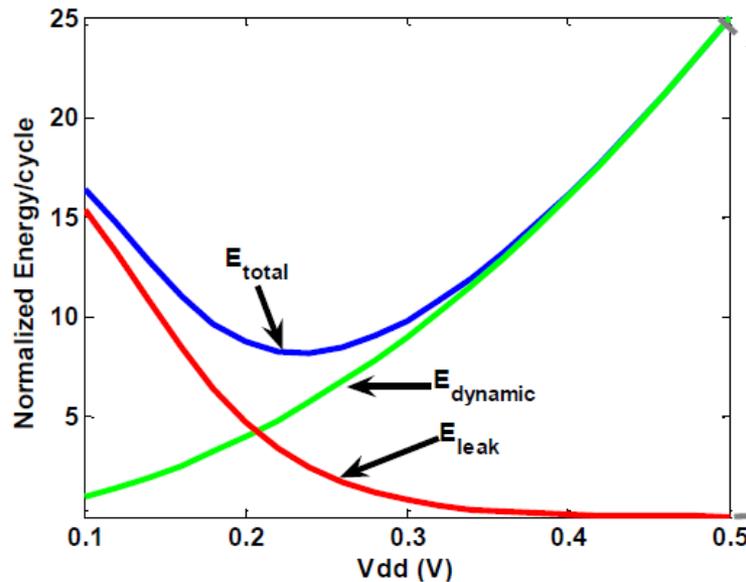
- 2-D model of the Intel 80 core
- 3-D2 model
- 3-D4 model

A 4-layer partitioned Intel 80 Core @ 90 nm achieves similar GOPS/Watt as 2-D implementation in 65 nm

Intel 80 Core @ 65 nm

GOPS/Watt

Technology Node (nm)

University of BRISTOL

# Modelling Hierarchy

**microelectronics research group**
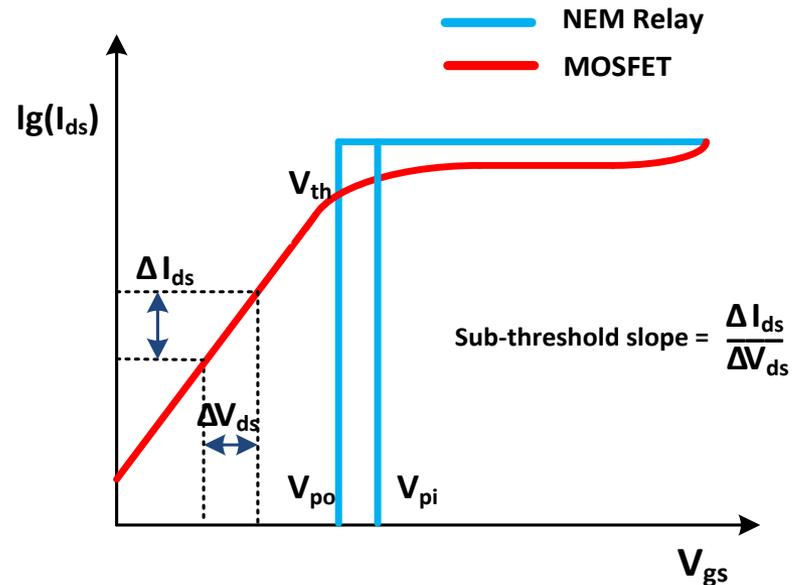
University of BRISTOL

# NEM Relay Based Computation

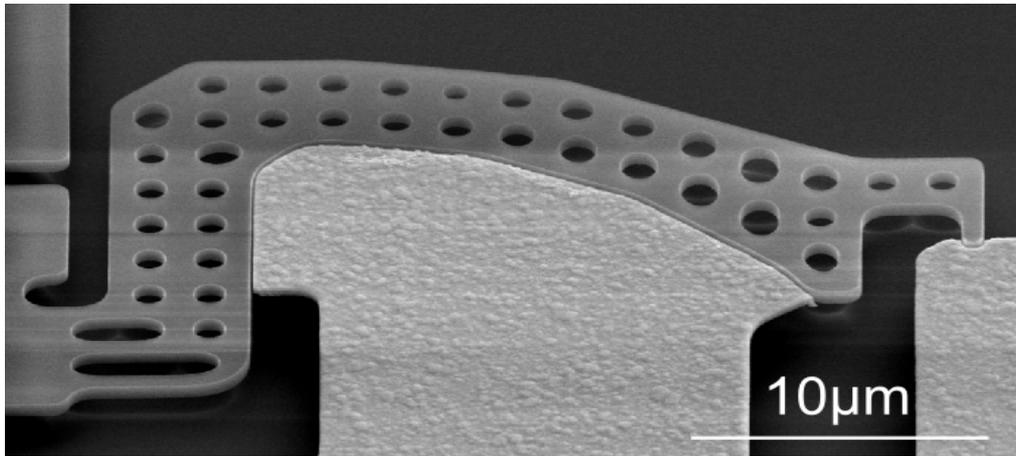## ❑ Limitation of CMOS Energy Efficiency



(**source: F. Chen et al, ICCAD 2008**)

## ❑ NEM relay advantages

– practically zero leakage

– Very steep slope for turn-on/-off transient

– high on-current

University of BRISTOL
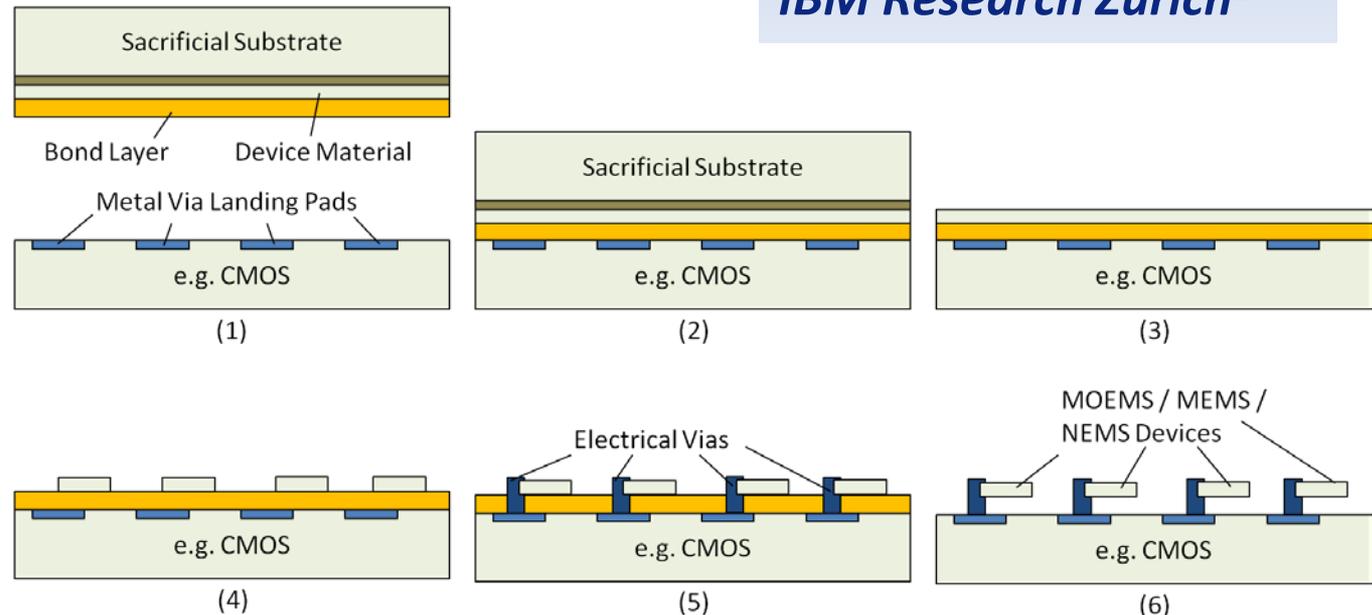
# NEM Relay Technology



10μm

**In-plane switch, fabricated using standard lithography in NEMIAC project; nm gap using sacrificial layer**
*courtesy D. Grogg et al. IBM Research Zurich[1]*

**Long term integration plan**
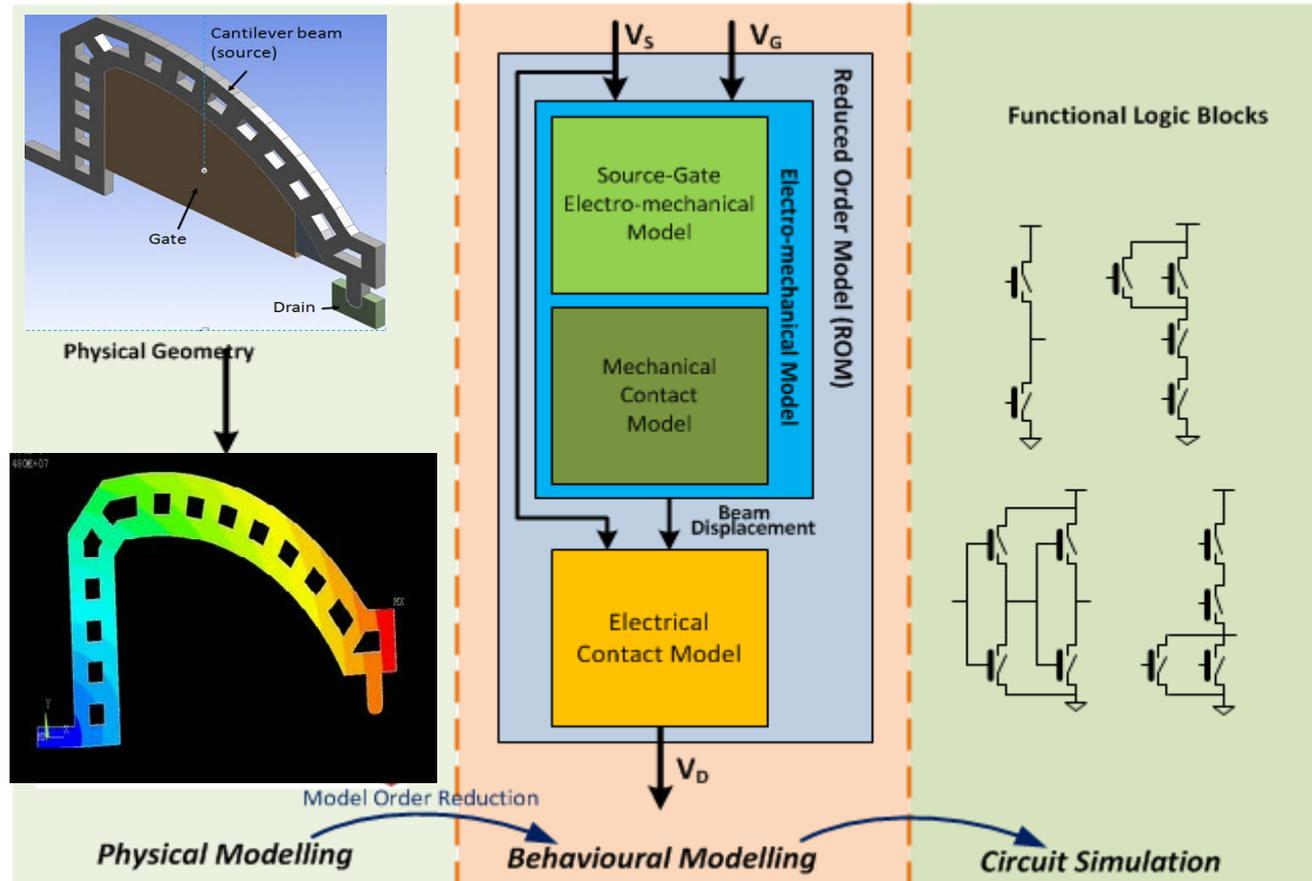*courtesy F. Niklaus et al. KTH, Sweden*



1. D. Grogg et al., "Curved cantilever design for a robust and scalable microelectromechanical switch," *Int. Conf. on Electron, Ion, and Photon Beam Technology and Nanofabrication*, Waikoloa, Hawaii, 2012.

University of BRISTOL

# NEM Relay Modelling

## ❑ FEA to Reduced Order Model to Circuit Model

**microelectronics research group**

University of BRISTOL

# ICE with NEM Logic

- **NEM "tech node" is much larger than CMOS**
  - Devices fabricated at 17 µm and 5 µm cantilever length

- **Miniaturisation increases speed and reduces energy**
  - No straightforward analogue to scaling for CMOS

- **Trajectory is promising**
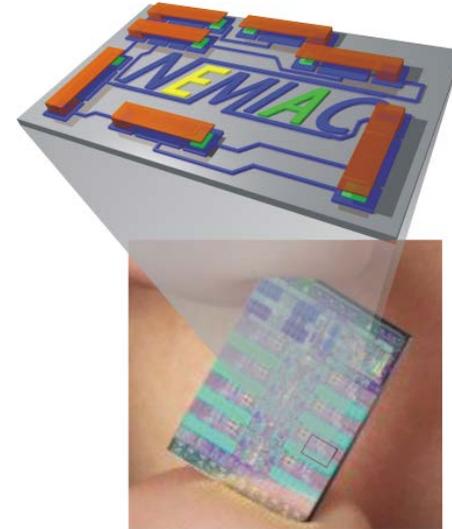  - Ultra-low power technology for low latency, low throughput applications

**Intrinsic Computational Efficiency (ICE)**



- ICE-T.Claasen
- ICE-This work
- + NEM 4 bit
- + NEM 8 bit
- + NEM 32 bit (approx)
- NEM 32 bit (approx)
- + Tilera 100
- + Tilera 36

MOPS / Watt vs Technology Node (nm)

> Ripple-carry architecture and worst-case energy
> Device models at 17µm and 5µm silicon qualified
> 4 and 8- bit adder energy at 17 µm based on accurate circuit model
> 32-bit adders based on scaling

University of BRISTOL

# Acknowledgements

## ❑ NEMIAC project

– EU FP7 Strep: Grant No. 288670



## ❑ ELITE Project

– EU FP7 Strep: Grant No. 215030

**microelectronics research group**

University of BRISTOL