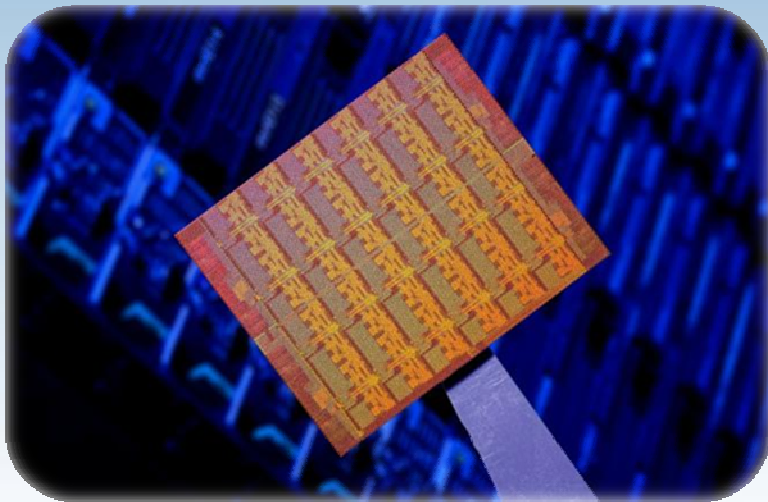


“A 48-Core Single Cloud Computer”

Intel’s experimental many-core processor



Jason Howard
Advanced Microprocessor Research
Intel Labs

Intel Labs
June 8, 2010



Content

- Feature set
- Architecture overview
 - Core
 - Memory Model
 - Interconnect Fabric
 - I/O and System Overview
- Design Overview
 - Clocking
 - Message Passing
 - Power management
- Silicon Results
- Programming
- Summary

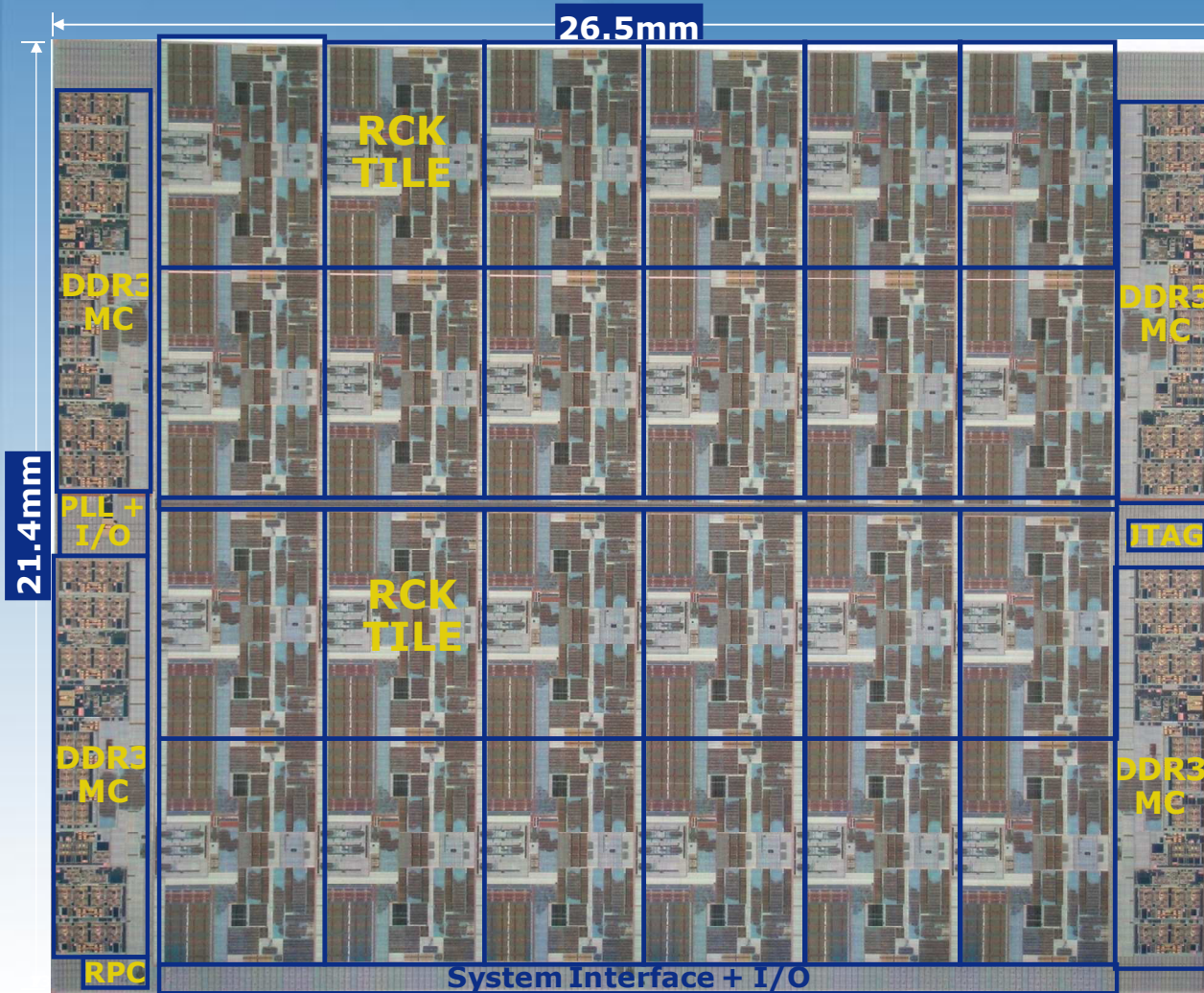


SCC Feature set

- First Si with 48 iA cores on a single die
- Next generation 2D mesh interconnect
 - Bisection B/W 1.5Tb/s to 2Tb/s, avg. power 6W to 12W
- Power envelope 125W
 - Core @ 1GHz, Mesh @ 2GHz
- Message passing architecture
 - No coherent shared memory
 - Proof of Concept for scalable solution for many core
- Fine grain dynamic power management
 - On die controller for on-package VRs
 - Frequency modulation

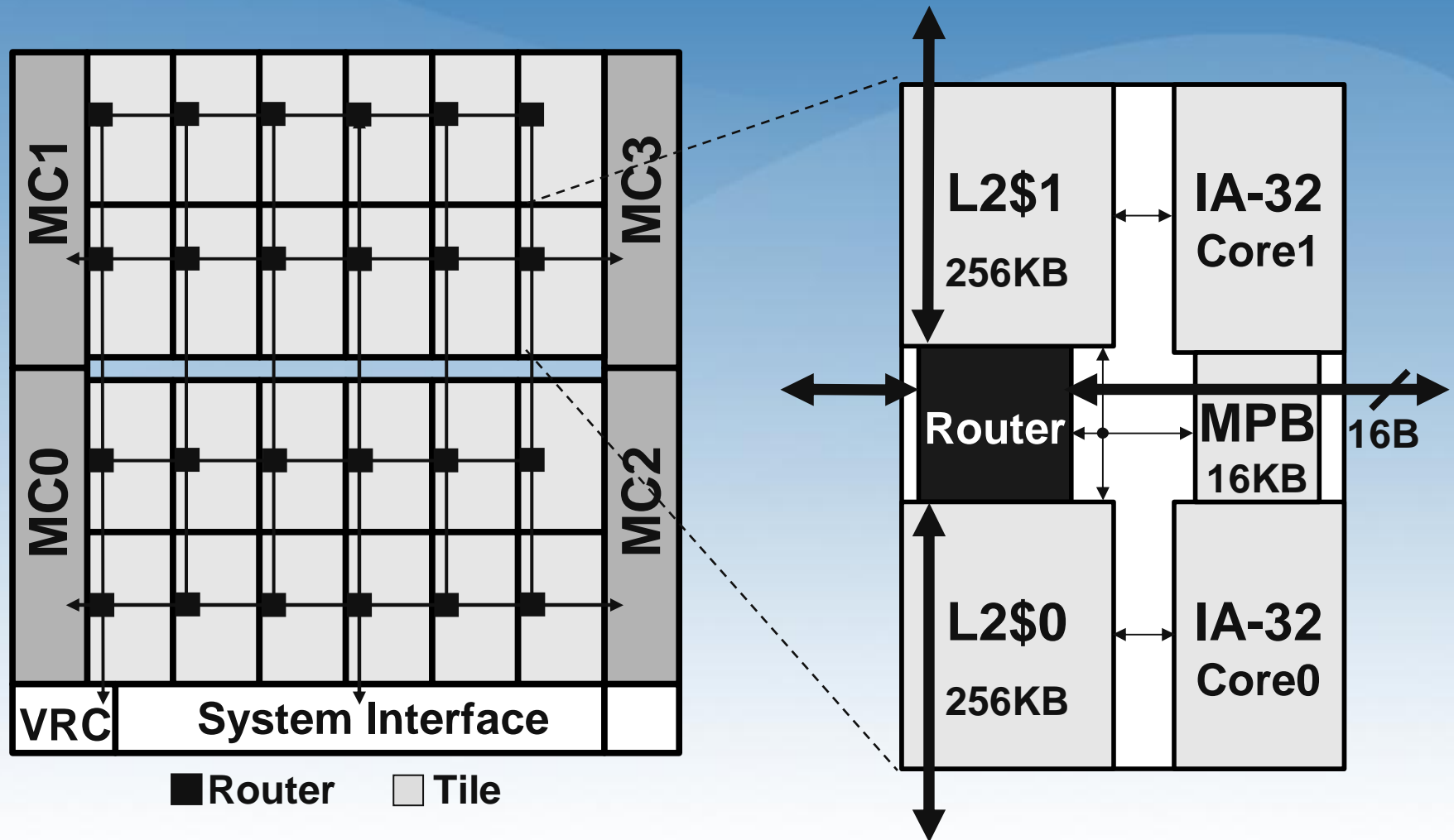


SCC Fullchip



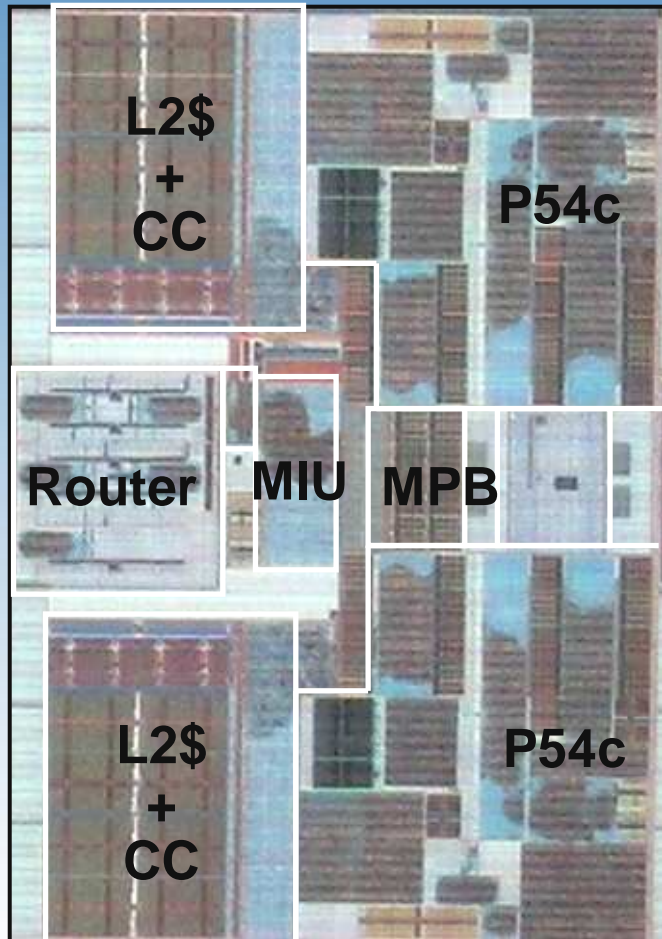
Technology	45nm Hi-K CMOS Process
Interconnect	1 Poly, 9 Metal (Cu)
Transistors	Die: 1.3B, Tile: 48M
Tile Area	18.7mm ²
Die Area	567.1mm ²

Die Architecture



2 core clusters in 6x4 2-D mesh

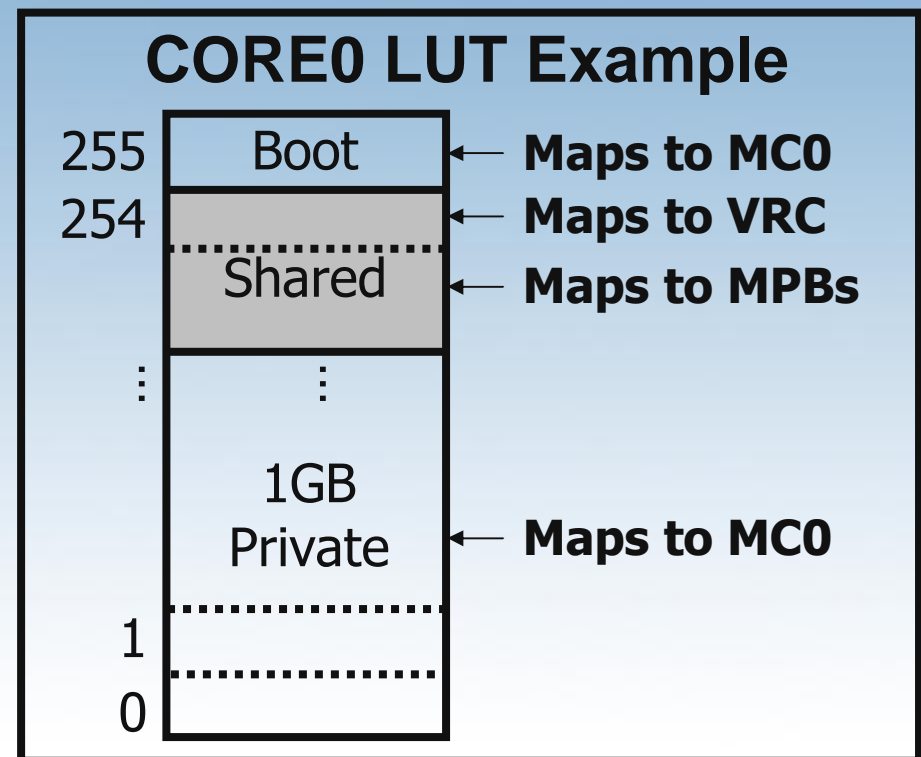
SCC Tile



- 2 P54C cores (16K L1\$/core)
- 256K L2\$ per core
- 16K Message passing buffer
- Mesh Interface Unit
 - Imbedded configuration registers
- Router
- Tile area 18.7mm²
- Core area 3.9mm²

Core Memory Management

- Each cores is allocated independent, private memory
- Core cache coherency is restricted to private memory space
 - Maintaining cache coherency for shared memory space is under software control
- Each core has an address Look Up Table (LUT) extension
 - Provides address translation and routing information
- LUT values must fit within the core and memory controller constraints
- LUT boundaries are dynamically programmable

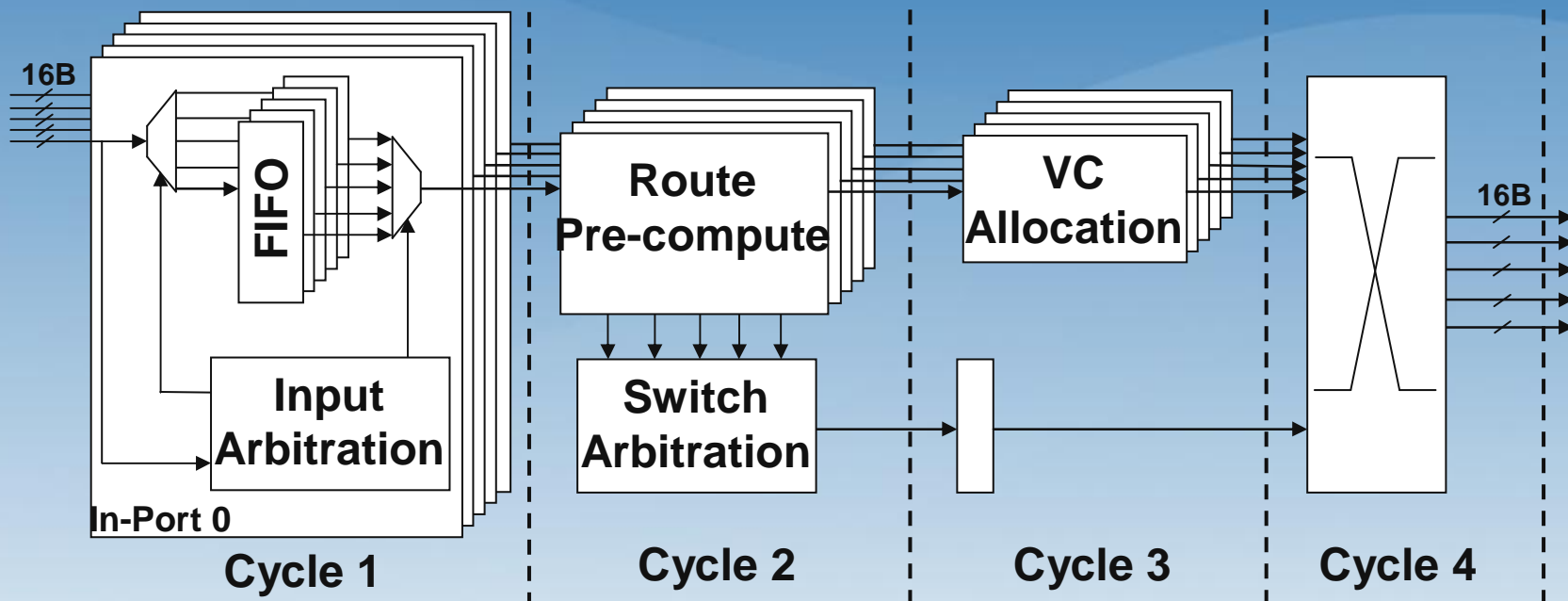


On-Die 2D Mesh

- 16B wide data links + 2B sideband
 - Target frequency: 2GHz
 - Bisection bandwidth: 2 Tb/s
 - Latency: 4 cycles (2ns)
- 2 message classes and 8 virtual channels
 - VC6 for request MCs
 - VC7 for response MCs
- Low power circuit techniques
 - Sleep, clock gating, voltage control, low power RF
 - Low power 5 port crossbar design
- Speculative VC allocation
- Route pre-computation
- Single cycle switch allocation

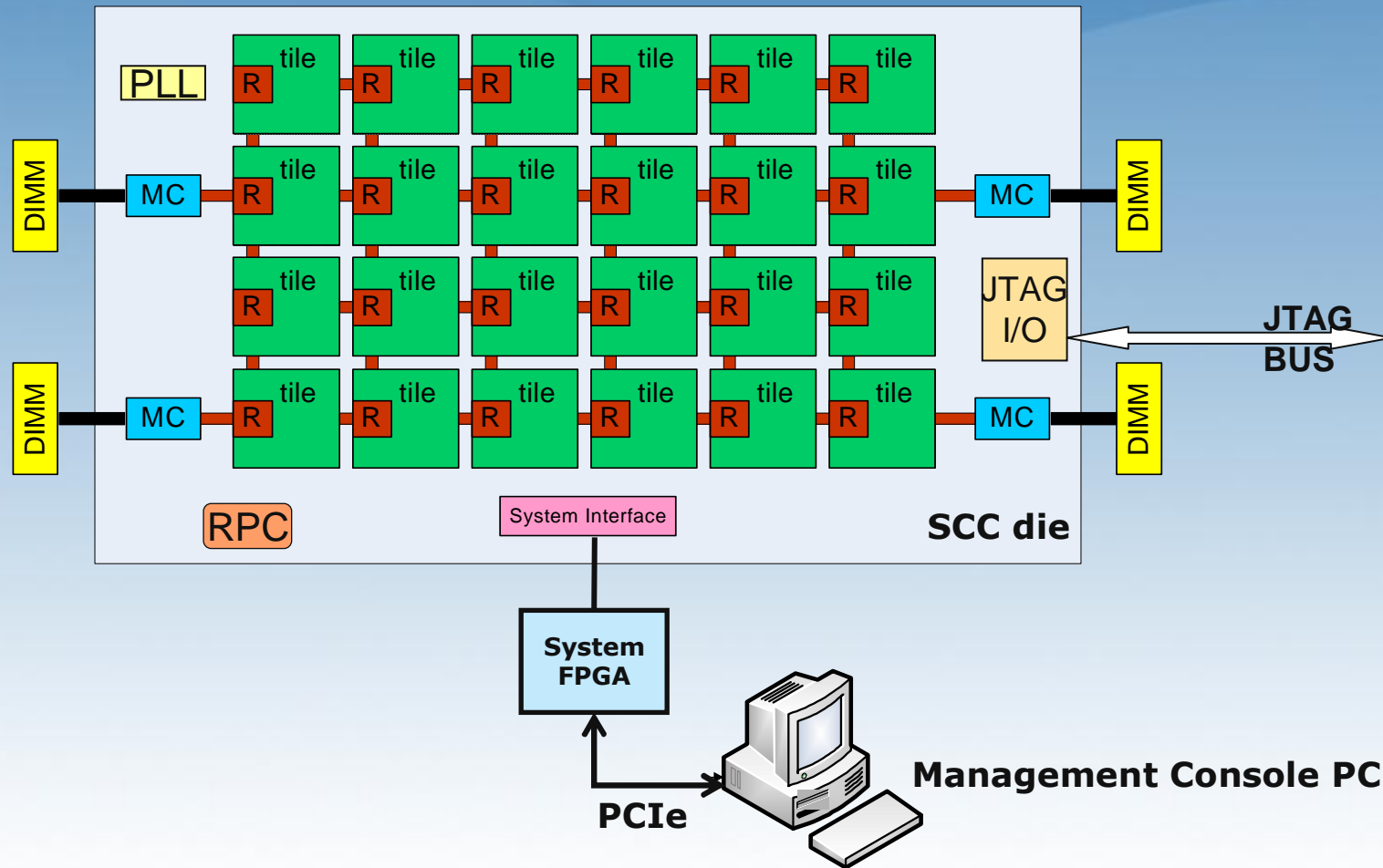


Router Architecture



Frequency	2GHz @ 1.1V
Latency	4 cycles
Link Width	16 Bytes
Bandwidth	64GB/s per link
Architecture	8 VCs over 2 MCs
Power Consumption	500mW @ 50°C

SCC system overview

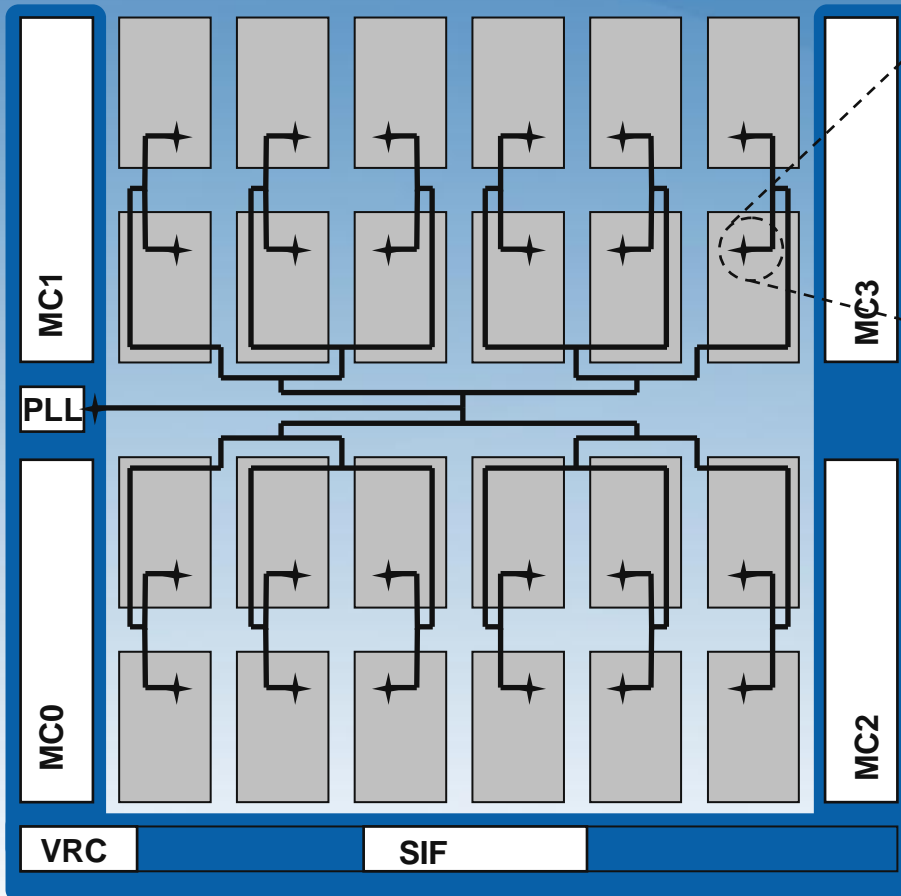


System Interface

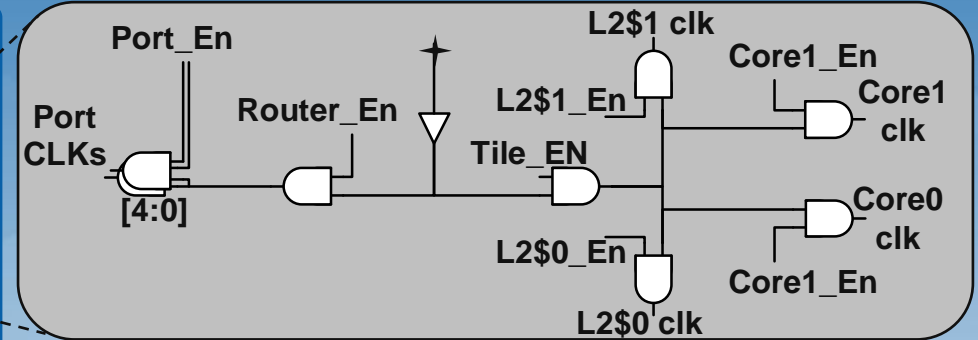
- JTAG access to config system while in reset/debug
 - Done on Power Reset from Host
 - Configuring memory controller etc.
 - Reset cores with default configuration
- Management Console PC can use Mem-mapped registers to modify default behavior
 - Configuration and voltage control registers
 - Message passing buffers
 - Memory mapping
- Preload image and reset rather than PC bootstrap
 - BIOS & firmware a work in progress



Clock Distribution



- Router
- Tile
- + Clock Gating



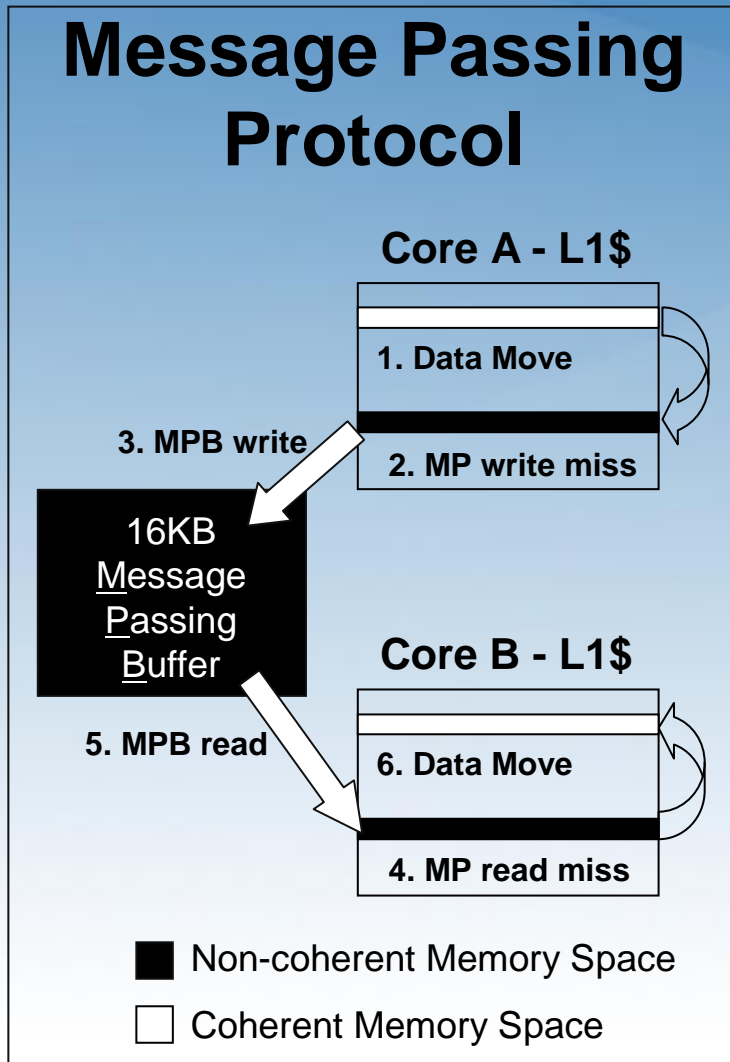
- **Balanced H-tree clock distribution**
- **Designed to provide 4GHz clock to tile entry points**
- **Simulated skew for adjacent tiles – 5ps**
- **Cross die skew irrelevant**

Message Passing on RC

- Message passing is done through shared memory space
- Two classes of shared memory:
 - Off-die, DRAM: Uncacheable shared memory ... results in high latency message passing
 - On-die, message passing buffers (MPB) ... low latency message passing
- MPB performance
 - Message Passing Buffers see a 15x improved latency as compared to off die DDR3-800



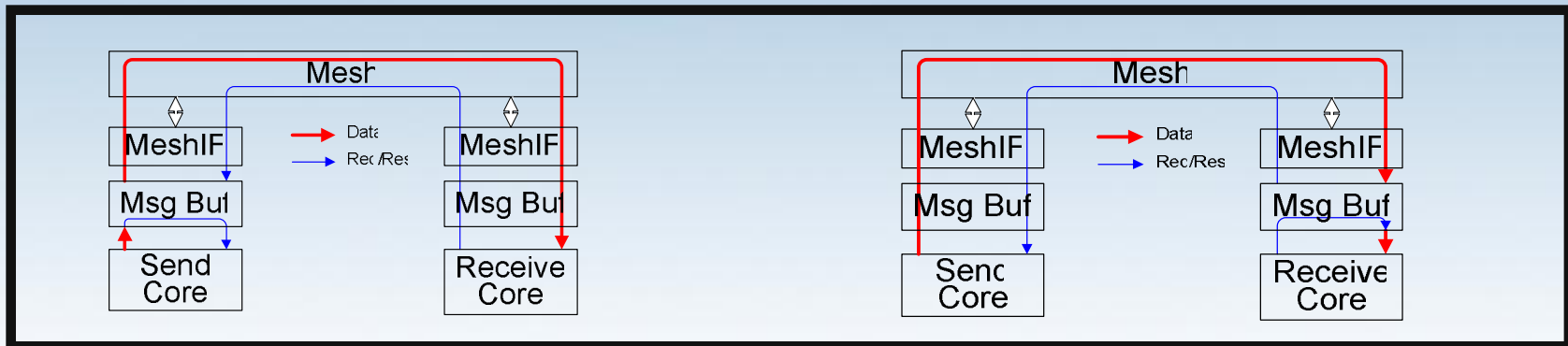
Message Passing Protocol



- Cores communicate through small fast messages
 - L1 to L1 data transfers
 - New Message Passing Data Type (MPDT)
- Message passing Buffer (MPB) – 16KB
 - 1 MPB per tile for 384KB of on-die shared memory
 - MPB size coincides with L1 caches

Dedicated Message Buffers

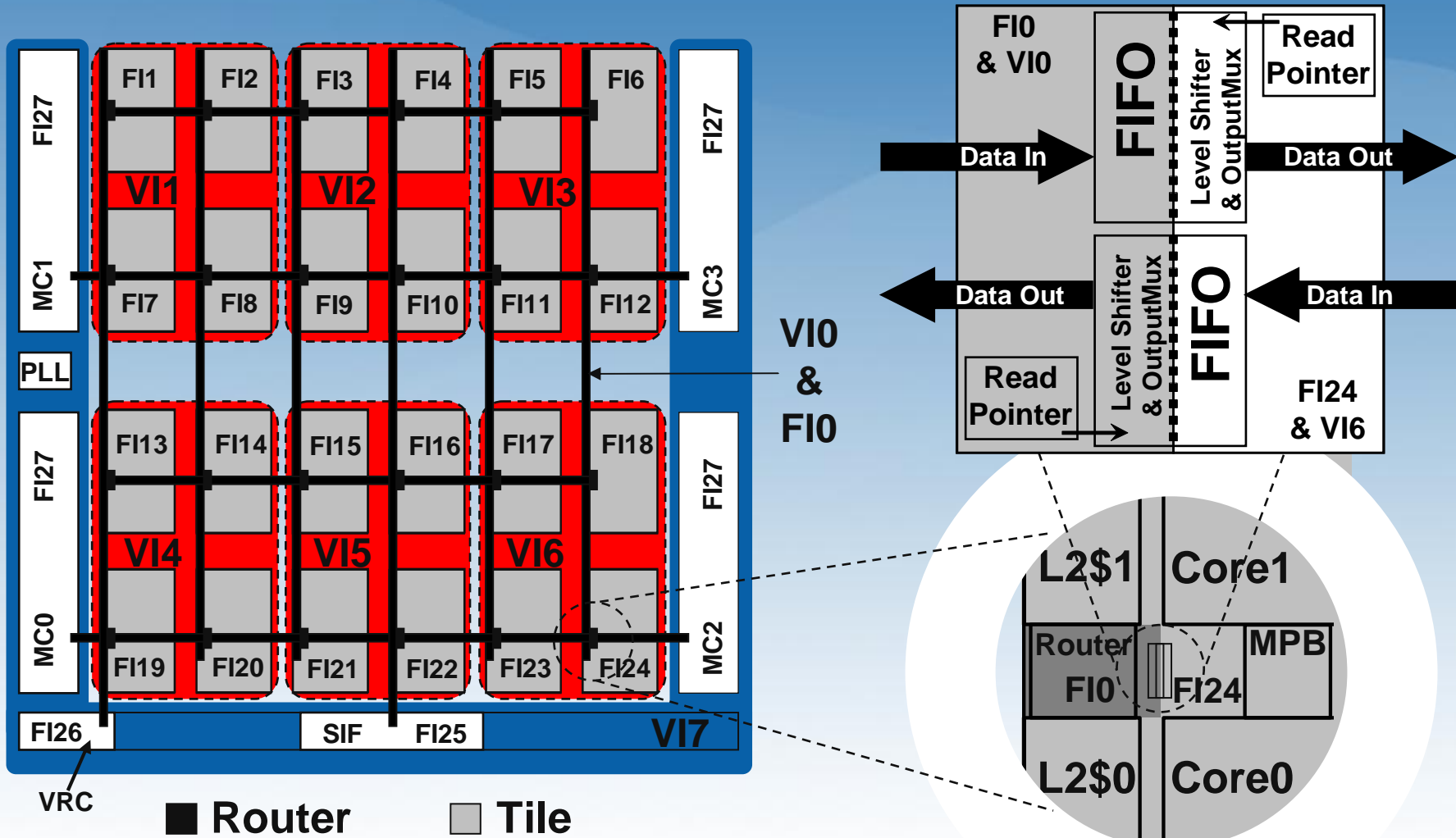
- Messages can be read from / written to one of three locations.
 - A message buffer locally in a core's tile.
 - A message buffer remotely in another tile
 - Off die to main memory
- We believe a remote write, local read has the best performance



Local write, remote read

Remote write, local read

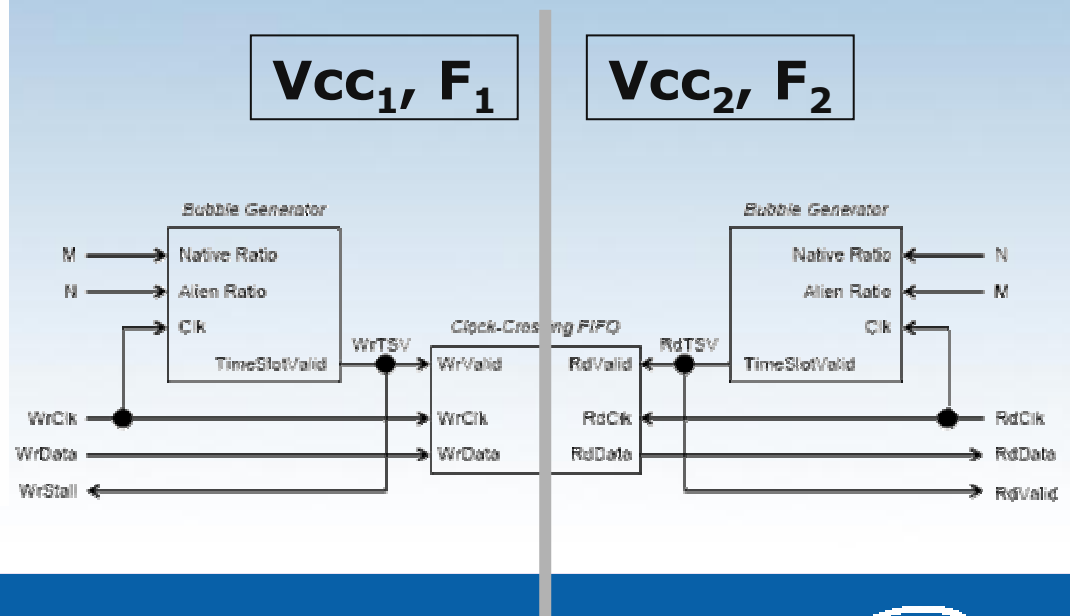
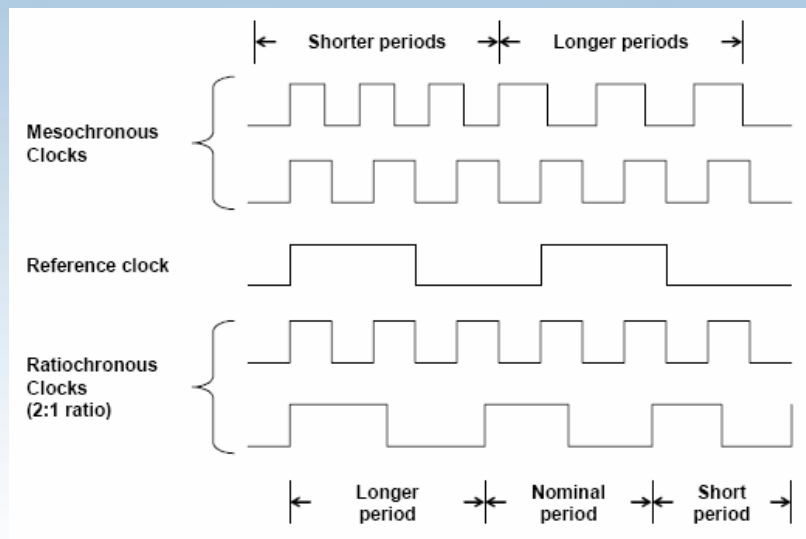
Voltage and Frequency islands



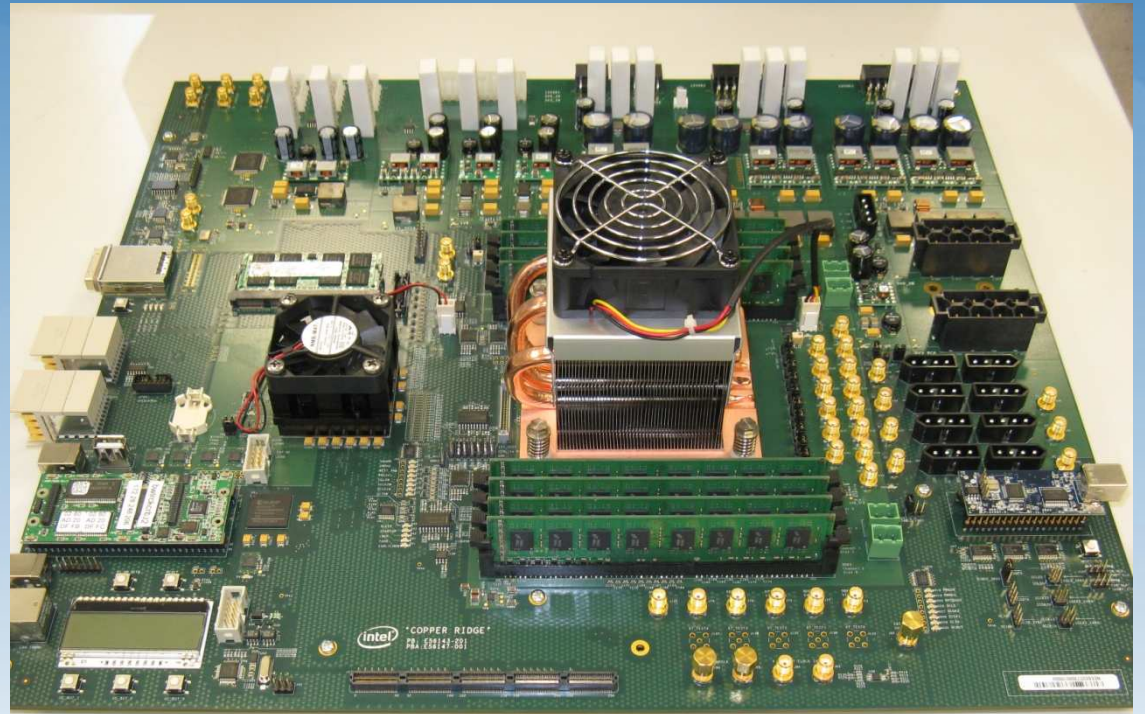
28 Frequency Islands (FI) 8 Voltage Islands (VI)

SCC Clock Crossing FIFO (CCF)

- 6 entry deep FIFO, 144-bits wide
- Built-in Voltage translation: 1:N ratios and pointer separation scanned in
- Key Benefit: independent mesh & tile frequency

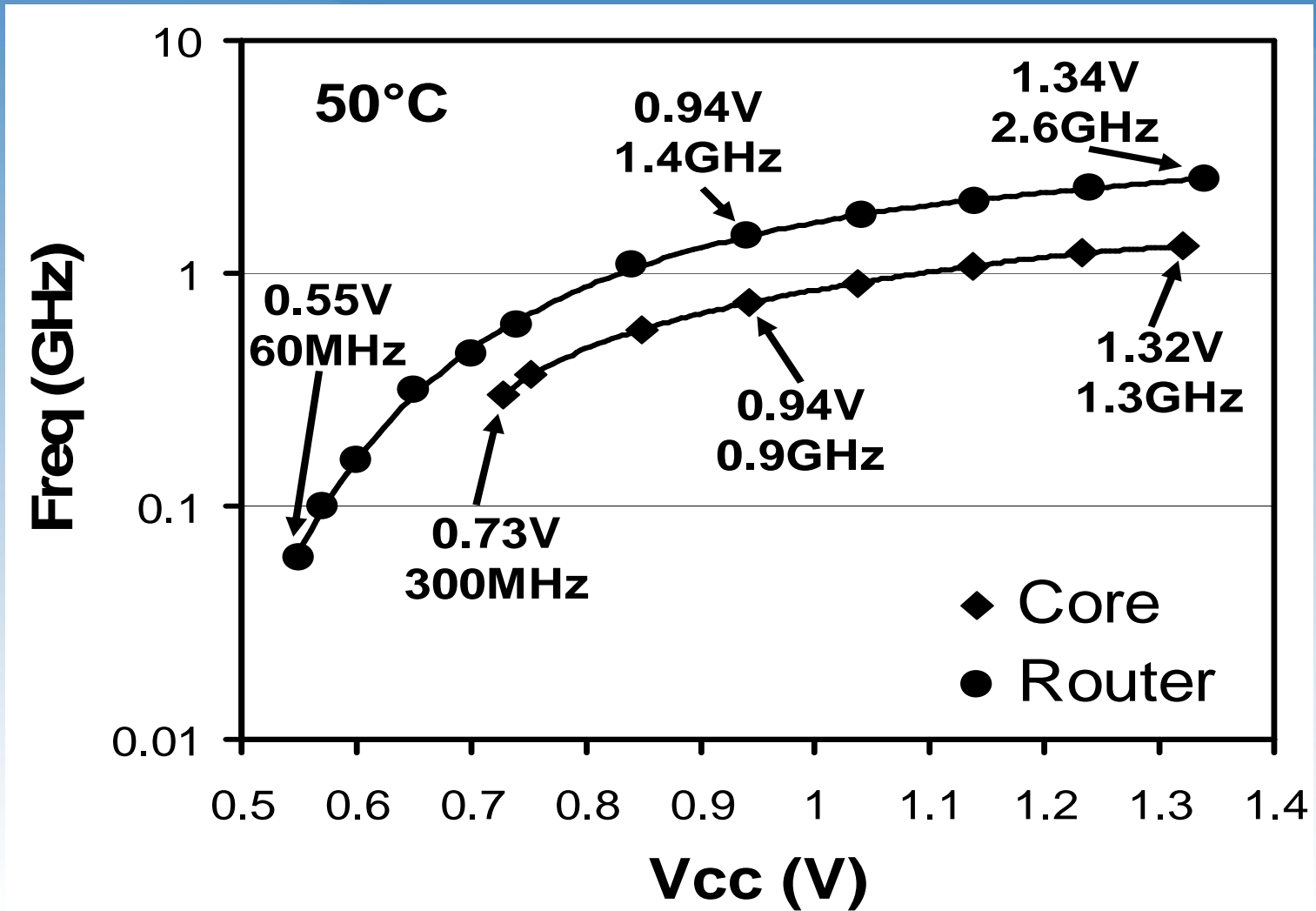


Package and Test Board

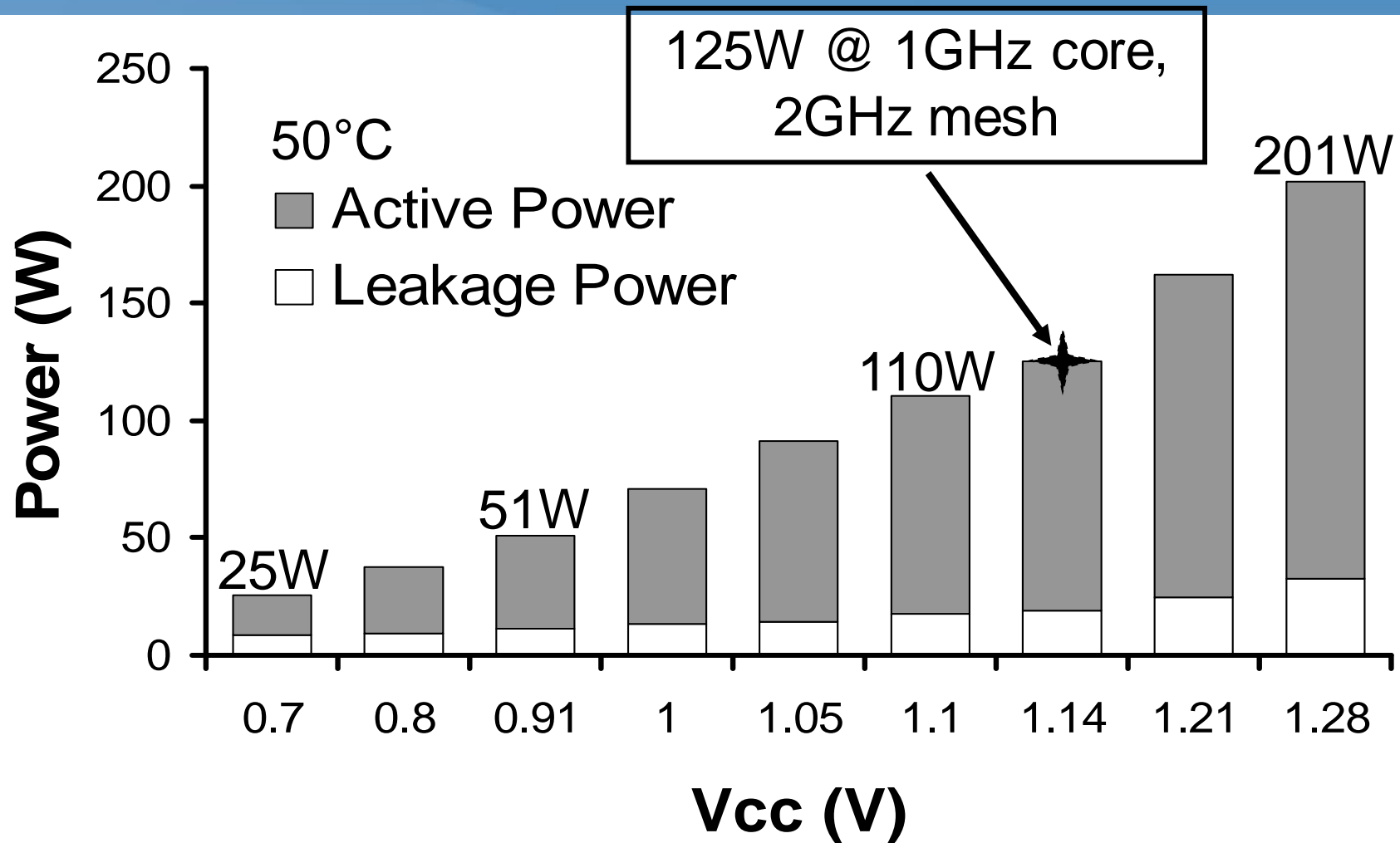


Technology	45nm Hi-K CMOS Process
Package	1567 pin LGA package
	14 layers (5-4-5)
Signals	970 pins

Core & Router Fmax

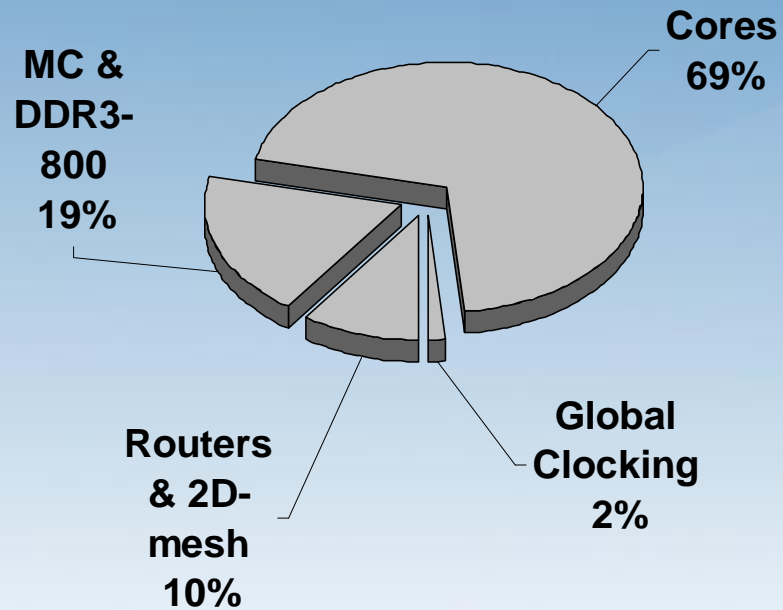


Measured full chip power



Power breakdown

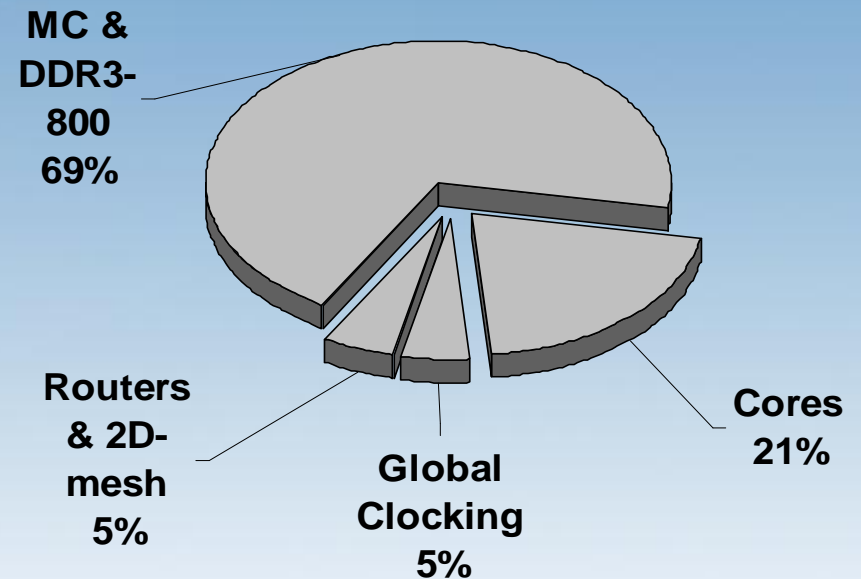
Full Power Breakdown Total -125.3W



Clocking: 1.9W Routers: 12.1W
Cores: 87.7W MCs: 23.6W

Cores-1GHz, Mesh-2GHz, 1.14V, 50°C

Low Power Breakdown Total - 24.7W



Clocking: 1.2W Routers: 1.2W
Cores: 5.1W MCs: 17.2W

Cores-125MHz, Mesh-250MHz, 0.7V, 50°C



Linux

A small Linux build was created targeted at the SCC specific features. The Linux memory driver was modified to enable Linux control of the on-die message passing buffers. Included in this Linux build is a TCP/IP driver for the 2-D mesh, connecting the host-PC and all 48 cores at the software layer.

Also, a TTY driver is included to allow the 48 cores to perform IO commands via memory mapped IO. This enables regular xterm type connections as shown.

```
telnet: Core 1 of Tile x=0, y=0 (localhost:5012)
-rwxr-x--- 1 root root 650280 Jan 10 23:00 stencil_synth_12
-rwxr-x--- 1 root root 650248 Jan 10 23:00 stencil_synth_4
-rwxr-x--- 1 root root 650568 Jan 10 23:00 stencil_synth_48
root@rck(x=0,y=0,Core=1):/ rce/stencil_synth_4
-sh: rce/: Permission denied
root@rck(x=0,y=0,Core=1):/ rce/stencil_synth_4
1.309680 1.312007 1.314311 1.316559 1.318740 1.320854 1.322914 1.324933 1.326930
1.328923 1.330935 1.332989 1.335105 1.337303 1.339585 1.341933
1.344285 1.346548 1.348793 1.350980 1.353108 1.355184 1.357216 1.359217 1.361201
1.363187 1.365201 1.367266 1.369420 1.371591 1.374107 1.376465
1.378731 1.380921 1.383052 1.385136 1.387180 1.389190 1.391173 1.393135 1.395091
1.397063 1.399076 1.401183 1.403438 1.405940 1.408221 1.410371
1.412572 1.414717 1.416762 1.418736 1.420673 1.422593 1.424503 1.426410 1.428329
1.430279 1.432313 1.434500 1.436999 1.439119 1.441118 1.443781
1.446491 1.448916 1.451031 1.452949 1.454784 1.456600 1.458421 1.460260 1.462128
1.464061 1.465822 1.468469 1.470490 1.472410 1.475770 1.479606
1.483253 1.486142 1.489353 1.492616 1.496127 1.499815 1.493328 1.497080 1.498888
1.500788 1.502888 1.504884 1.506799 1.510113 1.514518 1.519644
1.523842 1.526826 1.528947 1.530626 1.532157 1.533686 1.535270 1.536921 1.538651
1.540502 1.542414 1.544325 1.546811 1.550413 1.555716 1.562796
1.568056 1.567865 1.569259 1.570610 1.571893 1.573238 1.574682 1.576225 1.577865
1.579603 1.581428 1.583283 1.585582 1.589388 1.593880 1.598585 root@rck(x=0,y=0
,Core=1):/

telnet: Core 1 of Tile x=0, y=1 (localhost:5016)
root@rck(x=0,y=1,Core=1):/ rce/csimshift_4
Mark 03: Initial sum on UE 003 equals 322800,000000
Mark 14: Final sum on UE 003 equals 322000,000000
root@rck(x=0,y=1,Core=1):/ rce/stencil_synth_4
1.779153 1.780743 1.782249 1.783702 1.785120 1.786516 1.787905 1.789301 1.790721
1.792178 1.793687 1.795254 1.796880 1.798555 1.800260 1.801977
1.803756 1.805462 1.807075 1.808603 1.810060 1.811464 1.812839 1.814208 1.815595
1.817025 1.818521 1.820101 1.821777 1.823548 1.825397 1.827292
1.829192 1.831012 1.832712 1.834285 1.835747 1.837127 1.838455 1.839765 1.841092
1.842468 1.843927 1.845503 1.847221 1.849093 1.851105 1.853205
1.855296 1.857328 1.859389 1.861626 1.864056 1.866370 1.868432 1.869825 1.87052
1.868337 1.869725 1.871266 1.873006 1.874984 1.877201 1.879591
1.881979 1.884134 1.885982 1.887547 1.888894 1.890090 1.891197 1.892266 1.893348
1.894493 1.895758 1.897212 1.898932 1.901005 1.903486 1.906327
1.909232 1.911646 1.913527 1.914997 1.916188 1.917202 1.918116 1.918988 1.919871
1.920818 1.921890 1.923173 1.924784 1.926686 1.928773 1.933250
1.937213 1.939922 1.941706 1.942944 1.943871 1.944624 1.945283 1.945905 1.946534
1.947217 1.948011 1.949003 1.950336 1.952266 1.955242 1.959948
1.966593 1.969258 1.970563 1.971330 1.971856 1.972262 1.972610 1.972934 1.973263
1.973623 1.974050 1.974604 1.975402 1.976703 1.979176 1.984781
2.000000 2.000000 2.000000 2.000000 2.000000 2.000000 2.000000 2.000000 2.000000
2.000000 2.000000 2.000000 2.000000 2.000000 2.000000 2.000000 2.000000 root@rck(x=0,y=1
,Core=1):/

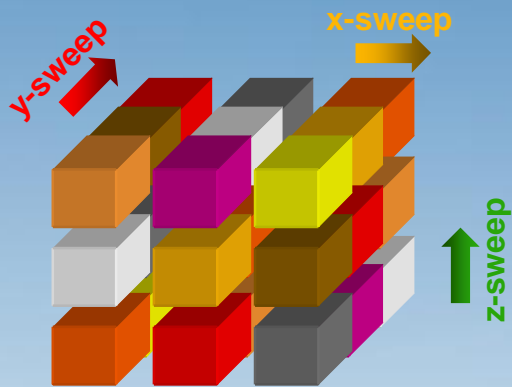
telnet: Core 0 of Tile x=0, y=0 (localhost:5010)
Mark 14: Final sum on UE 000 equals 322800,000000
root@rck(x=0,y=0,Core=0):/ rce/stencil_synth_4
Executing 1001 iterations
1.000000 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000
1.000000 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000
1.021070 1.028812 1.032219 1.034006 1.035100 1.035861 1.036449 1.036948 1.037406
1.037861 1.038350 1.038820 1.039695 1.040722 1.042529 1.044623
1.053390 1.061874 1.069865 1.068611 1.070437 1.071797 1.072890 1.073835 1.074713
1.075584 1.076511 1.077567 1.078864 1.080586 1.083056 1.086795
1.092284 1.097172 1.100994 1.103870 1.106704 1.107832 1.109308 1.110618 1.111850
1.113073 1.114358 1.115783 1.117451 1.119498 1.122104 1.125436
1.129432 1.133330 1.136728 1.139560 1.141910 1.143900 1.145641 1.147227 1.148736
1.150237 1.151796 1.153481 1.155373 1.157559 1.160124 1.163103
1.166389 1.169666 1.172700 1.176400 1.177776 1.179884 1.181794 1.183573 1.185287
1.186993 1.188746 1.190604 1.192621 1.194852 1.197333 1.200059
1.202945 1.205833 1.208592 1.211150 1.213494 1.215649 1.217656 1.219562 1.221415
1.223261 1.225144 1.227106 1.229186 1.231418 1.233817 1.236367
1.239010 1.241645 1.244199 1.246630 1.248921 1.251081 1.253135 1.255114 1.257054
1.258987 1.260947 1.262965 1.265070 1.267284 1.269616 1.272054
1.274560 1.277020 1.279426 1.281751 1.283984 1.286126 1.288194 1.290208 1.292192
1.294170 1.296168 1.298211 1.300322 1.302517 1.304806 1.307187
Total time: 0.543697
root@rck(x=0,y=0,Core=0):/

telnet: Core 0 of Tile x=0, y=1 (localhost:5014)
Cleaning out MPI... Start w/o argument for interactive mode...
Done! Bye...
root@rck(x=0,y=1,Core=0):/ rce/csimshift_4
Mark 03: Initial sum on UE 002 equals 322000,000000
Mark 14: Final sum on UE 002 equals 321000,000000
root@rck(x=0,y=1,Core=0):/ rce/stencil_synth_4
1.602526 1.603934 1.605060 1.606051 1.607090 1.608263 1.609580 1.611022 1.612568
1.614221 1.615902 1.617579 1.619677 1.622477 1.626366 1.629587
1.631156 1.632233 1.633079 1.633893 1.634613 1.635891 1.637127 1.638496 1.639971
1.641540 1.643121 1.644665 1.646461 1.648598 1.651055 1.653119
1.654447 1.655383 1.656097 1.656861 1.657748 1.658796 1.659997 1.661325 1.662754
1.664261 1.665780 1.667254 1.668808 1.670470 1.672172 1.673638
1.674739 1.675594 1.676363 1.677176 1.678107 1.679178 1.680384 1.681703 1.683110
1.684580 1.686067 1.687518 1.688964 1.690403 1.691792 1.693026
1.694468 1.694921 1.695660 1.696705 1.697606 1.698393 1.700173 1.701509 1.702912
1.704368 1.705845 1.707303 1.708733 1.710121 1.711441 1.712653
1.713755 1.714795 1.715825 1.716892 1.718024 1.719235 1.720521 1.721875 1.723286
1.724741 1.726222 1.727702 1.729163 1.730586 1.731947 1.733227
1.734435 1.735628 1.736821 1.738032 1.739277 1.740565 1.741898 1.743277 1.744698
1.746159 1.747651 1.749159 1.750669 1.752160 1.753605 1.754975
1.756270 1.757628 1.758981 1.760324 1.761667 1.763020 1.764391 1.765787 1.767214
1.768677 1.770180 1.771719 1.773286 1.774861 1.776413 1.777887 root@rck(x=0,y=1
,Core=0):/
```



Linpack and NAS Parallel benchmarks

1. **Linpack (HPL): solve dense system of linear equations**
 - Synchronous comm. with "MPI wrappers" to simplify porting

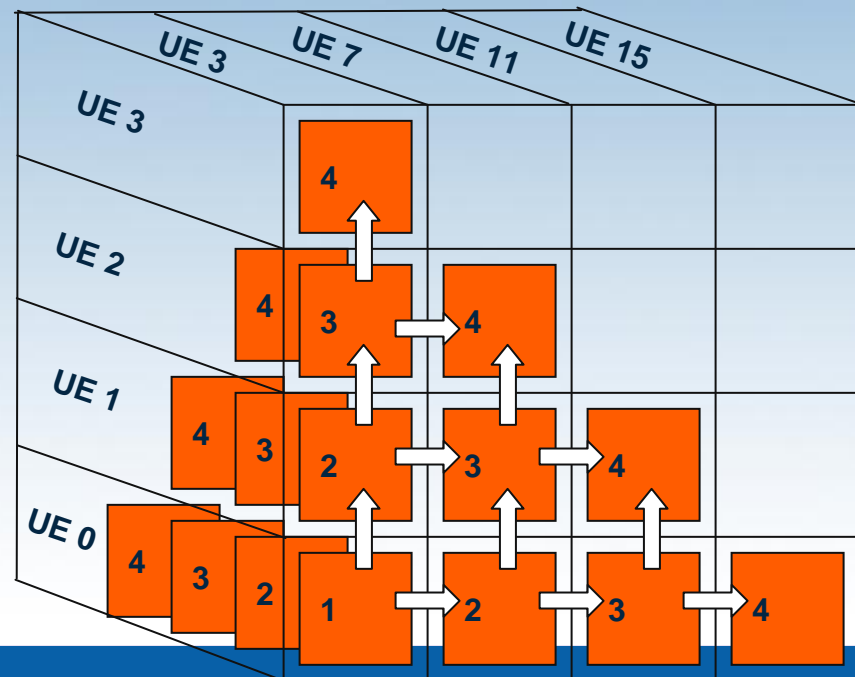


2. **BT: Multipartition decomposition**

- Each core owns multiple blocks (3 in this case)
- update all blocks in plane of 3x3 blocks
- send data to neighbor blocks in next plane
- update next plane of 3x3 blocks

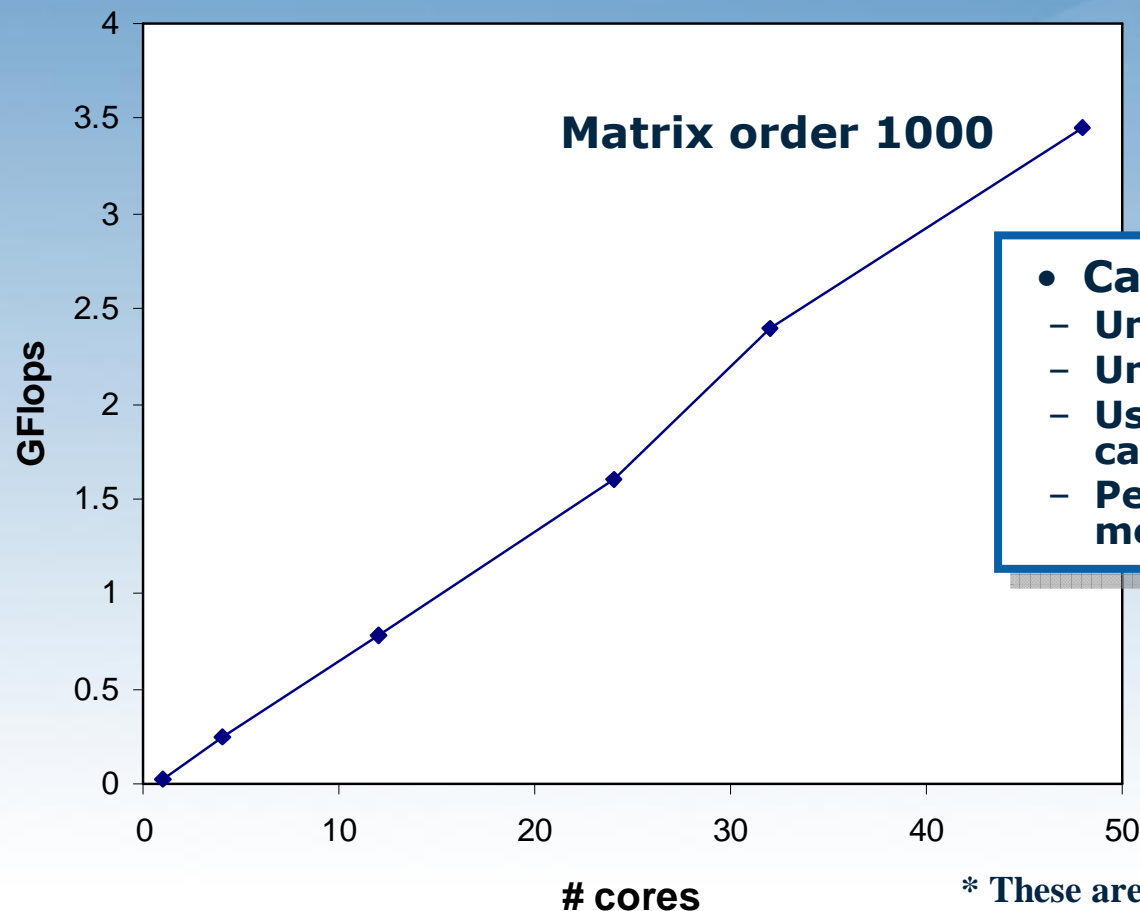
3. **LU: Pencil decomposition**
Define 2D-pipeline process

- await data (bottom+left)
- compute new tile
- send data (top+right)



Linpack, on the Linux SCC platform

- **Linpack (HPL)* strong scaling results:**
 - **GFLOPS vs. # of cores for a fixed size problem (1000).**
 - **This is a tough test ... scaling is easier for large problems.**



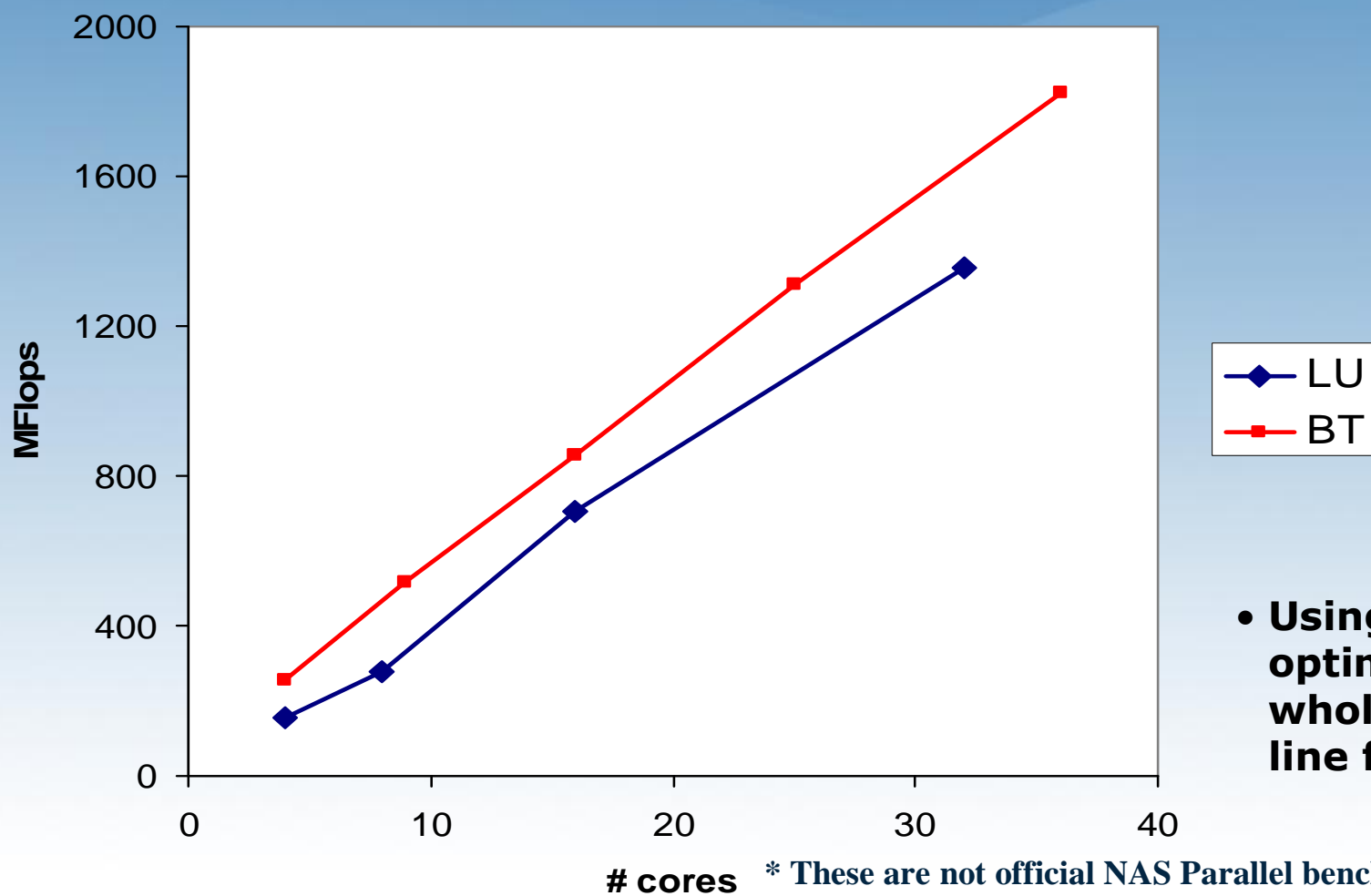
- **Calculation Details:**
 - **Un-optimized C-BLAS**
 - **Un-optimized block size (4x4)**
 - **Used latency-optimized whole cache line flags**
 - **Performance dropped ~10% with memory optimized 1-bit flags**

* These are not official LINPACK benchmark results.



LU/BT NAS Parallel Benchmarks, SCC

Problem size: Class A, 64 x 64 x 64 grid*



- Using latency optimized, whole cache line flags

* These are not official NAS Parallel benchmark results.



Power Management Demo

The screenshot displays a web browser window showing the SCC Power Management interface. The interface includes the Intel logo, the text "SCC Power Management", and "Advanced Workload Aware Power Management Technology". It features a "Power Management" toggle switch currently set to "OFF", with "ON" also visible. Below the text is a detailed image of a processor die. The status "ready" is shown in the bottom right corner of the interface.

To the right, a terminal window shows a series of shell commands and their outputs, including the execution of `/shared/DEMOS/ECO_Q/new_pwr_app.exe`.

Below the terminal is the "Rock Creek performance meter" window. It displays "Individual CPU usage..." with a grid of 24 small gauges, each showing a green needle. Below this grid are radio buttons for "Cockpit style" (selected), "Taskmanager style", and "Combined (overlay)". At the bottom, there is a section for "Over-all CPU usage of enabled cores..." featuring a large gauge with a green needle and a small graph titled "Over-all CPU usage over time".



Summary

- A 48 IA-32 core processor in 45nm CMOS
 - Second generation 2D-mesh network
 - 4 DDR3 channels in a 6×4
 - Highest level of IA-32 integration
- New message passing HW for increased performance
 - 384KB of on-die shared memory
 - Message passing memory type
- Power management employs 8VIs and 28FIs for DVFS
- Chip dissipates between 25W and 125W as performance scales
 - 25W at 0.7, 125MHz core, 250MHz mesh and 50°C
 - 125W at 1.14V, 1GHz core, 2GHz mesh and 50°C

