

L2 to Off-Chip Memory Interconnects for CMPs

Daniel Killebrew
University of California, Berkeley
Computer Science 258
dank@eecs.berkeley.edu

Allen Lee
University of California, Berkeley
Computer Science 258
leeallen@eecs.berkeley.edu

ABSTRACT

In recent years, chip multiprocessors (CMPs) have seen an increasing asymmetry between the number of cores and the number of memory access points on a single die, prompting a new study of network topologies that efficiently connect many nodes to few nodes. In this paper, we evaluate the latency and power efficiency of the tapered-fat tree (TFT) topology on a cycle-accurate simulator intended to model the TILE64 multiprocessor. We replace the original mesh network with two TFT networks (one for memory requests, one for memory responses) and run four synthetic benchmarks modeled after those in the PARSEC suite. Because several connections in the TFT network require global wires, we also modeled the multi-cycle latencies using a wire-delay model. Our simulator keeps track of activity factors in each of the routers, which we combine with the Orion power models to determine the switching power of the routers. We determined that for applications with large amounts of sharing and little off-chip traffic, the TFT topology offers negligible advantage over the mesh. However, the benchmarks that exhibited large amounts of off-chip traffic completed a workload up to 33.6% faster on the TFT than the same workload on the mesh. The TFT topology also has better power-efficiency in all cases except `canneal`, which dissipates a large amount of power in the inter-router interconnect due to the high volume of stores. Although this benchmark dissipated more average power on the TFT, the total energy consumed is less than on the mesh because it completed faster. At the cost of extra design complexity, our results suggest that the TFT topology offers better latency and energy-efficiency overall than the mesh for memory-intensive applications running on multicore processors.

1. INTRODUCTION

Chip multiprocessors (CMPs), single silicon dies with many simple cores, have gained increasing attention as the computing industry has the combined power and ILP wall[1]. Manycores present new engineering challenges and opportunities due to their unique new architecture. The proximity

of cores makes interprocessor communication much cheaper, but providing enough memory bandwidth for all these cores to function becomes a more serious problem.

As CMPs have evolved from single core architectures, we can observe the evolution of the network used to connect cores to both off-chip and on-chip memory. All designs have a private L1 cache, but the L2 cache, which is shared when there are a few cores, becomes more private as cores are added. In chips with less cores, interprocessor communication occurs through a shared memory model, while those with more cores add explicit interprocessor networks for message passing.

Besides offering a possible solution to the problems that have stymied the uniprocessors of the past, CMPs will make interesting new workloads feasible. Applications such data mining and advanced video encoding will place an increased strain on the memory system. Some of these workloads demonstrate an increasing per-core performance as more cores are added to a chip[2]. This translates into an exponential growth in memory bandwidth demand. Clearly this exponential growth has significant implications for future CMPs.

In this paper we examine several modern multicore chips, analyzing design trends in the off-chip memory interconnect (OCMI). We will summarize the results of the PARSEC benchmark suite as they apply to OCMI. We look at a modern 64-core chip, the Tile64, and use the results of the PARSEC suite to motivate a redesign of the Tile64's OCMI. We provide a comparison of our proposal and the Tile64's current solution using detailed power analysis and simulation results.

1.1 Manycore Bandwidth Challenge

The SPLASH2 benchmarks [3] have been the standard parallel benchmarking suite for many years now, but they are now somewhat dated. Some of the algorithms it uses are now inferior to newer methods for solving the same problems. Additionally, the benchmarks it contains are aimed at HPC and graphical simulations. CMPs will be used by a much more diverse set of users with a correspondingly wider array of applications. Thus, we use the PARSEC benchmarking suite [2] as the basis for our analysis.

These new applications will have a large appetite for memory bandwidth. There are 12 different applications in the

PARSEC suite, and of these 12, all but three have a working set size larger than 8 MB. This is much larger than any current on-chip cache, distributed or otherwise. Even if on-chip caches were to grow, these working set sizes grow along with the data set the algorithm is working on. Applications such as data mining become more useful as the data set is increased, so we expect data set sizes to grow with time. As a result, we expect the working set of at least eight of these applications to exceed future on-chip cache sizes, creating a distinct need for off-chip memory bandwidth.

The amount of off-chip memory bandwidth for half of the applications scales linearly with the number of cores. This means that as cores are added, the bandwidth demand increases because of the core added, as well as because the cores already on the chip increase their individual demands. This results in an exponentially increasing demand for off-chip bandwidth.

The PARSEC applications do not exhibit the same demand for interprocessor communication that they do for off-chip bandwidth. Four of the 12 exhibit an insignificant amount of sharing. Of those with true sharing, the data sets are so large that shared accesses between different processors are temporally spaced far apart. The effect is that by the time the second processor makes a shared access, the data has already been moved off-chip because of a capacity conflict. Two of the applications are parallelized in such a way that data is passed between different computational stages of the algorithm, allowing efficient data sharing. If cores were added, this on-chip traffic would scale proportionally, but so would the interprocessor network as well, presumably. None of the applications demonstrate an increasing usage of shared (or private) cache space as the number of processors scales.

The result of this will be a divergence of the requirements of the interprocessor communication network and the OCMI. Interprocessor communication networks beyond simple buses are usually chosen such that network bandwidth scales as processors are added. This will be sufficient to handle the traffic due to shared memory access and message passing, as this traffic volume grows linearly along with the number of cores. The off-chip memory bandwidth, on the other hand, will increase exponentially as cores are added. Working set sizes will only get larger as the utility of many of these applications increases as the data set grows (and the working set grows with the data set). The additional cores will increase the efficiency of the cores already present, resulting in the exponential demand for off-chip memory bandwidth. The off-chip bandwidth requirements will outstrip those of the on-chip communication network.

2. CASE STUDIES

On-chip caches have evolved as cores have been added to the chip, and this has had a direct effect on the off-chip memory interconnect. The Intel Core2 Duo has two cores on a die, which access a shared L2 cache and off-chip memory controller. AMD's Barcelona possesses four cores, each with a private L2 cache. A crossbar connects the cores to a shared L3 cache, which interfaces with the two memory controllers. Intel's Nehalem is similar to the Barcelona, except that there are eight cores. Sun's Niagara uses a full crossbar to connect

eight cores to a four way banked, shared L2[4]. Sun's Rock uses a hierarchical scheme, with the 16 cores divided into four clusters[5]. A cluster shares one instruction cache and two data caches, along with a crossbar that provides access to the central crossbar. This central crossbar connects to a four way banked, shared L2. See Table 1.

These processors do not have an explicit interprocessor communication network, instead relying on coherence of the shared caches to share data and synchronize. In all cases, there is a shared L2, or shared L3 if the cores have private L2's. This L3 accesses the memory upon a miss. Thus, the interprocessor communication network is also the processor to off-chip memory interconnect.

IBM's Cell processor is unique in both its cache structure (or lack thereof) and its interconnection network. The Power Processing Element (PPE) is a standard Power core with a cache hierarchy, while the 8 Synergistic Processing Elements (SPEs) possess no caches, only a 256 KB local store each. These local stores are accessed using DMA transfer initiated by the PPE or any of the SPEs. Only the owning SPE has direct access to its local store. The PPE, SPEs, and memory controller are all connected to four rings. The rings are divided into two pairs of unidirectional links, each pair traveling in the opposite direction from the other. The rings serve as both the interprocessor network and the off-chip memory network. A key difference is that there is no shared cache structure between the processors and the memory controller, as a result of the streaming model the Cell has adopted.

2.1 Tile64

Tile64 is a CMP with 64 cores[6]. As it has the most cores of any modern commercial CMP and has enough cores to make the scaling properties of the networks start to become evident, we have chosen it. The 64 cores are connected in an 8x8 mesh network. Mesh networks are popular because they are simple, have an efficient two-dimensional VLSI layout, and support tiling, which makes design and verification easier as well. Each tile contains a processor, the routers that support the networks, and an L2 cache.

The Tile64 has five physically separate 32 bit mesh networks that connect the processors. The mesh networks use dimension ordered routing (DOR) and credit based flit-level flow control. The routers require one cycle to route flits which are turning, and zero cycles for flits continuing in the same direction. Interrouter links take one cycle to traverse. Each router input has three flits of buffering, the minimum required to support full throughput on the interrouter links.

Each mesh serves a different purpose, the two that we are concerned with are the Tile Dynamic Network (TDN) and Memory Dynamic Network (MDN). The TDN delivers interprocessor cache-to-cache transfer requests. The actual bulk data transfer returns over the MDN to prevent protocol deadlock. An alternative would be to provide another virtual channel on the TDN to avoid deadlock, but this would require buffer space beyond the minimum used for full throughput flow control. The designers of the Tile64 specifically stated that they were attempting to minimize the chip area consumed by the networks, so they chose to multiplex

Table 1: Multicore Survey

Name	Cores	L2	Interconnect	Off-chip Bandwidth
Intel Core2 Duo	2 (on a single die)	Shared	Crossbar	10.7 GB/s
AMD Barcelona	4	Private	Crossbar	10.7 GB/s
Intel Nehalem	1-8	Private	Crossbar	31.9 GB/s
Sun Niagara	8	Shared	Crossbar	25.6 GB.s
IBM Cell	8	Local Store	4 Rings	25.6 GB/s
Sun Rock	16	Shared	Hierarchical Crossbar	54.4 read, 32.64 write GB/s
Tilera Tile64	64	Private	Mesh	25.6 GB/s

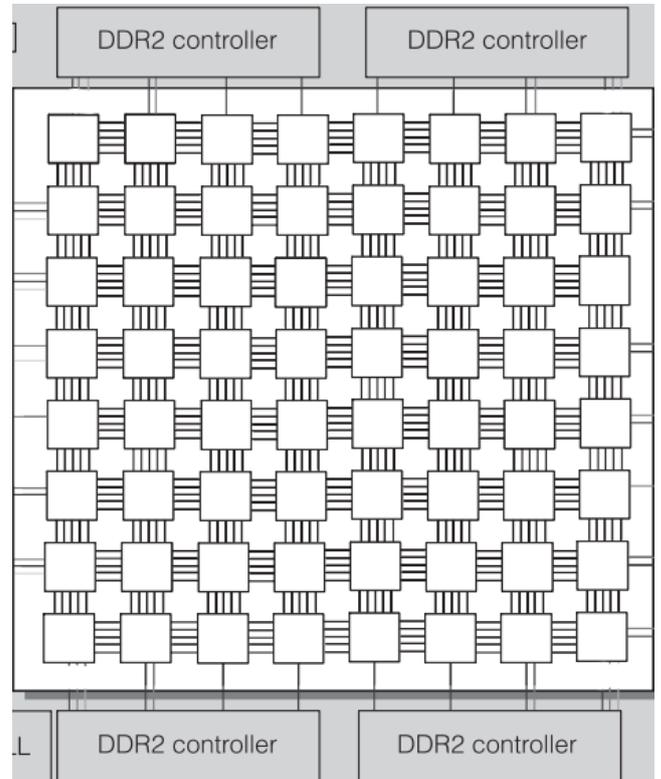
the MDN. As demonstrated later by the simulation, this results in an under-utilized TDN and an over-utilized MDN.

All loads and stores that are destined for off-chip memory use the MDN. Protocol deadlock is prevented because each DRAM reserves enough buffer space to buffer a certain number of requests from each core. The cores, in turn, do not have more than this predetermined number of outstanding requests to any single DRAM. If the responses traveled a different physical or logical network, this precaution would not be necessary. Removing the limit on outstanding requests would improve throughput and save on buffer space at the DRAM.

The connection between the memory controllers and the tiles is a little odd because of the mismatch between the number of memory controllers (4) and tiles (64). There are two memory controllers on the north side of the chip, and two on the south side. Each is connected to the three innermost adjacent tiles. Tiles route to the closest of the three access points (see 1). The memory controller uses round robin arbitration choose from among the three input tiles. Under a uniform traffic distribution, the three access points each see a different amount of traffic. The innermost of the three sees five times as much traffic as the middle one, and 2.5 times as much as the outer access point. This results in global unfairness because the five columns of tiles that use the innermost access point will see less throughput than the other columns of tiles.

The alternative, though, is much worse. The mesh employs DOR, with routing in the X direction followed by the Y direction. If tiles were to spread their traffic over the three access points instead of going to the closest, when the response returned from DRAM and entered the mesh, it would immediately travel in the X direction until it reached the appropriate column. Consider an example where two packets arrive, one at the middle and inner access points. Both are destined for a tile somewhere to the east, so they must take turns exiting the inner router's eastern output port. Add in the outer access point, and this results in horrible throughput for packets entering the mesh. Changing the routing to Y dimension followed by X would not solve anything, only cause the problem to occur as packets reached the edge and started traveling in the X dimension towards an access point. The problem is a result of the routing algorithm and the topology mismatch (four memory controllers talking to 64 cores on a mesh).

The four memory controllers support a total bandwidth of 256 bits/cycle, and the MDN has a bisection bandwidth of

**Figure 1: Tile64 memory controller connections**

1024 bits/cycle. As the MDN is shared by requests returning from memory as well as intertile bulk data transfers, this seeming 'overprovisioning' of the mesh is required to support full throughput from off-chip memory.

3. PROPOSAL

We propose to replace the Tile64's mesh MDN with two tapered fat trees (TFTs). TFTs used as processor to processor interconnection networks have low chip area usage and power dissipation[7]. They are also superior to meshes in good workload completion times, giving good efficiency in terms of area-delay and energy-delay. Replacing the mesh MDN with TFTs means that traffic that previously used the MDN for interprocessor communication (loads from remote L2 caches) now must transit the TDN.

3.1 Tapered Fat-Tree (TFT)

The topology shown in Figure 2 chosen also has benefits because the top four routers map to the four memory con-

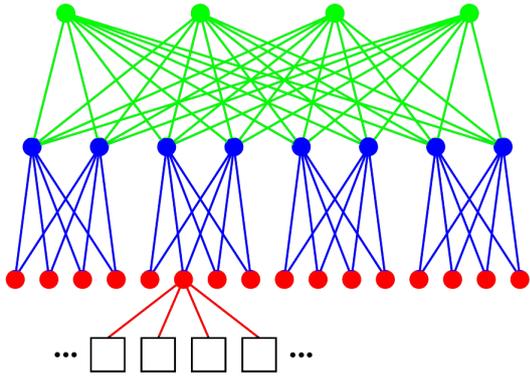


Figure 2: Tapered Fat-Tree topology

trollers. Because all traffic on the MDN is either traveling from a processor/L2 cache to a memory controller, or visa versa, there are never any U-turns in the network. There are no packets that travel up to a middle level router, then turn around and go back down. This reduces routing complexity, eliminates crossbar connections, and as a result, lowers power consumption and area consumed. Hop count is also reduced, lowering packet latency and reducing router energy because each packet visits fewer routers. Splitting the network up into separate physical networks, one in each direction does not increase router area (it lowers it because of the simplified crossbars)[7].

All of these benefits come at a cost, however. Wires no longer travel to the nearest neighbor, meaning wires must now cross over each other. This means using upper metal layers. Depending upon the specifics of the layout, it could require another metal layer to be added. Interrouter links are now too long to traverse in a single cycle, and so repeaters must be inserted, taking up additional area, and possibly increasing place and route complexity.

3.1.1 Topology

Figure 2 illustrates the tapered-fat tree (TFT) topology. Each of the bottom level routers (circles) connects to four processor/L2 nodes. Each router of the top level connects directly to a single memory controller.

3.1.2 Routing Function

The routing function is very simple. Packet travel to the bottom level of routers, randomly choose one of the two possible intermediate level routers, then move on to the router that connects to the desired memory controller. Routing on the way back down is just as simple, with one of the two possible middle level routers being chosen randomly. This routing simplicity also has latency advantages because the routing function takes less time.

3.1.3 Tapered Fat-Tree on TILE64

Unlike the original mesh network which has a very regular placement of routers and wires, place-and-route for a non-mesh network is fairly difficult. As a starting point, we use the router placements suggested in [7] for the TFT. Figure 3 illustrates the placement of the TFT routers and their logical connectivity mapped onto the TILE64. Due to the

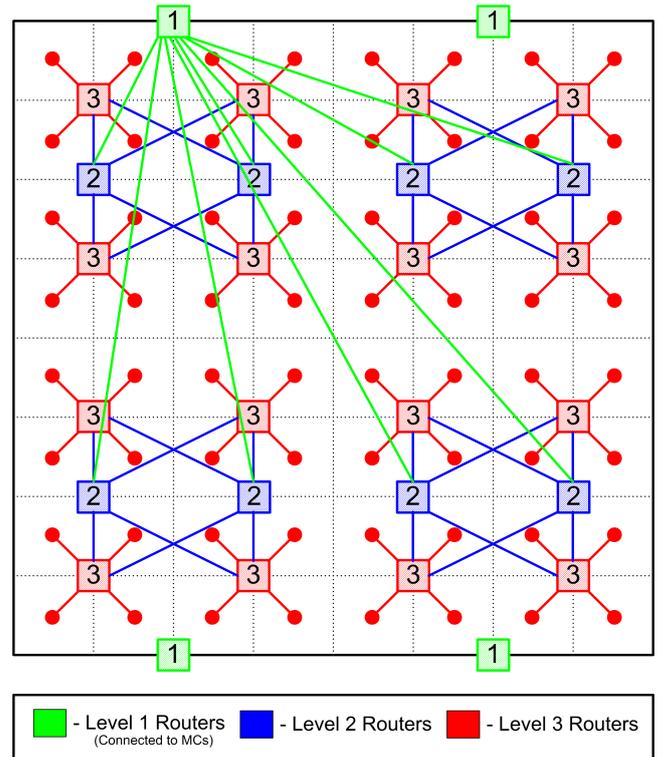


Figure 3: Tapered Fat-Tree mapped onto TILE64. The connections only represent logical connectivity. A possible physical layout is shown in Figure 5.

asymmetry, density, and lengths of the wire connections, the inter-router connections add significant power and delay. We discuss these issues in Section 4.

4. MODELS

4.1 Wire Model

The propagation delay of a signal across a wire is proportional to $R_{wire}C_{wire}$. Because both R_{wire} and C_{wire} grow linearly with the length of the wire, the delay across a wire increases quadratically with wire length. By inserting a chain of repeaters (inverters) along the wire, designers can make the delay grow only linear with wire length. To achieve the minimal delay, the repeaters must be sized and placed along the wire such that the fanouts the repeaters drive form a geometric series.

Achieving minimal wire delay is a very difficult problem to solve, not only mathematically, but also in terms of physical design. For a given wire length, designers must first determine the number of repeaters that achieves the best tradeoff in terms of power, area, and delay. Then, they must size and place the repeaters such that the fanouts form a geometric series. Since most semiconductor companies use standard cells, the optimal buffer sizes are likely not part of the standard cell library, so the closest ones must be chosen. Even after choosing the standard cell inverters and determining the optimal placements, it may not be possible to place them at the ideal spots. Globally repeated wires require many via cuts from the upper layer down to the substrate, and the

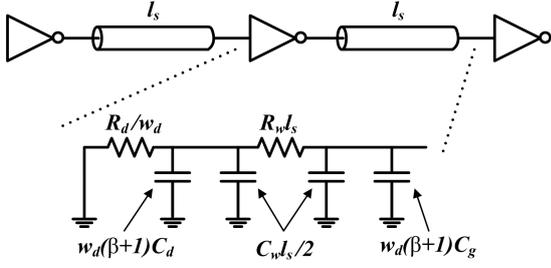


Figure 4: First-order repeater model

repeaters may be too large to fit into the substrate if the chip area is heavily utilized.

For these reasons, we use the first-order repeated wire delay model described in [7][8] and illustrated in Figure 4. Namely,

1. Repeaters are identically sized and evenly spaced.
2. The effects of inductance are not modeled.

We assume that each repeater has a PMOS/NMOS ratio of 2 and has twice the drive strength of a minimum-sized inverter.

The delay across a single repeater stage (T_s) is calculated as:

$$K_0 = R_d(1 + \beta)(C_d + C_g) \quad (1)$$

$$K_1 = \frac{R_d}{w_d}C_w + w_d R_w(1 + \beta)C_g \quad (2)$$

$$K_2 = \frac{1}{2}R_w C_w \quad (3)$$

$$T_s = 0.69(K_0 + K_1 l_s + K_2 l_s^2) \quad (4)$$

Here, R_d is the equivalent resistance of the repeater, β is the PMOS/NMOS ratio, C_d is the diffusion capacitance, C_g is the gate capacitance, w_d is the width of the repeater, C_w is the capacitance per millimeter, R_w is the resistance per millimeter, and l_s is the length of the segment. The total wire delay across n stages (T_w) is:

$$T_w = nT_s \quad (5)$$

4.1.1 Technology Model

Because our router power models are based on the 100nm process, we obtained and derived physical parameters for the wires corresponding to this technology node from the 2003 ITRS Roadmap [9]. We discuss the router power models in more detail in Section 4.2. Tables 2 and 3 present the relevant parameters.

4.1.2 Wire Layout

Figure 5 illustrates the physical layout of the wires for the routers at each level. We approximate the length of each wire as the Manhattan distance between the endpoints. We assume that Level 3 routers use M1 wires, Level 2 routers use intermediate wires, Level 1 routers use global wires, and repeaters are placed every 0.5mm.

Parameter	Symbol	Value	Units
Supply Voltage	V_{DD}	1.0	[V]
Clock Frequency	f	750	[MHz]
Tilera Tile Width	W_t	2.6025	[mm]
Gate Capacitance	C_g	1.95E-16	[F/ μ m]
Diffusion Capacitance	C_d	1.77E-15	[F/ μ m]
Resistivity of Copper	ρ_{Cu}	1.72E-8	[Ω -m]
Repeater Eq. Resistance	R_d	1138.4	[Ω - μ m]
NMOS Leakage Current	$I_{off,N}$	10	[nA/ μ m]
PMOS Leakage Current	$I_{off,P}$	10	[nA/ μ m]
PMOS/NMOS Ratio	β	2	

Parameter	M1	Int.	Global	Units
Pitch	240	320	475	[nm]
Width	120	160	320	[nm]
Thickness	75	94	113	[nm]
A/R	1.6	1.7	2.1	
R_w	1.91E3	1.14E3	4.75E2	[Ω /mm]
C_w	2.00E-13	1.85E-13	1.31E-13	[F/mm]

Because Tilera has not revealed the die size of the Tile64, we must estimate its dimensions. One source [10] postulated that the die size is approximately four times that of the quad-core Intel Core 2 Extreme QX6850, but this approximation results in such a large area that the yield would be virtually zero. Instead of using this estimate, we compare the transistor count of the Tile64 to the transistor count of NVIDIA's G80 chip, which was fabricated with the same process. The G80 GPU has 681 million transistors with a die size of 480mm² (21.9mm \times 21.9mm) [11]. Because the Tilera has 615 million transistors [12], we approximate the Tile64 die area as

$$\frac{615\text{M transistors}}{681\text{M transistors}} * 480\text{mm}^2 = 433.48\text{mm}^2$$

This area corresponds to dimensions of 20.82mm \times 20.82mm. At eight cores to a side, the width of one core (W_t) is thus 2.6025mm.

4.1.3 Multi-Cycle Wire Delay

As Figure 5 shows, the tapered fat-tree topology requires connectivity between points that are located at nearly opposite sides of the chip. Using the wire delay model described earlier in this section, we determined that these long paths require multiple cycles. Table 4 summarizes the results of

Table 4: Multi-Cycle Wire Delays ($t_{cycle} = 1.33$ ns)

Router Level	Paths	Length	n	T_w (ns)	Cycles
Level 1	P_1-P_4	W_t	5	1.22	1
Level 2	P_1-P_2	W_t	5	1.06	1
	P_3-P_4	$3 \cdot W_t$	15	3.18	3
Level 3	P_1-P_2	$3 \cdot W_t$	15	2.17	2
	P_3	$5 \cdot W_t$	26	3.61	3
	P_4-P_6	$7 \cdot W_t$	36	5.06	4
	P_7	$9 \cdot W_t$	46	6.50	5
	P_8	$11 \cdot W_t$	57	7.95	6

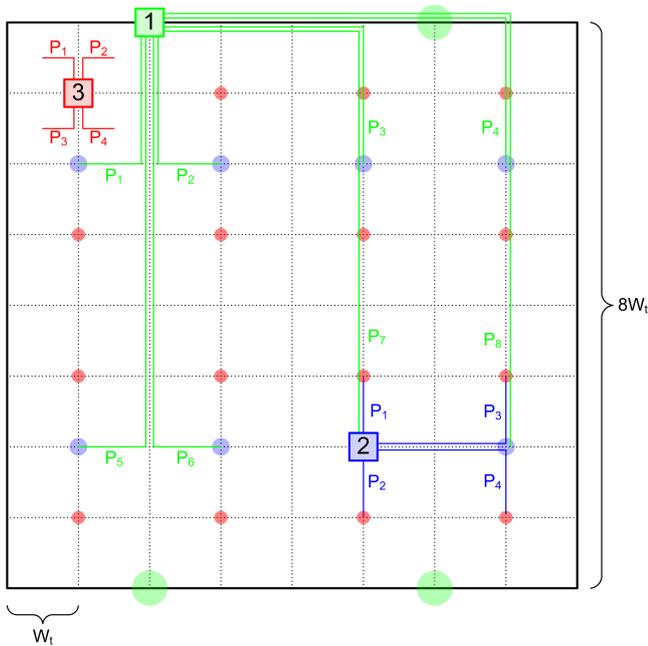


Figure 5: Physical layout of TFT on TILE64

our analysis.

4.2 Power Model

Because power is becoming the limiting constraint in modern multiprocessors, we model power dissipation in our network simulations to compare the energy efficiency of the different topologies. We consider the static and dynamic power dissipated both within the routers and in the inter-router interconnect. Our intra-router energy measurements are based on the Orion energy models[13] using the 100nm technology node.

We assume a clock frequency of 750MHz and supply voltage of 1.0V, the same parameters used by the actual TILE64[12].

Although the actual TILE64 multiprocessor uses the 90nm process, we were not able to obtain accurate power models for the routers at this technology node. However, because we are only interested in comparing the relative energy efficiency of the topologies, applying the same model to the different schemes is sufficient for comparison purposes.

4.2.1 Router Power

As a flit traverses a router, switching power ($P_{router,dyn}$) is consumed at each of the functional units: the input buffer (P_{buf}), the matrix crossbar (P_{xbar}), and the arbiter (P_{arb}):

$$P_{router,dyn} = P_{buf} + P_{xbar} + P_{arb} \quad (6)$$

Each functional unit is made up of a set of atomic components, i.e., elements consisting of one or more capacitors that have identical switching behavior. For an atomic component x with switching capacitance C_x and activity factor α , the average switching power $P_{x,sw}$ is calculated as:

$$P_{x,sw} = \frac{1}{2} \alpha C_x V_{DD}^2 f \quad (7)$$

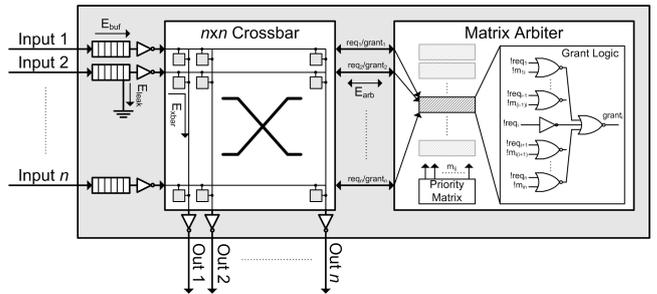


Figure 6: Input-buffered on-chip router with matrix arbiter.

The activity factor α is calculated individually for each atomic component of each functional unit from the load on the router. Here, load means the probability that a flit arrives at a single input buffer. To determine the appropriate load for each router, our simulations keep track of the total flit activity at each of the routers.

As we scale to deep submicron processes, subthreshold leakage power becomes a significant component of overall power dissipation. To account for leakage power in our simulations, we adopt the model detailed in [14]. The leakage power (P_{leak}) is calculated as:

$$P_{leak} = I_{leak} V_{DD} \quad (8)$$

Here, I_{leak} is the leakage current of the router buffers. Figure 6 illustrates the router we modeled with the energy components labeled. More details about the breakdown of the individual energy components can be found in [13].

4.2.2 Interconnect Power

Although the Orion models provide us with accurate power models within the routers, they do not account for power dissipated in the inter-router interconnect. Because the TFT requires many global connections, interconnect power consumption may be significant. We calculate the static ($P_{w,stat}$) and dynamic ($P_{w,dyn}$) power dissipated driving a wire segment (l_s) based on the equations disclosed in [7]:

$$\begin{aligned} P_{w,dyn} &= \alpha (w_d (C_g + C_d) + C_w l_s) f V_{DD}^2 \\ P_{w,stat} &= \frac{1}{2} w_d (I_{off,N} + \beta I_{off,P}) V_{DD} \\ P_w &= P_{w,dyn} + P_{w,stat} \end{aligned}$$

Here, α is the activity factor on the wire obtained from our simulations.

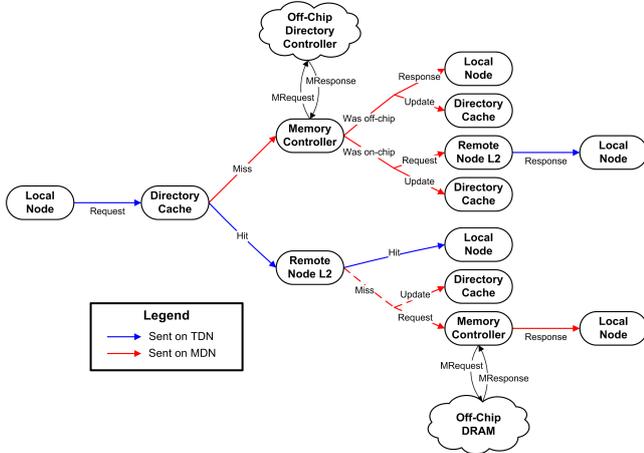
5. SIMULATION ENVIRONMENT

We wrote a detailed cycle-accurate flit-level simulator, consisting of 6000 lines of code, to provide completion times and activity factors for every router and link. These activity factors were then used to provide power estimates. The Tile64 was modeled after the description in 2.1. The TFT networks used the wire delays listed in Table 4.

The simulator consists of a number of nodes that are the sources and sinks for packets, representing the L2 cache injecting cache requests onto the TDN and MDN. These packets enter the appropriate network, and get buffered at the

Table 5: PARSEC Applications Simulation Parameters

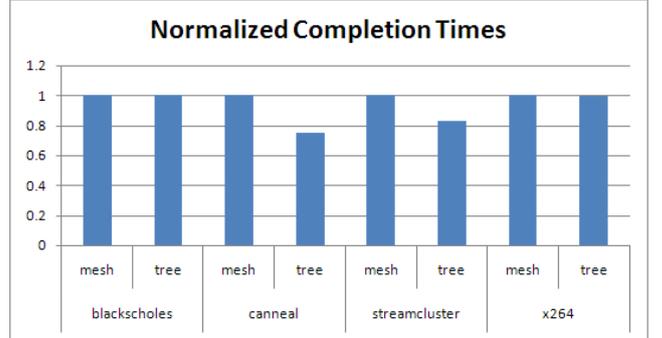
Benchmark	Working Set Size	Amount of Sharing	Per Proc. Injection Rate	% Shared Access	% Loads
streamcluster	small	small	.02705	.1	99
blackscholes	small	large	.0001373	18.5	33.3
cannal	large	small	.02341	.1	66.1
x264	large	large	.00408	50.9	68.7

**Figure 7: Protocol for TFT memory traffic****Table 6: Simulation Packet Sizes**

Packet Type	Size (bits)
Directory Load Query	32
Load Request	32
Load Response	288
Store Request	288
Store Response	32

input buffer to a router. The router examines the destination field of the first flit to decide which output port to request arbitration for. The arbiters use round-robin arbitration per output port to decide which packet gets to go; a hold line is invoked until the tail flit of a packet passes. The flit is then placed into a pipeline that simulates the intra-router and inter-router delays. The flit makes its way through the network until it arrives at the directory cache (inside a sink) or a memory. The memory changes the request into the appropriate response, placing it in a pipeline to delay the response for 52 cycles (our memory latency). The response then makes its way back to the source that originally issued the request.

We modeled a cache coherent directory based protocol (Figure 7). All stores are assumed to be writebacks destined for off-chip memory. Loads first traveled to an on-chip directory cache. If the request was a hit in the directory cache, the request was forwarded to the appropriate remote node, or off-chip memory as appropriate. A miss in the directory cache resulted in the packet being forwarded to the off-chip directory. If the data was off-chip, it was handled immediately. Otherwise, it was sent back on-chip to the appropriate remote node. Packet sizes were as according to Table 6 (flits on the MDN and TDN are 32 bits).

**Figure 8: Completion times of benchmarks normalized to the completion time with the Mesh topology.**

We performed synthetic benchmarks based on a selected set of the PARSEC benchmarks. We chose four that would give us a full range of small and large working sets, and low and high amounts of sharing. This would give us the full range of stressing the OCMI, as well and the interprocessor interconnection network. See 5 for the full simulation parameters. Per processor injection rate is the probability that a given processor would inject a request into the network on a cycle. Percent shared accesses refers to what percent of requests must go off-chip to be satisfied. We assumed that any shared requests could be satisfied by a node on-chip. To simulate a real processor, each processor is limited to a maximum of eight outstanding requests at any given time. After this point, a processor is considered to be stalling, and not producing any useful work.

6. RESULTS

6.1 Benchmark Completion Time

To compare the relative performance of the different topologies, we chose a workload for each benchmark and executed the same workload on each of the topologies. We normalize the completion time of each benchmark to the time it takes to finish executing when the MDN is comprised of the original mesh topology. The results are shown in Figure 8. It should be noted that when we refer to mesh or tree topology, or mesh or tree network, we are referring to the full simulation setup with a mesh TDN (always), and either a mesh MDN or a tapered fat tree TDN. All energy and power numbers include the activity of the TDN.

Our simulations indicate that for programs with large amounts of sharing, such as `blackscholes` and `x264`, the TFT offers little advantage over the original Mesh network. These applications require only a small amount off-chip bandwidth, so most of the memory traffic stays on the TDN, which is common to both networks we tested. The fact that the TFT

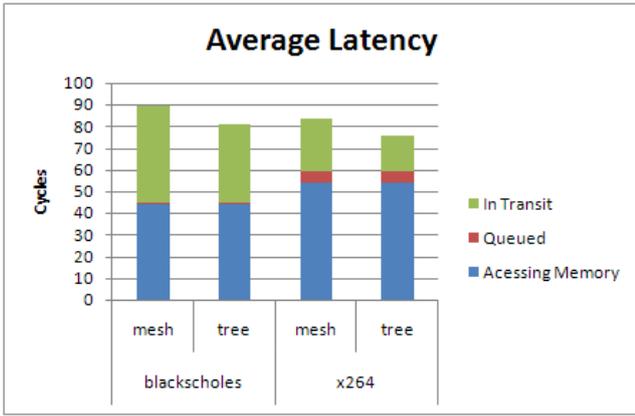


Figure 9: Average latency of packets for the blackscholes and x264 benchmarks.

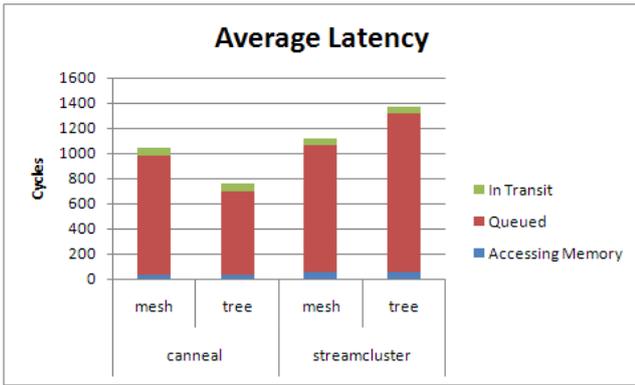


Figure 10: Average latency of packets for the canneal and streamcluster benchmarks.

network has a slightly faster completion time is due to packets spending less time in transit, as shown in Figure 9. As expected, the higher traffic injection rate of **x264** results in a small amount of queuing, of which there is almost none in **blackscholes**.

For applications with small amounts of sharing, such as **streamcluster** and **canneal**, most of the memory traffic goes off-chip. As demonstrated in Figure 10, the packets are spending the majority of their time in queues, waiting to make forward progress. At this point, the primary reason that the tree network completes faster is not due to the lower latency of the tree network, it is due to the increased throughput, shown in Table 7. In these benchmarks where congestion on the MDN becomes an issue, the TFT topology speeds up the completion time by 33.6% for **canneal** and 20.6% for **streamcluster**.

6.2 Power

The TFT topology shows a power savings for all of the benchmarks, as shown in Figure 11. Interestingly, a majority of the power dissipation appears to be from static leakage power because the routers on the MDN are fairly under-utilized. Even for memory intensive applications, the average load for routers in the Mesh network remain below

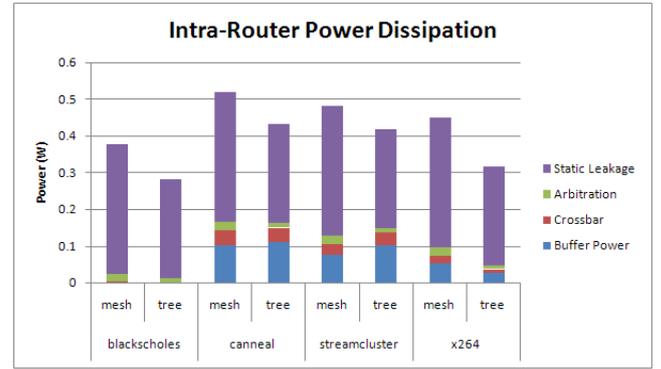


Figure 11: Power dissipated within the routers of the network.

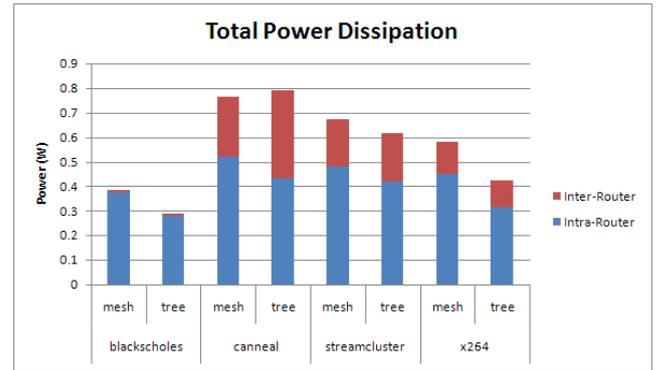


Figure 12: Total power dissipation in routers and interconnect.

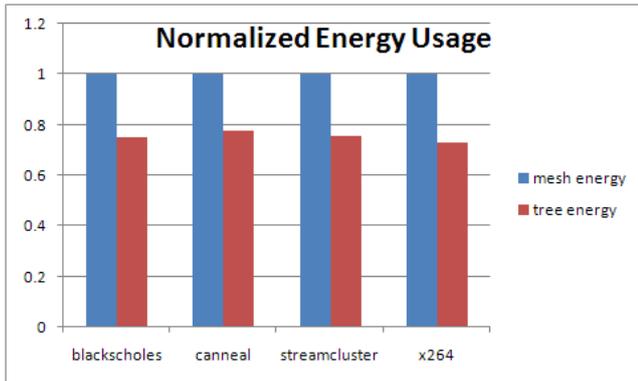
0.09. In the TFT network, the load for the routers was below 0.08 for the routers going toward the memory controller. We observed larger activity factors for the TFT routers returning from the memory controller (as much as 0.565 for the Level 1 routers in the **canneal** benchmark, but these routers have a small number of input ports, so their dynamic power dissipation remains low.

These power measurements suggest that a majority of our power savings comes from reducing static power dissipation. Because the original Mesh network used 64 routers in the MDN, and our TFT topology uses only 56 routers with fewer input ports on average, the TFT appears to be a more energy-efficient topology.

The effects of the inter-router interconnect can be seen in Figure 12. The tree network dissipates less power than the mesh network for all benchmarks except for **canneal**. The reason for this is due to the higher operating throughput of the tree network than the mesh network, especially for this specific workload. This higher throughput is achieved due to the nature of the traffic distribution for the **canneal**: 33% are stores, compared to 1% stores for **streamcluster**. In **streamcluster**, the large number of loads makes a bottleneck at the memories on the return path. The 288 bit load responses that are returning on the MDN completely congest that part of the network. **canneal**, on the other hand, balances traffic types better. As store requests are nine times

Table 7: Throughput of Selected MDN Routers

Benchmark	Router	Topology	Throughput (bits/cycle)
streamcluster	memory controller, outgoing	mesh	3.31
streamcluster	memory controller, incoming	mesh	28.98
streamcluster	top level, outgoing	tree	4
streamcluster	top level, incoming	tree	35
canneal	memory controller, outgoing	mesh	13.96
canneal	memory controller, incoming	mesh	27.07
canneal	top level, outgoing	tree	18.64
canneal	top level, incoming	tree	36.17

**Figure 13: Total power dissipated in routers and interconnect.**

larger than load requests on the request network, and load responses are nine times larger than store responses, traffic is distributed more evenly among the request and response networks, allowing higher network utilization.

This higher power usage does not cause an increase in energy required to complete the workload. To the contrary, the tree network uses less energy to finish all of the benchmarks. The greater power usage in **canneal** was due solely to the fact that the network was operating at a sufficiently higher throughput. It did more work per cycle, enough to result in higher power requirements despite lower per unit of work energy required. The total energy used to complete each workload is shown in Figure 13, which includes both intra-route energy and inter-router energy. The tree network consistently uses between 70%-80% of the energy consumed by the mesh network.

7. RELATED WORK

Interconnection networks have been studied for decades now. An excellent reference for all aspects of interconnection networks is [15]. The mesh networks used by the Tile64 are instances of k-ary n-cubes, which have been popular interconnection networks for years. Hypercubes were used for many parallel computers before Dally pointed out the benefits of lower radix k-ary n-cubes when under realistic constraints such as wiring and pin density [16]. Low dimensional meshes have been very popular interconnects since. Fat trees were first introduced in [17]. Generalized fat trees, of which the tapered fat tree used in the paper is a variant of, were described in detail in [18]. The tapered variant that

we employ was analyzed in terms of power, performance, and area usage in [7], prompting us to examine it for our special-purpose L2 to memory interconnect. Virtual channel flow control, employed in our modified design to allow the TDN to run without deadlock, was introduced in [19] and further developed [20].

Many different benchmarks for parallel workloads have been proposed over the years, the most commonly used (until recently) were the SPLASH and SPLASH2 benchmarks [3]. The authors of the PARSEC benchmark suite believe that the SPLASH2 benchmarks are now outdated and not representative of future CMP workloads [2]. The increased demands of the PARSEC benchmarks will be a primary motivating factor for re-examining the off-chip memory interconnect.

8. CONCLUSION

In this work, we proposed a detailed solution for providing an energy-efficient topology for off-chip processor-memory interconnect to improve the performance of memory-intensive parallel applications. We have demonstrated that the tapered fat-tree topology improves performance by up to 33.6% for applications with high off-chip bandwidth requirements without reducing the performance of those with mostly on-chip memory traffic. By reducing the number of total number of routers and simplifying the routers' crossbars, we have also shown that TFT offers substantial energy savings in all of the benchmarks, regardless of working set size or the amount of cache-line sharing. This energy savings translates into lower power requirements for all of the benchmarks but the most intensive.

As the number of cores on a single chip continues to grow and the number of memory access points scales less aggressively, we feel the study of similar "many-to-few" topologies warrants additional attention. Although design and physical layout of non-mesh networks is more difficult, the potential power savings and performance improvements may soon outweigh the added design complexity.

9. ACKNOWLEDGMENTS

We would like to thank Professor Krste Asanovic, Professor John Kubiawicz, and Christopher Batten for their insight and guidance throughout the life of this project.

10. REFERENCES

- [1] Asanovic et al., "The Landscape of Parallel Computing Research: A View from Berkeley", Tech. Rep., UC Berkeley, Dec 2006.

- [2] Bienia et al., “The PARSEC Benchmark Suite: Characterization and Architectural Implications”, Tech. Rep., Princeton University, Jan 2008.
- [3] “Stanford Parallel Applications for Shared Memory”, <http://www-flash.stanford.edu/SPLASH/>.
- [4] Leon et al., “A Power-Efficient High-Throughput 32-Thread SPARC Processor”, in *International Solid-State Circuits Conference (ISSCC)*, Feb 2006.
- [5] Konstadinidis et al., “Implementation of a Third Generation 16-Core 32-Thread Chip-Multithreading SPARC Processor”, in *International Solid-State Circuits Conference (ISSCC)*, Feb 2008.
- [6] Wentzlaff et al., “On-Chip Interconnection Architecture of the Tile Processor”, in *Micro, IEEE*, 2007.
- [7] James Balfour and William J. Dally, “Design Tradeoffs for Tiled CMP On-Chip Networks”, in *International Conference on Supercomputing (ICS)*, 2006.
- [8] Ron Ho, Kenneth W. Mai, and Mark A. Horowitz, “The Future of Wires”, in *Proceedings of the IEEE*, April 2001.
- [9] *International Technology Roadmap for Semiconductors (ITRS)*, 2003.
- [10] Guy Macon, “TILE64 Embedded Multicore Processors”, <https://www.gelato.unsw.edu.au/archives/comp-arch/2007-October/013123.html>.
- [11] “GeForce 8 Series”, <http://en.wikipedia.org/wiki/GeForce-8-Series>.
- [12] Bell et al., “TILE64 Processor: A 64-Core SoC with Mesh Interconnect”, in *IEEE International Solid-State Circuits Conference (ISSCC)*, Feb 2008.
- [13] Hangsheng Wang, “A Detailed Architectural-Level Power Model for Router Buffers, Crossbars and Arbiters”, Tech. Rep., Princeton University, Jan 2004.
- [14] Xuning Chen and Li-Shiuan Peh, “Leakage Power Modeling and Optimization in Interconnection Networks”, in *International Symposium on Low Power Electronics and Design (ISLPED)*, August 2003.
- [15] William James Dally and Brian Towles, *Principles and Practices of Interconnection Networks*, Morgan Kaufmann, 2004.
- [16] William J. Dally, “Performance analysis of k-ary n-cube interconnection networks”, in *IEEE Transactions on Computers*, June 1991.
- [17] C. E. Leiserson, “Fat-trees: Universal networks for hardware-efficient supercomputing”, in *IEEE Transactions on Computers*, Oct 1985.
- [18] Sabine Ohring, Maximilian Ibel, and Sajal Das, “On Generalized Fat Trees”, in *Parallel Processing Symposium*, 1995.
- [19] William J. Dally, “Virtual channel flow control”, in *International Symposium on Computer Architecture (ISCA)*, May 1990.
- [20] William J. Dally, “Virtual-channel flow control”, in *IEEE Transactions on Parallel and Distributed Systems*, March 1992.
- [21] Amit Kumar, Li-Shiuan Peh, Partha Kundu, and Niraj K. Jha, “Express Virtual Channels: Towards the Ideal Interconnection Fabric”, in *International Symposium on Computer Architecture (ISCA)*, June