



CACTI 5.1

Shyamkumar Thoziyoor, Naveen Muralimanohar, Jung Ho Ahn, and Norman P. Jouppi
HP Laboratories, Palo Alto
HPL-2008-20
April 2, 2008*

cache,
memory, area,
power, access
time, DRAM

CACTI 5.1 is a version of CACTI 5 fixing a number of small bugs in CACTI 5.0. CACTI 5 is the latest major revision of the CACTI tool for modeling the dynamic power, access time, area, and leakage power of caches and other memories. CACTI 5 includes a number of major improvements over CACTI 4. First, as fabrication technologies enter the deep-submicron era, device and process parameter scaling has become non-linear. To better model this, the base technology modeling in CACTI 5 has been changed from simple linear scaling of the original CACTI 0.8 micron technology to models based on the ITRS roadmap. Second, embedded DRAM technology has become available from some vendors, and there is interest in 3D stacking of commodity DRAM with modern chip multiprocessors. As another major enhancement, CACTI 5 adds modeling support of DRAM memories. Third, to support the significant technology modeling changes above and to enable fair comparisons of SRAM and DRAM technology, the CACTI code base has been extensively rewritten to become more modular. At the same time, various circuit assumptions have been updated to be more relevant to modern design practice. Finally, numerous bug fixes and small feature additions have been made. For example, the cache organization assumed by CACTI is now output graphically to assist users in understanding the output generated by CACTI.

CACTI 5.1

Shyamkumar Thoziyoor, Naveen Muralimanohar, Jung Ho Ahn, and Norman P. Jouppi
sthoziyo@cse.nd.edu, naveen@cs.utah.edu, jung-ho.ahn@hp.com, norm.jouppi@hp.com

April 1, 2008

Abstract

CACTI 5.1 is a version of CACTI 5 fixing a number of small bugs in CACTI 5.0. CACTI 5 is the latest major revision of the CACTI tool for modeling the dynamic power, access time, area, and leakage power of caches and other memories. CACTI 5 includes a number of major improvements over CACTI 4. First, as fabrication technologies enter the deep-submicron era, device and process parameter scaling has become non-linear. To better model this, the base technology modeling in CACTI 5 has been changed from simple linear scaling of the original CACTI 0.8 micron technology to models based on the ITRS roadmap. Second, embedded DRAM technology has become available from some vendors, and there is interest in 3D stacking of commodity DRAM with modern chip multiprocessors. As another major enhancement, CACTI 5 adds modeling support of DRAM memories. Third, to support the significant technology modeling changes above and to enable fair comparisons of SRAM and DRAM technology, the CACTI code base has been extensively rewritten to become more modular. At the same time, various circuit assumptions have been updated to be more relevant to modern design practice. Finally, numerous bug fixes and small feature additions have been made. For example, the cache organization assumed by CACTI is now output graphically to assist users in understanding the output generated by CACTI.

Contents

1	Introduction	5
2	Changes and Enhancements in Version 5	5
2.1	Organizational Changes	5
2.2	Circuit and Sizing Changes	5
2.3	Technology Changes	6
2.4	DRAM Modeling	7
2.5	Miscellaneous Changes	8
2.5.1	Optimization Function Change	8
2.5.2	New Gate Area Model	8
2.5.3	Wire Model	8
2.5.4	ECC and Redundancy	8
2.5.5	Display Changes	9
3	Data Array Organization	9
3.1	Mat Organization	10
3.2	Routing to Mats	11
3.3	Organizational Parameters of a Data Array	13
3.4	Comments about Organization of Data Array	14
4	Circuit Models and Sizing	16
4.1	Wire Modeling	16
4.2	Sizing Philosophy	17
4.3	Sizing of Mat Circuits	17
4.3.1	Predecoder and Decoder	17
4.3.2	Bitline Peripheral Circuitry	19
4.4	Sense Amplifier Circuit Model	21
4.5	Routing Networks	22
4.5.1	Array Edge to Bank Edge H-tree	22
4.5.2	Bank Edge to Mat H-tree	22
5	Area Modeling	23
5.1	Gate Area Model	24
5.2	Area Model Equations	25
6	Delay Modeling	29
6.1	Access Time Equations	30
6.2	Random Cycle Time Equations	30
7	Power Modeling	31
7.1	Calculation of Dynamic Energy	31
7.1.1	Dynamic Energy Calculation Example for a CMOS Gate Stage	31
7.1.2	Dynamic Energy Equations	32
7.2	Calculation of Leakage Power	33
7.2.1	Leakage Power Calculation for CMOS gates	33
7.2.2	Leakage Power Equations	34

8 Technology Modeling	35
8.1 Devices	35
8.2 Wires	37
8.3 Technology Exploration	38
9 Embedded DRAM Modeling	38
9.1 Embedded DRAM Modeling Philosophy	38
9.1.1 Cell	39
9.1.2 Destructive Readout and Writeback	39
9.1.3 Sense Amplifier Input Signal	39
9.1.4 Refresh	39
9.1.5 Wordline Boosting	39
9.2 DRAM Array Organization and Layout	40
9.2.1 Bitline Multiplexing	40
9.2.2 Reference Cells for V_{DD} Precharge	40
9.3 DRAM Timing Model	40
9.3.1 Bitline Model	40
9.3.2 Multisubbank Interleave Cycle Time	42
9.3.3 Retention Time and Refresh Period	42
9.4 DRAM Power Model	43
9.4.1 Refresh Power	43
9.5 DRAM Area Model	43
9.5.1 Area of Reference Cells	43
9.5.2 Area of Refresh Circuitry	43
9.6 DRAM Technology Modeling	44
9.6.1 Cell Characteristics	44
10 Cache Modeling	45
10.1 Organization	45
10.2 Delay Model	46
10.3 Area Model	47
10.4 Power Model	47
11 Quantitative Evaluation	47
11.1 Evaluation of New CACTI 5 Features	47
11.1.1 Impact of New CACTI Solution Optimization	48
11.1.2 Impact of Device Technology	49
11.1.3 Impact of Interconnect Technology	52
11.1.4 Impact of RAM Cell Technology	53
11.2 Version 4.2 vs Version 5.1 Comparisons	55
12 Validation	61
12.1 Sun SPARC 90nm L2 cache	61
12.2 Intel Xeon 65nm L3 cache	61
13 Commodity DRAM Technology and Main Memory Chip Modeling	65
14 Future Work	66
15 Conclusions	67
A Additional CACTI Validation Results for 90nm SPARC L2	68

1 Introduction

CACTI 5 is the latest major revision of the CACTI tool [30, 38, 40, 47] for modeling the dynamic power, access time, area, and leakage power of caches and other memories. CACTI 5.1 is a version of CACTI 5 fixing a number of small bugs in CACTI 5.0. CACTI has become widely used by computer architects, both directly and indirectly through other tools such as Wattch.

CACTI 5 includes a number of major improvements over CACTI 4.0. First, as fabrication technologies enter the deep-submicron era, device and process parameter scaling has become non-linear. To better model this, the base technology modeling in CACTI 5 has been changed from simple linear scaling of the original 0.8 micron technology to models based on the ITRS roadmap. Second, embedded DRAM technology has become available from some vendors, and there is interest in 3D stacking of commodity DRAM with modern chip multiprocessors. As another major enhancement, CACTI 5 adds modeling support of DRAM memories. Third, to support the significant technology modeling changes above and to enable fair comparisons of SRAM and DRAM technology, the CACTI code base has been extensively rewritten to become more modular. At the same time, various circuit assumptions have been updated to be more relevant to modern design practice. Finally, numerous bug fixes and small feature improvements have been made. For example, the cache organization assumed by CACTI is now output graphically by the web-based server, to assist users in understanding the output generated by CACTI.

The following section gives an overview of these changes, after which they are discussed in detail in subsequent sections.

2 Changes and Enhancements in Version 5

2.1 Organizational Changes

Earlier versions of CACTI (up to version 3.2) made use of a single row predecoder at the center of a memory bank with the row predecoded signals being driven to the subarrays for decoding. In version 4.0, this centralized decoding logic was implicitly replaced with distributed decoding logic. Using H-tree distribution, the address bits were transmitted to the distributed sinks where the decoding took place. However, because of some inconsistencies in the modeling, it was not clear at what granularity the distributed decoding took place - whether there was one sink per subarray or 2 or 4 subarrays. There were some other problems with the CACTI code such as the following:

- The area model was not updated after version 3.2, so the impact on area of moving from centralized to distributed decoding was not captured. Also, the leakage model did not account for the multiple distributed sinks. The impact of cache access type (normal/sequential/fast) [40] on area was also not captured;
- Number of address bits routed to the subarrays was being computed incorrectly;
- Gate load seen by NAND gate in the 3-8 decode block was being computed incorrectly; and
- There were problems with the logic computing the degree of muxing at the tristate subarray output drivers.

In version 5, we resolve these issues, redefine and clarify what the organizational assumptions of memory are and remove ambiguity from the modeling. Details about the organization of memory can be found in Section 3.

2.2 Circuit and Sizing Changes

Earlier versions of CACTI made use of row decoding logic with two stages - the first stage was composed of 3-8 predecode blocks (composed of NAND3 gates) followed by a NOR decode gate and wordline driver. The number of gates in the row decoding path was kept fixed and the gates were then sized using the method of logical effort [39] for an effective fanout of 3 per stage. In version 5, in addition to the row decoding logic, we also model the bitline mux decoding logic and the sense-amplifier mux decoding logic. We use the same circuit structures to model all decoding logic and we base the modeling on the effort described in [3]. We use the sizing heuristic described in [3] that has been shown to be good from an energy-delay perspective. With the new circuit structures and modeling that we use, the limit

on maximum number of signals that can be decoded is increased from 4096 (in version 4.2) to 262144 (in version 5). While we do not expect the number of signals that are decoded to be very high, extending the limit from 4096 helps with exploring area/delay/power tradeoffs in a more thorough manner for large memories, especially for large DRAMs. Details of the modeling of decoding logic are described in Section 4.

There are certain problems with the modeling of the H-tree distribution network in version 4.2. An inverter-driver is placed at branches of the address, datain, and dataout H-tree. However, the dataout H-tree does not model tristate drivers. The output data bits may come from a few subarrays and so the address needs to be distributed to a few subarrays, however, dynamic power spent in transmitting address is computed as if all the data comes from a single subarray. The leakage in the drivers of the datain H-tree is not modeled.

In version 5, we model the H-tree distribution network more rigorously. For the dataout H-tree we model tristate buffers at each branch. For the address and datain H-trees, instead of assuming inverters at the branches of the H-tree we assume the use of buffers that may be gated to allow or disallow the passage of signals and thereby control the dynamic power. We size these drivers based on the methodology described in [3] which takes the resistance and capacitance of intermediate wires into account during sizing. We also model the use of repeaters in the H-tree distribution network which are sized according to equations from [11].

2.3 Technology Changes

Earlier versions of CACTI relied on a complicated way of obtaining device data for the input technology-node. Computation of access/cycle time and dynamic power were based off device data of a 0.8-micron process that was scaled to the given technology-node using simple linear scaling principles. Leakage power calculation, however, made use of I_{off} (subthreshold leakage current) values that were based off device data obtained through BSIM3 parameter extractions. In version 4.2, BSIM3 extraction was carried out for a few select technology nodes (130/100/70nm); as a result leakage power estimation was available only for these select technology nodes.

There are several problems with the above approach of obtaining device data. Using two sets of parameters, one for computation of access/cycle time/dynamic power and another for leakage power, is a convoluted approach and is hard to maintain. Also, the approach of basing device parameter values off a 0.8-micron process is not a good one because of several reasons. Device scaling has become quite non-linear in the deep-submicron era. Device performance targets can no longer be achieved through simple linear scaling of device parameters. Moreover, it is well-known that physical gate-lengths (according to the ITRS, physical gate-length is the final, as-etched length of the bottom of the gate electrode) have scaled much more aggressively [4, 35] than what would be projected by simple linear scaling from the 0.8 micron process.

In version 5, we adopt a simpler, more evolvable approach of obtaining device data. We use device data that the ITRS [35] uses to make its projections. The ITRS makes use of the MASTAR software tool (Model for Assessment of CMOS Technologies and Roadmaps) [36] for computation of device characteristics of current and future technology nodes. Using MASTAR, device parameters may be obtained for different technologies such as planar bulk, double gate and Silicon-On-Insulator. MASTAR includes device profile and result files of each year/technology-node for which the ITRS makes projections and we incorporate the data from these files into CACTI. These device profiles are based off published industry process data and industry-consensus targets set by historical trends and system drivers. While it is not necessary that these device numbers match or would match process numbers of various vendors in an exact manner, they do come within the same ball-park as can be seen by looking at the Ion-Ioff cloud graphic within the MASTAR software which shows a scatter plot of various published vendor Ion-Ioff numbers and corresponding ITRS projections. With this approach of using device data from the ITRS, it also becomes possible to incorporate device data corresponding to different device types that the ITRS defines such as high performance (HP), LSTP (Low Standby Power), and Low Operating Power (LOP). More details about the device data used in CACTI can be found in Section 8.

There are some problems with interconnect modeling of version 4.2 also. Version 4.2 utilizes 2 types of wires in the delay model, ‘local’ and ‘global’. The local type is used for wordlines and bitlines, while the global type is used for all other wires. The resistance per unit length and capacitance per unit length for these two wire types are also calculated in a convoluted manner. For a given technology, the resistance per unit length of the local wire is calculated by assuming ideal scaling in all dimensions and using base data of a 0.8-micron process. The base resistance per unit length for the 0.8-micron process is itself calculated by assuming copper wires in the base 0.8-micron process and readjusting the

sheet resistance value of version 3.2 which assumed aluminium wires. As the resistivity of copper is about 2/3rd that of aluminium, the sheet resistance of copper was computed to be 2/3rd that of aluminium. However, this implies that the thickness of metal assumed in versions 3.2 and 4.2 are the same which turns out to be not true. When we compute sheet resistance for the 0.8-micron process with the thickness of local wire assumed in version 4.2 and assuming a resistivity of $2.2 \mu\text{ohm-cm}$ for copper, the value comes out to be a factor of 3.4 smaller than that used in version 3.2. In version 4.2, resistance per unit length for the global wire type is calculated to be smaller than that of local wire type by a factor of 2.04. This factor of 2.04 is calculated based on RC delays and wire sizes of different wire types in the 2004 ITRS but the underlying assumptions are not known. Another problem is that even though the delay model makes use of two types of wires, local and global, the area model makes use of just the local wire type and the pitch calculation of all wires (local type and global type) are based off the assumed width and spacing for the local wire type; this results in an underestimation of pitch (and area) occupied by the global wires.

Capacitance per unit length calculation of version 4.2 also suffers from certain problems. The capacitance per unit length values for local and global wire types are assumed to remain constant across technology nodes. The capacitance per unit length value for local wire type was calculated for a 65nm process as $(2.9/3.6)*230 = 185 \text{ fF/m}$ where 230 is the published capacitance per unit length value for an Intel 130nm process [42], 3.6 is the dielectric constant of the 130 nm process and 2.9 is the dielectric constant of an Intel 65nm process [4]. Computing the value of capacitance per unit length in this manner for a 65nm process ignores the fact that the fringing component of capacitance remains almost constant across technology-nodes and scales very slowly [11, 31]. Also, assuming that the dielectric constant remains fixed at 2.9 for future technology nodes ignores the possibility of use of lower-k dielectrics. Capacitance per unit length of the global type wire of version 4.2 is calculated to be smaller than that of local type wires by a factor of 1.4. This factor of 1.4 is again calculated based on RC delays and wire sizes of different wire types in the 2004 ITRS but the underlying assumptions again are not known.

In version 5, we remove the ambiguity from the interconnect modeling. We use the interconnect projections made in [11, 13] which are based off well-documented simple models of resistance and capacitance. Because of the difficulty in projecting the values of interconnect properties in an exact manner at future technology nodes the approach employed in [11, 13] was to come up with two sets of projections based on aggressive and conservative assumptions. The aggressive projections assume aggressive use of low-k dielectrics, insignificant resistance degradation due to dishing and scattering, and tall wire aspect ratios. The conservative projections assume limited use of low-k dielectrics, significant resistance degradation due to dishing and scattering, and smaller wire aspect ratios. We incorporate both sets of projections into CACTI. We also model 2 types of wires inside CACTI - semi-global and global with properties identical to that described in [11, 13]. More details of the interconnect modeling are described in Section 8.2. Comparisons of area, delay, and power of caches obtained using versions 4.2 and 5 are presented in Section 11.2.

2.4 DRAM Modeling

One of the major enhancements of version 5 is the incorporation of embedded DRAM models for a logic-based embedded DRAM fabrication process [19, 24, 27]. In the last few years, embedded DRAM has made its way into various applications. The IBM POWER4 made use of embedded DRAM in its L3 cache [41]. The main compute chip inside the Blue Gene/L supercomputer also makes use of embedded DRAM [14]. Embedded DRAM has also been used in the graphics synthesizer unit of Sony's PlayStation 2 [28].

In our modeling of embedded DRAM, we leverage the similarity that exists in the global and peripheral circuitry of embedded SRAM and DRAM and model only their essential differences. We use the same array organization for embedded DRAM that we used for SRAM. By having a common framework that, in general, places embedded SRAM and DRAM on an equal footing and emphasizes only their essential differences, we are able to compare relative tradeoffs between embedded SRAM and DRAM. We describe the modeling of embedded DRAM in Section 9.

2.5 Miscellaneous Changes

2.5.1 Optimization Function Change

In version 5, we follow a different approach in finding the optimal solution with CACTI. Our new approach allows users to exercise more control on area, delay, and power of the final solution. The optimization is carried out in the following steps: first, we find all solutions with area efficiency that is within a certain percentage (user-supplied value) of the area efficiency of the solution with best area efficiency. We refer to this area constraint as `max_area_constraint`. Next, from this reduced set of solutions that satisfy the `max_area_constraint`, we find all solutions with access time that is within a certain percentage of the best access time solution (in the reduced set). We refer to this access time constraint as `max_acc_time_constraint`. To the subset of solutions that results after the application of `max_acc_time_constraint`, we apply the following optimization function:

$$\begin{aligned} \text{optimization-func} = & \frac{\text{dynamic-energy}}{\text{min-dynamic-energy}} \text{flag-opt-for-dynamic-energy} + \\ & \frac{\text{dynamic-power}}{\text{min-dynamic-power}} \text{flag-opt-for-dynamic-power} + \\ & \frac{\text{leak-power}}{\text{min-leak-power}} \text{flag-opt-for-leak-power} + \\ & \frac{\text{rand-cycle-time}}{\text{min-rand-cycle-time}} \text{flag-opt-for-rand-cycle-time} \end{aligned}$$

where `dynamic-energy`, `dynamic-power`, `leak-power`, and `rand-cycle-time` are the dynamic energy, dynamic power, leakage power, and random cycle time of a solution respectively and `min-dynamic-energy`, `min-dynamic-power`, `min-leak-power`, and `min-rand-cycle-time` are their minimum (best) values in the subset of solutions being considered. `flag-opt-for-dynamic-energy`, `flag-opt-for-dynamic-power`, `flag-opt-for-leak-power`, and `flag-opt-for-rand-cycle-time` are user-specified boolean variables. The new optimization process allows exploration of the solution space in a controlled manner to arrive at a solution with user-desired characteristics.

2.5.2 New Gate Area Model

In version 5, we introduce a new analytical gate area model from [49]. With the new gate area model it becomes possible to make the areas of gates sensitive to transistor sizing so that when transistor sizing changes, the areas also change. With the new gate area model, transistors may get folded when they are subject to pitch-matching constraints and the area is calculated accordingly. This feature is useful in capturing differences in area caused due to different pitch-matching constraints that may have to be satisfied, particularly between SRAM and DRAM.

2.5.3 Wire Model

Version 4.2 models wires using the equivalent circuit model shown in Figure 1(a). The Elmore delay of this model is $RC/2$, however this model underestimates the wire-to-gate component ($R_{\text{wire}}C_{\text{gate}}$) of delay. In version 5, we replace this model with the Π RC model, shown in Figure 1(b), which has been used in more recent SRAM modeling efforts [2].

2.5.4 ECC and Redundancy

In order to be able to check and correct soft errors, most memories of today have support for ECC (Error Correction Code). In version 5, we capture the impact of ECC by incorporating a model that captures the ECC overhead in memory cell and data bus (datain and dataout) area. We incorporate a variable that specifies the number of data bits per ECC bit. By default, we fix the value of this variable to 8.

In order to improve yield, many memories of today incorporate redundant entities even at the subarray level. For example, the data array of the 16MB Intel Xeon L3 cache [7] which has 256 subarrays also incorporates 32 redundant subarrays. In version 5, we incorporate a variable that specifies the number of mats per redundant mat. By default, we fix the value of this variable to 8.

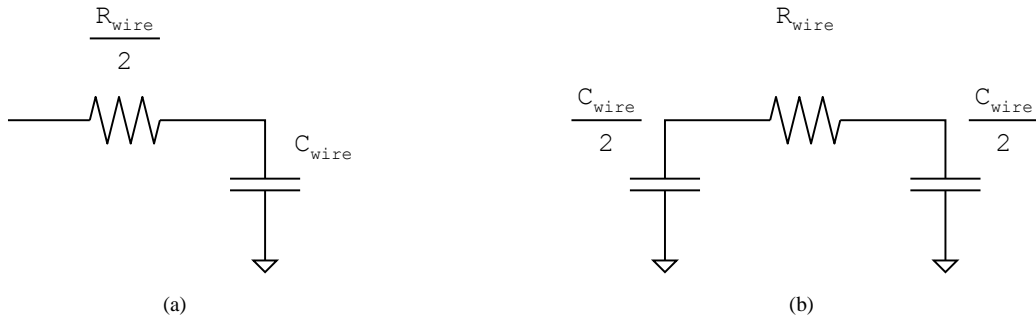


Figure 1: (a) L-model of wire used in version 4.2, (b) Π RC model of wire used in version 5.

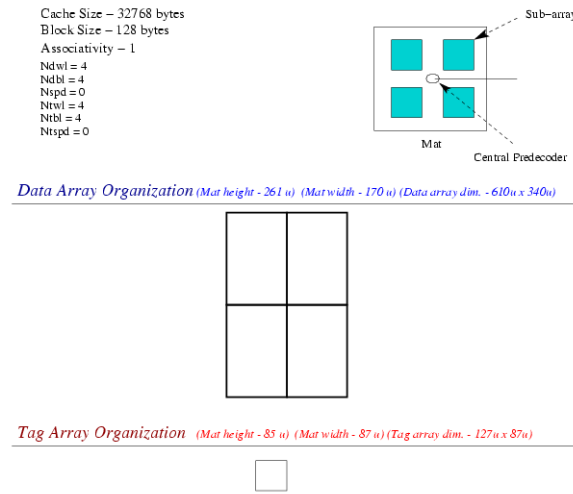


Figure 2: Example of the graphical display generated by version 5.

2.5.5 Display Changes

To facilitate better understanding of cache organization, version 5 can output data/tag array organization graphically. Figure 2 shows an example of the graphical display generated by version 5. The top part of the figure shows a generic mat organization assumed by CACTI. It is followed by the data and tag array organization plotted based on array dimensions calculated by CACTI.

3 Data Array Organization

At the highest level, a data array is composed of multiple identical banks (N_{banks}). Each bank can be concurrently accessed and has its own address and data bus. Each bank is composed of multiple identical subbanks (N_{subbanks}) with one subbank being activated per access. Each subbank is composed of multiple identical mats ($N_{\text{mats-in-subbank}}$). All mats in a subbank are activated during an access with each mat holding part of the accessed word in the bank. Each mat

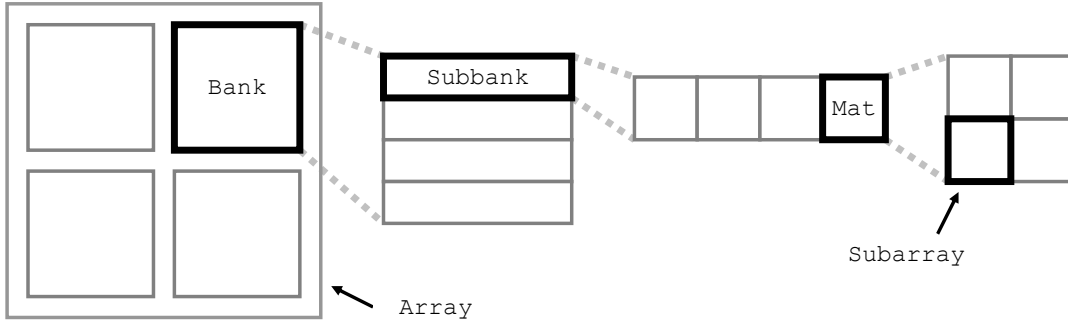


Figure 3: Layout of an example array with 4 banks. In this example each bank has 4 subbanks and each subbank has 4 mats.

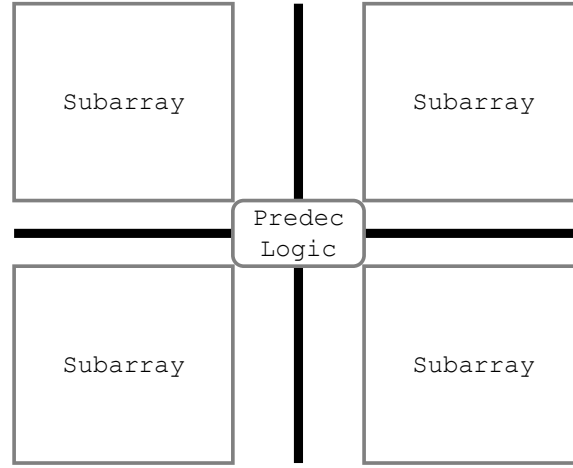


Figure 4: High-level composition of a mat.

itself is a self-contained memory structure composed of 4 identical subarrays and associated predecoding logic. Each subarray is a 2D matrix of memory cells and associated peripheral circuitry. Figure 3 shows the layout of an array with 4 banks. In this example each bank is shown to have 4 subbanks and each subbank is shown to have 4 mats. Not shown in Figure 3, address and data are assumed to be distributed to the mats on H-tree distribution networks.

The rest of this section further describes details of the array organization assumed in CACTI. Section 3.1 describes the organization of a mat. Section 3.2 describes the organization of the H-tree distribution networks. Section 3.3 presents the different organizational parameters associated with a data array.

3.1 Mat Organization

Figure 4 shows the high-level composition of all mats. A mat is always composed of 4 subarrays and associated predecoding/decoding logic which is located at the center of the mat. The predecoding/decoding logic is shared by all 4 subarrays. The bottom subarrays are mirror images of the top subarrays and the left hand side subarrays are mirror images of the right hand side ones. Not shown in this figure, by default, address/datain/dataout signals are assumed to enter the mat in the middle through its sides; alternatively, under user-control, it may also be specified to assume that they traverse over the memory cells.

Figure 5 shows the high-level composition of a subarray. The subarray consists of a 2D matrix of the memory cells and associated peripheral circuitry. Figure 6 shows the peripheral circuitry associated with bitlines of a subarray. After a wordline gets activated, memory cell data get transferred to bitlines. The bitline data may go through a level of bitline multiplexing before it is sensed by the sense amplifiers. Depending on the degree of bitline multiplexing, a single sense amplifier may be shared by multiple bitlines. The data is sensed by the sense amplifiers and then passed to tristate output

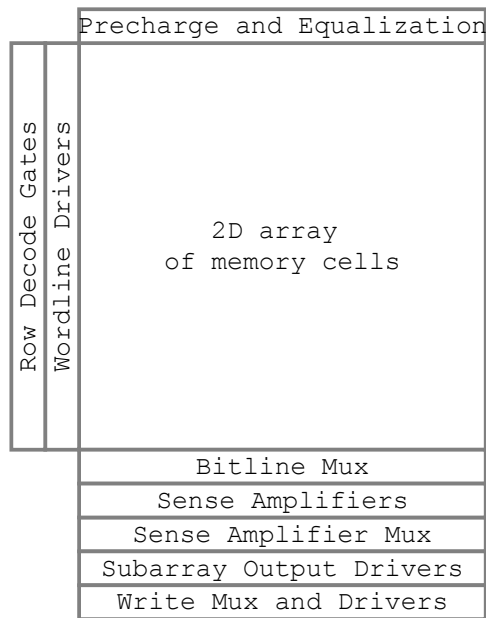


Figure 5: High-level composition of a subarray.

drivers which drive the dataout vertical H-tree (described later in this section). An additional level of multiplexing may be required at the outputs of the sense amplifiers in organizations in which the bitline multiplexing is not sufficient to cull out the output data or in set-associative caches in which the output word from the correct way needs to be selected. The select signals that control the multiplexing of the bitline mux and the sense amp mux are generated by the bitline mux select signals decoder and the sense amp mux select signals decoder respectively. When the degree of multiplexing after the outputs of the sense amplifiers is simply equal to the associativity of the cache, the sense amp mux select signal decoder does not have to decode any address bits and instead simply buffers the input way-select signals that arrive from the tag array.

3.2 Routing to Mats

Address and data are routed to and from the mats on H-tree distribution networks. H-tree distribution networks are used to route address and data and provide uniform access to all the mats in a large memory.¹ Such a memory organization is interconnect-centric and is well-suited for coping with the trend of worsening wire delay with respect to device delay. Rather than shipping a bunch of predecoded address signals to the mats, it makes sense to ship the address bits and decode them at the sinks (mats) [34]. Contemporary divided wordline architectures which make use of broadcast of global signals suffer from increased wire delay as memory capacities get larger [2]. Details of a memory organization similar to what we have assumed may also be found in [1]. For ease of pipelining multiple accesses in the array, separate request and reply networks are assumed. The request network carries address and datain from the edge of the array to the mats while the reply network carries dataout from the mats to the edge of the array. The structure of the request and reply networks is similar; here we discuss the high-level organization of the request network.

The request H-tree network is divided into two networks:

1. The H-tree network from the edge of the array to the edge of a bank; and,
2. The H-tree network from the edge of the bank to the mats.

Figure 7 shows the layout of the request H-tree network between the array edge and the banks. Address and datain are routed to each bank on this H-tree network and enter each bank at the middle from one of its sides. The H-tree

¹Non-uniform cache architectures (NUCA) are currently beyond the scope of CACTI 5 but may be supported by future versions of CACTI.

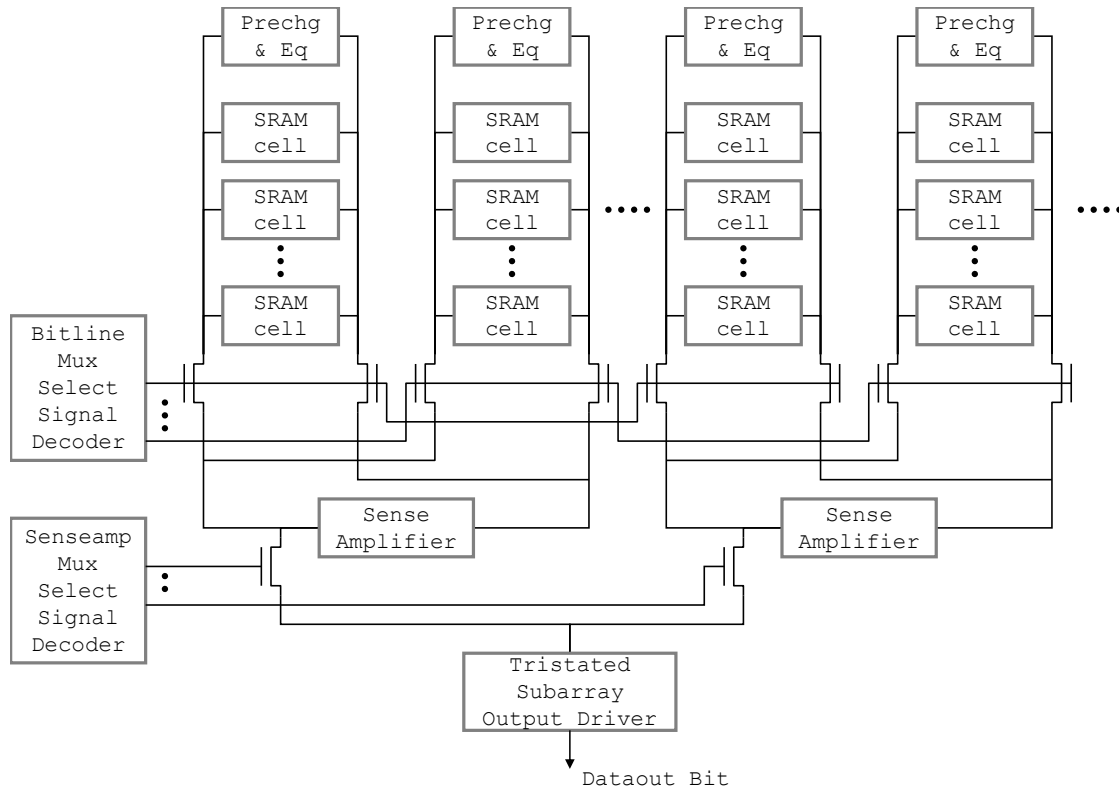


Figure 6: Peripheral circuitry associated with bitlines. Not shown in this figure, but the outputs of the muxes are assumed to be precharged high.

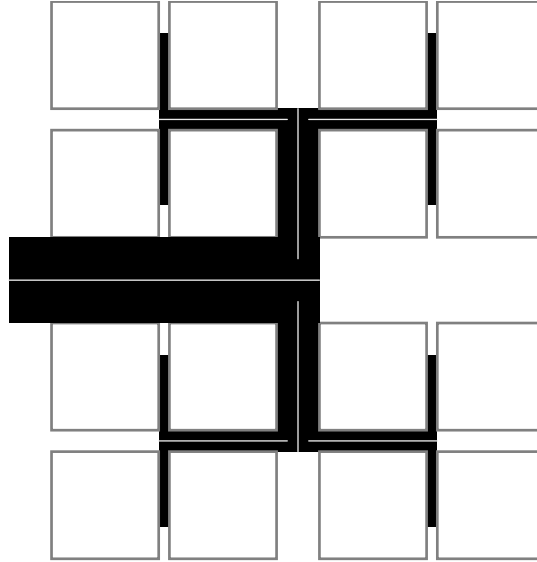


Figure 7: Layout of edge of array to banks H-tree network.

network from the edge of the bank to the mats is further divided into two 1-dimensional horizontal and vertical H-tree networks. Figure 8 shows the layout of the horizontal H-tree within a bank which is located at the middle of the bank while Figure 9 shows the layout of the vertical H-trees within a bank. The leaves of the horizontal H-tree act as the parent nodes (marked as V0) of the vertical H-trees. In order to understand the routing of signals on the H-tree networks within a bank, we use an illustrative example. Consider a bank with the following parameters: 1MB capacity, 256-bit

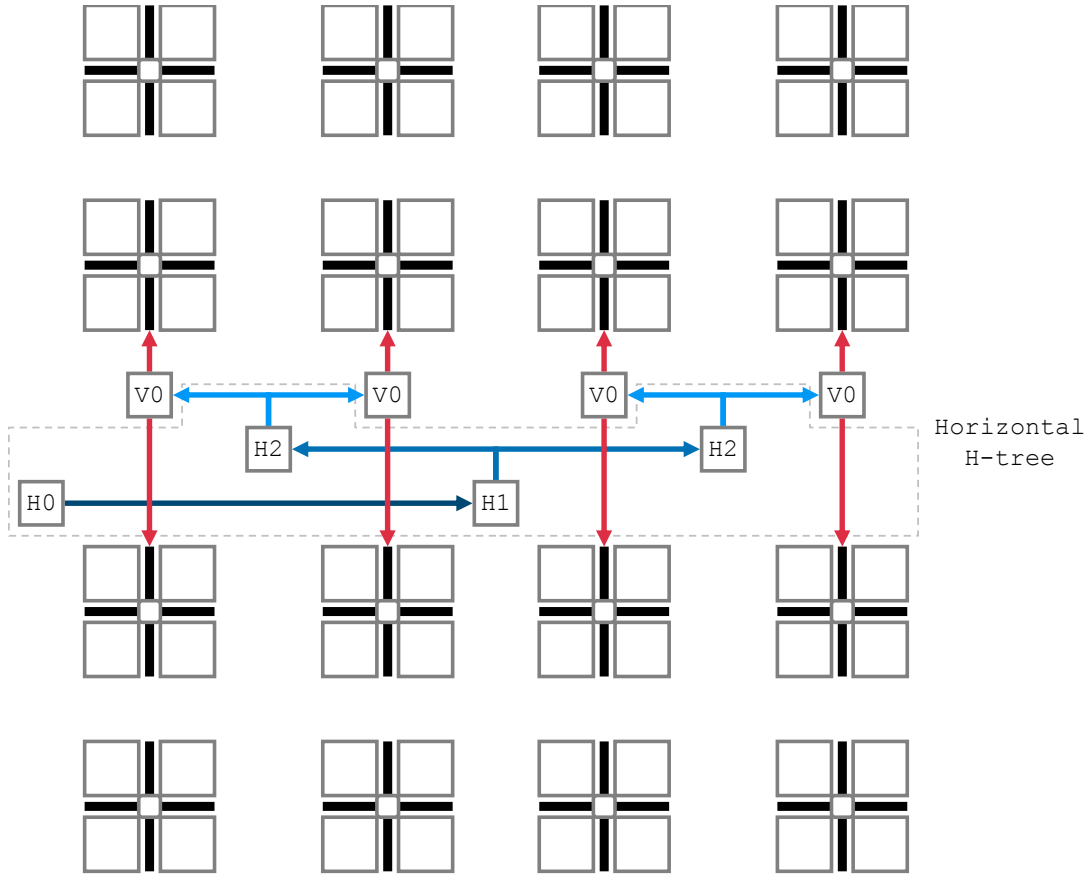


Figure 8: Layout of the horizontal H-tree within a bank.

output word, 4 subbanks, 4 mats in each subbank. Looked at together, Figures 8 and 9 can be considered to be the horizontal and vertical H-trees within such a bank. The number of address bits required to address a word in this bank is 15. As there are 4 subbanks and because each mat in a subbank is activated during an access, the number of address bits that need to be distributed to each mat is 13. Because each mat in a subbank produces 64 out of the 256 output bits, the number of datain signals that need to be distributed to each mat is 64. Thus 15 bits of address and 256 bits of datain enter the bank from the left side driven by the H0 node. At the H1 node, the 15 address signals are redriven such that each of the two nodes H1 receive the 15 address signals. The datain signals split at node H1 and 128 datain signals go to the left H2 node and the other 128 go to the right H2 node. At each H2 node, the address signals are again redriven such that all of the 4 V0 nodes end up receiving the 15 address bits. The datain signals again split at each H2 node so that each V0 node ends up receiving 64 datain bits. These 15 address bits and 64 datain bits then traverse to each mat along the 4 vertical H-trees. In the vertical H-trees, address and datain may either be assumed to be broadcast to all mats or alternatively, it may be assumed that these signals are appropriately gated so that they are routed to just the correct subbank that contains the data; by default, we assume the latter scenario.

The reply network H-trees are similar in principle to the request network H-trees. In case of the reply network vertical H-trees, dataout bits from each mat of a subbank travel on the vertical H-trees to the middle of the bank where they sink into the reply network horizontal H-tree, and are carried to the edge of the bank.

3.3 Organizational Parameters of a Data Array

In order to calculate the optimal organization based on a given objective function, like earlier versions of CACTI [30,38,40,47], each bank is associated with partitioning parameters N_{dwl} , N_{dbl} and N_{spd} , where N_{dwl} = number of segments in a bank wordline, N_{dbl} = number of segments in a bank bitline, and N_{spd} = number of sets mapped to each bank wordline.

Unlike earlier versions of CACTI, in CACTI 5 N_{spd} can take on fractional values less than one. This is useful for

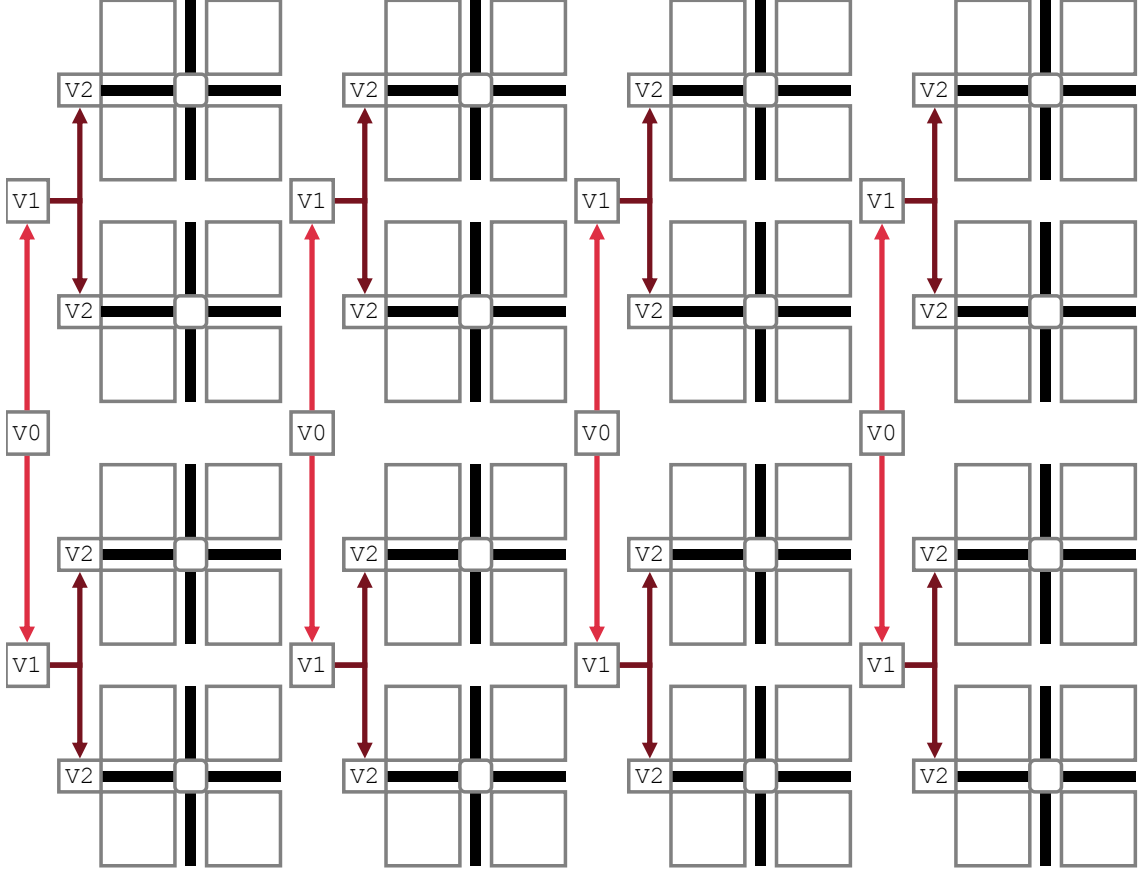


Figure 9: Layout of the vertical H-trees within a bank.

small highly-associative caches with large line sizes. Without values of N_{spd} less than one, memory mats with huge aspect ratios with only a few word lines but hundreds of bits per word line would be created. For a pure scratchpad memory (not a cache), N_{spd} is used to vary the aspect ratio of the memory bank.

N_{subbanks} and $N_{\text{mats-in-subbank}}$ are related to N_{dwl} and N_{dbl} as follows:

$$N_{\text{subbanks}} = \frac{N_{\text{dbl}}}{2} \quad (1)$$

$$N_{\text{mats-in-subbank}} = \frac{N_{\text{dwl}}}{2} \quad (2)$$

Figure 10 shows different partitions of the same bank. The partitioning parameters are labeled alongside. Table 1 lists various organizational parameters associated with a data array.

3.4 Comments about Organization of Data Array

The cache organization chosen in the CACTI model is a compromise between many possible different cache organizations. For example, in some organizations all the data bits could be read out of a single mat. This could reduce dynamic power but increase routing requirements. On the other hand, organizations exist where all mats are activated on a request and each produces part of the bits required. This obviously burns a lot of dynamic power, but has the smallest routing requirements. CACTI chooses a middle ground, where all the bits for a read come from a single subbank, but multiple mats. Other more complicated organizations, in which predecoders are shared by two subarrays instead of four, or in which sense amplifiers are shared between top and bottom subarrays, are also possible, however we try to model a simple common case in CACTI.

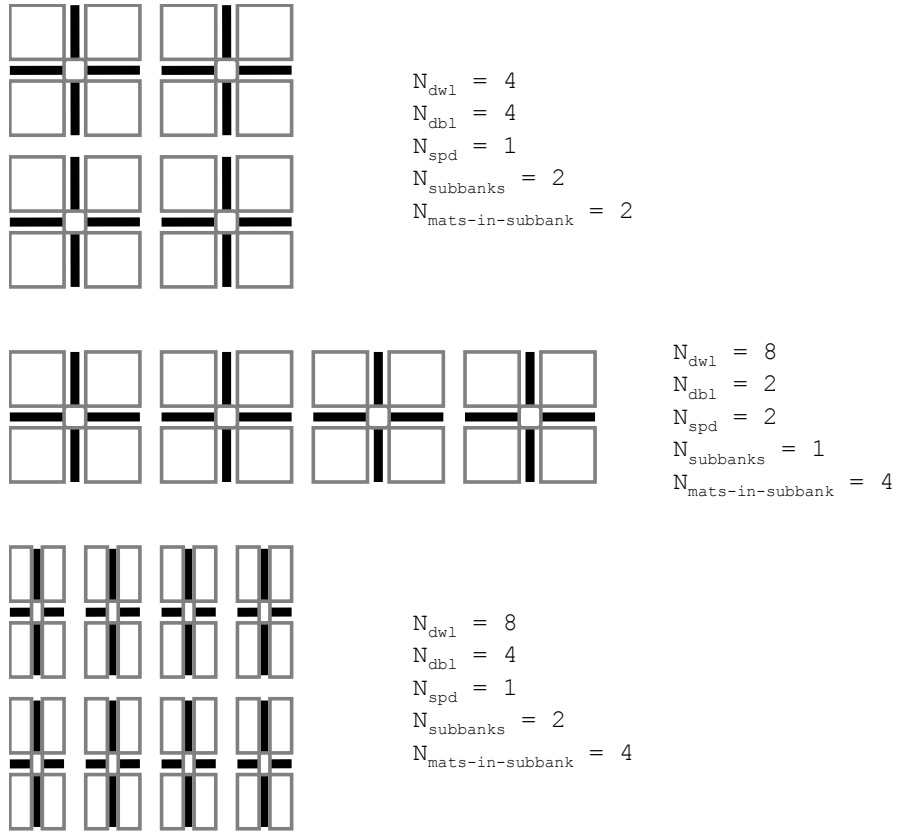


Figure 10: Different partitions of a bank.

Parameter Name	Meaning	Parameter Type
N_{banks}	Number of banks	User input
N_{dwl}	Number of divisions in a bank wordline	Degree of freedom
N_{dbl}	Number of divisions in a bank bitline	Degree of freedom
N_{spd}	Number of sets mapped to a bank wordline	Degree of freedom
$D_{\text{bitline-mux}}$	Degree of muxing at bitlines	Degree of freedom
$D_{\text{senseamp-mux}}$	Degree of muxing at sense amp outputs	Degree of freedom
N_{subbanks}	Number of subbanks	Calculated
$N_{\text{mats-in-subbank}}$	Number of mats in a subbank	Calculated
$N_{\text{subarr-rows}}$	Number of rows in a subarray	Calculated
$N_{\text{subarr-cols}}$	Number of columns in a subarray	Calculated
$N_{\text{subarr-senseamps}}$	Number of sense amplifiers in a subarray	Calculated
$N_{\text{subarr-out-drivers}}$	Number of output drivers in a subarray	Calculated
$N_{\text{bank-addr-bits}}$	Number of address bits to a bank	Calculated
$N_{\text{bank-datain-bits}}$	Number of datain bits to a mat	Calculated
$N_{\text{bank-dataout-bits}}$	Number of dataout bits from a mat	Calculated
$N_{\text{mat-addr-bits}}$	Number of address bits to a mat	Calculated
$N_{\text{mat-datain-bits}}$	Number of datain bits to a mat	Calculated
$N_{\text{mat-dataout-bits}}$	Number of dataout bits from a mat	Calculated
$N_{\text{mat-way-select}}$	Number of way-select bits to a mat (for data array of cache)	Calculated

Table 1: Organizational parameters of a data array.

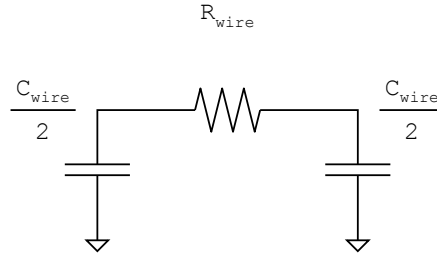


Figure 11: One-section Π RC model that we have assumed for non-ideal wires.

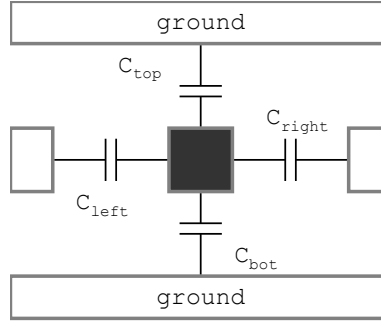


Figure 12: Capacitance model from [11].

4 Circuit Models and Sizing

In Section 3, the high-level organization of an array was described. In this section, we delve deeper into logic and circuit design of the different entities. We also present the techniques adopted for sizing different circuits. The rest of this section is organized as follows: First, in Section 4.1, we describe the circuit model that we have assumed for wires. Next in Section 4.2, we describe the general philosophy that we have adopted for sizing circuits. Next in Section 4.3, we describe the circuit models and sizing techniques for the different circuits within a mat, and in Section 4.5, we describe them for the circuits used in the different H-tree networks.

4.1 Wire Modeling

Wires are considered to belong to one of two types: ideal or non-ideal. Ideal wires are assumed to have zero resistance and capacitance. Non-ideal wires are assumed to have finite resistance and capacitance and are modeled using a one-section Π RC model shown in Figure 11. In this figure, R_{wire} and C_{wire} for a wire of length L_{wire} are given by the following equations:

$$R_{\text{wire}} = L_{\text{wire}} R_{\text{unit-length-wire}} \quad (3)$$

$$C_{\text{wire}} = L_{\text{wire}} C_{\text{unit-length-wire}} \quad (4)$$

For computation of $R_{\text{unit-length-wire}}$ and $C_{\text{unit-length-wire}}$ wires, we use the equations presented in [11, 13] which are reproduced below. Figure 12 shows the accompanying picture for the capacitance model from [11].

$$R_{\text{unit-length-wire}} = \alpha_{\text{scatter}} \frac{\rho}{(\text{thickness} - \text{barrier} - \text{dishing})(\text{width} - 2 * \text{barrier})} \quad (5)$$

$$C_{\text{unit-length-wire}} = \epsilon_0 \left(2M\epsilon_{\text{horiz}} \frac{\text{thickness}}{\text{spacing}} + 2\epsilon_{\text{vert}} \frac{\text{width}}{ILD_{\text{thick}}} \right) + \text{fringe}(\epsilon_{\text{horiz}}, \epsilon_{\text{vert}}) \quad (6)$$

4.2 Sizing Philosophy

In general the sizing of circuits depends on various optimization goals: circuits may be sized for minimum delay, minimum energy-delay product, etc. CACTI's goal is to model simple representative circuit sizing applicable to a broad range of common applications. As in earlier SRAM modeling efforts [2, 3, 20], we have made extensive use of the method of logical effort [39] in sizing different circuit blocks. Explanation of the method of logical effort may be found in [39].

4.3 Sizing of Mat Circuits

As described earlier in Section 3.1, a mat is composed of entities such as the predecoding/decoding logic, memory cell array, and bitline peripheral circuitry. We present circuits, models, and sizing techniques for these entities.

4.3.1 Predecoder and Decoder

As discussed in Section 2, new circuit structures have been adopted for the decoding logic. The same decoding logic circuit structures are utilized for producing the row-decode signals and the select signals of the bitline and sense amplifier muxes. In the discussion here, we focus on the row-decoding logic. In order to describe the circuit structures assumed within the different entities of the row-decoding logic, we use an illustrative example. Figure 13 shows the structure of the row-decoding logic for a subarray with 1024 rows. The row-decoding logic is composed of two row-predecode blocks and the row-decode gates and drivers. The row-predecode blocks are responsible for predecoding the address bits and generating predecoded signals. The row-decode gates and drivers are responsible for decoding the predecoded outputs and driving the wordline load. Each row-predecode block can predecode a maximum of 9 bits and has a 2-level logic structure. With 1024 rows, the number of address bits required for row-decoding is 10. Figure 14 shows the structure of each row predecode block for a subarray with 1024 rows. Each row predecode block is responsible for predecoding 5 address bits and each of them generates 32 predecoded output bits. Each predecode block has two levels. The first level is composed of one 2-4 decode unit and one 3-8 decode unit. At the second level, the 4 outputs from the 2-4 decode unit and the 8 outputs from the 3-8 decode unit are combined together using 32 NAND2 gates in order to produce the 32 predecoded outputs. The 32 predecoded outputs from each predecode block are combined together using the 1024 NAND2 gates to generate the row decode signals.

Figure 15 shows the circuit paths in the decoding logic for the subarray with 1024 rows. One of the paths contains the NAND2 of the 2-4 decode unit and the other contains the NAND3 gate of the 3-8 decode unit. Each path has 3 stages in its path. The branching efforts at the outputs of the first two stages are also shown in the figure. The predecode output wire is treated as a non-ideal wire with its $R_{\text{predec-out-wire}}$ and $C_{\text{predec-out-wire}}$ computed using the following equations:

$$R_{\text{predec-output-wire}} = L_{\text{predec-output-wire}} R_{\text{unit-length-wire}} \quad (7)$$

$$C_{\text{predec-output-wire}} = L_{\text{predec-output-wire}} C_{\text{unit-length-wire}} \quad (8)$$

where $L_{\text{predec-output-wire}}$ is the maximum length amongst lengths of predecode output wires.

The sizing of gates in each circuit path is calculated using the method of logical effort. In each of the 3 stages of each circuit path, minimum-size transistors are assumed at the input of the stage and each stage is sized independent of each other using the method of logical effort. While this is not optimal from a delay point of view, it is simpler to model and has been found to be a good sizing heuristic from an energy-delay point of view [3].

In this example that we considered for decoding logic of a subarray with 1024 rows, there were two different circuit paths, one involving the NAND2 gate and another involving the NAND3 gate. In the general case, when each predecode block decodes different number of address bits, a maximum of four circuit paths may exist. When the degree of decoding is low, some of the circuit blocks shown in Figure 13 may not be required. For example, Figure 16 shows the decoding logic for a subarray with 8 rows. In this case, the decoding logic simply involves a 3-8 decode unit as shown.

As mentioned before, the same circuit structures used within the row-decoding logic are also used for generating the select signals of the bitline and sense amplifier muxes. However, unlike the row-decoding logic in which the NAND2 decode gates and drivers are assumed to be placed on the side of subarray, the NAND2 decode gates and drivers are

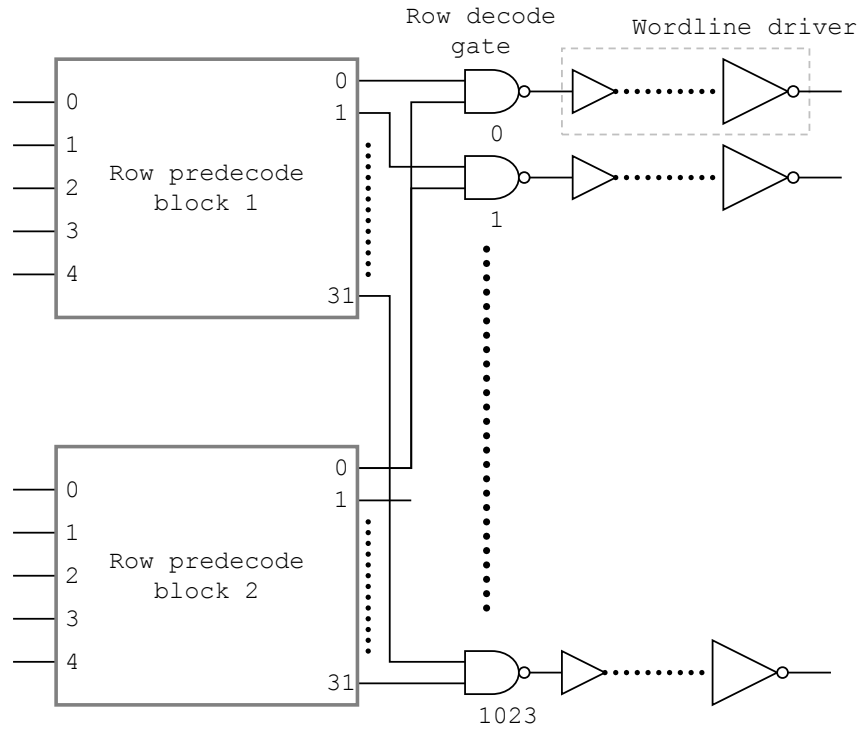


Figure 13: Structure of the row decoding logic for a subarray with 1024 rows.

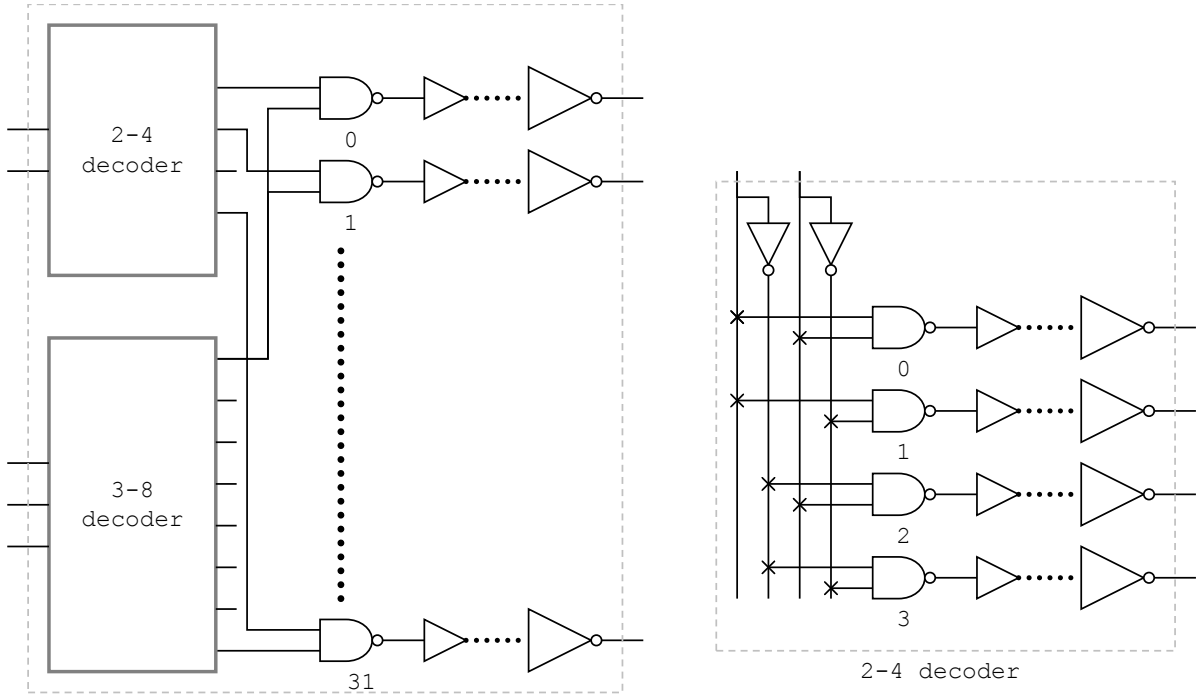


Figure 14: Structure of the row predecode block for a subarray with 1024 rows.

assumed to be placed at the center of the mat near their corresponding predecode blocks. Also, the resistance/capacitance of the wires between the predecode blocks and the decode gates are not modeled and are assumed to be zero.

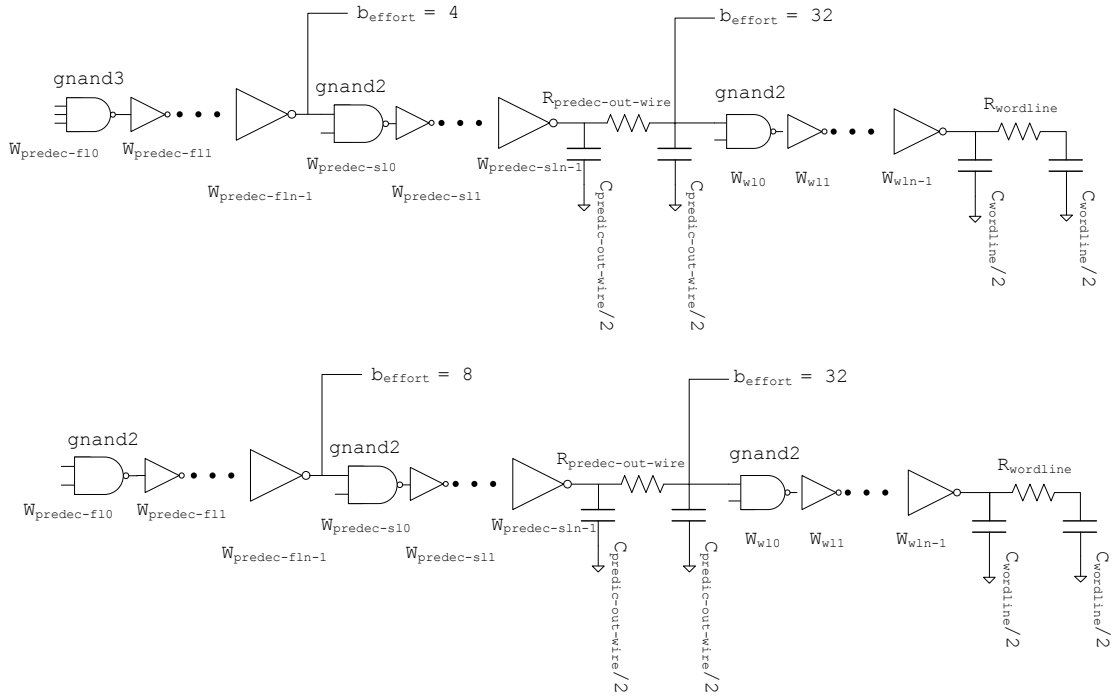


Figure 15: Row decoding logic circuit paths for a subarray with 1024 rows. One of the circuit paths contains the NAND2 gate of the 2-4 decode unit while the other contains the NAND3 gate of the 3-8 decode unit.

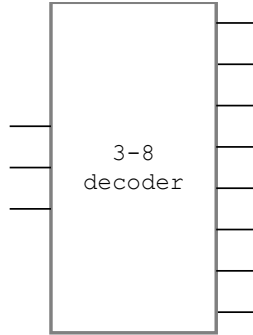


Figure 16: Structure of the row-decoding logic for a subarray with 8 rows. The row-decoding logic is simply composed of 8 decode gates and drivers.

4.3.2 Bitline Peripheral Circuitry

Memory Cell Figure 17 shows the circuit assumed for a 1-ported SRAM cell. The transistors of the SRAM cell are sized based on the widths specified in [14] and are presented in Section 8.

Sense Amplifier Figure 18 shows the circuit assumed for a sense amplifier - it's a clocked latch-based sense amplifier. When the ENABLE signal is not activated, there is no flow of current through the transistors of the latch. When the ENABLE signal is activated the sensing begins. The isolation transistors are responsible for isolating the high capacitance of the bitlines from the sense amplifier nodes during the sensing operation. The small-signal circuit model and analysis of this latch-based sense amplifier is presented in Section 4.4.

Bitline and Sense Amplifier Muxes Figure 19 shows the circuit assumed for the bitline and sense amplifier muxes. We assume that the mux is implemented using NMOS pass transistors. The use of NMOS transistors implies that the

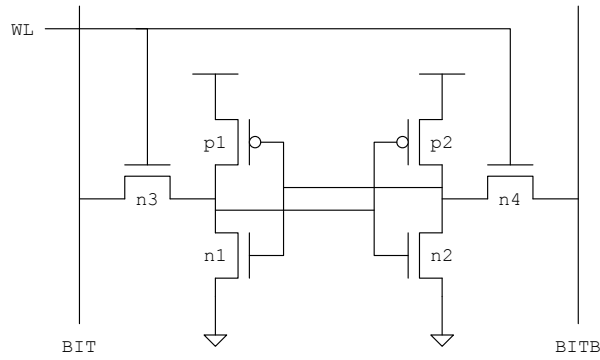


Figure 17: 1-ported 6T SRAM cell

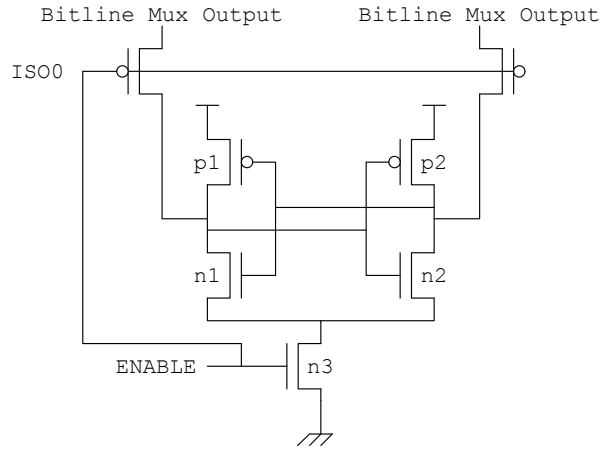


Figure 18: Clocked latch-based sense amplifier

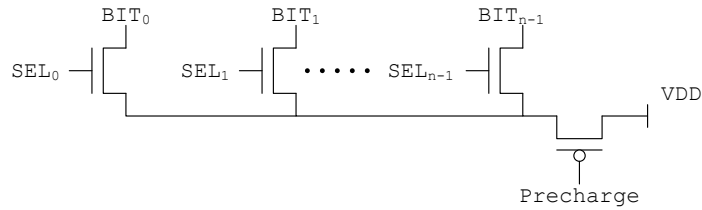


Figure 19: NMOS-based mux. The output is assumed to be precharged high.

output of the mux needs to be precharged high in order to avoid degraded ones. We do not attempt to size the transistors in the muxes and instead assume (as in [2]) fixed widths for the NMOS transistors across all partitions of the array.

Precharge and Equalization Circuitry Figure 20 shows the circuit assumed for precharging and equalizing the bitlines. The bitlines are assumed to be precharged to V_{DD} through the PMOS transistors. Just like the transistors in the bitline and sense amp muxes, we do not attempt to size the precharge and equalization transistors and instead assume fixed-width transistors across different partitions of the array.

Bitlines Read Path Circuit Model Figure 21 shows the circuit model for the bitline read path between the memory cell and the sense amplifier mux.

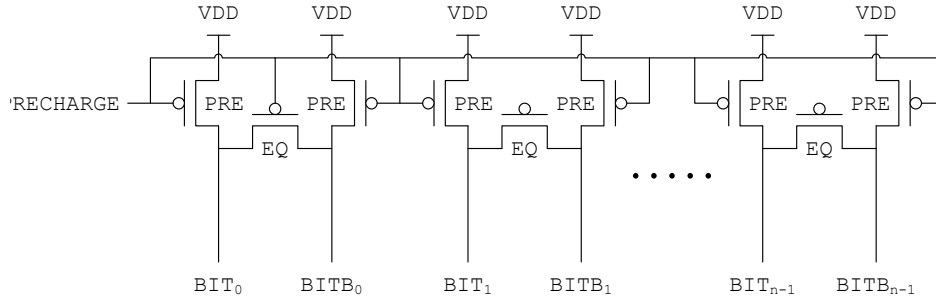


Figure 20: Bitline precharge and equalization circuitry.

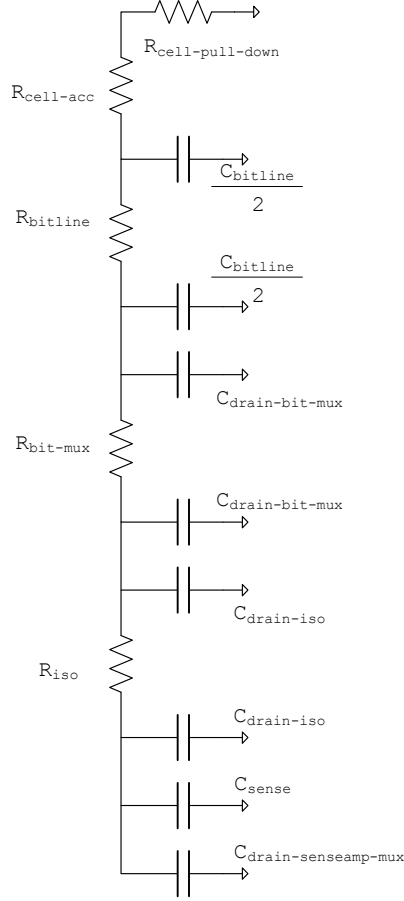


Figure 21: Circuit model of the bitline read path between the SRAM cell and the sense amplifier input.

4.4 Sense Amplifier Circuit Model

Figure 18 showed the clocked latch-based sense amplifier that we have assumed. [10] presents analysis of this circuit and equations for sensing delay under different assumptions. Figure 22 shows one of the small-signal models presented in [10]. Use of this small-signal model is based on two assumptions:

1. Current has been flowing in the circuit for a sufficiently long time; and
2. The equilibrating device can be modeled as an ideal switch.

For the small-signal model of Figure 22, it has been shown that the delay of the sensing operation is given by the following equation:

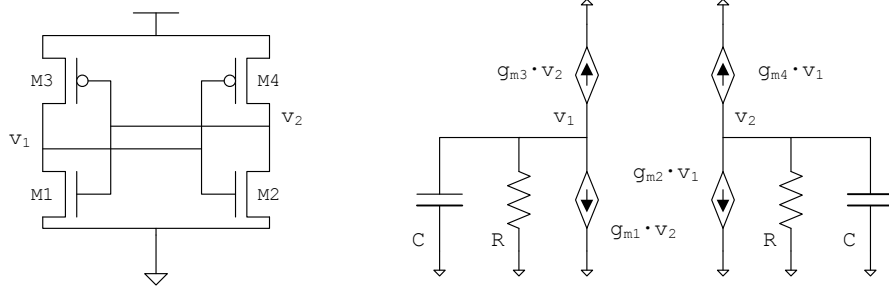


Figure 22: Small-signal model of the latch-based sense amplifier [10].

$$T_{\text{sense}} = \frac{C_{\text{sense}}}{G_m} \ln\left(\frac{V_{DD}}{V_{\text{sense}}}\right) \quad (9)$$

$$G_m = g_{mn} + g_{mp} \quad (10)$$

Use of Equation 9 for calculation of sense amplifier delay requires that the values of g_{mn} (NMOS transconductance) and g_{mp} (PMOS transconductance) be known. We assume that the transistors in the sense amplifier latch exhibit short-channel effects. For a transistor that exhibits short-channel effect, we use the following typical current equation [29] for computation of saturation current:

$$I_{\text{dsat}} = \frac{\mu_{\text{eff}}}{2} C_{\text{ox}} \frac{W}{L} (V_{\text{GS}} - V_{\text{TH}}) V_{\text{dsat}} \quad (11)$$

Differentiating the above equation with respect to V_{GS} gives the equation for g_m of the transistor. It can be seen that because of short-channel effect, g_m comes out to be independent of V_{GS} .

$$g_m = \frac{\mu_{\text{eff}}}{2} C_{\text{ox}} \frac{W}{L} V_{\text{dsat}} \quad (12)$$

4.5 Routing Networks

As described earlier in Section 3.2, address and data are routed to and from the mats on H-tree distribution networks. First address/data are routed on an H-tree from array edge to bank edge and then on another H-tree from bank edge to the mats.

4.5.1 Array Edge to Bank Edge H-tree

Figure 7 showed the layout of H-tree distribution of address and data between the array edge and the banks. This H-tree network is assumed to be composed of inverter-based repeaters. The sizing of the repeaters and the separation distance between them is determined based on the formulae given in [11]. In order to allow for energy-delay tradeoffs in the repeater design, we introduce an user-controlled variable “maximum percentage of delay away from best repeater solution” or `max_repeater_delay_constraint` in short. A `max_repeater_delay_constraint` of zero results in the best delay repeater solution. For a `max_repeater_delay_constraint` of 10%, the delay of the path is allowed to get worse by a maximum of 10% with respect to the best delay repeater solution by reducing the sizing and increasing the separation distance. Thus, with the `max_repeater_delay_constraint`, limited energy savings are possible at the expense of delay.

4.5.2 Bank Edge to Mat H-tree

Figures 8 and 9 showed layout examples of horizontal and vertical H-trees within a bank, each with 3 nodes. We assume that drivers are placed at each of the nodes of these H-trees. Figure 23 shows the circuit path and driver circuit structure

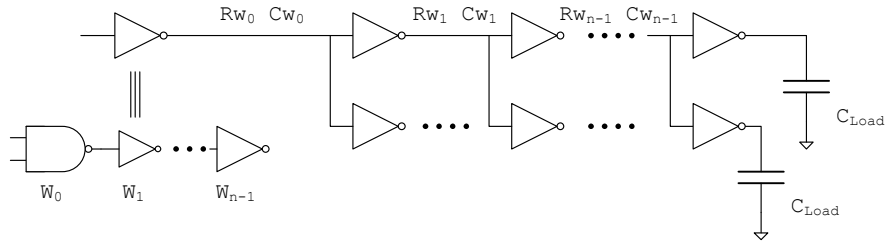


Figure 23: Circuit path of address/datain H-trees within a bank.

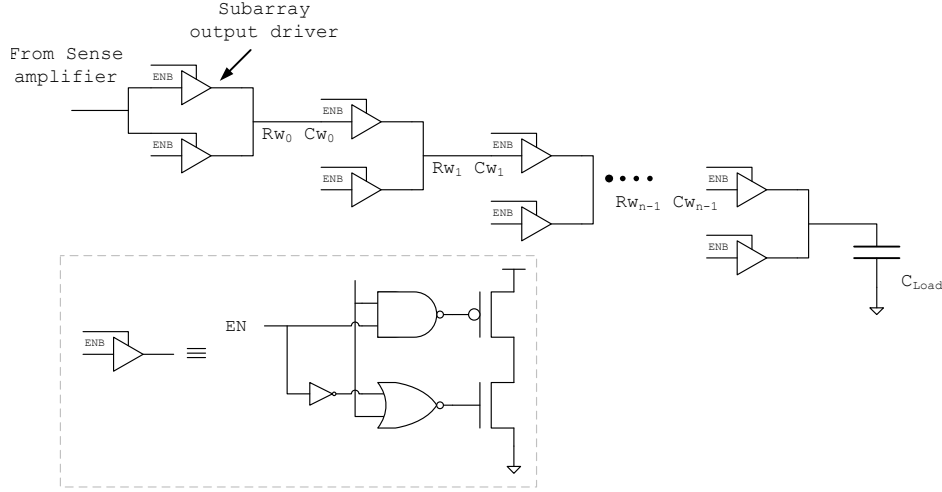


Figure 24: Circuit path of vertical dataout H-trees.

of the address/datain H-trees, and Figure 24 shows the circuit path and driver circuit structure of the vertical dataout H-tree. In order to allow for signal-gating in the address/datain H-trees we consider multi-stage buffers with a 2-input NAND gate as the input stage. The sizing and number of gates at each node of the H-trees is computed using the methodology described in [3] which takes into account the resistance and capacitance of the intermediate wires in the H-tree.

One problem with the circuit paths of Figures 23 and 24 is that they start experiencing increased wire delays as the wire lengths between the drivers start to get long. This also limits the maximum random cycle time that can be achieved for the array. So, as an alternative to modeling drivers only at H-tree branching nodes, we also consider an alternative model in which the H-tree circuit paths within a bank are composed of buffers at regular intervals (i.e. repeaters). With repeaters, the delay through the H-tree paths within a bank can be reduced at the expense of increased power consumption. Figure 25 shows the different types of buffer circuits that have been modeled in the H-tree path. At the branches of the H-tree, we again assume buffers with a NAND gate in the input stage in order to allow for signal-gating whereas in the H-tree segments between two nodes, we model inverter-based buffers. We again size these buffers according to the buffer sizing formulae given in [11]. The `max_repeater_delay_constraint` that was described in Section 4.5.1 is also used here to decide the sizing of the buffers and their separation distance so that delay in these H-trees also may be traded off for potential energy savings.

5 Area Modeling

In this section, we describe the area model of a data array. In Section 5.1, we describe the area model that we have used to find the areas of simple gates. We then present the equations of the area model in Section 5.2.

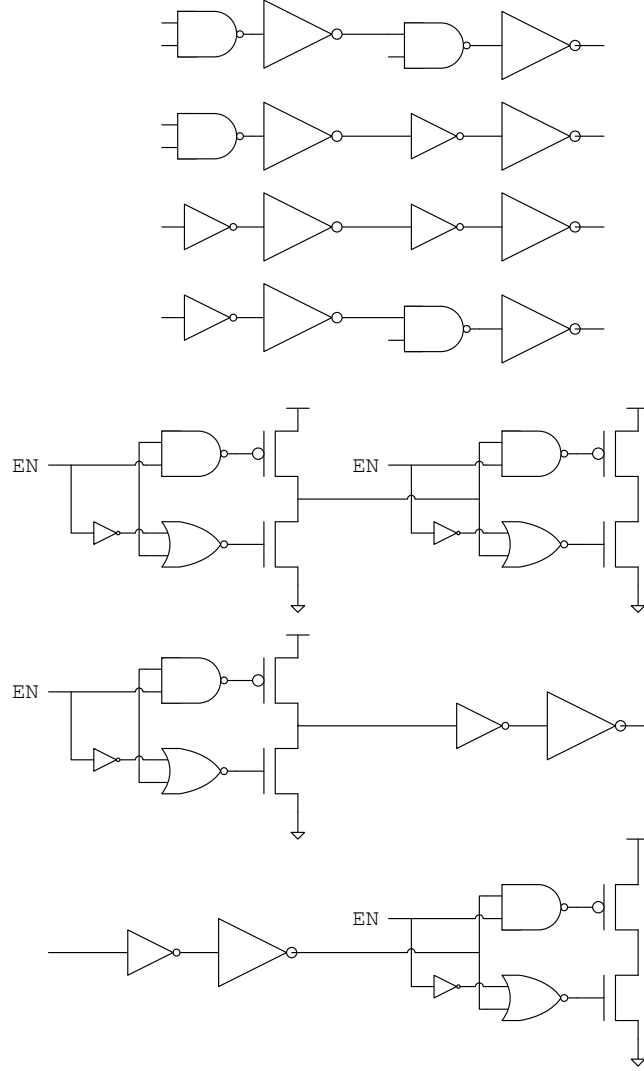


Figure 25: Different types of buffer circuit stages that have been modeled in the H-trees within a bank.

5.1 Gate Area Model

A new area model has been used to estimate the areas of transistors and gates such as inverter, NAND, and NOR gates. This area model is based off a layout model from [49] which describes a fast technique to estimate standard cell characteristics before the cells are actually laid out. Figure 26 illustrates the layout model that has been used in [49]. Table 2 shows the process/technology input parameters required by this gate area model. For a thorough description of the technique, please refer to [49]. Gates with stacked transistors are assumed to have a layout similar to that described in [47]. When a transistor width exceeds a certain maximum value ($H_{n\text{-diff}}$ for NMOS and $H_{p\text{-diff}}$ for PMOS in Table 2), the transistor is assumed to be folded. This maximum value can either be process-specific or context-specific. An example of when a context-specific width would be used is in case of memory sense amplifiers which typically have to be laid out at a certain pitch.

Given the width of an NMOS transistor, $W_{\text{before-folding}}$, the number of folded transistors may be calculated as follows:

$$N_{\text{folded-transistors}} = \left\lceil \frac{W_{\text{before-folding}}}{H_{n\text{-diff}}} \right\rceil \quad (13)$$

The equation for total diffusion width of N_{stacked} transistors when they are not folded is given by the following equation:

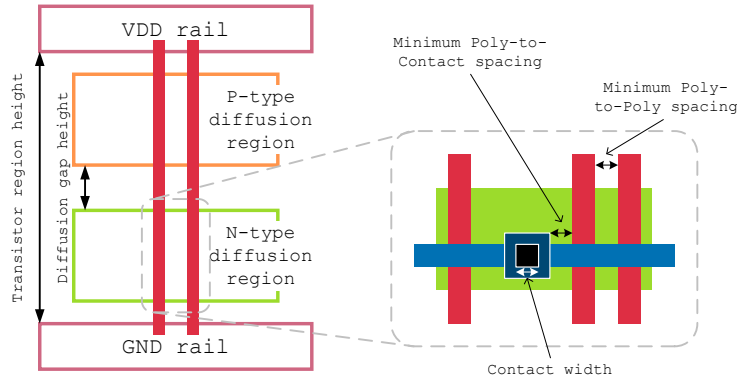


Figure 26: Layout model assumed for gates [49].

Parameter name	Meaning
$H_{n\text{-diff}}$	Maximum height of n diffusion of a transistor
$H_{p\text{-diff}}$	Maximum height of p diffusion for a transistor
$H_{\text{gap-bet-same-diffs}}$	Minimum gap between diffusions of the same type
$H_{\text{gap-bet-opp-diffs}}$	Minimum gap between n and p diffusions
$H_{\text{power-rail}}$	Height of V_{DD} (GND) power rail
W_p	Minimum width of poly (poly half-pitch or process feature size)
S_{p-p}	Minimum poly-to-poly spacing
W_c	Contact width
S_{p-c}	Minimum poly-to-contact spacing

Table 2: Process/technology input parameters required by the gate area model.

$$\text{total-diff-width} = 2(W_c + 2S_{p-c}) + N_{\text{stacked}}W_p + (N_{\text{stacked}} - 1)S_{p-p} \quad (14)$$

The equation for total diffusion width of N_{stacked} transistors when they are folded is given by the following equation:

$$\text{total-diff-width} = N_{\text{folded-transistors}}(2(W_c + 2S_{p-c}) + N_{\text{stacked}}W_p + (N_{\text{stacked}} - 1)S_{p-p}) \quad (15)$$

Note that Equation 15 is a generalized form of the equations used for calculating diffusion width (for computation of drain capacitance) in the original CACTI report [47]. Earlier versions of CACTI assumed at most two folded transistors; in version 5, we allow the degree of folding to be greater than 2 and make the associated layout and area models more general. Note that drain capacitance calculation in version 5 makes use of equations similar to 14 and 15 for computation of diffusion width.

The height of a gate is calculated using the following equation:

$$H_{\text{gate}} = H_{n\text{-diff}} + H_{p\text{-diff}} + H_{\text{gap-bet-opp-diffs}} + 2H_{\text{power-rail}} \quad (16)$$

5.2 Area Model Equations

The area of the data array is estimated based on the area occupied by a single bank and the area spent in routing address and data to the banks. It is assumed that the area spent in routing address and data to the bank is decided by the pitch of the routed wires. Figures 27 and 28 show two example arrays with 8 and 16 banks respectively; we present equations for the calculation of the areas of these arrays.

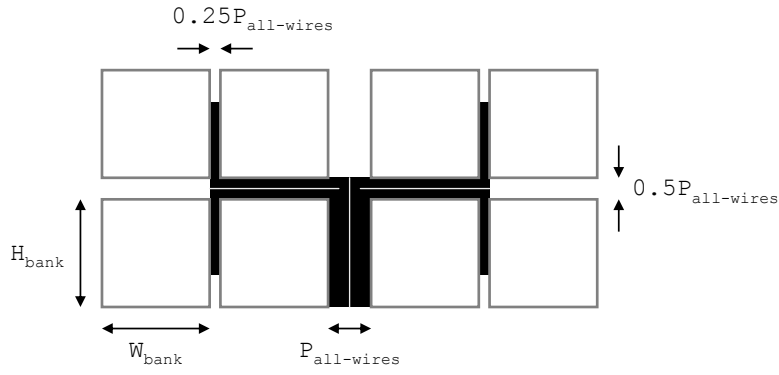


Figure 27: Supporting figure for example area calculation of array with 8 banks.

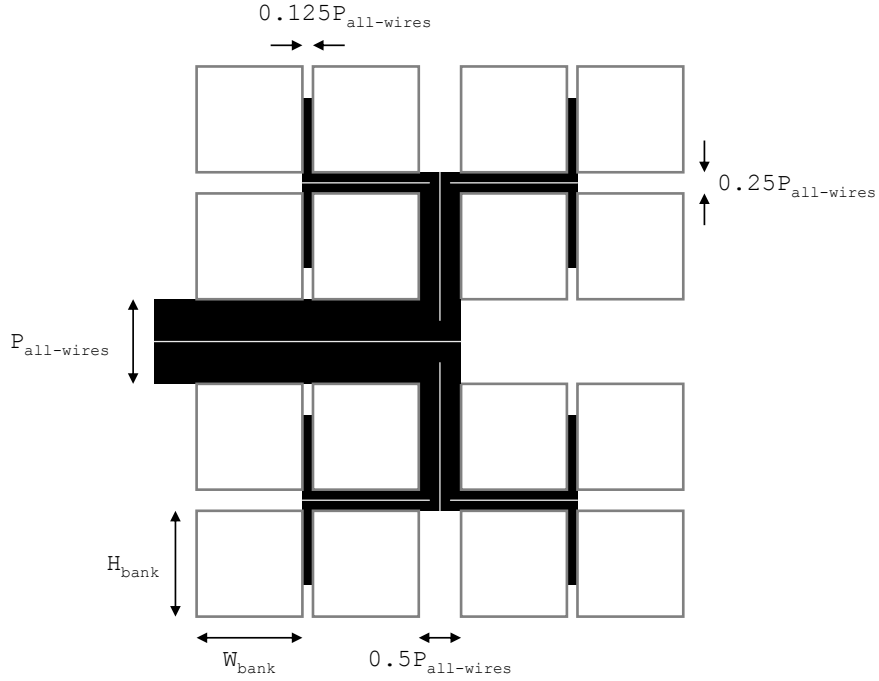


Figure 28: Supporting figure for example area calculation of array with 16 banks.

$$A_{\text{data-arr}} = H_{\text{data-arr}} W_{\text{data-arr}} \quad (17)$$

The pitch of wires routed to the banks is given by the following equation:

$$P_{\text{all-wires}} = P_{\text{wire}} N_{\text{wires-routed-to-banks}} \quad (18)$$

For the data array of Figure 27 with 8 banks, the relevant equations are as follows:

$$W_{\text{data-arr}} = 4W_{\text{bank}} + P_{\text{all-wires}} + 2 \frac{P_{\text{all-wires}}}{4} \quad (19)$$

$$H_{\text{data-arr}} = 2H_{\text{bank}} + \frac{P_{\text{all-wires}}}{2} \quad (20)$$

$$N_{\text{wires-routed-to-banks}} = 8(N_{\text{bank-addr-bits}} + N_{\text{bank-datain-bits}} + N_{\text{bank-dataout-bits}} + N_{\text{way-select-signals}}) \quad (21)$$

For the data array of Figure 28 with 16 banks, the relevant equations are as follows:

$$W_{\text{data-arr}} = 4W_{\text{bank}} + \frac{P_{\text{all-wires}}}{2} + 2\frac{P_{\text{all-wires}}}{8} \quad (22)$$

$$H_{\text{data-arr}} = 4H_{\text{bank}} + P_{\text{all-wires}} + 2\frac{P_{\text{all-wires}}}{4} \quad (23)$$

$$N_{\text{wires-routed-to-banks}} = 16(N_{\text{bank-addr-bits}} + N_{\text{bank-datain-bits}} + N_{\text{bank-dataout-bits}} + N_{\text{way-select-signals}}) \quad (24)$$

The banks in a data array are assumed to be placed in such a way that the number of banks in the horizontal direction is always either equal to or twice the number of banks in the vertical direction. The height and width of a bank is calculated by computing the area occupied by the mats and the area occupied by the routing resources of the horizontal and vertical H-tree networks within a bank. We again use an example to illustrate the calculations. Figures 8 and 9 showed the layouts of horizontal and vertical H-trees within a bank. The horizontal and vertical H-trees were each shown to have three branching nodes (H0, H1, and H2; V0, V1, and V2). Combined together, these horizontal and vertical H-trees may be considered as H-trees within a bank with 4 subbanks and 4 mats in each subbank. We present area model equations for such a bank.

$$A_{\text{bank}} = H_{\text{bank}}W_{\text{bank}} \quad (25)$$

In version 5, as described in Section 4.5, for the H-trees within a bank we assume that drivers are placed either only at the branching nodes of the H-trees or that there are buffers at regular intervals in the H-tree segments. When drivers are present only at the branching nodes of the vertical H-trees within a bank, we consider two alternative models in accounting for area overhead of the vertical H-trees. In the first model, we consider that wires of the vertical H-trees may traverse over memory cell area; in this case, the area overhead caused by the vertical H-trees is in terms of area occupied by drivers which are placed between the mats. In the second model, we do not assume that the wires traverse over the memory cell area and instead assume that they occupy area besides the mats. The second model is also applicable when there are buffers at regular intervals in the H-tree segments. The equations that we present next for area calculation of a bank assume the second model i.e. the wires of the vertical H-trees are assumed to not pass over the memory cell area. The equations for area calculation under the assumption that the vertical H-tree wires go over the memory cell area are quite similar. For our example bank with 4 subbanks and 4 mats in each subbank, the height of the bank is calculated to be equal to the sum of heights of all subbanks plus the height of the routing resources of the horizontal H-tree.

$$H_{\text{bank}} = 4H_{\text{mat}} + H_{\text{hor-htree}} \quad (26)$$

The width of the bank is calculated to be equal to the sum of widths of all mats in a subbank plus the width of the routing resources of the vertical H-trees.

$$W_{\text{bank}} = 4(W_{\text{mat}} + W_{\text{ver-htree}}) \quad (27)$$

The height of the horizontal H-tree is calculated as the height of the area occupied by the wires in the H-tree. These wires include the address, way-select, datain, and dataout signals. Figure 29 illustrates the layout that we assume for the wires of the horizontal H-tree. We assume that the wires are laid out using a single layer of metal. The height of the area occupied by the wires can be calculated simply by finding the total pitch of all wires in the horizontal H-tree. Figure 30 illustrates the layout style assumed for the vertical H-tree wires, and is similar to that assumed for the horizontal H-tree wires. Again the width of the area occupied by a vertical H-tree can be calculated by finding the total pitch of all wires in the vertical H-tree.

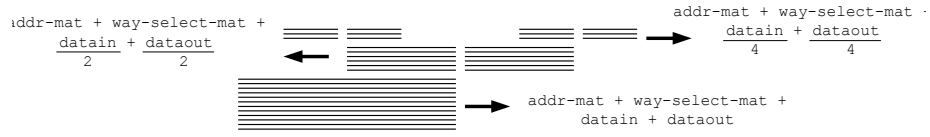


Figure 29: Layout assumed for wires of the horizontal H-tree within a bank.

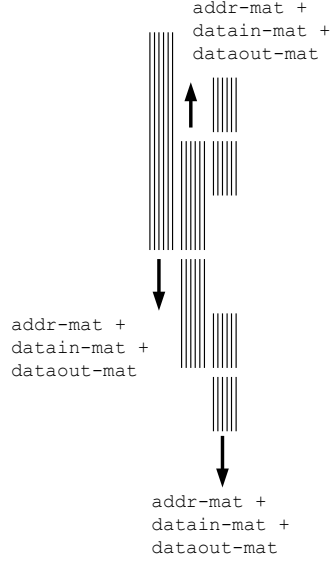


Figure 30: Layout assumed for wires of the vertical H-tree within a bank.

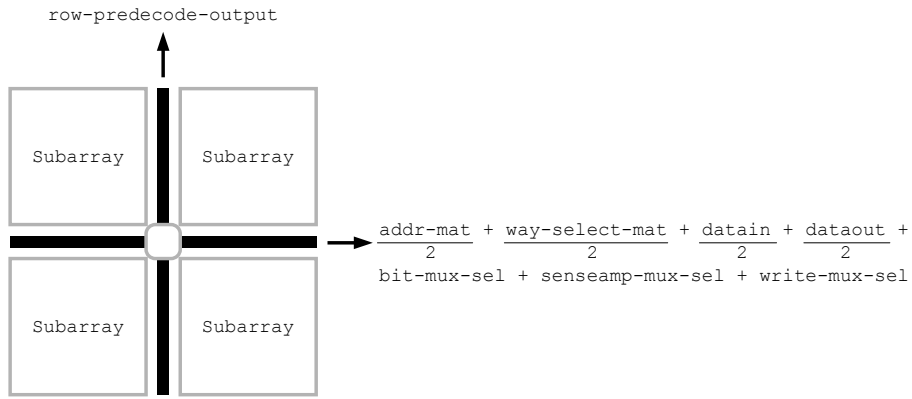


Figure 31: Layout of a mat.

$$H_{\text{hor-htree}} = P_{\text{hor-htree-wires}} \quad (28)$$

$$W_{\text{ver-htree}} = P_{\text{ver-htree-wires}} \quad (29)$$

The height and width of a mat are estimated using the following equations. Figure 31 shows the layout of a mat and illustrates the assumptions made in the following equations. We assume that half of the address, way-select, datain, and dataout signals enter the mat from its left and the other half enter from the right.

$$W_{\text{mat}} = \frac{H_{\text{mat}} W_{\text{initial-mat}} + A_{\text{mat-center-circuitry}}}{W_{\text{initial-mat}}} \quad (30)$$

$$H_{\text{mat}} = 2H_{\text{subarr-mem-cell-area}} + H_{\text{mat-non-cell-area}} \quad (31)$$

$$W_{\text{initial-mat}} = 2W_{\text{subarr-mem-cell-area}} + W_{\text{mat-non-cell-area}} \quad (32)$$

$$\begin{aligned} A_{\text{mat-center-circuitry}} = & A_{\text{row-predec-block-1}} + A_{\text{row-predec-block-2}} \\ & + A_{\text{bit-mux-predec-block-1}} + A_{\text{bit-mux-predec-block-2}} \\ & + A_{\text{senseamp-mux-predec-block-1}} + A_{\text{senseamp-mux-predec-block-2}} + \\ & A_{\text{bit-mux-dec-drivers}} + A_{\text{senseamp-mux-dec-drivers}} \end{aligned} \quad (33)$$

$$H_{\text{subarr-mem-cell-area}} = N_{\text{subarr-rows}} H_{\text{mem-cell}} \quad (34)$$

$$\begin{aligned} W_{\text{subarr-mem-cell-area}} = & N_{\text{subarr-cols}} W_{\text{mem-cell}} + \lfloor \frac{N_{\text{subarr-cols}}}{N_{\text{mem-cells-per-wordline-stitch}}} \rfloor W_{\text{wordline-stitch}} + \\ & \lceil \frac{N_{\text{subarr-cols}}}{N_{\text{bits-per-ecc-bit}}} \rceil W_{\text{mem-cell}} \end{aligned} \quad (35)$$

$$H_{\text{mat-non-cell-area}} = 2H_{\text{subarr-bitline-peri-circ}} + H_{\text{hor-wires-within-mat}} \quad (36)$$

$$\begin{aligned} H_{\text{hor-wires-within-mat}} = & H_{\text{bit-mux-sel-wires}} + H_{\text{senseamp-mux-sel-wires}} + H_{\text{write-mux-sel-wires}} + \\ & \frac{H_{\text{number-mat-addr-bits}}}{2} + \frac{H_{\text{number-way-select-signals}}}{2} + \\ & \frac{H_{\text{number-mat-datain-bits}}}{2} + \frac{H_{\text{number-mat-dataout-bits}}}{2} \end{aligned} \quad (37)$$

$$W_{\text{mat-non-cell-area}} = \max(2W_{\text{subarr-row-decoder}}, W_{\text{row-predec-out-wires}}) \quad (38)$$

$$H_{\text{subarr-bitline-peri-cir}} = H_{\text{bit-mux}} + H_{\text{senseamp-mux}} + H_{\text{bitline-pre-eq}} + H_{\text{write-driver}} + H_{\text{write-mux}} \quad (39)$$

Note that the width of the mat is computed as in Equation 31 because we optimistically assume that the circuitry laid out at the center of the mat does not lead to white space in the mat. The areas of lower-level circuit blocks such as the bitline and sense amplifier muxes and write drivers are calculated using the area model that was described in Section 5.1 while taking into account pitch-matching constraints.

When redundancy in mats is also considered, the following area contribution due to redundant mats is added to the area of the data array computed in Equation 17.

$$A_{\text{redundant-mats}} = N_{\text{redundant-mats}} A_{\text{mat}} \quad (40)$$

$$N_{\text{redundant-mats}} = \lfloor \frac{N_{\text{banks}}}{N_{\text{mats}}} N_{\text{mats-per-redundant-mat}} \rfloor \quad (41)$$

where $N_{\text{mats-per-redundant-mat}}$ is the number of mats per redundant mat that and is set to 8 by default. The final height of the data array is readjusted under the optimistic assumption that the redundant mats do not cause any white space in the data array.

$$H_{\text{data-arr}} = \frac{A_{\text{data-arr}}}{W_{\text{data-arr}}} \quad (42)$$

6 Delay Modeling

In this section we present equations used in CACTI to calculate access time and random cycle time of a memory array.

6.1 Access Time Equations

$$T_{\text{access}} = T_{\text{request-network}} + T_{\text{mat}} + T_{\text{reply-network}} \quad (43)$$

$$T_{\text{request-network}} = T_{\text{arr-edge-to-bank-edge-htree}} + T_{\text{bank-addr-din-hor-htree}} + T_{\text{bank-addr-din-ver-htree}} \quad (44)$$

$$T_{\text{mat}} = \max(T_{\text{row-decoder-path}}, T_{\text{bit-mux-decoder-path}}, T_{\text{senseamp-decoder-path}}) \quad (45)$$

$$T_{\text{reply-network}} = T_{\text{bank-dout-ver-htree}} + T_{\text{bank-dout-hor-htree}} + T_{\text{bank-edge-to-arr-edge}} \quad (46)$$

The critical path in the mat usually involves the wordline and bitline access. However, Equation 45 also must include a max with the delays of the bitline mux decoder and sense amp mux decoder paths as these circuits operate in parallel with the row decoding logic, and in general may act as the critical path for certain partitions of the data array. Usually when that happens, the number of rows in the subarray would be too few and the partitions would not get selected.

$$T_{\text{row-decoder-path}} = T_{\text{row-predec}} + T_{\text{row-dec-driver}} + T_{\text{bitline}} + T_{\text{senseamp}} \quad (47)$$

$$T_{\text{bit-mux-decoder-path}} = T_{\text{bit-mux-predec}} + T_{\text{bit-mux-dec-driver}} + T_{\text{senseamp}} \quad (48)$$

$$T_{\text{senseamp-mux-decoder-path}} = T_{\text{senseamp-mux-predec}} + T_{\text{senseamp-mux-dec-driver}} \quad (49)$$

$$T_{\text{row-predec}} = \max(T_{\text{row-predec-blk-1-nand2-path}}, T_{\text{row-predec-blk-1-nand3-path}}, T_{\text{row-predec-blk-2-nand2-path}}, T_{\text{row-predec-blk-2-nand3-path}}) \quad (50)$$

$$T_{\text{bit-mux-sel-predec}} = \max(T_{\text{bit-mux-sel-predec-blk-1-nand2-path}}, T_{\text{bit-mux-sel-predec-blk-1-nand3-path}}, T_{\text{bit-mux-sel-predec-blk-2-nand2-path}}, T_{\text{bit-mux-sel-predec-blk-2-nand3-path}}) \quad (51)$$

$$T_{\text{senseamp-mux-sel-predec}} = \max(T_{\text{senseamp-mux-sel-predec-blk-1-nand2-path}}, T_{\text{senseamp-mux-sel-predec-blk-1-nand3-path}}, T_{\text{senseamp-mux-sel-predec-blk-2-nand2-path}}, T_{\text{senseamp-mux-sel-predec-blk-2-nand3-path}}) \quad (52)$$

The calculation for bitline delay is based on the model described in [46]. The model considers the effect of the wordline rise time by considering the slope (m) of the wordline signal.

$$T_{\text{bitline}} = \begin{cases} \sqrt{2T_{\text{step}} \frac{V_{\text{DD}} - V_{\text{th}}}{m}} & \text{if } T_{\text{step}} \leq 0.5 \frac{V_{\text{DD}} - V_{\text{th}}}{m} \\ T_{\text{step}} + \frac{V_{\text{DD}} - V_{\text{th}}}{2m} & \text{if } T_{\text{step}} > 0.5 \frac{V_{\text{DD}} - V_{\text{th}}}{m} \end{cases} \quad (53)$$

$$\begin{aligned} T_{\text{step}} = & (R_{\text{cell-pull-down}} + R_{\text{cell-acc}})(C_{\text{bitline}} + 2C_{\text{drain-bit-mux}} + C_{\text{iso}} + C_{\text{senseamp}} + C_{\text{drain-senseamp-mux}}) + \\ & R_{\text{bitline}} \left(\frac{C_{\text{bitline}}}{2} + 2C_{\text{drain-bit-mux}} + C_{\text{iso}} + C_{\text{senseamp}} + C_{\text{drain-senseamp-mux}} \right) + \\ & R_{\text{bit-mux}} (C_{\text{drain-bit-mux}} + C_{\text{iso}} + C_{\text{senseamp}} + C_{\text{drain-senseamp-mux}}) + R_{\text{iso}} (C_{\text{iso}} + C_{\text{senseamp}} + \\ & C_{\text{drain-senseamp-mux}}) \end{aligned} \quad (54)$$

The calculation of sense amplifier delay makes use of the model described in [10].

$$T_{\text{senseamp}} = \tau \ln \left(\frac{V_{\text{DD}}}{V_{\text{senseamp}}} \right) \quad (55)$$

$$\tau = \frac{C_{\text{senseamp}}}{G_{\text{m}}} \quad (56)$$

6.2 Random Cycle Time Equations

Typically, the random cycle time of an SRAM would be limited by wordline and bitline delays. In order to come up with an equation for lower bound on random cycle time, we consider that the SRAM is potentially pipelineable with placement of latches at appropriate locations.

$$\begin{aligned}
T_{\text{random-cycle}} = & \max(T_{\text{row-dec-driver}} + T_{\text{bitline}} + T_{\text{senseamp}} + T_{\text{wordline-reset}} + \\
& \max(T_{\text{bitline-precharge}}, T_{\text{bit-mux-out-precharge}}, T_{\text{senseamp-mux-out-precharge}}), \\
& T_{\text{between-buffers-bank-hor-htree}}, T_{\text{between-buffers-bank-ver-dataout-htree}}, T_{\text{row-predec-blk}}, \\
& T_{\text{bit-mux-predec}} + T_{\text{bit-mux-dec-driver}}, \\
& T_{\text{senseamp-mux-predec}} + T_{\text{senseamp-mux-dec-driver}})
\end{aligned} \tag{57}$$

We come up with an estimate for the wordline reset delay by assuming that the wordline discharges through the NMOS transistor of the final inverter in the wordline driver.

$$T_{\text{wordline-reset}} = \ln\left(\frac{V_{\text{DD}} - 0.1V_{\text{DD}}}{V_{\text{DD}}}\right)(R_{\text{final-inv-wordline-driver}}C_{\text{wordline}} + \frac{R_{\text{final-inv-wordline-driver}}C_{\text{wordline}}}{2}) \tag{58}$$

$$\begin{aligned}
T_{\text{bitline-precharge}} &= \ln\left(\frac{V_{\text{DD}} - 0.1V_{\text{bitline-swing}}}{V_{\text{DD}} - V_{\text{bitline-swing}}}\right)(R_{\text{bit-pre}}C_{\text{bitline}} + \frac{R_{\text{bitline}}C_{\text{bitline}}}{2}) \\
T_{\text{bit-mux-out-precharge}} &= \ln\left(\frac{V_{\text{DD}} - 0.1V_{\text{bitline-swing}}}{V_{\text{DD}} - V_{\text{bitline-swing}}}\right)(R_{\text{bit-mux-pre}}C_{\text{bit-mux-out}} + \frac{R_{\text{bit-mux-out}}C_{\text{bit-mux-out}}}{2})
\end{aligned} \tag{59}$$

$$\begin{aligned}
T_{\text{senseamp-mux-out-precharge}} &= \ln\left(\frac{V_{\text{DD}} - 0.1V_{\text{bitline-swing}}}{V_{\text{DD}} - V_{\text{bitline-swing}}}\right)(R_{\text{senseamp-mux-pre}}C_{\text{senseamp-mux-out}} + \frac{R_{\text{senseamp-mux-out}}C_{\text{senseamp-mux-out}}}{2})
\end{aligned} \tag{60}$$

7 Power Modeling

In this section, we present the equations used in CACTI to calculate dynamic energy and leakage power of a data array. We present equations for dynamic read energy; the equations for dynamic write energy are similar.

7.1 Calculation of Dynamic Energy

7.1.1 Dynamic Energy Calculation Example for a CMOS Gate Stage

We present a representative example to illustrate how we calculate the dynamic energy for a CMOS gate stage. Figure 32 shows a CMOS gate stage composed of a NAND2 gate followed by an inverter which drives the load. The energy consumption of this circuit is given by:

$$E_{\text{dyn}} = E_{\text{dyn-nand2}} + E_{\text{dyn-inv}} \tag{61}$$

$$E_{\text{dyn-nand2}} = 0.5(C_{\text{intrinsic-nand2}} + C_{\text{gate-inv}})V_{\text{DD}}^2 \tag{62}$$

$$E_{\text{dyn-inv}} = 0.5(C_{\text{intrinsic-inv}} + C_{\text{gate-load-next-stage}} + C_{\text{wire-load}})V_{\text{DD}}^2 \tag{63}$$

$$C_{\text{intrinsic-nand2}} = \text{draincap}(\text{nand2}, W_{\text{nand-pmos}}, W_{\text{nand-nmos}}) \tag{64}$$

$$C_{\text{gate-inv}} = \text{gatecap}(\text{inv}, W_{\text{inv-pmos}}, W_{\text{inv-nmos}}) \tag{65}$$

$$C_{\text{drain-inv}} = \text{draincap}(\text{inv}, W_{\text{inv-pmos}}, W_{\text{inv-nmos}}) \tag{66}$$

The multiplicative factor of 0.5 in the equations of $E_{\text{dyn-nand2}}$ and $E_{\text{dyn-inv}}$ assumes consecutive charging and discharging cycles for each gate. Energy is consumed only during the charging cycle of a gate when its output goes from low to high.

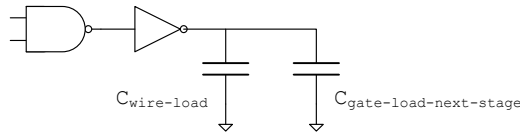


Figure 32: A simple CMOS gate stage composed of a NAND2 followed by an inverter which is driving a load.

7.1.2 Dynamic Energy Equations

The dynamic energy per read access consumed in the data array is the sum of the dynamic energy consumed in the mats and that consumed in the request and reply networks during a read access.

$$E_{\text{dyn-read}} = E_{\text{dyn-read-request-network}} + E_{\text{dyn-read-mats}} + E_{\text{dyn-read-reply-network}} \quad (67)$$

$$E_{\text{dyn-read-mats}} = (E_{\text{dyn-predec-blks}} + E_{\text{dyn-decoder-drivers}} + E_{\text{dyn-read-bitlines}} + E_{\text{senseamps}}) N_{\text{mats-in-subbank}} \quad (68)$$

$$E_{\text{dyn-predec-blks}} = E_{\text{dyn-row-predec-blks}} + E_{\text{dyn-bit-mux-predec-blks}} + E_{\text{dyn-senseamp-mux-predec-blks}} \quad (69)$$

$$E_{\text{dyn-row-predec-blks}} = E_{\text{dyn-row-predec-blk-1-nand2-path}} + E_{\text{dyn-row-predec-blk-1-nand3-path}} + E_{\text{dyn-row-predec-blk-2-nand2-path}} + E_{\text{dyn-row-predec-blk-2-nand3-path}} \quad (70)$$

$$E_{\text{dyn-bit-mux-predec-blks}} = E_{\text{dyn-bit-mux-predec-blk-1-nand2-path}} + E_{\text{dyn-bit-mux-predec-blk-1-nand3-path}} + E_{\text{dyn-bit-mux-predec-blk-2-nand2-path}} + E_{\text{dyn-bit-mux-predec-blk-2-nand3-path}} \quad (71)$$

$$E_{\text{dyn-senseamp-mux-predec-blks}} = E_{\text{dyn-senseamp-mux-predec-blk-1-nand2-path}} + E_{\text{dyn-senseamp-mux-predec-blk-1-nand3-path}} + E_{\text{dyn-senseamp-mux-predec-blk-2-nand2-path}} + E_{\text{dyn-senseamp-mux-predec-blk-2-nand3-path}} \quad (72)$$

$$E_{\text{dyn-decoder-drivers}} = E_{\text{dyn-row-decoder-drivers}} + E_{\text{dyn-bitmux-decoder-driver}} + E_{\text{dyn-senseamp-mux-decoder-driver}} \quad (73)$$

$$E_{\text{dyn-row-decoder-drivers}} = 4E_{\text{dyn-mat-row-decoder-driver}} \quad (74)$$

$$E_{\text{dyn-read-bitlines}} = N_{\text{subarr-cols}} E_{\text{dyn-read-bitline}} \quad (75)$$

$$E_{\text{dyn-read-bitline}} = C_{\text{bitline}} V_{\text{bitline-swing}} V_{\text{DD}} \quad (76)$$

$$V_{\text{bitline-swing}} = 2V_{\text{senseamp}} \quad (77)$$

$$E_{\text{dyn-read-request-network}} = E_{\text{dyn-read-arr-edge-to-bank-edge-request-htree}} + E_{\text{dyn-read-bank-hor-request-htree}} + E_{\text{dyn-read-bank-ver-request-htree}} \quad (78)$$

$$E_{\text{dyn-read-reply-network}} = E_{\text{dyn-read-bank-ver-reply-htree}} + E_{\text{dyn-read-bank-hor-reply-htree}} + E_{\text{dyn-read-bank-edge-to-arr-edge-reply-htree}} \quad (79)$$

Equation 78 assumes that the swing in the bitlines rises up to twice the signal that can be detected by the sense amplifier [2]. $E_{\text{dyn-read-request-network}}$ and $E_{\text{dyn-read-reply-network}}$ are calculated by determining the energy consumed in the wires/drivers/repeaters of the H-trees. The energy consumption in the horizontal and vertical H-trees of the request network within a bank for the example 1MB bank discussed in Section 4.5 with 4 subbanks and 4 mats in each subbank may be written as follows (referring to Figures 8 and 9 in Section 3.2):

$$E_{\text{dyn-read-bank-hor-request-htree}} = E_{\text{dyn-read-req-network-H0-H1}} + E_{\text{dyn-read-req-network-H1-H2}} + E_{\text{dyn-read-req-network-read-H2-V0}} \quad (80)$$

$$E_{\text{dyn-read-bank-ver-request-htree}} = E_{\text{dyn-read-req-network-V0-V1}} + E_{\text{dyn-read-req-network-V1-V2}} \quad (81)$$

The energy consumed in the H-tree segments depends on the location of the segment in the H-tree and the number of signals that are transmitted in each segment. In the request network, during a read access, between nodes H0 and H1, a total of 15 (address) signals are transmitted; between node H1 and both H2 nodes, a total of 30 (address) signals are transmitted; between all H2 and V0 nodes, a total of 60 (address) signals are transmitted. In the vertical H-tree, we assume signal-gating so that the address bits are transmitted to the mats of a single subbank only; thus, between all V0 and V1 nodes, a total of 56 (address) signals are transmitted; between all V1 and V2 nodes, a total of 52 (address) signals are transmitted.

$$E_{\text{dyn-read-req-network-H0-H1}} = (15)E_{\text{H0-H1-1-bit}} \quad (82)$$

$$E_{\text{dyn-read-req-network-H1-H2}} = (30)E_{\text{H1-H2-1-bit}} \quad (83)$$

$$E_{\text{dyn-read-req-network-H2-V0}} = (60)E_{\text{H2-V0-1-bit}} \quad (84)$$

$$E_{\text{dyn-read-req-network-V0-V1}} = (56)E_{\text{V0-V1-1-bit}} \quad (85)$$

$$E_{\text{dyn-read-req-network-V1-V2}} = (52)E_{\text{V1-V2-1-bit}} \quad (86)$$

The equations for energy consumed in the H-trees of the reply network are similar in form to the above equations. Also, the equations for dynamic energy per write access are similar to the ones that have been presented here for read access. In case of write access, the data bits are written into the memory cells at full swing of the bitlines.

7.2 Calculation of Leakage Power

We estimate the standby leakage power consumed in the array. Our leakage power estimation does not consider the use of any leakage control mechanism in the array. We make use of the methodology presented in [20, 21] to simply provide an estimate of the drain-to-source subthreshold leakage current for all transistors that are off with V_{DD} applied across their drain and source.

7.2.1 Leakage Power Calculation for CMOS gates

We illustrate our methodology of calculation of leakage power for the CMOS gates that are used in our modeling. Figure 33 illustrates the leakage power calculation for an inverter. When the input is low and the output is high, there is subthreshold leakage through the NMOS transistor whereas when the input is high and the output is low, there is subthreshold leakage current through the PMOS transistor. In order to simplify our modeling, we come up with a single average leakage power number for each gate. Thus for the inverter, we calculate leakage as follows:

$$P_{\text{leak-inv}} = \frac{W_{\text{inv-pmos}}I_{\text{off-pmos}} + W_{\text{inv-nmos}}I_{\text{off-nmos}}}{2} V_{DD} \quad (87)$$

where $I_{\text{off-pmos}}$ is the subthreshold current per unit width for the PMOS transistor and $I_{\text{off-nmos}}$ is the subthreshold current per unit width for the NMOS transistor.

Figure 34 illustrates the leakage power calculation for a NAND2 gate. When both inputs are high, the output is low and for this condition there is leakage through the PMOS transistors as shown. When either of the inputs is low, the output is high and there is leakage through the NMOS transistors. Because of the stacked NMOS transistors [20, 21], this leakage depends on which input(s) is low. The leakage is least when both inputs are low. Under standby operating conditions, for NAND2 and NAND3 gates in the decoding logic within the mats, we assume that the output of each NAND is high (deactivated) with both of its inputs low. Thus we attribute a leakage number to the NAND gate based on the leakage through its stacked NMOS transistors when both inputs are low. We consider the reduction in leakage due to the effect of stacked transistors and calculate leakage for the NAND2 gate as follows:

$$P_{\text{leak-nand2}} = W_{\text{inv-nmos}}I_{\text{off-nmos}}SF_{\text{nand2}} \quad (88)$$

where SF_{nand2} is the stacking fraction for reduction in leakage due to stacking.

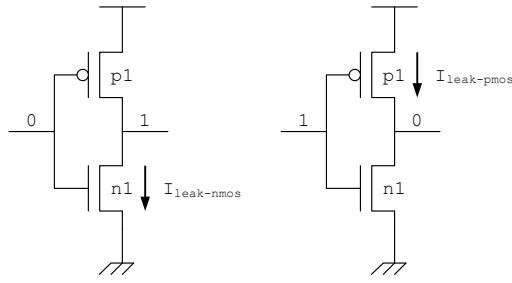


Figure 33: Leakage in an inverter.

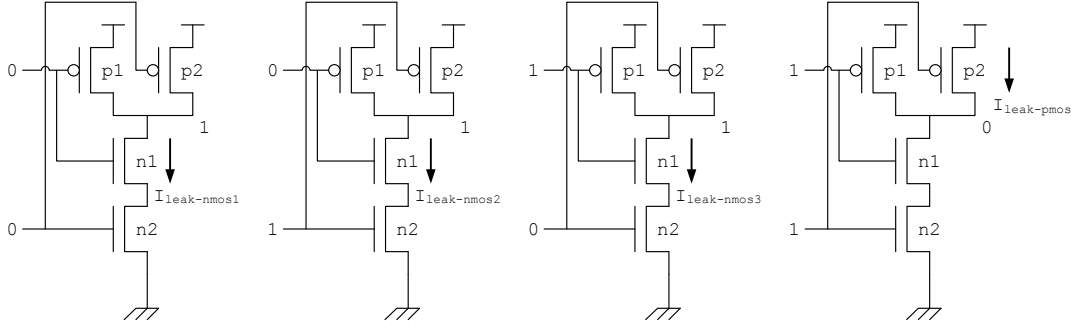


Figure 34: Leakage in a NAND2 gate.

7.2.2 Leakage Power Equations

Most of the leakage power equations are similar to the dynamic energy equations in form.

$$P_{\text{leak}} = P_{\text{leak-request-network}} + P_{\text{leak-mats}} + P_{\text{leak-reply-network}} \quad (89)$$

$$P_{\text{leak-mats}} = (P_{\text{leak-mem-cells}} + P_{\text{leak-predec-blks}} + P_{\text{leak-decoder-drivers}} + P_{\text{leak-senseamps}}) N_{\text{banks}} N_{\text{subbanks}} N_{\text{mats-in-subbank}} \quad (90)$$

$$P_{\text{leak-mem-cells}} = N_{\text{subarr-rows}} N_{\text{subarr-cols}} P_{\text{mem-cell}} \quad (91)$$

$$P_{\text{leak-decoder-drivers}} = P_{\text{leak-row-decoder-drivers}} + P_{\text{leak-bitmux-decoder-driver}} + P_{\text{leak-senseamp-mux-decoder-driver}} \quad (92)$$

$$P_{\text{leak-row-decoder-drivers}} = 4 N_{\text{subarr-rows}} P_{\text{leak-row-decoder-driver}} \quad (93)$$

$$P_{\text{leak-request-network}} = P_{\text{leak-arr-edge-to-bank-edge-request-htree}} + P_{\text{leak-bank-hor-request-htree}} + P_{\text{leak-bank-ver-request-htree}} \quad (94)$$

$$P_{\text{leak-reply-network}} = P_{\text{dyn-ver-reply-htree}} + P_{\text{dyn-hor-reply-htree}} + P_{\text{dyn-bank-edge-to-arr-edge-reply-htree}} \quad (95)$$

Figure 35 shows the subthreshold leakage paths in an SRAM cell when it is in idle/standby state [20, 21]. The leakage power contributed by a single memory cell may be given by:

$$P_{\text{mem-cell}} = V_{\text{DD}} I_{\text{mem-cell}} \quad (96)$$

$$I_{\text{mem-cell}} = I_{p1} + I_{n2} + I_{n3} \quad (97)$$

$$I_{p1} = W_{p1} I_{\text{off-pmos}} \quad (98)$$

$$I_{n2} = W_{n2} I_{\text{off-nmos}} \quad (99)$$

$$I_{n3} = W_{n2} I_{\text{off-nmos}} \quad (100)$$

Figure 36 shows the subthreshold leakage paths in a sense amplifier during an idle/standby cycle [20, 21].

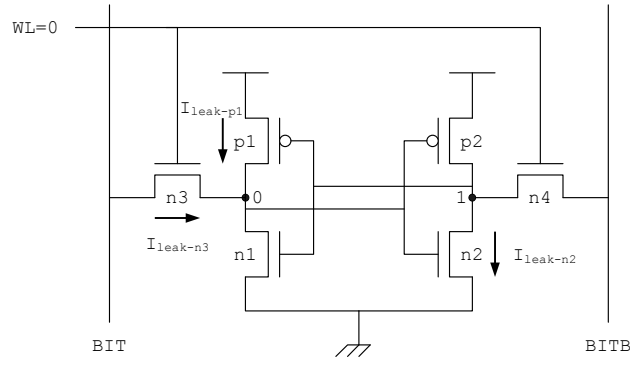


Figure 35: Leakage paths in a memory cell in idle state. BIT and BITB are precharged to V_{DD} .

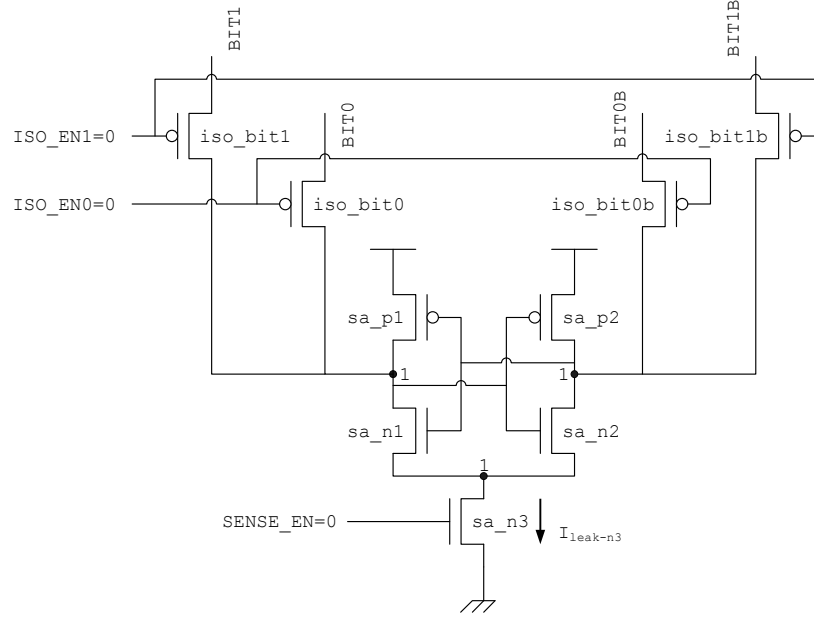


Figure 36: Leakage paths in a sense amplifier in idle state.

8 Technology Modeling

Version 5 makes use of technology projections from the ITRS [35] for device data and projections from [11, 13] for wire data. Currently we look at four ITRS technology nodes (we use MPU/ASIC metal 1 half-pitch to define the technology node) – 90, 65, 45, and 32 nm – which cover years 2004 to 2013 in the ITRS. Section 8.1 gives more details about the device data and modeling and Section 8.2 gives more details about the wire data and modeling.

8.1 Devices

Table 3 shows the characteristics of transistors modeled by the ITRS that are incorporated within CACTI. We include data for the three device types that the ITRS defines - High Performance (HP), Low Standby Power (LSTP) and Low Operating Power (LOP). The HP transistors are state-of-the-art fast transistors with short gate lengths, thin gate oxides, low V_{th} , and low V_{DD} whose CV/I is targeted to improve by 17% every year. As a consequence of their high on-currents, these transistors tend to be very leaky. The LSTP transistors on the other hand are transistors with longer gate lengths, thicker gate oxides, higher V_{th} , and higher V_{DD} . The gate-lengths of the LSTP transistors lag the HP transistors by 4 years. The LSTP transistors trade off high on-currents for maintenance of an almost constant low leakage of 10 pA across the technology nodes. The LOP transistors have performance that lie in between the HP and LSTP transistors.

Parameter	Meaning	Units
V_{DD}	Voltage applied between drain and source, gate and source	V
L_{gate}	Physical length of the gate	micron
V_{th}	Saturation threshold voltage	V
M_{eff}	Effective mobility	cm^2/Vs
V_{dsat}	Drain saturation voltage	V
$C_{ox-elec}$	Capacitance of gate-oxide in inversion	F/μ^2
$C_{gd-overlap}$	Gate to drain overlap capacitance	F/μ
$C_{gd-fringe}$	Gate to drain fringing capacitance	F/μ
$C_{j-bottom}$	Bottom junction capacitance	F/μ^2
$C_{j-sidewall}$	Sidewall junction capacitance	F/μ^2
I_{on}	On-current (saturation)	A/μ
I_{off}	Channel leakage current (for $V_{gate} = 0$ and $V_{drain} = V_{DD}$)	A/μ

Table 3: Technology characteristics of transistors used in the model.

Technology-node	90nm	65nm	45nm	32nm
L_{gate} (nm)	37/75/53	25/45/32	18/28/22	13/20/16
EOT (Equivalent oxide thickness) (nm)	1.2/2.2/1.5	1.1/1.9/1.2	0.65/1.4/0.9	0.5/1.1/0.8
V_{DD} (V)	1.2/1.2/0.9	1.1/1.2/0.8	1/1.1/0.7	0.9/1/0.7
V_{th} (mV)	237/525/318	195/554/315	181/532/256	137/513/242
I_{on} ($\mu A/\mu$)	1077/465/550	1197/519/573	2047/666/749	2496/684/890
I_{off} (nA/ μ)	32.4/0.008/2.0	196/0.009/4.9	280/0.01/4.0	1390/0.021/65
$C_{ox-elec}$ (fF/ μ^2)	17.9/12.2/16.0	18.8/13.6/18.7	37.7/20.1/28.2	45.8/22.9/31.2
τ (Intrinsic switching delay) (ps)	1.01/2.98/1.78	0.64/1.97/1.17	0.4/1.33/0.79	0.25/0.9/0.53
FO1 delay (ps)	7.3/25.1/19.9	4.8/18.1/10.0	2.75/11.5/6.2	1.63/7.13/3.51

Table 4: Values of key technology metrics of HP, LSTP, and LOP NMOS transistors for four technology-nodes from the 2005 ITRS [35].

They use the lowest V_{DD} to control the operating power and their gate-lengths lag those of HP transistors by 2 years. The CV/I of the LSTP and LOP transistors improves by about 14% every year.

Table 4 shows values of key technology metrics of the HP, LSTP, and LOP NMOS transistors for four technology nodes. The data is obtained from MASTAR [36] files. According to the 2003 ITRS², the years 2004, 2007, 2010, and 2013 correspond to 90, 65, 45, and 32nm technology-nodes. Because the 2005 ITRS does not include device data for the 90 nm technology-node (year 2004), we obtain this data using MASTAR and targeting the appropriate CV/I. Note that all values shown are for planar bulk devices. The ITRS actually makes the assumption that planar high-performance bulk devices reach their limits of practical scaling in 2012 and therefore includes multiple parallel paths of scaling for SOI and multiple-gate MOS transistors such as FinFETs starting from the year 2008 which run in parallel with conventional bulk CMOS scaling. We however use MASTAR to meet the target CV/I of the 32 nm node with planar bulk devices. For all technology nodes, the overlap capacitance value has been assumed to be 20% of ideal (no overlap) gate capacitance. The bottom junction capacitance value for the planar bulk CMOS transistors has been assumed to be $1fF/\mu^2$, which is the value that MASTAR assumes. As MASTAR does not model sidewall capacitance, we compute values for sidewall capacitance in the following manner: we use process data provided at the MOSIS website [37] for TSMC and IBM 130/180/250nm processes and compute average of the ratios of sidewall-to-bottom junction capacitances for these processes. We observe that average error in using this average value for projecting sidewall capacitance given bottom junction capacitance is less than 10%. We use this average value in projecting sidewall capacitances for the ITRS processes.

²Because of ambiguity associated with the “technology-node” term, the 2005 ITRS has discontinued the practice of using the term, however, for the sake of convenience, we continue to use it in CACTI.

Parameter	Value	Reference
$A_{\text{sram-cell}}$ (Area of an SRAM cell) (μ^2)	$146F^2$	[14]
$W_{\text{sram-cell-acc}}$ (Width of SRAM cell access transistor) (μ)	$1.31F$	[14]
$W_{\text{sram-cell-pd}}$ (Width of SRAM cell pull-down transistor) (μ)	$1.23F$	[14]
$W_{\text{sram-cell-pu}}$ (Width of SRAM cell pull-up transistor) (μ)	$2.08F$	[14]
$AR_{\text{sram-cell}}$ (Aspect ratio of the cell)	1.46	[14]

Table 5: Technology data assumed for an SRAM cell.

We calculate the drive resistance of a transistor during switching as follows:

$$R_{\text{on}} = \frac{V_{\text{DD}}}{I_{\text{eff}}} \quad (101)$$

The effective drive current is calculated using the following formula described in [26] [48]:

$$I_{\text{eff}} = \frac{I_{\text{H}} + I_{\text{L}}}{2} \quad (102)$$

where $I_{\text{H}} = I_{\text{DS}} (V_{\text{GS}} = V_{\text{DD}}, V_{\text{DS}} = \frac{V_{\text{DD}}}{2})$ and $I_{\text{L}} = I_{\text{DS}} (V_{\text{GS}} = \frac{V_{\text{DD}}}{2}, V_{\text{DS}} = V_{\text{DD}})$.

For PMOS transistors, we find the width of the transistor that produces the same I_{off} as a unit-width NMOS transistor. Using this width, we compute the PMOS effective drive current ($I_{\text{eff-pmos}}$) and the PMOS-to-NMOS sizing ratio that is used during the application of the method of logical effort:

$$S_{\text{pmos-to-nmos-logical-effort}} = \frac{I_{\text{eff-nmos}}}{I_{\text{eff-pmos}}} \quad (103)$$

Table 5 shows technology data that we have assumed for an SRAM cell.

It may be useful to know that while currently we provide device data for just the three ITRS device types, it is not difficult to incorporate device data from other sources into CACTI. Thus, published data of various industrial fabrication processes or data from sources such as [50] may also be utilized. Also, by making use of MASTAR, it is possible to obtain device data for scaling models and assumptions that are different from those of the ITRS. As an example, while the ITRS device data for its High Performance device type is based on an improvement in device CV/I of 17 % every year, one may obtain alternative device data by targeting a different CV/I improvement and/or I_{off} . Another example is to start off with the ITRS High Performance device type and use MASTAR to come up with higher V_t or longer channel variations of the base device.

8.2 Wires

Wire characteristics in CACTI are based on the projections made in [11, 13]. The approach followed in [11, 13] is to consider both aggressive (optimistic) and conservative (pessimistic) assumptions regarding interconnect technology. The aggressive projections assume aggressive use of low-k dielectrics, insignificant resistance degradation due to dishing and scattering, and tall wire aspect ratios. The conservative projections assume limited use of low-k dielectrics, significant resistance degradation due to dishing and scattering, and smaller wire aspect ratios. For these assumptions, [11, 13] looks at two types of wires, semi-global and global. Wires of semi-global type have a pitch of $4F$ (F = Feature size) whereas wires of global type have a pitch of $8F$. We incorporate the properties of both these wire types into CACTI. The values of the semi-global and global wire characteristics under aggressive and conservative assumptions are presented in Table 6 for 90/65/45/32 technology nodes. The resistance per unit length and capacitance per unit length values are calculated based off Equations 5 and 6 respectively. For the capacitance per unit micron calculation, we assume a Miller factor of 1.5 as a “realistic worst-case” value [31]. For material strength, we assume that low-k dielectrics are not utilized between wire layers as suggested in [31].

Technology-node	90nm	65nm	45nm	32nm
Common wire characteristics (aggressive/conservative)				
$\rho(m\Omega.\mu)$	0.022/0.022	0.018/0.022	0.018/0.022	0.018/0.022
$\epsilon_r for C_c$	2.709/3.038	2.303/2.734	1.958/2.46	1.664/2.214
Semi-global wire properties (aggressive/conservative)				
Pitch(nm)	360	280	180	128
Aspect ratio	2.4/2.0	2.7/2.0	3.0/2.0	3.0/2.0
Thickness (nm)	432/400	351/280	270/200	192/140
ILD (nm)	480/480	405/405	315/315	210/210
Miller factor	1.5/1.5	1.5/1.5	1.5/1.5	1.5/1.5
Barrier (nm)	10/8	0/6	0/4	0/3
Dishing (%)	0/0	0/0	0/0	0/0
$\alpha_{scatter}$	1/1	1/1	1/1	1/1
Resistance per unit length (Ω/μ)	0.33/0.38	0.34/0.73	0.74/1.52	1.46/3.03
Capacitance per unit length (fF/ μ)	0.314/0.302	0.302/0.282	0.291/0.265	0.269/0.254
Global wire properties (aggressive/conservative)				
Pitch(nm)	800	560	400	280
Aspect ratio	2.7/2.2	2.8/2.2	3.0/2.2	3.0/2.2
Thickness (nm)	1080/880	784/616	600/440	420/308
ILD (nm)	960/1100	810/770	630/550	420/385
Miller factor	1.5	1.5	1.5	1.5
Barrier (nm)	10/8	0/6	0/4	0/3
Dishing (%)	0/10	0/10	0/10	0/10
$\alpha_{scatter}$	1/1	1/1	1/1	1/1
Resistance per unit length (Ω/μ)	0.067/0.09	0.095/0.17	0.19/0.36	0.37/0.72
Capacitance per unit length (fF/ μ)	0.335/0.315	0.308/0.298	0.291/0.281	0.269/0.267

Table 6: Aggressive and conservation wire projections from [11].

8.3 Technology Exploration

As an additional feature in version 5, we allow the user to map different device and wire types to different parts of the array. We divide the devices in the array into two parts: one, devices used in the memory cells and wordline drivers, and two, the rest of the peripheral and global circuitry. Different device types such as the ITRS HP, LSTP, LOP or other user-added device types may be mapped to the devices in the two parts of the array³. We divide the wires in the array also into two parts, wires inside mats and wires outside mats. Different wire types such as the semi-global or global wire types or other user-defined wire types may be mapped to the wires inside and outside mats.

9 Embedded DRAM Modeling

In this section, we describe our modeling of embedded DRAM.

9.1 Embedded DRAM Modeling Philosophy

We model embedded DRAM and assume a logic-based embedded DRAM fabrication process [19, 24, 27]. A logic-based embedded DRAM process typically means that DRAM has been embedded into the logic process without affecting the characteristics of the original process much [23]. In our modeling of embedded DRAM, we leverage the similarity that

³It is important to note that in reality, SRAM cell functionality and design does depend on device type [8, 12, 22], however, we do not model different SRAM cell designs for the different device types.

exists in the global and peripheral circuitry of embedded SRAM and DRAM and model only their essential differences. We also use the same array organization for embedded DRAM that we used for SRAM. By having a common framework that, in general, places embedded SRAM and DRAM on an equal footing and emphasizes only their essential differences, we would be able to compare relative tradeoffs involving embedded SRAM and DRAM.

We capture the following essential differences between embedded DRAM and SRAM in our area, delay and power models:

9.1.1 Cell

The most essential difference between SRAM and DRAM is in their storage cell. While SRAM typically makes use of a 6T cell and the principle of positive feedback to store data, DRAM typically makes use of a 1T-1C cell and relies on the charge-storing capability of a capacitor. Because it makes use of only one transistor, a DRAM cell is usually laid out in a much smaller area compared to an SRAM cell. For instance the embedded DRAM cells presented in [45] for four different technology nodes – 180/130/90/65nm have areas in the range of $19\text{--}26F^2$ where F is the feature size of the process. In contrast, a typical SRAM cell would have an area of about $120\text{--}150F^2$.

9.1.2 Destructive Readout and Writeback

When data is read out from a DRAM cell, the charge stored in the cell gets destroyed because of charge redistribution between the cell and its capacitive bitline. Because of the destructive readout, there is a need for data to be written back into the cell after every read access. This writeback takes time and increases the random cycle time of a DRAM array. In an SRAM there is no need for writeback because the data is not destroyed during a read.

9.1.3 Sense Amplifier Input Signal

In a DRAM, the maximum differential signal developed on the bitlines is limited by the amount of charge transferred between the DRAM cell and the bitline which in turn depends on the capacitance of the DRAM cell and the bitline. The lower the differential signal, the greater the sense amplifier delay. In an SRAM, there is no charge-based limit on the differential signal developed on the bitlines. In any case, in modern technologies the sense amplifiers of SRAMs or DRAMs are operating at signal level inputs of more or less the same amplitude [23], so the delay of the sense amplifier in either SRAM or DRAM can come out to have similar values.

9.1.4 Refresh

In a DRAM cell, charge cannot be stored for an infinite time in the capacitor and the charge leaks out because of various leakage components. If charge from a DRAM cell is allowed to leak out for a sufficient period of time, the differential voltage developed on the bitline pair becomes so small that the data stored in the cell can no longer be detected by the sense amplifier. Thus there is an upper bound on the time for which data may be retained in a DRAM cell without it being refreshed, and this time is known as the retention time. Because of a finite retention time, the DRAM cell needs to be refreshed periodically.

9.1.5 Wordline Boosting

In a DRAM cell, because the access takes place through an NMOS pass transistor, there is a V_{th} drop during the write/writeback of a 1 into the cell. In order to prevent this V_{th} drop, DRAM wordlines are usually boosted to a voltage, $V_{PP} = V_{DD} + V_{th}$. In commodity DRAMs, V_{th} is relatively high (typically around 1V for 65nm) in order to maintain the high refresh period (64ms) that requires extremely low leakage. This means that V_{PP} is also high (typically around 3.3V for 65nm) and forces the use of high voltage (thicker gate-oxide) slower transistors in the wordline driver. For the embedded DRAMs that we have modeled, however, V_{th} is not very high (0.44V for 65nm), consequently V_{PP} is also not very high (1.6V for 65nm).

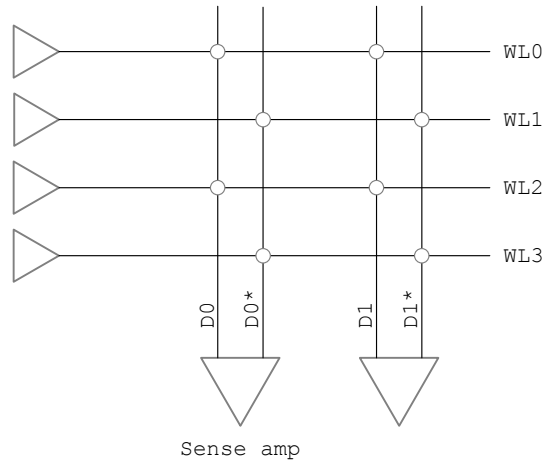


Figure 37: Folded array architecture from [17].

9.2 DRAM Array Organization and Layout

For DRAM, we assume a folded array architecture [17] in the subarray, shown in Figure 37. In the folded array architecture, the bitline being read (true bitline) and its complement are laid out next to each other, similar to the dual bitlines of an SRAM cell. The difference here is that the true and complement bitlines connect to alternate rows of the array and not to the same row as in SRAM. This has an impact on bitline capacitance calculation. Assuming drain contacts are shared, the bitline capacitance for DRAM may be given by the following equation:

$$C_{\text{bitline}} = \frac{N_{\text{subarr-rows}}}{2} C_{\text{drain-cap-acc-transistor}} + N_{\text{subarr-rows}} C_{\text{bit-metal}} \quad (104)$$

9.2.1 Bitline Multiplexing

In DRAM, the read access is destructive. This means that during a read access after data is read from a DRAM cell, it needs to be written back into the cell. This writeback is typically accomplished by using the sense amplifier which detects the data stored in the cell during a read. During a read access, because each cell connected to a wordline is read out through its associated bitline, this means that there needs to be a sense-amplifier associated with each cell connected to a wordline. Hence bitline multiplexing, which is common in SRAMs to connect multiple bitlines to a single sense amplifier, is not feasible in DRAMs. Thus in DRAMs, there needs to be a sense amplifier associated with every bitline that can carry out the writeback. With respect to the bitline peripheral circuitry shown in Figure 6 this means that DRAM arrays do not have a bitline mux between the bitlines and sense amplifiers.

9.2.2 Reference Cells for V_{DD} Precharge

We assume that the bitlines are precharged to V_{DD} (GND) just like the DRAM described in [5, 27]. As in [5], we assume the use of reference cells that store $V_{DD}/2$ and connect to the complement bitline during a read. Figure 38 shows the bitline peripheral circuitry with the reference cells. For each subarray, we assume an extra two rows of reference cells that store $V_{DD}/2$. One of the rows with reference cells is activated during read of even-numbered rows in the subarray and the other row is activated during read of odd-numbered rows in the subarray.

9.3 DRAM Timing Model

9.3.1 Bitline Model

In DRAM the differential voltage swing developed on a bitline pair that acts as input to the sense amplifier is limited by the ratio of charge transferred between the bitline and DRAM cell, and given by the following equation:

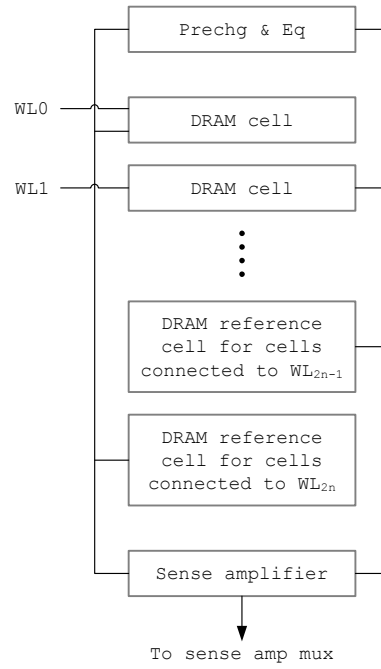


Figure 38: DRAM bitline circuitry showing reference cells for V_{DD} precharge.

$$V_{\text{sense-max}} = \frac{V_{DD}}{2} \frac{C_{\text{dram}}}{C_{\text{dram}} + C_{\text{bitline}}} \quad (105)$$

The delay for the above differential signal to develop may be given by the following equation [6] (ignoring the effect of wordline rise time):

$$T_{\text{step}} = 2.3R_{\text{dev}} \frac{C_{\text{dram}}C_{\text{bitline}}}{C_{\text{dram}} + C_{\text{bitline}}} \quad (106)$$

where R_{dev} is the resistance in series with the storage capacitor of the DRAM cell and may be given by the following equation:

$$R_{\text{dev}} = \frac{V_{DD}}{I_{\text{cell-on}}} \quad (107)$$

It is important to note that use of Equations 106 and 107 assumes that the impact of bitline resistance on signal development time is negligible. This approximation works well for contemporary logic-based embedded DRAM processes. When bitline resistance becomes significant, as in the case of commodity DRAM processes that do not make use of copper bitlines, more sophisticated models need to be used.

Equation 106 assumes that 90% of the data stored in the cell is read out and corresponds to the development of approximately $V_{\text{sense-max}}$ (given by Equation 105) on the bitline pair. In order to improve the random cycle time of a DRAM macro further, nowadays less than 90% of the data stored in a cell is read out [15], just enough to generate the required input signal of the sense amplifier ($V_{\text{senseamp-input}}$). To accommodate this case, Equation 106 may be generalized as follows:

$$T_{\text{step-generalized}} = 2.3R_{\text{dev}} \frac{C_{\text{dram}}C_{\text{bitline}}}{C_{\text{dram}} + C_{\text{bitline}}} \frac{V_{\text{senseamp-input}}}{V_{\text{sense-max}}} \quad (108)$$

When $V_{\text{senseamp-input}}$ is equal to $V_{\text{sense-max}}$, Equation 108 reduces to Equation 106. In CACTI, we assume a certain value for $V_{\text{senseamp-input}}$ (such as 80mV) and use Equation 108 to compute the signal development delay.

When rise time of the wordline is also considered, the bitline delay (T_{bitline}) of DRAM may be calculated using the same methodology that was used for SRAM (Equation 53 in Section 6).

The time taken to write data back into a DRAM cell after a read depends on the time taken for the charge transfer to take place between the bitline and the DRAM and may be given by the following equation:

$$T_{\text{writeback}} = T_{\text{step}} \quad (109)$$

9.3.2 Multisubbank Interleave Cycle Time

For a DRAM array, we consider three timing characteristics: random access time, random cycle time and multibank interleave cycle time. Calculation of random access time makes use of the same equations that were used for calculation of random access time of an SRAM array (in Section 6). For a DRAM array, typically there are two kinds of cycle time: random cycle time and multibank interleave cycle time. Random cycle time has the same meaning as the random cycle time of an SRAM viz. it is the time interval between two successive random accesses. This time interval is typically limited by the time it takes to activate a wordline, sense the data, write back the data, and then precharge the bitlines. Random cycle time can thus be calculated using the following equation:

$$T_{\text{random-cycle}} = T_{\text{row-dec-driver}} + T_{\text{bitline}} + T_{\text{senseamp}} + T_{\text{writeback}} + T_{\text{wordline-reset}} + \max(T_{\text{bitline-precharge}}, T_{\text{bit-mux-out-precharge}}, T_{\text{senseamp-mux-out-precharge}}) \quad (110)$$

In order to improve the rate at which a DRAM array is accessed so that it is not limited by the random cycle time of the array, DRAM arrays usually employ the concept of multibank interleaving. Multibank interleaving takes advantage of the fact that while random access to a particular bank is limited by the random cycle time, accesses to other banks need not be. With multibank interleaving, accesses to multiple DRAM banks that are on the same address/data bus are interleaved at a rate defined by the *multibank interleave cycle time*. In our terminology, each *bank* in an array has its own address and data bus and may be concurrently accessed. For our array organization, the concept of multibank interleaved mode is relevant to subbank access and not bank access, so in the rest of this discussion we use the terminology of multisubbank interleave mode and multisubbank interleave cycle. Thus, the multisubbank interleave cycle time is the rate at which accesses may be interleaved between different subbanks of a bank. The multisubbank interleave cycle time depends on the degree of pipelining employed in the request and reply networks of a subbank, and is limited by the pipelining overhead. We assume minimum pipeline overhead and use the following simple equation to calculate multisubbank-interleave cycle time:

$$T_{\text{multisubbank-interleave}} = \max(T_{\text{request-network}} + T_{\text{row-predec}}, T_{\text{reply-network}}) \quad (111)$$

9.3.3 Retention Time and Refresh Period

An equation for the retention time of a DRAM array may be written as follows [16]:

$$T_{\text{retention}} = \frac{C_{\text{dram-cell}} \Delta V_{\text{cell-worst}}}{I_{\text{worst-leak}}} \quad (112)$$

where $\Delta V_{\text{cell-worst}}$ is the worst-case change in the voltage stored in a DRAM cell which leads to a read failure, and $I_{\text{cell-worst-leak}}$ is the worst-case leakage in a DRAM cell.

We assume that $\Delta V_{\text{cell-worst}}$ is limited by $V_{\text{min-sense}}$, the minimum input signal that may be detected by the bitline sense amplifier. Thus, for a given array organization, $\Delta V_{\text{cell-worst}}$ may be calculated by solving the following equation for $\Delta V_{\text{cell-worst}}$:

$$V_{\text{min-sense}} = \frac{C_{\text{dram-cell}}}{C_{\text{dram-cell}} + C_{\text{bitline}}} \left(\frac{V_{\text{DD}}}{2} - \Delta V_{\text{cell-worst}} \right) \quad (113)$$

If we assume that the differential voltage detected by the sense amplifier is independent of array organization, then this means that different array partitions would have different retention times depending on the charge transfer ratio between the DRAM cell and the bitlines. For each array organization, it's thus possible to calculate the value for $\Delta V_{\text{cell-worst}}$ using Equation 113, which may then be plugged into Equation 112 to find the retention time for that array organization.

The upper bound on the refresh period of a DRAM cell would be equal to its retention time. We assume that a safety margin of 10% with respect to the retention time is built into the refresh period and thus calculate the refresh period using the following equation:

$$T_{\text{refresh}} = 0.9T_{\text{retention}} \quad (114)$$

9.4 DRAM Power Model

During the read of a 0 from a DRAM cell, the true bitline is pulled down to GND during the writeback. Energy is then consumed in restoring the bitline to V_{DD} during the precharge operation. During the read of a 1 from a DRAM cell, because of our assumption of V_{DD} -precharge, the voltage of the true bitline does not change but the voltage of the complementary bitline gets pulled down to GND and needs to be restored to V_{DD} . So for DRAM, the power consumed in a bitline during a read may be approximated by the following equation:

$$E_{\text{dyn-read-bitline}} = C_{\text{bitline}} V_{\text{DD}}^2 \quad (115)$$

9.4.1 Refresh Power

Refreshing the data in each cell of the array consumes power. In order to carry out refresh, every cell in the array needs to be accessed, its data read out, and then written back.

$$P_{\text{refresh}} = \frac{E_{\text{refresh}}}{T_{\text{refresh}}} \quad (116)$$

$$E_{\text{refresh}} = E_{\text{refresh-predec-blks}} + E_{\text{refresh-row-dec-drivers}} + E_{\text{refresh-bitlines}} \quad (117)$$

$$E_{\text{refresh-predec-blks}} = N_{\text{banks}} N_{\text{subbanks}} N_{\text{mats-in-subbank}} E_{\text{dyn-mat-predec-blks}} \quad (118)$$

$$E_{\text{refresh-row-dec-drivers}} = 4N_{\text{banks}} N_{\text{subbanks}} N_{\text{mats-in-subbank}} E_{\text{dyn-mat-row-dec-drivers}} \quad (119)$$

$$E_{\text{refresh-bitlines}} = 4N_{\text{banks}} N_{\text{subbanks}} N_{\text{mats-in-subbank}} N_{\text{subarr-cols}} E_{\text{dyn-read-bitline}} \quad (120)$$

9.5 DRAM Area Model

9.5.1 Area of Reference Cells

As mentioned earlier in Section 9.2.2, the use of V_{DD} -precharge leads to the use of reference cells in the array [5]. For our array organization, this means that there are two additional wordlines per subarray.

9.5.2 Area of Refresh Circuitry

In order to enable continued scaling of a logic-based embedded DRAM cell in terms of performance and cell area, [45] describes a new scalable embedded DRAM cell that makes use of an access transistor with an intermediate gate-oxide of moderate thickness (2.2nm for 90/65nm). This transistor is a standard offering in the logic process which incorporates the embedded DRAM. Conventional cells [14] in earlier technologies made use of access transistors with much thicker gate-oxides. An effect of the scalable embedded DRAM cell described in [45] is that it results in the cell having a lower retention time and a lower refresh period (because of higher worst-case leakage - 10s of pAs compared to 1 fA for commodity DRAM). The macro discussed in [18] that makes use of the cell described in [45] has a refresh period of 64 μs compared to conventional macros which have refresh period of 64 ms. This low refresh period required innovation

Parameter	Meaning	Unit
C_{dram}	Storage capacitance of a DRAM cell	F
$A_{\text{dram-cell}}$	Area occupied by the DRAM cell	mm^2
$AR_{\text{dram-cell}}$	Aspect ratio of the DRAM cell	
$V_{\text{DDdram-cell}}$	Voltage representing a 1 in a DRAM cell	V
$V_{\text{th-dram-acc-transistor}}$	Threshold voltage of DRAM cell access transistor	mV
$L_{\text{dram-acc-transistor}}$	Length of DRAM cell access/wordline transistor	nm
$W_{\text{dram-acc-transistor}}$	Width of DRAM cell access transistor	nm
$EOT_{\text{dram-acc-transistor}}$	Equivalent oxide thickness of DRAM access transistors	nm
$I_{\text{on-dram-cell}}$	DRAM cell on-current under nominal conditions	μA
$I_{\text{off-dram-cell}}$	DRAM cell off-current under nominal conditions	pA
$I_{\text{worst-off-dram-cell}}$	DRAM cell off-current under worst-case conditions	A/μ
V_{pp}	Boosted wordline voltage applied to gate of access transistor	V
$I_{\text{on-dram-wordline-transistor}}$	On-current of wordline transistor	$\mu\text{A}/\mu$

Table 7: Characteristics of the DRAM cell and wordline driver.

Parameter	90nm	65nm
C_{dram} (F)	20	20
$A_{\text{dram-cell}}$ (F - feature size)	$20.7F^2$	$25.6F^2$
$V_{\text{DDdram-cell}}$	1.2	1.2
$V_{\text{th-dram-acc-transistor}}$	350	350
$L_{\text{dram-acc-transistor}}$ (nm)	120	120
$W_{\text{dram-acc-transistor}}$	140	90
$I_{\text{on-dram-cell}}$ (μA)	45	36
$I_{\text{off-dram-cell}}$ (pA)	2	2
V_{pp}	1.5	1.5

Table 8: DRAM technology data for 90nm and 65nm from [18, 45].

at the circuit level through the development of a concurrent refresh scheme described in [18] in order to guarantee high availability of the DRAM macro. This concurrent refresh scheme adds an extra bank select port to each bank (subbank in our terminology) thereby allowing for concurrent memory access and bank refresh operations in different banks. Each bank is equipped with a row address counter that contains the address of the row to be refreshed. A concurrent refresh scheduler composed of an up-shift-register and a down-shift-register is required in order to generate the bank select signals.

Because we loosely base the parameters of our logic-based embedded DRAM technology on information presented in [5, 18, 45], we model the overhead of the concurrent refresh scheme on area. For our organization in which each subbank has multiple mats, we assume that each mat incurs overhead of a row address counter placed at the center of the mat. Because of the requirements of the design of the concurrent refresh scheme, for our organization, we assume $N_{\text{subbanks-in-mat}}$ number of concurrent refresh schedulers per bank.

9.6 DRAM Technology Modeling

9.6.1 Cell Characteristics

Similar to the SRAM technology assumptions, we assume two types of transistors in the DRAM array. One transistor type is used in the DRAM cell and wordline driver, while the other is used in the rest of the peripheral and global circuitry. Table 3 showed a list of transistor characteristics used in CACTI. Table 7 shows characteristics of the DRAM cell and wordline driver that we consider in our model.

We obtain embedded DRAM technology data for four technology nodes – 90, 65, 45, and 32 nm – by using an

Parameter	90nm	65nm	45nm	32nm
C_{dram} (F)	20	20	20	20
$A_{\text{dram-cell}}$ (F - Feature size)	$20.7F^2$	$25.6F^2$	$30.4F^2$	$30.6F^2$
$V_{\text{DDdram-cell}}$ (V)	1.2	1.2	1.1	1.1
$V_{\text{th-dram-acc-transistor}}$ (mV)	455	438	446	445
$L_{\text{dram-acc-transistor}}$ (nm)	120	120	78	56
$W_{\text{dram-acc-transistor}}$ (nm)	140	90	79	56
$I_{\text{on-dram-cell}}$ (μA)	45	36	36	36
$I_{\text{off-dram-cell}}$ (pA)	2	2	2	2
$I_{\text{worst-off-dram-cell}}$ (pA)	21.1	19.6	19.5	18.9
V_{PP} (V)	1.6	1.6	1.5	1.5

Table 9: Characteristics of DRAM cell for the four technology nodes.

approach that makes use of published data, transistor characterization using MASTAR, and our own scaling projections. For 90 nm and 65 nm, we use technology data from [18, 45]; Table 8 shows this data. In order to obtain technology data for the 45nm and 32nm technology nodes, we make the following scaling assumptions:

1. Capacitance of the DRAM cell is assumed to remain fixed at 20fF;
2. The nominal off-current is assumed to remain fixed at 2pA for the cell;
3. Gate oxide thickness is scaled slowly in order to keep gate leakage low and subthreshold current as the dominant leakage current. It has a value of 2.1 nm for 45 nm and 2 nm for 32 nm;
4. $V_{\text{DDdram-cell}}$ is scaled such that the electric field in the dielectric of the DRAM ($V_{\text{PP}}/EOT_{\text{dram-acc-transistor}}$) access transistor remains almost constant;
5. There is excellent correlation in the 180–130nm (for conventional thick-oxide device) and 90–65nm (for the intermediate-oxide device) scaling-factors for width and length of the DRAM cell access transistor. We assume that there would be good correlation in the 130–90nm and 65–45nm scaling-factors as well. For 32 nm, we assume that the width and length are scaled in the same proportion as feature size;
6. We calculate area of the DRAM cell using the equation $A_{\text{dram-cell}} = 10W_{\text{dram-acc-transistor}}L_{\text{dram-acc-transistor}}$. This equation has good correlation with the actual cell area of the 90 and 65 nm cells that made use of the intermediate-oxide based devices; and
7. We simply assume that nominal on-current of the cell can be maintained at the 65 nm value. This would require aggressive scaling of the series parasitic resistance of the transistor.

With these scaling assumptions, we use MASTAR to model the transistors. It is assumed that the resulting channel doping concentrations calculated by MASTAR would be feasible. Table 9 summarizes the characteristics of the DRAM cell for the four technology nodes.

10 Cache Modeling

In this section we describe how a cache has been modeled in version 5. The modeling methodology is almost identical to earlier versions of CACTI with a few changes that simplify the code.

10.1 Organization

As described in [47], a cache has a tag array in addition to a data array. In earlier versions of CACTI the data and tag arrays were modeled separately with separate code functions even though the data and tag arrays are structurally very

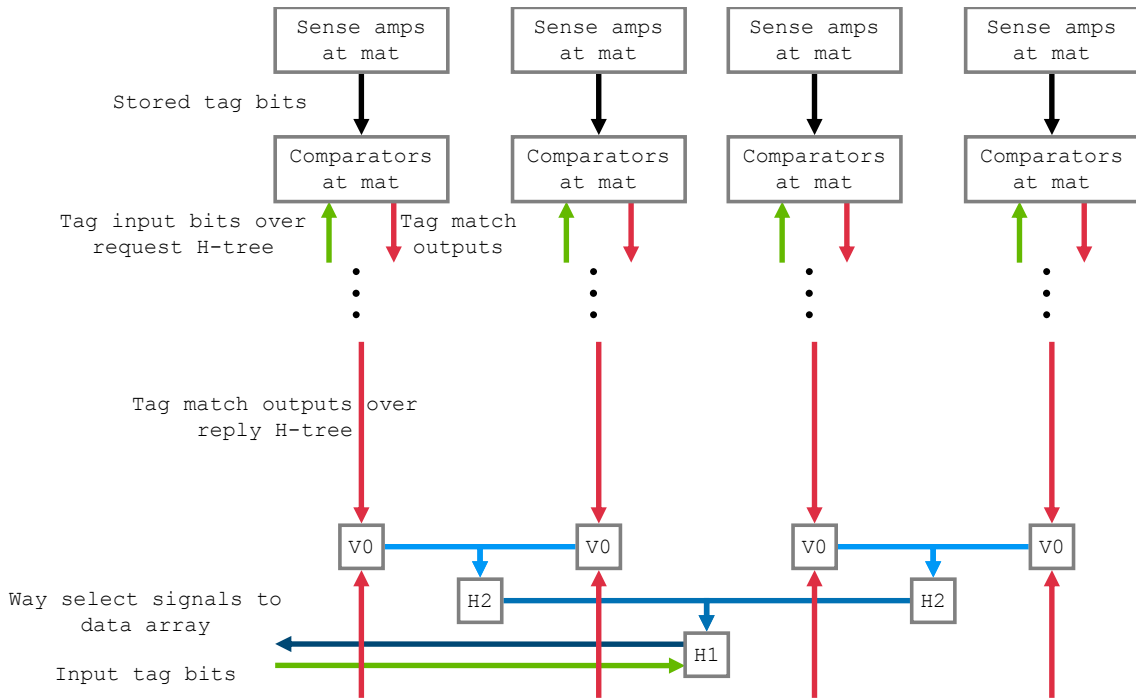


Figure 39: Organization of a set-associative tag array.

similar. The essential difference between the tag array and the data array is that the tag array includes comparators that compare the input tag bits with the stored tags and produce the tag match output bits. Apart from the comparators, the rest of the peripheral/global circuitry and memory cells are identical for data and tag arrays. In version 5, we leverage this similarity between the data and tag arrays and use the same set of functions for their modeling. For the tag array, we reuse the comparator area, delay, and power models.

Figure 39 illustrates the organization of a set-associative tag array. Each mat includes comparators at the outputs of the sense amplifiers. These comparators compare the stored tag bits with the input tag bits and produce the tag match output bits. These tag match output signals are the way-select signals that serve as inputs to the data array. The way-select signals traverse over the vertical and horizontal H-trees of the tag array to get to the edge of the tag array from where they are shipped to the data array. For a cache of normal access type, these way-select signals then enter the data array where, like the address and datain signals, they travel along the horizontal and vertical H-tree networks to get to mats in the accessed subbank. At the mats, these way-select signals are ‘anded’ with sense amplifier mux decode signals (if any) and the resultant signals serve as select signals for the sense amplifier mux which generates the output word from the mat.

10.2 Delay Model

We present equations for access and cycle times of a cache. The access time of a cache depends on the type of cache access (normal, sequential, or fast [40]).

The equation for access time of a normal cache which is set-associative is as follows:

$$T_{\text{access-normal-set-associative}} = \max(T_{\text{tag-arr-access}} + T_{\text{data-arr-request-network}} + T_{\text{data-arr-senseamp-mux-decode}}, T_{\text{data-arr-request-network}} + T_{\text{data-arr-mat}}) + T_{\text{data-arr-reply-network}} \quad (121)$$

$$T_{\text{tag-arr-access}} = T_{\text{tag-arr-request-network}} + T_{\text{tag-arr-mat}} + T_{\text{tag-arr-reply-network}} \quad (122)$$

In the above equation, $T_{\text{tag-arr-access}}$, the access time of the tag array, is calculated using the following equation.

$$T_{\text{tag-arr-access}} = T_{\text{tag-arr-request-network}} + T_{\text{tag-arr-mat}} + T_{\text{tag-arr-reply-network}} + T_{\text{comparators}} \quad (123)$$

The equation for access time of a normal cache which is direct-mapped is as follows:

$$T_{\text{access-normal-direct-mapped}} = \max(T_{\text{tag-arr-access}}, T_{\text{data-arr-access}}) \quad (124)$$

The equation for access time of a sequentially accessed (tag array is accessed first before data array access begins) cache is as follows:

$$T_{\text{access-sequential}} = T_{\text{tag-arr-access}} + T_{\text{data-arr-access}} \quad (125)$$

The equation for access time of a ‘fast’ cache (late way-select multiplexing) cache is as follows:

$$T_{\text{access-fast}} = \max(T_{\text{tag-arr-access}}, T_{\text{data-arr-access}}) + T_{\text{way-select-mux}} \quad (126)$$

where $T_{\text{way-select-mux}}$ is the delay through the way-select mux. The way-select mux is assumed to be placed at the edge of the data array and selects the appropriate output word corresponding to the correct way (which is selected based on the way-select signals from the tag array).

10.3 Area Model

Total area of the cache is calculated by simply adding the areas occupied by the tag and data arrays.

$$A_{\text{cache}} = A_{\text{data-array}} + A_{\text{tag-array}} \quad (127)$$

$A_{\text{tag-array}}$ is calculated using the equations presented in Section 5 with the area of the comparators also added.

10.4 Power Model

The dynamic energy consumed in the cache and its standby leakage power are calculated by simply adding their values for the data and tag arrays. For the tag array, the leakage in the comparators is also considered.

$$E_{\text{dyn-energy-cache}} = E_{\text{dyn-energy-data-array}} + E_{\text{dyn-energy-tag-array}} \quad (128)$$

$$P_{\text{leak-cache}} = P_{\text{leak-data-array}} + P_{\text{leak-tag-array}} \quad (129)$$

11 Quantitative Evaluation

In this section we evaluate the impact of new CACTI 5 features. We also compare results from CACTI 5.1 with version 4.2 in order to give users an idea of what to expect when upgrading to CACTI 5.1.

11.1 Evaluation of New CACTI 5 Features

Table 10 shows the default parameters that we have used while carrying out the evaluation of new CACTI 5 features. For this evaluation we use plain RAMs instead of caches. In Section 12, we show results for caches. For each study, we present charts that show the following metrics: access time, random cycle time, area, dynamic energy per read access, and standby leakage power.

Parameter	Value
Capacity (MB)	16
Output width (bits)	512
Number of banks	1
Number of read/write ports	1
Number of exclusive read ports	0
Number of exclusive write ports	0
Technology-node (nm)	65
DRAM	No
max_area_constraint	40
max_acc_time_constraint	10
max_repeater_delay_constraint	10
Optimize for dynamic energy	No
Optimize for dynamic power	No
Optimize for leakage power	No
Optimize for cycle time	Yes
Temperature (K)	360
SRAM cell/wordline technology flavor	ITRS HP
Peripheral/Global circuitry technology flavor	ITRS HP
Interconnect projection type	Conservative
Wire type inside mat	Semi-global
Wire type outside mat	Semi-global

Table 10: CACTI input parameters

11.1.1 Impact of New CACTI Solution Optimization

Figure 40 shows the impact of varying max_area_constraint (that was described in Section 2.5.1) for a 16MB SRAM. As max_area_constraint is increased, the number of subarrays in the SRAM is allowed to grow, and so the area grows steadily. As the number of subarrays increases, the components of delay within a mat decreases and the access time falls up to a point after which it starts to increase again. The random cycle time keeps decreasing as the number of subarrays increases because the wordline and bitline delays keep getting smaller. The trend for dynamic read energy per access shows some up-and-down variation. For our assumed CACTI cache organization, an increase in Ndw1 typically increases the dynamic energy consumption because more wordlines are activated per read access, while an increase in Ndbl typically decreases the dynamic energy consumption because of reduction in bitline power. The standby leakage power keeps growing as the area of the RAM increases.

Figure 41 shows the impact of varying max_acc_time_constraint. For the assumed set of SRAM parameters, it can be seen that the solution with best access time (corresponding to max_acc_time_constraint of 0) also has the best dynamic read energy per access. So further relaxation of the max_acc_time_constraint does not lead to any energy reduction benefits in this case.

Figure 42 shows the impact of varying max_repeater_delay_constraint. The max_repeater_delay_constraint changes the separation distance and sizing of repeaters/buffers in the H-tree networks and is useful for trading off delay for energy benefits. It can be seen here that varying max_repeater_delay_constraint does not lead to energy savings much unless the access time is allowed to degrade heavily. Initially, as max_repeater_delay_constraint is increased from 0 to 20%, it can be seen that the access time does not change and there are no energy savings. This is because of the maximum limit that we have imposed on transistor size in CACTI. For values of max_repeater_delay_constraint between 0 and 20%, the sizing of the repeater comes out to be larger than the maximum allowed transistor size and is therefore being fixed at the maximum allowed transistor size (the maximum allowed transistor size was fixed at $100F$ (F = feature size) for NMOS transistors). For a max_repeater_delay_constraint of 400% there is significant energy savings but with a disproportionate degradation of access time.

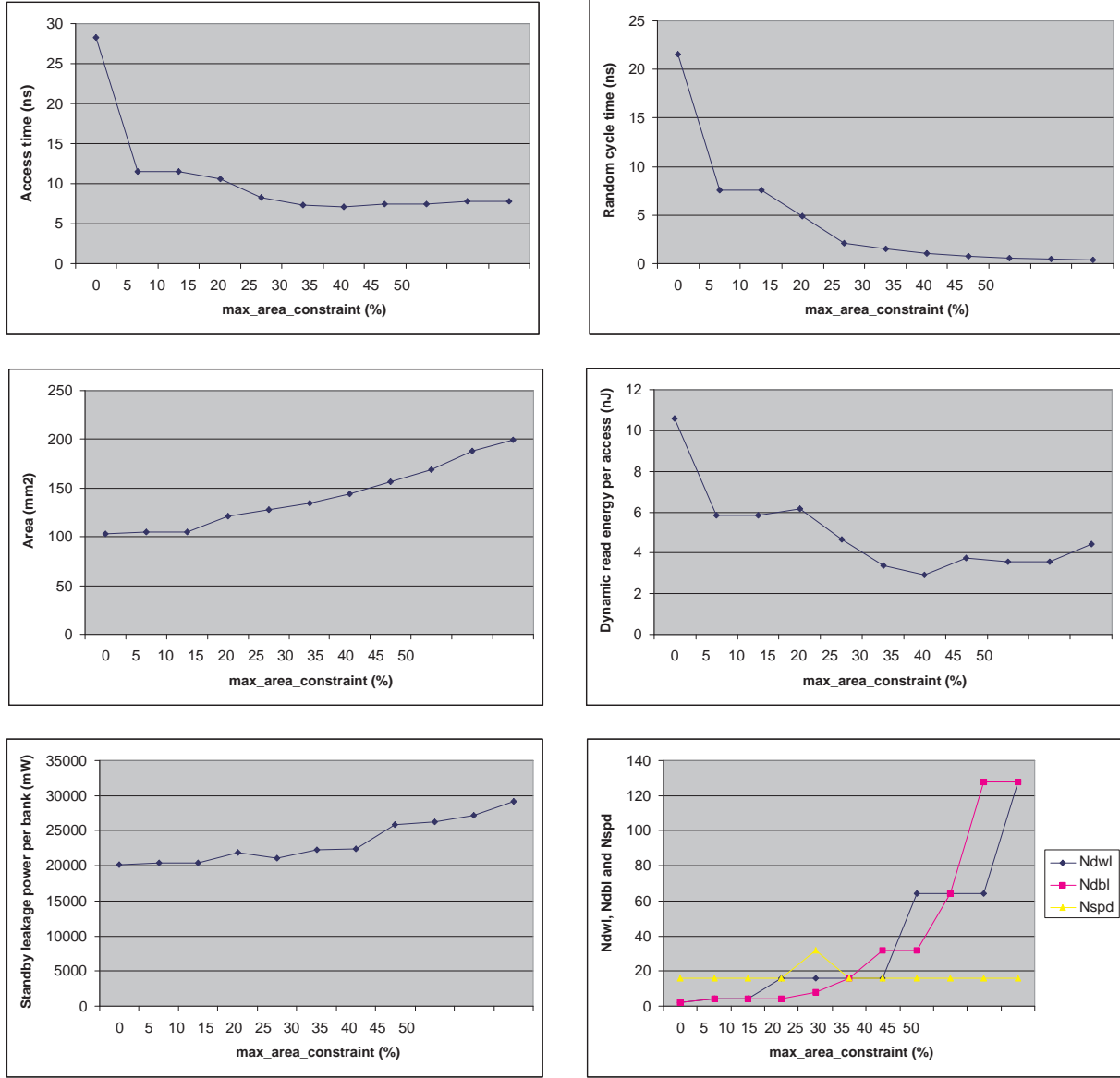


Figure 40: Access time, random cycle time, area, dynamic energy, and leakage power of a 16MB SRAM as `max_area_constraint` is varied.

Figure 43 shows the impact of optimizing the solution generated by CACTI for a 16MB SRAM in different ways. Table 11 shows the different optimization scenarios targeting metrics of random cycle time, dynamic read energy per access, dynamic power, and standby leakage power. The percentage variation between the worst and best values for each metric shown in Figure 43 is as follows: access time (4%), random cycle time (273%), area (28%), dynamic read energy per access (38%), and standby leakage power (24%). These variations illustrate the dependence of RAM and cache performance estimation on the kind of optimization that is applied.

11.1.2 Impact of Device Technology

Figure 44 illustrates the tradeoffs associated with assuming different types of devices in the memory cells/wordline drivers and the rest of the peripheral/global circuitry. Three scenarios are considered:

1. ITRS HP only;

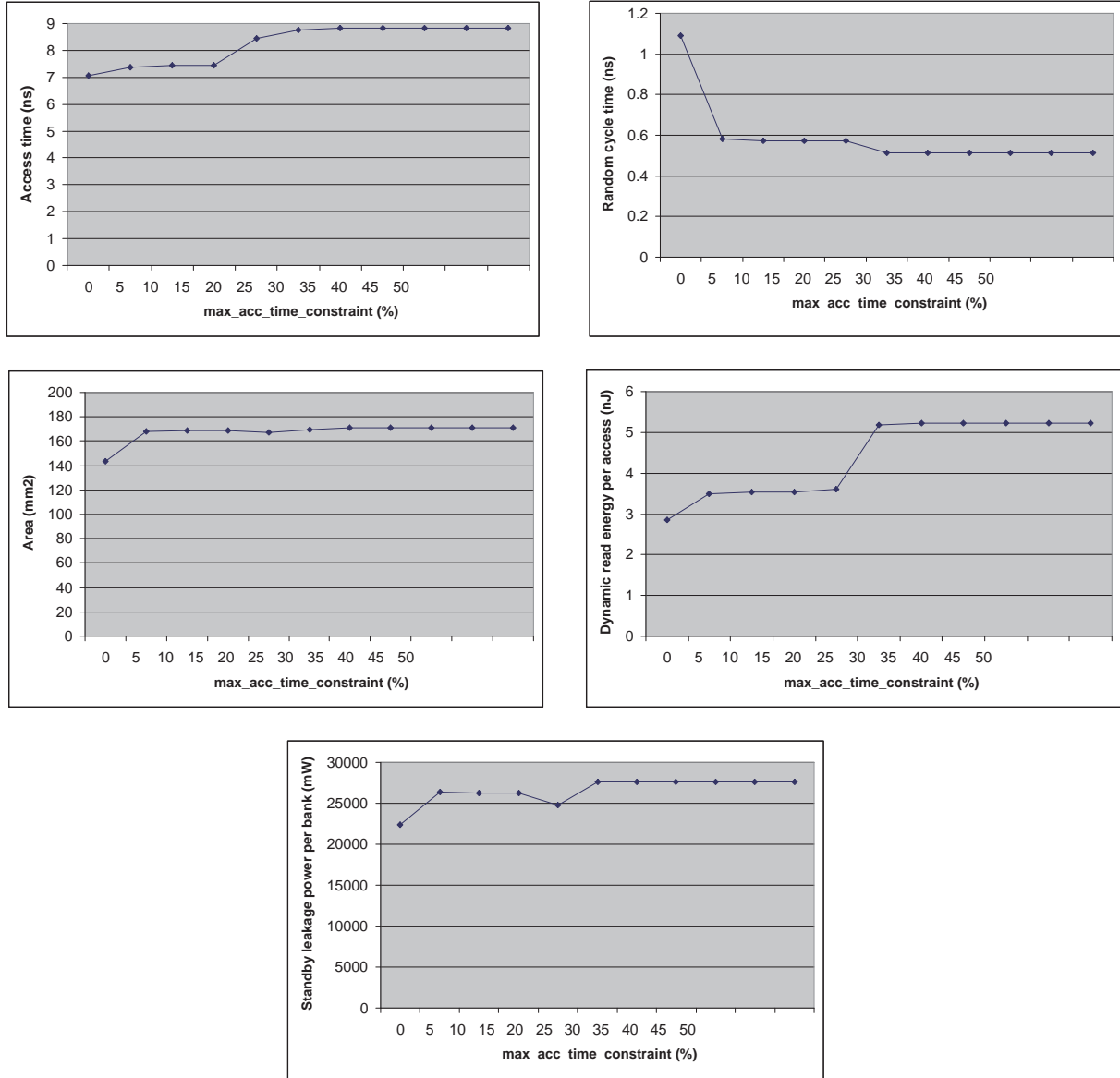


Figure 41: Access time, random cycle time, area, dynamic energy, and leakage power of a 16MB SRAM as max_acc_time_constraint is varied.

Optimization Scenario	Optimize for random cycle time	Optimize for dynamic energy	Optimize for dynamic power	Optimize for leakage power
A	Yes	No	No	No
B	Yes	Yes	No	No
C	No	Yes	No	No
D	No	No	Yes	No
E	No	No	No	Yes
F	Yes	Yes	Yes	Yes

Table 11: Different solution optimization scenarios targeting metrics of random cycle time, dynamic read energy per access, dynamic power, and standby leakage power.

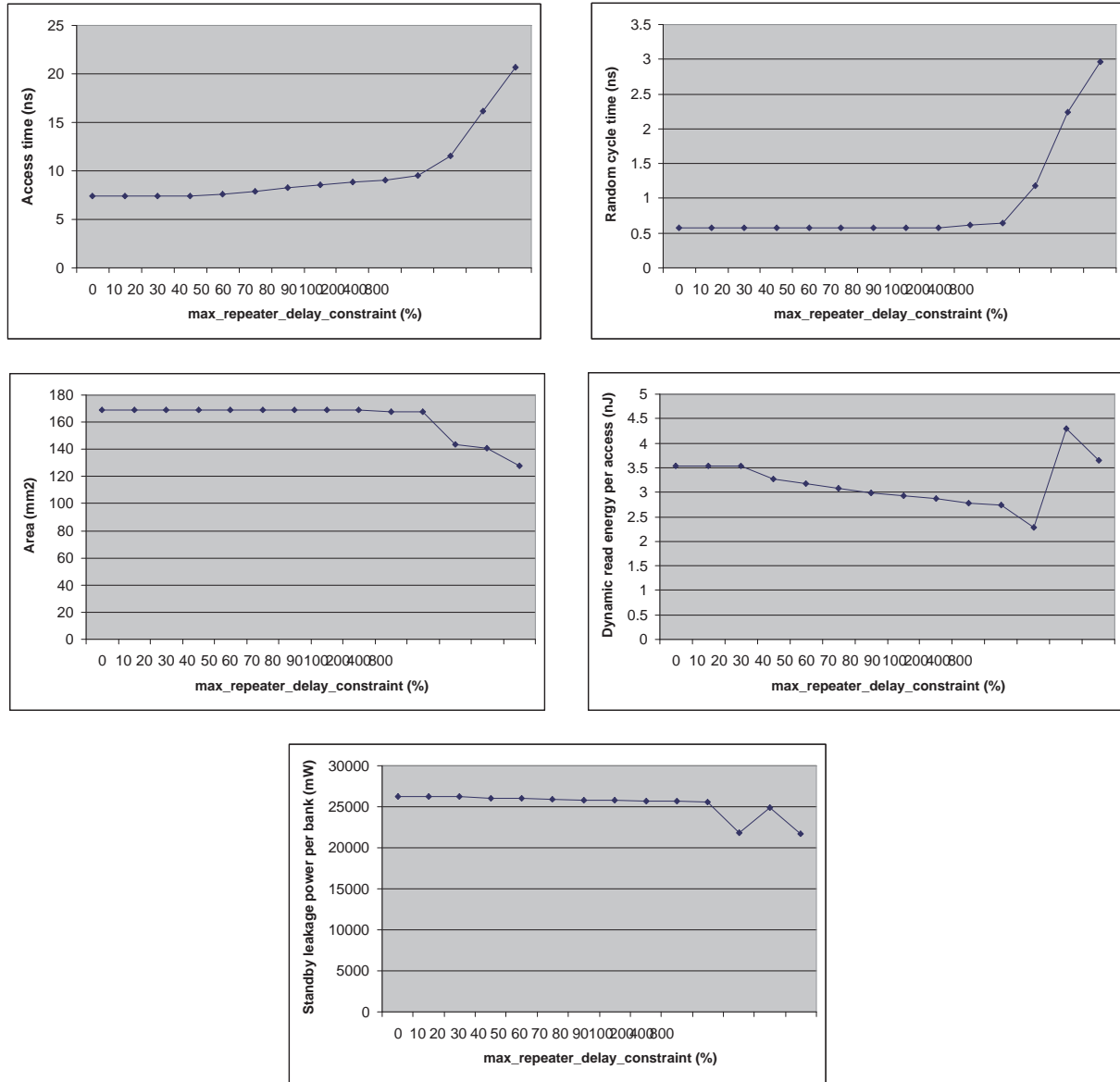


Figure 42: Access time, random cycle time, area, dynamic energy, and leakage power of a 16MB SRAM as max_repeater_delay_constraint is varied.

2. ITRS LSTP (memory cells/wordline drivers) + ITRS HP (peripheral/global circuitry);
3. ITRS LSTP only;

It can be seen that the areas of RAMs for the 3 considered scenarios remain more or less the same. With respect to “ITRS HP only”, on average over the considered capacities, “ITRS LSTP + ITRS HP” exhibits an improvement of 76% in the standby leakage power. This improvement comes at the cost of 11% worse access time and 39% worse random cycle time. “ITRS LSTP only” shows an improvement of almost 100% in standby leakage power with respect to “ITRS HP only”, and this improvement comes at a cost of 159% worse access time and 246% worse random cycle time.

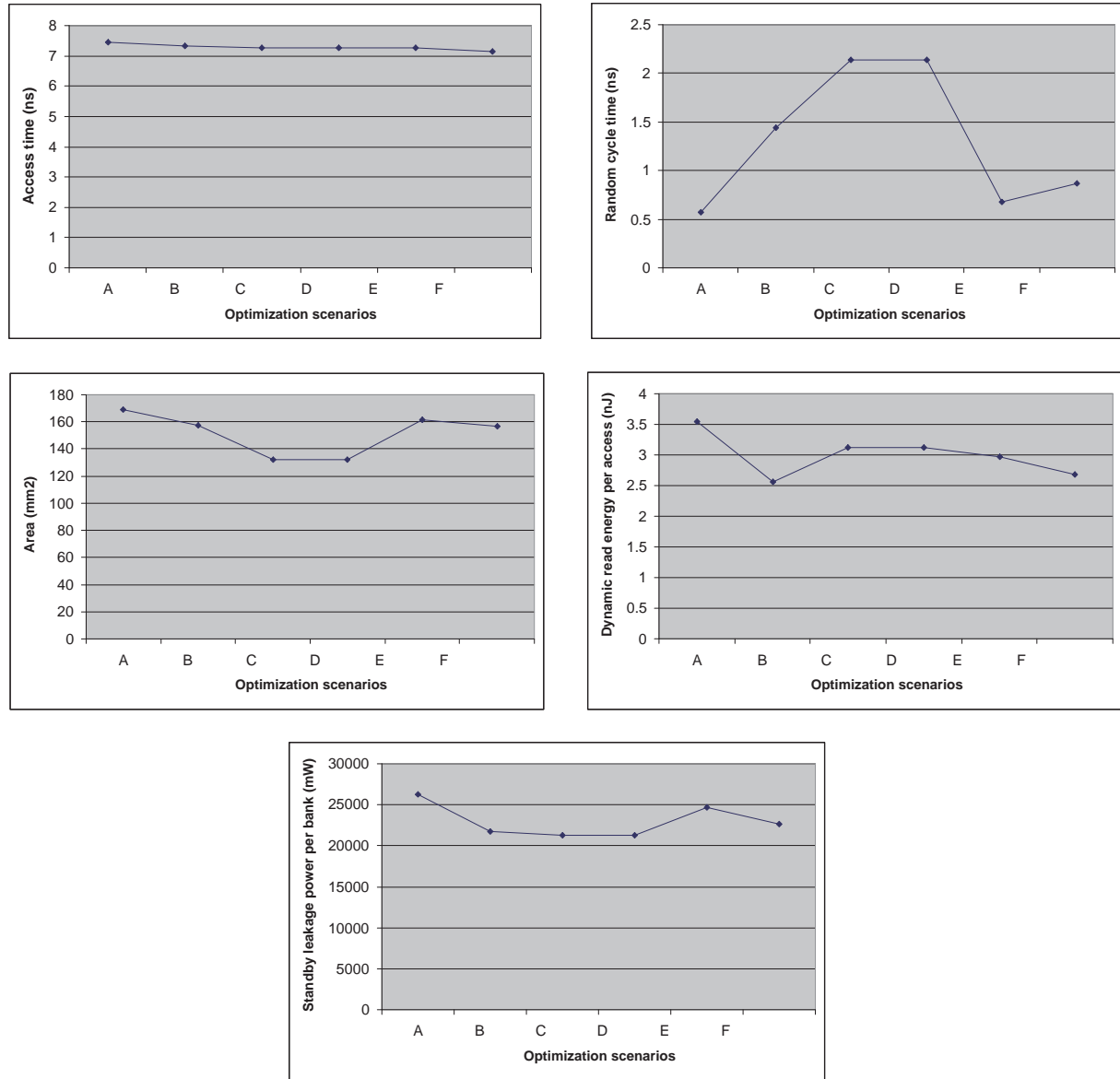


Figure 43: Access time, random cycle time, area, dynamic energy, and leakage power of a 16MB SRAM under different optimization function scenarios.

11.1.3 Impact of Interconnect Technology

Figure 45 illustrates the dependence of RAM/cache performance on interconnect technology assumptions. As described in Section 8 on “Technology Modeling”, instead of assuming a single set of scaling assumptions for interconnect technology, we consider aggressive and conservative scaling projections as in [11, 13]. From Figure 45, it can be seen that for the SRAM capacities and technology (65nm) considered, the cache performance is not very different under either conservative or aggressive interconnect technology assumptions. The lower resistance per unit length of the aggressive projections leads to lowering of the access time by about 11% on an average. For smaller technologies, the impact of interconnect technology assumptions would be more.

Figure 46 shows the impact of wire type on cache performance. As described in Section 8 on “Technology Modeling”, wires outside a mat can be of either ‘semi-global’ or ‘global’ type. With respect to semi-global type, global type wires outside mats lead to an improvement in access time of 22% on an average with about 3% increase in area. The

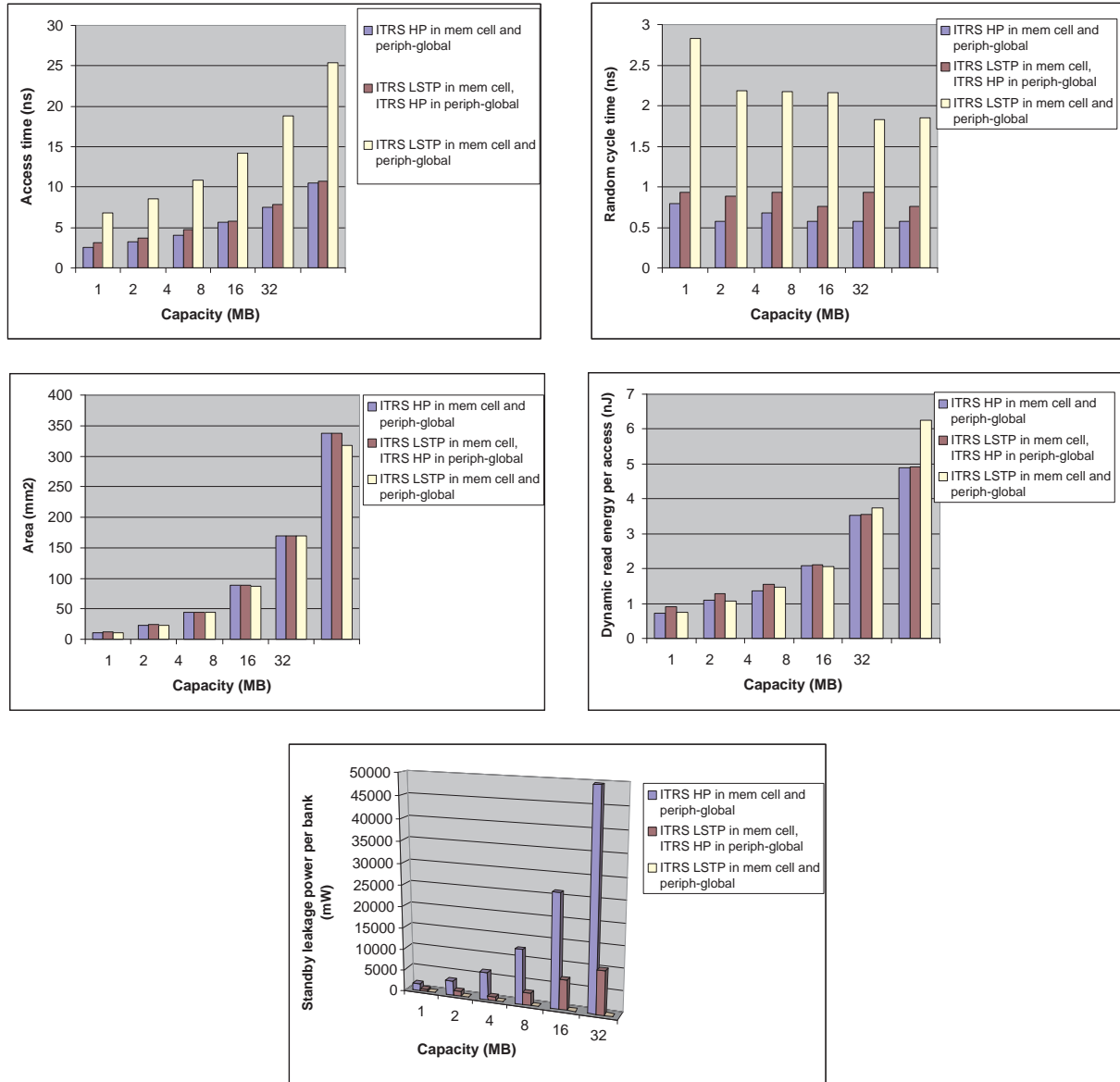


Figure 44: Access time, random cycle time, area, dynamic energy, and leakage power of SRAMs for different 65nm device technology assumptions.

global type wires take up greater area than the semi-global type wires, so the number of mats in the bank with global wire type is fewer than that with semi-global wire type. This leads to increase in the random cycle time because of greater wordline and bitline delays.

11.1.4 Impact of RAM Cell Technology

Figure 47 illustrates the dependence of cache performance on the type of RAM cell technology – SRAM or logic-based embedded DRAM. For each capacity, the properties of SRAMs and DRAMs shown in Figure 47 are for CACTI solutions with the best access time amongst all solutions. It can be seen from this figure that up to about 2 MB capacity, the access time of SRAM is lower after which the access time of DRAM becomes lower. This is because of the decreased wire delay experienced in DRAMs which occupy a much smaller area compared to the SRAMs. For the larger-capacity RAMs (≥ 1 MB), on average over all capacities, the area of the DRAMs is about a factor of 2.6 smaller than that of

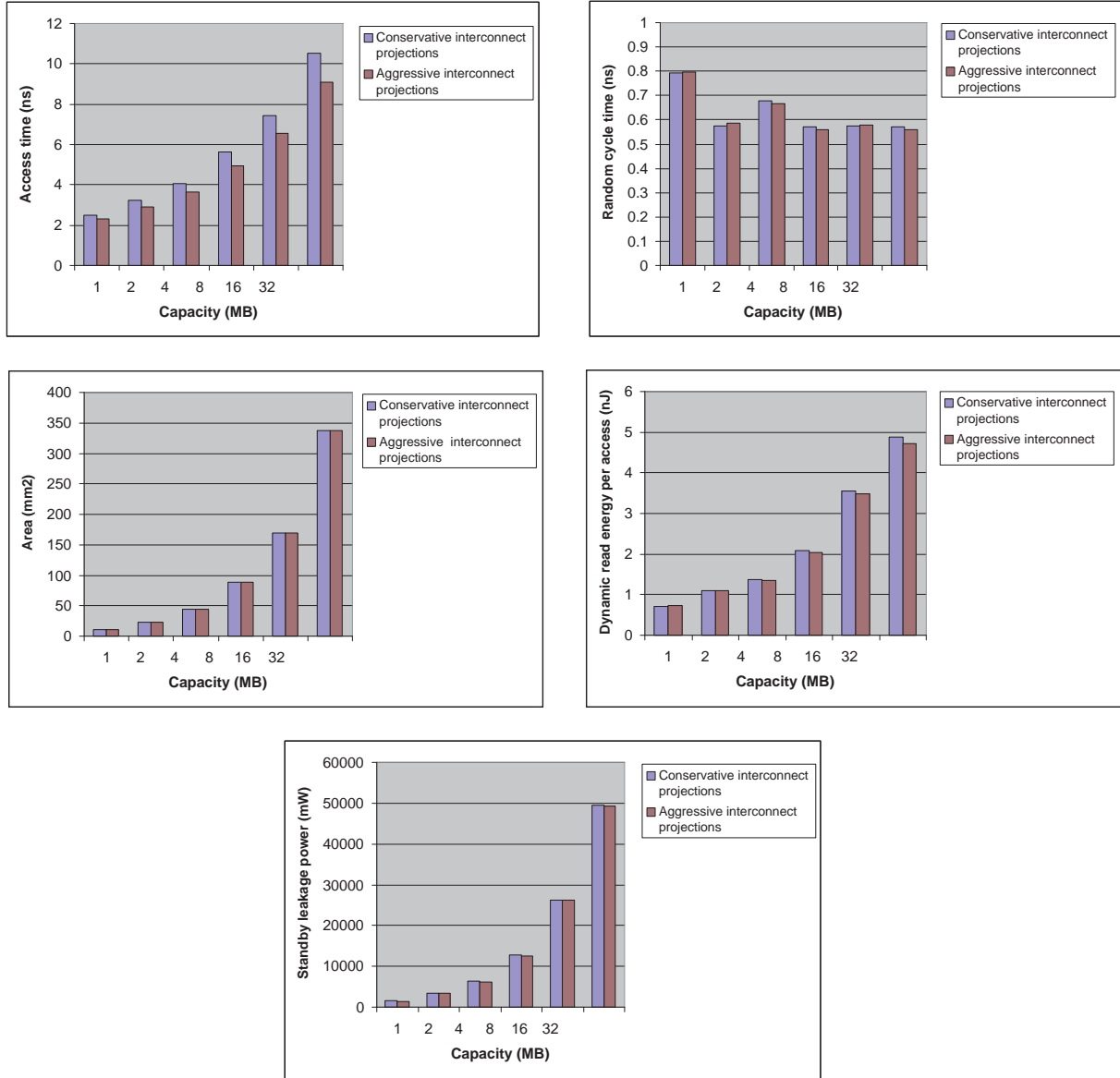


Figure 45: Access time, random cycle time, area, dynamic energy, and leakage power of SRAMs under aggressive and conservative interconnect technology assumptions.

the SRAMs. Because of the destructive read out from a capacitor and the subsequent writeback, it can be seen that the random cycle time of the DRAM is much higher than that of the SRAM. On an average over all capacities it is higher by a factor of about 2.2. However, it can be seen that the multisubbank interleave cycle time of the DRAMs can come close to the random cycle time of SRAMs up to about 4 MB. With further pipelining in the request and reply networks, the multisubbank interleave cycle time can be improved further.

The standby leakage power of the DRAM is much lower than that of the SRAMs because of the use of the low-leakage 1T cell. For the larger-capacity RAMs ($\geq 1\text{MB}$), on average the standby leakage power of DRAMs is lower than that of SRAMs by a factor of about 6. For the larger-capacity RAMs, on average the dynamic read energy per access of the DRAMs is more or less the same as that of the SRAMs. It is important to note that the comparisons presented here were for SRAMs and DRAMs with the best access time amongst all solutions.

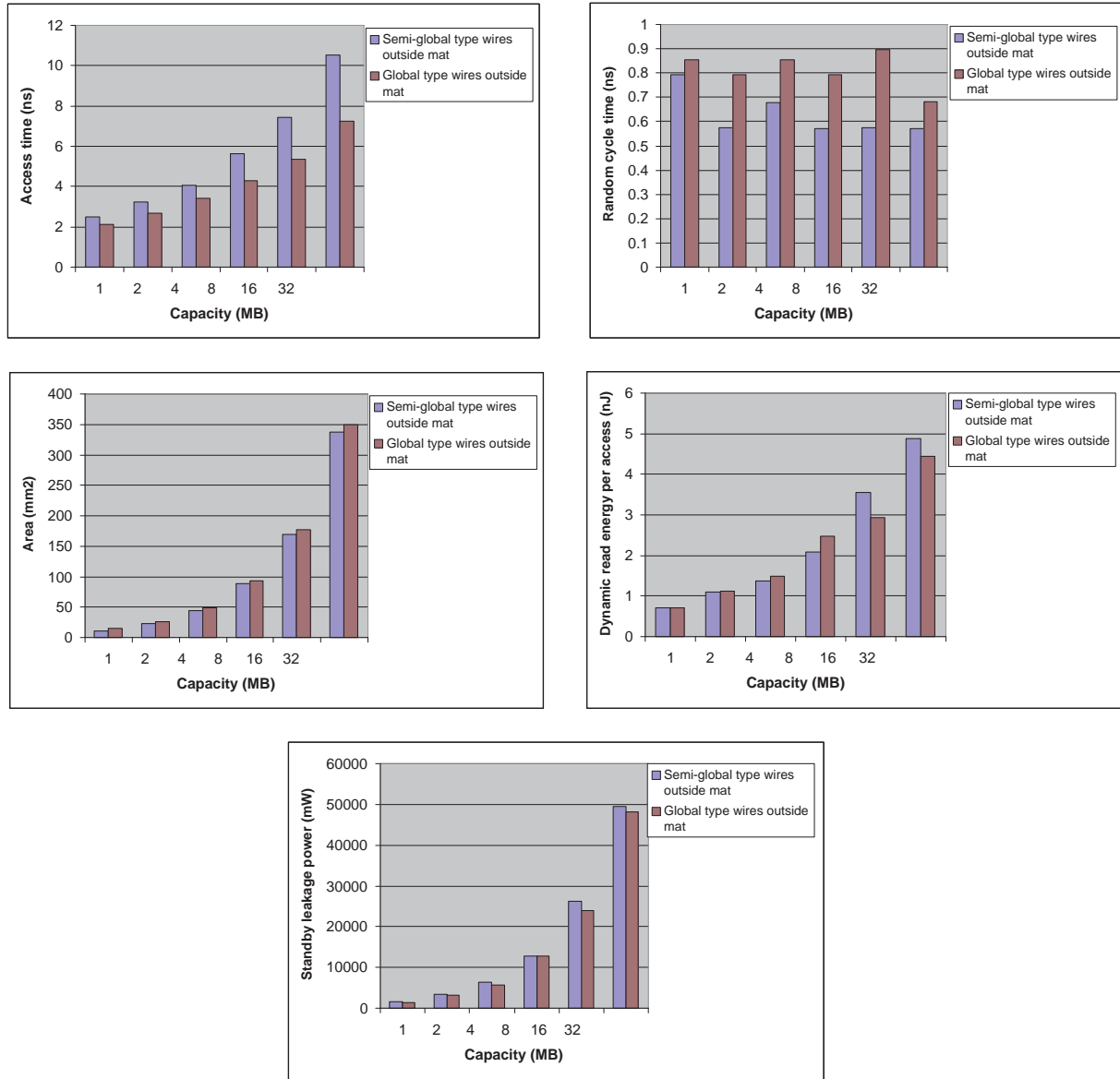


Figure 46: Access time, random cycle time, area, dynamic energy, and leakage power of SRAMs with “global” interconnect type used for wires outside the mat.

11.2 Version 4.2 vs Version 5.1 Comparisons

We first present the differences in the technology metrics of versions 4.2 and 5.1. Figure 48 shows FO4 delays for versions 4.2 and 5.1. For version 5, the FO4 delay has been shown for the ITRS HP device type. It can be seen that, surprisingly, there is good agreement between the FO4 delays of versions 4.2 and 5.1, particularly at the 65nm and 45nm nodes.

Figures 49, 50 and 51 shows the trends for resistance per unit micron, capacitance per unit micron and unbuffered delay for a wire of length 1-micron for versions 4.2 and 5.1. For version 4.2, trends are shown for both ‘local’ and ‘global’ wire types which were discussed in Section 2.3. For version 5.1, trends are shown for both ‘semi-global’ and ‘global’ wire types and for both aggressive and conservative assumptions of interconnect technology as discussed in Section 8.2. It can be seen that the unbuffered wire delay of the local wire type of version 4.2 is greater than that of the semi-global wire type of version 5.1 by a factor of about 3.5.

Figure 52 compares results generated by CACTI version 4.2 and version 5.1. For version 5.1, we show three solutions – solutions with best access time, best area and best dynamic read energy per access. We also show results from a modified version of version 4.2. In the modified version, we increase the limits defined for N_{dwl} , N_{dbl} , N_{spd} , number of subarrays, and number of segments allowed in the H-trees. Increasing these set limits allows a richer exploration of the search space of array partitions, particularly for large SRAMs. We call this version, version 4.2-limits-removed.

Because of significant modeling differences between versions 4.2 and 5.1, it is not easy to compare and analyze the behavior of the two versions, but we make general observations and present high-level plausible reasons for the trends. Firstly, regarding area, we observe that the areas of solutions produced by version 5.1 are much greater than that of the version 4.2 solution. Version 5.1 has made a major update to the area model and has introduced new features such as inclusion of ECC and redundancy overheads which increase area overhead. Also version 5.1 makes use of a bigger SRAM cell at $146F^2$ compared to the $120F^2$ cell used in version 4.2. For the 32 MB SRAM, the solution generated by version 4.2-limits-removed has an area of 343 mm^2 which is greater than the areas of all version 5.1 solutions.

Regarding access time, it can be seen that the access times for the best access time and best dynamic energy solutions of version 5.1 are much lower and scale in a much better way compared to the version 4.2 solution. The access time of the 32 MB SRAM gets better with version 4.2-limits-removed, however, it's still much worse than the access time of the version 5.1 best access time and best dynamic energy solutions. The main reason for the larger access times of version 4.2 is because of the high resistance per unit length of the local wires in version 4.2.

Regarding dynamic energy, it can be seen that the dynamic energy per read access of the version 5.1 best access time and best dynamic energy solutions are greater than that of version 4.2 by a factor of about 5 on an average. This is mainly because of the organization that has been assumed in version 5.1 in which wordlines and bitlines in multiple mats are activated per access. Also, as described in Section 2.2, bugs in version 4.2 with respect to dynamic energy of routing in the H-trees causes the dynamic energy per read access to be underestimated.

Leakage power is heavily dependent on the underlying technology data. The standby leakage of version 4.2 is greater than that of the best access time and best dynamic energy solutions of version 5.1 by a factor of about 3 on an average.

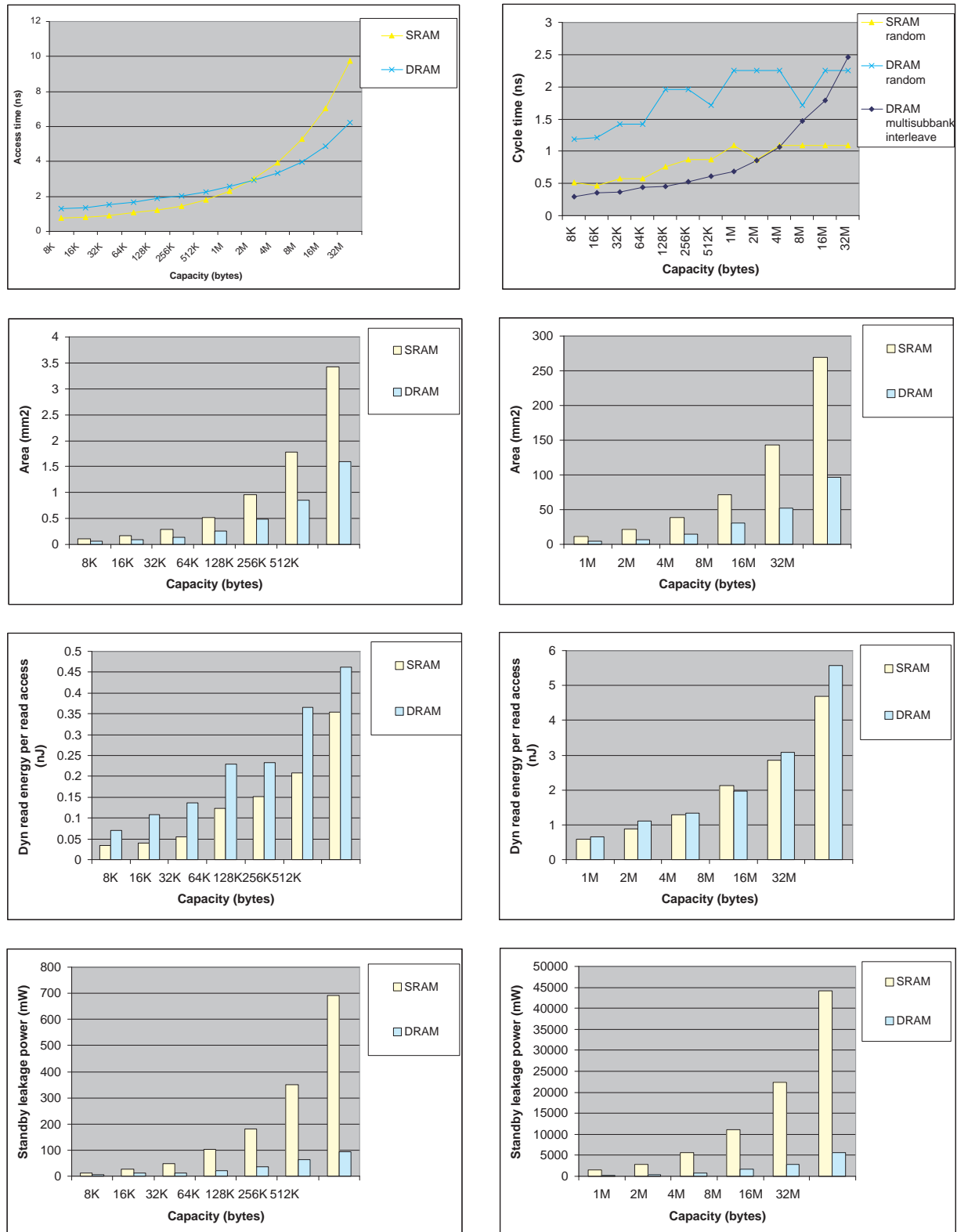


Figure 47: Access time, cycle time, area, dynamic read energy per access, and standby leakage power of SRAM and logic-based embedded DRAM for 65nm technology. The area, dynamic read energy per access, and standby leakage power trends are split up into two charts based on capacities.

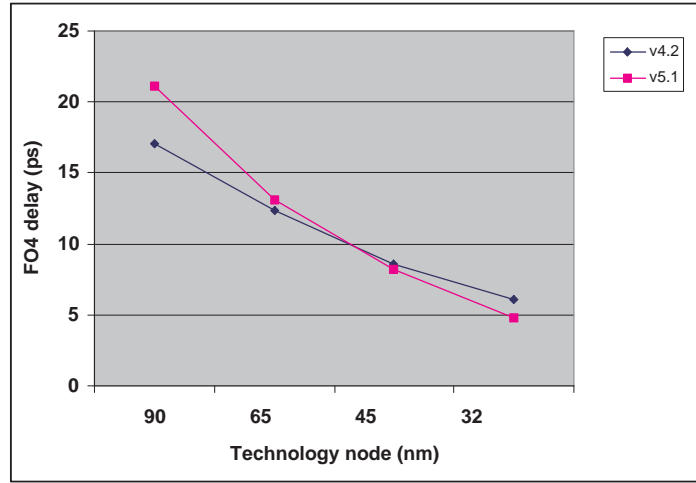


Figure 48: FO4 delays for various technology nodes in versions 4.2 and 5.1. The version 5.1 FO4 delays are for ITRS high-performance device type.

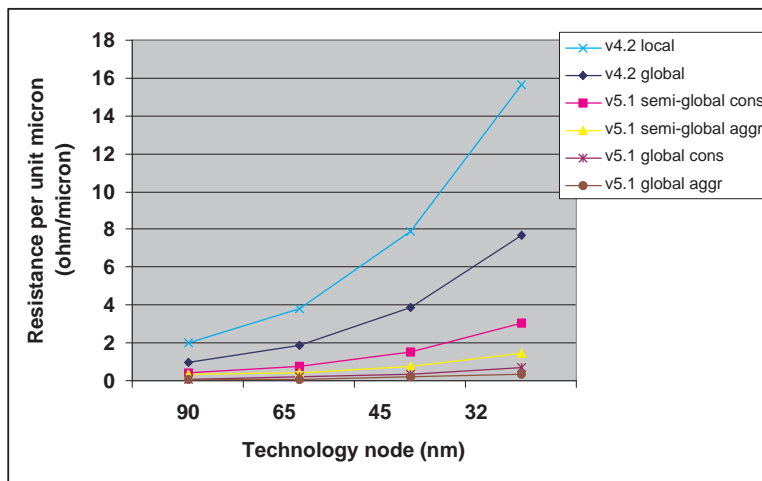


Figure 49: Resistance per unit micron of wire at various technology nodes for versions 4.2 and 5.1.

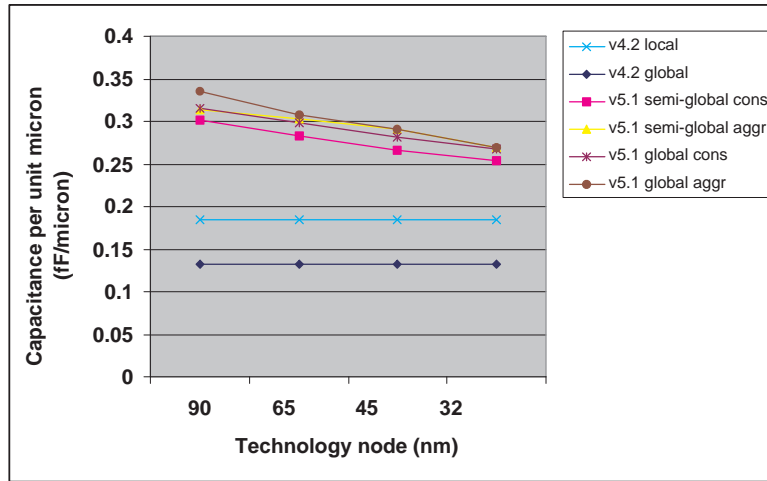


Figure 50: Capacitance per unit micron of wire at various technology nodes for versions 4.2 and 5.1.

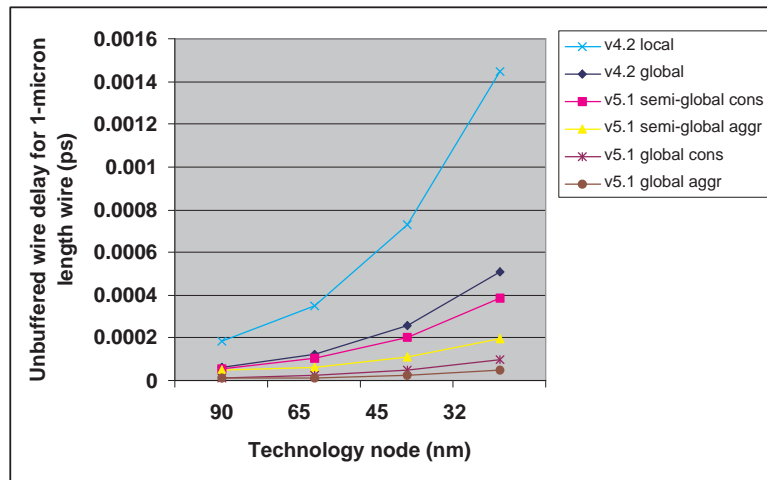


Figure 51: Unbuffered delays through a wire of length 1-micron for various technology nodes in versions 4.2 and 5.1. The wire delays for version 5.1 are for both aggressive and conservative projections.

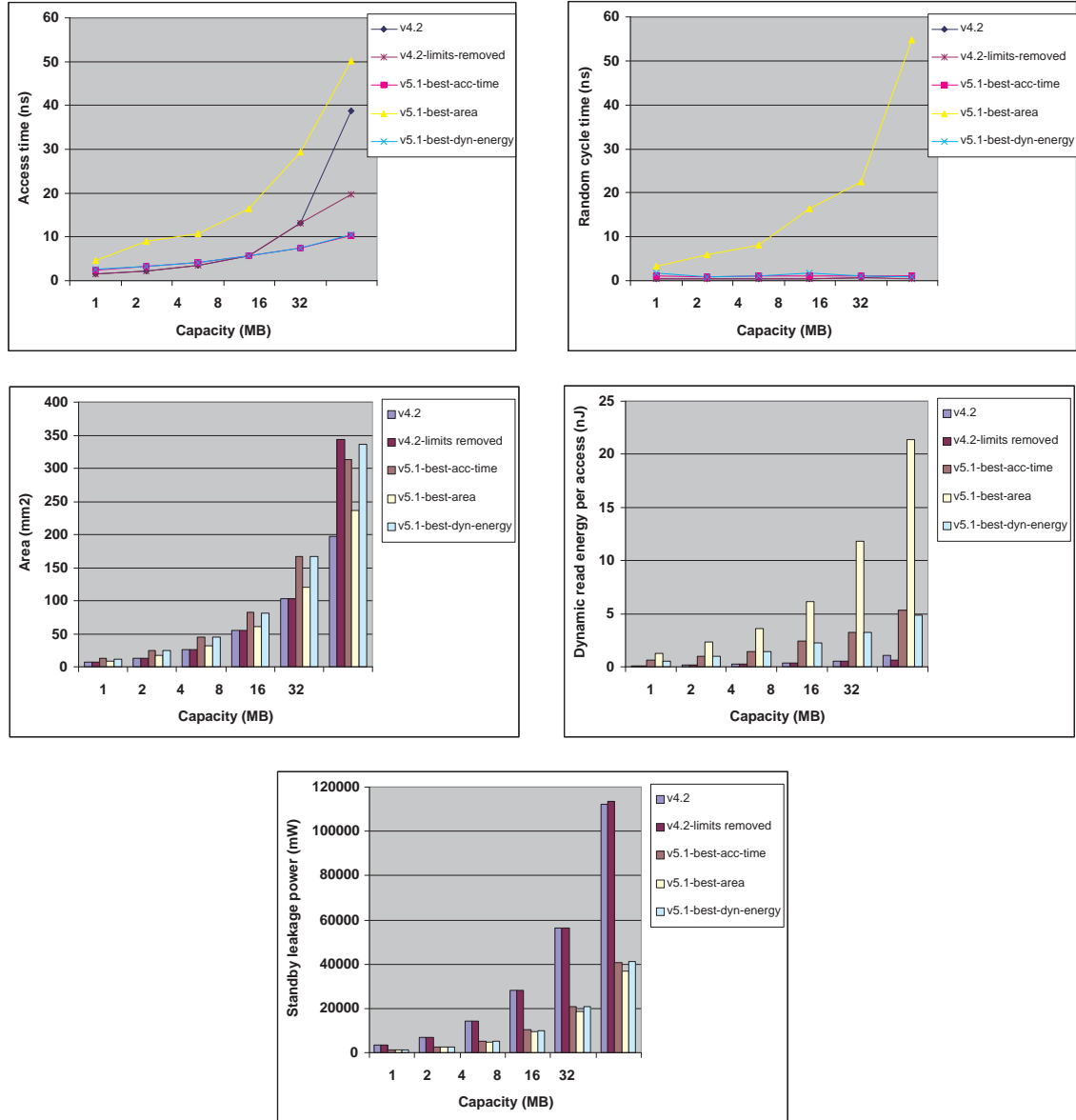


Figure 52: Access time, random cycle time, area, dynamic energy, and leakage power obtained from versions 4.2 and 5.1 for SRAMs in 70nm technology.

Parameter	Value
Area (mm ²)	128
Access time (ns)	5
Clock frequency (MHz)	800
Total power (W)	8

Table 12: Characteristics of Sun’s SPARC 90nm L2 cache.

12 Validation

In this section, we compare the published values of area, delay, and power of real cache designs with the projections of area, delay, and power produced by CACTI. The area, delay, and power of a real cache or RAM design are influenced by various factors. The design process inherently makes certain area, delay, and power tradeoffs based on budgets and requirements. Area, delay, and power are also influenced by design methodology, human bias and other practical considerations such as availability of IP from past designs etc. CACTI is based on generic assumptions of cache organization, circuits, design methodology, layout, design rules, and technology, whereas a real cache design is based on specific choices of all these. With CACTI 5, however, as was shown in the previous section, we provide a number of knobs that can be turned in order to try to emulate a real cache design in a better way. So it is interesting to see how the projections produced by CACTI would compare with real designs.

We use information from real cache specifications to fix as many of the input parameters required by CACTI as possible, such as capacity, associativity, line size, technology-node, etc. In order to understand and explore area/delay/power tradeoffs, we vary parameters such as “maximum percentage away from best area solution” and “maximum percentage away from best access time solution” within reasonable bounds.

12.1 Sun SPARC 90nm L2 cache

[25] describes the implementation of a 90nm SPARC 4MB L2 cache. Table 12 shows the area, access time, random cycle time, and power of the SPARC L2 cache. The clock frequency of the CPU core itself is 1.6 GHz but the L2 cache has a throughput of two clock cycles, so we fix the random cycle time of the L2 cache as 800 MHz.

Table 13 presents the input parameters used with CACTI to model this cache. From the description of the cache, we could not be sure whether the cache access mode is ‘normal’ or ‘fast’, so we try out both scenarios. In order to explore a range of area/delay/power tradeoffs, we vary `max_area_constraint` between 0 and 50 and `max_delay_constraint` between 0 and 30. In order to meet the aggressive random cycle time of the cache, we optimize the solution for random cycle time only. Because we do not have information about the interconnect properties of the fabrication process, we consider both aggressive and conservative projections for interconnect. Also, for ‘wire type outside mat’, we try out both ‘semi-global’ and ‘global’ wire types.

Figure 53 shows bubble charts showing access time, area, and power of the SPARC L2 cache and the various solutions generated by CACTI. The charts shown in Figure 53 are for CACTI solutions with ‘fast’ access mode, ‘conservative’ interconnect projections and ‘semi-global’ wire type outside mat. We believe that these values for the parameters are likely to be closest to the actual cache design. Results of the validation exercise with other configuration values are presented in the appendix. As we do not know the operating conditions corresponding to the published value for power of the SPARC L2, we compute dynamic power for the CACTI solutions for three activity factors – 0.1, 0.5, and 1.0. Also, we assume that the ratio of read to write accesses is 3. Note that the solutions shown in Figure 53 are the ones that can meet the random cycle time of the SPARC L2. It can be seen that many solutions have access time, area, and power that are quite similar to that of the L2 cache. Table 14 shows error percentages of prominent CACTI solutions with respect to the SPARC L2.

12.2 Intel Xeon 65nm L3 cache

Table 15 shows the area, access time, dynamic power, and leakage power of an Intel Xeon L3 cache in a 65nm process. [9] mentions that the access time of the cache is less than 9 ns; we assume the access time to be 9 ns. The clock

Parameter	Value
Capacity (MB)	4
Line size (bytes)	32
Associativity	4
Number of read/write ports	1
Number of exclusive ports	0
Number of banks	1
Technology-node (nm)	90
Output width (bits)	256
Specific tag	Yes
Tag width	34
Access mode	Normal/Fast
Pure RAM	No (cache)
DRAM	No
Repeaters in bank H-trees	Yes
max_area_constraint	0 - 70
max_acc_time_constraint	0 - 30
max_repeater_delay_constraint	10
Optimize for dynamic energy	No
Optimize for dynamic power	No
Optimize for leakage power	No
Optimize for cycle time	Yes
Temperature (K)	360
SRAM cell/wordline technology flavor	ITRS HP
Peripheral/Global circuitry technology flavor	ITRS HP
Interconnect projection type	Conservative
Wire type inside mat	Semi-global
Wire type outside mat	Semi-global/Global

Table 13: CACTI input parameters used for modeling 90nm SPARC L2 cache

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	8/8/8	-31/-31/-31	-67/-52/-34	35/30/24
Best % error in area	17/17/17	-33/-33/-33	-69/-52/-31	39/34/27
Best % error in power	17/8/8	-33/-23/-23	-69/-54/-36	39/28/22
Best average of area, acc time, and power % errors	8/8/8	-23/-23/-23	-68/-54/-36	33/28/22
Best average of area and acc time % errors	8/8/8	-23/-23/-23	-68/-54/-36	33/28/22
Best average of acc time and power % errors	8/8/8	-31/-31/-31	-67/-52/-34	35/30/24
Best acc time	8/8/8	-31/-31/-31	-67/-52/-34	35/30/24
Best area	17/17/17	-33/-33/-33	-69/-52/-31	39/34/27
Best power	17/8/8	-33/-23/-23	-69/-54/-36	39/28/22

Table 14: Error percentages of some prominent solutions generated by CACTI with respect to a 90nm SPARC L2 cache. The CACTI solutions assume ‘fast’ access mode, ‘conservative’ interconnect projections, and ‘semi-global’ wire type outside mat. We have used 3 activity factors of 0.1, 0.5, and 1, and so each entry in the table has 3 values.

frequency of the core itself is given to be 3.4 GHz and the clock frequency of the L3 is given to be half that of the core [33]. Also, because the output bus width of the L3 is 256 bits, it would require two cycles to transmit a 64-byte

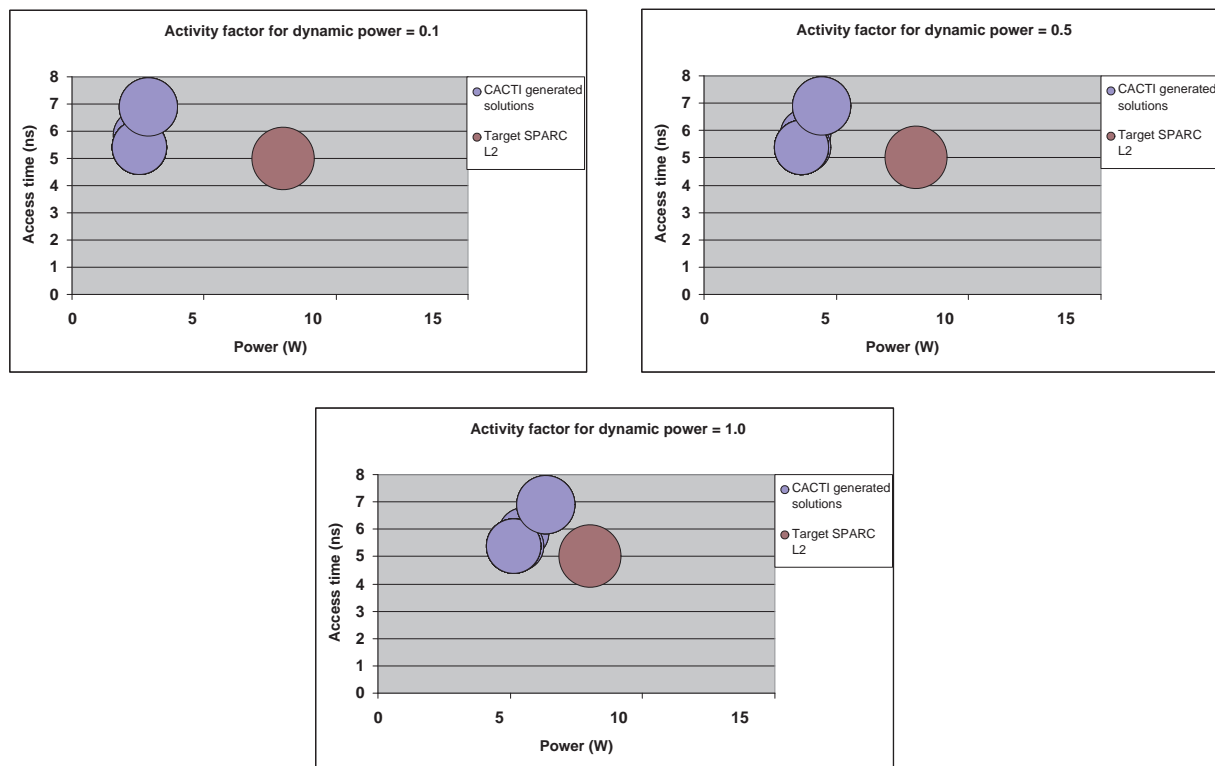


Figure 53: Access time, area, and power of the 90nm SPARC L2 cache and of solutions generated by CACTI. The CACTI solutions assume ‘fast’ access mode, ‘conservative’ interconnect projections, and ‘semi-global’ wire type outside mat. The 3 plots correspond to 3 activity factors assumed while computing dynamic power for the CACTI-generated solutions.

Area (mm^2)	200	Measured from die photo [32]
Access time (ns)	9 ns	[9]
Clock frequency (GHz)	850 MHz	[33]
Dynamic power (W)	1.7/5.4	[7, 33, 44]
Leakage power (W)	6.6	[33, 44]

Table 15: Characteristics of Intel Xeon’s 65nm L3 cache.

line, so we fix the random cycle frequency of the L3 to be one-fourth that of the CPU, i.e. 850 MHz. The dynamic power of the cache comes out to be 5.4W based on information from [33, 44], however [7] mentions that the cache consumes about 1.7W for “average applications”. We speculate that these differences in dynamic power numbers that have been quoted are because of different activity factors in the cache caused due to measurements or simulations of applications with different characteristics. While carrying out comparisons of the area, delay, and power of the Intel cache with those of the solutions generated by CACTI, we use both values of power.

The 65nm process offers transistors with 35nm gate-lengths. The cache itself, however, makes use of longer-channel devices with lengths that are about 10% longer than the nominal. The longer-channel devices have on-currents that are about 10% less than the nominal devices but have leakage that is lower by a factor of 3. The cache operates in a voltage domain different from that of the cores. The cores can operate at 1.25V while the cache operates at 1.1V.

In order to control leakage in the cache, the Intel cache implements n and p sleep transistors at the level of ‘blocks’ within subarrays (Each subarray within the Intel cache is composed of multiple ‘blocks’ with one block within a subarray activated per access. The ‘subarray’ of the Intel cache is not the same as that of CACTI). The impact of these sleep transistors is that leakage power in all blocks that are not activated during an access is cut down by half.

Parameter	Value
Capacity (MB)	16
Line size (bytes)	64
Associativity	16
Number of read/write ports	1
Number of exclusive ports	0
Number of banks	2
Technology-node (nm)	65
Output width (bits)	512
Specific tag	No
Access mode	Serial
Pure RAM	No (cache)
DRAM	No
max_area_constraint	0 - 50
max_delay_constraint	0 - 30
max_repeater_delay_constraint	10
Optimize for dynamic energy	No
Optimize for dynamic power	No
Optimize for leakage power	No
Optimize for cycle time	Yes
Temperature (K)	360
SRAM cell/wordline technology flavor	ITRS HP
Peripheral/Global circuitry technology flavor	ITRS HP
Interconnect projection type	Conservative
Wire type inside mat	Semi-global
Wire type outside mat	Semi-global/Global

Table 16: CACTI input parameters used for modeling 65nm Intel Xeon L3 cache.

Table 16 shows the input parameters used with CACTI to model the Intel L3 cache. In order to compare the power numbers produced by CACTI with those of the Intel cache in a fair manner, we assume the use of leakage control mechanisms within CACTI similar to that used in the Intel cache. To model the longer-channel devices that have been used in the Intel cache which reduce leakage by a factor of 3, we also reduce the leakage of the CACTI 65 nm high-performance transistors by a factor of 3. Also, we assume the use of sleep transistors that cut down the leakage of all mats that are not activated during an access by half.

Figure 54 shows bubble charts of access time, area, and power of the Intel cache and the various solutions generated by CACTI. The charts shown in Figure 54 are for CACTI solutions with ‘conservative’ interconnect projections and ‘semi-global’ wire type outside mat as we believe that these values for the parameters are likely to be closest to the actual cache design. Results of the validation exercise with other configuration values are again presented in the appendix. Again, as we do not know the operating conditions corresponding to the published value for power, we compute dynamic power for the CACTI solutions for three activity factors – 0.1, 0.5, and 1.0. Again, we assume that the ratio of read to write accesses is 3. There are two targets for the Intel cache corresponding to the two values of dynamic power shown in Table 15. It can be seen from Figure 54 that many CACTI solutions have area, access time and power that are quite similar to that of the Xeon L3. Tables 17 and 18 show error percentages of prominent CACTI solutions with respect to the Xeon L3.

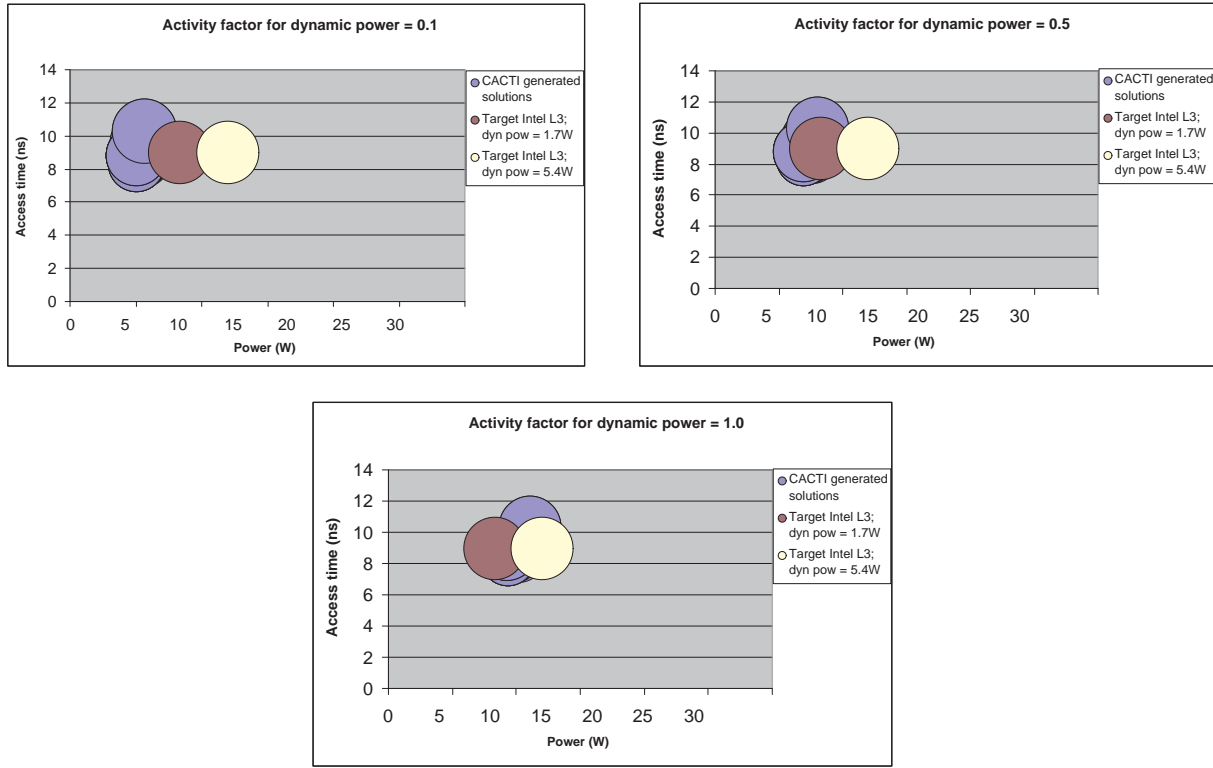


Figure 54: Access time, area, and power of the 65nm Xeon L3 cache and of solutions generated by CACTI. The CACTI solutions assume ‘conservative’ interconnect projections and ‘semi-global’ wire type outside mat. The 3 plots correspond to 3 activity factors assumed while computing dynamic power for the CACTI-generated solutions.

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	-7/-7/-7	-23/-23/-23	-39/-16/12	23/15/14
Best % error in area	6/6/6	-25/-25/-25	-32/-3/35	21/11/22
Best % error in power	-2/-2/-2	-6/-6/-6	-39/-17/12	16/8/6
Best average of area, acc time, and power % errors	-2/14/-2	-6/2/-6	-39/-3/12	16/6/6
Best average of area and acc time % errors	-2/-2/-2	-6/-6/-6	-39/-17/12	16/8/6
Best average of acc time and power % errors	6/6/-2	-25/-25/-6	-32/-3/12	21/11/6
Best acc time	-7/-7/-7	-23/-23/-23	-39/-16/12	23/15/14
Best area	6/6/6	-25/-25/-25	-32/-3/35	21/11/22
Best power	-2/-2/-2	-6/-6/-6	-39/-17/12	16/8/6

Table 17: Error percentages of some prominent solutions generated by CACTI with respect to a 65nm Intel Xeon L3 cache when we assume that the dynamic power consumed by the cache is 1.7W. The CACTI solutions assume ‘conservative’ interconnect projections and ‘semi-global’ wire type outside mat. We have used 3 activity factors of 0.1, 0.5, and 1, and so each entry in the table has 3 values.

13 Commodity DRAM Technology and Main Memory Chip Modeling

As an additional enhancement to CACTI we are adding support for modeling commodity DRAM technology and main memory chip organization and operation.

Commodity DRAM technology offers a cell that is much denser than that offered by logic process based DRAM

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	-7/-7/-7	-23/-23/-23	-58/-42/-22	29/24/18
Best % error in area	6/6/6	-25/-25/-25	-53/-33/-7	28/21/13
Best % error in power	-2/-2/-2	-6/-6/-6	-58/-42/-23	22/17/10
Best average of area, acc time, and power % errors	-2/14/14	-6/2/2	-58/-33/-8	22/16/8
Best average of area and acc time % errors	-2/-2/-2	-6/-6/-6	-58/-42/-23	22/17/10
Best average of acc time and power % errors	-5/6/6	-20/-25/-25	-54/-33/-7	26/21/13
Best acc time	-7/-7/-7	-23/-23/-23	-58/-42/-22	29/24/18
Best area	6/6/6	-25/-25/-25	-53/-33/-7	28/21/13
Best power	-2/-2/-2	-6/-6/-6	-58/-42/-23	22/17/10

Table 18: Error percentages of some prominent solutions generated by CACTI with respect to a 65nm Intel Xeon L3 cache when we assume that the dynamic power consumed by the cache is 5.4W. The CACTI solutions assume ‘conservative’ interconnect projections and ‘semi-global’ wire type outside mat. We have used 3 activity factors of 0.1, 0.5, and 1, and so each entry in the table has 3 values.

technology (the typical cell area of a modern commodity DRAM technology is $6F^2$). Commodity DRAM technology is also characterized by a higher refresh period (64ms for commercial applications). However commodity DRAM technology has transistors that are slower than those found in a logic process. We are incorporating technology data for commodity DRAM technology into CACTI.

The organization and operation of a main memory chip is different from an embedded memory. Main memory is usually organized as DIMMs which are typically composed of 8 or 16 main memory chips. Each main memory chip on a DIMM has 4 or 8 data output pins which is much lower than what is typical of embedded memories. Because of their limited number of data pins, main memory chips operate in burst mode with a *burst length* of 4 or 8. Modern main memory chips incorporate the concept of *internal prefetch width* which determines the number of bits that are prefetched internally inside the DRAM core. For DDR, the internal prefetch width is 2, for DDR2 it’s 4, for DDR3 it’s 8, and so on. Main memory chips are also characterized by *page size* which is equal to the number of sense amplifiers that are activated per access. Because of their limited number of address pins, a main memory chip also operates differently from an embedded memory. In an embedded memory all the address bits may be decoded in parallel. In a main memory chip, row and column address are multiplexed on the address bus. The row address is first applied and is used to latch data from a row into the sense amplifiers. After the data gets latched, the column address is then applied and the read or write access is carried out. We are incorporating the impact of these organizational and operational features of a main memory DRAM chip into the CACTI models.

With these new technology parameters and organizational and operational models users will be able to obtain area, delay, and power projections for memories and caches that are based on commodity DRAM technology in addition to those based on SRAM and logic process based DRAM technologies. Further details about support for commodity DRAM technology and main memory modeling in CACTI may be found in [43]. We expect the source code corresponding to this enhancement will be released in CACTI version 5.2.

14 Future Work

Non Uniform Cache Access (NUCA) is an interesting architecture that CACTI could support in the future. Incorporation of models for low-swing interconnect into CACTI could also be an interesting enhancement.

15 Conclusions

In this technical report, we have described the various enhancements carried out in CACTI 5 while also providing a comprehensive overview of the CACTI area, access time and power modeling. CACTI 5 includes a number of major improvements over CACTI 4.0. The base technology modeling has been changed from simple linear scaling of the original 0.8 micron technology to models based on the ITRS roadmap. Data for different ITRS device types has been incorporated into CACTI. Interconnect technology data has also been updated so that it is now based off well-documented models and data. CACTI 5 has also added support for modeling embedded DRAM in such a way that it now becomes possible to compare tradeoffs involving the use of embedded SRAM or DRAM for identical input specifications. This has been achieved by an extensive rewrite of the CACTI code base. Various organizational and circuit assumptions have been clarified and updated. The modeling has also been restructured in such a way that it is now more modular and easier to extend and evolve.

In the studies shown in this report, the impact of technology assumptions on cache performance has been pointed out and emphasized. The importance of solution optimization techniques on memory and cache performance has also been highlighted. Area, delay, and power results obtained from version 5.1 have been compared against published data available for two prominent caches. Taking into account the extremely generic nature of CACTI, it was found that there is reasonable agreement between the results produced by CACTI and the published data.

Finally, as in the original CACTI report, we would like to caution users against making too many conclusions based on results shown in this report. It is important to know that CACTI is a simplified model for memories and caches with various limitations at the various levels of modeling, so appropriate caution and judgement needs to be exercised with its use. In general, it is best to use CACTI for studies that involve relative optimization.

Acknowledgments

We would like to thank the following people for many helpful discussions and comments: Al Davis from the University of Utah, and Jay Brockman and Christian Poellabauer from the University of Notre Dame.

A Additional CACTI Validation Results for 90nm SPARC L2

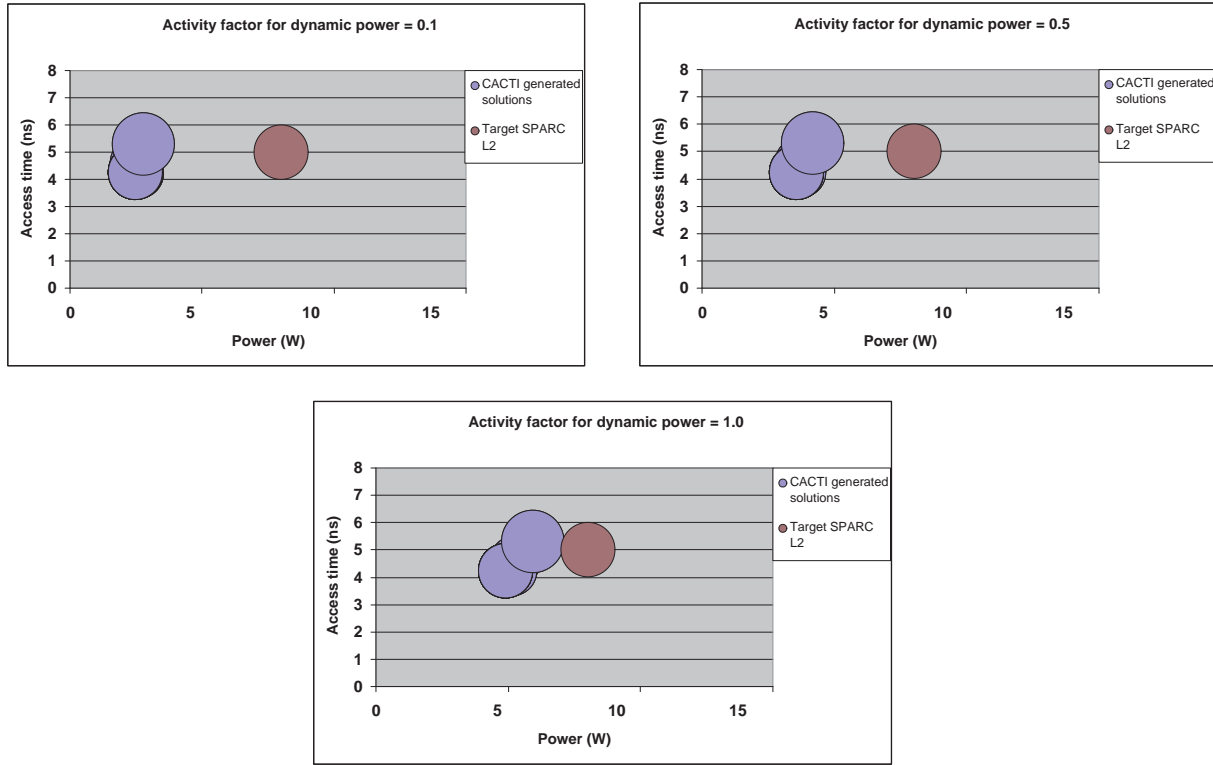


Figure 55: Access time, area and power of the 90nm SPARC L2 cache and of solutions generated by CACTI. The CACTI solutions are for assumptions of ‘fast’ access mode, ‘conservative’ interconnect projections and ‘global’ wire type outside mat. The 3 plots correspond to 3 activity factors assumed while computing dynamic power for the CACTI-generated solutions.

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	-15/-15/-15	-13/-13/-13	-68/-54/-36	32/27/21
Best % error in area	-8/-8/-8	-16/-16/-16	-69/-54/-34	31/26/20
Best % error in power	-8/-15/-15	-16/1/1	-69/-56/-39	31/24/18
Best average of area, acc time, and power % errors	-15/-15/-15	1/1/1	-69/-56/-39	28/24/18
Best average of area and acc time % errors	-15/-15/-15	1/1/1	-69/-56/-39	28/24/18
Best average of acc time and power % errors	6/6/6	26/26/26	-65/-48/-26	32/27/19
Best acc time	-15/-15/-15	-13/-13/-13	-68/-54/-36	32/27/21
Best area	-8/-8/-8	-16/-16/-16	-69/-54/-34	31/26/20
Best power	-8/-15/-15	-16/1/1	-69/-56/-39	31/24/18

Table 19: Error percentages of some prominent solutions generated by CACTI with respect to a 90nm SPARC L2 cache under assumptions of ‘fast’ access mode, ‘conservative’ interconnect projections and ‘global’ wire type outside mat. We have used 3 activity factors of 0.1, 0.5 and 1, and so each entry in the table has 3 values.

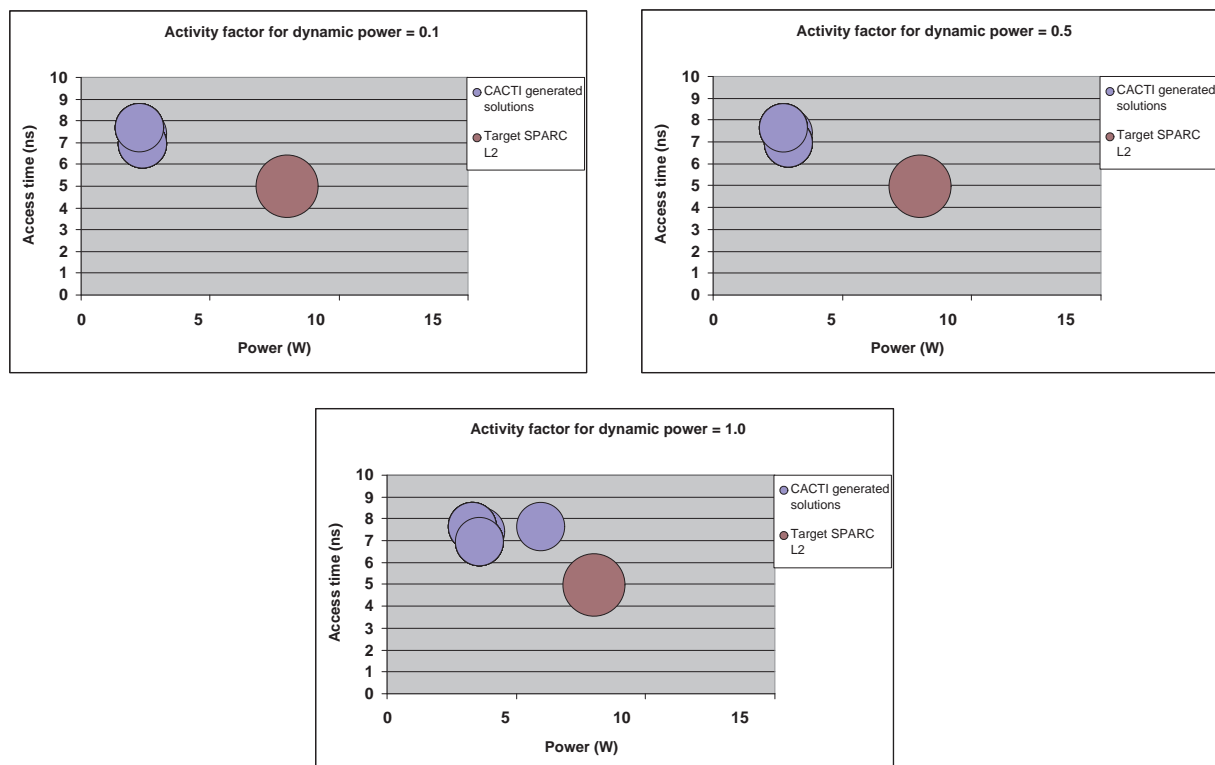


Figure 56: Access time, area and power of the 90nm SPARC L2 cache and of solutions generated by CACTI. The CACTI solutions are for assumptions of ‘normal’ access mode, ‘conservative’ interconnect projections and ‘semi-global’ wire type outside mat. The 3 plots correspond to 3 activity factors assumed while computing dynamic power for the CACTI-generated solutions.

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	39/39/39	-39/-39/-39	-70/-64/-56	49/47/45
Best % error in area	39/39/39	-39/-39/-39	-70/-64/-56	49/47/45
Best % error in power	53/53/53	-37/-37/-37	-72/-66/-59	54/52/50
Best average of area, acc time and power % errors	39/39/39	-39/-39/-39	-70/-64/-56	49/47/45
Best average of area and acc time % errors	39/39/39	-39/-39/-39	-70/-64/-56	49/47/45
Best average of acc time and power % errors	39/39/39	-39/-39/-39	-70/-64/-56	49/47/45
Best acc time	39/39/39	-39/-39/-39	-70/-64/-56	49/47/45
Best area	39/39/39	-39/-39/-39	-70/-64/-56	49/47/45
Best power	53/53/53	-37/-37/-37	-72/-66/-59	54/52/50

Table 20: Error percentages of some prominent solutions generated by CACTI with respect to a 90nm SPARC L2 cache for a normal cache under assumptions of ‘normal’ access mode, ‘conservative’ interconnect projections and ‘semi-global’ wire type outside mat. The CACTI solutions and errors depends on the activity factor. We have used 3 activity factors of 0.1, 0.5 and 1, and so each entry in the table has 3 values.

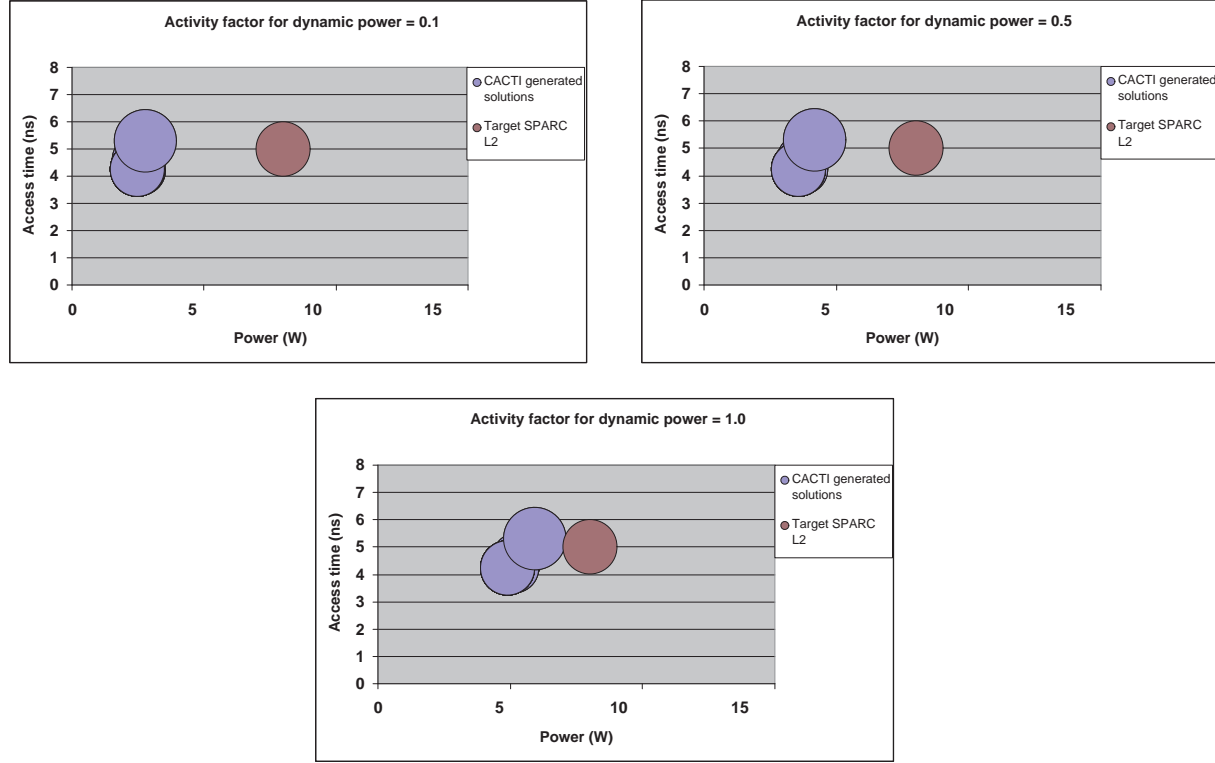


Figure 57: Access time, area and power of the 90nm SPARC L2 cache and of solutions generated by CACTI. The CACTI solutions are for assumptions of ‘normal’ access mode, ‘conservative’ interconnect projections and ‘global’ wire type outside mat. The 3 plots correspond to 3 activity factors assumed while computing dynamic power for the CACTI-generated solutions.

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	12/12/12	-31/-31/-31	-70/-64/-56	38/36/33
Best % error in area	23/23/23	-33/-33/-33	-72/-66/-57	43/41/38
Best % error in power	23/22/22	-33/-27/-27	-72/-67/-60	43/39/37
Best average of area, acc time and power % errors	12/12/12	-31/-31/-31	-70/-64/-56	38/36/33
Best average of area and acc time % errors	12/12/12	-31/-31/-31	-70/-64/-56	38/36/33
Best average of acc time and power % errors	12/12/12	-31/-31/-31	-70/-64/-56	38/36/33
Best acc time	12/12/12	-31/-31/-31	-70/-64/-56	38/36/33
Best area	23/23/23	-33/-33/-33	-72/-66/-57	43/41/38
Best power	23/22/22	-33/-27/-27	-72/-67/-60	43/39/37

Table 21: Error percentages of some prominent solutions generated by CACTI with respect to a 90nm SPARC L2 cache for a normal cache under assumptions of ‘normal’ access mode, ‘conservative’ interconnect projections and ‘global’ wire type outside mat. We have used 3 activity factors of 0.1, 0.5 and 1, and so each entry in the table has 3 values.

B Additional CACTI Validation Results for 65 nm Xeon L3

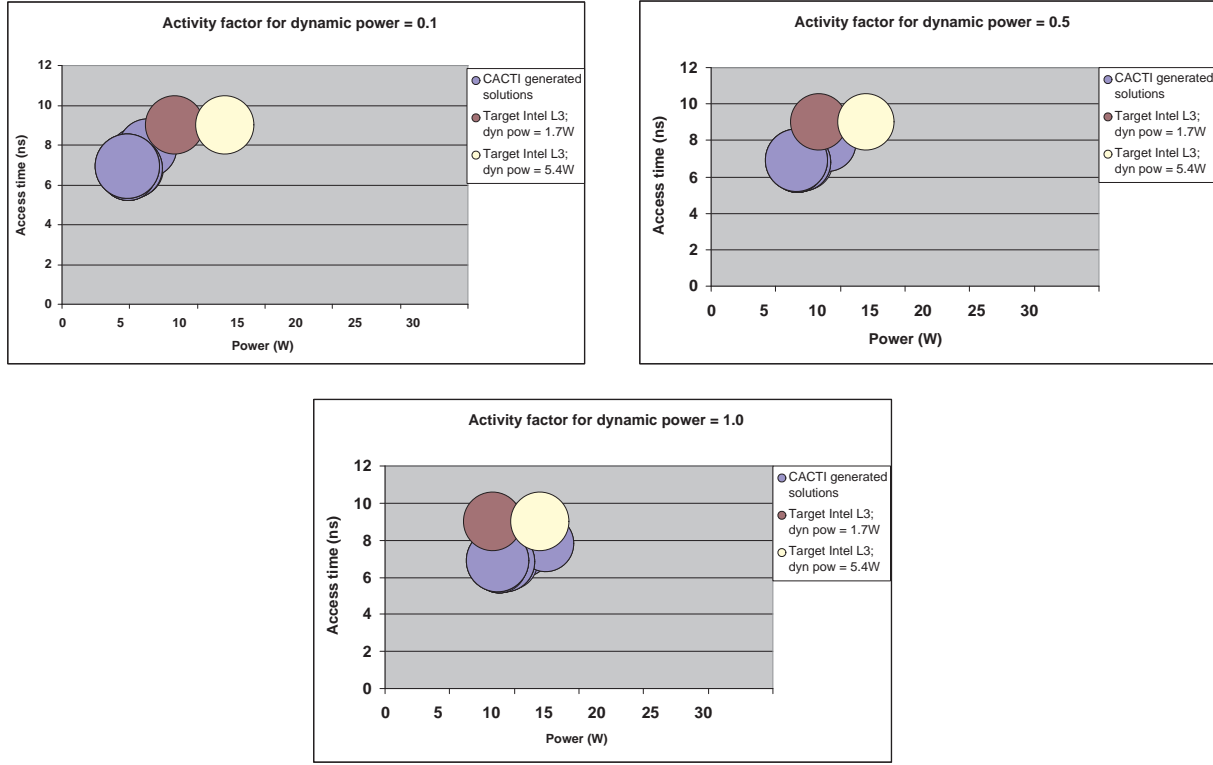


Figure 58: Access time, area and power of the 65 nm Xeon L3 cache and of solutions generated by CACTI. The CACTI solutions are for assumptions of ‘conservative’ interconnect projections and ‘global’ wire type outside mat. The 3 plots correspond to 3 activity factors assumed while computing dynamic power for the CACTI-generated solutions.

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	-26/-26/-26	-7/-7/-7	-41/-44/7	25/26/13
Best % error in area	-17/-17/-17	-11/-11/-11	-34/-35/30	20/21/19
Best % error in power	-23/-23/-23	18/18/18	-41/-45/5	28/29/15
Best average of area, acc time, and power % errors	-13/-13/-26	-1/-1/-7	-24/-25/7	13/13/13
Best average of area and acc time % errors	-13/-13/-13	-1/-1/-1	-24/-25/49	13/13/21
Best average of acc time and power % errors	-13/-13/-23	-1/-1/18	-24/-25/5	13/13/15
Best acc time	-26/-26/-26	-7/-7/-7	-41/-44/7	25/26/13
Best area	-17/-17/-17	-11/-11/-11	-34/-35/30	20/21/19
Best power	-23/-23/-23	18/18/18	-41/-45/5	28/29/15

Table 22: Error percentages of some prominent solutions generated by CACTI with respect to a 65nm Intel Xeon L3 cache when we assume that the dynamic power consumed by the cache is 1.7W. The CACTI solutions are for assumptions of ‘conservative’ interconnect projections and ‘global’ wire type outside mat. We have used 3 activity factors of 0.1, 0.5 and 1, and so each entry in the table has 3 values.

Solution	% error in acc time	% error in area	% error in power	Avg of acc time, area and power % errors
Best % error in acc time	-26/-26/-26	-7/-7/-7	-59/-44/-26	31/26/20
Best % error in area	-17/-17/-17	-11/-11/-11	-54/-35/-10	27/21/13
Best % error in power	-23/-23/-23	18/18/18	-60/-45/-27	34/29/23
Best average of area, acc time and power % errors	-13/-13/-13	-1/-1/-1	-48/-25/3	21/13/6
Best average of area and acc time % errors	-13/-13/-13	-1/-1/-1	-48/-25/3	21/13/6
Best average of acc time and power % errors	-13/-13/-13	-1/-1/-1	-48/-25/3	21/13/6
Best acc time	-26/-26/-26	-7/-7/-7	-59/-44/-26	31/26/20
Best area	-17/-17/-17	-11/-11/-11	-54/-35/-10	27/21/13
Best power	-23/-23/-23	18/18/18	-60/-45/-27	34/29/23

Table 23: Error percentages of some prominent solutions generated by CACTI with respect to a 65nm Intel Xeon L3 cache when we assume that the dynamic power consumed by the cache is 5.4W. The CACTI solutions are for assumptions of ‘conservative’ interconnect projections and ‘global’ wire type outside mat. We have used 3 activity factors of 0.1, 0.5 and 1, and so each entry in the table has 3 values.

References

- [1] S. Alexenian, "Multi-level Semiconductor Memory Architecture and Method of Forming the Same," *US Patent 20050041513*, Feb 2005.
- [2] B. S. Amrutur and M. A. Horowitz, "Speed and Power Scaling of SRAM's," *IEEE Journal of Solid-State Circuits*, vol. 35, no. 2, pp. 175–185, Feb 2000.
- [3] —, "Fast Low Power Decoders for RAMs," *IEEE Journal of Solid-State Circuits*, vol. 36, no. 10, pp. 1506–1515, Oct 2001.
- [4] P. Bai, *et al.*, "A 65nm Logic Technology Featuring 35nm Gate Lengths, Enhanced Channel Strain, 8 Cu Interconnect Layers, Low-k ILD and 0.57 micron² SRAM cell," in *IEEE International Electron Devices Meeting*, Dec 2004.
- [5] J. Barth, J. Dreibelbis, E. Nelson, D. Anand, G. Pomichter, P. Jakobsem, M. Nelms, J. Leach, and G. Belansek, "Embedded DRAM design and architecture for the IBM 0.11-micron ASIC offering," *IBM Journal of Research and Development*, vol. 46, no. 6, Nov 2002.
- [6] J. Barth, *et al.*, "A 500MHz random cycle 1.5ns-latency SOI embedded DRAM macro featuring a 3T micro sense amplifier," in *International Solid-State Circuits Conference*, San Francisco, CA, Feb 2007.
- [7] J. Chang, *et al.*, "The 65-nm 16-MB Shared On-Die L3 Cache for the Dual-Core Intel Xeon Processor 7100 Series," *IEEE Journal of Solid-State Circuits*, vol. 42, no. 4, Apr 2007.
- [8] L. Chang, *et al.*, "Stable SRAM cell design for the 32 nm node and beyond," in *Symposium on VLSI*, 2005.
- [9] J. Gilbert, S. Hunt, D. Gunadi, and G. Srinivasa, "TULSA, A Dual P4 Core Large Shared Cache Intel Xeon Processor for the MP Server Market Segment," in *Hot Chips*, Aug 2006.
- [10] A. Hajimiri and R. Heald, "Design Issues in Cross-coupled Sense Amplifier," in *International Symposium on Circuits and Systems*, May/Jun 1998, pp. 149–152.
- [11] R. Ho, "On-chip Wires: Scaling and Efficiency," Ph.D. dissertation, Stanford University, 2003.
- [12] T. Hook, M. Breitwisch, J. Brown, P. Cottrell, D. Hoyniak, C. Lam, and R. Mann, "Noise margin and leakage in ultra-low leakage SRAM cell design," *IEEE Transactions on Electron Devices*, vol. 49, no. 8, Aug 2002.
- [13] M. Horowitz, R. Ho, and K. Mai, "The Future of Wires," In Invited Workshop Paper for SRC Conference., Tech. Rep., 1999.
- [14] S. S. Iyer, J. J. E. Barth, P. C. Parries, J. P. Norum, J. P. Rice, L. R. Logan, and D. Hoyniak, "Embedded DRAM: Technology platform for the Blue Gene/L chip," *IBM Journal of Research and Development*, vol. 49, no. 2/3, pp. 333–350, Mar/May 2005.
- [15] J. J. Barth, *et al.*, "A 500-MHz multi-banked compilable DRAM macro with direct write and programmable pipelining," *IEEE Journal of Solid-State Circuits*, vol. 40, no. 1, Jan 2005.
- [16] S. Jin, J. Yi, J. Choi, D. Kang, Y. Park, and H. Min, "Prediction of data retention time distribution of DRAM by physics-based statistical simulation," *IEEE Transactions on Electron Devices*, vol. 52, no. 11, Nov 2005.
- [17] B. Keeth and R. Baker, *DRAM Circuit Design: A Tutorial*. IEEE Press, 2000.
- [18] T. Kirihaata, *et al.*, "An 800-MHz embedded DRAM with a concurrent refresh mode," *IEEE Journal of Solid-State Circuits*, vol. 40, no. 6, Jun 2005.
- [19] V. Klee, *et al.*, "A 0.13-micron logic-based embedded DRAM technology with electrical fuses, Cu interconnect in SiLKTM, sub-7ns random access time and its extension to the 0.10-micron generation," in *IEEE International Electron Devices Meeting*, 2001.
- [20] M. Mamidipaka and N. Dutt, "eCACTI: An Enhanced Power Estimation Model for On-chip Caches," Center for Embedded Computer Systems (CECS), Tech. Rep. TR-04-28, 2004.
- [21] M. Mamidipaka, K. Khouri, N. Dutt, and M. Abadir, "Analytical Models for Leakage Power Estimation of Memory Array Structures," in *Proceedings of the 2nd IEEE/ACM/IFIP international conference on Hardware/software codesign and system synthesis*, Sep 2004.
- [22] R. Mann, *et al.*, "Ultralow-power SRAM technology," *IBM Journal of Research and Development*, vol. 47, no. 5/6, Sep/Nov 2003.
- [23] R. Matick and S. Schuster, "Logic-based eDRAM: Origins and rationale for use," *IBM Journal of Research and Development*, vol. 49, no. 1, Jan 2005.
- [24] Y. Matsubara, *et al.*, "Fully integration of high density embedded DRAM with 65 nm CMOS technology (CMOS5)," in *IEEE International Electron Devices Meeting*, 2003.

- [25] H. McIntyre, *et al.*, “A 4-MB On-Chip L2 Cache for a 90-nm 1.6-GHz 64-bit Microprocessor,” *IEEE Journal of Solid-State Circuits*, vol. 40, no. 1, Jan 2005.
- [26] M. Na, E. Nowak, W. Haensch, and J. Cai, “The effective drive current of CMOS inverters,” in *IEEE International Electron Devices Meeting*, Dec 2002.
- [27] K. Noh, Y. Choi, J. Joo, M. Kim, J. Jung, J. Lim, C. Lee, G. Kim, and M. Kim, “A 130nm 1.1V 143 MHz SRAM-like embedded DRAM compiler with dual asymmetric bit line sensing scheme and quiet unselected IO scheme,” in *Symposium on VLSI Circuits*, 2004.
- [28] M. Oka and M. Suzuoki, “Designing and Programming the Emotion Engine,” *IEEE Micro*, vol. 19, no. 6, pp. 20–28, 1999.
- [29] J. M. Rabaey, A. Chandrakasan, and B. Nikolic, *Digital Integrated Circuits*, 2nd ed. Prentice Hall, 2003.
- [30] G. Reinman and N. P. Jouppi, “CACTI 2.0: An Integrated Cache Timing and Power Model,” DEC WRL, Tech. Rep. technical report number 2000/7, 2000.
- [31] Ron Ho, “Tutorial: Dealing with issues in VLSI interconnect scaling,” in *International Solid-State Circuits Conference*, 2007.
- [32] S. Rusu, S. Tam, H. Muljono, D. Ayers, and J. Chang, “A Dual-Core Multi-Threaded Xeon Processor with 16MB L3 Cache,” in *International Solid-State Circuits Conference*, San Francisco, CA, Feb. 2006.
- [33] S. Rusu, *et al.*, “A 65-nm Dual-Core Multithreaded Xeon Processor With 16-MB L3 Cache,” *IEEE Journal of Solid-State Circuits*, vol. 42, no. 1, Jan 2007.
- [34] P. Saxena, N. Menezes, P. Cocchini, and D. Kirkpatrick, “The scaling challenge: can correct-by-construction help?” in *International Symposium on Physical Design*, Apr 2003.
- [35] Semiconductor Industries Association, “International Technology Roadmap for Semiconductors,” 2005, <http://www.itrs.net/>.
- [36] —, “Model for Assessment of CMOS Technologies and Roadmaps (MASTAR),” 2005, <http://www.itrs.net/models.html>.
- [37] T. M. Service, “<http://www.mosis.org>.”
- [38] P. Shivakumar and N. P. Jouppi, “CACTI 3.0: An Integrated Cache Timing, Power, and Area Model,” Hewlett Packard Laboratories, Tech. Rep. WRL-2001-2, 2001.
- [39] I. E. Sutherland, R. F. Sproull, and D. Harris, *Logical Effort: Designing Fast CMOS Circuits*, 1st ed. San Mateo, CA: Morgan Kaufmann, 1999.
- [40] D. Tarjan, S. Thoziyoor, and N. P. Jouppi, “CACTI 4.0,” HP Labs, Tech. Rep. HPL-2006-86, 2006.
- [41] J. M. Tendler, J. S. Dodson, J. J. S. Fields, H. Le, and B. Sinharoy, “POWER4 system microarchitecture,” *IBM Journal of Research and Development*, vol. 46, no. 1, pp. 5–26, Jan. 2002.
- [42] S. Thompson, M. Alavi, M. Hussein, P. Jacob, C. Kenyon, P. Moon, M. Prince, S. Sivakumar, S. Tyagi, and M. Bohr, “130nm Logic Technology Featuring 60nm Transistors, Low-K Dielectrics, and Cu Interconnects,” *Intel Technical Journal*, vol. 6, no. 2, May 2002.
- [43] S. Thoziyoor, Ph.D. dissertation, University of Notre Dame, 2008.
- [44] R. Varada, M. Sriram, K. Chou, and J. Guzzo, “Design and Integration Methods for a Multi-threaded Dual Core 65nm Xeon Processor,” in *International Conference on Computer-Aided Design*, Nov 2006.
- [45] G. Wang, P. Parries, B. Khan, J. Liu, Y. Otani, J. Norum, N. Robson, T. Kirihaata, and S. Iyer, “A 0.168/0.11 micron² highly scalable high performance embedded DRAM cell for 90/65 nm logic applications,” in *Symposium on VLSI Circuits*, 2004.
- [46] S. Wilton and N. Jouppi, “CACTI: An enhanced cache access and cycle time model,” *IEEE Journal of Solid-State Circuits*, vol. 31, no. 5, pp. 677–688, May 1996.
- [47] S. Wilton and N. P. Jouppi, “An Enhanced Access and Cycle Time Model for On-Chip Caches,” DEC WRL, Tech. Rep. technical report number 93/5, 1994.
- [48] P. Yeh, D. Nayak, and D. Gitlin, “Improved CV/I methodology to accurately predict CMOS technology performance,” *IEEE Transactions on Electron Devices*, vol. 54, no. 7, Jul 2007.
- [49] H. Yoshida, K. De, and V. Boppana, “Accurate pre-layout estimation of standard cell characteristics,” in *Proceedings of the 41st annual conference on Design automation*. New York, NY, USA: ACM Press, 2004, pp. 208–211.
- [50] W. Zhao and Y. Cao, “New Generation of Predictive Technology Model for Sub-45nm Design Exploration,” in *IEEE International Symposium on Quality Electronic Design*, Mar 2006.