# An SRAM Compiler for Monolithic-3-D Integrated Circuit With Carbon Nanotube Transistors

**DAEHYUN KIM** (Graduate Student Member, IEEE),
**EDWARD LEE** (Student Member, IEEE), **JAMIN SEO** (Student Member, IEEE),
**JINWOO KIM** (Graduate Student Member, IEEE), **SUNG KYU LIM** (Senior Member, IEEE),
and **SAIBAL MUKHOPADHYAY** (Fellow, IEEE)

Department of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA 30332 USA

CORRESPONDING AUTHOR: D. KIM (daehyun.kim@gatech.edu)

**ABSTRACT**   This article presents monolithic-3-D (M3D) SRAM arrays using multiple tiers of carbon nanotube (CNT) transistors. The compiler automatically generates single-tier 2-D SRAM subarrays and multitier 3-D SRAM subarrays with different tiers for cells and peripheral logic. Moreover, the compiler can integrate multiple subarrays of different dimensions to generate larger capacity SRAM arrays. The compiler is demonstrated in a commercial-grade M3D process design kit (PDK) with two tiers of carbon nanotube transistors (CNFETs). Simulations show that the M3D CNT SRAM design can improve the properties of memory compared to the 2-D CNT SRAM design. In a 32-kB memory implementation, the M3D design can reduce footprint, latency, and energy by 33%, 10%, and 19%, respectively. The compiler is used to show the feasibility of fine-grain logic and SRAM stacking in M3D technology.

**INDEX TERMS**   CNFET, memory compiler, monolithic-3-D (M3D), electronic design automation (EDA), SRAM.

## I. INTRODUCTION

**M**ONOLITHIC 3-D (M3D) integrated circuits (ICs) using fine-grain nanoscale interlayer vias (ILVs) promise significant energy-efficiency improvements over 2-D ICs [1]. However, M3D requires sequential fabrication of multiple layers of transistors in one substrate, where the circuit components in different tiers are interconnected via high-density 3-D vias. The need for high-temperature processing of silicon-based MOSFET (1000 °C) during sequential fabrication of one tier can degrade the reliability and performance of devices in previously fabricated tiers [1]–[5]. The recent works on silicon-based M3D processes are exploring techniques that address the high-temperature processing challenges [6], [7].

Carbon nanotube FET-based M3D process has emerged as an attractive alternative to the silicon-based M3D process [2], [4], [8]–[11]. This is because CNFET can be fabricated at the temperature below 425 °C, which eliminates potential defects of devices and interconnections on previously fabricated tiers [2]. Hence, in a CNT-based M3D IC, on/off currents of CNFET transistors in different tiers are close

to each other. CNFET also promises high energy efficiency in designing logic and memory circuits [8], [9]. Shulaker *et al.* [12], [14] and Hills *et al.* [13] have demonstrated applications of CNFET-based M3D processes. Srimani *et al.* [2] have demonstrated commercial-grade M3D process design kits (PDKs) and the operation of logic and SRAM.

In this article, we present an SRAM compiler for CNFET-M3D using the PDK developed by Srimani *et al.* [2]. Although there are few prior works on SRAM cell design in M3D silicon [15] or CNFET [2], to the best of our knowledge, there is no SRAM compiler for CNFET-based M3D. Fig. 1 shows a schematic of the M3D stack that includes two CNFETs layers and six metal layers. Our compiler leverages similar CNFET performance in different tiers in two ways. We first exploit this observation to generate a large SRAM array by combining *stackable single-tier* subarrays designed in individual tiers. These *stackable single-tier* subarrays are referred to as the **STF**, i.e., **S**ingle **t**ier of **F**ETs contains bit-cells and peripherals, all in a single tier. The *stackable single-tier* design is achieved by separating metal layers used in each SRAM tier, as shown in Fig. 1. As subarrays in
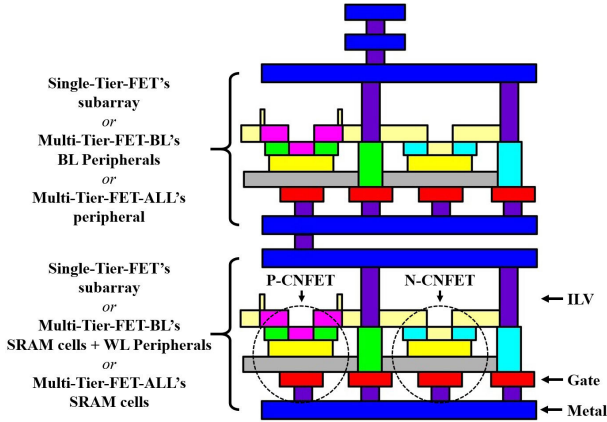
**FIGURE 1. Schematic layer organization of M3D PDK used in our design (redrawn after [2]). The two tiers of CNFETs are included in this process.**



**FIGURE 2. Schematic of SRAM subarray.**

different tiers have similar performance, we can efficiently integrate them to generate a large array.

Second, we present 3-D SRAM subarrays composed of multiple tiers of the transistor where bit-cells in one tier and peripheral circuits in other tiers are connected using fine-grain ILVs. As transistors in different tiers have similar performance, the folding of the peripheral circuits into multiple tiers allows reducing footprint while maintaining (or improving) performance. We present two different types of multitier SRAM subarrays. An **MTF-BL** subarray contains **M**ultiple **t**iers of **F**ETs with bitline (**BL**) peripherals in a second layer of transistors. An **MTF-ALL** subarray contains **M**ultiple **t**iers of **F**ETs with **all** peripherals in wordline (WL) and BL peripherals in one tier of transistors and bit-cells in a different tier.

Our memory compiler flow generates the layout files [library exchange format (LEF) and graphical data system (GDS)] and the timing file [liberty timing file (LIB)] for the SRAM memory with different types of subarrays with varying dimensions. Our compiler supports the scaling capacity of the generated SRAM by integrating multiple subarrays. First, the compiler can generate the physical design of an **SRAM block** by integrating smaller subarrays using an H-tree architecture. Second, the further scaling of memory capacity is achieved by connecting multiple **SRAM blocks** using Network-on-Chip (NoC) to generate an **SRAM array**. We demonstrate the application of the CNFET M3D SRAM compiler for generating SRAM subarrays, blocks, and arrays of varying capacities. We apply the compiler to generate the physical design of a system architecture with a multicore processor in one tier integrated with an SRAM array in the other tier. Each core locally connects to a smaller capacity SRAM block, but all SRAM blocks are connected via an NoC to create a large capacity array but with multiple distributed access ports.

We observe that M3D subarray designs can improve the properties of the memory compared to 2-D subarray designs. The combination of WL and BL, used for the comparison, is 64WL × 64BL, 64WL × 128BL, 128WL × 64BL,

and 128WL × 128BL. The footprint, read energy, write energy, and read latency can be reduced by 27.8%–39.8%, 1.7%–2.8%, 4.7%–8%, and 9.7%–11.3%, respectively. In addition, the properties of the SRAM block also can be improved by using MTF designs instead of STF designs. 32 kB of SRAM block can achieve a 24.8% lower footprint and 9.5% lower energy.

## II. M3D SRAM ARCHITECTURES AND COMPILER
### A. SUBARRAY ARCHITECTURES
The SRAM subarrays include bit-cells, address decoders and drivers, sense amplifiers (SAs), and write-driver circuits (see Fig. 2). We use the single-ended dynamic SA (see Fig. 2). Before WL is raised high, the signal EN is high and pre-discharges the node X, thereby precharging OUT to high. During reading, the signal EN is made ''low,'' which turns the pMOS (P1) ON. If the SRAM bit-cell is storing a ''0,'' the BL discharges, thereby turning on the pMOS P2, which charges the node X and discharges OUT to low. If the SRAM cell is storing an ''1,'' the BL remains high, which ensures the node X and OUT remains low and high, respectively. All the logics are custom-designed and automatically placed by the SKILL code. We assume all BLs are read/written in parallel, i.e., no BL interleaving or column multiplexing inside the subarray. All the columns are read out in parallel and multiplexed outside of the subarray to create the designed data width (32 bit).

Fig. 3 shows the STF, MTF-BL, and MTF-ALL subarray microarchitectures. All these microarchitectures follow the same schematic design (see Fig. 2), but different placements. For both STF and MTF designs, the bit-cell arrays remain in a 2-D arrangement. The key difference resides in the placement and arrangement of peripheral circuitry. The physical implementation of bit-cells and all peripherals limit layer usage to one of the two tiers (layers) of CNFETs and the two immediate layers of metal (one above and below the FET, each) (see Fig. 1) to prevent metal usage overlap between top tier and bottom tier designs. This allows M3D stacking of bit-cells and peripherals without the need for reimplementing physical layouts.

In the STF subarray, peripheral circuits are placed in the same tier as bit-cells following conventional 2-D subarray arrangements. Our compiler supports the design of STF subarrays in both tiers of transistors provided by the PDK.
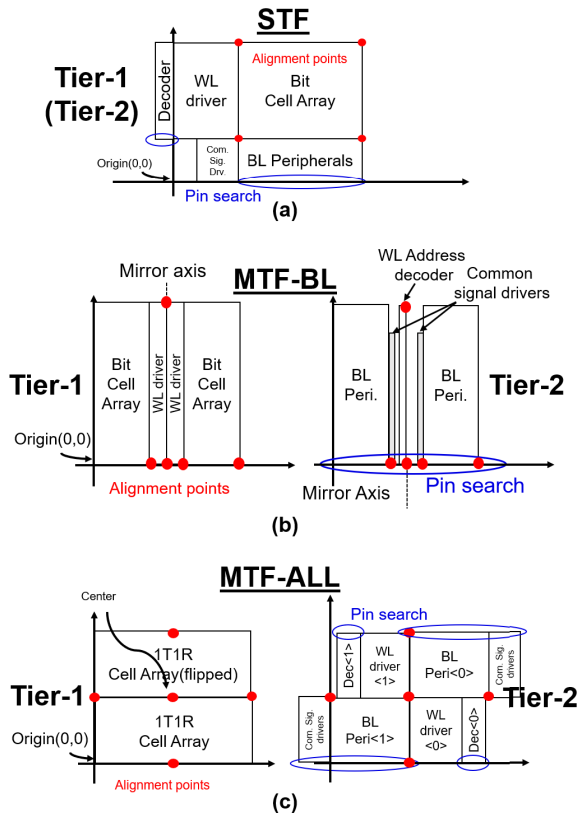
**FIGURE 3.** Subarray peripheral/bit-cell arrangement for (a) 2-D STF, (b) MTF-BL, and (c) MTF-ALL subarrays.



**FIGURE 4.** Centered placement of peripherals reduces worst case mismatch.



**FIGURE 5.** Multiple ILV path for writing with MTF stacked peripherals.

From post-PEX simulations, Tier-2 STF subarrays demonstrate only 4% and 6% reductions in read energy and latency, respectively, compared to Tier-1 STF designs. The detailed analysis shows that the difference is mainly contributed by different parasitics involved when changing tiers. However, for comparison with MTF designs, Tier-1 STF designs are used as the baseline due to the fact that Tier-1 STF and MTF designs all use the same tier (Tier-1) for implementing bit-cells.

Fig. 3(b) shows the MTF-BL subarray architecture. The WL drivers are placed in Tier-1 along with bit-cells, but WL address decoders and BL peripherals are in Tier-2. As BL peripherals and WL drivers reside in different tiers, additional functionality is added to scale/distribute the BL driver output stage as subarray dimensions change. This allows trimming the BL driver for small subarrays to reduce footprint/performance overhead. Fig. 3(c) shows MTF-ALL architecture where all peripherals are in Tier-2. The arrangement is such that, with peripherals and bit-cells rotated/flipped, it can ensure that output edges of peripherals are aligned along the center axes of the bottom bit-cell arrays. This structure prevents BL peripherals from blocking WL access and vice versa by segmenting all peripheral instantiating to twice. However, the central location of peripherals also separates the access pins to different edges, which needs to be considered during automated generation for multiple interconnected subarrays.
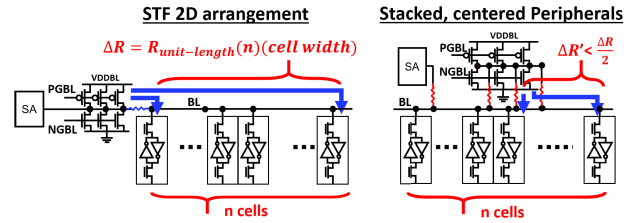
A unique property of the CNFET M3D process is the use of back-gated CNFET (see Fig. 1). A key advantage of back-gate FET geometries is the reduction in gate-to-plug capacitance [16], which reduces overall WL capacitance and BL capacitance in M3D SRAM. Our compiler utilizes the back-gate structure by placing bit-cells in the bottom tier and BL peripherals in the top tier, where BLs occupy the metal layers between tiers. This arrangement allows multiple ILVs to be constructed for resistance reduction and minimizes parasitic capacitance applied from surrounding devices/metals to WLs that occupy the bottommost metal layer.

MTF subarrays created by stacking peripherals directly on top of the bit-cell array allow several advantages. The first advantage is reduced mismatch seen by peripherals (see Fig. 4). MTF designs allow flexible placement of peripherals compared to traditional 2-D structures. By centering the stacked peripherals, up to $2\times$ reduction in worst case mismatch can be achieved. As worst case conditions have been cut short, improvement in read/write performance can also be observed. The second advantage is a benefit enabled by the ability to manufacture dense ILVs in M3D technologies. Instead of utilizing ILVs for digital signals, MTF subarray structures use multi-ILV paths to create low-resistance connections between WLs/BLs and their corresponding drivers (see Fig. 5). Reduced resistance not only allows faster BL and WL switching but also reduces I-R drop, which can affect write stability for far-end bit-cells.

### B. SRAM BLOCK AND ARRAY ARCHITECTURE
We compile a large capacity SRAM array by integrating multiple smaller capacity subarrays. First, we use an H-tree architecture to integrate the smaller subarrays to compile an **SRAM block** (see Fig. 6). Next, we combine multiple memory blocks using an NoC to design an **SRAM array** (see Fig. 7).
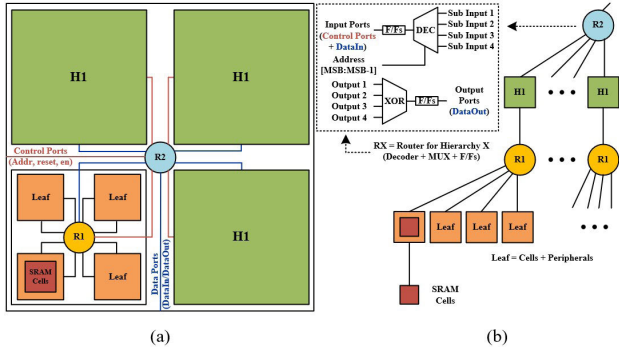
**FIGURE 6.** SRAM block design: (a) high-level placement and (b) H-tree architecture.
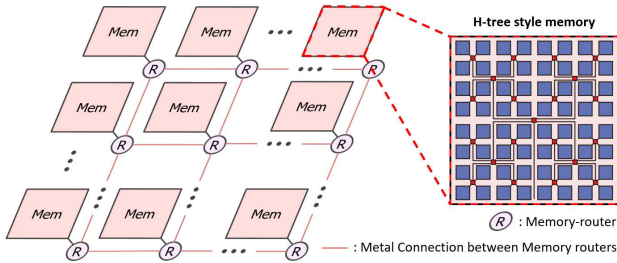


**FIGURE 7.** Multiport memory network architecture.

## 1) ARCHITECTURE OF THE SRAM BLOCK

An H-tree is a hierarchical design where each node of the tree accumulates data from lower level nodes. We refer to the logic necessary to combine the subarrays within the H-tree as the **top module**. The compiler can integrate any size and type of subarrays as the leaf nodes of an H-tree. We place the H-tree router in the empty space among the submodules. During reading, the H-tree router multiplexes multiple incoming data ports from a lower level in the hierarchy to a single output port to the higher level in the hierarchy. During writing, the H-tree router demultiplexes the incoming data port from the higher level in the hierarchy to one of the output ports to the lower level in the hierarchy.

## 2) ARCHITECTURE OF THE SRAM ARRAY

An SRAM array is implemented by connecting multiple SRAM blocks using a memory NoC (see Fig. 7). The current compiler generates a mesh NoC where each SRAM block is connected to a router. We use the open-source NoC router, which has a virtual channel (VC) for the deadlock-free algorithm [17]. The router controls the read/write access to individual SRAM blocks and manages the data movement within the memory NoC.

## 3) DATAFLOW OF SRAM ACCESS

Our design creates a distributed multiport memory where each router acts as an I/O port. Individual logic blocks can be connected to one router and access the entire memory. When there is a memory access request from the logic blocks, the memory network automatically calculates the index of the target SRAM block from the input address. Depending on the physical distance between the target SRAM block and the

requesting logic block, the access request can traverse zero, one, or multiple hops in the network.

### C. EVALUATION AND DESIGN SPACE EXPLORATION

The latency, area, and energy of the subarrays are evaluated using SPICE simulations of the extracted netlists. The SRAM block, including the H-tree network and the routers of the memory NoC in the entire array, is synthesized to meet the target memory frequency while ensuring single cycle communication between successive levels in H-tree and nearest routers in the NoC.

To evaluate the energy consumption of the array, we synthesize and perform the PNR of each level of the hierarchy independently. To calculate the dynamic energy of an SRAM block, we recognize that only a single subarray and only one node in each level of the H-tree are active every cycle. The static energy for the entire SRAM block is computed. The energy for accessing the entire array includes the energy of one SRAM block and the number of active routers (i.e., hops in NoC). As the number of hops for access can vary and is only available from a detailed architectural simulation, this article considers all the routers are active (the worst case scenario).

The compiler allows exploring various dimensions of subarrays, the capacity of SRAM blocks (i.e., number of levels in the H-tree hierarchy), and the number of SRAM blocks (i.e., number of nodes in the mesh NoC) for a given memory capacity to meet a design goal, such as the minimum footprint or minimum energy. The subarray options include the type of the subarray (STF, MTF-BL, and MTF-ALL) and different dimensions of a given subarray type, all of which determine the access latency/energy of individual subarrays. Generating a memory array with smaller capacity SRAM blocks will lead to more hops while accessing a distant address. However, as a smaller capacity of SRAM block have lower read/write latency/energy, the cost of local accesses will be reduced.

## III. SIMULATION RESULTS
### A. RUN TIME OF COMPILER

The run time of the compiler is measured on a desktop with an i7-9700 core and 16-GB memory. The subarray layout generation takes less than 1 min. The analysis takes 30 and 120 min for a 64WL × 64BL and a 128WL × 128BL subarray, respectively. The run time to compile an SRAM block depends on the number of levels in the H-tree, where each level requires 20–30 min. The generation of a 2-MB SRAM array with 4 × 4 SRAM blocks (128 kB) and mesh NoC requires 120 min.

### B. SUBARRAY COMPILATION RESULTS

Fig. 8 visualizes the layout of different subarray types.

#### 1) FOOTPRINT ANALYSIS

Fig. 9(a) shows the 2-D footprint area for different subarray dimensions normalized to 64WL × 64BL STF subarray
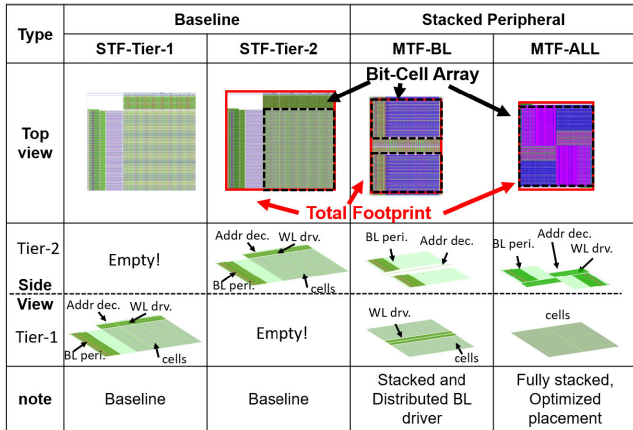
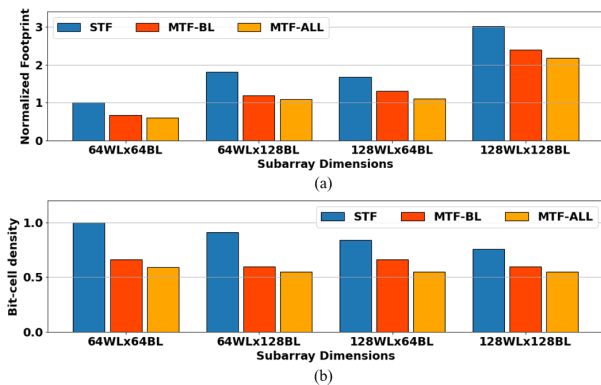**FIGURE 8. Visualization of layouts of different subarray (64WL × 64BL) structures.**



**FIGURE 9. Subarray area analysis: (a) footprint and (b) bit-cell density.**



**FIGURE 10. Normalized subarray's comparison for different subarrays on (a) read energy per bit, (b) read latency, and (c) write energy per bit.**

subarrays for all subarray dimensions. MTF-BL structures show lower read latency for small subarrays but higher read latency for large subarrays due to the additional parasitics introduced when the driver output stage is distributed across longer BLs for large subarrays.

### 4) WRITE ENERGY ANALYSIS

As the write energy is dominated by the parasitic resistance and capacitance along the BL, we observe a reduction in write energy for both MTF designs at smaller subarray dimensions. However, similar to read access metrics, a marginal increase in write energy for MTF-BL due to the higher parasitics introduced to/by distributed drivers is also observed for larger subarray dimensions. MTF-ALL subarrays, on the other hand, demonstrate a reduction in write energy for all subarray dimensions.

### C. SRAM BLOCK COMPILATION RESULTS

Fig. 11 shows the top view and the side view of an SRAM block. The STF-Stack design represents the stacking of Tier-1 and Tier-2 STF arrays to create the entire capacity. In this example, 32 kB of SRAM block is designed with 128WL × 128BL subarrays. Sixteen subarrays are used in MTF designs, while STF designs in each tier use eight subarrays. A 32-bit bus width is assumed for simulation. The frequency target for the PNR follows the subarrays' maximum frequency (see Fig. 10).

### 1) FOOTPRINT ANALYSIS

Fig. 12(a) shows the footprint comparison results (normalized to footprint generated by 64WL × 64BL STF subarray). As expected, the subarray with a larger dimension reduces the total footprint. The STF-Stack shows a lower footprint as the STF subarrays are stacked using both Tier-1 and Tier-2.

footprint. MTF designs show a lower footprint than the STF design. However, the benefit reduces for larger subarrays where peripheral circuits have relatively lower contributions. MTF-ALL structures suffer from peripheral overhead at small subarray dimensions. This is because peripherals placed in Tier-2 extend beyond the bit-cell array boundaries in Tier-1. This overhead is addressed by scaled drivers in MTF-BL, which provides the lowest footprint for 64WL × 64BL subarray.

### 2) CELL-DENSITY ANALYSIS

Fig. 9(b) shows the bit-cell density for different subarray architectures. In traditional 2-D structures, the reduced cell density limits the usage of small subarrays. The cell density can be increased by 1.32× by changing STF subarray dimensions from 64WL × 64BL to 128WL × 128BL. MTF subarrays with 64WL × 64BL can achieve superior bit-cell density than 128WL × 128BL STF.

### 3) READ LATENCY AND ENERGY ANALYSIS

Fig. 10 compares the read performance and read/write energy per bit of different types of subarrays. The data are normalized to the energy and latency of STF subarrays at the **same dimensions**. We observe that the MTF-ALL designs show lower read latency and read energy compared to the STF
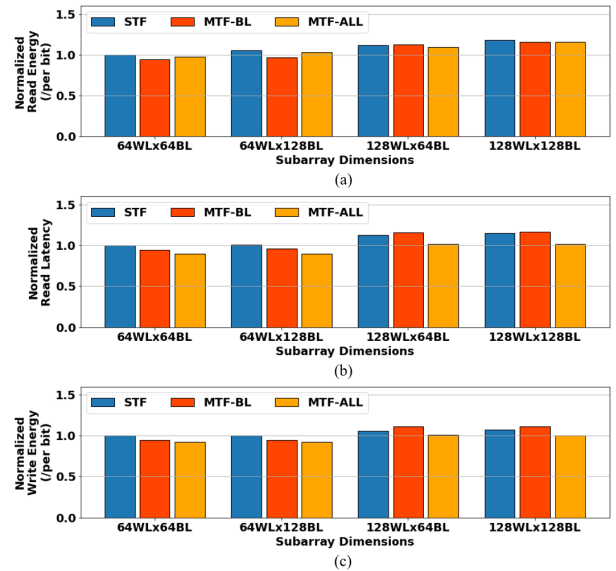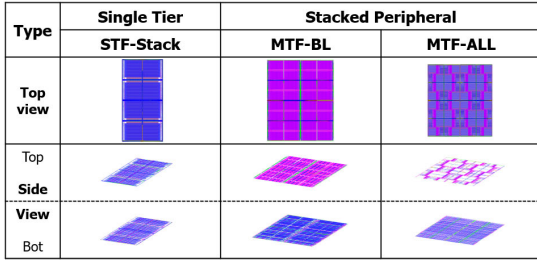
**FIGURE 11. Layout of compiled arrays (32 kB with 128WL × 128BL subarrays).**
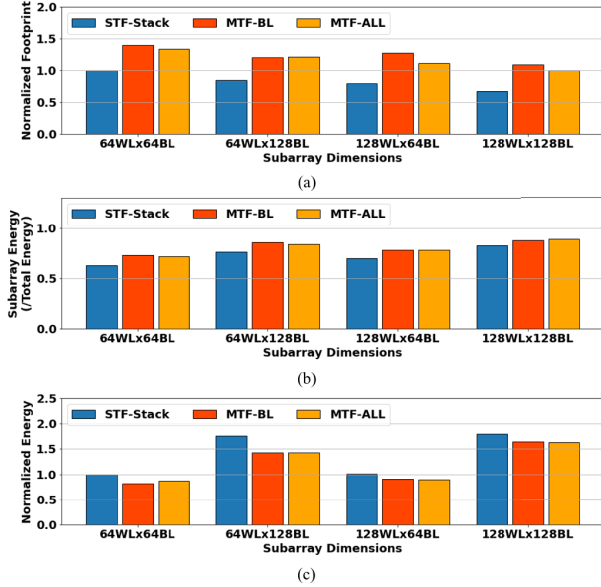


**FIGURE 12. SRAM block (32 kB) analysis. (a) Normalized footprint. (b) Subarray's energy per total energy. (c) Normalized total energy.**

The MTF-ALL shows a smaller footprint than the MTF-BL as all peripherals are in 3-D.

*2) ENERGY ANALYSIS*

The minimum energy design of an SRAM block depends on the tradeoff between subarray and top module energy. A larger dimension increases subarray energy but reduces the top module's energy as fewer subarrays are used. Also, using large subarrays reduces footprint and parasitic capacitance in the top module. Hence, we observe that larger subarrays increase the ratio of the subarray's energy to the SRAM block's energy, as shown in Fig. 12(b). As the majority of the energy consumption is from the subarrays, using larger subarray dimensions increases the energy of an SRAM block [see Fig. 12(c)].

*3) SCALABILITY OF THE COMPILER*

The developed compiler is used to generate SRAM blocks of varying capacities (see Fig. 13). All the designs use the 128WL × 128BL subarray to implement the SRAM block. As only one subarray consumes read/write energy irrespective of the total capacity, normalized read energy only increases marginally for larger capacity, mainly due to the
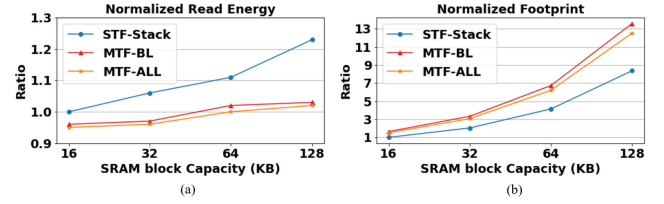


**FIGURE 13. Scalability of SRAM block capacity: (a) read energy and (b) footprint.**



|  | Minimum Footprint | Minimum Latency | Minimum Energy |
|---|---|---|---|
| Top view | | | |
| Type | STF-Stack | MTF-ALL | MTF-BL |
| Dimension | 128WLx128BL | 64WLx64BL | 64WLx64BL |
| Normalized Footprint | 0.67 | 1.34 | 1.40 |
| Normalized Latency | 1.15 | 0.9 | 0.95 |
| Normalized Energy | 1.8 | 0.86 | 0.81 |

**FIGURE 14. Memory compiler generates the best design for different objectives. All the data are normalized by STF-Stack 64WL × 64BL subarray-based design.**

increased complexity of the top module [see Fig. 13(a)]. On the other hand, the footprint advantage of STF-Stack over MTF designs increases with the larger capacity of the SRAM block [see Fig. 13(b)].

*4) BLOCK OPTIMIZATION*

Fig. 14 shows the memory compiler output for 32-kB SRAM block but varying design targets. We explored the three types of subarrays (STF-Stack, MTF-BL, and MTF-ALL) and four dimensions (64WL × 64BL, 64WL × 128BL, 128WL × 64BL, and 128WL × 128BL) while generating final designs. We observe that designs with the largest possible STF-Stack result in the minimum footprint. On the other hand, the MTF-ALL subarray with the smallest dimension shows the minimum latency. Among the 12 cases that we compared, the 32-kB SRAM block designed with 64WL × 64BL MTF-BL subarray shows the minimum energy. However, when we generate the 128-kB SRAM block, MTF-ALL with 128WL × 64BL subarray-based design shows minimum energy.

*D. MULTIPORT SRAM ARRAY ANALYSIS*

Fig. 15 shows the layout of a 2-MB SRAM array designed with 4 × 4 SRAM blocks and routers. Each 128-kB SRAM block is designed with STF-64WL × 64BL subarrays. Fig. 16 shows the footprint and power of the design considering various SRAM block and subarray capacities. All the results are normalized to the SRAM array generated by 16-kB SRAM block and 64WL × 64BL subarray. As expected, using a higher capacity SRAM block shows a smaller footprint and lower worst case power (assuming all routers
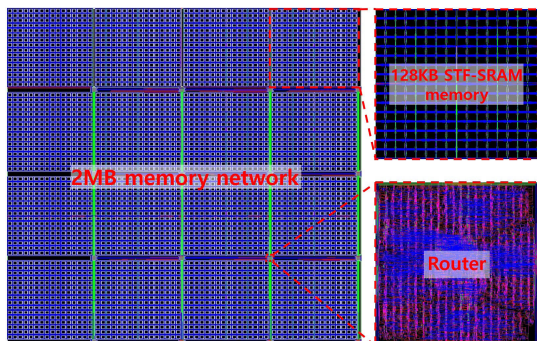
**FIGURE 15.** 2-MB memory network consisting of 4 × 4 128-kB SRAM blocks.



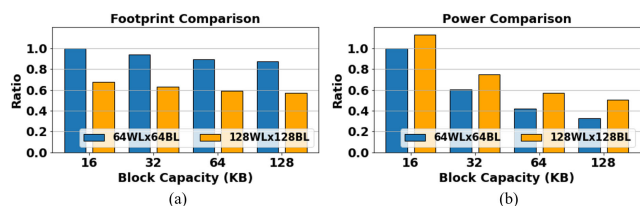**FIGURE 16.** Analyze (a) footprint and (b) power ratio about 2-MB memory network, which consists of different SRAM blocks and subarray capacities.



**FIGURE 17.** 3-D multicore processor and multiport memory network stacked architecture.



**FIGURE 18.** Analyze (a) local memory access latency and (b) worst case hop count about 2-MB memory network, which consists of different SRAM block capacities.

are active) because the number of SRAM blocks and memory routers is decreased. Although the power of individual SRAM blocks is higher with a larger capacity, the power reduction due to fewer routers dominates. Note that, if only a few routers are active (i.e., local memory access) and/or router power is more optimized, the smaller capacity blocks will be more power-efficient. Among different subarray choices, as expected from Fig. 12, the design with 128WL × 128BL subarray shows a lower footprint but higher power than the design with 64WL × 64BL subarray.

### E. LOGIC AND MEMORY STACKING
We use the memory compiler to implement the physical design of a 3-D multicore design (see Fig. 17). We use the OpenPiton as our multicore architecture [18] connected to a shared on-chip cache organized as multiple networked SRAM blocks. Fig. 17 shows the layout of a design with 4 × 4 OpenPiton cores connected to 2 MB of shared cache organized as 4 × 4 128 kB SRAM blocks (same as in Fig. 15).

The multiple cores in M3D tier-2 are connected via a core NoC, and each core vertically connects to a memory router in tier-1. The core-to-core communication occurs via the NoC in the core tier, while the NoC in the memory tier supports the core-to-memory communication. The distributed memory access enables fast and low-latency access between a core and its local memory due to direct vertical access (reduced wire-length) using ILVs. The multiple access ports also avoid the memory congestion associated with using a single port. On the other hand, the memory network enables any core to access any memory location enabling a large shared cache. The absence of the memory network will limit the available memory capacity of each core and lead to data duplication and cache coherence challenges.
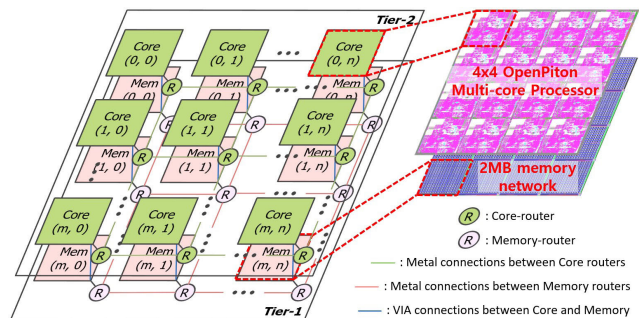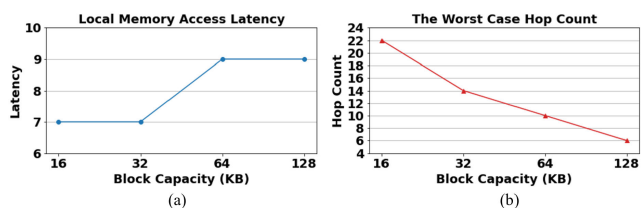
We study the impact of SRAM block capacity on the memory access latency (see Fig. 18). The local read or write access from a core to its local SRAM block is performed directly without using a memory router. Hence, the latency of the "local" memory access depends on the access latency of the SRAM block, i.e., the number of levels in the H-tree hierarchy. Increasing capacity of SRAM block increases the local access latency [see Fig. 18(a)]. Due to the H-tree organization, the latency of the local access is constant for any address. The access to a distant memory block, referred to as the "global memory access," is performed by using the memory network. Hence, the global access latency is variable depending on the target address. The worst case latency of the global access depends on the maximum number of hops in the mesh network. Using larger SRAM block capacity reduces worst case number of hops [see Fig. 18(b)]. Note that the global access latency of a specific request depends on the memory access patterns of *all* cores, which determines the congestion in the network.

### IV. DISCUSSION
This article has studied an M3D memory compiler that can exploit the multiple layers of transistors on a single substrate and the high-density (3-D via) interconnection between tiers. The design concepts and compiler methodologies developed in this article can be adopted for alternative M3D IC technologies. In particular, as the array generation and the multiport memory network generation follow logic design automation methods, they can be easily ported to other M3D processes. The circuit design concepts within the subarray generation are also technology agnostic. However, due to the inherently analog nature of the circuits, the exact topologies, sizing, and layout are technology-specific.

The 3-D SRAM compiler needs additional considerations compared to a traditional 2-D SRAM compiler. The goal of the M3D compiler is to maximize the benefits of the M3D PDK by allowing: 1) logic to be placed above (below) STF subarrays in bottom (top) tiers and 2) external routing to cross MTF subarrays when needed. The compiler also needs to consider the availability of only two metal layers between the top and bottom tiers of CNFETs. Hence, the M3D SRAM compiler must restrict the use of metal layers (to only one layer above and one layer below the transistors) while implementing bit-cells and peripherals. The above restriction creates several challenges. For example, the M3D compiler can use only a single metal layer to design a power delivery network (PDN) of SRAM subarrays. Therefore, compared to traditional 2-D SRAMs that use multiple layers of metals for PDN, the M3D SRAMs can experience higher IR drop. This IR drop reduces the robustness of the memory operation.

The orientation and organization of BLs/WLs metal wires and VDD/VSS grids while compiling MTF subarrays must consider the constraints on metal layer usage. For example, multiple routing vacancies are required to allow ILVs for cross-tier access and intratier routing. These routing vacancies are aligned to avoid possible connectivity hazards during the automated generation when subarray dimensions scale. The cross-tier access is the connections for signal/power nets from M3D peripherals to nets in the bottom bit-cells. The intratier routing is the transistor-bypass path to allow access between the horizontal-routing layer and the vertical-routing layers that are below and above transistors. Moreover, the routing connections passing through a transistor layer need to be carefully distributed to place ILV stacks in MTF subarrays. We place routing connections for WL and BL in the top and bottom tier metal layers, respectively. The ILVs required for BL connection use the vacancies in the top tier, and the WL connection uses the vacancies in the bottom tier to improve area efficiency.

The parasitics of individual tiers are important factors while optimizing M3D designs. For example, although CNFETs in different tiers have similar on current, there is still a 4%–6% difference in the performance of STF SRAMs in different tiers, which can be important for high-frequency designs. The differences in parasitics also play a key role in optimizing the design/placement of peripheral circuits for MTF subarrays.

We use an H-tree-based connection of multiple subarrays of small dimensions to create a single-port memory block of higher capacity. The H-tree-based connectivity simplifies the router design and wiring density within a memory block but can only support one memory access per read/write cycle. On the other hand, we connect multiple memory blocks using a mesh-style NoC to create a high-capacity memory array. As each memory block is vertically connected to a logic core, the NoC-based design allows multiport access to the overall large capacity memory allowing a high degree of access concurrence (and, hence, bandwidth) for local core-to-cache communications, such as a distributed cache architecture. However, this approach also preserves the logically shared structure of the cache as all memory blocks can be accessed from any core, thereby avoiding the complex cache coherency protocol required in traditional distributed caches. In other words, the proposed approach uses the M3D integration to create a physically distributed but logically shared cache that harnesses the bandwidth advantage of M3D for local access while maintaining globally shared models of the memory to reduce the burdens on memory managements of distributed caches. Furthermore, connecting memory blocks with an NoC in the memory tier reduces the communication burden on the core routers and separates the core-to-core and core-to-memory communications.

M3D technology can improve the latency and the energy of the memory access on the multicore system with the additional level of cache and routing resources. In the 2-D system design, the memory is placed next to the edge of the processor. Therefore, to read/write function, data are transmitted across all the NoC routers from/to the edge of the processor to/from the target core. This transmission increases the memory access latency, energy, and NoC congestion. However, in M3D designs, we can place the multiport memory above/below the processor. In this design, we can directly connect the memory block to the core and do not require the data transmission through NoC routers in the multicore processor.

## V. CONCLUSION

This article demonstrates an SRAM memory compiler for CNFET-based M3D technologies with multiple tiers of transistors. The compiler exploits the ability to have identical device performance in all tiers to generate scalable subarray using single and multiple tiers of peripherals. Our compiler efficiently integrates the subarrays in a single or multiple tiers to generate an SRAM block and integrates the blocks using an NoC to generate an SRAM array. Ultimately, our compiler automatically generates the layout (LEF/GDS) and timing (LIB) files of SRAM blocks that can be used in the full-chip design. The simulation results show that multitier SRAMs with peripheral and bit-cells in different tiers show lower energy and latency, while 3-D stacking single-tier SRAMs show a better footprint. We also show the feasibility of using the memory compiler to generate an M3D stack of logic and memory. The future work will be improving the compiler and exploring new system architectures enabled by efficient M3D memory compilers.

## REFERENCES

[1] S. Wong, A. El-Gamal, P. Griffin, Y. Nishi, F. Pease, and J. Plummer, "Monolithic 3D integrated circuits," in *Proc. Int. Symp. VLSI Technol., Syst. Appl. (VLSI-TSA)*, Apr. 2007, pp. 1–4.

[2] T. Srimani *et al.*, "Heterogeneous integration of BEOL logic and memory in a commercial foundry: Multi-tier complementary carbon nanotube logic and resistive RAM at a 130 nm node," in *Proc. IEEE Symp. VLSI Technol.*, Jun. 2020, pp. 1–2.

[3] D. Akinwande, S. Yasuda, B. Paul, S. Fujita, G. Close, and H. S. P. Wong, "Monolithic integration of CMOS VLSI and carbon nanotubes for hybrid nanotechnology applications," *IEEE Trans. Nanotechnol.*, vol. 7, no. 5, pp. 636–639, Sep. 2008.

[4] M. M. Shulaker, K. Saraswat, H.-S.-P. Wong, and S. Mitra, "Monolithic three-dimensional integration of carbon nanotube FETs with silicon CMOS," in *Symp. VLSI Technol. Dig. Tech. Papers*, Jun. 2014, pp. 1–2.

[5] P. Batude, T. Ernst, J. Arcamone, G. Arndt, P. Coudrain, and P. E. Gaillardon, "3-D sequential integration: A key enabling technology for heterogeneous co-integration of new function with CMOS," *IEEE J. Emerg. Sel. Topics Circuits Syst.*, vol. 2, no. 4, pp. 714–722, Dec. 2012.

[6] C. Fenouillet-Beranger *et al.*, "New insights on bottom layer thermal stability and laser annealing promises for high performance 3D VLSI," in *IEDM Tech. Dig.*, Dec. 2014, pp. 27.5.1–27.5.4.

[7] L. Pasini *et al.*, "nFET FDSOI activated by low temperature solid phase epitaxial regrowth: Optimization guidelines," in *Proc. SOI-3D-Subthreshold Microelectron. Technol. Unified Conf. (S3S)*, Oct. 2014, pp. 1–2.

[8] M. D. Bishop *et al.*, "Fabrication of carbon nanotube field-effect transistors in commercial silicon manufacturing facilities," *Nature Electron.*, vol. 3, no. 8, pp. 492–501, Aug. 2020.

[9] G. Hills *et al.*, "Understanding energy efficiency benefits of carbon nanotube field-effect transistors for digital VLSI," *IEEE Trans. Nanotechnol.*, vol. 17, no. 6, pp. 1259–1269, Nov. 2018.

[10] M. M. Shulaker, T. F. Wu, M. M. Sabry, H. Wei, H.-S. P. Wong, and S. Mitra, "Monolithic 3D integration: A path from concept to reality," in *Proc. Design, Autom. Test Eur. Conf. Exhib. (DATE)*, 2015, pp. 1197–1202.

[11] T. F. Wu *et al.*, "Hyperdimensional computing exploiting carbon nanotube FETs, resistive RAM, and their monolithic 3D integration," *IEEE J. Solid-State Circuits*, vol. 53, no. 11, pp. 3183–3196, Nov. 2018.

[12] M. Shulaker *et al.*, "Experimental demonstration of a fully digital capacitive sensor interface built entirely using carbon-nanotube FETs," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2013, pp. 112–113.

[13] G. Hills *et al.*, "Modern microprocessor built from complementary carbon nanotube transistors," *Nature*, vol. 572, no. 7771, pp. 595–602, 2019.

[14] M. M. Shulaker *et al.*, "Monolithic 3D integration of logic and memory: Carbon nanotube FETs, resistive RAM, and silicon FETs," in *IEDM Tech. Dig.*, Dec. 2014, pp. 27.4.1–27.4.4.

[15] S. Srinivasa, X. Li, M. Chang, J. Sampson, S. K. Gupta, and V. Narayanan, "Compact 3-D-SRAM memory with concurrent row and column data access capability using sequential monolithic 3-D integration," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 26, no. 4, pp. 671–683, Apr. 2018.

[16] T. Srimani, G. Hills, M. D. Bishop, and M. M. Shulaker, "30-nm contacted gate pitch back-gate carbon nanotube FETs for sub-3-nm nodes," *IEEE Trans. Nanotechnol.*, vol. 18, pp. 132–138, 2018.

[17] H. Matsutani, M. Koibuchi, H. Amano, and T. Yoshinaga, "Prediction router: Yet another low latency on-chip router architecture," in *Proc. IEEE 15th Int. Symp. High Perform. Comput. Architecture*, Feb. 2009, pp. 367–378.

[18] J. Balkind *et al.*, "OpenPiton: An open source manycore research framework," *ACM SIGARCH Comput. Archit. News*, vol. 44, no. 2, pp. 217–232, 2016.