

# TDRAM: Tag-enhanced DRAM for Efficient Caching

Maryam Babaie<sup>1</sup>, Ayaz Akram<sup>1</sup>, Wendy Elsasser<sup>2</sup>, Brent Haukness<sup>2</sup>, Michael Miller<sup>2</sup>,  
Taeksang Song<sup>2</sup>, Thomas Vogelsang<sup>2</sup>, Steven Woo<sup>2</sup>, Jason Lowe-Power<sup>1</sup>

<sup>1</sup>Department of Computer Science, University of California, Davis

<sup>2</sup>Rambus Labs, Rambus Inc.

**Abstract**—As SRAM-based caches are hitting a scaling wall, manufacturers are integrating DRAM-based caches into system designs to continue increasing cache sizes. While DRAM caches can improve the performance of memory systems, existing DRAM cache designs suffer from high miss penalties, wasted data movement, and interference between misses and demand requests. In this paper, we propose TDRAM, a novel DRAM microarchitecture tailored for caching. TDRAM enhances HBM3 by adding a set of small low-latency mats to store tags and metadata on the same die as the data mats. These mats enable fast parallel tag and data access, on-DRAM-die tag comparison, and conditional data response based on comparison result (reducing wasted data transfers) akin to SRAM caches mechanism. TDRAM further optimizes the hit and miss latencies by performing opportunistic early tag probing. Moreover, TDRAM introduces a flush buffer to store conflicting dirty data on write misses, eliminating turnaround delays on data bus. We evaluate TDRAM using a full-system simulator and a set of HPC workloads with large memory footprints showing TDRAM provides at least  $2.6\times$  faster tag check,  $1.2\times$  speedup, and 21% less energy consumption, compared to the state-of-the-art commercial and research designs.

## I. INTRODUCTION

Today’s computers, equipped with significant processing capabilities and memory capacities, aim to fulfill the requirements of high-performance computing (HPC) tasks such as machine learning and artificial intelligence. To find the right balance between performance and capacity, manufacturers have embraced heterogeneous memory systems. These systems combine high-performance memories like HBM with high-capacity memories with lower performance. The introduction of interconnect technologies like Compute Express Link (CXL) is increasing memory heterogeneity with local and remote memory pools. Intel’s Sapphire Rapids CPU demonstrates this strategy by employing on-package HBM DRAMs as cache [25], [61]. This approach effectively boosts cache capacities and tackles the scaling limitations encountered by SRAM [5]. The expanded cache capacity and enhanced bandwidth provided by HBM DRAMs offer the potential for improved data locality, without programmers intervention.

However, the potential benefits of DRAM caches have not borne fruit. Previous studies of DRAM caches have shown that using standard DRAM devices as a cache can slow down applications with large memory footprints and high miss rates [38], [59]. The current designs of DRAM caches, such as those in Intel’s Cascade Lake DRAM cache [8], [38], store tag and metadata together with the cache line data in the same DRAM. Storing tags and data together reduces hit time for read demands [58], but significantly increases miss penalties since a separate DRAM read is necessary to retrieve tag and metadata for hit/miss and status information, causing contention with demand reads. Also, *all write requests, including those hitting on the cache, require a DRAM read* to fetch tag and data to ensure dirty data is not overwritten, further exacerbating the contention and causing expensive turnaround bubbles on the data bus [19]. These extra accesses for read-misses and write demands increase: (i) latency for demand misses, (ii) contention in the read buffer which extends the queue occupancy time, and (iii) wasted data movement and energy consumption.

Today’s applications that require large capacity memories have high miss rates in DRAM caches which cause these issues to significantly affect the workload’s performance. SRAM caches cannot scale to the capacities required by today’s applications, and thus it becomes imperative to enhance existing DRAM cache designs to address these challenges.

In this work, we introduce TDRAM (Tag-enhanced DRAM), a DRAM microarchitecture specifically tailored for caching purposes.<sup>1</sup> TDRAM enhances HBM3 by adding a set of small low-latency mats to store tags and metadata on the same die as the data mats. These mats enable faster access than the data mats through the reduction of wordline and bitline lengths. The additional on-die storage is sufficiently large to accommodate the tag and metadata for all DRAM cache lines; thus, the tag store scales with the data capacity. By placing the tags in separate mats on-die, TDRAM enables rapid on-die tag checking which reduces the latency for demand misses, mitigates contention on the DRAM read queue, and decreases wasted data transfers (and thus energy).

This work is licensed under the Creative Commons Attribution 4.0 International (CC-BY-SA 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate websites with the appropriate attribution.

<sup>1</sup>In the recent past, it was uneconomical for DRAM manufacturers to modify the core DRAM microarchitecture. However, specialty DRAMs are becoming increasingly common (e.g., Samsung’s Aquabolt [4], Micron’s Automata Processor [67] and RLDRAM [11]).

TDRAM extends HBM3’s interface in three ways: (1) It adds a unidirectional hit-miss (HM) bus to its interface to transfer the tag check result and metadata to the controller. (2) It adds two new commands to the HBM3’s protocol: *ActRd* and *ActWr*, which access both tag and data mats in lockstep. These commands check the tag for the block and only send data to the controller when it is needed. (3) We add a *flush buffer* to store conflicting dirty data from write misses which eliminates costly turnaround delays on the data bus and immediate cache line data transfer to the controller for write requests. The overhead of this new design compared to HBM3 is 192 pins (out of 1,972 existing pins, a 10% increase) and 8.24% total die area.

TDRAM further improves performance by implementing *early tag probing*, which opportunistically performs tag checks (without data access) in otherwise unused command and HM bus slots. Tag probing returns early hit-miss and status indication of a demand access, allowing certain operations (e.g., main memory access for read demand misses) to start earlier. Early tag probing also reduces request queue occupancy time by removing misses from the queue early, allowing other demands to proceed with fewer stalls.

TDRAM is orthogonal to many prior works that focus on improving DRAM caching performance by adding predictors, prefetchers, tag caches, modifying coherence protocols, bypass policies, and other application-specific mechanisms [23], [24], [28], [36], [68]. TDRAM is designed in a way that all of these techniques can be applied on top of it to further improve its performance. Overall, TDRAM enables a perfectly scalable HBM-based cache with a cohesive caching paradigm akin to processors’ SRAM-based caches. Thus, TDRAM focuses on optimizing the core and fundamental DRAM cache operations by eliminating inefficiencies in existing designs.

We have extensively modeled TDRAM in the gem5 simulator [51], for a detailed full-system cycle-level timing analysis. Our evaluations using scientific and graph analytics applications with large memory footprints, have shown TDRAM provides 2.6× faster tag check and 1.2× speedup and at least 21% less energy consumption, compared to the commercial and research designs such as Intel’s Cascade Lake and Alloy.

In this paper, we make the following contributions:

#### Microarchitecture

- We propose a new *HBM3*-based DRAM microarchitecture, TDRAM, designed for caching with on-DRAM-die tag management, to enable *perfectly scalable* DRAM caching where the tag storage scales with data capacity.
- We extend the HBM3 interface with a unidirectional *Hit-Miss bus* to transfer tag check results and metadata from DRAM to the controller, decoupling them from data transfer.
- We add a *flush buffer* to hold conflicting dirty data from write misses which eliminates both costly turnaround delays on the data bus and immediate cache line data transfer to the controller for write requests. TDRAM opportunistically sends them to the controller when data bus is idle or in read-state.

#### Protocol

- We integrate two new commands to HBM3’s protocol, *ActRd* and *ActWr*, which enable parallel independent access to both tag and data banks. The protocol *selectively streamlines data* to the controller only when necessary based on tag comparison, reducing bandwidth bloat.
- We propose *opportunistic early tag check* mechanism in unused HM and command bus slots, to minimize queueing delay for tag check, reducing buffer contention and optimizing miss latencies. This mechanism does not access data banks. If this tag check results in miss, then TDRAM initiates backing store access immediately, if necessary; and avoids future cache line data access, if not necessary.

#### Evaluation

- We demonstrate DRAM caching using existing designs cause *slowdown* while TDRAM provides 1.11× speedup.
- We show TDRAM reduces energy consumption by 21%, since its conditional data response to the controller removes wasted data transfers.
- We analyze performance of DRAM caching in disaggregated systems with remote main memories and show TDRAM provides 1.14× speedup compared to the state-of-the-art commercial DRAM cache, for low miss ratio applications.

## II. BACKGROUND AND MOTIVATION

### A. HBM3 Architecture

This section provides an overview of HBM DRAMs as the basis of TDRAM. HBM provides the highest bandwidth and capacity of any single DRAM package in mass production. HBM3 DRAMs stack multiple DRAM die into a single package, and can support up to 64 GiB capacity using 12 to 16-high stacks [47]. These devices provide up to 1024 GiB/s of bandwidth when running at 8 Gbps across 16 independent channels with 64b data (DQ) and 10b Row command (R) and 8b Column command (C) buses. Each DQ channel can be split into two 32-bit pseudo-channels (PCs) that share the same R and C buses, with each PC providing 32B access granularity [9], [14]–[16]. Each channel includes 38 additional signals for clocks, strobes, ECC, redundancy, and other functions. The wires connecting the high pin count interface between the DRAM and host (1024 DQs, 288 command/address (CA) buses, and more than 650 pins for additional channel and global functions), are implemented in technologies such as TSMC’s InFO or silicon (e.g., a silicon interposer) to support the high pin and trace densities required for this technology.

HBM DRAMs are hierarchically organized, storing data bits in arrays of capacitors. DRAM bit cells are grouped into *rows* or *pages*, with multiple *rows* aggregated into *mats*, and *mats* organized into 2D structures called *banks*. Decoded row and column addresses identify bits within a bank. Multiple banks form a *bank group* that share some resources. Back-to-back accesses to the same bank group require longer latencies to allow these shared resources to free up, while accesses to different bank groups enable lower latencies.

TABLE I: Comparison of TDRAM with related work. Tag storage can be on the processor die (e.g., in SRAM or eDRAM) or off die (e.g., on DDR or HBM device). The tag check can occur before the off-die memory controller, within the memory controller, or in the off-chip data storage. The area overhead on the processor die depends on the tag storage location. Some designs require extra hardware (e.g., tag cache, prefetchers) or significant changes to the coherence protocol (Extra hardware row). The tag capacity can either scale with the amount of data stored in the cache or not. Some designs require multiple (tag and/or off-chip) device accesses to complete read and write hits. TDRAM has low are overhead, no extra hardware, scales tag capacity with data capacity, and requires only one device access for read and write hits.

Tag storage maintained on:	On the processor die		Off the processor die			
			Tag & data in the same row	Tag & data in separate storage		
Tag storage type	SRAM	eDRAM	DRAM	RRAM	DRAM	
Examples	[33,41,52,72]	eTag[68]	[23,26-28,30,34,35,40,48,49,53,54,56,58,62,66,69]	R-Cache [24]	[36]	<b>TDRAM</b>
Tag check	Before MC <sup>1</sup>	Before MC	In MC	In RRAM	In DRAM	<b>In DRAM</b>
Processor die area	High <sup>2</sup>	High	Low	Low	Low	<b>Low</b>
No Extra HW	✓	×	×	×	×	✓
Tags scale with data	×	×	✓	✓	✓	✓
Low hit/miss latency	✓	✓	×	×	×	✓

Notes: <sup>1</sup> MC: memory controller, <sup>2</sup> Some prior work propose 3D-based solutions.

A command decoder receives commands and addresses from the memory controller over a CA bus. When data is read from a DRAM, an activate command provides a row address to move all bits in one row of a bank to sense amplifiers (or sense amps). A separate read command provides a column address to select a subset of the bits from the sense amps to be returned across the DQ bus. Write commands work similarly, providing data to be written into the DRAM.

#### B. Tag Management in Existing DRAM Caches

Numerous studies have investigated the management of tag and metadata (referred to collectively as *tag*) in hardware-managed DRAM caches and Table I compares some of these prior works to TDRAM. Also, DRAM cache products, like Intel’s Xeon series, offering gigabytes of DRAM cache, are available in the market. In terms of storage size, a 64 GiB block-based DRAM cache requires 3 GiB of storage for 3B tag per 64B blocks. This is far beyond the cache sizes in high-end CPUs by AMD (384 MiB in EPYC 9654P [1]) and Intel (105 MiB in Xeon Platinum 8468H [7]) today. While SRAM caches are hitting a scaling wall, tags-in-SRAM solutions (e.g., on processor die) will add to the area overhead and price [33], [41], [52], [72]. Moreover, solutions that put tags on the processor die, e.g., eTag [68], severely limits scalability of DRAM cache capacity by tying it to the tag capacity that processor chip can provide. Where Aurora supercomputer [2] offers over 64 GiB HBM per CPU, moving towards highly-scalable HBM-based DRAM caches is the direction of future.

Previous studies suggest storing tags in the same cache line that data resides [23], [26]–[28], [30], [34], [35], [40], [48], [49], [53], [54], [56], [58], [62], [66], [69]. For instance, in Alloy cache instead of accessing 64B block, 72B (plus 8B ignored) must be accessed, causing misalignment in column layout within DRAM rows and leaving unused bits that causes scalability overhead. In Intel’s Xeon Series (e.g., Cascade Lake), tags are stored in the unused bits of ECC in commodity DRAM devices [38]. However, ECC bits are not designed for this purpose. These designs depend on a DRAM read to access

tag which can create a serialization of tag and data access (e.g., in write-hits), increasing bandwidth bloat.

Some prior work proposed to store tags in separate storage on DRAM die. R-Cache uses resistive RAM for tag storage [24]. Since tag access is on the critical path (i.e., data access in DRAM cache depends on the tag comparison result), the tag read and update latencies must be minimized. Resistive RAM cannot provide the required speed. Moreover, tag and metadata are subject to frequent updates, which can wear out resistive RAM quickly. Other works suggested DRAM-based tag storage [36], [68]. These works optimize tag management and data layout in DRAM rows for associative caches which require multiple tag comparisons, and both activate tag and data regions in parallel. However, they have to delay the start of column operations till tag comparison logic finds the corresponding column, which in fact internally ties the data access to the tag access. Most importantly, they fall short in providing an efficient mechanism to handle write misses to dirty cache lines that necessitates a data read before write for correctness. As a result, these solutions have to rely on either speculative mechanisms (e.g., predictors and DRAM bypassing unit with application-specific hardcoded designs [36]), or require deep cache coherence protocol changes (e.g., for the LLC to send clean writeback messages to the memory controller) [50]. Changes to the cache coherence protocol ties the designs of coherence protocol, the memory controller, and the DRAM, which we avoid in this work.

#### C. Opportunities to Improve DRAM Cache Designs

DRAM cache hit/miss latencies vary based on system design, influenced by factors like tag storage and access mechanisms. Tag check latency is always on the critical path of servicing a memory demand. The prior designs storing tag in the DRAM cache line [42], [58] or in the ECC bits of DRAM (e.g., Intel’s Cascade Lake), must perform a DRAM read that retrieves the cache line’s tag and data, simultaneously. These designs aimed for parallelizing the tag and data access

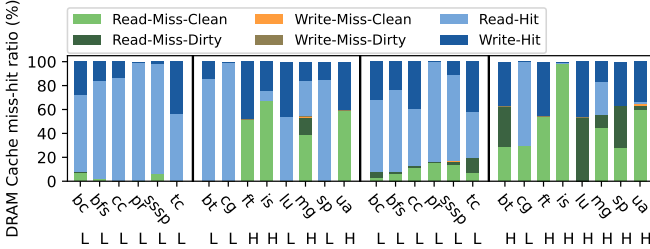


Fig. 1: The breakdown of hit and miss ratios of DRAM cache. The letters show high or low miss ratio.

to improve hit latency. Through some experiments we show their inefficiencies.

For the experiments we consider Intel’s Xeon Max series [25], rounded up to 64 cores and 64 GiB of HBM (1 GiB per core, in DRAM cache mode). Using gem5 simulator [51], we modeled  $\frac{1}{8}$  of the target system. *In this work, for the baseline to represent existing designs storing tags in DRAM, we choose Intel Xeon series DRAM cache (e.g., Cascade Lake), recognized as the state-of-the-art real commercial product in this domain which implements a direct-mapped insert-on-miss cache [38].* We executed 28 HPC multithreaded workloads from GAPBS [22] and NPB [21]. These workloads’ memory footprints are 0.1–80 GiB, while the total DRAM cache size remains at 8 GiB. We conducted full-system simulation and we employed LoopPoint technique [60] for precise control over workloads execution phases. Notably, our experimental setup differs from previous works, allowing us to uncover pathological pathways not observed in prior studies. §IV provides more details about our methodology.

**1) DRAM Cache’s Increased Hit Latency:** In DRAM caches that use standard DRAM devices with tags stored in the device, the cache hit latency almost equals the DRAM read latency for LLC read misses. For LLC writebacks, i.e., evicting dirty data from LLC, the hit latency involves a DRAM read latency (to retrieve the tag) and then a DRAM write is issued (for writing incoming data into the cache). Previous efforts aimed to parallelize tag and data accesses for each memory request [58]. However, this parallelization is compromised for all LLC writebacks, even those hitting on the DRAM cache. The reason is the controller must read the tag (in which the data is also read) for tag comparison and it will not issue the incoming data write into the cache until after DRAM read is ready in the controller and tag check is done. Note that the data read happens because in the commodity DRAMs the only mechanism to fetch tag is to read the whole cache line data, regardless of incoming write data size. In this case, the DRAM read remains on the critical path of write demand and causes access amplification (bandwidth bloat). Prior work reported DRAM caches’ access amplification can reach to as high as 5 accesses [38]. Figure 1 illustrates the miss ratio percentage of the DRAM cache and its breakdown in our experiments. As shown in dark blue color, in the majority of workloads, the number of LLC writebacks hitting on the DRAM cache is significant, potentially affecting the hit latency of the DRAM

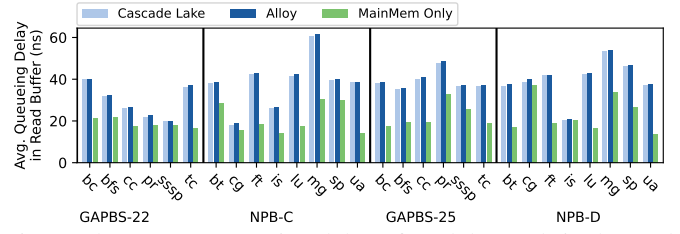


Fig. 2: The average queueing delay of read demands in the read buffer of controller, in Intel’s Cascade Lake and Alloy DRAM caches, compared to the system having a main memory only (no DRAM cache). This time marks the waiting time requests spend in the buffer before accessing the memory.

cache.

**2) DRAM Cache’s Increased Miss Latency:** The process time of a memory demand is sum of two components: *QueueingDelay* + *MemoryAccessLatency*. Memory access latency is the time it takes since a read/write command is issued for a demand, until the data is available on the DQ bus. The queueing delay marks the waiting time requests spend in a buffer before a read/write command is issued. In current DRAM caches, all read and write requests (i.e., LLC’s read misses and writebacks) must undergo a DRAM read to fetch the tag. The controller handles these DRAM reads, including those for LLC writebacks, in the same read buffer. This arrangement heightens contention in the buffer, increasing queueing delay of all demands.

Figure 2 displays the average queueing delay of all DRAM reads in state-of-the-art DRAM caches, compared to a system solely equipped with main memory (no DRAM cache). As depicted, the bars are significantly higher in the DRAM cache system compared to system relying solely on main memory. The main reason is that every read and write demand has to start by reading a tag in DRAM cache, which increases contention in the read buffer and bank conflicts when the tag reads occur.

This extended delay directly impacts the tag check latency for all read demands that miss in the DRAM cache, leading to a delay in fetching the missing line from the main memory where the response to LLC resides. In simpler terms, it extends the miss latency of the DRAM cache. This latency is crucial for LLC read misses, as their processing time in the DRAM cache contributes to the miss penalty of LLC, which is observed by the CPU. This directly influences the overall system throughput. As Figure 1 shows, the number of read misses (in dark/light green) in the DRAM cache is significant.

**3) Increased Bandwidth Bloat and Energy Consumption:** The read data during tag access in the DRAM cache is only beneficial for read demands hitting the cache or for demands missing to a dirty cache line. *In cases of read/write misses to a clean line and write hits, the controller discards the read data immediately after tag comparison, serving no purpose.* In such cases, existing DRAM cache designs introduce data movement overheads by: (i) keeping DRAM’s command bus and banks busy, and (ii) occupying data bus for unnecessary data transfers. These extra communications between the DRAM





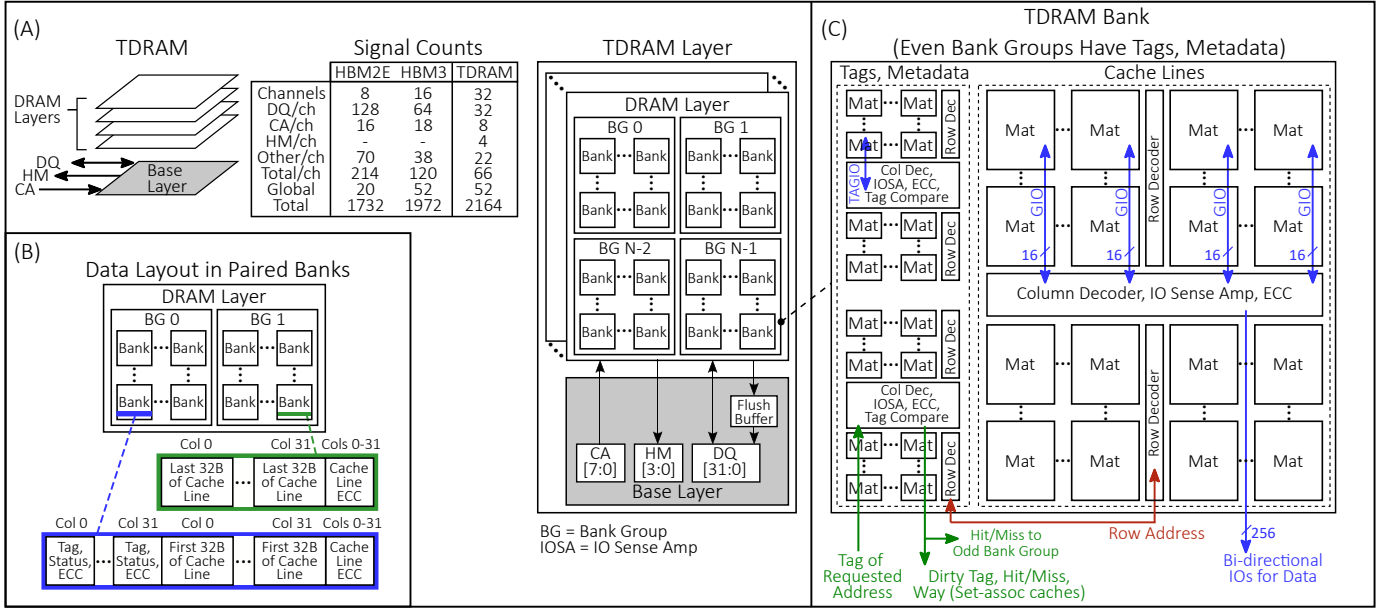


Fig. 4: TDRAM's architecture and bank organization.

## B. TDRAM's Internal Architecture

1) **Data Storage and Access Granularity:** TDRAM maintains the standard bank microarchitecture found in HBM3 devices for data storage. Since server CPUs from Intel and AMD operate on 64B cache-lines, but HBM DRAMs are designed to provide 32B granularity, TDRAM pairs banks in different bank groups and staggers accesses to them to achieve 64B granularity at lower latencies. Figure 4B shows the layout of these paired banks. The controller views the paired banks as one logical bank and schedules accesses accordingly. To simplify the management of paired banks, the controller issues a single command (e.g., activate, read, etc.) and the logic on the base die replicates it, staggering it in time across the bank pair. Pairing banks across bank groups in this way simplifies the controller management since the design eliminates the back-to-back accesses to the same bank group.

2) **On-Die Tag Storage:** TDRAM stores tags and metadata (referred collectively as tag), and their ECC in a separate structure on the same die as the cache line data. TDRAM uses a set of small low-latency mats to store tags to provide fast access. The size of tag storage is much smaller than the size of data storage (about  $\frac{3}{64}$ , i.e., 3B tag per each 64B cache line). The smaller size allows these mats to have shorter wordline and bitline lengths (than the data mats), similar to the design of Reduced Latency DRAM (RLDRAM) [11]. Latency improvements with scaled mats are discussed in [64]. Our design scales the tag mats by  $\frac{1}{2}$  in each direction, reducing the wordline delay time and bitline charge sharing completion time. Centralized decoder and IOSA structures further improve the latency. These low-latency mats allow parallel tag and data lookup, with hit/miss being determined before the data becomes available in the data mats. The special mats are placed at the edge of each bank (Figure 4C).

As an alternative design, the reduced latency tag arrays can also be implemented on a separate die within the TDRAM stack. However, an advantage of placing tags on the same die as the cache-line data is that tag storage scales with data storage. For the remainder of this paper we assume the tags are on the same die as the data.

3) **Metadata Access and Tag Comparison:** We add two new DRAM commands to the HBM3 command set: *activate-read* (ActRd) and *activate-write* (ActWr). When a ActRd or ActWr command is issued to a bank, the tag mats are activated in parallel with the data mats. To avoid sending tags and metadata back to the controller, TDRAM uses on-die tag comparators implemented in the IOSA area of the tag mats to determine whether an access is a hit or a miss. Then, the HM result is routed to the periphery of the chip for output on the HM pins. The HM result is also sent to the column decoders of the data mats where it is used to gate the column decode logic using the same hardware as SALP [44]. If the tag comparison results in a hit or in a miss to a dirty cache-line for read demands, the data is transferred through the DQ bus. If the tag comparison results in a miss to a clean cache-line, the column decode does not happen and no data is transferred on the DQ bus, saving energy. Furthermore, to improve reliability, there are ECC bits for the tag which are analyzed and corrected if needed by on-DRAM-die circuitry as in the baseline HBM3 [9]. Figures 5, 6, 7, and 8 show the timing transactions of read and write operations in TDRAM and are discussed in detail in §III-C.

4) **Direct-Mapped & Set-Associative TDRAM:** The architecture of TDRAM applies equally well to direct-mapped and set-associative caches. On-die tag comparison can provide greater benefits for set-associative caches. In Figure 4, if pairs of bank groups (e.g., 0 and 1, 2 and 3, etc) form two ways of a

set, tag comparisons can be performed in parallel if each way has its own comparator. A signal from the matching way is sent to the internal control logic to select the proper column in the data mats. Implementations without on-die tag comparators send all tags in the set back to the controller, incurring additional latency and energy, and the controller subsequently sends a request for the proper column to the DRAM, again incurring additional latency and energy consumption [48].

5) **Tag Mats Timing Values:** TDRAM architecture minimizes the tags access latencies using small low-latency mats as discussed in §III-B2. *In our evaluation, we use timing parameters for the tag mats that are loosely based on RL-DRAM technology.* We choose to base our timings on public datasheets as the timing parameters for RLDRAM are close to the *proprietary analysis we conducted*. Through discussions with DRAM designers, we validated the internal timing values. Table III shows a list of timing values we used. The RLDRAM spec values (e.g.,  $t_{RL}=15\text{ns}$  and  $t_{RC}=8\text{ns}$ ) match, or are more optimistic than, our values (e.g.,  $t_{RCD\_TAG}+t_{HM}=15\text{ns}$  and  $t_{RC\_TAG}=12\text{ns}$ ). Furthermore, these values and internal TDRAM timings were also correlated with prior work analyzing the use of smaller mats [64].

Additionally,  $t_{HM\_int} = t_{CCD\_L}+t_{HM\_detect}$  (which is a fast equal comparison). Address comparisons are already done in DRAMs today to quickly determine if every row or column address that the DRAM receives is a repaired row or column. We set  $t_{HM\_detect}$  to 0.5ns (one 2GHz clock cycle) for the fast equal comparison based on discussions with expert DRAM designers. The use of  $t_{HM\_int}$  depends on  $t_{RCD}$  since a read operation cannot occur until  $t_{RCD}$  is met. In our design,  $t_{RCD} = 12\text{ns}$  which is longer than  $t_{RCD\_TAG}+t_{HM\_int}=10\text{ns}$ , effectively hiding the tag access and hit/miss detection latency.  $t_{HM\_int}$  was also correlated with prior work [64], which breaks ACT-to-data delay into: 47%-sensing, 26%-address-decode, 20%-MUXing (transfer+rate-conversion), and 7%-IO. The column decode is done in parallel to sensing ( $t_{RCD}$ ) with our ActRd/ActWr commands. Finally, the I/O delay is not relevant to internal timing. Only a portion of the MUXing delay, the delay to move data out of the IOSA, is relevant to internal timing, and this delay is optimized with smaller mats to achieve  $t_{HM\_int}=2.5\text{ns}$ , including HM detect logic.

To ensure dirty data is not overwritten in the SA,  $t_{RL\_CORE}$  (used in write operations as illustrated in Figure 7) needs to be less than or equal to  $\text{intRD-to-WR\_data\_Delay}+t_{BURST}/2=9\text{ns}$ . We performed our evaluation using  $t_{RL\_CORE}=t_{CCD\_L}=2\text{ns}$ .

6) **Tag Storage Area Overhead:** A 64 GiB direct-mapped DRAM cache can support 1 petabyte address space using a 14-bit tag. We assume 3B of tag and metadata for each 64B cache line. Tags are stored only in one bank group of the pair (the even-numbered bank groups), and the result of the tag comparison is communicated to the other (odd-numbered) bank group through an internal bus.

We estimate the die size impact of tag storage, control logic, and on-die comparison as follows. HBM3 DRAMs store an

TABLE II: TDRAM’s cache operations on different accesses.

Cache Access	CMD	DQ Activity	HM Bus	Later Actions
Read hit to clean	ActRd	Hit Data	Hit	None
Read hit to dirty		Hit Data	Hit	None
Read to invalid		None	Miss	Read main mem & fill
Read miss to clean		None	Miss	Read main mem & fill
Read miss to dirty	ActWr	Dirty Data	Miss, Dirty Tag	Read main mem & fill Writeback dirty data
Write to invalid		Wr Data	Miss	None
Write miss to clean		Wr Data	Miss	None
Write miss to dirty		Wr data	Miss, Dirty Tag	Dirty data to flush buffer
Write hit to clean		Wr Data	Hit	None
Write hit to dirty		Wr Data	Hit	None

additional 6B of information (2B metadata and 4B parity) for every 32B of data (i.e., total column size is 38B) across 19 mats as shown by Park et al. [57]. The HBM3 die photo shows that banks (including mats, BLSAs, and Sub-WL drivers) occupy about 66% of die area. The remaining 34% of die area includes shared resources like through silicon vias (TSVs), IOSAs, per-bank group ECC, and column decoders.

For the tag mats, we use four smaller tag mats per data mat to reduce the row cycle time by reducing the bit line and word line lengths. Son et al. show that the overhead of changing the aspect ratio by a factor of 4 is 19% [64]; however, we estimate a more pessimistic 24.3% when we scale by 1/2 in each dimension, based on our discussions with DRAM designers. Additionally, we only need tag mats in the even banks, reducing our overall area overhead. Thus, even banks require 24.3% additional area for the tags and the data banks occupy 66% of the die. So the overall impact on die size is  $24.3\% \times 0.5$  (only even banks)  $\times 0.66$  (area for banks) = 8.02%. We also add additional area for wire routing (for example, to route Hit/Miss signals from the even bank to the odd bank), resulting in our estimate of 8.24% die area impact.

7) **Tag Storage Power Overhead:** On-die tag storage and tag checking increase the power of TDRAM over standard HBM3 DRAMs. However, *TDRAM provides power benefits at the system level*, as (i) tags are not communicated back to the controller to perform the tag check, and (ii) the protocol minimizes data movement amplification, reducing the number of commands and cache line packets sent between the TDRAM and the controller. In the HBM2 generation, 62.6% of the power is spent moving data between the DRAM and the controller [12]. Avoiding data transfers associated with the tag and cache-lines, which are potentially not needed if the tag comparison result is a read miss clean for example, can be beneficial for overall memory system power. §V-C provides an extensive power analysis of TDRAM.

### C. Protocol

TDRAM’s command protocol is similar to traditional DRAM protocols, with modifications to minimize access latency of tags and bandwidth bloat. TDRAM provides new combined *ActRd* and *ActWr* commands that activate a row and read/write a column at both tag and data banks with an auto-precharge for close-page policy. These combined commands include the row and column addresses, bank group, bank, and tag address needed to determine cache hit/miss. Internal state machines in TDRAM handle sequencing and timing of the

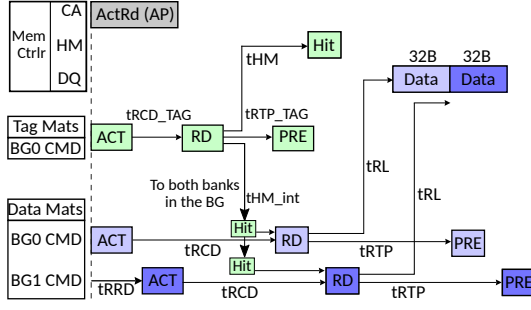


Fig. 5: Timing transactions of a read operations in TDRAM. The timing is the same for a read miss dirty.

activate and column operations to the banks and sense amps. Read and write data appear at fixed offsets on the DQ bus from these commands, as is done in modern DRAMs.

Having a single command to access tag and data banks reduces command amplification and saves energy [6], [10], [11]. Moreover, it simplifies the memory controller since the tag and data banks have the same number of rows and columns. A single address is decoded for both tag and data, allowing their banks to be activated by a single command in lockstep. Regarding the scheduling policy of read/write requests, TDRAM’s controller can adopt any policy such as first-ready first-come first-serve (FR-FCFS). Table II shows the operations the cache performs for each type of access.

1) **Read Operations:** The low-latency tag mats allow hit-miss determination to occur before cache-line data is available. Figure 5 shows the timing transaction of commands involved in read operations of TDRAM. For reads, the HM response will precede the DQ bus transfer, allowing a *conditional response* based on the hit/miss result: (1) on a *read-hit*, cache-line data is returned to the controller. (2) On a *read-miss-clean*, no read command is issued and no cache-line data is returned to the controller. The unused DQ slot can be used to transfer data from the flush buffer (more details in §III-C2) to the controller. (3) On a *read-miss-dirty*, the dirty data is returned to the controller with the same timing on the DQ bus that would have been used to return data on a cache hit. The dirty tag is returned on the HM bus along with the indication that the transaction is a dirty miss. Figure 6 shows the timing transactions of pipelined read accesses in different cases. When the controller receives miss indicator for read requests on the HM bus, it can initiate a backing store read to access the data needed (for the cache line fill and LLC response) before data (if any) arrives at the controller. Early tag probing optimizes this further (§III-D).

2) **Write Operations:** Writes are more complicated than reads since they must avoid overwriting a dirty cache-line with the new data on a write-miss – a rare occurrence but one that needs to be handled correctly. All existing DRAM caches have to serialize the cache-line data read (sending it back to the controller) and the incoming write data, for all write demands. TDRAM avoids this serialization by implementing a flush buffer. The flush buffer is shared among all banks, and operates similar to how a write buffer in a controller stores

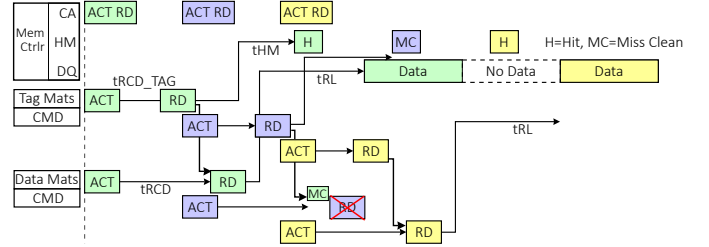


Fig. 6: Timing transactions of pipelined read accesses.

data to be written to the DRAMs. The flush buffer (along with additional logic to support caching) is placed on the existing base layer which already contains logic to support HBM protocol, etc. The base layer is not area limited, thus can support the needed buffer and logic.

TDRAM issues an ActWr command that initiates an internal tag and data access. Once the tag comparison result arrives to the data banks, in case of hit and miss-clean, only an internal write command is issued. If the tag check indicates a miss-dirty, an internal read command followed by an internal write command is issued. Figure 7 shows the sequence of these commands. TDRAM places the read dirty data into the flush buffer and then writes the new data to the DRAM. The flush buffer needs to be sized large enough such that the controller does not need to interrupt a sequence of cache writes for the sole purpose of emptying a full flush buffer, which would require insertion of a full DQ bus turnaround from write to read and then back to write direction. When this occurs, the turnaround delays are optimized since the flush buffer is not in the DRAM core and traditional internal resource conflicts are avoided. Since write-dirty-miss is expected to be a relatively rare event, the flush buffer can be sized modestly (16 entries or less in our simulations) to eliminate virtually any need to require a forced emptying of the flush buffer. It is true that there still will be a small read-to-write turnaround internally to support moving the dirty data from the DRAM bank to the flush buffer, but the much larger turnaround on the DQ bus and on to the controller, can be avoided. A sequence of cache writes from the controller would not experience any delay on the DQ bus due to the write-dirty miss. Direct RDRAM uses a similar approach implementing a Write Buffer and a Write/Retire mechanism to minimize turnarounds due to resource conflicts in the DRAM core [3]. Next we explain how TDRAM opportunistically returns the dirty data in flush buffer when DQ bus is idle or in read-state.

**Unloading the Flush Buffer:** TDRAM transmits the dirty data in the flush buffer to the controller opportunistically or on-demand, as follows: (i) when the DQ bus is idle, such as during *refresh* operations, (ii) in *read-miss-clean* accesses in which DQ is in read-state and is not used for data transfer, and (iii) if the flush buffer becomes full, the controller sends explicit read from flush buffer commands, transmitting multiple entries as a group to amortize any bus turnarounds. The controller has a global knowledge of the addresses in the flush buffer. If the DRAM cache receives a read request to any of



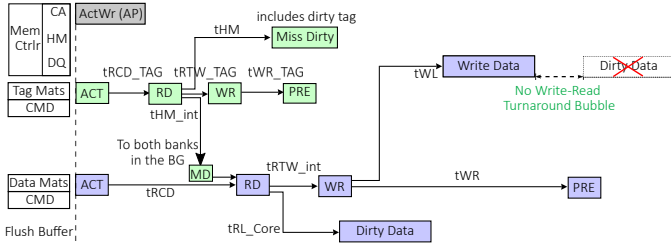


Fig. 7: Timing transactions of write operations in TDRAM.

the addresses in the flush buffer, the controller will get the data from the buffer. In case of a write demand to an address in the flush buffer, the incoming write demand will proceed into the DRAM cache and the controller removes the older data from the flush buffer. Our analysis (§V-D) has shown if we assume 16 entries for the buffer, transferring during read-miss-cleans and refresh cycles, prevents its overflow and any explicit read command to empty the buffer.

#### D. Early Tag Probing Optimization

The TDRAM architecture and command bus have additional unused bandwidth, because: (i) the timing parameters of the tag banks is shorter than the data banks; thus, the tag banks's busy-time is less than data banks, and (ii) the size of the packets transferred on HM bus (few bytes) is much smaller than DQ bus (64B), while both buses work at the same frequency. We use this unused bandwidth for early tag probing, in which the controller can query the status of a cache-line and get an *earlier hit/miss determination* so that following actions (e.g., read main memory for read misses) can begin earlier. Figure 8 illustrates a set of pipelined read transactions. While the data bus is fully occupied by back-to-back data transfers, the CA bus and HM bus are not. Tag probing allows a request to be sent over unused CA bus cycles to perform a tag access and comparison.

1) **Probing Mechanism:** Tag probing only involves accessing the tag mats and returning a result on the HM bus to the controller and does not access the cache-line data. We refer to commands that can access both tag and data, and transferring them on the HM and DQ bus as MAIN slot commands, and tag probe requests that only access the tags and transmit status on the HM bus as PROBE slot commands.

While TDRAM without probing accelerates the tag check through fast tag bank access, the probing mechanism aims to reduce the tag check latency by *minimizing the queue occupancy time of the requests waiting to be scheduled for tag access*. For instance, if the probing indicates a miss-clean for a read demand, the request can be removed from the read queue as soon as the tag check result arrives to the controller on HM bus. The early tag probing lowers the contention in the read buffer, requiring fewer entries and reducing the average queueing delay. Moreover, if a tag probe for the read requests results in a miss, the main memory access starts earlier than if the system waited for the MAIN slot for tag check. A future MAIN slot can then be used for the cache-line fill using the data returning from the main memory.

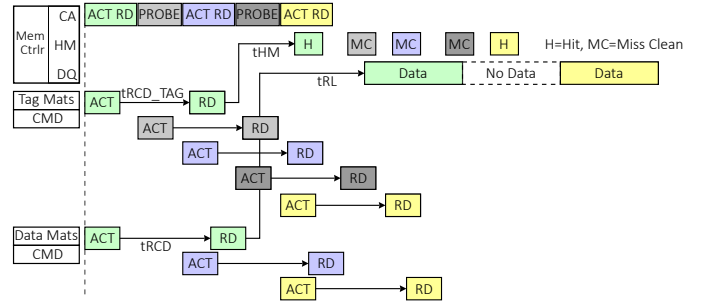


Fig. 8: The timing transaction of early tag probing. The tag check results transfer from tag banks to data banks is left out to make figure clear.

2) **Selection Policy:** Once the controller finds a PROBE slot, amongst all tag check requests that can be issued at that time (i.e., no bank conflict), it picks the *youngest* request to minimize the average queueing delay in the controller. Even though the write packets can also use probing, TDRAM focuses on using these slots for read requests to reduce potential bank conflicts induced by early tag probing. Our analysis has shown that the probing-induced bank conflicts are not common (less than 1% of total demands).

### IV. EVALUATION METHODOLOGY

#### A. Modeled System for Evaluation

Many of the previous studies on DRAM caches [28], [42], [43], [58], [70] rely on trace-based or functional-first simulators, which might not faithfully simulate the behavior of applications that take different paths depending on I/O or thread timings [29]. In contrast, we use an execute-in-execute simulator *gem5* and use full-system simulation. Notably, prior DRAM cache research often omits full-system simulations, failing to capture OS effects. Bin et al.'s work demonstrated that OS kernel bottlenecks can degrade memory access latency in DRAM cache systems [32].

*gem5* implements off-package memory systems via two event-driven components: (i) memory controller responsible for receiving demands from the CPU/LLC, enqueueing them into appropriate queues, and scheduling them to access the memory device, and (ii) a memory interface that handles device-specific timings and operations and communicates with the memory controller [31], [37]. We extended *gem5*'s memory system and implemented TDRAM device and DRAM cache controller [20]. This memory interface uses the DRAM timing parameters listed in Table III. We explained the timing values setup for the tag banks of TDRAM in detail in Section III-B5.

We have integrated alternative DRAM cache designs into *gem5* to assess the performance of TDRAM cache: **Cascade Lake:** As the state-of-the-art commercial hardware-managed DRAM cache, we use Intel's Cascade Lake model to establish a baseline for our evaluation. This is a block-granule direct-mapped insert-on-miss cache storing tag and metadata in DRAM. **Alloy:** This is a research proposal DRAM cache designed specifically to reduce hit latency [28], [58]. We chose Alloy since has the most similar design principles to

TDRAM. We do not consider predictor parts of Alloy, as they are orthogonal to our work (explained in §II-B). To model Alloy’s 80B burst size, we have proportionally increased the corresponding timing parameters of TDRAM (e.g., tBURST, etc.). **TDRAM-NP**: TDRAM no probing (NP), implements all the microarchitectures we described for TDRAM in this work, except the early tag probing. **TDRAM**: implements all the optimizations described for TDRAM, including early tag probing. **Ideal**: We consider an ideal cache which has an architecture similar to TDRAM and knows hit/miss and metadata status in zero latency and overhead.

In order to hold a fair comparison between our approach and other DRAM caching protocols, *we use the same timing parameters for modeling the DRAM device unless a parameter does not apply to a caching protocol (e.g., tRCD\_TAG in Table III is only used for TDRAM cache)*. We modeled  $\frac{1}{8}$  of a target system similar to Intel’s Xeon Max series [25], rounded up to 64 cores and 64 GiB of HBM (1 GiB per core, in DRAM cache mode). Table III shows the detailed parameters of the modeled system.

TABLE III: System Configurations

Processors	
Number of cores	8
Frequency	5 GHz
On-chip Caches	
Private Inst.	32 KB
Private Data	512 KB
Shared LLC	8 MB
DRAM Cache Controller	
Read & Write Buffers	64 entries each
Writeback Buffer	64 entries
Controller latency	20ns round-trip
Shed. Policy	FR-FCFS
DRAM Cache (TDRAM)	
Capacity	8 GiB (8 channels)
Peak BW	32 GiB/s per channel
Read/Write Buffer	64 entries each
Main Memory (DDR5)	
Capacity	128 GiB (2 channels)
Peak BW	32 GiB/s per channel
Read/Write Buffer	64 entries each
Timing Parameters (ns) (same for all evaluated DRAM cache designs)	
Clk=2 GHz, data rate = 8Gbps, close page, RoCoRaBaCh, tBURST = 2, tRCD = 12, tRCD_WR = 6, tCCD_L = 2, tRP = 14, tRAS = 28, tCL = 18, tCWL = 7, tRRD = 2, tFAW = 16, tRL_core = 2, For Tag Banks in TDRAM Architecture only: tHM = 7.5, tHM_int=2.5, tRCD_TAG = 7.5, tRTP_TAG = 2.5, tRRD_TAG = 2, tWR_TAG = 1, tRTW_TAG = 1, tRC_TAG = 12	

## B. Benchmarks

Since we focus on large-scale computing systems for our DRAM cache design, we picked multithreaded workloads with memory footprints larger than the DRAM cache from NPB [21] and GAPBS [22] that are used to evaluate high-performance computing systems. Many past studies often use copies of benchmarks across multiple cores, neglecting inter-thread dependencies apparent in real-world workloads. In contrast, we leverage multithreaded workloads to fully utilize simulated cores, enhancing realism. We utilize the C and D class of NPB workloads and a synthetic graph for the GAPBS workloads with 22 and 25 vertices as inputs. Note: the performance of same workload at different classes or inputs must not be compared together, as the workload has different execution phases in different classes (§IV-C) that remains

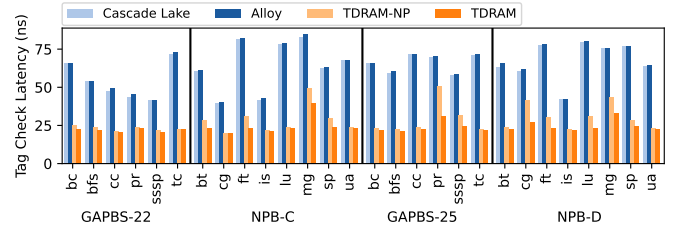


Fig. 9: Tag check latency comparison, the time it takes once controller receives a demand until the tag check result is ready. TDRAM is 2.6× and 2.65× faster than Cascade Lake and Alloy, respectively.

the same across different microarchitectures we are analyzing. Thus, they should be seen as 28 separate workloads.

The working-set sizes of these workloads range from a few hundred megabytes to tens of gigabytes, giving different miss ratios in the 8 GiB DRAM cache. Figure 1 shows the miss ratio of these workloads. We grouped our applications based on their miss ratios: (i) below 30% miss ratio naming them low-miss-ratio, and (ii) above 50% miss ratio calling them high-miss-ratio. There are no workloads in the middle range (i.e., no mid-miss-ratio group).

## C. Methodology for Experiments

Simulating large scale applications like NPB and GAPBS to completion is impractical, necessitating focused execution segments for each benchmark. Measuring work in such applications is complex due to thread interference and extended spin loops. Traditional metrics like instruction count can cause misleading inaccurate performance measurements [18]. Instead, we employ a technique similar to LoopPoint, a sampling technique for multithreaded applications, tracking work progress via global loop instruction counts [60]. LoopPoint ensures that we capture the same phases of execution on different target configurations while comparing their performance.

Following is a summary of our evaluation methodology. Per each workload, first, Linux kernel boots on the target system in *gem5*, the program starts and continues until the start of the region of interest (ROI) of the workload using KVM CPU. Beginning at the ROI, we warm up the system including the CPU caches and the DRAM cache to ensure that the cold miss ratios stay a small fraction of the overall miss count (less than 1%). At the end of warmup, we take a checkpoint. This process is done once per workload. Later, we restore from the checkpoint to run all simulations using an out-of-order CPU with different DRAM cache configurations. Using a checkpoint ensures that all runs start at the same system state for a fair comparison across different configurations.

## V. RESULTS AND DISCUSSION

### A. Impact of Optimizing Tag Check Mechanism

Figure 9 illustrates the average tag check latency for TDRAM and TDRAM-NP in comparison to Intel’s Cascade Lake and Alloy caches. Tag check latency is the time once the controller issues a tag check request to the cache until

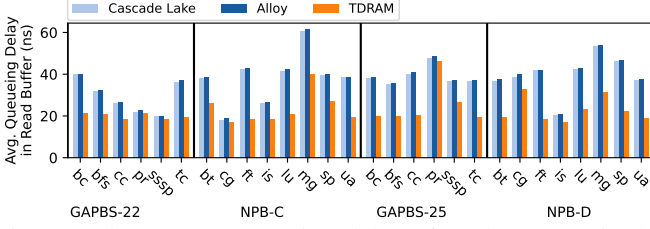


Fig. 10: The average queueing delay of read requests in the controller’s read buffer. TDRAM significantly reduces the queueing delay compared to the other two designs.

the result is ready at the controller. The reported numbers are measured in the controller during simulation and include the queue occupancy time, DRAM cache tag access time, tag compare latency, bus latency, etc. The tag access time in Cascade Lake and Alloy designs consist a read from cache line data while for TDRAM is an access to the separate tag storage of the DRAM cache. All designs use the same timing parameters (Table III) for cache line data access and TDRAM uses validated timings (based on RLDRAM) for tag storage as explained in §III-B5. In baseline design tRCD, tRL(tCL), tBURST timing parameters and for TDRAM tRCD\_TAG and tHM have the most impact in tag check latency. Based on the simulation measurements, in both Cascade Lake and Alloy, the tag check latency falls into 40–85 ns interval, as shown in Figure 9. TDRAM-NP reduces this to 20–50 ns, and TDRAM (with early tag probing), further improves it to 19–39 ns. TDRAM-NP achieves faster hit/miss indication across all applications compared to Cascade Lake and Alloy by parallelizing tag and data access and employing conditional data response. TDRAM which also incorporates early tag probing, further expedites this process by opportunistically performing tag checks. On a geomean, TDRAM’s tag check is 2.6× faster than Cascade Lake and Alloy.

Tag check latency is on the critical path of the hit and miss latencies. Specifically, for read demands that miss on DRAM cache, tag check latency directly impacts the LLC miss penalty, thereby affecting CPU throughput. *Improving the tag check latency accelerates the fetch of missing line from the main memory (response to the LLC), thus, reduces LLC miss penalty.* Figure 9 shows how much faster this main memory read can be issued.

In both Cascade Lake and Alloy, for all demands, the controller issues a DRAM read for tag check, placing them in the controller’s read buffer. I.e., all read and write demands compete in the same queue for DRAM read access in their tag check process. This causes contention in this buffer, increasing queue occupancy time and extends the process time of read requests. Figure 10 shows the average queueing delay of these read requests: the time taken since a read request enters the queue, until the read command for that demand is issued.

Figure 10 shows that the read queueing delay is significantly shorter in TDRAM compared to two other designs, thanks to TDRAM’s early tag probing mechanism. By employing opportunistic tag probing, TDRAM allows a read request to

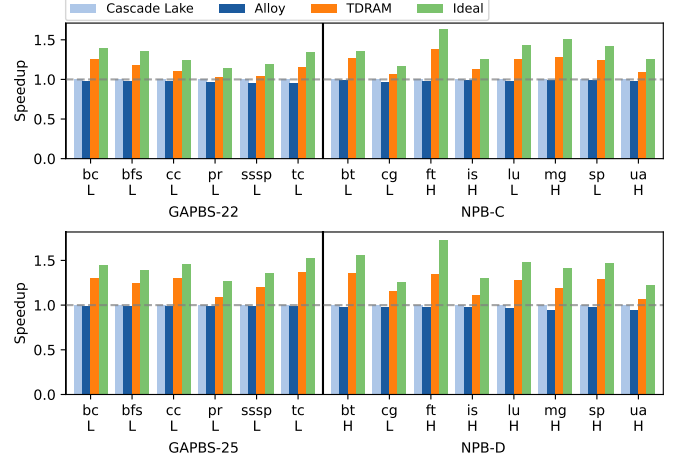


Fig. 11: System’s speedup normalized to the Intel’s Cascade Lake. TDRAM gives 1.2× speedup on a geomean.

leave the read queue as soon as the hit-miss indicator arrives on the HM bus in the case of a miss-clean, without even activating the data bank. This leads to fewer bank conflicts in the system and significantly impacts bank availability, resulting in reduced access time for future demands. Note that in all read-misses in TDRAM, the controller issues a fetch from main memory immediately after seeing the miss indicator on HM bus.

The early-tag-probing benefit is totally workload dependent. Read-misses get the most benefit from this mechanism as it accelerates the off-package memory access with 0 overhead. If we have a smaller DRAM cache size, or workloads with larger memory footprints with the current DRAM cache size (i.e., increasing miss rates), TDRAM will benefit more from early-tag-probing. In other words, TDRAM allows misses to happen while minimizing the miss penalty.

## B. Overall Performance

Figure 11 shows the speedup of TDRAM compared to Cascade Lake and Alloy and Ideal designs. In all workloads TDRAM outperforms Cascade Lake and Alloy, providing a geometric mean speedup of 1.20× and 1.23×, respectively.

As discussed in §V-A, TDRAM effectively reduces tag check latency and read demands queueing delay. This improvement positively impacts the hit and miss latency of the DRAM cache and the miss penalty of LLC. Consequently, the overall performance of the system is enhanced compared to existing designs, as evident in Figure 11. The ideal cache offers tag check results with zero latency, eliminating the need to endure queueing delay and DRAM access latency for tag checks. In essence, the ideal cache sets a performance upper bound for caching, and Figure 11 demonstrates that TDRAM closely approaches this ideal.

Figure 12 compares the speedup of the aforementioned four systems to the same system that has only a main memory (no DRAM cache). As the figure shows for applications with lower miss ratios DRAM caching improves system overall throughput. This improvement decreases as the miss ratio

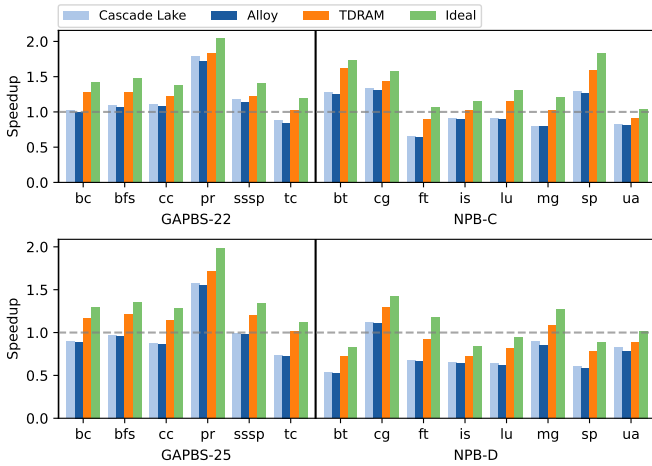


Fig. 12: Speedup of systems with DRAM cache normalized to system without DRAM cache. TDRAM provides 1.11 $\times$  speedup, while Cascade Lake and Alloy cause 8% and 10% slowdown, respectively.

increases due to the miss penalty of DRAM cache that involves main memory access. Analyzing the data, Intel’s Cascade Lake and Alloy caches cause a geomean *slowdown* of 8% and 10%, respectively. In contrast, TDRAM provides an overall *speedup* of 1.11 $\times$ , primarily due to optimizations in hit latency and miss penalty within its protocol.

### C. TDRAM’s Energy Improvement

Prior work has defined bandwidth bloat factor as: total number of bytes moved, over total useful bytes moved [28]. Table IV shows the bandwidth bloat factor for the evaluated designs. Our results shows that TDRAM reduces the bandwidth bloat factor upto 39.9% and 25.1% compared to Alloy and Cascade Lake, respectively.

TABLE IV: Bandwidth Bloat Factor

Design	Low-Miss Ratio	High-Miss Ratio
Alloy	1.68	3.43
Cascade Lake	1.35	2.75
TDRAM	1.13	2.06
<b>TDRAM Reductions</b>		
Over Alloy	32.7%	39.9%
Over Cascade Lake	16.3%	25.1%

As the bandwidth bloat factor increases, more energy is consumed since more data is transferred. TDRAM eliminates unnecessary data transfers, saving energy by reducing the total data movement, while servicing the same number of memory demands as the other two designs. To analyze the energy consumption of TDRAM, we developed an HBM3 power model using HBM2 power data in [55] and scaled it for HBM3 speeds and timings (Table III). Processor interface power is calculated from our validated HBM3 PHY design. Compared to a standard HBM3 DRAM, TDRAM’s power is increased to account for on-die tag storage and associated operations, and both DRAM cache and processor interface power are increased for the additional signals and HM buses and associated logic.

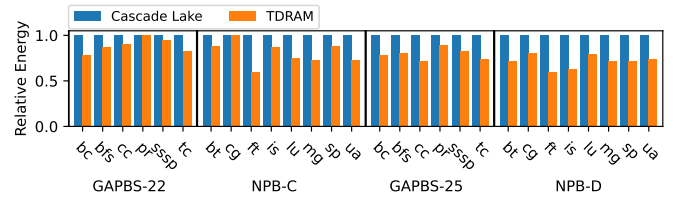


Fig. 13: Relative energy consumption of TDRAM compared to baseline. TDRAM reduces energy by 21% on average compared to the baseline.

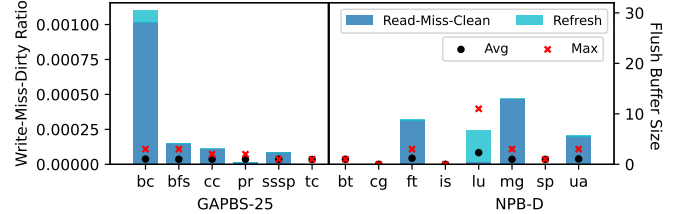


Fig. 14: Flush buffer size sensitivity test. The left axis is the ratio of write-miss-dirty demands out of total requests that DRAM cache received. Each bar is broken into different colors depicting the ratio of dirty data in the flush buffer that was unloaded during read-miss-cleans or refresh cycles. The right axis shows the maximum and average occupancy of flush buffer when it had 32 entries total.

Figure 13 shows the relative energy consumption (power  $\times$  benchmark runtime) of TDRAM and Cascade Lake DRAM cache. We are not showing Alloy’s energy consumption as it is much higher than Cascade Lake. On average, TDRAM saves 21% energy compared to the Cascade Lake due to reducing bandwidth bloat. Applications with high number of write-hits or read/write miss-cleans demands, such as *ft* and *is*, show more energy savings with TDRAM’s.

### D. Flush Buffer Size Sensitivity Analysis

We assessed the sensitivity to the flush buffer size with 8, 16, 32, and 64 entries. The results indicate that the flush buffer consistently avoids becoming full, preventing TDRAM stalls to empty the buffer, except for a minor exception with *lu* in NPB-D with a buffer size of 8. In this case, TDRAM stalled only 13 times, resulting in negligible performance overhead. Figure 14 shows the total write-miss-dirty accesses ratio in NPB-D (which has the largest memory footprint and highest miss rate among all workloads we tested) with a flush buffer size of 32. The figure represents how the dirty data was transferred from the flush buffer to the controller, via read-miss-cleans, refresh cycles, or explicit command to read from flush buffer (which is zero).

Write-miss-dirty accesses in the DRAM cache should be rare, since write-misses in general are the LLC write-backs, that previously required a read in the DRAM cache. This expectation is confirmed by our results, as shown in the left Y-axis the write-miss-dirty ratios are extremely low. The results also shows that most applications heavily rely on read-miss-clean accesses to unload the flush buffer. Notably, *lu* and *bc* efficiently use refresh cycles to unload the flush buffer. This



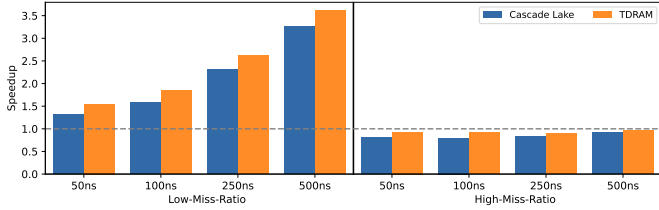


Fig. 15: Speedup of systems with DRAM cache in disaggregated systems normalized to the same system having a remote main memory only (no DRAM cache).

data confirms the effectiveness of TDRAM’s opportunistic behavior in unloading the flush buffer to minimize data transfer overhead. The flush buffer’s average occupancy is 5, with a maximum of 12. Setting the buffer size to 16 prevents TDRAM stalls. Thus, the overhead of flush buffer is minimal.

#### E. Link Latency Case Study

We assume incorporating DRAM caches in disaggregated memory systems where it accesses the main memory through an interconnect such as CXL. We consider round-trip link latencies of 50 ns, 100 ns, 250 ns and 500 ns. Figure 15 shows the speedup of the DRAM cache systems compared to the same system having a remote main memory only (no DRAM cache) and the link latency of this main memory also varies from 50–500 ns.

As Figure 15 shows, in all cases TDRAM outperforms Cascade Lake DRAM cache. We are not showing the Alloy cache data as it was worse performant than Cascade Lake in all cases. For low miss ratio applications, Cascade Lake has an overall speedup of 1.99 $\times$  while TDRAM increases this to 2.27 $\times$ , based on a geometric mean of all speedups. For high miss ratio applications, the speedup is challenging due to the extended miss penalty that includes the interconnect’s latency. Cascade Lake DRAM cache system causes 17% slowdown; however, TDRAM reduces it to 7%.

#### F. Set-Associative TDRAM

TDRAM’s on-die tag comparison is beneficial for set-associative caches since it avoids transferring multiple tags in a set to the controller for comparison. In general, set-associativity is helpful for applications with high number of miss conflicts. However, our analysis showed the HPC workloads we tested barely have miss conflicts on DRAM cache. Thus, they did not get significant performance improvement from set-associativity compared to direct-mapped TDRAM. Our results showed both direct-mapped and 16-way set-associative TDRAMs have similar speedup (over the same system having a main memory only) for tested workloads.

### VI. RELATED WORK

Table I and II-B provides a comparison of TDRAM with prior work. Loh and Hill [48] proposed one of the earliest block-based DRAM cache where tag and data access were stored in the same row with a MissMap to avoid accessing the DRAM cache on predicted misses. Alloy [58] reduces latency

by streaming data and tags together in a single burst. Furthermore, they introduced a memory access predictor, which incurred less overhead compared to the MissMap technique. Retagger [26] uses tags in the controller to mitigate the DRAM row buffer miss cost. RedCache [23] adapts at runtime to start and stop caching for individual blocks. While these works have explored different approaches for storing tags and data in DRAM caches, they all require the tags to be moved to the controller for tag comparison and checks to be performed. R-Cache proposed to use RRAM memory for on-die tag storage [24]. Due to longer latency of RRAM compared to DRAM, it can extend the tag check latency, exacerbating the hit and miss latencies of DRAM cache. In contrast, our work modifies the DRAM microarchitecture to enable tag checks to be performed inside the DRAM, thereby reducing the data movement overhead and improving overall cache efficiency.

The Footprint Cache [43] and Unison Cache [42] blend block and page-based designs in a hybrid architecture. This lowers off-chip traffic compared to page-based designs, with high hits, low latency, and minimal tag overhead. Coarse-grain tracking leads to bandwidth waste and poor utilization of cache capacity in contrast to block-based caches like TDRAM.

Several works [45], [46], [71] combine software and hardware techniques for DRAM caching. Also, Hong et al. proposed a DRAM cache specifically for GPUs working with storage-class memories [39]. We envision potential software/OS integration benefits for TDRAM, as well as GPU-specific changes which are directions we plan to explore.

To our knowledge, the only work to leverage HBM’s embedded logic die for cache management is Stockdale et al.’s [65]. They enhance the base HBM DRAM layer, introducing a cache result signal and reserving one pseudo-channel for tags. The eTag DRAM cache uses eDRAM storage on the processor die, with tag comparison preceding DRAM cache access [68], but eTag cannot scale with increased off-chip capacity as eDRAM size limits the data cache capacity. Hameed et al. proposed a DRAM cache with a separate tag and data storage that relies on a predictor and a Data-Absence-Table [36]. These speculation-based designs are orthogonal to TDRAM. TDRAM minimizes amplification, using the HM bus for cache outcomes and tag transfers, and employs tag probing to mitigate read miss impact. TDRAM puts tags and data on each channel which scales with cache capacity.

### CONCLUSION

In this paper we introduced TDRAM, tag-enhanced energy-efficient DRAM for caching, to optimize caches hit and miss latencies. We showed TDRAM’s 1.2 $\times$ speedup and 21% energy saving over commercial designs (Intel’s Cascade Lake) and research proposal (Alloy). TDRAM can bridge performance gaps between LLC DRAMs and remote memories in heterogeneous/disaggregated systems.

### REFERENCES

- [1] “Amd epyc 9654p.” [Online]. Available: <https://www.techpowerup.com/cpu-specs/epyc-9654p.c2934>

- [2] "Aurora supercomputer." [Online]. Available: <https://wccftch.com/intel-unveils-aurora-supercomputer-specifications-21248-xeon-cpus-63744-gpus-for-over-2-exaflops/>
- [3] "Direct rldram." [Online]. Available: <https://datasheetspdf.com/pdf-file/623377/HynixSemiconductor/HY5R288HC745/1/>
- [4] "Hbm aquabolt: New potential breakthrough memory." [Online]. Available: <https://semiconductor.samsung.com/us/dram/hbm/hbm2-aquabolt/>
- [5] "Iedm 2022: Did we just witness the death of sram?" [Online]. Available: <https://fuse.wikichip.org/news/7343/iedm-2022-did-we-just-witness-the-death-of-sram/>
- [6] "Intel. rldram ii and rldram 3 features." [Online]. Available: <https://www.intel.com/content/www/us/en/docs/programmable/710283/17-0/rldram-ii-and-rldram-3-features.html>
- [7] "Intel xeon platinum 8468h." [Online]. Available: <https://www.itcreations.com/product/140851>
- [8] "Intel's cascade lake: 2nd generation intel® xeon® scalable processors." [Online]. Available: <https://www.intel.com/content/www/us/en/products/platforms/details/cascade-lake.html>
- [9] "Jedec. high bandwidth memory dram (hbm3), jedec standard jesd238, jan 2022." [Online]. Available: <https://www.jedec.org/standards-documents/docs/jesd238a>
- [10] "Micron. async/page/burst cellularram@ 1.0 memory mt45w2mw16gbg." [Online]. Available: [https://www.digchip.com/datasheets/parts/datasheet/301/MT45W2MW16BGB-701\\_IT-pdf.php](https://www.digchip.com/datasheets/parts/datasheet/301/MT45W2MW16BGB-701_IT-pdf.php)
- [11] "Micron rldram 3 specifications." [Online]. Available: [https://media-www.micron.com/-/media/client/global/documents/products/data-sheet/dram/1,-d-,125gb\\_x18\\_x36\\_rldram3.pdf](https://media-www.micron.com/-/media/client/global/documents/products/data-sheet/dram/1,-d-,125gb_x18_x36_rldram3.pdf)
- [12] "Power consumption in hbm." [Online]. Available: <https://semiengineering.com/where-power-is-spent-in-hbm/>
- [13] "Rambus inc. hbm3 controller." [Online]. Available: <https://www.rambus.com/interface-ip/hbm/hbm3-controller/>
- [14] "Rambus inc. hbm3: Everything you need to know." [Online]. Available: <https://www.rambus.com/blogs/hbm3-everything-you-need-to-know/Oct2023>
- [15] "Samsung. hbm3 icebolt: Powering the next frontier." [Online]. Available: <https://www.semiconductor.samsung.com/us/dram/hbm/hbm3-icebolt/>
- [16] "Synopsys. what is high bandwidth memory 3 (hbm3)?" [Online]. Available: <https://www.synopsys.com/glossary/what-is-high-bandwidth-memory-3.html>
- [17] "[tech day 2022] dram solutions to advance data intelligence." [Online]. Available: <https://semiconductor.samsung.com/news-events/tech-blog/dram-solutions-to-advance-data-intelligence/>
- [18] A. R. Alameldeen and D. A. Wood, "Ipc considered harmful for multiprocessor workloads," *IEEE Micro*, vol. 26, no. 4, pp. 8–17, 2006.
- [19] M. Arafat, B. Fahim, S. Kottapalli, A. Kumar, L. P. Looi, S. Mandava, A. Rudoff, I. M. Steiner, B. Valentine, G. Vedaraman, and S. Vora, "Cascade lake: Next generation intel xeon scalable processor," *IEEE Micro*, vol. 39, no. 2, pp. 29–36, 2019.
- [20] M. Babaie, A. Akram, and J. Lowe-Power, "Enabling design space exploration of dram caches for emerging memory systems," in *2023 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 2023, pp. 340–342.
- [21] D. H. Bailey, E. Barszcz, J. T. Barton, D. S. Browning, R. L. Carter, L. Dagum, R. A. Fatoohi, P. O. Frederickson, T. A. Lasinski, R. S. Schreiber *et al.*, "The nas parallel benchmarks," *The International Journal of Supercomputing Applications*, vol. 5, no. 3, pp. 63–73, 1991.
- [22] S. Beamer, K. Asanović, and D. Patterson, "The gap benchmark suite," *arXiv preprint arXiv:1508.03619*, 2015.
- [23] P. Behnam and M. N. Bojnordi, "Adaptively reduced dram caching for energy-efficient high bandwidth memory," *IEEE Transactions on Computers*, vol. 71, no. 10, pp. 2675–2686, 2022.
- [24] P. Behnam, A. P. Chowdhury, and M. N. Bojnordi, "R-cache: A highly set-associative in-package cache using memristive arrays," in *2018 IEEE 36th International Conference on Computer Design (ICCD)*. IEEE, 2018, pp. 423–430.
- [25] A. Biswas and S. Kottapalli, "Next-Gen Intel Xeon CPU - Sapphire Rapids," in *Hot Chips 33*, 2021.
- [26] M. N. Bojnordi and F. Nasrullah, "Retagger: An efficient controller for dram cache architectures," in *Proceedings of the 56th Annual Design Automation Conference 2019*, 2019, pp. 1–6.
- [27] C. C. Chou, A. Jaleel, and M. K. Qureshi, "Cameo: A two-level memory organization with capacity of main memory and flexibility of hardware-managed cache," in *2014 47th Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE, 2014, pp. 1–12.
- [28] C. Chou, A. Jaleel, and M. K. Qureshi, "Bear: Techniques for mitigating bandwidth bloat in gigascale dram caches," *ACM SIGARCH Computer Architecture News*, vol. 43, no. 3S, pp. 198–210, 2015.
- [29] L. Eeckhout, "Computer architecture performance evaluation methods," *Synthesis Lectures on Computer Architecture*, vol. 5, no. 1, pp. 1–145, 2010.
- [30] M. El-Nacouzi, I. Atta, M. Papadopoulou, J. Zebchuk, N. E. Jerger, and A. Moshovos, "A dual grain hit-miss detector for large die-stacked dram caches," in *2013 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2013, pp. 89–92.
- [31] W. Elsasser and N. Nikoleris, "Memory controller updates for new DRAM technologies, NVM interfaces and flexible memory topologies," in *3rd gem5 Users' Workshop with ISCA 2020*, 2020.
- [32] B. Gao, H.-W. Tee, A. Sanaee, S. B. Jun, and D. Jevdjic, "Os-level implications of using dram caches in memory disaggregation," in *2022 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 2022, pp. 153–155.
- [33] M. Ghosh and H.-H. S. Lee, "Smart refresh: An enhanced memory controller design for reducing energy in conventional and 3d die-stacked drams," in *40th Annual IEEE/ACM international symposium on microarchitecture (MICRO 2007)*. IEEE, 2007, pp. 134–145.
- [34] N. Gulur, M. Mehendale, R. Manikantan, and R. Govindarajan, "Bi-modal dram cache: Improving hit rate, hit latency and bandwidth," in *2014 47th Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE, 2014, pp. 38–50.
- [35] F. Hameed, L. Bauer, and J. Henkel, "Simultaneously optimizing dram cache hit latency and miss rate via novel set mapping policies," in *2013 International Conference on Compilers, Architecture and Synthesis for Embedded Systems (CASES)*. IEEE, 2013, pp. 1–10.
- [36] F. Hameed, A. A. Khan, and J. Castrillon, "Improving the performance of block-based dram caches via tag-data decoupling," *IEEE Transactions on Computers*, vol. 70, no. 11, pp. 1914–1927, 2020.
- [37] A. Hansson, N. Agarwal, A. Kolli, T. Wenisch, and A. N. Udupi, "Simulating dram controllers for future system architecture exploration," in *2014 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 2014, pp. 201–210.
- [38] M. Hildebrand, J. T. Angeles, J. Lowe-Power, and V. Akella, "A case against hardware managed dram caches for nvram based systems," in *2021 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 2021, pp. 194–204.
- [39] J. Hong, S. Cho, G. Park, W. Yang, Y.-H. Gong, and G. Kim, "Bandwidth-effective dram cache for gpu s with storage-class memory," in *2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2024, pp. 139–155.
- [40] C.-C. Huang and V. Nagarajan, "Atcache: Reducing dram cache latency via a small sram tag cache," in *Proceedings of the 23rd international conference on Parallel architectures and compilation*, 2014, pp. 51–60.
- [41] K. Inoue, S. Hashiguchi, S. Ueno, N. Fukumoto, and K. Murakami, "3d implemented sram/dram hybrid cache architecture for high-performance and low power consumption," in *2011 IEEE 54th International Midwest Symposium on Circuits and Systems (MWSCAS)*. IEEE, 2011, pp. 1–4.
- [42] D. Jevdjic, G. H. Loh, C. Kaynak, and B. Falsafi, "Unison cache: A scalable and effective die-stacked dram cache," in *2014 47th Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE, 2014, pp. 25–37.
- [43] D. Jevdjic, S. Volos, and B. Falsafi, "Die-stacked dram caches for servers: Hit ratio, latency, or bandwidth? have it all with footprint cache," *ACM SIGARCH Computer Architecture News*, vol. 41, no. 3, pp. 404–415, 2013.
- [44] Y. Kim, V. Seshadri, D. Lee, J. Liu, and O. Mutlu, "A case for exploiting subarray-level parallelism (salp) in dram," in *2012 39th Annual International Symposium on Computer Architecture (ISCA)*, 2012, pp. 368–379.
- [45] Y. Kim, H. Kim, and W. J. Song, "Nomad: Enabling non-blocking os-managed dram cache via tag-data decoupling," in *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2023, pp. 193–205.
- [46] J. B. Kotra, H. Zhang, A. R. Alameldeen, C. Wilkerson, and M. T. Kandemir, "Chameleon: A dynamically reconfigurable heterogeneous memory system," in *2018 51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 2018, pp. 533–545.

- [47] J. Lee and et al., "A 48-gb 16-high 1280-gb/s hbm3e dram with all-around power tsv and a 6-phase rdqs scheme for tsv area optimization," International Solid State Circuits Conference (ISSCC) 2024.
- [48] G. Loh and M. D. Hill, "Supporting very large dram caches with compound-access scheduling and missmap," *IEEE Micro*, vol. 32, no. 3, pp. 70–78, 2012.
- [49] G. H. Loh, N. Jayasena, K. Mcgrath, M. O'Connor, S. Reinhardt, and J. Chung, "Challenges in heterogeneous die-stacked and off-chip memory systems," in *3rd Workshop on SoCs, Heterogeneous Architectures and Workloads*, vol. 20, 2012, p. 12.
- [50] J. Lowe-Power, *On Heterogeneous Compute and Memory Systems*. The University of Wisconsin-Madison, 2017.
- [51] J. Lowe-Power et al., "The gem5 simulator: Version 20.0+," 2020.
- [52] N. Madan, L. Zhao, N. Muralimanohar, A. Udipi, R. Balasubramonian, R. Iyer, S. Makineni, and D. Newell, "Optimizing communication and capacity in a 3d stacked reconfigurable cache hierarchy," in *2009 IEEE 15th International Symposium on High Performance Computer Architecture*. IEEE, 2009, pp. 262–274.
- [53] M. R. Meswani, S. Blagodurov, D. Roberts, J. Slice, M. Ignatowski, and G. H. Loh, "Heterogeneous memory architectures: A hw/sw approach for mixing die-stacked and off-package memories," in *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 2015, pp. 126–136.
- [54] J. Meza, J. Chang, H. Yoon, O. Mutlu, and P. Ranganathan, "Enabling efficient and scalable hybrid memories using fine-granularity dram cache management," *IEEE Computer Architecture Letters*, vol. 11, no. 2, pp. 61–64, 2012.
- [55] M. O'Connor, N. Chatterjee, D. Lee, J. Wilson, A. Agrawal, S. W. Keckler, and W. J. Dally, "Fine-grained dram: Energy-efficient dram for extreme bandwidth systems," in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, 2017, pp. 41–54.
- [56] K. H. Park, S. K. Park, H. Seok, W. Hwang, D.-J. Shin, J. H. Choi, and K.-W. Park, "Efficient memory management of a hierarchical and a hybrid main memory for mn-mate platform," in *Proceedings of the 2012 International Workshop on Programming Models and Applications for Multicores and Manycores*, 2012, pp. 83–92.
- [57] M.-J. Park, J. Lee, K. Cho, J. Park, J. Moon, S.-H. Lee, T.-K. Kim, S. Oh, S. Choi, Y. Choi et al., "A 192-gb 12-high 896-gb/s hbm3 dram with a tsv auto-calibration scheme and machine-learning-based layout optimization," *IEEE Journal of Solid-State Circuits*, vol. 58, no. 1, pp. 256–269, 2022.
- [58] M. K. Qureshi and G. H. Loh, "Fundamental latency trade-off in architecting dram caches: Outperforming impractical sram-tags with a simple and practical design," in *2012 45th Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE, 2012, pp. 235–246.
- [59] A. Raybuck, T. Stamler, W. Zhang, M. Erez, and S. Peter, "Hemem: Scalable tiered memory management for big data applications and real nvm," in *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, 2021, pp. 392–407.
- [60] A. Sabu, H. Patil, W. Heirman, and T. E. Carlson, "Loopoint: Checkpoint-driven sampled simulation for multi-threaded applications," in *2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2022, pp. 604–618.
- [61] D. D. Sharma, "System on a package innovations with universal chiplet interconnect express (ucie) interconnect," *IEEE Micro*, vol. 43, no. 2, pp. 76–85, 2023.
- [62] J. Sim, G. H. Loh, V. Sridharan, and M. O'Connor, "Resilient die-stacked dram caches," *ACM SIGARCH Computer Architecture News*, vol. 41, no. 3, pp. 416–427, 2013.
- [63] A. Sodani, R. Gramunt, J. Corbal, H.-s. Kim, K. Vinod, S. Chinthamani, S. Hutsell, R. Agarwal, and Y.-C. Liu, "Knights Landing: Second-Generation Intel Xeon Phi Product," *IEEE Micro*, vol. 36, no. 2, pp. 34–46, mar 2016. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=7453080>
- [64] Y. H. Son, O. Seongil, Y. Ro, J. W. Lee, and J. H. Ahn, "Reducing memory access latency with asymmetric dram bank organizations," in *Proceedings of the 40th annual international symposium on computer architecture*, 2013, pp. 380–391.
- [65] T. Stocksdales, M.-T. Chang, H. Zheng, and F. Mueller, "Architecting hbm as a high bandwidth, high capacity, self-managed last-level cache," in *Proceedings of the 2nd Joint International Workshop on Parallel Data Storage & Data Intensive Scalable Computing Systems*, 2017, pp. 31–36.
- [66] H. Sun, J. Liu, R. Anigundi, N. Zheng, J. Lu, R. Ken, and T. Zhang, "Design of 3d dram and its application in 3d integrated multi-core computing systems," *IEEE Design and Test of Computers*, pp. 36–47, 2009.
- [67] K. Wang, K. Angstadt, C. Bo, N. Brunelle, E. Sadredini, T. Tracy, J. Wadden, M. Stan, and K. Skadron, "An overview of micron's automata processor," in *2016 International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*, 2016, pp. 1–3.
- [68] K.-H. Yang, H.-J. Tsai, C.-Y. Li, P. Jendra, M.-F. Chang, and T.-F. Chen, "etag: Tag-comparison in memory to achieve direct data access based on edram to improve energy efficiency of dram cache," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 64, no. 4, pp. 858–868, 2016.
- [69] S. Yin, J. Li, L. Liu, S. Wei, and Y. Guo, "Cooperatively managing dynamic writeback and insertion policies in a last-level dram cache," in *2015 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2015, pp. 187–192.
- [70] V. Young, C. Chou, A. Jaleel, and M. Qureshi, "Accord: Enabling associativity for gigascale dram caches by coordinating way-install and way-prediction," in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2018, pp. 328–339.
- [71] X. Yu, C. J. Hughes, N. Satish, O. Mutlu, and S. Devadas, "Banshee: Bandwidth-efficient dram caching via software/hardware cooperation," in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, 2017, pp. 1–14.
- [72] L. Zhao, R. Iyer, R. Illikkal, and D. Newell, "Exploring dram cache architectures for cmp server platforms," in *2007 25th International Conference on Computer Design*. IEEE, 2007, pp. 55–62.