
Bulk Bitwise Accumulation in Commercial DRAM

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Processing-in-memory (PIM) is a promising paradigm for addressing data transfer
2 bottlenecks in data-intensive workloads, particularly in machine learning. Among
3 PIM techniques, Processing-using-Commercial-DRAM (PuCD) offers a practical
4 approach for enabling in-memory computing by employing widely available
5 DRAM modules without hardware modifications. With its massive bit-level par-
6 allelisms, PuCD has a high-performance capability of bulk bit logic operation.
7 However, implementing *accumulation* operations, crucial for machine learning
8 tasks, remains challenging in PuCD. The need for multiple consecutive opera-
9 tions in accumulation leads to increased latency and error propagation. To address
10 these challenges, we propose a novel method for bulk bitwise accumulation us-
11 ing PuCD. As a fundamental building block for our accumulation method, we
12 introduce a novel implementation of the *population-count-of-3* (POPCNT3) opera-
13 tion tailored for commercial DRAM. On top of this, we present a POPCNT3-based
14 bitwise accumulation method that efficiently handles large input sizes, enabling
15 scalable bitwise accumulation for various input sizes. We evaluate the through-
16 put and errors of our approach using commercial DDR4 DRAM modules with an
17 FPGA. The experiments indicate that the throughput improvement is up to 348
18 times over A100 GPU across various input sizes with negligible errors to main-
19 tain the accuracy of machine learning applications. These results demonstrate that
20 PuCD can provide a practical pathway for accelerating machine learning tasks
21 without requiring specialized memory chips.

22 1 Introduction

23 Emerging data-intensive workloads, particularly in machine learning (ML), are increasingly con-
24 strained by data transfer costs rather than computation. Deep neural networks, especially large
25 language models, exemplify this trend with their demand for processing vast amounts of data, often
26 requiring tens of gigabytes of parameters [Brown et al., 2020]. The predominance of multiply-
27 accumulate (MAC) operations of general matrix-vector multiplication (GeMV) in these models,
28 coupled with their low on-chip data reuse, makes off-chip data transfer the primary performance
29 bottleneck [Choi et al., 2023, Wu et al., 2024].

30 To address this challenge, processing-in-memory (PIM) has reemerged as a promising solution. PIM
31 is a computing paradigm that integrates computational capabilities directly within memory devices,
32 utilizing high internal bandwidth and reducing off-chip data movement. While digital PIM [Devaux,
33 2019, Kwon et al., 2021, Lee et al., 2022] incorporates processing units within memory devices,
34 analog PIM [Chi et al., 2016, Seshadri et al., 2017, Eckert et al., 2018] transforms memory cells into
35 computational units. Compared to digital PIM, analog PIM offers higher computational density and
36 energy efficiency, while digital PIM provides greater arithmetic precision and flexibility.

37 While most workloads require high-precision computations, some DNN models, including LLMs,
38 have inherent redundancy that allows them to tolerate lower precision [Ma et al., 2024]. These mod-
39 els can quantize their weights to just a few bits without sacrificing accuracy. This trend towards
40 low-bit arithmetic presents a unique opportunity for analog PIM, which operates with bit-level par-
41 allelism, to efficiently accelerate DNNs. Even more intriguing is the discovery of a technology
42 called Processing-using-Commercial-DRAM (PuCD), which transforms standard DRAM into an
43 analog PIM device [Gao et al., 2019, Olgun et al., 2021, Gao et al., 2022, Yuksel et al., 2024, Yuksel
44 et al., 2024]. PuCD enables massive parallelism and substantial computational power by leverag-
45 ing the density of DRAM, without requiring additional circuitry or changes to the cost-optimized,
46 low-margin DRAM design. By making use of the widespread availability of DRAM and enabling in-
47 memory computation without hardware modifications, PuCD provides a practical solution to bring
48 PIM capabilities into mainstream computing systems.

49 Despite the potential of PuCD, implementing *accumulation* using PuCD presents significant chal-
50 lenges. Accumulation is a fundamental operation in MAC computations, which form the backbone
51 of many ML processes. Unlike multiplication, which typically involves two inputs, accumulation
52 requires a number of inputs proportional to the matrix size to sum up the multiple results. This char-
53 acteristic poses two main difficulties for PuCD implementation: performance issues and accuracy
54 concerns. Executing accumulation with many inputs necessitates a proportional number of logical
55 operations, resulting in slow performance when implemented in PuCD. Additionally, PuCD inher-
56 ently carries computational errors, which can accumulate significantly when multiple operations are
57 performed consecutively, as required in accumulation.

58 To address the challenges of implementing fast and accurate accumulation operations in PuCD, we
59 propose a novel method for bulk bitwise accumulation. Our approach consists of two key com-
60 ponents: First, we introduce an optimized implementation of the *population-count-of-3* (POPCNT3)
61 operation specifically designed for commercial DRAM. POPCNT3 serves as a fundamental building
62 block for our accumulation method, efficiently counting the number of '1' bits in three input bits.
63 Our implementation leverages PuCD to perform POPCNT3 with minimal latency and error propaga-
64 tion. Building upon this optimized POPCNT3 operation, we develop a scalable bitwise accumulation
65 technique capable of handling larger input sizes. This method iteratively applies POPCNT3 operations
66 to groups of inputs, progressively computing higher-order bit positions. By leveraging the efficiency
67 of our POPCNT3 implementation and the parallelism inherent in PuCD, we achieve high-throughput
68 accumulation across various input sizes.

69 The key contributions of this work are as follows:

- 70 1. As a fundamental building block for our accumulation method, we introduce a new imple-
71 mentation of the POPCNT3 operation tailored for commercial DRAM.
- 72 2. We present a POPCNT3-based bitwise accumulation method that efficiently handles large
73 input sizes, enabling scalable bitwise accumulation for various input sizes.
- 74 3. Using DDR4 DRAM modules, we demonstrate up to 348 times higher throughput than
75 A100 GPU across various input sizes with negligible errors to maintain the accuracy of
76 machine learning applications.

77 **2 DRAM Structure and PuCD Technology**

78 Modern computing systems widely employ dynamic random-access memory (DRAM) as their main
79 memory due to its high density and low cost characteristics. DRAM modules utilize a hierarchical
80 structure to enable efficient data management as shown in Figure 1. At the highest level, channels
81 offer independent data paths. Each channel encompasses multiple ranks, with a rank comprising a
82 collection of DRAM chips. Manufacturers divide each chip into multiple banks, which function as
83 independently operable memory arrays. Banks further subdivide into subarrays, with each subarray
84 consisting of a two-dimensional array of memory cells arranged in rows and columns. The memory
85 cell, the fundamental unit of DRAM, incorporates a single transistor and a capacitor. This cell stores
86 data as the presence or absence of an electrical charge. Each column of memory cells is connected
87 to a bitline, which serves as a communication for reading from and writing to the cells. At the end
88 of each bitline is a sense amplifier, a crucial component that detects and amplifies the small voltage
89 differences on the bitline operations. A memory controller governs data access in DRAM. This

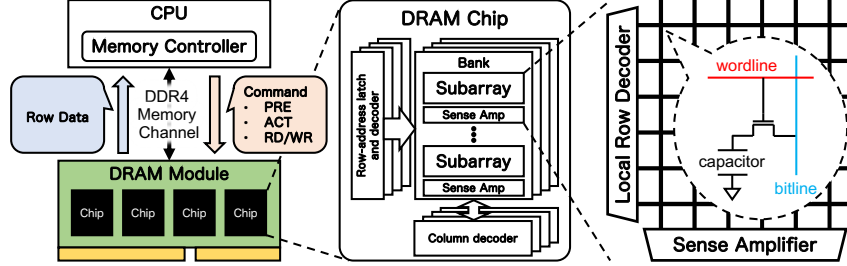


Figure 1: DRAM hierarchy and basic operation. The figure illustrates the structural organization of DRAM from the highest level (memory channel) down to the individual memory cell.

90 controller issues commands in carefully timed sequences. Through this orchestration of commands,
 91 the memory controller achieves efficient data access.

92 Processing-using-Commercial-DRAM (PuCD) is an innovative technique that enables in-memory
 93 computation on existing DRAM modules without hardware modifications by exploiting the charge-
 94 sharing effect between simultaneously activated rows to perform logical operations. Gao et al.
 95 [2019] discovered that logical operations could be performed directly on commercial DDR3 DRAM
 96 modules without circuit modifications. Building on this insight, researchers have demonstrated the
 97 ability to execute RowCopy [Gao et al., 2019, Yuksel et al., 2024], AND/OR [Gao et al., 2022], and
 98 NOT [Yuksel et al., 2024] operations on DDR4 DRAM modules. The charge-sharing effect occurs
 99 when multiple DRAM rows are activated simultaneously, causing the electrical charges stored in the
 100 capacitors of these rows to redistribute across the bitlines. PuCD implements this charge-sharing
 101 effect through simultaneous multi-row activation (SiMRA), a method that involves issuing DRAM
 102 commands in a way that violates conventional timing constraints. This technique activates multi-
 103 ple DRAM rows simultaneously on commercial DRAM modules, inducing charge sharing on the
 104 bitlines. PuCD then utilizes this charge sharing effect to execute *majority-of-X* (MAJX) operations.
 105 MAJX operations, which determine the majority value among input bits, serve as the basic logical
 106 units in PuCD. Prior works build upon MAJX to implement fundamental logical operations such as
 107 AND and OR. For instance, they can achieve AND/OR operations by fixing certain inputs of MAJX to spe-
 108 cific values. By combining these basic operations, PuCD can perform more complex computations.
 109 However, PuCD faces challenges in operational reliability. The reliability of MAJX operations heav-
 110 ily depends on input patterns, with performance degrading significantly when the numbers of '0's
 111 and '1's in the input are closely balanced. Furthermore, complex operations often require numerous
 112 consecutive MAJX operations, potentially leading to cumulative errors and performance deteriora-
 113 tion.

114 3 Bulk Bitwise Accumulation

115 In this section, we introduce a novel approach to perform efficient bulk bitwise accumulation using
 116 PuCD. We present our method in two stages: first, we propose an optimized implementation of
 117 POPCNT3, a fundamental operation that counts the number of '1's in three input bits. Then, we
 118 build upon this POPCNT3 operation to develop a scalable bitwise accumulation technique capable
 119 of handling larger input sizes. By combining these two components, we achieve an efficient and
 120 accurate method for bulk bitwise accumulation for various input sizes.

121 3.1 POPCNT3 in Commercial DRAM

122 We propose a novel method for efficiently implementing POPCNT3 in PuCD. POPCNT3 represents a
 123 bit operation that counts the number of '1's in three input bits and expresses the result in two bits.
 124 POPCNT3 is defined for three input bits A, B, and C using the following logical expressions:

$$125 \text{MSB}(A, B, C) = (A \wedge B) \vee (B \wedge C) \vee (C \wedge A)$$

$$126 \text{LSB}(A, B, C) = (A \wedge \neg B \wedge \neg C) \vee (\neg A \wedge B \wedge \neg C) \vee (\neg A \wedge \neg B \wedge C) \vee (A \wedge B \wedge C)$$

126 where MSB denotes the most significant bit and LSB the least significant bit. Naive PuCD ap-
 127 proaches attempting to implement these logical expressions directly require numerous consecutive

128 AND, OR, and NOT operations. Specifically, LSB calculation necessitates 11 consecutive AND/OR
 129 operations, leading to high latency and reduced accuracy. In PuCD, multiple consecutive logical
 130 operations increase the overall execution time. Moreover, each logical operation carries a success
 131 probability. For instance, even if a DRAM module can perform AND/OR operations with an exact
 132 99% success rate, the probability of all 11 consecutive operations succeeding drops to approxi-
 133 mately 90%. These challenges highlight the inefficiency and accuracy issues of directly implement-
 134 ing POPCNT3 in PuCD.

135 To achieve fast and high-precision POPCNT3 in PuCD, we propose a novel implementation method
 136 that leverages MAJX operations in PuCD. This approach introduces reference rows and performs
 137 MAJX operations combining input and reference rows to directly and efficiently calculate the
 138 POPCNT3 output. Figure 2 illustrates the execution timeline of POPCNT3. Our method prepares
 139 three input rows, two reference rows, and two output rows within the same subarray. The process
 140 unfolds as follows: ① We store input values in the input rows. We can move these values from
 141 other rows within the same subarray using RowCopy. ② We execute MAJ3 on the input rows to
 142 calculate the MSB value and store the result in the MSB row. ③ We reload the input values into
 143 the input rows and copy the NOT of the previous MSB value to the reference rows using NOT. ④ We
 144 perform a MAJ5 operation combining input and reference rows to calculate the LSB value. Table 1
 145 presents the truth table for POPCNT3. The essence of this approach lies in recursively comparing the
 146 number of '1's in the input to identify the output bits in binary representation from the most signifi-
 147 cant bit. We can consider this process as a binary search. Our proposed method executes POPCNT3
 148 with significantly fewer logical operations (two MAJX operations). This reduction in operations sub-
 149 stantially decreases the overall execution time. Moreover, by minimizing the number of consecutive
 150 PuCD logical operations, we also mitigate output accuracy degradation.

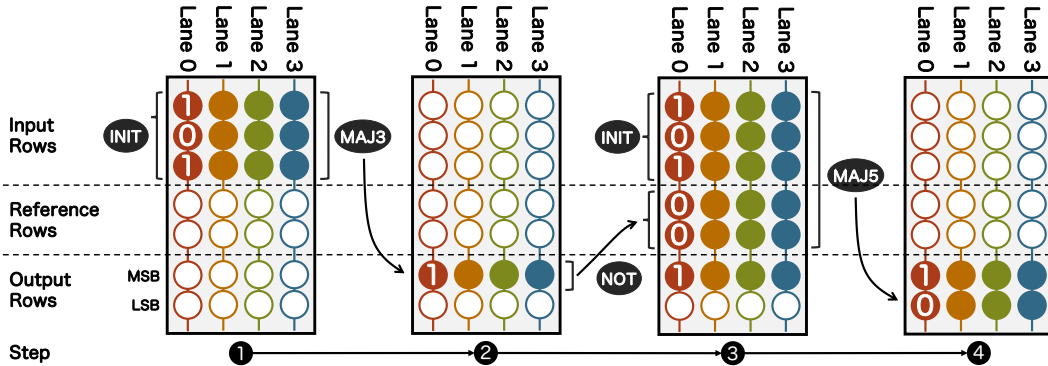


Figure 2: Timeline illustrating data changes in a DRAM segment during a POPCNT3 operation. The figure shows the step-by-step process of executing a POPCNT3 operation using PuCD.

# of 1s in Input Rows		0	1	2	3
Reference Row	Value	1	1	0	0
		1	1	0	0
Output Bits	MSB	0	0	1	1
	LSB	0	1	0	1

Table 1: Truth table for the POPCNT3 operation in PuCD. This table illustrates the relationship between the number of '1' bits in the input rows, the reference row bit, and the resulting output bit pattern in the POPCNT3 operation.

151 3.2 POPCNT3-based Bitwise Accumulation

152 Building upon the POPCNT3 operation, we propose POPCNT3-based bitwise accumulation, a method
 153 for executing bitwise accumulation on a larger number of inputs. This approach enables counting for
 154 a larger number of inputs by iteratively applying POPCNT3. While POPCNT3 generates a 2-bit output

155 from three inputs, we can continue to calculate higher bit position values by repeatedly executing
 156 POPCNT3 on each output bit position. Figure 3 illustrates the procedure for bitwise accumulation
 157 with 15 inputs. First, we initially execute POPCNT3 five times on groups of three inputs from the
 158 15 inputs, yielding five 1st bits and five 2nd bits as outputs. Then, we apply POPCNT3 to three
 159 of the five 1st bits, generating an additional 1st bit and 2nd bit. Finally, we repeat this process,
 160 ultimately obtaining a 4-bit output. Consequently, our approach offers an efficient solution for
 161 bitwise accumulation across various input sizes, leveraging the simplicity and effectiveness of the
 162 POPCNT3 operation.

163 Our POPCNT3-based bitwise accumulation inherits the bulk processing capabilities and bank par-
 164 allelism of PuCD, enabling simultaneous computation across all columns in an array and parallel
 165 execution across multiple banks. This parallelism allows our method to achieve high throughput for
 166 bitwise accumulation operations. However, it's constrained by PuCD's requirement that all inputs
 167 and intermediate results must reside within the same subarray. Given that typical DRAM subar-
 168 rays contain only a few hundred rows, this limits the size of accumulations that can be performed
 169 solely using PuCD. For larger accumulations, alternative or hybrid approaches may be necessary,
 170 highlighting the importance of careful data layout consideration in practical applications of this
 171 technique.

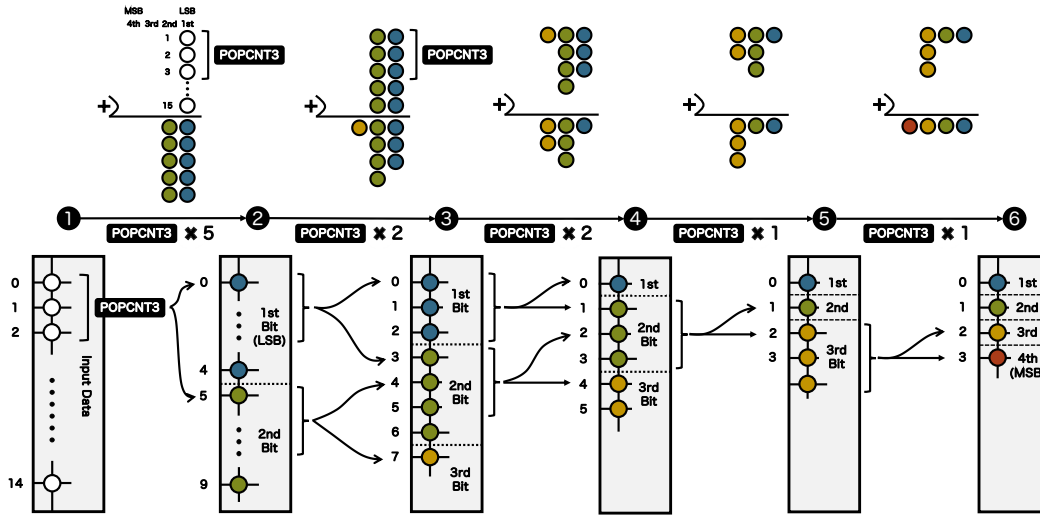


Figure 3: Timeline of a 15-input bitwise accumulation operation in a single DRAM bitline using POPCNT3-based method. This figure illustrates the step-by-step process of performing bitwise accumulation on 15 inputs using our proposed POPCNT3-based approach.

172 4 Results and Discussion

173 We evaluate the performance and accuracy of our proposed bulk bitwise accumulation method using
 174 DDR4 DRAM for different accumulation sizes. We refer to a single bulk bitwise accumulation
 175 computation as a *kernel*. We express the kernel size as $K \times N$, where K represents the input length
 176 and N denotes the bulk data width. To implement our proposed method, we use an FPGA connected
 177 to four DDR4 DRAM modules. We compare the performance of our method with an NVIDIA A100
 178 GPU as a baseline. Further details about our evaluation setup are in Appendix A.

179 Table 2 presents the throughput and latency of our proposed method for each kernel size. Through-
 180 put indicates the amount of data that can be processed per second. For our method, we calculate
 181 throughput when executing in parallel across 64 banks of 4 modules. Latency of our method repre-
 182 sents the time required from the execution of the first DRAM command to the completion of the last
 183 DRAM command for a single kernel execution which comprises multiple POPCNT3 operations. As
 184 shown in Table 2, our proposed method consistently outperforms the GPU across all kernel sizes.
 185 Specifically, our method achieves approximately 348 times higher throughput for the smallest ker-
 186 nel size (7x65536) and still maintains about 27 times higher throughput for the largest kernel size
 187 (127x65536). This significant improvement in throughput demonstrates that our proposed method

188 effectively leverages the bit-level parallelism of DRAM to accelerate bulk bitwise accumulation oper-
 189 ations. The ability to process data at such high rates is particularly beneficial for large-scale MAC
 190 operations, potentially enabling much faster execution of machine learning algorithms and neural
 191 network computations.

Table 2: Comparison of throughput and latency for bulk bitwise accumulation using our method and GPU (A100) for different kernel sizes. Our evaluation covers four dimensions: $K = 7, 15, 31, 63$ and $N = 65536$. We chose these K values to represent realistic accumulations within a subarray, with output bit counts of exactly 3, 4, 5, and 6, respectively. N is equal to the number of columns in the DRAM.

Kernel Size K x N	Our Method		GPU (A100)	
	Throughput (TB/s)	Latency (us)	Throughput (TB/s)	Latency (us)
7x65536	47.3	(4.97)	0.136	(7.06)
15x65536	37.0	(13.6)	0.352	(7.36)
31x65536	32.2	(32.3)	0.609	(5.95)
63x65536	30.0	(70.8)	0.834	(6.09)
127x65536	28.6	(149)	1.04	(6.48)

192 Due to the inherent error in PuCD operations, our proposed bulk bitwise accumulation also has
 193 error. Table 3 presents the normalized mean square error (NMSE) of our DRAM-based calculations
 194 compared to ideal computational results (details in Appendix B). As shown in Table 3, the error
 195 in our proposed method is three to four orders of magnitude smaller than the quantization error
 196 typically used in machine learning [Wei et al., 2024]. This implies that adopting our method for
 197 neural network inference would have a negligible impact on inference accuracy.

Table 3: NMSE error of our bulk bitwise accumulation for different kernel sizes.

Kernel Size	NMSE error
7x65536	9.1e-08
15x65536	1.1e-07
31x65536	2.9e-07
63x65536	6.9e-07
127x65536	1.5e-06

198 5 Conclusion

199 This paper has introduced an efficient method for bulk bitwise accumulation using PuCD. Our
 200 approach addresses the challenges of implementing accumulation operations in PuCD, which are crucial
 201 for machine learning tasks. We presented a novel implementation of the POPCNT3 operation
 202 optimized for commercial DRAM, serving as a fundamental building block for our scalable bit-
 203 wise accumulation technique. Evaluation using DDR4 DRAM modules shows that our approach
 204 consistently outperforms GPU-based methods across all tested kernel sizes, with throughput im-
 205 provements ranging from 27 to 348 times, while maintaining high accuracy. These results highlight
 206 the effectiveness of leveraging bit-level parallelism in DRAM for accelerating bulk bitwise accu-
 207 mulation operations. This substantial increase in throughput is crucial for constructing large matrix
 208 multiplications, which form the cornerstone of many machine learning inference tasks. The ability
 209 to process vast amounts of data rapidly can lead to significant advancements in real-time predictions
 210 and efficient handling of large datasets. Future work could explore integrating this approach into
 211 broader machine learning frameworks and extending it to other types of operations. By enabling
 212 efficient PIM capabilities in standard DRAM, our work contributes to addressing the data transfer
 213 bottleneck in modern computing systems.

214 References

- 215 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
216 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
217 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
218 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,
219 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,
220 Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proceedings of the*
221 *34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook,
222 NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 223 Jaewan Choi, Jaehyun Park, Kwanhee Kyung, Nam Sung Kim, and Jung Ho Ahn. Unleashing
224 the Potential of PIM: Accelerating Large Batched Inference of Transformer-Based Generative
225 Models. *IEEE Computer Architecture Letters*, 22(2):113–116, July 2023. ISSN 1556-6064. doi:
226 10.1109/LCA.2023.3305386. URL [https://ieeexplore.ieee.org/document/10218731/](https://ieeexplore.ieee.org/document/10218731/?arnumber=10218731)
227 ?arnumber=10218731. Conference Name: IEEE Computer Architecture Letters.
- 228 Yuting Wu, Ziyu Wang, and Wei D. Lu. PIM GPT a hybrid process in memory accelerator
229 for autoregressive transformers. *npj Unconventional Computing*, 1(1):4, July 2024. ISSN
230 3004-8672. doi: 10.1038/s44335-024-00004-2. URL [https://www.nature.com/articles/](https://www.nature.com/articles/s44335-024-00004-2)
231 [s44335-024-00004-2](https://www.nature.com/articles/s44335-024-00004-2).
- 232 Fabrice Devaux. The true processing in memory accelerator. In *2019 IEEE Hot Chips 31 Symposium*
233 *(HCS)*, pages 1–24, 2019. doi: 10.1109/HOTCHIPS.2019.8875680.
- 234 Young-Cheon Kwon, Suk Han Lee, Jaehoon Lee, Sang-Hyuk Kwon, Je Min Ryu, Jong-Pil Son,
235 O Seongil, Hak-Soo Yu, Haesuk Lee, Soo Young Kim, Youngmin Cho, Jin Guk Kim, Jongy-
236 oon Choi, Hyun-Sung Shin, Jin Kim, BengSeng Phuah, HyoungMin Kim, Myeong Jun Song,
237 Ahn Choi, Daeho Kim, SooYoung Kim, Eun-Bong Kim, David Wang, Shinhaeng Kang, Yuhwan
238 Ro, Seungwoo Seo, JoonHo Song, Jaeyoun Youn, Kyomin Sohn, and Nam Sung Kim. 25.4
239 a 20nm 6gb function-in-memory dram, based on hbm2 with a 1.2tflops programmable com-
240 puting unit using bank-level parallelism, for machine learning applications. In *2021 IEEE In-*
241 *ternational Solid-State Circuits Conference (ISSCC)*, volume 64, pages 350–352, 2021. doi:
242 10.1109/ISSCC42613.2021.9365862.
- 243 Seongju Lee, Kyuyoung Kim, Sanghoon Oh, Joonhong Park, Gimoon Hong, Dongyoon Ka,
244 Kyudong Hwang, Jeongje Park, Kyeongpil Kang, Jungyeon Kim, Junyeol Jeon, Nahsung Kim,
245 Yongkee Kwon, Kornijcuk Vladimir, Woojae Shin, Jongsoon Won, Minkyu Lee, Hyunha Joo,
246 Haerang Choi, Jaewook Lee, Donguc Ko, Younggun Jun, Keewon Cho, Ilwoong Kim, Choungki
247 Song, Chunseok Jeong, Daehan Kwon, Jieun Jang, Il Park, Junhyun Chun, and Joochwan
248 Cho. A 1ynm 1.25v 8gb, 16gb/s/pin gddr6-based accelerator-in-memory supporting 1tflops
249 mac operation and various activation functions for deep-learning applications. In *2022 IEEE*
250 *International Solid-State Circuits Conference (ISSCC)*, volume 65, pages 1–3, 2022. doi:
251 10.1109/ISSCC42614.2022.9731711.
- 252 Ping Chi, Shuangchen Li, Cong Xu, Tao Zhang, Jishen Zhao, Yongpan Liu, Yu Wang, and Yuan
253 Xie. PRIME: A Novel Processing-in-Memory Architecture for Neural Network Computation
254 in ReRAM-Based Main Memory. In *2016 ACM/IEEE 43rd Annual International Symposium*
255 *on Computer Architecture (ISCA)*, pages 27–39, June 2016. doi: 10.1109/ISCA.2016.13. URL
256 <https://ieeexplore.ieee.org/document/7551380/?arnumber=7551380>. ISSN: 1063-
257 6897.
- 258 Vivek Seshadri, Donghyuk Lee, Thomas Mullins, Hasan Hassan, Amirali Boroumand, Jeremie
259 Kim, Michael A. Kozuch, Onur Mutlu, Phillip B. Gibbons, and Todd C. Mowry. Ambit: in-
260 memory accelerator for bulk bitwise operations using commodity DRAM technology. In *Pro-*
261 *ceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, pages
262 273–287, Cambridge Massachusetts, October 2017. ACM. ISBN 978-1-4503-4952-9. doi:
263 10.1145/3123939.3124544. URL <https://dl.acm.org/doi/10.1145/3123939.3124544>.
- 264 Charles Eckert, Xiaowei Wang, Jingcheng Wang, Arun Subramaniyan, Ravi Iyer, Dennis Sylvester,
265 David Blaauw, and Reetuparna Das. Neural Cache: Bit-Serial In-Cache Acceleration of Deep

- 266 Neural Networks. In *2018 ACM/IEEE 45th Annual International Symposium on Computer Archi-*
267 *tecture (ISCA)*, pages 383–396, Los Angeles, CA, June 2018. IEEE. ISBN 978-1-5386-
268 5984-7. doi: 10.1109/ISCA.2018.00040. URL [https://ieeexplore.ieee.org/document/](https://ieeexplore.ieee.org/document/8416842/)
269 [8416842/](https://ieeexplore.ieee.org/document/8416842/).
- 270 Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong,
271 Ruiping Wang, Jilong Xue, and Furu Wei. The Era of 1-bit LLMs: All Large Language Models are
272 in 1.58 Bits, February 2024. URL <http://arxiv.org/abs/2402.17764>. arXiv:2402.17764
273 [cs].
- 274 Fei Gao, Georgios Tziantzioulis, and David Wentzlaff. ComputeDRAM: In-Memory Compute Us-
275 ing Off-the-Shelf DRAMs. In *Proceedings of the 52nd Annual IEEE/ACM International Sym-*
276 *posium on Microarchitecture*, pages 100–113, Columbus OH USA, October 2019. ACM. ISBN
277 978-1-4503-6938-1. doi: 10.1145/3352460.3358260. URL [https://dl.acm.org/doi/10.](https://dl.acm.org/doi/10.1145/3352460.3358260)
278 [1145/3352460.3358260](https://dl.acm.org/doi/10.1145/3352460.3358260).
- 279 Ataberk Olgun, Minesh Patel, A. Giray Yağlıkcı, Haocong Luo, Jeremie S. Kim, Nisa Bostancı,
280 Nandita Vijaykumar, Oğuz Ergin, and Onur Mutlu. QUAC-TRNG: High-Throughput True Ran-
281 dom Number Generation Using Quadruple Row Activation in Commodity DRAM Chips, May
282 2021. URL <http://arxiv.org/abs/2105.08955>. arXiv:2105.08955 [cs].
- 283 Fei Gao, Georgios Tziantzioulis, and David Wentzlaff. FracDRAM: Fractional Values in Off-the-
284 Shelf DRAM. In *2022 55th IEEE/ACM International Symposium on Microarchitecture (MICRO)*,
285 pages 885–899, Chicago, IL, USA, October 2022. IEEE. ISBN 978-1-66546-272-3. doi: 10.1109/
286 MICRO56248.2022.00066. URL <https://ieeexplore.ieee.org/document/9923819/>.
- 287 Ismail Emir Yuksel, Yahya Can Tuğrul, F. Nisa Bostancı, Geraldo F. Oliveira, A. Giray Yaglıkcı,
288 Ataberk Olgun, Melina Soysal, Haocong Luo, Juan Gómez-Luna, Mohammad Sadrosadati, and
289 Onur Mutlu. Simultaneous Many-Row Activation in Off-the-Shelf DRAM Chips: Experimen-
290 tal Characterization and Analysis, May 2024. URL <http://arxiv.org/abs/2405.06081>.
291 arXiv:2405.06081 [cs].
- 292 İsmail Emir Yüksel, Yahya Can Tuğrul, Ataberk Olgun, F. Nisa Bostancı, A. Giray Yağlıkcı, Ger-
293 aldo F. Oliveira, Haocong Luo, Juan Gómez-Luna, Mohammad Sadrosadati, and Onur Mutlu.
294 Functionally-Complete Boolean Logic in Real DRAM Chips: Experimental Characterization
295 and Analysis. In *2024 IEEE International Symposium on High-Performance Computer Archi-*
296 *tecture (HPCA)*, pages 280–296, Edinburgh, United Kingdom, March 2024. IEEE. ISBN
297 9798350393132. doi: 10.1109/HPCA57654.2024.00030. URL [https://ieeexplore.ieee.](https://ieeexplore.ieee.org/document/10476435/)
298 [org/document/10476435/](https://ieeexplore.ieee.org/document/10476435/).
- 299 Jianyu Wei, Shijie Cao, Ting Cao, Lingxiao Ma, Lei Wang, Yanyong Zhang, and Mao Yang. T-
300 MAC: CPU Renaissance via Table Lookup for Low-Bit LLM Deployment on Edge, June 2024.
301 URL <http://arxiv.org/abs/2407.00088>. arXiv:2407.00088 [cs].
- 302 Ataberk Olgun, Hasan Hassan, A. Giray Yağlıkcı, Yahya Can Tuğrul, Lois Orosa, Haocong Luo,
303 Minesh Patel, Oğuz Ergin, and Onur Mutlu. DRAM Bender: An Extensible and Versatile FPGA-
304 based Infrastructure to Easily Test State-of-the-art DRAM Chips, September 2023. URL <http://arxiv.org/abs/2211.05838>.
305 [arXiv:2211.05838](http://arxiv.org/abs/2211.05838) [cs].

306 A Evaluation Setup

307 Table 4 provides a detailed overview of the hardware specifications used in our experiments. The
308 core of our setup consists of a Xilinx Alveo U200 FPGA board, which serves as the platform for
309 implementing our PuCD-based method. We paired this FPGA with four SK Hynix DDR4 DIMM
310 modules, each with a density of 4GB and operating at 2400MT/s. For precise control over DRAM
311 operations, we employed DRAM Bender [Olgun et al., 2023], an open-source memory controller
312 implemented on the FPGA. DRAM Bender is crucial to our setup as it enables us to issue DRAM
313 commands with arbitrary timings, a capability essential for implementing PuCD operations. We
314 controlled DRAM Bender through a host CPU program, allowing us to generate and execute the
315 specific command sequences required for our bulk bitwise accumulation method.

316 Figure 4 shows the setup consisting of the following key components: ① A Xilinx Alveo U200
 317 FPGA board programmed with DRAM Bender. This board serves as the central processing unit for
 318 our PuCD operations. ② Four DDR4 DIMM modules, each containing 8 DRAM chips, are installed
 319 in the DIMM slots of our setup. This configuration results in a total of 32 DRAM chips available for
 320 our experiments, allowing us to explore the full potential of bank-level parallelism in our method.
 321 ③ A separate computer, responsible for generating DRAM commands and controlling the overall
 322 experiment flow. This host machine runs the software that interfaces with DRAM Bender to execute
 323 our bulk bitwise accumulation operations.

324 Our measurement methodology involved conducting 1000 independent trials for each kernel size
 325 and configuration to ensure statistical significance and account for DRAM behavior variability. We
 326 focused on two primary metrics: throughput, calculated as the amount of data processed per second
 327 when executing our method in parallel across all 64 banks of the 4 DRAM modules, and latency,
 328 defined as the time elapsed from the first DRAM command to the last for a single kernel execution.

329 To provide a meaningful performance comparison, we implemented a baseline version of our bulk
 330 bitwise accumulation method on an NVIDIA A100-PCIe-40GB GPU using CUDA. GPU throughput
 331 was calculated based on the execution time of 8192 parallel kernel operations.

Table 4: Specifications of the FPGA and DRAM module used in our evaluation.

FPGA	DRAM Module	Memory Bandwidth (GB/s)
Xilinx Alveo U200	SK Hynix’s DDR4 DIMM ¹	76.8

¹ We used four DRAM modules HMA851U6CJR6N-UHN0 that have a density of 4GB and operates at 2400MT/s.

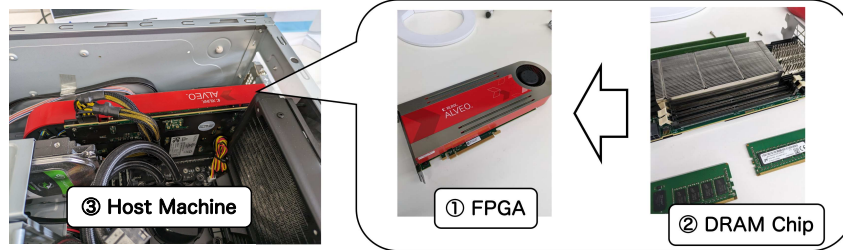


Figure 4: DRAM chips testing setup.

332 B NMSE Calculation

333 To quantify the accuracy of our bulk bitwise accumulation method, we employed the normalized
 334 mean square error (NMSE) metric. NMSE is a statistical measure that quantifies the relative differ-
 335 ence between predicted and actual values, normalized by the variance of the actual values. It is
 336 defined as:

$$\text{NMSE} = \frac{E[(Y - \hat{Y})^2]}{\text{Var}(Y)}$$

337 NMSE values provide a measure of our method’s accuracy relative to the inherent variability in the
 338 accumulation task. Lower NMSE values indicate better performance, with an NMSE of 0 represent-
 339 ing perfect prediction.

340 For this evaluation, we generated input bits randomly with a 50% probability for each bit state.
 341 Our NMSE calculation process began with performing an ideal accumulation using a high-precision
 342 software implementation for each set of randomly generated inputs. This provided our ”actual”
 343 values (Y). We then performed the accumulation using our proposed PuCD-based method, obtaining
 344 our ”predicted” values (\hat{Y}). For each column, we calculated the squared difference between the ideal
 345 and PuCD-based results: $(Y - \hat{Y})^2$. We then computed the mean of these squared differences across
 346 all columns, giving us $E[(Y - \hat{Y})^2]$. Next, we calculated the variance of the ideal accumulation

347 results across all columns: $\text{Var}(Y)$. Finally, we divided the mean squared error by the variance to
348 obtain the NMSE.