

### 37.6 A 22nm 60.81TFLOPS/W Diffusion Accelerator with Bandwidth-Aware Memory Partition and BL-Segmented Compute-in-Memory for Efficient Multi-Task Content Generation

Yiqi Jing<sup>1</sup>, Jiaqi Zhou<sup>1</sup>, Yiyang Sun<sup>1</sup>, Siyuan He<sup>1</sup>, Peiyu Chen<sup>2,3</sup>, Ru Huang<sup>1</sup>, Le Ye<sup>1,2</sup>, Tianyu Jia<sup>1</sup>

<sup>1</sup>Peking University, Beijing, China

<sup>2</sup>Advanced Institute of Information Technology of Peking University, Hangzhou, China

<sup>3</sup>Nano Core Chip Electronic Technology, Hangzhou, China

Initially applied for image synthesis [1], Diffusion Models (DMs) have been rapidly expanded into many content-generation tasks, e.g. 3D scenes [2-3] or video [4], and deliver exceptional performance. Figure 37.6.1 provides an overview of DM architecture, which typically processes random noisy input through multiple, i.e. 20-50, denoising steps to generate the desired output content. Each denoising step incorporates a U-Net structure with a down-sampling encoder and an up-sampling decoder, which contains repetitive transformer blocks. To support diverse content generation, multi-view [5] or temporal attention blocks [6] are integrated to enhance 3D scene or video-frame consistency. Due to the large number of denoising steps, generating a single piece of content consumes significant latency, e.g. ~70s for a 4 seconds of 6fps video on an A100 GPU. To improve the hardware performance and efficiency, compute-in-memory (CIM) accelerators have been developed for transformers [7-9] or DMs [10]. However, for several reasons, it is still challenging to use existing CIM-based accelerators for practical image or multi-task DMs. First, a significant operational intensity (Opl) variation exists along DM layers leading to dynamic memory bandwidth (BW) requirements. The DMs also require excessive data storage, i.e. ~10 $\times$  larger than VAEs and GANs. Second, the impressive CIM macro efficiency is often significantly degraded at the system-level due to the diminishing reuse rate for large AI models. Moreover, ~59% of CIM macro power is consumed on data access rather than computation. Third, emerging DM tasks for 3D or video require additional consistency operations, i.e. multi-view/frame attention, for smooth transitions across viewpoints or frames. Such consistency attention is quite computationally-intensive and comprises ~31% of the total operations.

To address these challenges, this paper presents a digital CIM-based accelerator for multi-task DMs with following key features: 1) a dynamic BW-aware memory partitioning scheme is developed with dense on-chip eDRAM storage to optimize CIM utilization and reduce EMA, 2) a bitline (BL)-segmented CIM cluster is designed with reuse-aware weight reordering to enhance system efficiency, and 3) a hierarchical consistency optimization flow is presented to minimize frame/pixel-level operations to improve performance. Overall, these innovations enable our chip to achieve a 60.81TFLOPS/W system efficiency, which has 1.4 $\times$  better performance than a prior image DM chip [10] and our chip also shows promising performance for more diverse content-generation tasks.

Figure 37.6.2 shows the overall architecture of our chip. It comprises a CIM acceleration subsystem with 9 CIM clusters, each contains four 24Kb digital-CIM macros, 3Mb on-chip eDRAM, a multi-frame consistency management (MFCM) unit, a host RISC-V CPU and peripheral circuits. Direct connections are implemented between CIM clusters to support pipeline parallelism and crossbar (Xbar) interfaces are used to connect eDRAM. A reuse-aware weight update scheduler is incorporated inside the CIM subsystem to improve the data reuse rate for all clusters. The eDRAM is designed into four splittable banks and each has 8 $\times$ 3 32Kb 3T gain-cell arrays. A leak-tracking reference column is added in each array to enhance Vref accuracy and extend eDRAM retention time (Fig. 37.6.3). A dynamic BW-aware partitioning module is designed with an interconnect coupler to support flexible memory bank partitioning to adapt to varying Opl of different DM layers. The MFCM performs the hierarchical consistency optimizations, which contains a progressive view extension unit for frame-level optimizations, along with a spatial-temporal compress unit and epipolar-attention sparse unit to reduce pixel-level computations.

Figure 37.6.3 illustrates the architecture and operation of our BW-aware memory partitioning. A two-step partitioning scheme is developed for each runtime subtask, i.e. tiled matrix workload in DM layers. First, the computation and memory resource demand for each subtask are evaluated based on operational intensity and recorded in a resource-aware subtask table. Second, the data arrangement within eDRAM, i.e. eDRAM bank partitioning plan, is determined by the required BW and stored in a memory-partition table. To maximize utilization, a time-multiplexed BW distribution method is adopted to decouple BW from memory capacity, e.g. data for subtask 5 is stored equally in ten eDRAM columns to increase available peak BW. Compared to a conventional fixed memory allocation (Case 1), which evenly distributes resources to subtasks, our BW-aware partitioning jointly optimizes CIM and BW utilization for subtasks with diverse Opl, e.g., subtasks 3-5 have improved capacity and BW utilization by 1.98 $\times$  and 2.76 $\times$  in Case 2. Based on the partitioning plan, area-efficient crossbar and interconnect couplers are reconfigured for proper topology and flow control. The coupler determines source and destination addresses through a loop decoder,

while a credit noter monitors eDRAM refreshes and CIM backpressure to ensure reliable data transfer. Overall, the BW-aware memory partitioning improves CIM and BW utilization by 1.27-9.54 $\times$  and 1.51-10.19 $\times$  across DM layers for a Wonder3D model, leading to a total 2.68 $\times$  performance gain with only 3.7% area overhead.

Figure 37.6.4 illustrates the implementation details of the CIM cluster and the weight-update scheduler. Each cluster integrates four BL-segmented CIM macros, an aggregator unit, and a local NoC. Each macro contains six 128 $\times$ 32 6T SRAM MAC-arrays, pre/post-processing units, a weight-alignment unit, and I/O buffers. The CIM array is designed with an architecture of 4 $\times$ 64 weight subarrays, which consist of 16 SRAM cells with a local stationary unit (LSU), and a LUT-bypass adder tree. To alleviate costly data access, a 2-stage BL-segmentation technique is incorporated, which segments BL based on operational addresses at both the MAC-array level and a finer-grained subarray level. This approach reduces the effective BL loading and unnecessary precharging during SRAM access to enhance CIM efficiency by 27%. A LUT-bypass adder tree is designed by leveraging a 4b sparse LUT as a multiplier and a first-stage adder to reduce high dynamic transitions. A bypass adder is used as the second-stage adder to skip zeros, resulting in 11% power reduction. During computation, the weight-update scheduler supports simultaneous computation with our reuse-aware weight reordering, in which a top-k module is used to generate a sparse attention pattern to indicate the weights to be reused in the CIM subarray. The pattern is first row-wise reordered via a reuse-driven activation sorter by similarity comparisons, forming the CIM's activation sequence. Then, a column-wise reordering based on weight lifespan decides the weight-update sequence with the least-recently-used queue-replacement policy, thereby reducing memory access for attention layers by 29%. Overall, the BL-S CIM macro achieves 1.31 $\times$  performance gain and 54% energy savings by the above techniques.

Figure 37.6.5 depicts our MFCM scheme for efficient content generation. Conventional multi-view DMs denoise frames from all viewpoints and apply multi-view attention on pixels to ensure consistency, resulting in significant overhead. We develop a hierarchical consistency computation flow with both frame-level and pixel-level optimizations with following three stages. Stage 1 adopts a progressive view extension technique, which utilizes fewer viewpoint frames at initial timesteps to reduce frame-level computations, e.g. only 2 frames at timestep 5. New frames from different viewpoints will be added by duplicating the previous frame once its similarity with reference frame (calculated by a frame compare unit) reaches predefined threshold. Stage 2 reduces pixel-level computation by skipping background and trivial pixels using a spatial-temporal compression unit. The target object is segmented from the background using RGBA values in salient object detection unit (SODU) based on a spiral search pattern. Unmodified pixels are further skipped in trivial-pixel sparsity unit (TPSU) by similarity assessment across denoising iterations. Stage 3 further leverages a pixel-level epipolar-attention mechanism to minimize irrelevant inter-frame interactions. The epipolar solver controls a 3 $\times$ 3 MAC array to generate the pixels and epipolar line on each view plane associated with the target light ray. A sparse attention is performed between the epipolar line and the pixel, e.g. pixel on the P3 frame and the epipolar line on the P0 frame, to mask out irrelevant regions and reduce computations. The above hierarchical consistency optimizations bring us a 3.71 $\times$  speedup and 68% energy saving.

Figure 37.6.6 shows the measurement results of the 22nm CIM processor. Multiple content-generation tasks are evaluated using SOTA DMs, i.e. SD-v1.5, Wonder3D, SVD, based on customized hybrid BF16-W4A8 quantization. Compared to a SOTA image DM accelerator [10], our chip achieves a 1.4 $\times$  performance improvement. Since there is no prior accelerator for 3D or video, we also provide the execution time for 3D and video models with high system efficiency. The BW-aware partitioning, BL-S macros and MFCM together contribute to 13.03 $\times$  performance and 3.69 $\times$  system-efficiency improvement, leading to a 49.74-60.81TFLOPS/W system efficiency (1.52 $\times$  better than [10]). Our CIM macro FoM, which considers both energy and area efficiency, is 1.28 $\times$  higher than the Booth8 CIM in [10]. The system FoM is 1.44 $\times$  and 3.82 $\times$  higher than SOTA CIM [8] and DM [11] accelerators, illustrating better optimized CIM computation and BW utilization. Figure 37.6.7 shows the die photo and more specifications.

#### Acknowledgement:

This work was supported in part by NSFC Grant 92164301, Grant 62225401, and Grant U23A6007; Zhejiang Provincial Key R&D program under Grant 2021C01035; Grant QYJS-2023-2401-B, and Grant QYJS-2023-2402-B. Corresponding authors: Tianyu Jia and Le Ye.



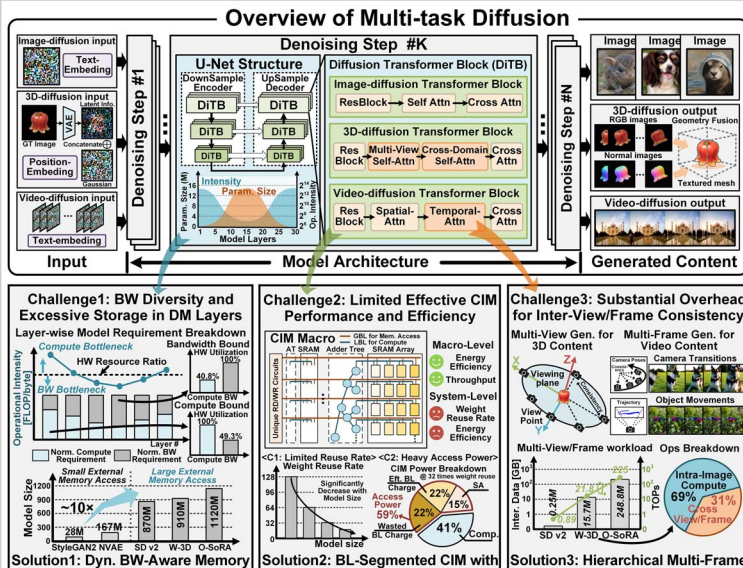


Figure 37.6.1: Overview of diffusion models for multi-task content generation and the deployment challenges.

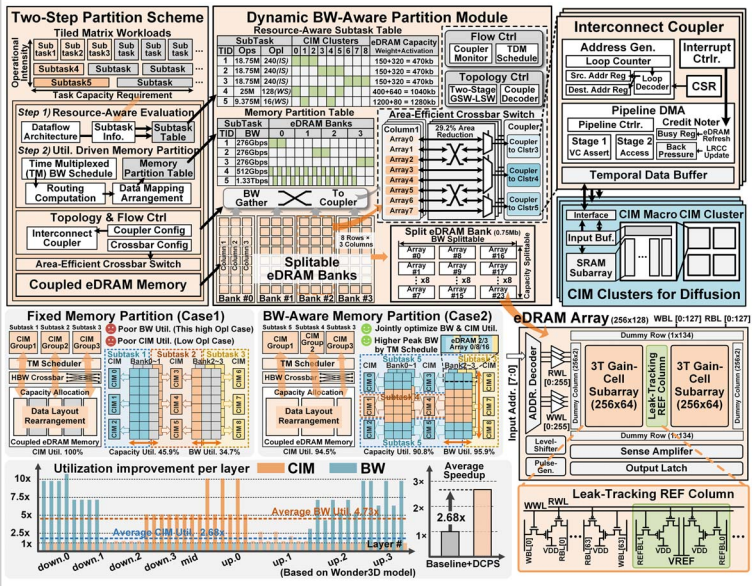


Figure 37.6.3: Dynamic BW-aware memory partitioning scheme with eDRAM to jointly improve CIM and BW utilization.

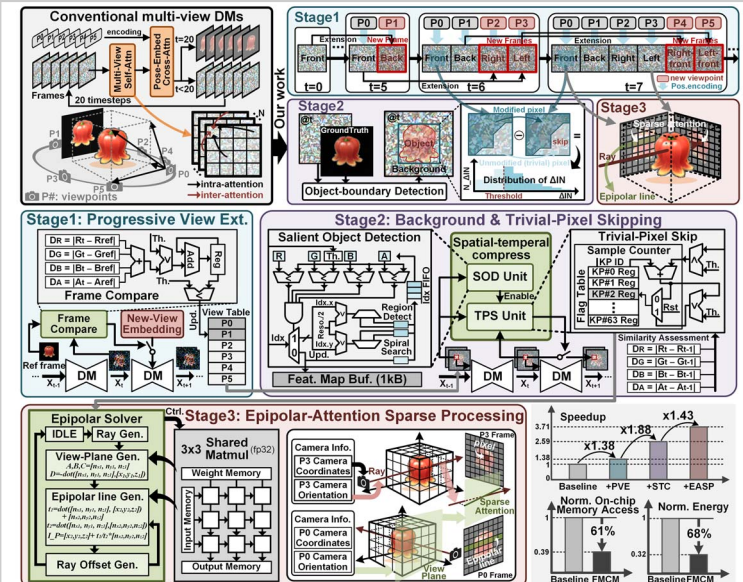


Figure 37.6.5: Multi-frame consistency management with hierarchical optimizations for diverse content generation.

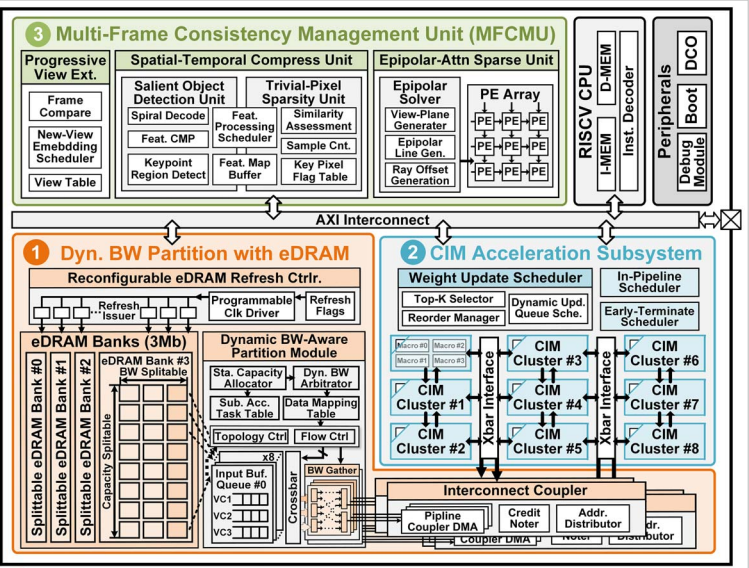


Figure 37.6.2: Overall chip architecture.

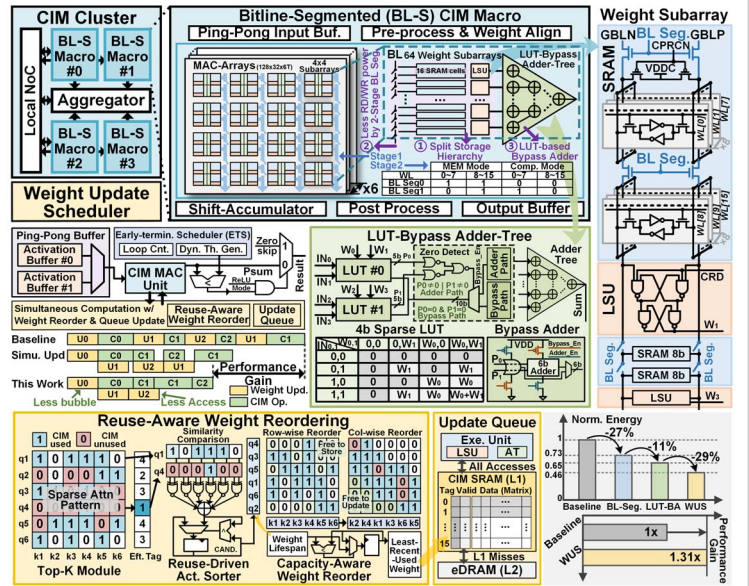


Figure 37.6.4: Bitline-segmented (BL-S) CIM macro and weight update scheduler with reuse-aware reordering update.

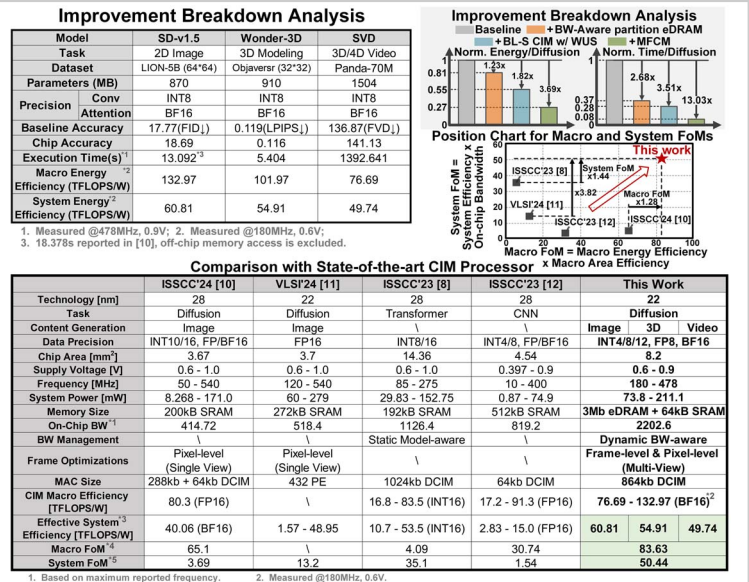


Figure 37.6.6: Measurement results and performance comparison table.



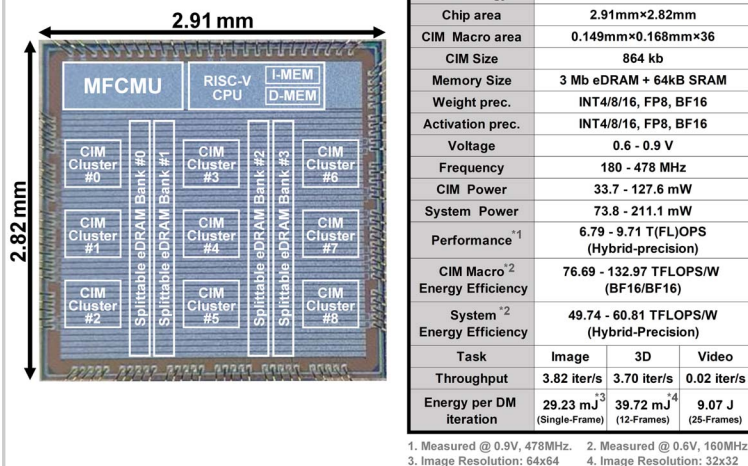


Figure 37.6.7: Chip micrograph and specifications.

## References:

- [1] R. Rombach et al., "High-Resolution Image Synthesis with Latent Diffusion Models," *IEEE CVPR*, pp. 10674-10685, 2022.
- [2] X.-Y. Zheng et al., "MVD<sup>2</sup>: Efficient Multiview 3D Reconstruction for Multiview Diffusion," *ACM SIGGRAPH*, 2024.
- [3] X. Long et al., "Wonder3D: Single Image to 3D using Cross-Domain Diffusion," *IEEE CVPR*, pp. 9970-9980, 2024.
- [4] H. Ni et al., "Conditional Image-to-Video Generation with Latent Flow Diffusion Models," *IEEE CVPR*, pp. 18444-18455, 2023.
- [5] Z. Deng et al., "MV-Diffusion: Motion-aware Video Diffusion Model," *ACM Multimedia*, pp. 7255-7263, 2023.
- [6] H. Jeong et al., "VMC: Video Motion Customization using Temporal Attention Adaption for Text-to-Video Diffusion Models," *IEEE CVPR*, pp. 9212-9221, 2024.
- [7] F. Tu et al., "A 28nm 15.59μJ/Token Full-Digital Bitline-Transpose CIM-Based Sparse Transformer Accelerator with Pipeline/Parallel Reconfigurable Modes," *ISSCC*, pp. 466-468, 2022.
- [8] F. Tu et al., "MultiTCIM: A 28nm 2.24μJ/Token Attention-Token-Bit Hybrid Sparse Digital CIM-Based Accelerator for Multimodal Transformers," *ISSCC*, pp. 248-250, 2023.
- [9] S. Liu et al., "A 28nm 53.8TOPS/W 8b Sparse Transformer Accelerator with In-Memory Butterfly Zero Skipper for Unstructured-Pruned NN and CIM-Based Local-Attention-Reusable Engine," *ISSCC*, pp. 250-252, 2023.
- [10] R. Guo et al., "A 28nm 74.34TFLOPS/W BF16 Heterogenous CIM-Based Accelerator Exploiting Denoising-Similarity for Diffusion Models," *ISSCC*, pp. 362-364, 2024.
- [11] Y. Qin et al., "A 52.01 TFLOPS/W Diffusion Model Processor with Inter-Time-Step Convolution-Attention-Redundancy Elimination and Bipolar Floating-Point Multiplication," *IEEE VLSI Technology and Circuits*, 2024.
- [12] J. Yue et al., "A 28nm 16.9-300TOPS/W Computing-in-Memory Processor Supporting Floating-Point NN Inference/Training with Intensive-CIM Sparse-Digital Architecture," *ISSCC*, pp. 252-254, 2023.