



ViDA: Video Diffusion Transformer Acceleration with Differential Approximation and Adaptive Dataflow

Li Ding¹, Jun Liu¹², Shan Huang¹, Guohao Dai^{12*}

¹Shanghai Jiao Tong University ²Infinigence AI

Equal Contribution *Corresponding to:<<u>daiguohao@sjtu.edu.cn</u>>

ASP-DAC 2025

2025.1.21



• Current video generation models excel due to the diffusion paradigm and the Video Diffusion Transformer backbone.





Backboldedulideood)ifooaioonTsaksformer

Sora^[2]: Clifftop Waves and Waterfalls

[1] https://www.vidu.studio/ [2] https://sora.com

2025/3/7

- The Paradigm: What is diffusion model ?
 - Frames are iteratively denoised, with noise predicted by VDiT.

Training: forward diffusion



• The Backbone: What is Video Diffusion Transformer (VDiT) ?

• A transformer-based noise predictor with a spatial and temporal patching structure.



• The Backbone: What is Video Diffusion Transformer (VDiT) ?

• A transformer-based noise predictor with a spatial and temporal patching structure.



- Video generation is slow due to large inference data.
 - It costs over 1 hour to generate a one-minute video.^[1]

Inference: reverse diffusion



VDiT is the main bottleneck! Take up over 90% of the inference time.

*Other processes include initialization and noise handling.



Noise predictor: VDiT

[1] https://www.wired.com/story/openai-sora-generative-ai-video/

Key Idea

How to accelerate VDiT ?

• We can learn from traditional video processing.

Video Processing^[1]



[1] Research on Key Technologies of H.264/HEVC Video Coding & Decoding Standard.

Key Idea

- How to accelerate VDiT ?
 - We can learn from traditional video processing.



Challenge 1

- VDiT exists large unaccelerated computations.
 - Existing token reduction methods^{[1][2]} fail to accelerate Act-Act operator due to accuracy loss caused by dimension mismatch.



Song Z, et al. CMC: Video Transformer Acceleration via CODEC Assisted Matrix Condensing [ASPLOS 2024]
 Wang X, et al. InterArch: Video Transformer Acceleration via Inter-Feature Deduplication with Cube-based Dataflow [DAC 2024]

Contribution 1

 Algorithm-hardware co-design for Act-Act operator acceleration.



Tech 1 (algorithm): Differential Approximation Method



Tech1: Differential Approximation Method

• What is differential computing?

• Reducing redundant computations by leveraging the **similarity** between activations.



Differential computing can be applied to the **Act-W operator** efficiently.

Tech1: Differential Approximation Method

How can differential computing be applied to Act-Act operator ?



Tech1: Differential Approximation Method

How can differential computing be applied to Act-Act operator ?



- How to perform SpMM introduced by differential computing efficiently ?
 - Observation: Column-split computing existing column-sparse pattern.



Red: Non-zero elements Blue: zero elements

Non-zero elements are mainly clustered in some columns.

- How to perform SpMM introduced by differential computing efficiently ?
 - Column-split computing based on column-sparse pattern.
 1) Allocation 2) Computing



 How to perform SpMM introduced by differential computing efficiently ?

1) Allocation: determine dense and sparse columns by thresholds.



 How to perform SpMM introduced by differential computing efficiently ?

2) Computing: Compute **outer products** of allocated columns using sparse and dense arrays.



1.10x / 1.56x higher area efficiency compared with dense-only and sparse-only architectures.

Area efficiency



Challenge 2

Large operational intensity (OI) difference among operators leads to inefficient hardware utilization.



Operational intensity^[1] = computations / memory accesses

Contribution 2

 Use intensity adaptive dataflow architecture design for dynamically allocate resources for different operators.



Tech 3 (dataflow): Intensity Adaptive Dataflow Architecture

Tech3: Adaptive Dataflow Architecture

Dataflow: Intensity adaptive dataflow for reasonable resource allocation.



Tech3: Adaptive Dataflow Architecture

 Architecture: Reconfigurable architecture including routing controller to support flexible resource allocation.



Evaluation

А •	<i>ccuracy evaluation</i> <i>Models:</i> STDiT.1, Latte, VDT		FVD↓	CLIPSIM↑	Average loss
• D	Datasets: UCF101, MSR-VTT	STDiT(dense)	477.97	0.264	
		VIDA-STDIT	479.56	0.262	0.55%
•	 Baselines NVIDIA A100 GPU 	CMC-STDIT	483.53	0.263	0.77%
	SOTA Vision Accelerators	Latte(dense)	505.27	0.294	
	 CMC^[1][ASPLOS 2024] InterArch^[2][DAC 2024] 	ViDA-Latte	514.33	0.293	1.07%
•	MetricsFVD, FID, CLIPSIM	CMC-Latte	518.59	0.293	1.45%
		*We furtheptest the ful(SJCFT)01/edataset forts av/iDA0ablcffne)SOTACMCaccuracy test.			

[1] Song Z, et al. CMC: Video Transformer Acceleration via CODEC Assisted Matrix Condensing [ASPLOS 2024]
[2] Wang X, et al. InterArch: Video Transformer Acceleration via Inter-Feature Deduplication with Cube-based Dataflow [DAC 2024]

Evaluation

Area evaluation

Simulation process library

- 32nm standard librabry
- Tools
 - Synopsys Design Compiler
 - CACTI 7
- Architecture
 - 4 x PE group
 - 8 x column-concentrate PEs
 - 4 x 4 DPUs
 - 1 x 4 SPUs
 - Allocation unit
 - Merge unit

Component			Area (mm^2)	Breakdown
		4×4 DPUs	1.04	51.39%
	$CCPE \times 8$	1×4 SPUs	0.52	25.60%
$PEG \times 4$		Allocation unit	0.01	0.70%
		Merge unit	0.07	3.27%
	Routing controller		0.09	4.38%
	Memory elements \times 8			3.38%
Global buffer			0.23	11.27%
Total (7nm)			2.03	100.00%

*Scale to 7nm with a frequency of 1GHz for comparing with other works

Evaluation

 ViDA achieves average 16.44x/2.18x speedup and 18.39x/2.35 area efficiency compared with A100 GPU/SOTA Vision Accelerator.



You H, et al. Vitcod: Vision transformer acceleration via dedicated algorithm and accelerator co-design Song Z, et al. CMC: Video Transformer Acceleration via CODEC Assisted Matrix Condensing [ASPLOS 2024] Wang X, et al. InterArch: Video Transformer Acceleration via Inter-Feature Deduplication with Cube-based Dataflow [DAC 2024]



Conclusion

Design Automation Innovation & Domain-specific Artificial Intelligence

Thank you for your attention!

Communication e-mail : <u>daiguohao@sjtu.edu.cn</u>

ViDA: Video Diffusion Transformer Acceleration with Differential Approximation and Adaptive Dataflow

 $A_2B_2 \approx A_2\Delta B_2 + C$ Video Algorithm Processing Similarity Hardware & PF ME 2 Video **Dataflow** Operator Operator Generation with with low Ol high OI

Lab homepage