



Systematic Approaches for Efficient and Scalable Deep Learning

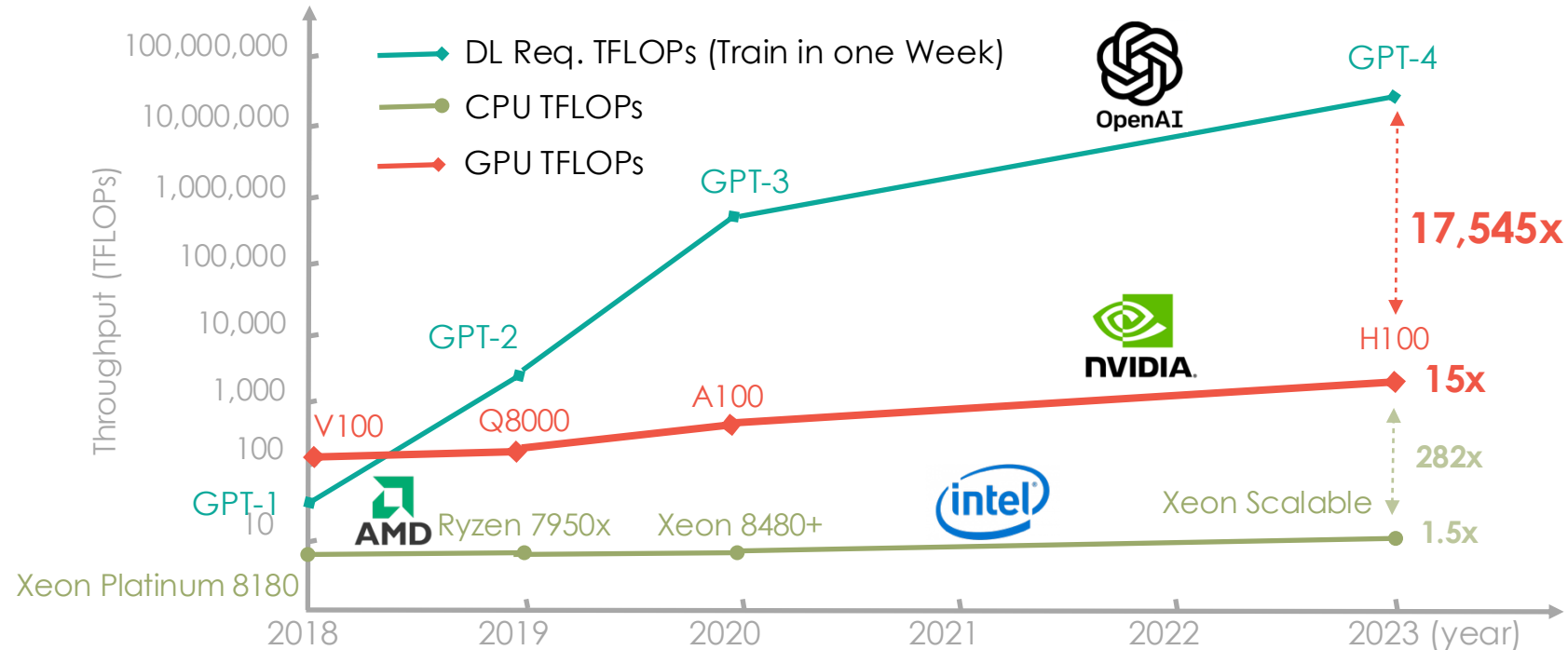
Yuke Wang

Computer Science

Rice University

The Trend of DL Algorithm and Hardware

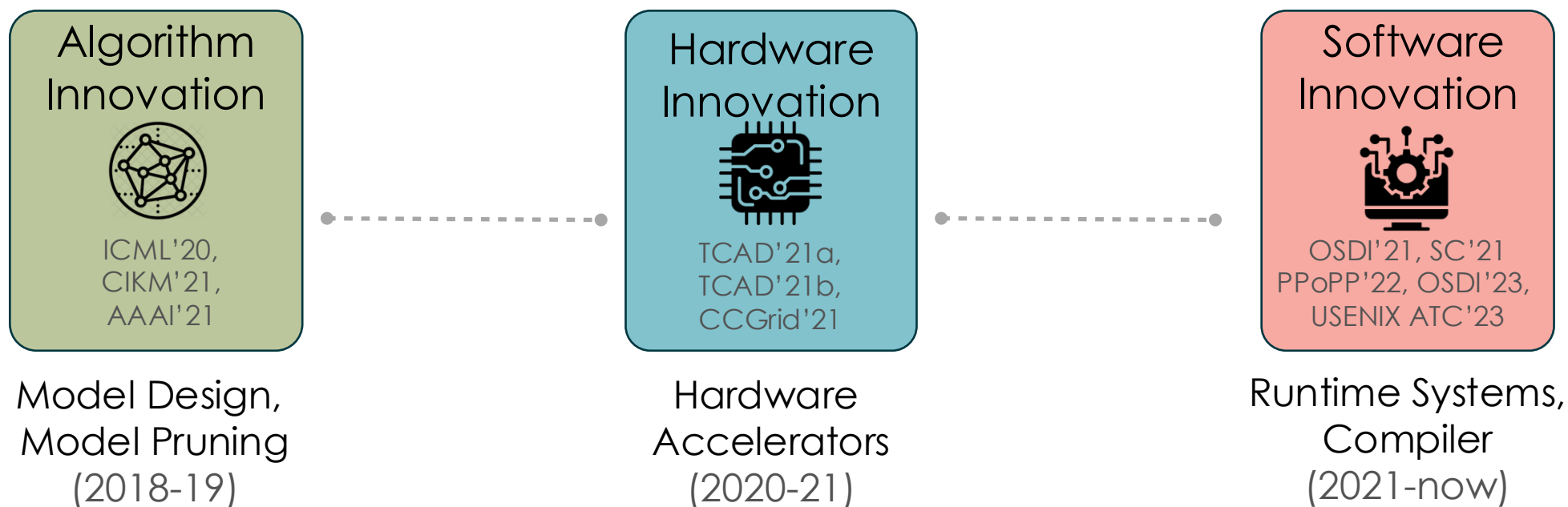
❖ Recap of DL algorithms and hardware performance scaling.



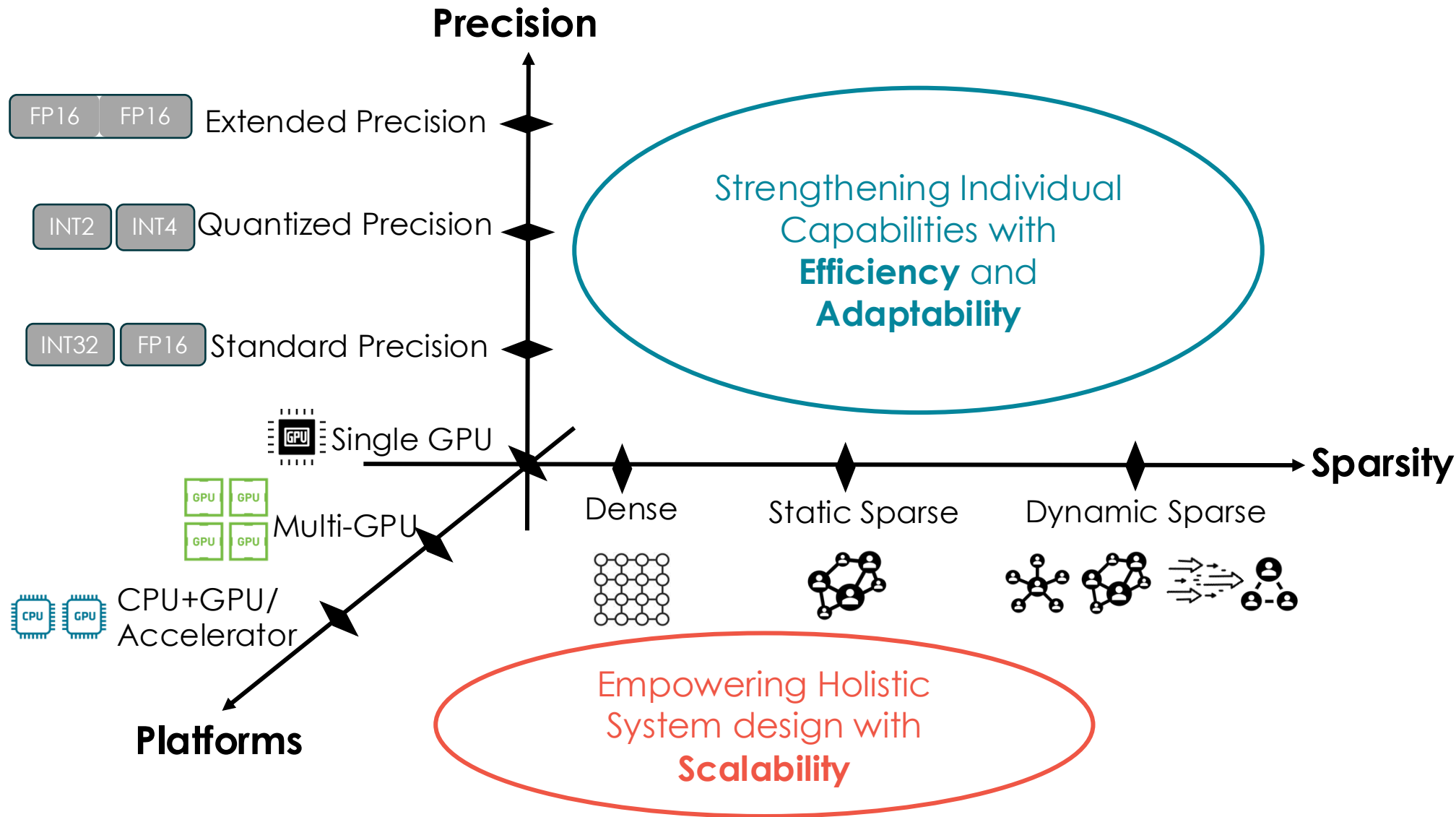
Huge Potential with GPUs! But it still has a **Large Gap!**

Deep Learning Drives Computing Innovations

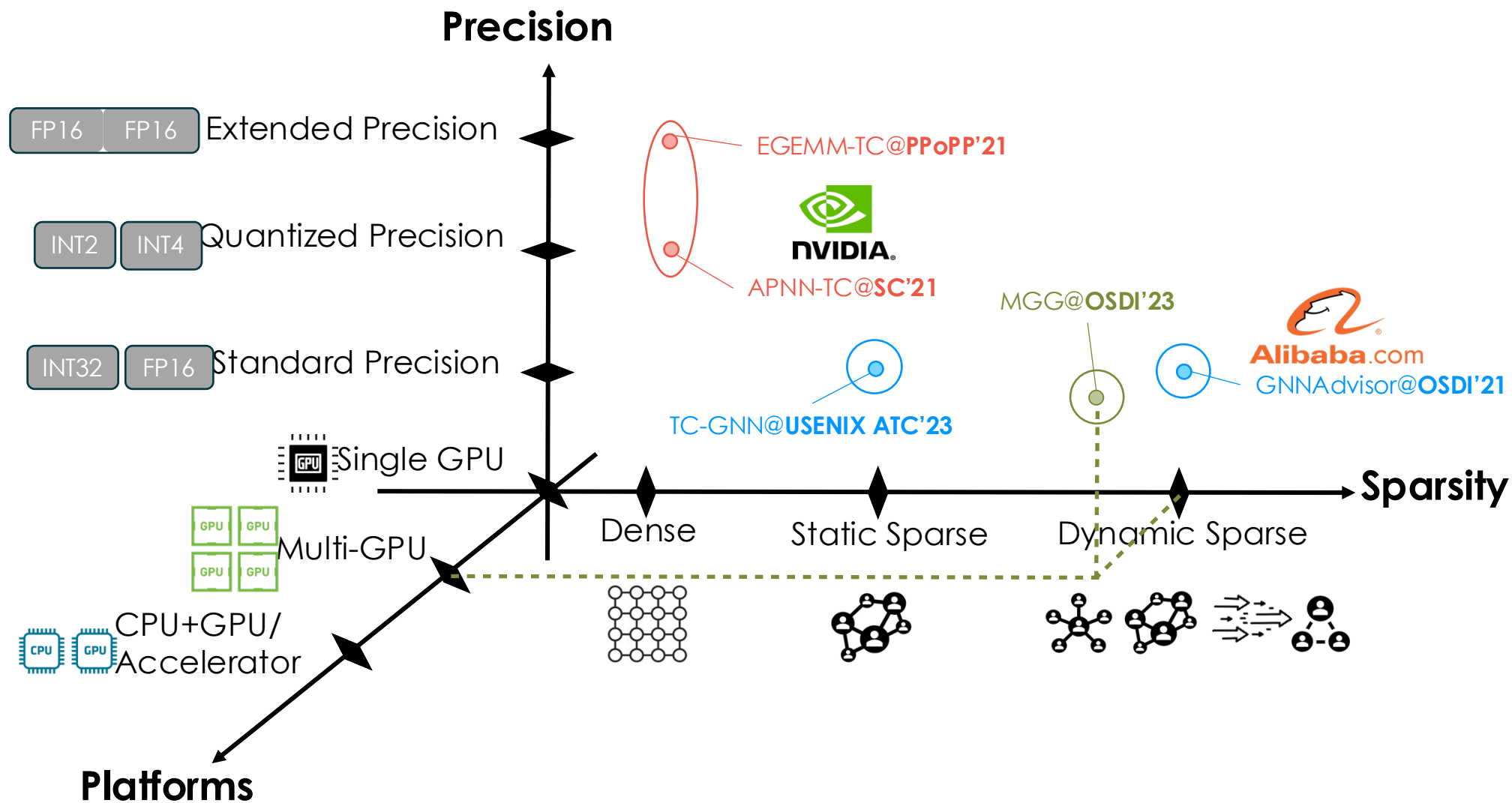
❖ Overview of my prior Ph.D. Research.



My Prior Ph.D. Research Recap



My Prior PhD Research Recap



Precision

Sparsity

Scalability

Diverse Precision Demands for DL Applications

❖ **Low-precision** quantized deep-learning applications.

Quantized Deep Learning	Precision Requirements
QNNs [JMLR'18]	1-bit Weight, 2-bit Activation for Vision Model, 3-bit Weight, 4-bit Activation for Language Model.
SGQuant [ICTAI'20]	Graph Attention Model: 2-bit Neighbor Attention , 4-bit Neighbor Aggregation .
LLM.int8() [NeurIPS'22]	8-bit Quantization for Transformers.
...	...

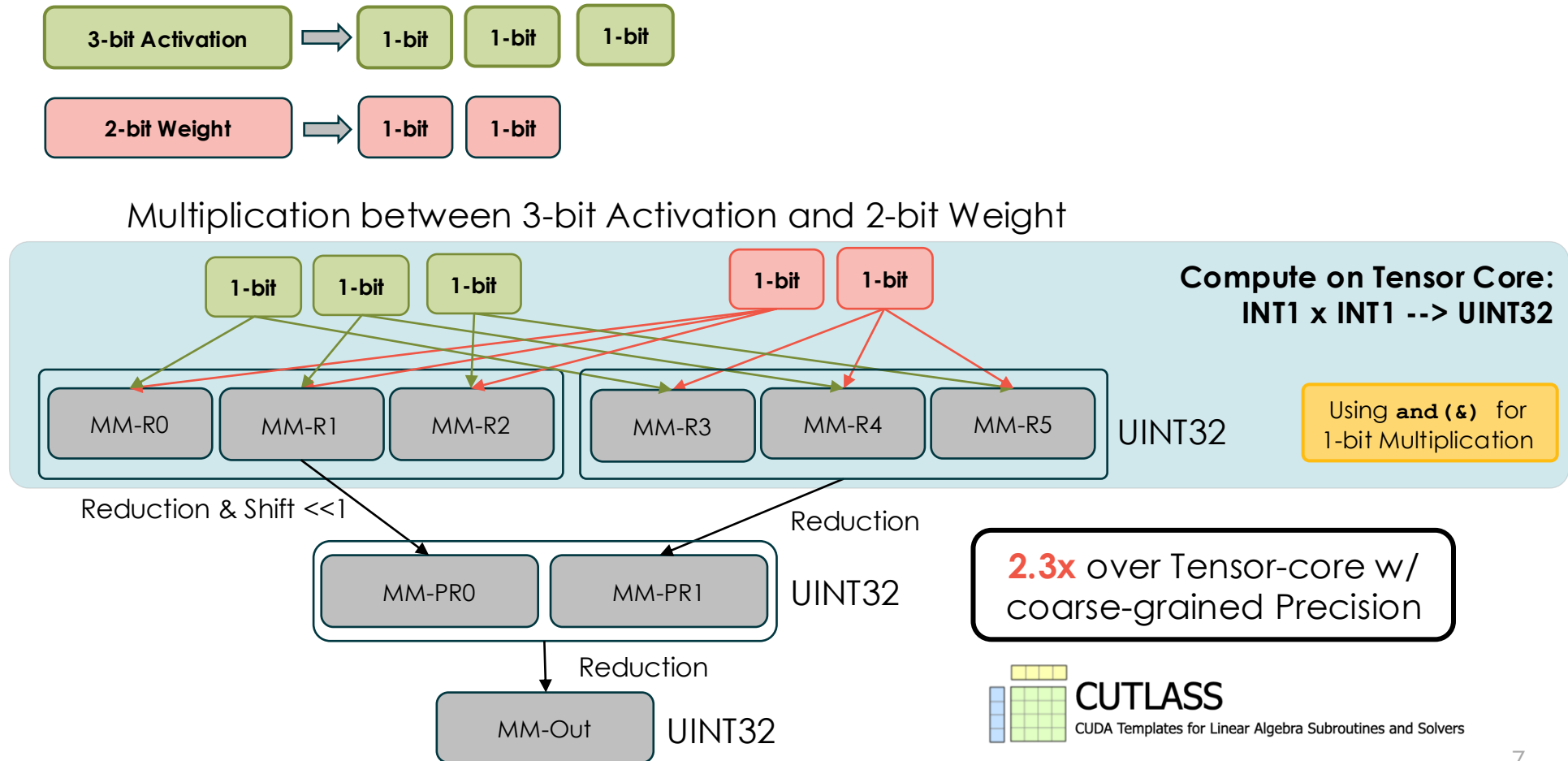


GPU Tensor cores would suffer from **low efficiency**.

Bit Composition for Quantized Deep Learning [SC'21]

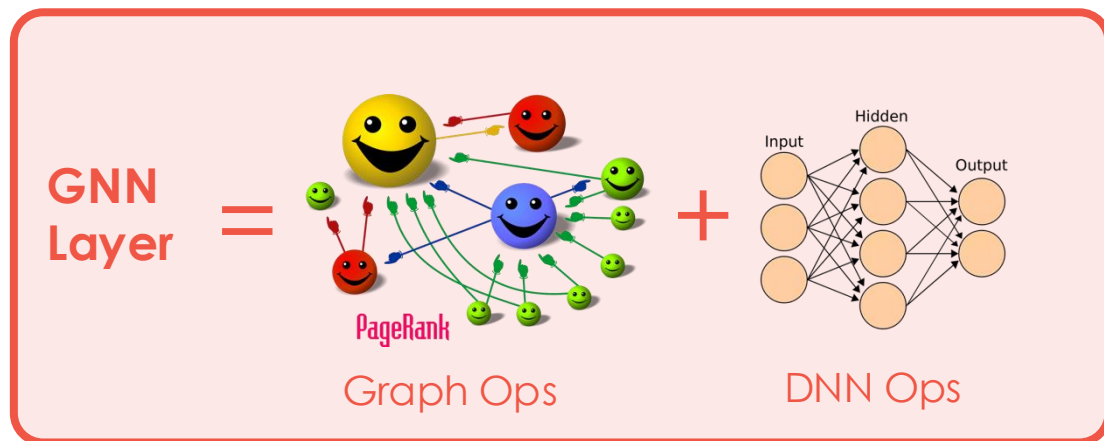
❖ Insight: **Quantized deep learning** can be **composed** with the **binary (1-bit) precision**.

Example of **2-bit** and **3-bit** Precision in Quantized DNN Computing.

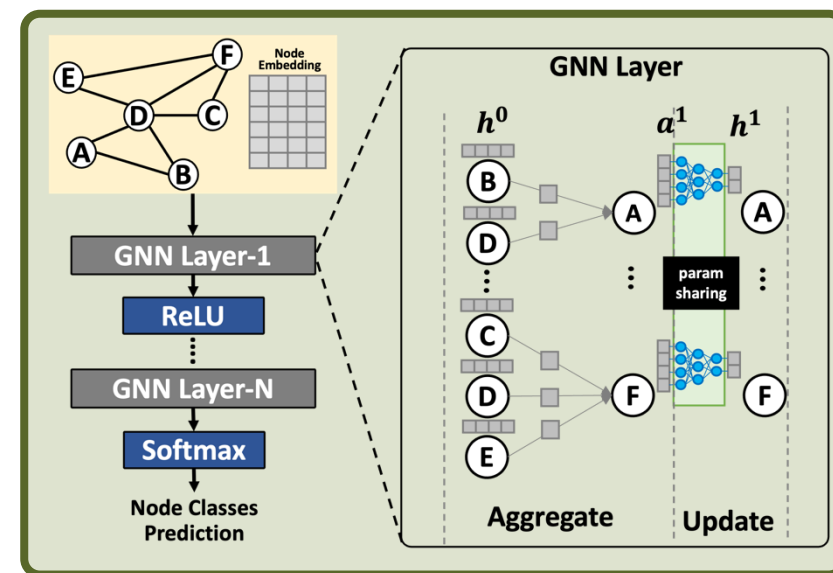


A Typical Paradigm of Graph Deep Learning

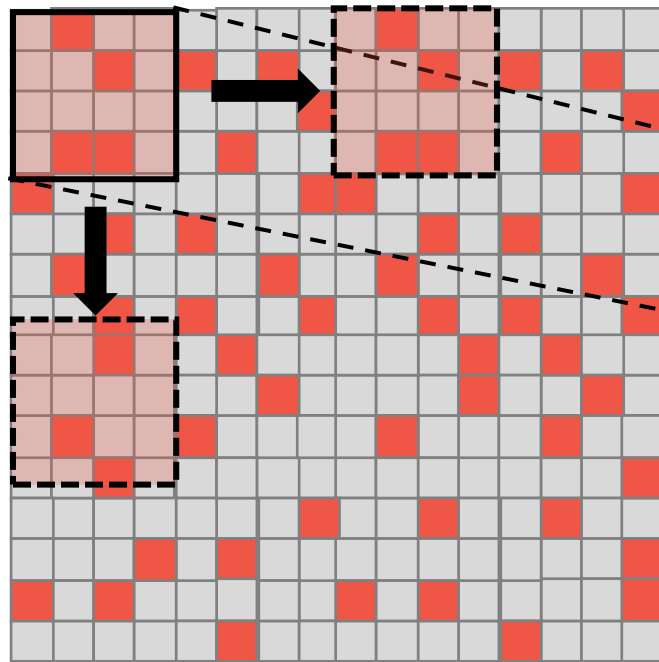
Operation view



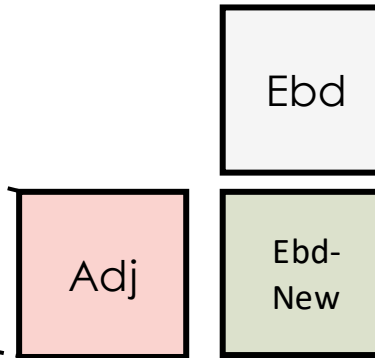
Model view



Challenge of Mapping Sparse Computing to Dense Units



Sparse Adjacent Matrix of Graph



Dataset	# Nodes	# Edges	Memory	Eff.Comp
OVCR-8H	1,890,931	3,946,402	14302.48 GB	0.36%
Yeast	1,714,644	3,636,546	11760.02 GB	0.32%
DD	334,925	1,686,092	448.70 GB	0.03%

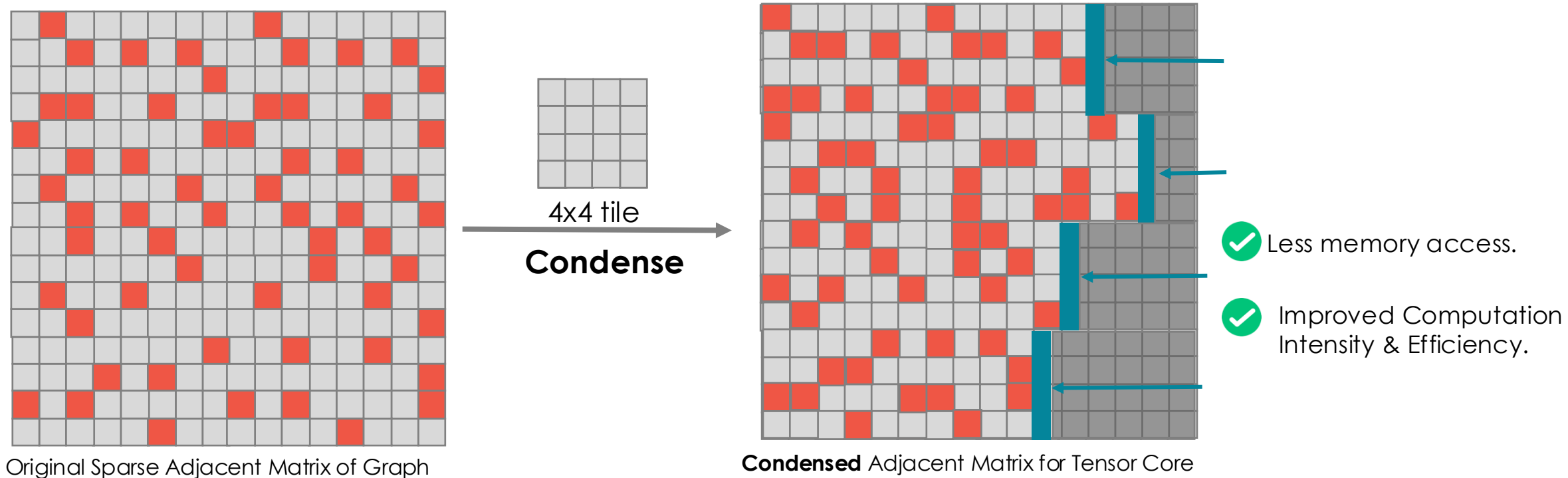
>> A100/H100
(80GB)

Largely **Wasted** Computation
and Memory Access

Direct mapping: **high memory** consumption
and **low computing** efficiency.

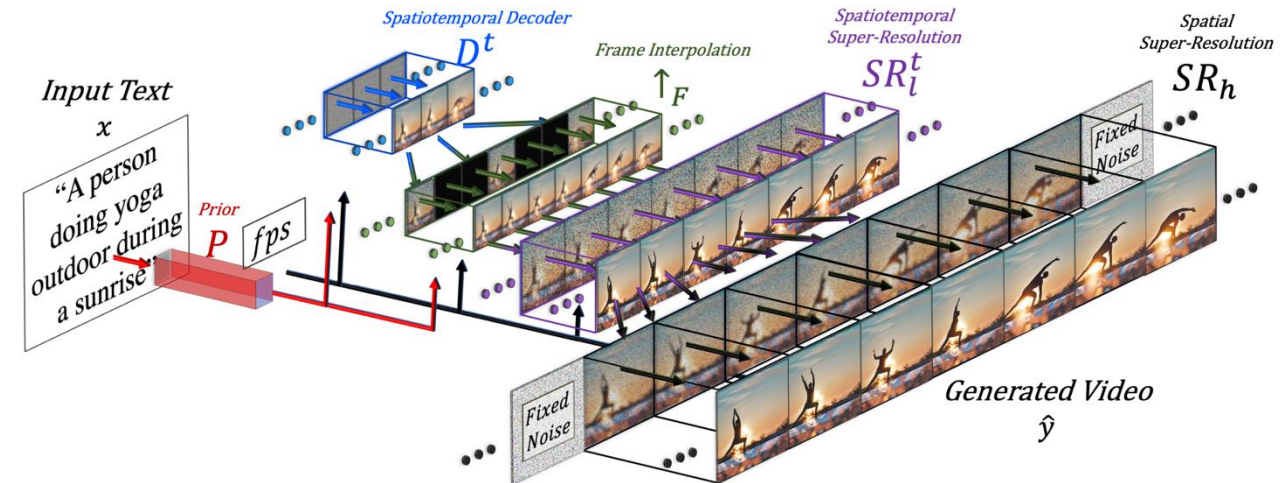
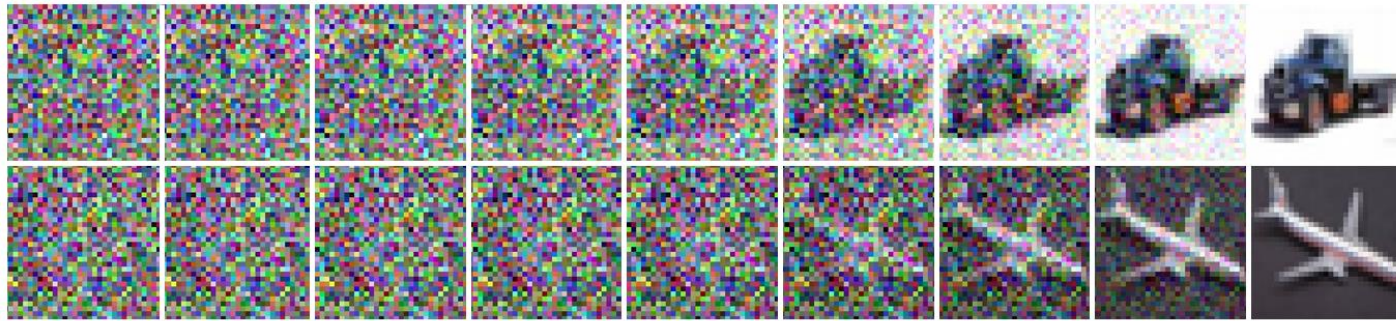
TC-GNN: Order-Invariant Transformation [ATC'23]

- Irregularly-scattered elements can be **condensed** to benefit high-performance dense GPU units.



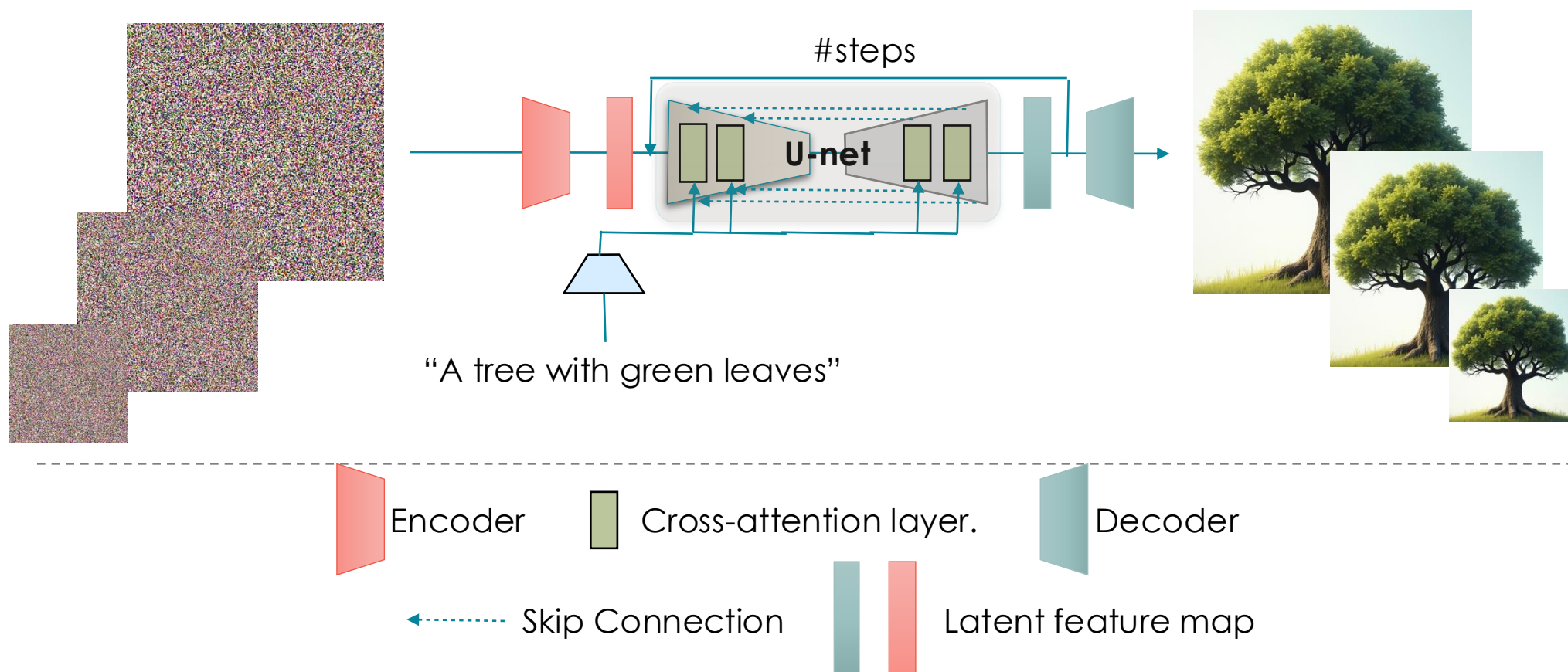
1.50x ~ 6.70x over DGL operators (cuSPARSE).
Incorporated by SparseTIR in [TVM](#) Project.

In The Era of Generative AI



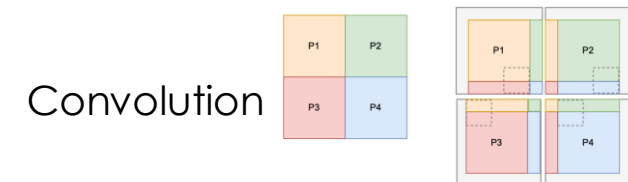
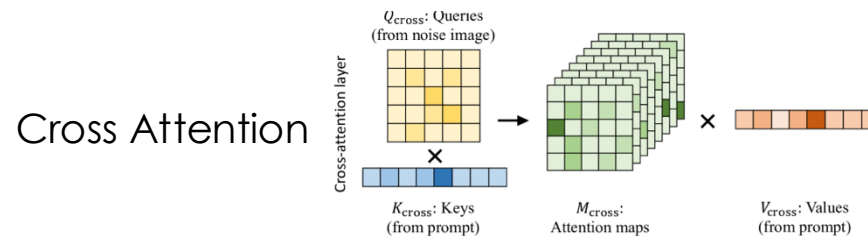
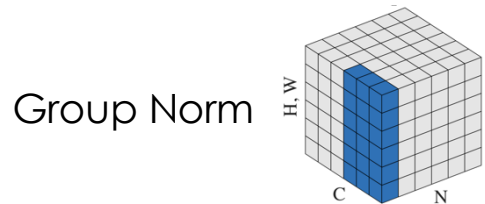
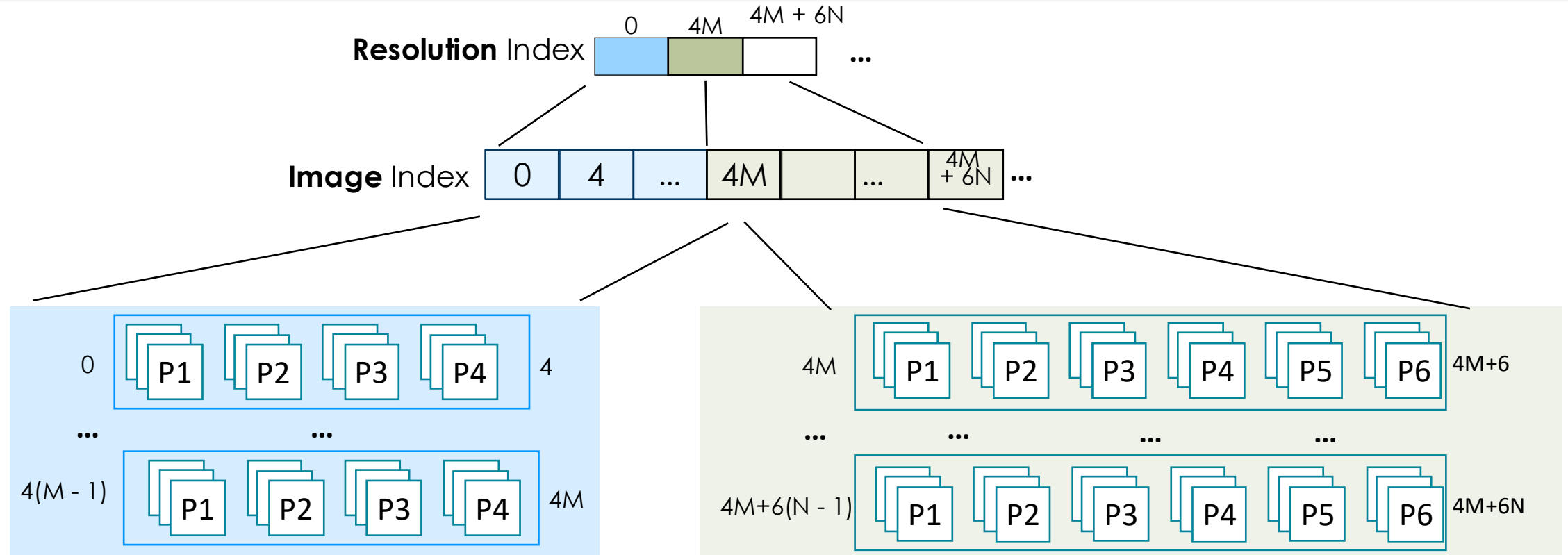
Patch-based Diffusion Serving [PPoPP'26]

- Architecture of **Image Diffusion** Model.



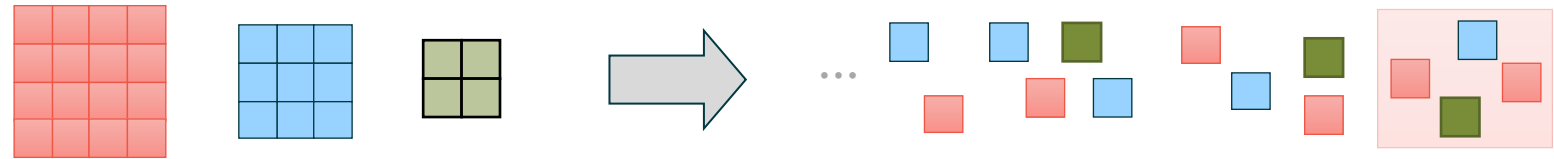
Note that Conv are omitted in Unit for simplicity.

Hierarchical Tile Storage and Indexing

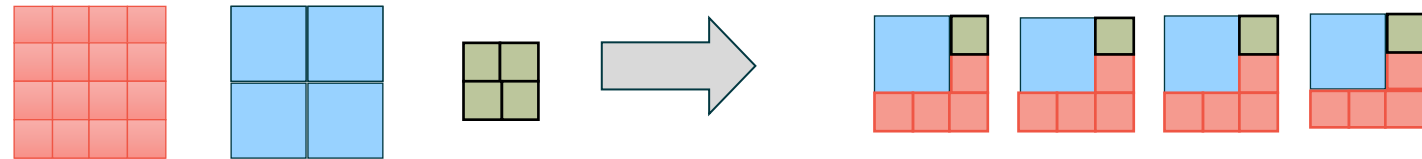


Tile-based Workload Scheduling

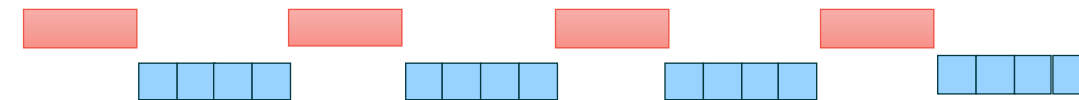
Unify the processing of different resolutions



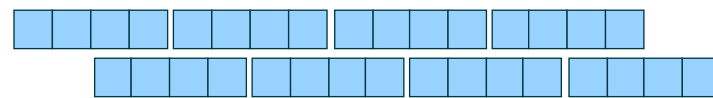
Control granularity for mixed workload composition



Unlock more fine-grained pipelining (FCFS)



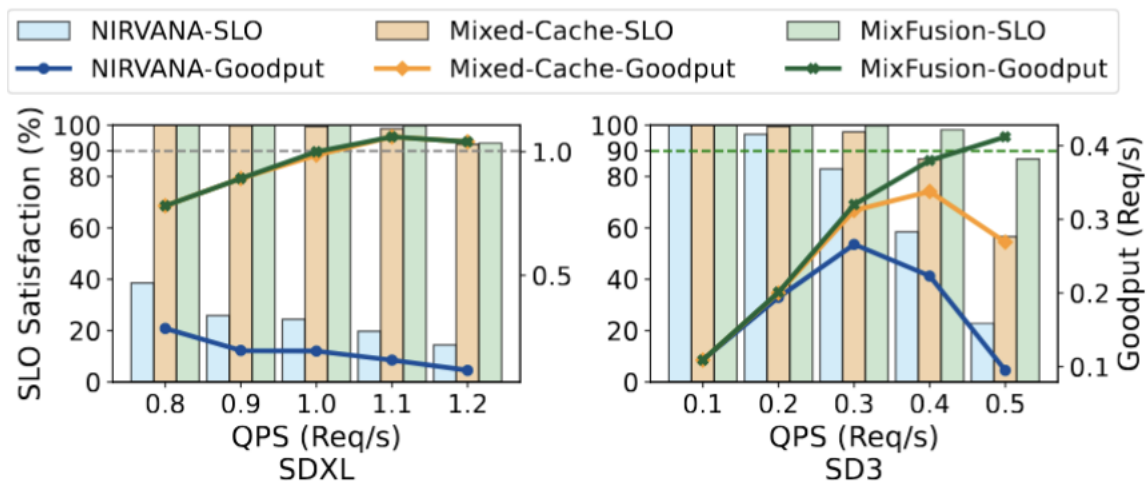
a. Resolution-mismatch for largely sequent processing.



b. Decomposed Resolution for batched processing.

Evaluation

- End-to-End SLO satisfaction Ratio



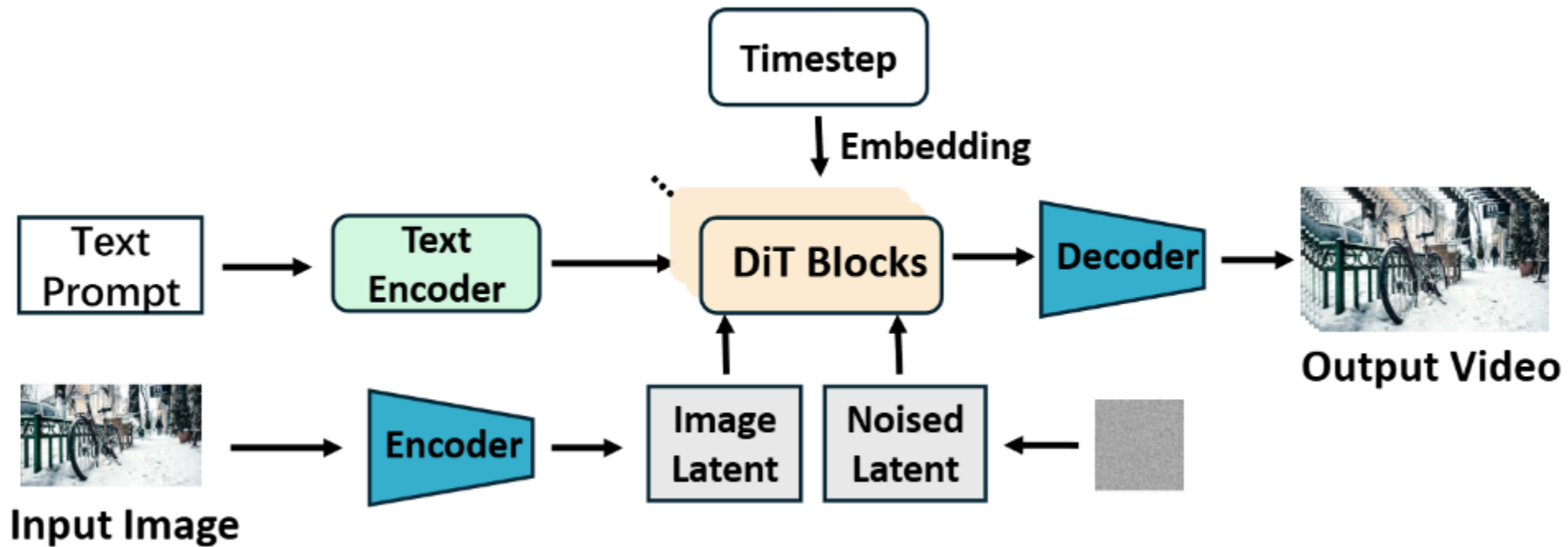
- Quality Score Comparison

Model	Method	SDXL		SD3	
		COCO	diffusiondb	COCO	diffusiondb
CLIP (↑)	Original	14.92	16.24	14.79	16.65
	MixFusion	15.43	16.62	15.13	17.06
FID (↓)	Original	31.92	35.56	28.94	32.38
	MixFusion	28.85	33.42	26.56	38.01

5.33× higher goodput when achieving **90 % SLO**
while over SOTA NIRVANA [1]

Ultra-Resolution Video Generation

- Architecture of Diffusion Transformer Model.



Observation of Ultra-high-resolution Video Generation

- Maximum supported resolution and running time

Model	Max Resolution	Frames	VRAM	Latency
CogVideoX1.5 5B [15]	1360 × 768	80	40GB	400s
HunyuanVideo 13B [20]	1280 × 720	128	70GB	1,800s

$$T_{4K} \approx 1800s \times \left(\frac{3840 \times 2160}{1280 \times 720} \right)^2 / 3600 \approx 40 \text{ h,}$$

Why only supports 720p

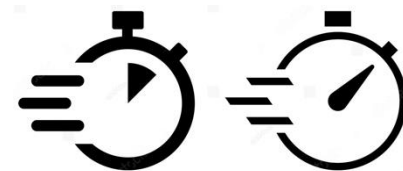
Not enough high-resolution training datasets. (Billions)



How to training-free and efficiently to generate ultra-high-resolution video?

Generation inefficiency

5-second 4K video directly could take a day and intensive memory cost.



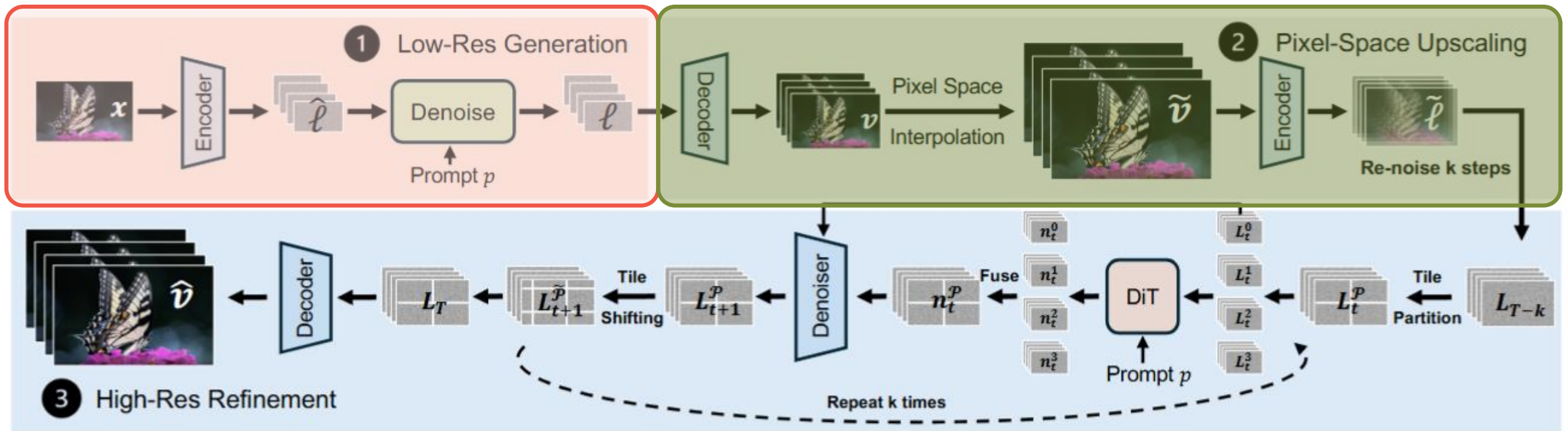
Current super-resolution video generation?

Only fixed scale factor, not training-free



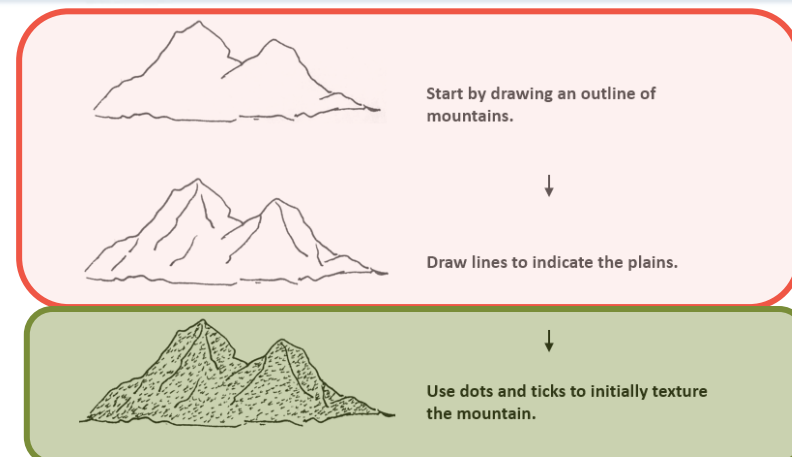
Opportunity: Sketching and Tiling

- Training-free Two-stage Generation



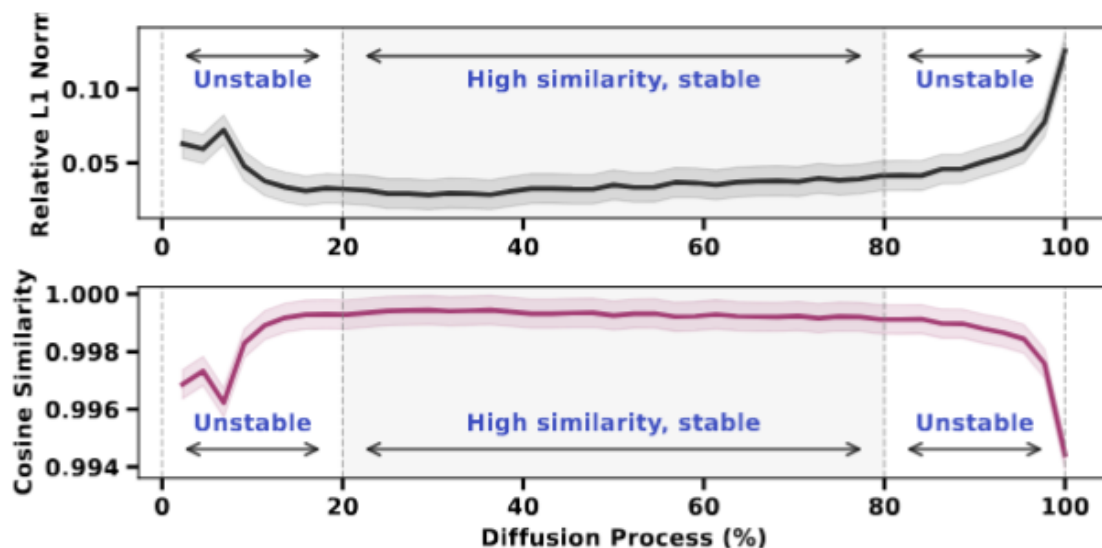
First stage: **Global Semantic Guidance**

Second stage: **Local Details Refinement**



System Support: Fine-grained Cache

- Observation: High similarity of predicted noise across timestep



$$L1_{\text{rel}}(\mathbf{O}, t) = \frac{\|\mathbf{O}_t - \mathbf{O}_{t+1}\|_1}{\|\mathbf{O}_{t+1}\|_1}, \text{CosSim}(\mathbf{O}, t) = \frac{\langle \mathbf{O}_t, \mathbf{O}_{t+1} \rangle}{\|\mathbf{O}_t\|_2 \|\mathbf{O}_{t+1}\|_2}$$

How to Utilize this Similarity?

System Support: Fine-grained Caching Strategy

- Fine-grained Region-aware Cache

Introduce cache residual $\delta_t \triangleq O_t - I_t$

Approximate $O_t \approx I_t + \delta_c$.

Error Accumulation $E_{c \rightarrow t} = \sum_{k=c+1}^t \|O_k - O_{k-1}\|.$

Involve transformation rate $k_t = \frac{\|O_t - O_{t-1}\|}{\|I_t - I_{t-1}\|}.$

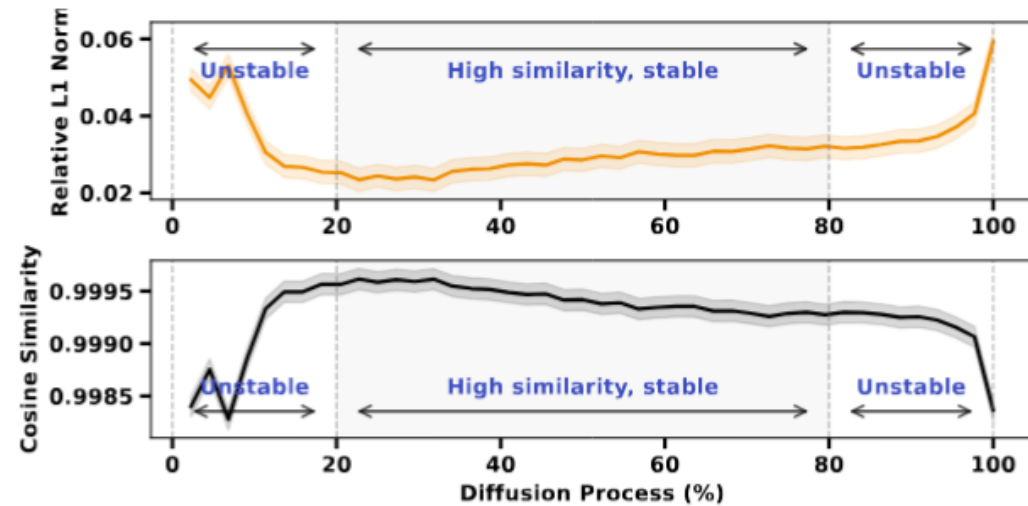
$$E_{c \rightarrow t} \approx \sum_{k=c+1}^t k_c \|I_k - I_{k-1}\| = k_c \sum_{k=c+1}^t \|I_k - I_{k-1}\| = k_c L_{c \rightarrow t},$$

Cache Decision
if $k_c L_{c \rightarrow t} < \tau \Rightarrow$ reuse cache at step t ,
else $k_c L_{c \rightarrow t} \geq \tau \Rightarrow$ recompute O_t and set $c \leftarrow t$.

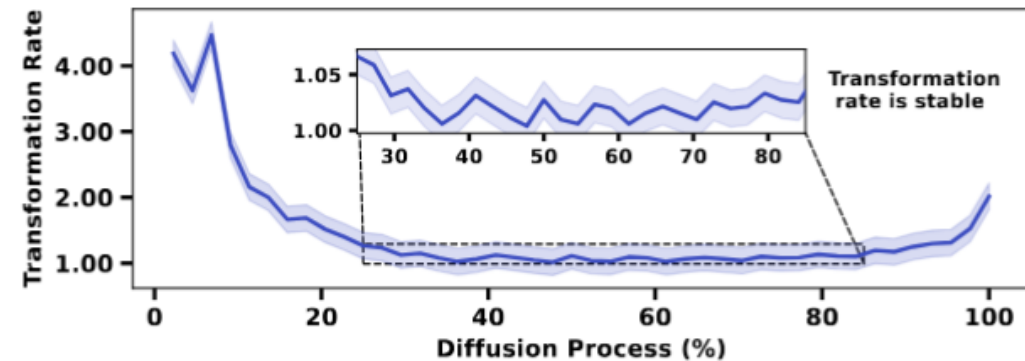
1) **Predict** next diffusion step **outcome** based on current outcome.

2) **Accumulate** the prediction **errors** to determine caching.

High similarity of Cache Residual

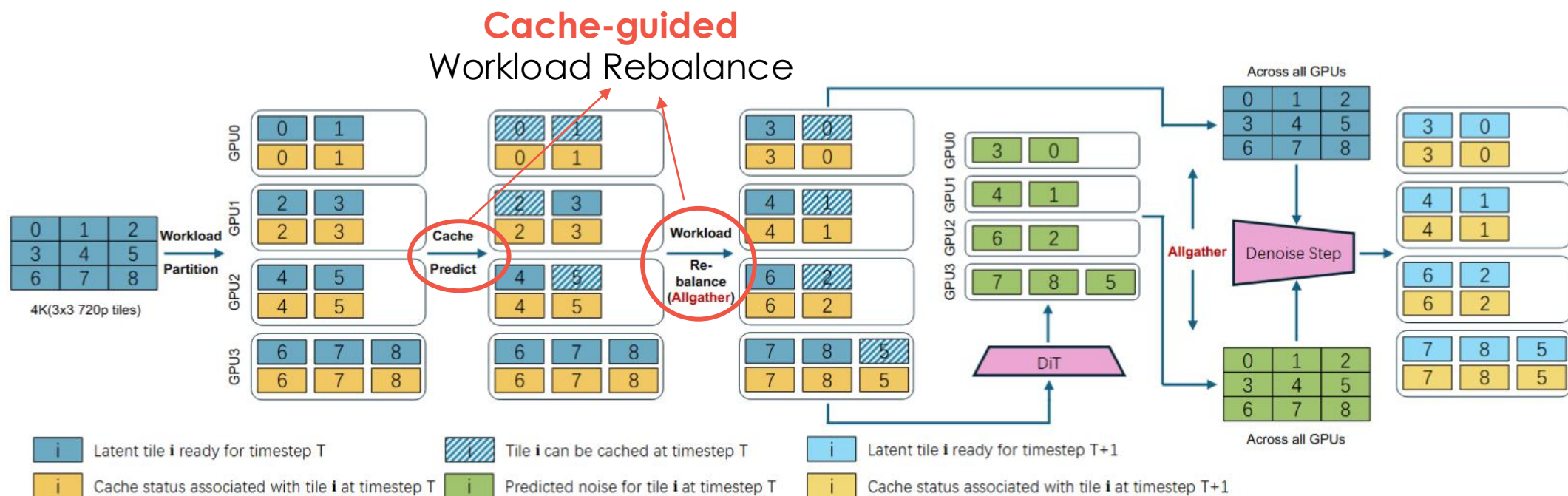


Transformation Rate K is Stable



System Support: Intelligent Cost-efficient parallelism

Workload imbalance: e.g., 9 tiles on 4 GPUs



KV in different tile are **independent**,
Minimum communication.

1) Reuse the computed tiles in **nearby GPUs**.

2) Reuse the computed tiles in from **prior iteration**.

Evaluation: Video Quality

Table 2: Quality results of SUPERGEN on VBench benchmark. V1–V5 denote the five evaluation metrics: **V1**: Subject Consistency, **V2**: Background Consistency, **V3**: Motion Smoothness, **V4**: Aesthetic Quality, and **V5**: Imaging Quality.

Model	Setting	V1(%)	V2(%)	V3(%)	V4(%)	V5(%)	Avg.
Cogvideo	720p	96.29	96.23	98.41	61.88	70.20	84.60
	2K w/o Cache	95.66	96.06	97.22	63.86	70.38	84.64
	2K w/ Cache	95.45	95.91	97.21	62.75	69.75	84.21
	4K w/o Cache	92.94	94.11	98.10	57.92	67.38	82.09
	4K w/ Cache	93.22	94.32	98.04	57.95	67.56	82.22
Hunyuanvideo	720p	98.55	97.86	99.53	64.54	70.83	86.26
	2K w/o Cache	98.02	97.31	99.42	66.47	69.62	86.17
	2K w/ Cache	98.30	97.48	99.44	66.16	70.26	86.33
	4K w/o Cache	97.76	97.24	99.43	63.07	69.68	85.44
	4K w/ Cache	98.12	97.58	99.51	62.57	70.31	85.62

Both 2K and 4K can achieve **high quality**.

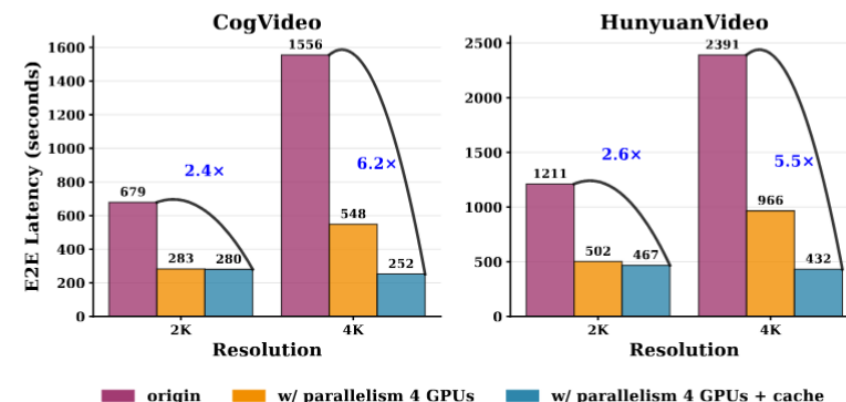
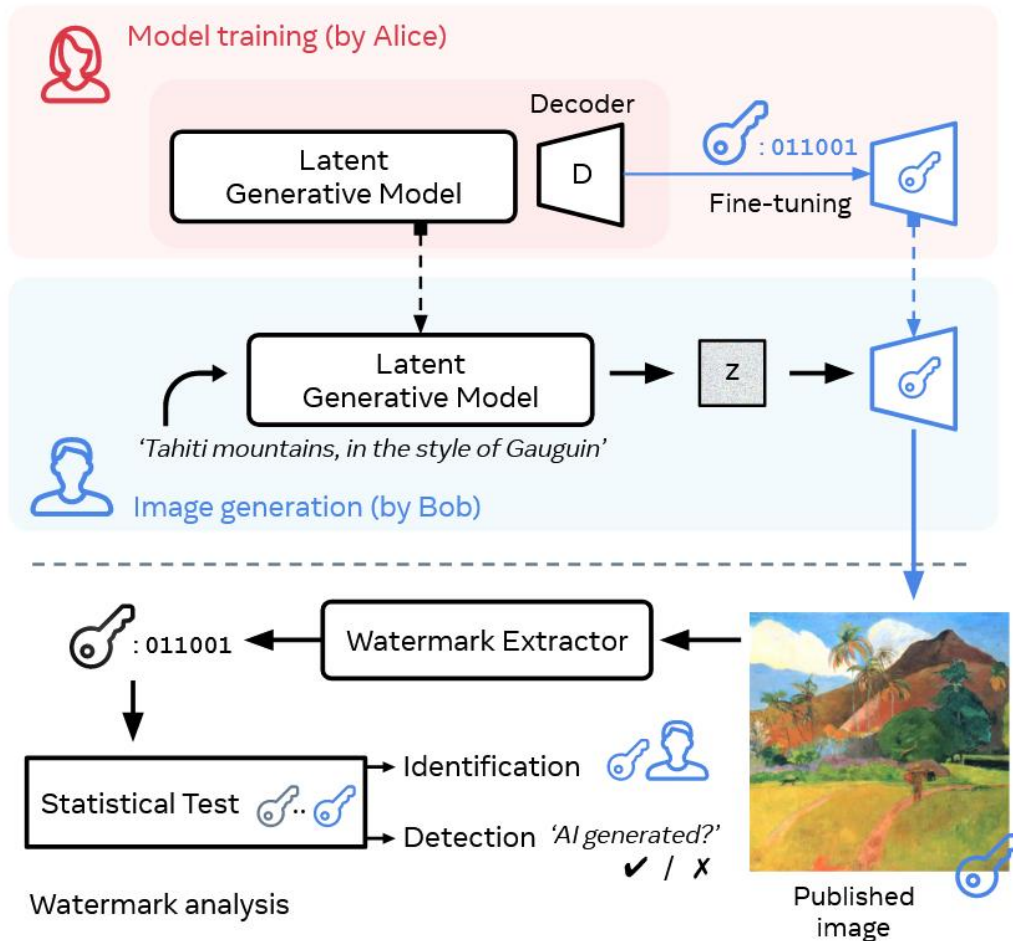


Figure 12: End-to-end latency. Origin setting is evaluated with 1 GPU without cache. The other two settings with parallelism are measured with workload rebalance.

End-to-end can achieve up to **6.2x** speedup.

Efficient and Adaptive Watermark Detection

- Overview of diffusion image watermarking



- Diffusion models now generate images nearly **indistinguishable** from real photos.
- Social platforms face **billions** of daily uploads (YouTube, TikTok, Facebook).
- Watermarking is critical for **verifying** AI-generated content and copyright attribution.

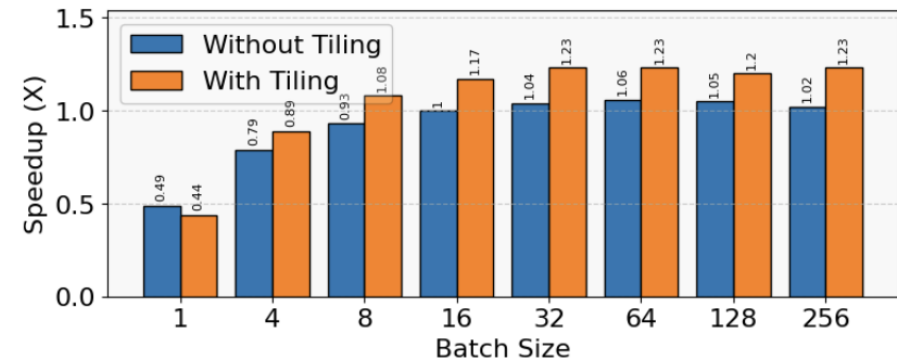
Challenges of Diffusion Watermarking Detection

- Validation accuracy of naive tile-based Stable Signature on different tile size.

None	128	96	80	64	48	32	16
0.997	0.960	0.933	0.897	0.875	0.804	0.714	0.624

Challenge-1: How to balance the trade-off between accuracy and efficiency?

- Naive tile-based design brings only limited improvement.



Challenge-2: How to optimize the tile-based watermark detection pipeline for higher throughput?

Motivation from Real-world Observations

Observation1: **One tile** of an entire image is **sufficient** for watermark detection.

None

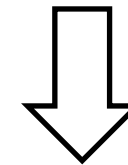


Crop 10%



Accuracy reduction: 99% -> 95%
(even without pretraining on tiles)

Observation2: A QR code can still be **accurately** and **quickly** even when part of it is **missing**.

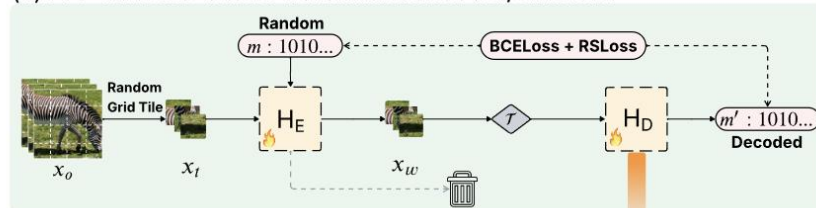


https://en.m.wikipedia.org/wiki/Main_Page

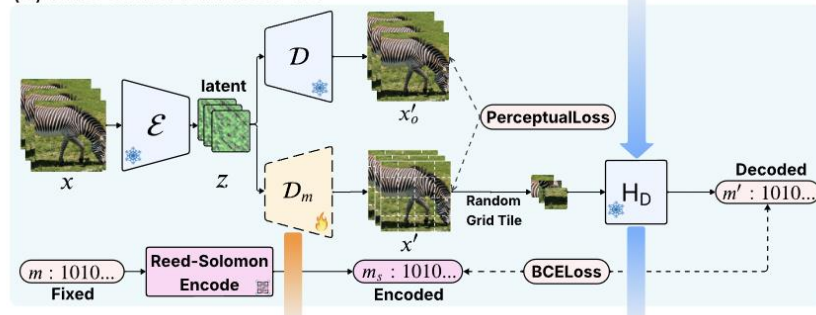
Algorithmic-System Co-Design Solution

- Tile-based detection with RS Correction

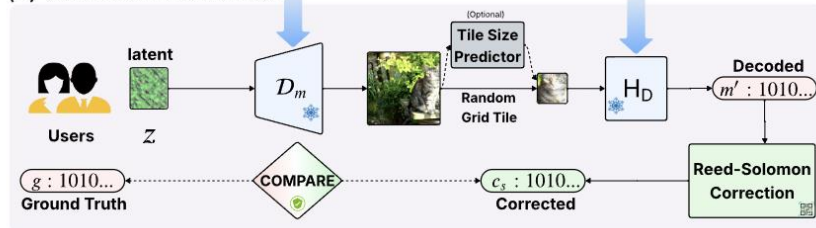
(a) Pre-train tile-based watermark encoder/extractor



(b) Fine-tune LDM decoder



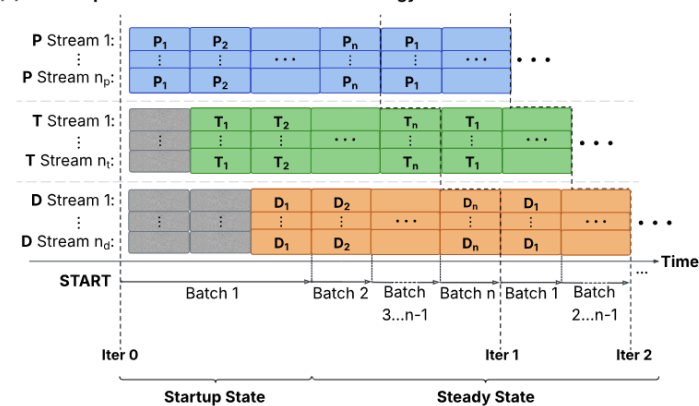
(c) Generate and Detect



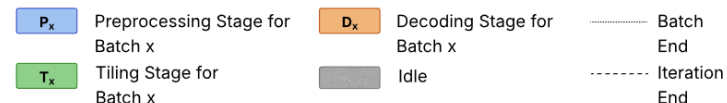
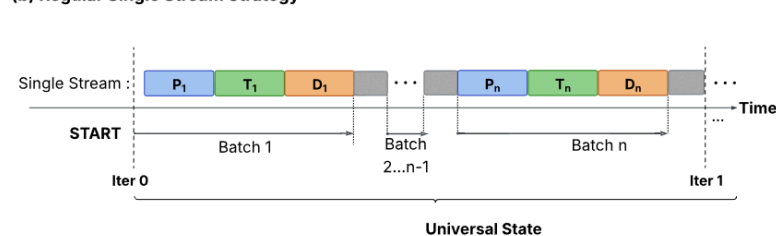
Leverage **RS correction** to offset the accuracy degradation caused by **tiling**

- Adaptively allocate CUDA stream for each stage

(a) Our Adaptive CUDA Stream Allocation Strategy



(b) Regular Single Stream Strategy



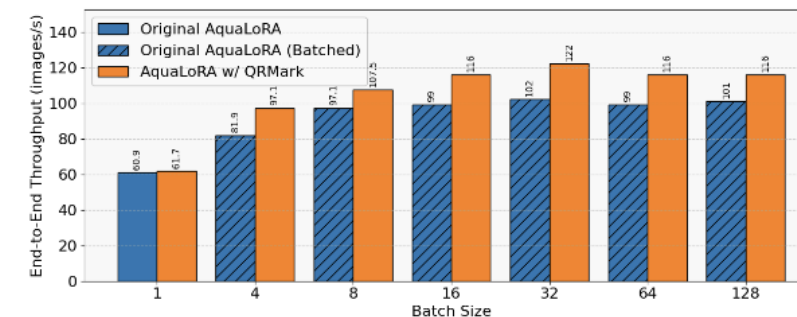
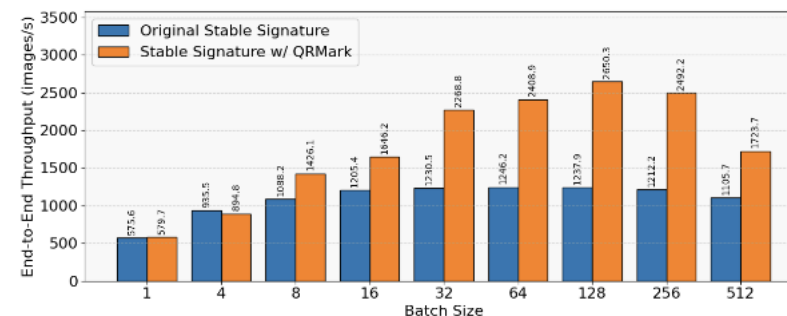
Equalize the **per-minibatch execution** time across stages

End-to-end Performance

- End-to-end accuracy and robustness

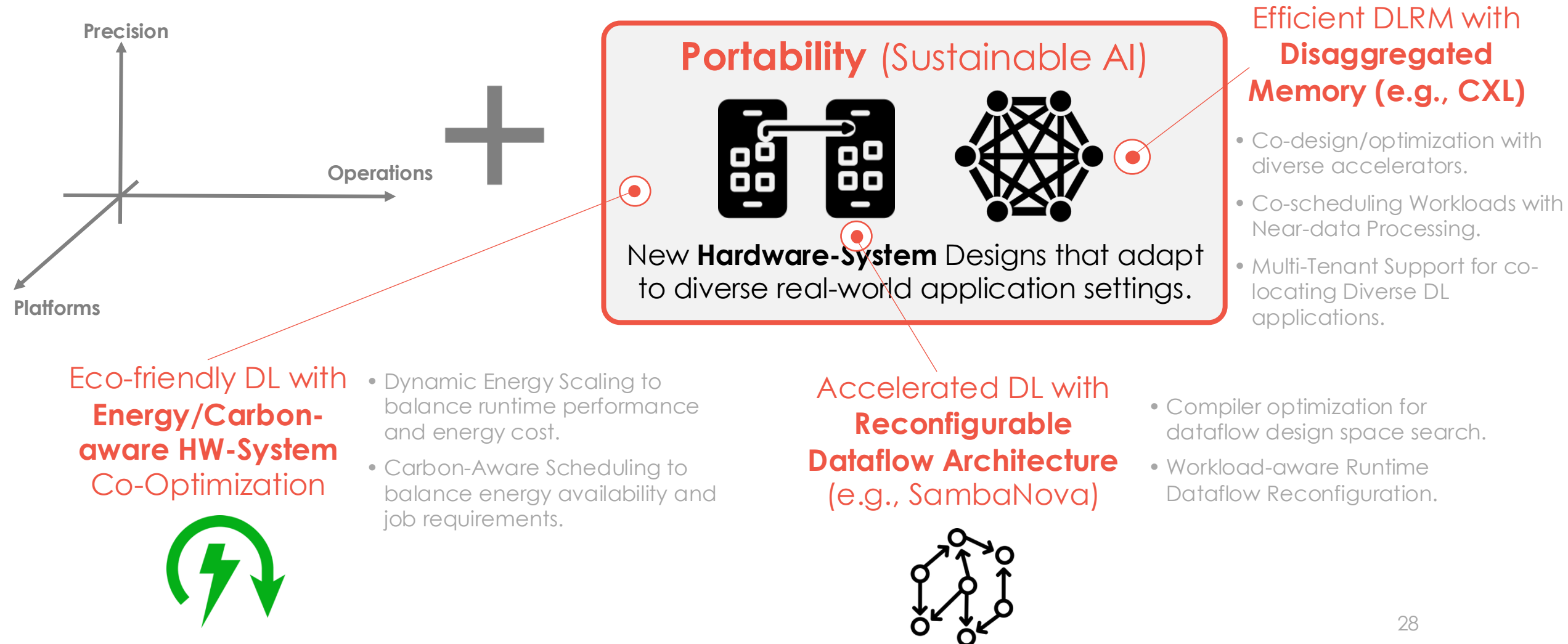
Model	TS	BitAcc. \uparrow	BitAcc. (ADV.) \uparrow	PSNR \uparrow	TPR \uparrow
StableQR	16	0.748	0.665	27.67	0.761
	32	0.989	0.907	29.47	0.993
	48	0.997	0.933	29.63	0.996
	64	0.999	0.945	30.35	0.998
	80	0.999	0.949	30.76	0.999
StableBL	—	0.999	0.974	30.05	0.993
AquaLoRA _{QR}	256	0.947	0.883	17.13	0.970
AquaLoRA _{BL}	—	0.958	0.912	17.65	0.985

- End-to-end Throughput

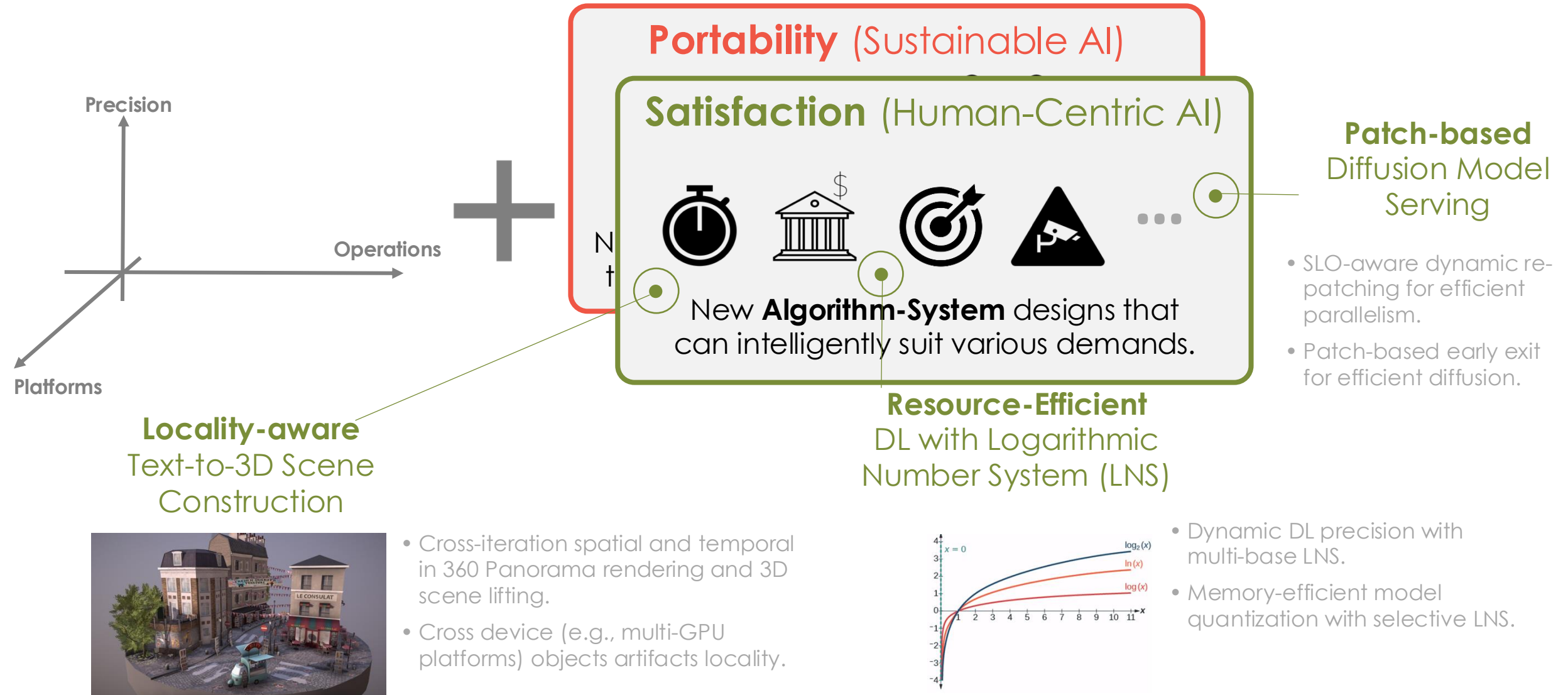


2.20x over Stable Signature [1]
and **2X** over AquaLoRA [2].

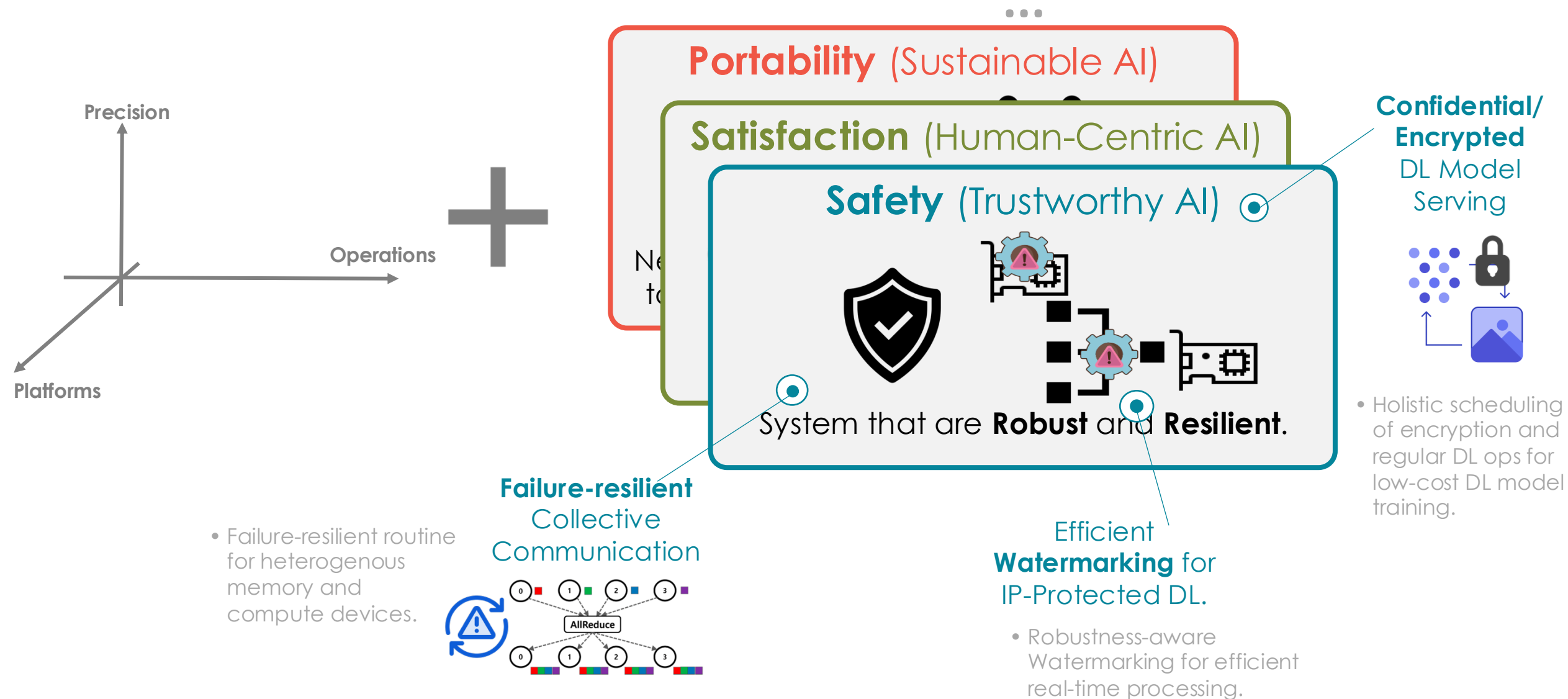
Future Research: New Hardware-System Optimization



Future Research: Exploring New DL Workloads



Future Future: Secure and Resilient DL System





RICE UNIVERSITY

Thank You

Q & A



yuke.wang@rice.edu



github.com/YukeWang96



wang-yuke.com