Efficient Inference for Large Language Models Algorithm, Model, and System

Xuefei Ning¹ Guohao Dai^{2,4}



Haoli Bai³



Lu Hou³



Yu Wang¹



Qun Liu³



¹Tsinghua University ²Shanghai Jiao Tong University ³Huawei ⁴Infinigence-Al foxdoraame@gmail.com, daiguohao@sjtu.edu.cn, baihaoli@huawei.com, houlu3@huawei.com, yu-wang@tsinghua.edu.cn, qun.liu@huawei.com

Tutorial Website: https://haolibai.github.io/emnlp-2025-tutorial-efficiency/



Background

Model-Level Optimization

2 Preliminary

- 6 System-Level Optimization
- Problem Definition & Conceptual Analysis
- Algo-Level Optimization

4 Practical Pipeline

8 Conclusion

- Background (Xuefei Ning)
- Preliminary (Xuefei Ning)
- Problem Definition & Conceptual Analysis (Xuefei Ning)
- Practical Pipeline (Xuefei Ning)

- Model-Level
 Optimization
 (Lu Hou)
- System-Level
 Optimization
 (Guohao Dai)
- Algo-Level
 Optimization
 (Haoli Bai)
- 8 Conclusion (Haoli Bai)

Background

Model-Level Optimization

2 Preliminary

- 6 System-Level Optimization
- 3 Problem Definition & Conceptual Analysis
- Algo-Level Optimization

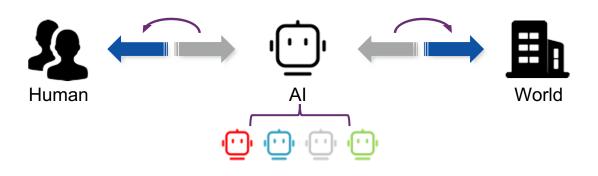
4 Practical Pipeline

8 Conclusion

Background

Towards general and generative intelligence, scaling up model / data / computation based on Transformer is a mainstream and effective pathway.

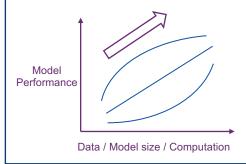
Focus Shift 1: Discriminative => Generative



Focus Shift 2: Specialized => General

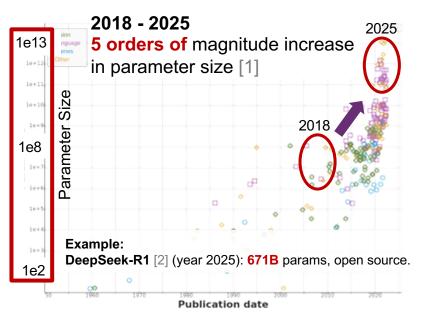
Current mainstream path

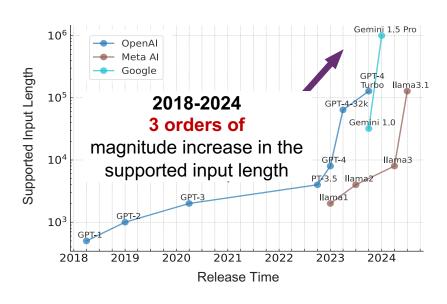
Based on the Transformer architecture, scaling model size, training data / computation, and test computation.



Background

The model scale and input/output length in generative intelligence research and applications have increased significantly.





- [1] Villalobos et al. "Machine Learning Model Sizes and the Parameter Gap." arXiv 2022.
- [2] Guo, Daya, et al. "Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning arXiv 2025.
- [3] Esser, Patrick et al., Scaling rectified flow transformers for high-resolution image synthesis, ICML 2024.

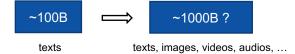
[4] Achiam, Josh, et al. "Gpt-4 technical report." arXiv 2023. [5] Reid, Machel, et al. "Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context." arXiv 2024.

Background

Application demands e.g., multimodal input, advanced tasks, may continually drive increases in model scale and input/output length.

Multimodal Input

Modeling multimodal data *may* push the scaling saturation point of model size higher.



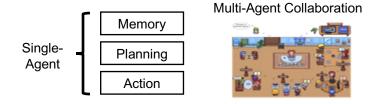
Supporting high-resolution images and longer videos requires extended input/output contexts.





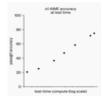
Advanced Tasks

Agentic pipelines for broad applications require stronger models and longer input/output contexts.

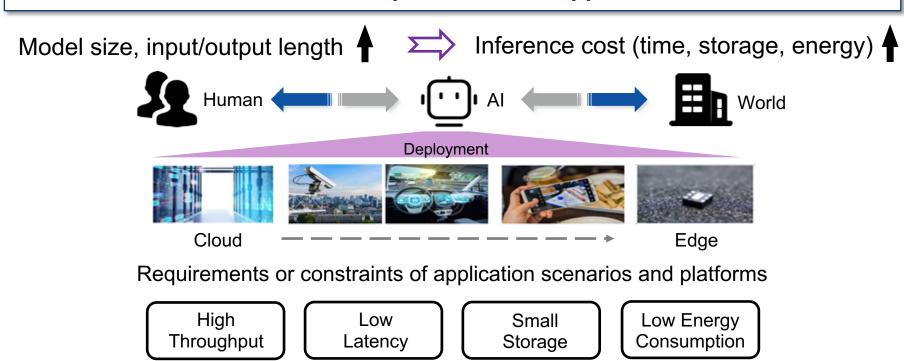


Test-time compute scaling (especially CoT) for reasoning require longer input/output contexts.

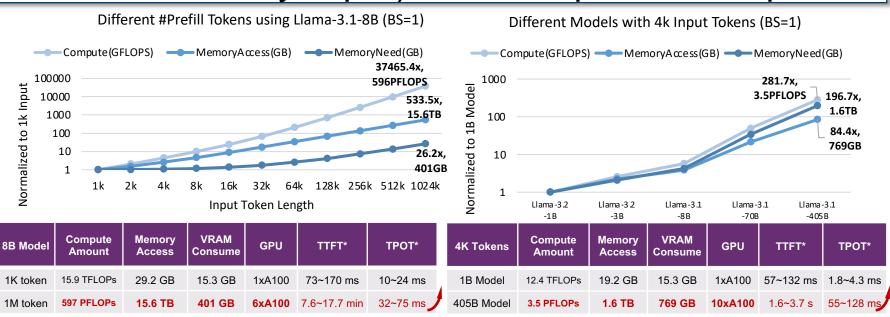




This scaling of the model size and input/output poses challenges for efficient inference across platforms and application scenarios.



<u>Time to first token (TTFT)</u> and <u>time per output token (TPOT)</u> w.r.t. model size and input/output lengths, <u>estimated</u> with three costs (compute, memory access and memory footprint) of the model spec and device spec.



^{*}TTFT is estimated using **Compute Amount** / (**Peak OPS** x **compute-util**), as prefill is computation-bounded. We assume a compute utilization range 30%~70% to report the estimation. TPOT is estimated using **Memory Access** / (**Bandwidth x bandwidth-util**), as decoding is memory-bounded. We assume a bandwidth utilization range 30%~70% to report the estimation. A100 Peak Compute Performance (FP16) = 312 TFLOPS; Peak Bandwidth = 2 TB/s

Per-request energy consumption w.r.t. model size and input/output lengths, estimated with actual latency measurement, device & schedule assumptions, and device spec.

| Model | Compa ny; Host (Device) | Date | Estimated energy consumption (100in-300out) (Wh) | Estimated energy consumption (10kin-1.5kout) (Wh) |
|--------------------|----------------------------------|--------------|---|--|
| GPT-4.1 nano | OpenAl; | | 0.10±0.04 | 0.45±0.21 |
| GPT-4.1 mini | Azure (H200& H100) | Apr, 2025 | 0.42±0.20 | 1.59±0.80 |
| GPT-4.1 | , | | 0.92±0.50 | 4.23±1.97 |
| LLaMA- 3.1-8B | Meta; | | 0.10±0.02 | 0.60±0.09 |
| LLaMA- 3.1-70B | AWS (H200& | Jul, 2024 | 1.10±0.13 | 11.63±1.39 |
| LLaMA- 3.1-405B | H100) | | 1.99±0.32 | 20.76±1.80 |



Jegham et al. made an attempt to estimate the environmental footprint of LLM inference at perprompt level of commercial AI providers, based on assumptions on the infrastructure and scheduling.

^[1] Jegham, Nidhal, et al. "How hungry is ai? benchmarking energy, water, and carbon footprint of Ilm inference." arXiv preprint arXiv:2505.09598 (2025).

Overall energy consumption, estimated with reported usage, device & schedule assumptions, and device spec.

A rough estimation result: Assuming a total of 772 billion queries (estimated with OpenAl 2024 report and the usage growth pattern, assume 80% short queries) annually in 2025, GPT-40 inference require approximately **4**×**10**¹¹**Wh**, exceeding the total electricity consumption of **35,000 U.S. residential households**.



Jegham et al. made an attempt to estimate the environmental footprint of LLM inference at perprompt level of commercial AI providers, based on assumptions on the infrastructure and scheduling.

1 Background

Model-Level Optimization

2 Preliminary

- 6 System-Level Optimization
- 3 Problem Definition & Conceptual Analysis
- Algo-Level Optimization

4 Practical Pipeline

8 Conclusion

2 Preliminary

- 2.1: LLMs' Model & Algorithm
- 2.2: Basics of Software and Hardware
- 2.3: NVIDIA GPU's Software and Hardware
- 2.4: Deployment

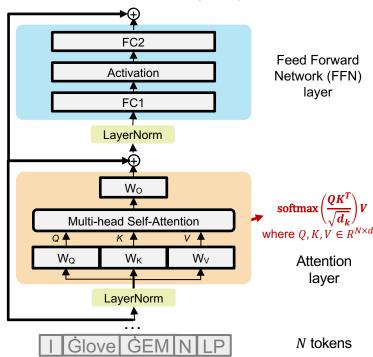
2 Preliminary

- 2.1: LLMs' Model & Algorithm
- 2.2: Basics of Software and Hardware
- 2.3: NVIDIA GPU's Software and Hardware
- 2.4: Deployment

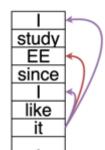
Model: Architecture of LLMs

- Most LLMs are based on Transformer architecture^[1], consisting of an input embedding layer, Transformer blocks, and a decoding layer.
- A Transformer block consists of:
 - Attention-Linear (transform for Q, K, V, O)
 - Multi-Head Self-Attention
 - Feed Forward Network (FFN)
 - Layer Norm

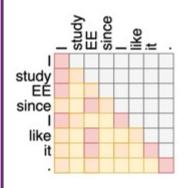
Residual Stream



Attention Layer

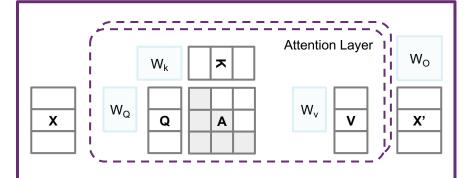


The idea of attention module is to establish **token-to-token** "attention" relationships within a sequence.



This relationship is modeled by an **attention matrix**, where each row represents one token's attention distribution to **previous** tokens (sum up to 1).

Causal attention mask: Each token only has positive attention to previous tokens



Calculation of each attention head:

- $Q = XW_Q$; $K = XW_K$; $V = XW_V$ to map $X \in \mathbb{R}^{N \times D}$ to $Q, K, V \in \mathbb{R}^{N \times d}$
- Calculate inner product $S = QK^T$, apply mask and softmax to get the attention matrix $A = \operatorname{softmax}(S + M) \in \mathbb{R}^{N \times N}$
- X' = AV to get the output

Generative Modeling Algo: Autoregressive Model

- Contemporary LLMs are generative models that uses the "autoregressive" generative modeling method.
- The **core task of generative modeling** is to learn a parametrized model $p_{\theta}(x)$ from observed data $\left\{x^{(i)}\right\}_{i=1}^{N}$, which in some sense capture the unknown real distribution of data $p_{\text{data}}(x)$, and can do stochastic sampling (i.e., sample generation) & probabilistic inference (e.g., likelihood estimation).
- Autoregressive models are a family of generative modeling methods that models the joint probability of a token sequence $x = [x_1, x_2, ..., x_n]$ as a product of conditional probability distributions, each conditioned on the preceding tokens:

$$p(x_1, x_2, ..., x_n) = p(x_1) \prod_{i=2}^{n} p(X_i = x_i | x_1, ..., x_{i-1})$$

Sampling Process of LLMs

LLMs are autoregressive models that uses one transformer to model the conditional distributions.

• Denoting the vocabulary set as |V|, the transformer maps the token sequence $x_{< i} \in V^{i-1}$ to a sequence of logits $l \in \mathbb{R}^{(i-1)\times |V|}$, where the logits $l[i-1] \in \mathbb{R}^{|V|}$ corresponding to x_{i-1} is regarded as the categorical distribution's parameter of x_i .

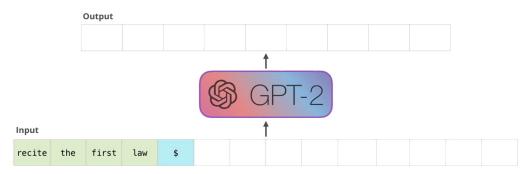
The KV cache design for sampling from causal LLMs:

- At first glance, in each step of autoregressive sampling: $x_i \sim p_{\theta}(x_{< i})$, the model forward process has $O((i-1)^2)$ computation complexity. Thus the overall sampling process (assume N steps) has $O(N^3)$ complexity.
- Most contemporary LLMs choose to use causal attention: The calculation of features of x_i only attends to $x_{< i}$. In this way, newly sampled tokens $x_{> i}$ don't influence the features corresponding to x_i .
- This enables us to "cache" already calculated features (specifically, the Key and Value) of preceding tokens to avoid recalculating their features, thus reduce overall computation complexity of sampling to $O(N^2)$.

Sampling Process of LLMs

After applying the KV cache technique, a typical LLM inference/sampling* process can be divided into two stages:

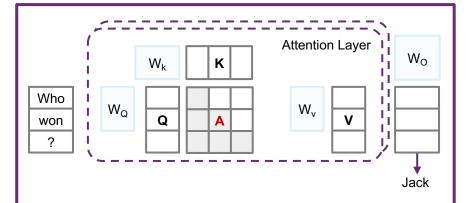
- Prefilling Stage: The LLM calculates and saves the KV cache of the initial input tokens, and samples the first output token.
- Decoding Stage: The LLM samples the output tokens one by one with the KV cache, and in the meantime updates the KV cache.



^{*}In this tutorial, the terms "inference" and "sampling" will be used interchangeably. Although they differ in the context of probabilistic modeling, here they both refer to either a single sampling step—i.e., a forward pass of the model—or the overall sampling process, depending on context.

Attention Computation in Prefill/Decode Stages

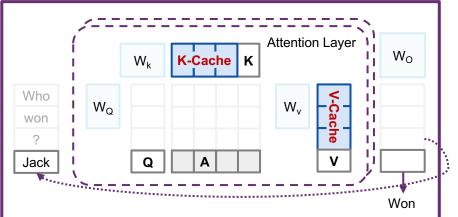
Prefill



The LLM calculates and saves the KV cache of the initial input tokens, and samples the first output token.

Calculate N query, key, value; calculate $N \times N$ -sized attention matrix A.

Decode



The LLM samples the output tokens one by one with the KV cache, and in the meantime updates the KV cache.

Calculate 1 query, key, value; read N key, value from KV cache; calculate $1 \times (N + 1)$ -sized attention matrix A.

2 Preliminary

- 2.1: LLMs' Model & Algorithm
- 2.2: Basics of Software and Hardware
- 2.3: NVIDIA GPU's Software and Hardware
- 2.4: Deployment

Software, Hardware, Instruction

Let's first have an overview of the most basic terms of interest:

- What is a computing hardware? Hardware consists of microelectronic components that transfer and transform electrical signals. Through these physical operations, hardware realizes high-level abstractions, that is executing functionalities described by instructions, such as storing, transferring, and transforming data.
- What is software? All software -- no matter it is an OS, a game, or an Al model's inference engine -- is a structured pack of instructions and data.
- Instructions are the interface between software and hardware. Instruction Set Architecture (ISA) defines this instruction interface: what instructions are available, how they are encoded and executed.

Instruction Set Architecture (ISA)

Define the instruction interface of a device

Example 1: CPU ISA defines instruction format and meaning. [1] Typical ISA for CPU includes MIPS, RISC-V, x86, x86-64 (amd64), etc.

MIPS32 Add Immediate Instruction

| 001000 | 00001 | 00010 | 0000000101011110 |
|----------------------|--------|--------|-------------------|
| OP Code | Addr 1 | Addr 2 | Immediate value |
| OP Code Equivaler | A | h | addi Sr1, Sr2, 35 |

Example 2: NVIDIA Ampere GPU architecture has compute capability sm80, which indicates its supported certain instructions (ISA) and some microarchitectural features

A Conceptual Layered Overview

Runtime Lib Toolchain OS / Driver Software Orchestrate resources or provide Interact with Compile instructions for host and device low-level resource-access API device at runtime **Hardware System Accelerator Device(s) Host Machine** (CPU & Memory & ...) **Off-Chip Memory** Al Chip **Accelerator Chip Microarchitecture** Define functional modules (controller, computation, Physical Design / memory) and how they "connect": (1) control path: how the controller parse the instruction and distribute control information; (2) data path: how modules parse and pass data based on control information **Functional Modules Circuits**

- Toolchain produces hardware-specific instructions.
- Runtime library manages program execution, send instruction to device, transfer data to/from device, and optionally call toolchain dynamically.
- A hardware system consists of host, accelerator devices, and their interconnects

Accelerator's Instruction Set Architecture (ISA)

Define the instruction interface of the AI accelerator device

- A microarchitecture design organizes modules to implement the accelerator's ISA
- A functional module combines circuits to perform certain functions
- Basic logic and memory elements

A Composition Overview

 A system primarily consists of the following components: Data cache PCIe Bandwidth: ~100GB/s during Host Al Chip computation Global control CPU On-chip cache SRAM Capacity: tens of MB Capacity: The total amount of data that a memory component can hold Bandwidth: 10-Weight cache Input/Output cache 100 TB/s Off-chip Large-capacity memory storage Processing unit DRAM (Use HBM as an example) regis Computation Vector/Scalar Capacity: tens of GB Matrix/Tenso Bandwidth: ~TB/s processing core HBM, DDR, r processing unit LPDDR. unit Bandwidth: The rate at which data can be transferred between two components of a system. Computation Interconnection Controller flow control of Connecting with other AI chips interface Al chips

Common Hardware Types

- Common hardware types: CPU, GPU, FPGA, ASIC
 - Chip-level metrics/characteristics: peak compute performance (OPS), energy efficiency (OPS/W), power (W), and area (mm²).
 - Also need to consider: generality and suitableness to algorithms; comprehensiveness of the software ecosystem.



CPU

Central Processing Unit High generality Low Al computing performance



GPU

Graphics Processing Unit
Strong parallel computing
capability
High-bandwidth memory



FPGA

Field-Programmable Gate Array
Hardware programmability
Flexibility and
reconfigurability



ASIC

Application Specific Accelerator
High specialization
High performance and low
power consumption

General More general and fine-grained ISA. Rely on software to implement coarse-grained operators, algorithms.

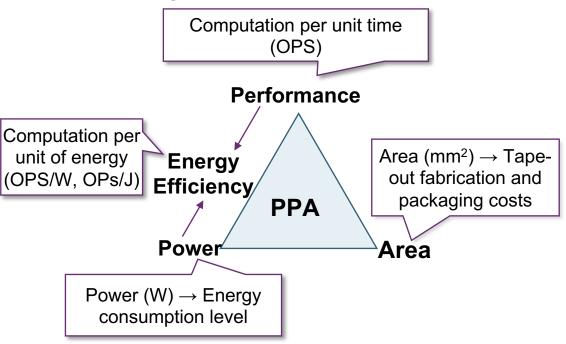
More specialized and coarse-grained ISA or even template-based design that hard code an algorithm.

Specialized

Chip-level Evaluation Metrics

Chip-level metrics/characteristics

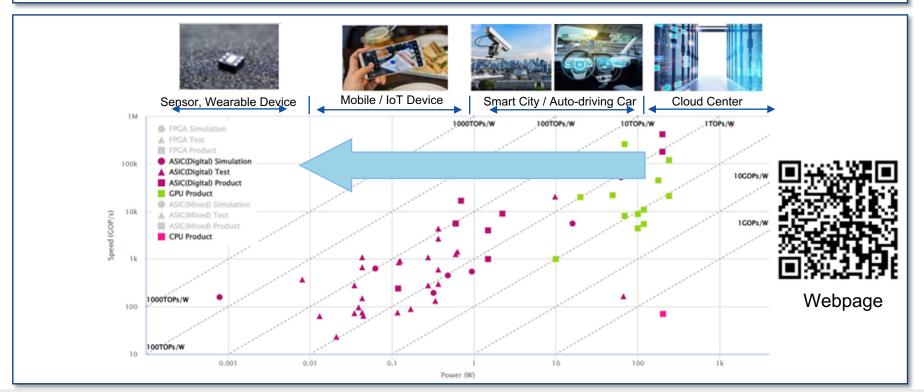
Chip-level evaluation metrics



PPT credit: Prof. Zhenhua Zhu@Tsinghua University

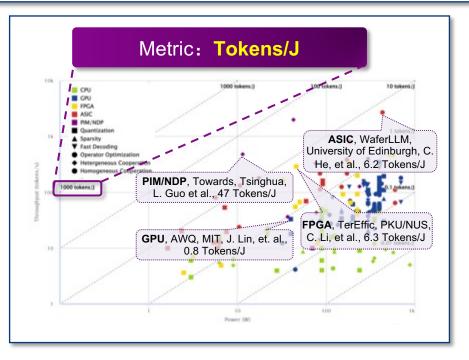
AI1.0 Accelerators in Different Scenarios

From cloud center to tiny edge device



From AI 1.0 to AI 2.0: Energy Efficiency Metric

Hardware energy efficiency → Inference-system energy efficiency TOPs/J → Tokens/J

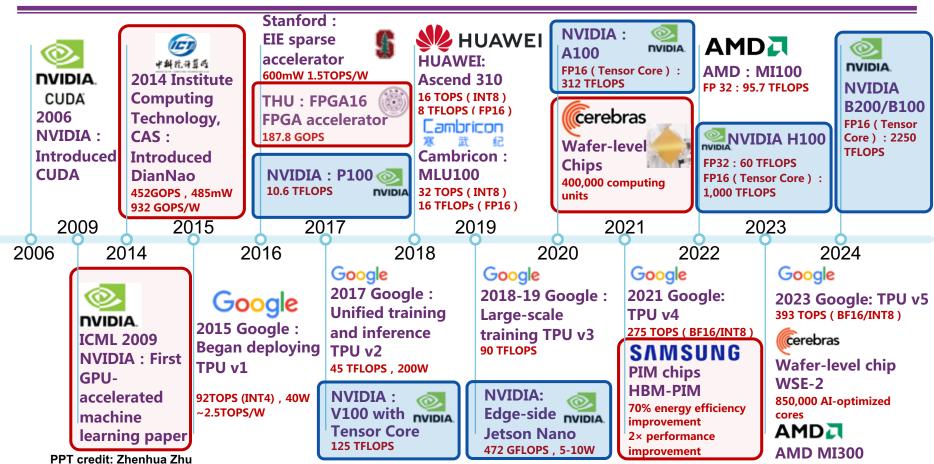






Paper

Basic Knowledge: Development of Chips



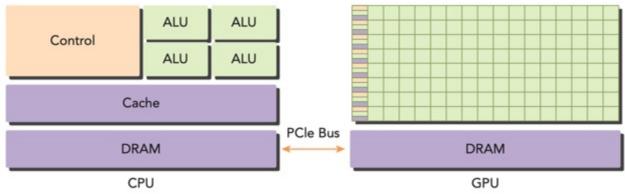
2 Preliminary

- 2.1: LLMs' Model & Algorithm
- 2.2: Basics of Software and Hardware
- 2.3: NVIDIA GPU's Software and Hardware
- 2.4: Deployment

Hardware System & GPU Device

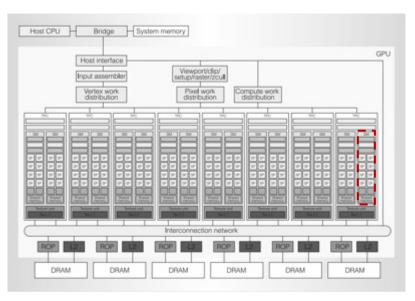
NVIDIA GPU

- GPU device:
 - GPUs use a Single Instruction Multiple Threads (SIMT) architecture.
 - Compared with CPUs, GPUs are better suited for programs featuring simple control logic and large-scale parallel computation.
- Hardware system: A GPU is not an independent computing platform, but rather a co-processor to the CPU.

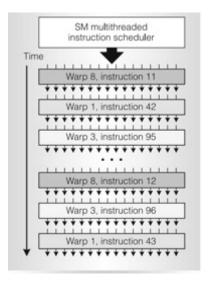


GPU Microarchitecture

- Design Concept of GPUs
 - Single Instruction Multiple Threads (SIMT): A single instruction is executed in parallel by multiple threads.



SM (Streaming Multiprocessor) serves as the basic hardware unit for parallel instructions.



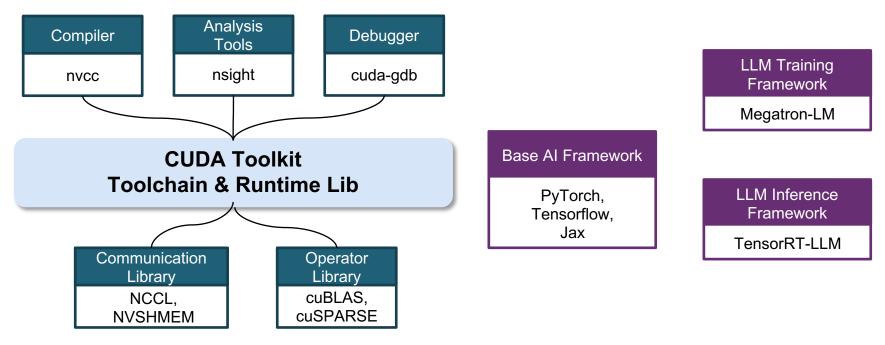
SIMT

GPU Architecture

GPU Software Stack

NVIDIA GPU Software Stack

- CUDA (Compute Unified Device Architecture) is the NVIDIA's GPU parallel programming platform and programming model, featuring a rich software ecosystem.
- AI & LLM frameworks is built on them.

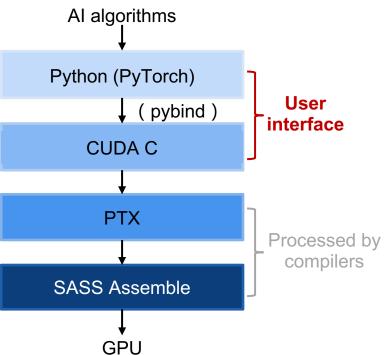


[1] M. Shoeybi, et al. " Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism.", arXiv, 2019.

GPU Programming Model

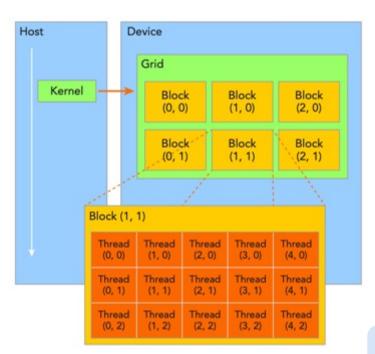
NVIDIA GPU Software Stack

 CUDA allows writing GPU code in high-level languages such as C/C++, reducing programming complexity.

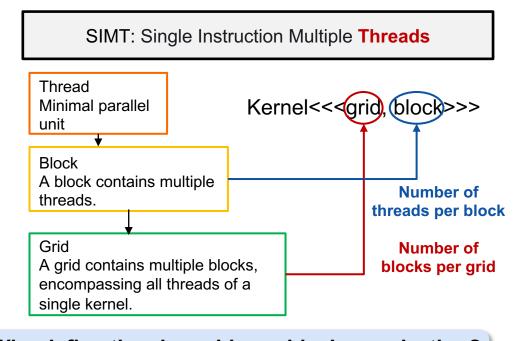


GPU Programming Model

Programming Organization: How to Parallelize?



Thread



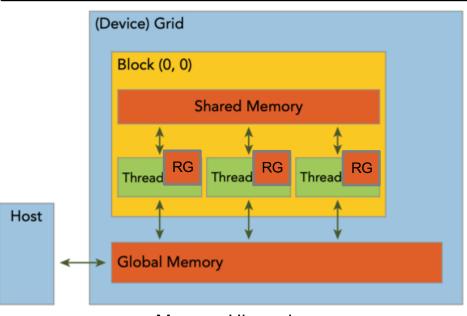
Why define the above hierarchical organization? It corresponds to the GPU memory hierarchy!

GPU Programming Model

GPU Memory Hierarchy

- Thread
 - Each thread has its own registers (RG).
 - Register contents are not shared between threads.
- Block
 - All threads within a block share Shared Memory. Can cooperate & communicate through it.
- Grid
 - All blocks access data from Global Memory (High Bandwidth Memory, HBM).

SIMT: Single Instruction Multiple Threads

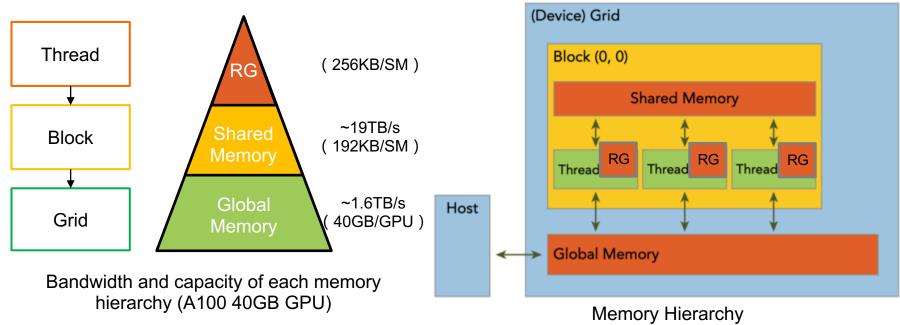


Memory Hierarchy

GPU Programming Model

GPU Memory Hierarchy

 From the perspective of memory access efficiency: how threads are organized significantly impacts kernel performance.



[1] C. John, et al. "Professional CUDA C Programming."

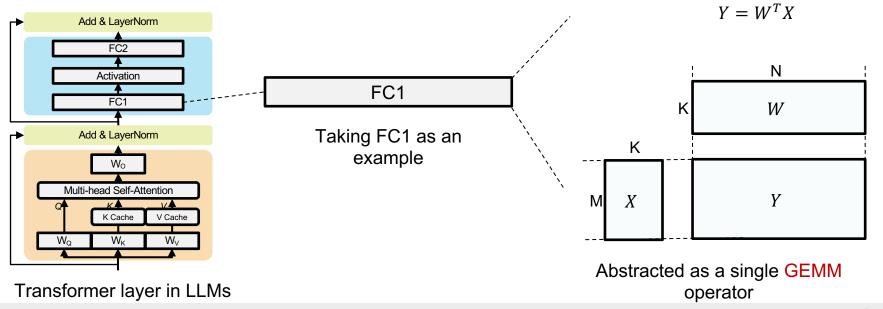
Contents

2 Preliminary

- 2.1: LLMs' Model & Algorithm
- 2.2: Basics of Software and Hardware
- 2.3: NVIDIA GPU's Software and Hardware
- 2.4: Deployment

Algorithm Deployment Process

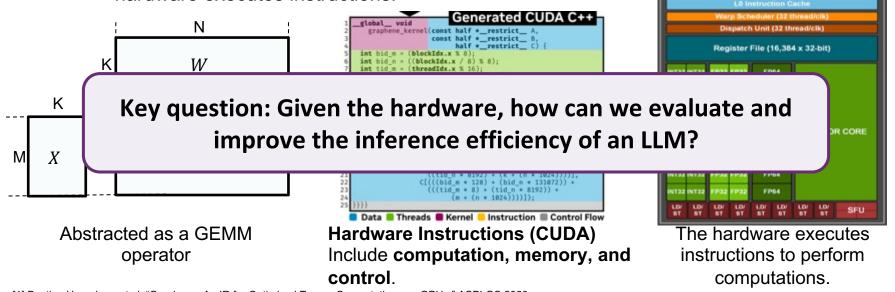
- How Al algorithms are deployed on hardware for actual computation
 - Al Inference can be seen as forwarding data on a **computational graph**, where each node represents a single operator, edge represents dependency.
 - During deployment, operators are translated into hardware instructions. In the runtime, hardware executes instructions.



Algorithm Deployment Process

- How Al algorithms are deployed on hardware for actual computation
 - Al Inference can be seen as forwarding data on a **computational graph**, where each node represents a single operator, edge represents dependency (**software level**).

 During deployment, operators are translated into hardware instructions. In the runtime, hardware executes instructions.



[1] Bastian Hagedorn, et al. "Graphene: An IR for Optimized Tensor Computations on GPUs." ASPLOS 2023.

目录 Contents

1 Background

Model-Level Optimization

2 Preliminary

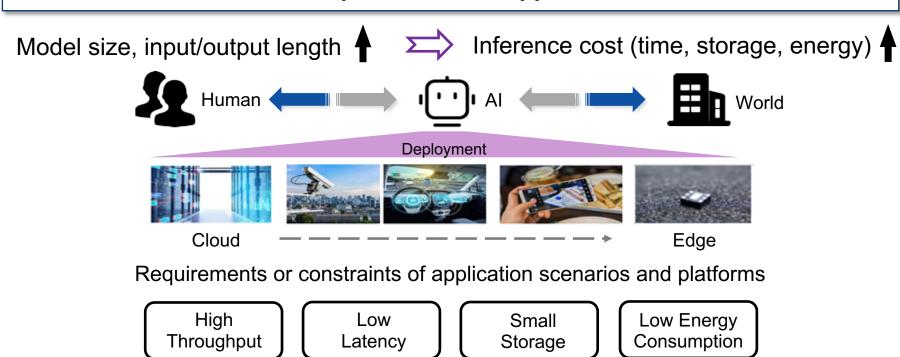
- 6 System-Level Optimization
- 3 Problem Definition & Conceptual Analysis
- Algo-Level Optimization

4 Practical Pipeline

8 Conclusion

Review: Application Challenge

The scale of the model size and input/output pose challenges for efficient inference across platforms and application scenarios.



Review: Application Challenge

The scale of the model size and input/output pose challenges for efficient inference across platforms and application scenarios.







To meet requirements or constraint of application scenarios and platforms, we need to optimize the resource consumption of Al inference.

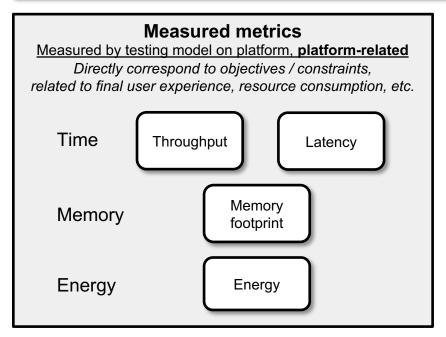


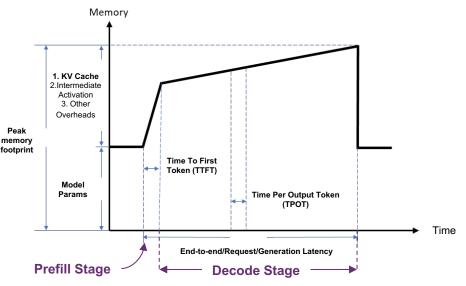
Requirements or constraints of application scenarios and platforms

High

Problem Definition: Objectives & Constraints

Optimization objective or constraint: Usually, latency, memory, energy consumption or throughput will be the ultimate objective or constraint on "efficiency". In the meantime, the intelligence level of AI needs to be retained.

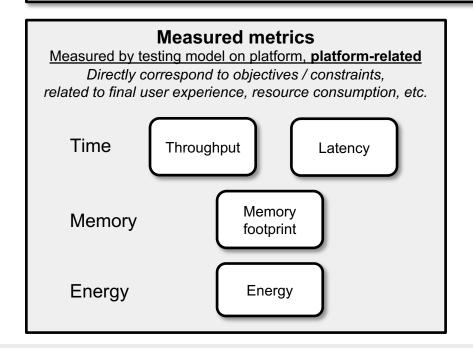


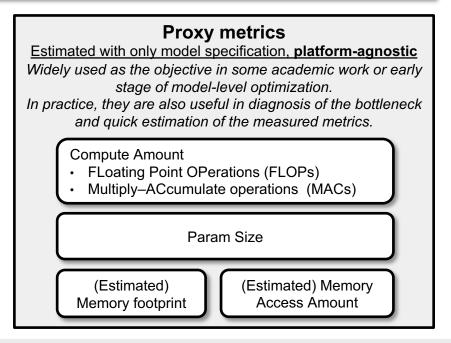


^{*}Note this is only a conceptual illustration. In actual serving framework, the KV cache pool is usually pre-allocated.

Problem Definition: Objectives & Constraints

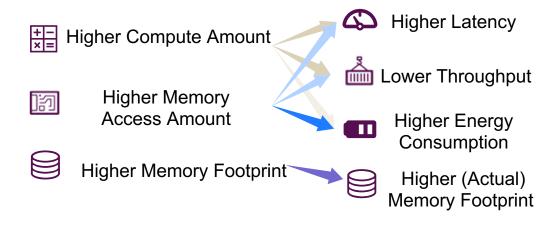
Optimization objective or constraint: Usually, latency, memory, energy consumption or throughput will be the ultimate objective or constraint on "efficiency". In the meantime, the intelligence level of AI needs to be retained.





How Proxy Metrics Relate with Measured Metrics

- We can use three important proxy metrics to analyze the efficiency
 - Compute amount: the amount of operations
 - Memory access amount: the amount of data that read or written between off-chip DRAM and GPU chip
 - Memory footprint: the occupied off-chip DRAM size to store parameters/KV cache/activation



Then, let's bring hardware specification into the picture.

How Proxy Metrics Relate with Measured Metrics

- Intuitively, $\frac{\text{Compute Amount}}{\text{Peak Compute Perf.}}$, $\frac{\text{Memory Access Amount}}{\text{Bandwidth}}$ are two lower bounds of latency.
- But the compute units might not be fully utilized, the bandwidth might not be fully utilized:

$$Compute_utilization = \frac{Achieved\ Compute\ Perf.}{Peak\ Compute\ Perf.} \qquad Bandwidth_utilization = \frac{Achieved\ Bandwidth}{Peak\ Bandwidth}$$

Review our previous estimation example:

| Llama-3.1-8B #Prefill tokens | Compute Amount | Memory Access | VRAM Consume | GPU | TTFT* |
|---------------------------------|-------------------|------------------|-----------------|--------|-----------|
| 1K token | 15.9 TFLOPs | 29.2 GB | 15.3 GB | 1xA100 | 73~170 ms |

1K token 15.9 TFLOPs 29.2 GB 15.3 GB 1xA100 73~170 ms estimate latency for the prefilling stage?

=> It's because prefilling stage is usually more "computation-bounded"

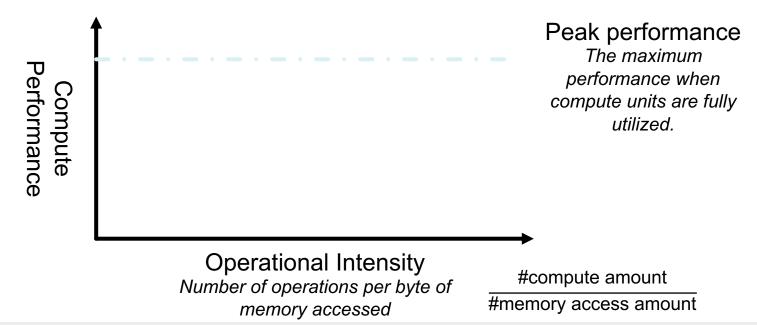
NVIDIA A100 80G's FP16 peak compute performance is 312 TFLOPS.

$$Latency = \frac{\textit{Compute Amount}}{\textit{Peak Compute Perf.} \times \textit{compute_utilization}} = \frac{15.9\,\text{TFLOPs}}{312\,\text{TFLOPS} \times \textit{compute_utilization}} = \frac{51\,\text{ms}}{\textit{compute_utilization}}$$

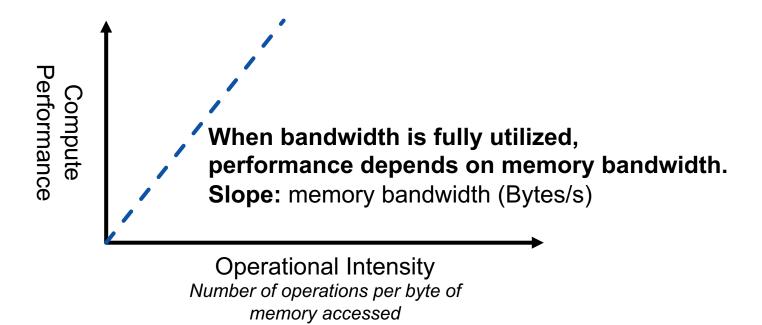
If we assume 30%~70% compute utilization, we can get 73ms~170ms.

Why do we use compute amount instead of memory access amount to

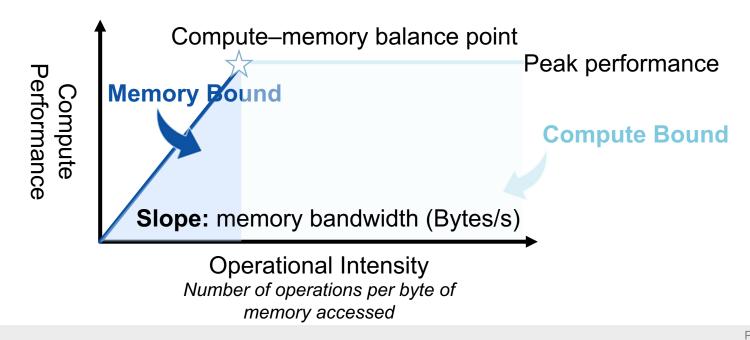
- Proposed by David Patterson in 2009
- Purpose: Uses an algorithm's compute and memory access characteristics (operational intensity) along with the chip's peak performance and memory bandwidth to roughly assess computational bottlenecks and guide subsequent optimization directions.



- Proposed by David Patterson in 2009
- Purpose: Uses an algorithm's compute and memory access characteristics (operational intensity) along with the chip's peak performance and memory bandwidth to roughly assess computational bottlenecks and guide subsequent optimization directions.

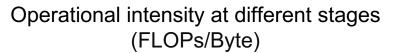


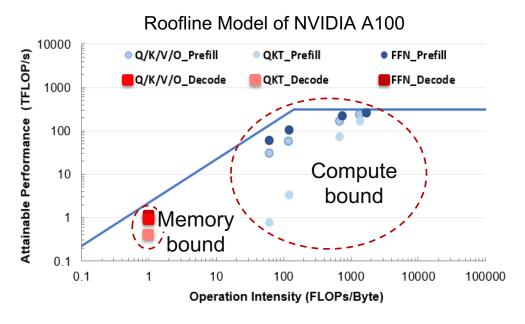
- Proposed by David Patterson in 2009
- Purpose: Uses an algorithm's compute and memory access characteristics (operational intensity) along with the chip's peak performance and memory bandwidth to roughly assess computational bottlenecks and guide subsequent optimization directions.



- Take LLM as an example:
 - Prefill stage: Compute bound; Decode stage: memory bound

| | Prefill | Decode | | |
|----------|------------------------------|-----------------------------|--|--|
| Q/K/V/O | $\frac{ld}{l+d}$ | $\frac{d}{d+1}$ | | |
| QK^{T} | <i>l</i> ≫ 1 | $\frac{l}{l+1}$ ~1 | | |
| AttenV | l l | $\frac{l}{l+1}$ | | |
| FFN | $\frac{ld_{FFN}}{l+d_{FFN}}$ | $\frac{d_{FFN}}{d_{FFN}+1}$ | | |

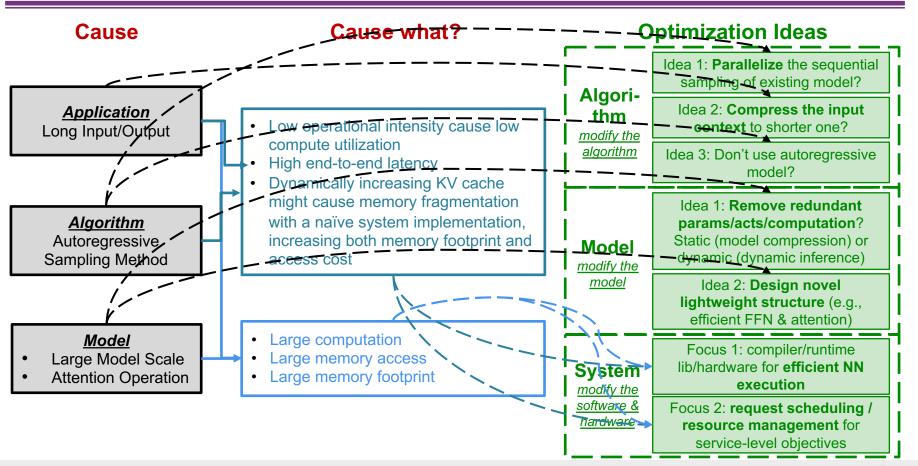


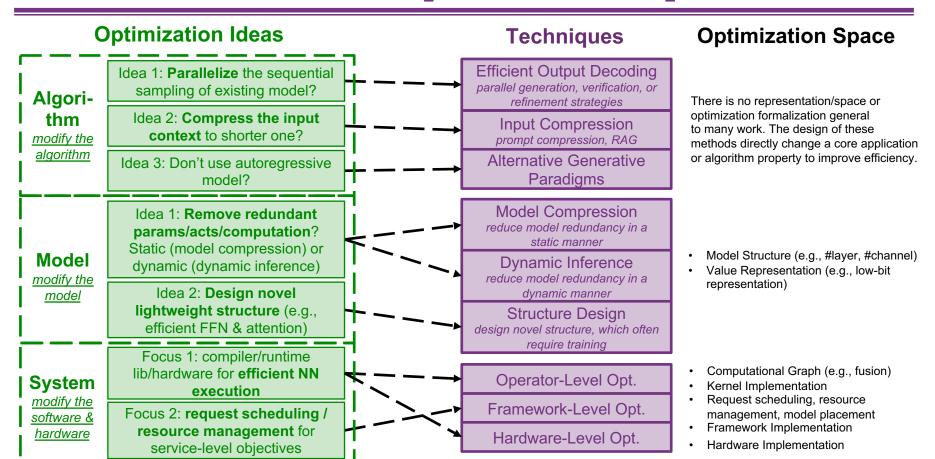


For designing the optimization space, we need to know the "bottleneck" modules or properties of the current algorithm, model, and software that hinders the efficiency of running them on the given hardware.

Root causes of LLM inference inefficiency

- **Application:** The input / output token length can be very long.
- Algorithm: Autoregressive model samples tokens one by one.
- Model: (i) The transformer model has a large number of weights and computations. (ii) Attention modules have quadratic complexity w.r.t. the input token length.





Contents

1 Background

Model-Level Optimization

2 Preliminary

- 6 System-Level Optimization
- 3 Problem Definition & Conceptual Analysis
- Algo-Level Optimization

4 Practical Pipeline

8 Conclusion

Practical Pipeline of Model/System-Level Method Design

Estimation According to Specifications include Application & Model & Hardware

Application Specification

- Objectives, e.g., memory footprint, latency, throughput
- Context lengths
- Batch size

• ..

Model Specification

- · Num hidden layers
- Num key-value heads
- Hidden size
- Intermediate size
- Num attention heads
- Head dim

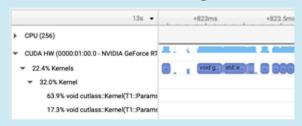
Hardware Specification

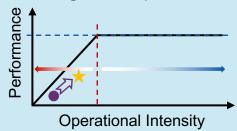
- · Peak Compute Performance
- · Memory Capacity
- Memory Bandwidth

Can estimate theoretically:

- Judge compute or memory bound of each module by roofline model
- Estimate bottleneck module
- Estimate overall objectives

Profile/Diagnosis Tool: Nsight System & Nsight Compute





Model-Level Design

- Analyze compressing which dimension / how to redesigning the module might help with the efficiency most
- Analyze algorithmic redundancy & property, how can we retain/restore performance

System-Level Design (Operator-Level)

- Operation fusion to reduce memory access and kernel launch overhead
- Reimplement some operations

NVIDIA GPU Diagnosis Tool

Nsight System

 A system-level analysis tool used for diagnosing performance bottlenecks from a global perspective and identifying key operators that require optimization.

nsys nvprof -o {output-file} python3 xxx.py

- Taking **LLaMA2-7B** as an example, Nsight Systems can visualize all operators and their calls across all inference stages. It enables researchers to quickly pinpoint efficiency bottlenecks from a global perspective.
 - Input length: 2048 tokens; decoding: 2 tokens



NVIDIA GPU Diagnosis Tool

NVIDIA Nsight Systems

Nsight System

 A system-level analysis tool used for diagnosing performance bottlenecks from a global perspective and identifying key operators that require optimization.

nsys nvprof -o {output-file} python3 xxx.py

- Taking LLaMA2-7B as an example, Nsight Systems can visualize all operators and their calls across all inference stages. It enables researchers to quickly pinpoint efficiency bottlenecks from a global perspective.
 - Nsight Systems analysis reveals that, during the decode stage, the General Matrix-Vector Multiplication (GEMV) operator in the linear layer is the primary performance bottleneck.



NVIDIA GPU Diagnosis Tool

Nsight Compute

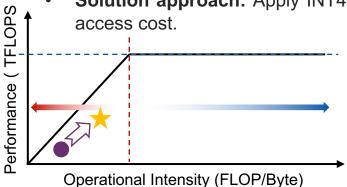


 A kernel-level performance analysis tool focused on deep optimization of individual CUDA kernels, used in conjunction with the Roofline Model for performance analysis.

ncu -set full -o {output-file} python3 xxx.py

Taking LLaMA2-7B as an example, using Nsight Compute to analyze the GEMV operator in the linear layer shows that the operator is severely memory-bound.

Solution approach: Apply INT4 quantization to the weights to reduce weight's memory



| Input channel | Output channel | FP16 (us) | INT4 (us) | |
|------------------|-------------------|-----------|-----------|--|
| 4096 | 11008 | 159.3 | 52.0 | |
| 11008 | 4096 | 45.6 | 37.6 | |
| 4096 | 4096 | 43.5 | 23.0 | |

RTX 3090 GPU

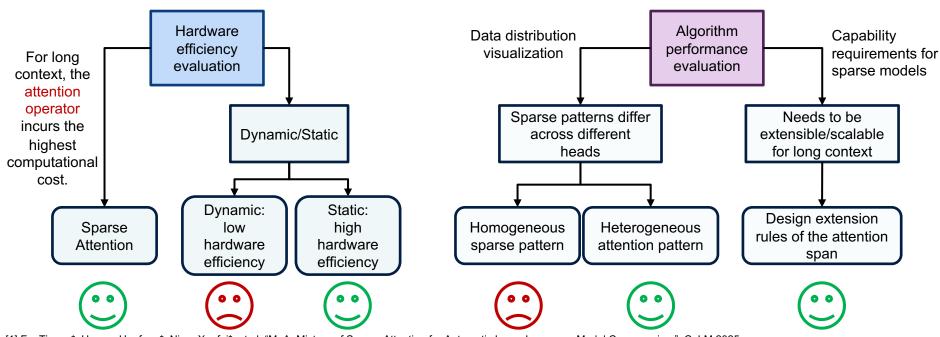
NVIDIA GPU Diagnosis Tool

- Nsight Compute
 - Taking LLaMA2-7B as an example
 - Implementation: To achieve acceleration, the dequantization operator and the GEMV computation operator need to be fused.

```
#pragma unroll
for (int ic_0 = 0; ic_0 < 4; ic_0++){
 // iterate over different uint32_t packed_weights in this loop
  uint32_t current_packed_weight = packed_weights[ic_0];
  half packed_inputs[PACK_FACTOR];
  // each thread load 8 inputs, starting index is packed_group_idx * 128 * 8 (because each iter loads 128*8)
  if (inputs_ptr_delta + ic_0 < IC / PACK_FACTOR) {</pre>
   *((float4*)packed_inputs) = *(inputs_ptr + ic_0);
    #pragma unroll
    for (int ic_1 = 0; ic_1 < PACK_FACTOR; ic_1++){</pre>
     // iterate over 8 numbers packed within each uint32_t number
                                                                                                Dequantization
     float current_single_weight_fp = (float)(current_packed_weight & 0xF);
                                                                                                                             Operator
     float dequantized_weight = scaling_factor * (current_single_weight_fp - current_zeros);
                                                                                                                              Fusion
      //if(blockTdv.v == 0 && blockTdv.v == 0 && threadTdv.v == 0 && ic 0 == 0 && ic 1 == 0 && na
                                                                                                            GEMV
     psum += dequantized_weight * __half2float(packed_inputs[ic_1]);
     current_packed_weight = current_packed_weight >> 4;
```

Example: Design Thought of Model Compression Method

- How to design a sparsification method for a given model and scenario
 - Example: MoA^[1]



Contents

1 Background

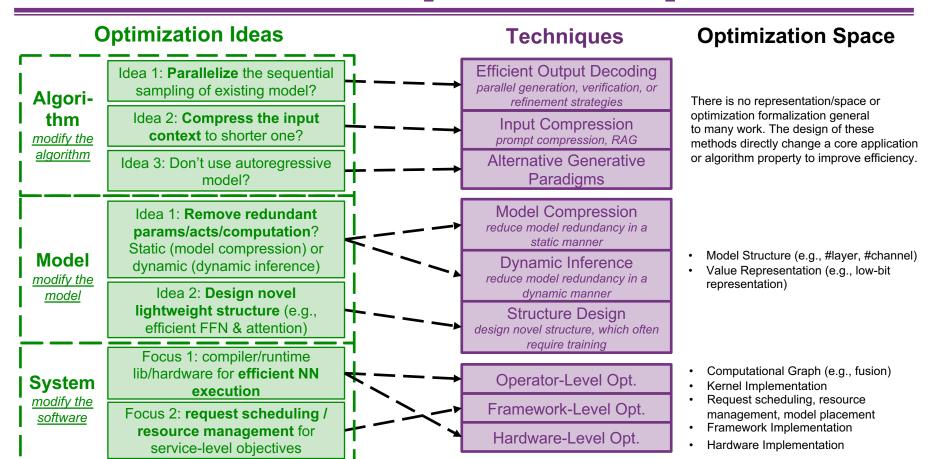
Model-Level Optimization

2 Preliminary

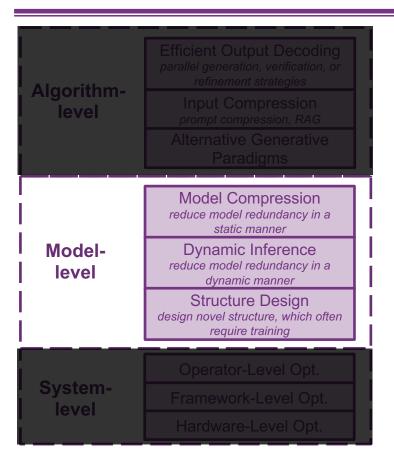
- 6 System-Level Optimization
- 3 Problem Definition & Conceptual Analysis
- Algo-Level Optimization

4 Practical Pipeline

8 Conclusion



Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

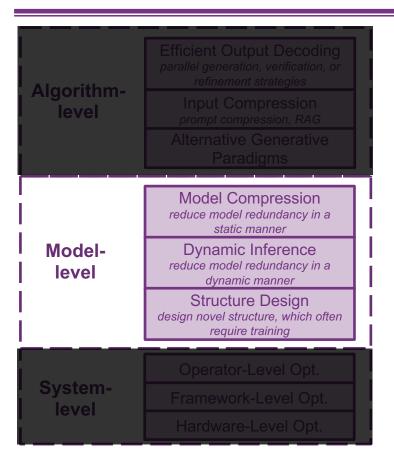
Dynamic Inference

- Module-granularity
- Model-granularity

Structure Design

- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

Dynamic Inference

- Module-granularity
- Model-granularity

Structure Design

- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

Definition

Motivation

Neural network weights and activations are typically represented using high-precision floating-point formats. However, there exists numerical redundancy in neural network computations, and using lower-precision arithmetic does not significantly affect the accuracy of the network [1].

Definition

Quantization: Represent weights & activations with low-bit numbers, thereby storing them or computing them with reduced numerical precision.



IEEE 754 standard 32-bit floating-point data 8-bit fixed-point data

图片来源: MIT EfficientML Course

^[1] Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, Han et al., ICLR 2016.

Two Procedures

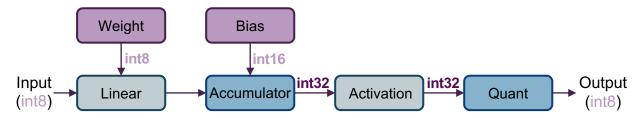
Offline stage: Model quantization procedure

- Convert the FP weights into low-bit-width weights
- Determine the quantization parameters for activations (optional, if using low-precision activation and not using online quant) Quantize floating-point weights (FP32/FP16) into low-bit weights (INT8/INT4).



Online stage: Quantized inference procedure

Low-Precision computation: LP arithmetic -> Requantization



High-Precision computation: Dequant -> HP arithmetic -> (optional) Quant

Model Quantization Procedure

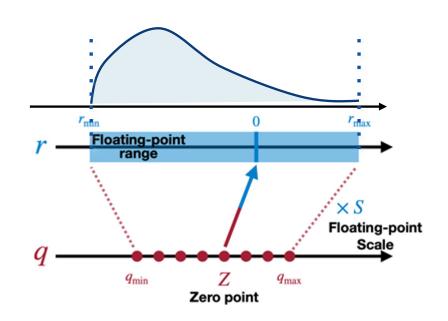
- Example: Uniform Fixed-Point Quantization
- A quantization method needs to answer: "How to convert floating-point numbers into fixed-point number representation?"

$$q = (\lceil r/S \rceil - Z)$$

Fixed-point Floating-point Rounding function

- Quantization parameters: scaling factor (S); zero point (Z)
- How to decide quantization parameters? Take asymmetric quantization ($Z \neq 0$) as an example, if we want the float-point range to cover $[r_{min}, r_{max}]$:

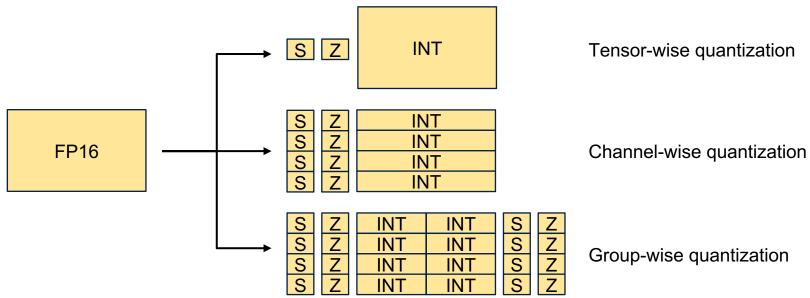
$$S = \frac{r_{max} - r_{min}}{2^{N-1}} \qquad Z = \frac{r_{max} + r_{min}}{2 * 2^{N-1}}$$



^{*} 图片来源: MIT EfficientML Course

^[1] Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, Han et al., ICLR 2016.

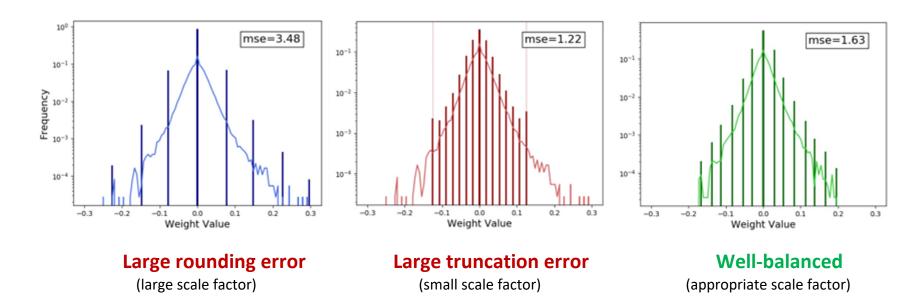
- Quantization Granularity
- For actual storage decrease and speed-up, a group of values needs to share the same quantization parameters (the group size is called quantization granularity, e.g., tensor-wisely, channel-wisely).



^{*} 图片来源: MIT EfficientML Course

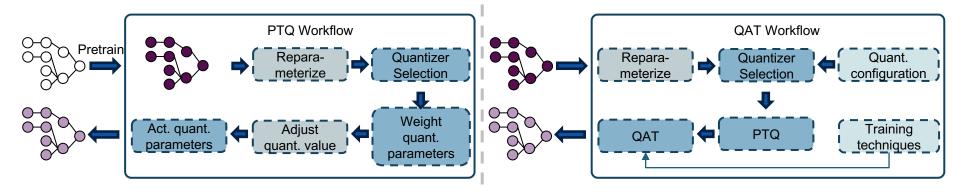
^[1] Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, Han et al., ICLR 2016.

• Core focus of quantization parameter decision: Appropriately balance representational range and precision (i.e., balancing truncation error and rounding error).



^[1] Zhao, Ritchie, et al. "Improving neural network quantization without retraining using outlier channel splitting." International conference on machine learning. PMLR, 2019.

- Two types of quantization process
 - PTQ v.s. QAT



Pros
No weight optimization,
No training cost,
Fast quantization
process

Cons
Lack of recovery,
Difficult to quantize
models to ultra-low
bitwidth

Pros
Quantization with
performance recovery,
Enabling lower bitwidth

Cons
High computational
cost and data
demand

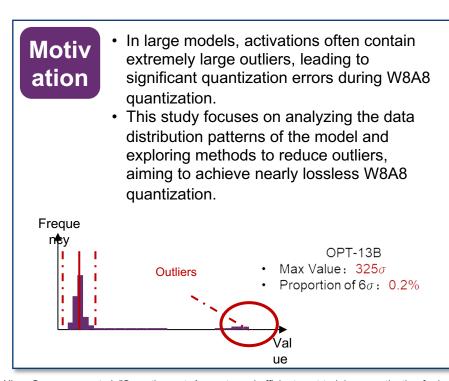
Representative studies

| | Quantized Tensor Type | | | Quantized | Quantized | Quantized Value |
|-------------|-----------------------|------------|--------------|-------------|-----------------|-----------------|
| | Weight | Activation | KV Cache | Format | Criterion | Update |
| GPTQ | | | | Uniform | Statistic-based | $\sqrt{}$ |
| AWQ | $\sqrt{}$ | | | Uniform | Search-based | $\sqrt{}$ |
| SqueezeLLM | $\sqrt{}$ | | | Non-uniform | Statistic-based | |
| GPT3.int8() | $\sqrt{}$ | $\sqrt{}$ | | Uniform | Statistic-based | |
| SmoothQuant | $\sqrt{}$ | $\sqrt{}$ | | Uniform | Statistic-based | $\sqrt{}$ |
| RPTQ | $\sqrt{}$ | $\sqrt{}$ | | Uniform | Statistic-based | |
| OminiQuant | $\sqrt{}$ | $\sqrt{}$ | | Uniform | Search-based | |
| FlexGen | $\sqrt{}$ | | \checkmark | Uniform | Statistic-based | |
| Atom | $\sqrt{}$ | $\sqrt{}$ | \checkmark | Uniform | Statistic-based | |
| KVQuant | | | \checkmark | Non-uniform | Statistic-based | |
| KIVI | | | $\sqrt{}$ | Uniform | Statistic-based | |

- [1] Frantar, Elias, et al. "Gptq: Accurate post-training quantization for generative pre-trained transformers." arXiv preprint arXiv:2210.17323 (2022).
- [2] Lin, Ji, et al. "Awq: Activation-aware weight quantization for on-device Ilm compression and acceleration." Proceedings of machine learning and systems 6 (2024): 87-100.
- [3] Kim, Sehoon, et al. "Squeezellm: Dense-and-sparse quantization." arXiv preprint arXiv:2306.07629 (2023).
- [4] Dettmers, Tim, et al. "Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale." Advances in neural information processing systems 35 (2022): 30318-30332.
- [5] Xiao, Guangxuan, et al. "Smoothquant: Accurate and efficient post-training quantization for large language models." International conference on machine learning. PMLR, 2023.
- [6] Yuan, Zhihang, et al. "Rptq: Reorder-based post-training quantization for large language models." arXiv preprint arXiv:2304.01089 (2023).
- [7] Shao, Wenqi, et al. "OmniQuant: Omnidirectionally Calibrated Quantization for Large Language Models." ICLR. 2024.
- [8] Sheng, Ying, et al. "Flexgen: High-throughput generative inference of large language models with a single gpu." International Conference on Machine Learning. PMLR, 2023.
- [9] Zhao, Yilong, et al. "Atom: Low-bit quantization for efficient and accurate Ilm serving." Proceedings of Machine Learning and Systems 6 (2024): 196-209.
- [10] Hooper, Coleman, et al. "Kvquant: Towards 10 million context length Ilm inference with kv cache quantization." Advances in Neural Information Processing Systems 37 (2024): 1270-1303.
- [11] Liu, Zirui, et al. "KIVI: a tuning-free asymmetric 2bit quantization for KV cache." Proceedings of the 41st International Conference on Machine Learning. 2024.

SmoothQuant

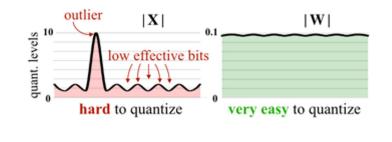
Typical PTQ Method: SmoothQuant^[1]



Insight

Analysis of weight and activation data distribution

 Outliers in large-model activations appear in specific channels. This property can be leveraged to balance the data distribution of weights and activations.



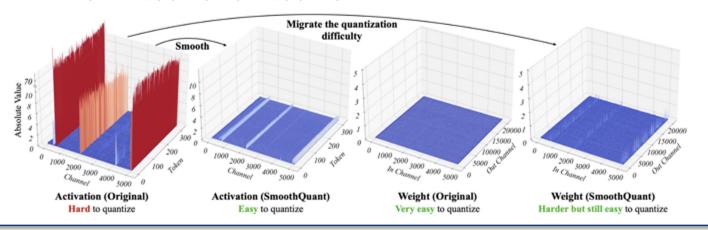
SmoothQuant

Method

- Introducing channel equalization to suppress outliers in activations.
 - Activation channels with large outliers are divided by an equalization factor s_j to reduce their magnitude.

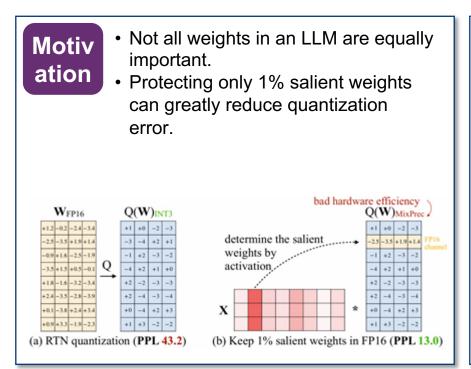
The corresponding weight channels are multiplied by the same factor s_j to appropriately increase the weight values.

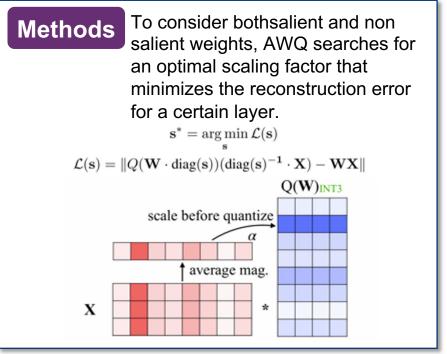
$$\mathbf{Y} = (\mathbf{X} \operatorname{diag}(\mathbf{s})^{-1}) \cdot (\operatorname{diag}(\mathbf{s})\mathbf{W}) = \hat{\mathbf{X}}\hat{\mathbf{W}} \qquad \mathbf{s}_j = \max(|\mathbf{X}_j|)^{\alpha}/\max(|\mathbf{W}_j|)^{1-\alpha}$$



AWQ

Typical PTQ Method: AWQ^[1]



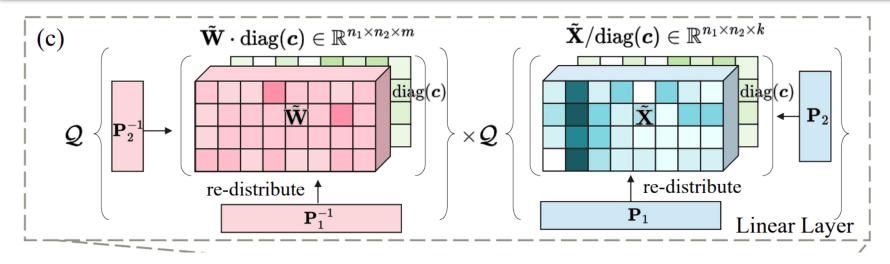


[1] Lin, Ji, et al. "Awq: Activation-aware weight quantization for on-device Ilm compression and acceleration." Proceedings of machine learning and systems 6 (2024): 87-100.

FlatQuant: Flatness Matters for LLM Ouantization

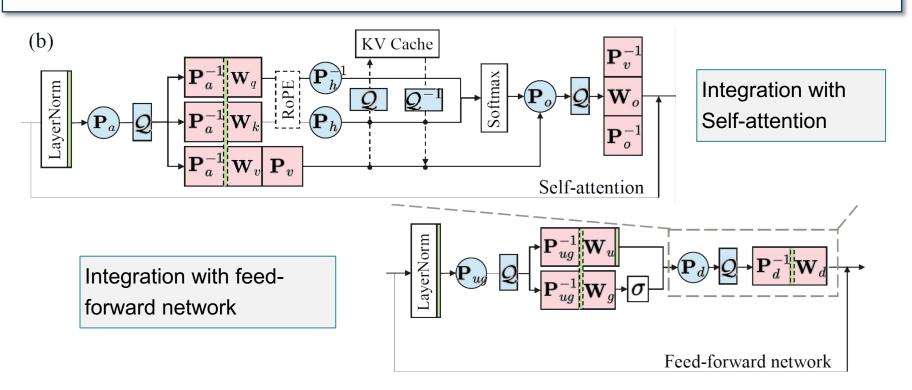
- Motivation: Affine transformations are more powerful to suppress outliers
- Methodology: Learning affine transformations for each linear layer
 - Reducing transformation overhead: Kronecker product & kernel fusion

$$\mathbf{P}^* = \arg\min_{\mathbf{P}} \|\mathbf{Y} - \mathcal{Q}(\mathbf{X}\mathbf{P})\mathcal{Q}(\mathbf{P}^{-1}\mathbf{W}^\top)\|_F^2, \qquad \mathbf{P} = \mathbf{P}_1 \otimes \mathbf{P}_2,$$



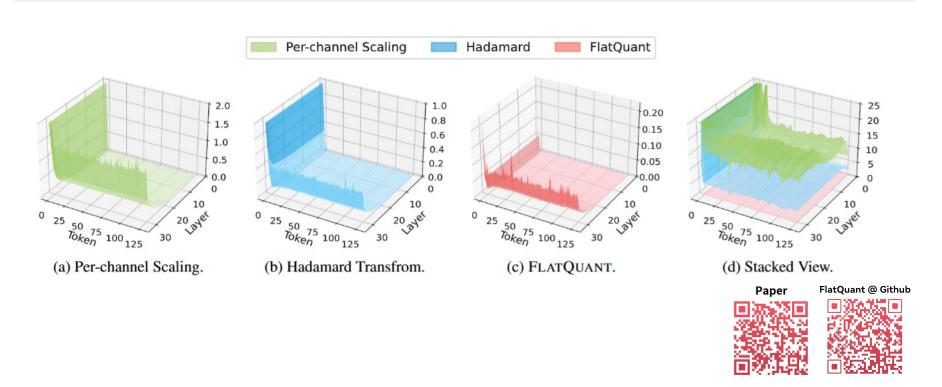
FlatQuant: Flatness Matters for LLM Ouantization

How to Integrate FlatQuant with the Transformer architecture?



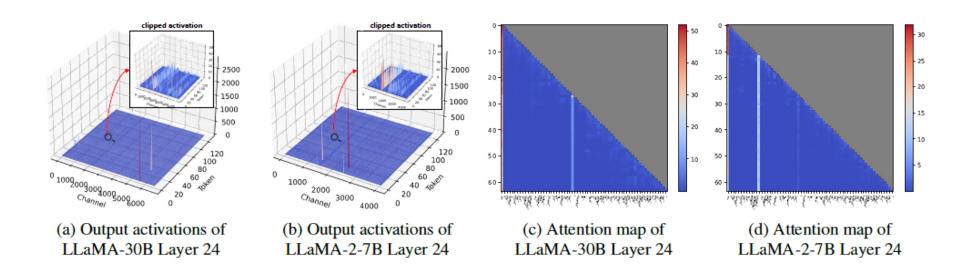
FlatQuant: Flatness Matters for LLM Ouantization

The mean square error of quantization along channels & tokens can be effectively reduced



IntactKV: Keeping Pivot Tokens Intact

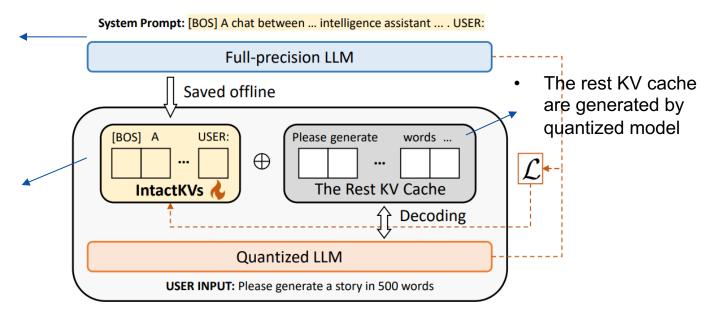
Observation: the pivot tokens exhibit attention sinks together with massive activation outliers.



IntactKV: Keeping Pivot Tokens Intact

Avoid the quantization error accumulated on pivot tokens that are critical to the performance.

- The system prompt contains most pivot tokens with massive outliers
- System KV cache are generated by the BF16 model
- They can be further trained like LLM parameters



Evaluating Quantized Models

[1] Li, S., et al. "Evaluating Quantized Large Language Models." ICML 2024.

- Evaluation Dimensions
 - Effects of quantization on 5 major categories of tasks
 - Effects of quantization on 11 model families
 - Effects of quantizing 3 tensor types on model performance
 - Application scope of SOTA quantization methods

| Section | Tasks & Ability | Benchmark | Size | Model Family |
|---------|-----------------------|-------------------------------------|-------|---------------------------------------|
| l I | Language Modeling | CHID (Zheng et al., 2019) | 2002 | |
| | Language Prodering | Winogrande (Sakaguchi et al., 2021) | 1267 | ĺ |
| | Understanding | RACE (Lai et al., 2017) | 3489 | OPT (125M-66B), LLaMA2 (7B-70B), |
| Sec. 3 | Cincinaining | LAMBADA (Paperno et al., 2016) | 5153 | Falcon (7B-180B), Bloomz (560M-176B) |
| | Reasoning | SIQA (Sap et al., 2019) | 1950 | Mistral(7B, 8×7B) |
| | Reasoning | PIQA (Bisk et al., 2020) | 1876 | ĺ |
| | In-Context Learning | MMLU (Hendrycks et al., 2021b) | 14079 | |
| | III-Context Examing | CEval (Huang et al., 2023) | 13948 | ĺ |
| Sec. 4 | Multi-Step Reasoning | GSM8K (Cobbe et al., 2021) | 1319 | ĺ |
| | Muiti-Step Reasoning | StrategyQA (Geva et al., 2021) | 2290 | ĺ |
| | Instruction-Following | Hellaswag (Zellers et al., 2019) | 10003 | LLaMA2 (7B-70B), Falcon (7B-180B), |
| | instruction-rottowing | ARC (Clark et al., 2018) | 7787 | ChatGLM (6B), Mistral (7B, 8×7B) |
| ì | Self-Calibration | MMLU (Hendrycks et al., 2021b) | 14079 | Gemma (2B, 7B), Mamba (2.8B) |
| | Ethics | ETHICS (Hendrycks et al., 2021a) | 15160 | I |
| App. D | Hallucination | TruthfulQA (Lin et al., 2021) | 817 | ĺ |
| Ĭ | Robustness | AdvGLUE (Wang et al., 2021) | 738 | ĺ |
| Sec. 5 | Dialogue | MT-bench (Zheng et al., 2023a) | 80 | (+ StableLM-3B) |
| Sec. 6 | Long-Context | Longeval (Li et al., 2023) | 3000 | Vicuna (7B, 13B), LongChat (7B, 13B), |
| 366.0 | Zong-Context | Multi-Doc QA (Liu et al., 2023a) | 700 | ChatGLM (6B), Mistral (7B, 8×7B) |



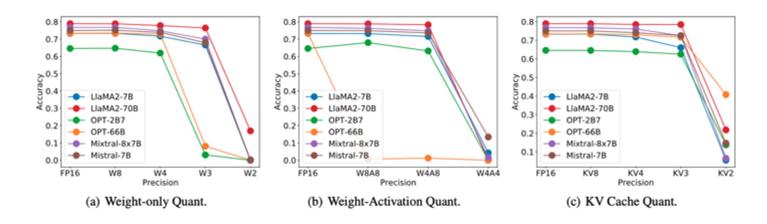
Paper link



Open source

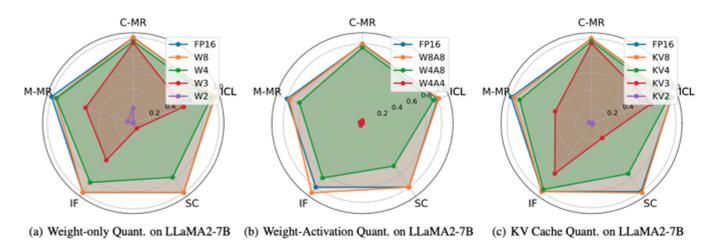
Evaluating Quantized Models

- Effects of Quantization on Tensor Types
 - The larger the model size, the higher the tolerance for Weight and KVcache Quantization, and the lower the tolerance for Activation Quantization.
 - The larger the model size, the fewer outliers in the Weight and KV Cache tensors, and the more outliers in the Activation tensors.



Evaluating Quantized Models

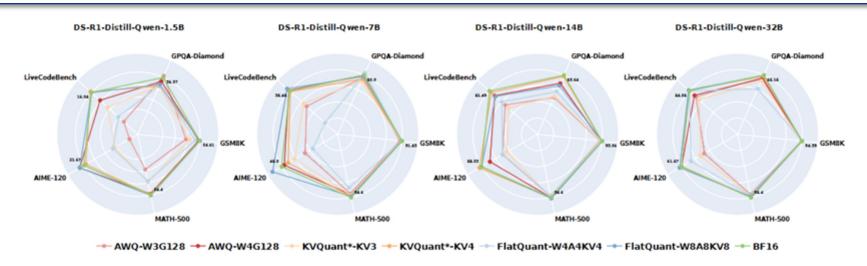
- Effects of Quantization on Emergent Abilities
 - The tolerance to quantization varies across the four abilities, listed in descending order of tolerance: In-context Learning ~ Instruction Following > Multi-Step Reasoning ~ Selfcalibration.



[1] Li, Shiyao, Ning, Xuefei, et al. "Evaluating Quantized Large Language Models." ICML2024.

Evaluating Quantized Reasoning Models

- Reasoning LLMs: Qwen 1.5B 32B distilled from DeepSeek-R1
- Hard tasks (e.g., AIME) suffer more than easier ones (e.g., GSM8K)
 - W8A8 and W4A16 is safe to use (<1% acc drop)
 - W4A4 and KV4 can be still risky in practice

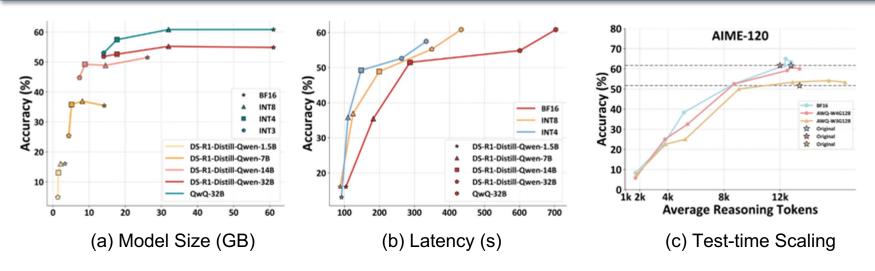


[1] Liu, R., Sun Y., Zhang M., Bai H., et al. "Quantization Hurts Reasoning? An Empirical Study on Quantized Reasoning Models". COLM 2025.

Evaluating Quantized Reasoning LLMs

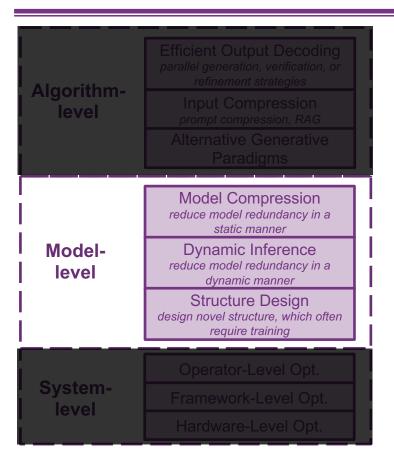
The scaling effect of quantized reasoning LLMs

- (a) & (b): Large quantized LLMs are preferred to small BF16 LLMs w.r.t. size and latency
- **(c) Test-time scaling**: higher accuracy with more reasoning tokens, but at a slower rate when compared to BF16 models



[1] Liu, R., Sun Y., Zhang M., Bai H., et al. "Quantization Hurts Reasoning? An Empirical Study on Quantized Reasoning Models". COLM 2025.

Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

Dynamic Inference

- Module-granularity
- Model-granularity

Structure Design

- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

Model-level: Sparse Attention

Sparse Attention

- Omit certain attention calculations
 - to enhance computational efficiency: saving computation on S and O

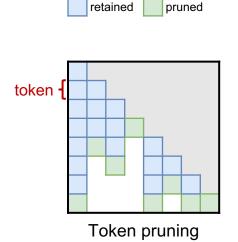
$$S = QK^T$$
, $A = softmax(S + M)$, $O = AV$

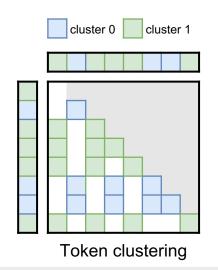
- Static vs Dynamic Mask
 - Static Mask: the attention mask is predefined and remains fixed.
 - Dynamic Mask: the attention mask is determined online based on the input.

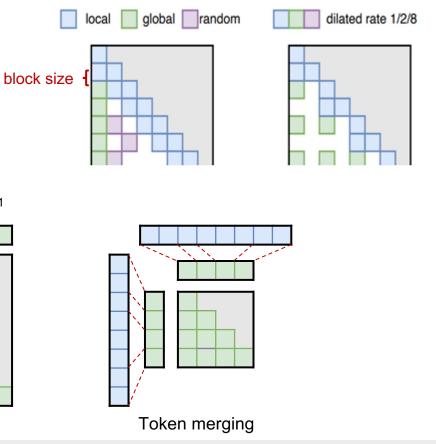
Model-level: Sparse Attention

Sparse Attention

- Sparse Pattern: local, global, random, dilated
 - Granularity: blockwise
- Token pruning vs clustering vs merging







Sparse Attention

Representative studies

| | Mask G | eneration | Sparse Pattern | | | | | | |
|------------------------|--------------|-----------|----------------|-----------|-----------|-----------|------------|-----------|-----------|
| | Static | Dynamic | Local | Global | Dilated | Random | Clustering | Pruning | Merging |
| Sparse Transformers | V | | $\sqrt{}$ | | $\sqrt{}$ | | | $\sqrt{}$ | |
| StreamingLLM | \checkmark | | $\sqrt{}$ | $\sqrt{}$ | | | | | |
| BigBird | \checkmark | | $\sqrt{}$ | $\sqrt{}$ | | $\sqrt{}$ | | | |
| Spatten | | $\sqrt{}$ | | | | | | $\sqrt{}$ | |
| Reformer | | $\sqrt{}$ | | | | | $\sqrt{}$ | | |
| H2O | | $\sqrt{}$ | $\sqrt{}$ | | | | | $\sqrt{}$ | |
| MoA | | $\sqrt{}$ | $\sqrt{}$ | $\sqrt{}$ | | | | | |
| NSA | | $\sqrt{}$ | $\sqrt{}$ | | | | | $\sqrt{}$ | $\sqrt{}$ |

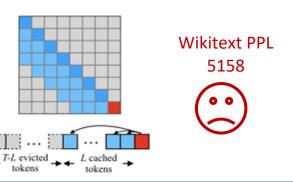
- [1] Child, Rewon, et al. "Generating long sequences with sparse transformers." arXiv preprint arXiv:1904.10509 (2019).
- [2] Xiao, Guangxuan, et al. "Efficient Streaming Language Models with Attention Sinks." The Twelfth International Conference on Learning Representations.
- [3] Zaheer, Manzil, et al. "Big bird: Transformers for longer sequences." Advances in neural information processing systems 33 (2020): 17283-17297.
- [4] Wang, Hanrui, Zhekai Zhang, and Song Han. "Spatten: Efficient sparse attention architecture with cascade token and head pruning." 2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA). IEEE, 2021.
- [5] Kitaev, Nikita, Lukasz Kaiser, and Anselm Levskaya. "Reformer: The Efficient Transformer." International Conference on Learning Representations.
- [6] Zhang, Zhenyu, et al. "H2o: Heavy-hitter oracle for efficient generative inference of large language models." Advances in Neural Information Processing Systems 36 (2023): 34661-34710.
- [7] Fu, Tianyu, et al. "Mixture of Attention Spans: Optimizing LLM Inference Efficiency with Heterogeneous Sliding-Window Lengths." Second Conference on Language Modeling.
- [8] Yuan, Jingyang, et al. "Native sparse attention: Hardware-aligned and natively trainable sparse attention." arXiv preprint arXiv:2502.11089 (2025).

StreamingLLM

Sparse Attention: StreamingLLM

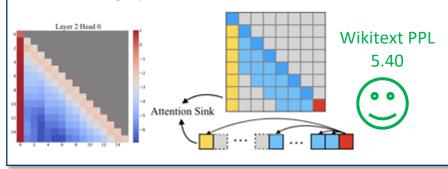
Motiv ation

 Sliding-window attention mechanism discards part of the long-term historical information that is rarely needed, but this often results in severe performance degradation.





- Key finding: LLMs tend to assign attention scores to the initial tokens (attention sink).
- Therefore, StreamingLLM not only retains sliding-window attention but also preserves the attention scores of the initial tokens.



MoA

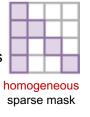
Sparse Attention: MoA

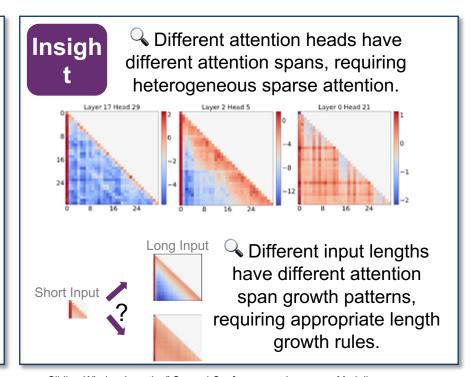


 Existing methods apply a homogeneous sparse mask to each attention head, which fails to capture the diverse attention patterns in LLMs and consequently leads to a significant drop in model performance.



StreamingLLM [2] with fixed-length local attention and global attention on the initial tokens





MoA

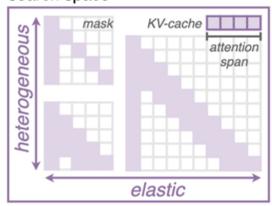


- Different attention heads: require searching for different sparse attention span.
- **Different input lengths:** require searching for suitable attention span growth rules for different attention heads.

Heterogeneous search space

Construct a search space of sparse schemes that includes different sparse attention patterns and their variation rules with sequence length.

search space

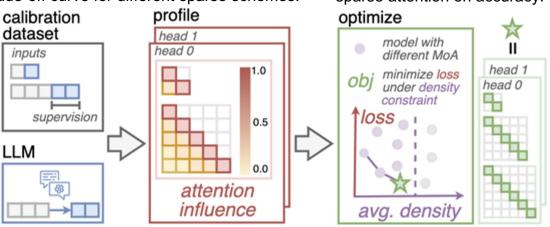


Analyze attention value importance

Based on gradients, analyze the impact of different attention values on the prediction results, and obtain the accuracy—sparsity trade-off curve for different sparse schemes.

portance Optimization

Formulate an optimization problem to select sparse patterns under sparsity constraints, minimizing the impact of sparse attention on accuracy.

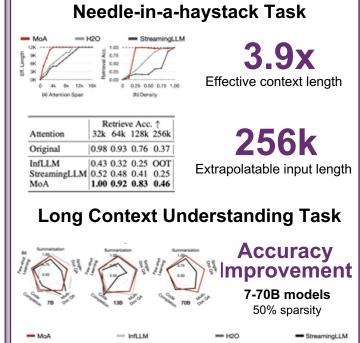


MoA



Increase the inference throughput by about 7×. Expand the effective context length by 3.9×.

Accuracy-Throughput MoA === H2O === StreamingLLM 1.00 0.75 0.50 0.25 0.00 350 1050 700 1400 Throughput (token/s) Vicuna-7B, 8K input length Code



Efficiency

1.7x - 1.9x

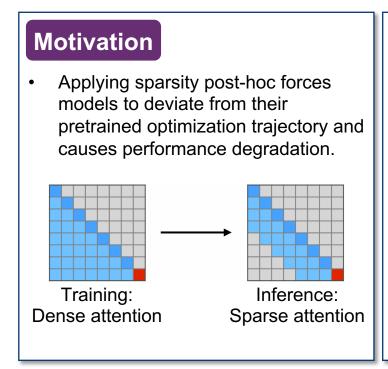
Throughput Improvement

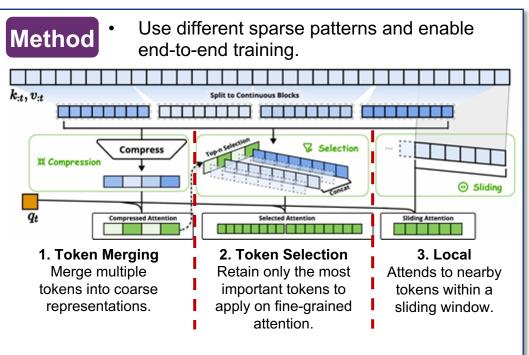
Compared with VLLM on 7B and 13B LLMs using 50% attention sparsity on A100-80GB GPUs

| Model | Framework | Attention | Batch | 4k Throughput | Batch | 8k Throughput | Batch | 16k Throughpu |
|-------|-------------|--------------------|-------|------------------|-------|------------------|-------|------------------|
| | vLLM | PagedAttention | 30 | 628.8 | 15 | 323.0 | 8 | 145.5 |
| | FlexGen | H2O | 20 | 754.9 | 6 | 296.3 | 1 | 51.7 |
| | HuggingFace | InfLLM | 15 | 62.0 | 10 | 37.5 | 6 | 19.2 |
| | HuggingFace | StreamingLLM | 50 | 945.1 | 25 | 467.3 | 12 | 232.0 |
| 7B | | FlashAttention2 | 30 | 134.6 | 15 | 66.9 | 8 | 32.9 |
| | | +Static KV-Cache | 30 | 496.1 | 15 | 219.5 | 8 | 91.6 |
| | HuggingFace | +Reduced Attention | 30 | 722.5 | 15 | 369.9 | 8 | 178.3 |
| | | +Increased Batch | 50 | 897.7 | 25 | 436.7 | 12 | 206.4 |
| | | +Kernel (=MoA) | 50 | 1099.0 | 25 | 535.7 | 12 | 257.3 |
| | vLLM | PagedAttention | 16 | 314.8 | - 8 | 160.5 | 4 | 71.1 |
| | FlexGen | H2O | 12 | 330.2 | 4 | 138.2 | 1 | 37.4 |
| | HuggingFace | InfLLM | 8 | 30.3 | 5 | 17.63 | 3 | 11.3 |
| | HuggingFace | StreamingLLM | 28 | 478.4 | 14 | 241.2 | 7 | 116.5 |
| 13B | | FlashAttention2 | 16 | 81.3 | 8 | 40.8 | 4 | 19.8 |
| | | +Static KV-Cache | 16 | 264.6 | 8 | 111.3 | 4 | 62.2 |
| | HuggingFace | +Reduced Attention | | 329.6 | 8 | 156.4 | 4 | 87.3 |
| | | +Increased Batch | 28 | 471.5 | 14 | 222.6 | 7 | 108.3 |
| | | +Kernel (=MoA) | 28 | 550.9 | 14 | 267.6 | 7 | 132.3 |

NSA

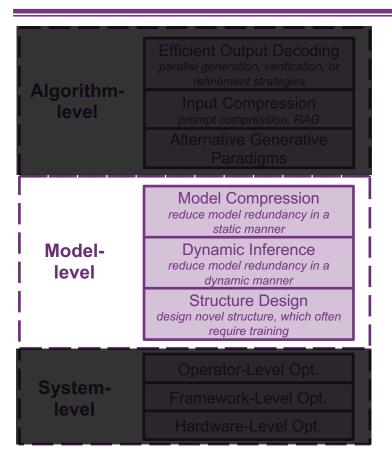
Sparse Attention: NSA





[1] Yuan, Jingyang, et al. "Native sparse attention: Hardware-aligned and natively trainable sparse attention." arXiv preprint arXiv:2502.11089 (2025).

Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

Dynamic Inference

- Module-granularity
- Model-granularity

Structure Design

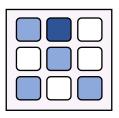
- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

Model-level: Weight Pruning

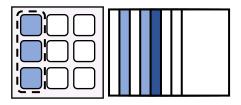
Weight Pruning

Remove less critical weights and structures from models

| | Unstructured Pruning | Structured Pruning |
|-----------------------------|--|--|
| Granularity | Individual weight values | Structural units, e.g., channels, layers, experts |
| Performance loss | Low | High |
| Actual speed-up on hardware | No | Yes |
| Research directions | (1) accelerate pruning process(2) design effective pruning strategies (e.g., pruning metrics, pruning ratios) | (1) decide structured pattern(2) design effective pruning metrics |
| Representative Studies | SparseGPT, Prune and Tune, ISC, BESA | LLM-Pruner , LLaMA-Sheard, ZipLM, LoRAPrune, EEP |



Unstructured Pruning Granularity: Weight



Structured Pruning
Granularity: Channel/Group/Layer

SparseGPT

Idea

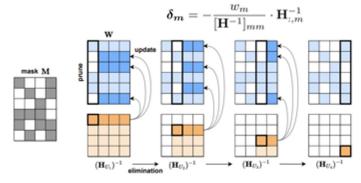
 Removing weights and updating the remaining ones to compensate for the error.

Category

- Type: Unstructured Pruning
- Granularity: Individual weight values

Method

 Incrementally prune weights in each column of the weight, using a sequence of Hessian inverses, and updating the remainder of the weights.



$$(\mathbf{H}_{U_{j+1}})^{-1} = \left(\mathbf{B} - \frac{1}{[\mathbf{B}]_{11}} \cdot \mathbf{B}_{:,1} \mathbf{B}_{1,:}\right)_{2:,2:} \quad \mathbf{B} = (\mathbf{H}_{U_j})^{-1}$$

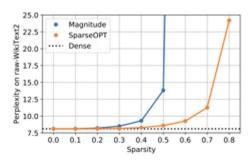
SparseGPT

Results

- **High Sparsity, Low Accuracy Loss**: Prunes OPT-175B to 60% sparsity in one shot with a negligible increase in perplexity
- **High Efficiency**: Process the 175-billion-parameter models in under 4.5 hours, removing more than 100 billion weights.

Table 2. ZeroShot results on several datasets for sparsified variants of OPT-175B.

| Method | Spars. Lamb. PIQA ARC-e ARC-c Story. | Avg. |
|-------------------------------------|---|-------|
| Dense | 0% 75.59 81.07 71.04 43.94 79.82 | 70.29 |
| Magnitude | 50% 00.02 54.73 28.03 25.60 47.10 | 31.10 |
| SparseGPT SparseGPT SparseGPT | 50% 78.47 80.63 70.45 43.94 79.12 4:8 80.30 79.54 68.85 41.30 78.10 2:4 80.92 79.54 68.77 39.25 77.08 | |



LLM-Pruner

Idea

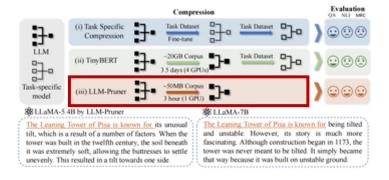
Identify and remove non-critical, coupled structures.

Category

- Type: Structured Pruning
- Granularity: Head, Channel

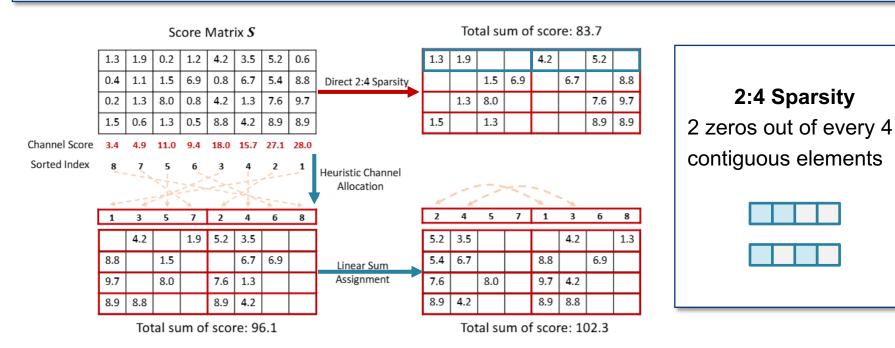
Methods

 LLM-Pruner automatically identifies and removes non-critical, coupled structures based on gradient information, and recovers performance using a LoRA with a small dataset.



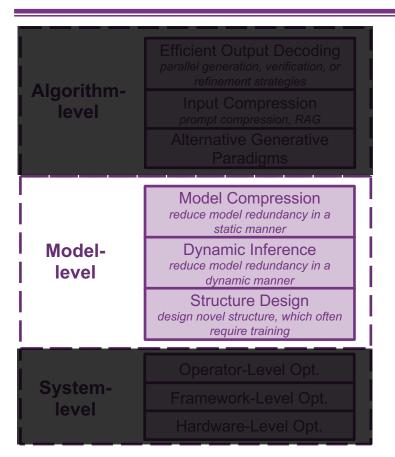
Channel Permutation for Better N:M Sparsity

Find a **permutation strategy** that preserves more important parameters under N:M sparsity



[1] Zhang, Y., Bai, H., et al. "Plug-and-play: An Efficient Post-training Pruning Method for Large Language Models". ICLR 2024.

Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

Dynamic Inference

- Module-granularity
- Model-granularity

Structure Design

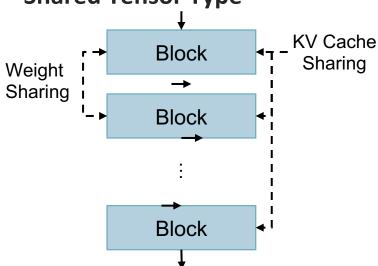
- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

Sharing

Definition

Reuse parameters, states, or intermediate results across different parts of the model.





| | Shared Tensor Type | | | |
|---------------------|--------------------|-----------|--|--|
| | Weight KV Cache | | | |
| Subformer | | | | |
| MobileLLM | $\sqrt{}$ | | | |
| Dynamic layer tying | $\sqrt{}$ | | | |
| LCKV | | $\sqrt{}$ | | |
| CLA | | V | | |

- [1] Reid, Machel, Edison Marrese-Taylor, and Yutaka Matsuo. "Subformer: Exploring weight sharing for parameter efficiency in generative transformers." arXiv preprint arXiv:2101.00234 (2021).
- [2] Liu, Zechun, et al. "Mobilellm: Optimizing sub-billion parameter language models for on-device use cases." Forty-first International Conference on Machine Learning. 2024.
- [3] Hay, Tamir David, and Lior Wolf. "Dynamic Layer Tying for Parameter-Efficient Transformers." The Twelfth International Conference on Learning Representations.
- [4] Wu, Haoyi, and Kewei Tu. "Layer-Condensed KV Cache for Efficient Inference of Large Language Models." Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2024.
- [5] Brandon, William, et al. "Reducing transformer key-value cache size with cross-layer attention." Advances in Neural Information Processing Systems 37 (2024): 86927-86957.

MobileLLM

MobileLLM

Idea

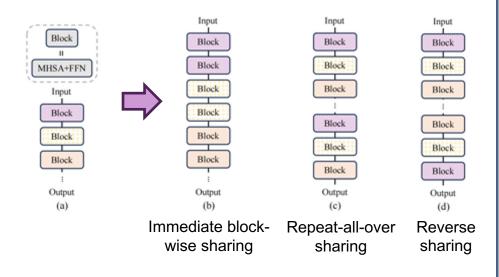
 Weight sharing between two adjacent blocks avoids weight movement, requiring only computing the block twice and incurring minimal latency overhead.

Category

Type: Weight Sharing

Methods

Design three different weight-sharing strategies:



LCKV

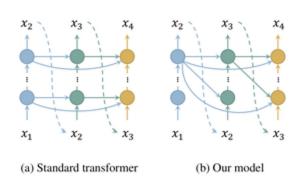
Layer-Condensed KV Cache

Motivation

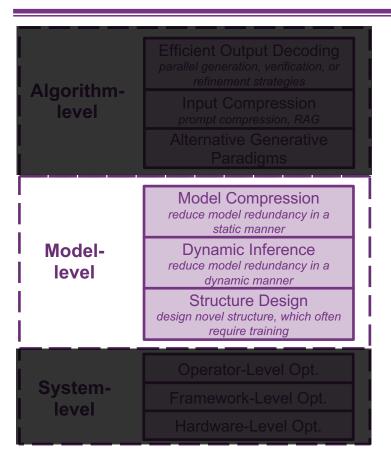
- Existing methods focus on compressing the KV cache sequence length.
- This approach reduces the number of cached layers, not just the sequence length.

Methods

- In LCKV, all layers attend to only the top layer's KVs.
- A few "warmup" layers with standard attention are kept to maintain performance.



Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

Dynamic Inference

- Module-granularity
- Model-granularity

Structure Design

- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

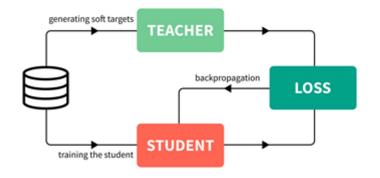
Knowledge Distillation

Motivation

Although the compressed lightweight model achieves better hardware efficiency, its accuracy is lower under conventional training methods. It is necessary to design training approaches to achieve better accuracy recovery.

Definition

Knowledge Distillation: Use a teacher model to guide the training of a student model, enabling the student model to learn the "knowledge" of the teacher model to help improve its accuracy.

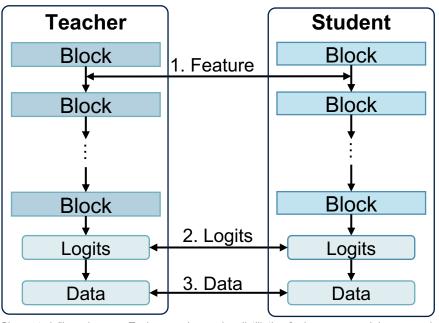




The teacher model (large) helps the student model (small) recover accuracy.

Knowledge Distillation

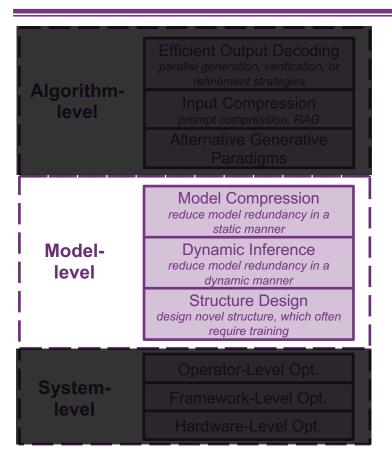
Aligning Objective



| | | Objective | | | | |
|-------------|---------------------|-----------|-----------|--|--|--|
| | Feature Logits Data | | | | | |
| TED | √ | | | | | |
| MiniLLM | | $\sqrt{}$ | | | | |
| GKD | | $\sqrt{}$ | | | | |
| DISCO | | | $\sqrt{}$ | | | |
| MCKD | | | $\sqrt{}$ | | | |
| DeepSeek-R1 | | | $\sqrt{}$ | | | |

- [1] Liang, Chen, et al. "Less is more: Task-aware layer-wise distillation for language model compression." International Conference on Machine Learning. PMLR, 2023.
- [2] Gu, Yuxian, et al. "MiniLLM: Knowledge Distillation of Large Language Models." The Twelfth International Conference on Learning Representations.
- [3] Agarwal, Rishabh, et al. "Gkd: Generalized knowledge distillation for auto-regressive sequence models." CoRR (2023).
- [4] Chen, Zeming, et al. "DISCO: Distilling Counterfactuals with Large Language Models." Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics. 2023.
- [5] Zhao, Jiachen, et al. "Multistage collaborative knowledge distillation from large language models." arXiv preprint arXiv:2311.08640 (2023).
- [6] Guo, Daya, et al. "Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning." arXiv preprint arXiv:2501.12948 (2025).

Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

Dynamic Inference

- Module-dimension
- Model-dimension

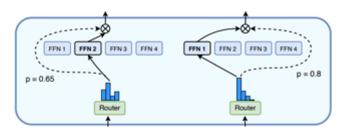
Structure Design

- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

Model-level: Dynamic Inference

Motivation & Definition

- During neural network inference, not all data needs to go through the same computation flow. The core design idea of dynamic inference algorithms is to determine the required computations based on the runtime input data.
 - What is the data granularity for dynamic inference (e.g., query-level, token-level, etc.)
 - Which dimensions are dynamically adjusted (e.g., layer, model)
 - How to dynamically adjust the corresponding dimensions based on input data (e.g., training a router)



```
(a) question: Compute 9999² - 9998×1000.

LLM: Okey, let's think step by step... 9999² is hard, rewrite it...

SLM: Okey, let us think step by step... 9999² is 999801...

identical ✓ neutral ✓ divergent X
```

MoE: Token-level dynamic module routing

R2R: Token-level dynamic model routing

Mixture-of-Depths

Mixture-of-Depths

Idea

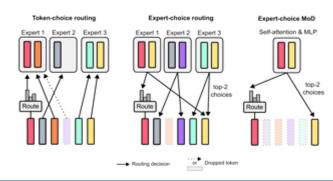
 In language modeling, not all tokens and sequences require the same time or effort to accurately make a prediction.

Category

- Data Granularity: token-level
- Dimension: layer
- Method: training a router

Methods

 At specific layers, a learned router selects the top-k most important tokens to be processed by the self-attention and MLP blocks, while other tokens bypass these computations through a simple residual connection.



RouteLLM

RouteLLM

Motivation

 Direct simpler queries to smaller models and more complex ones to larger models to balance response quality with cost efficiency.

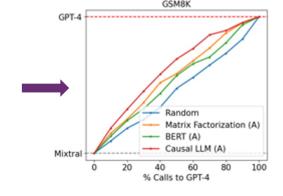
Category

- Granularity:query-level
- **Dimension**: model
- Method: similarity-based retrieval / training a router

Methods

 RouteLLM trains a router model on human preference data to intelligently direct queries to either a strong, expensive LLM or a weak, cheap one.

 routers outperform random baselines



Recent Work: SLM-LLM Mix Inference (R2R)

Use small language model (SLM) and LLM for different reasoning steps

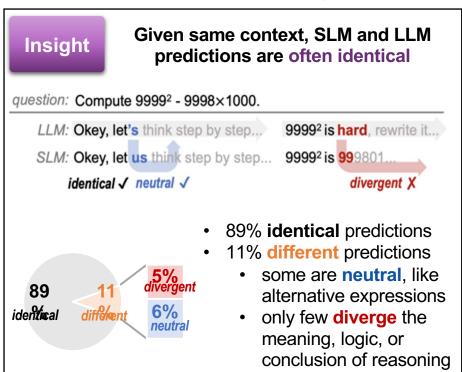
Motivation

Fast but weak SLM slow but strong LLM

Tested results on AIME'24-25

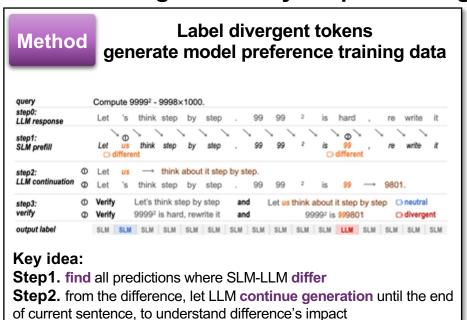
| Туре | Model | Accuracy | Latency (s / question) |
|------|---------|--------------|---------------------------|
| SLM | R1-1.5B | ⊗ 9% | © 199 |
| LLM | R1-32B | 9 45% | 8 498 |

We should selectively use SLM and LLM for different generation steps, constructing a **fast and strong** mix-inference method

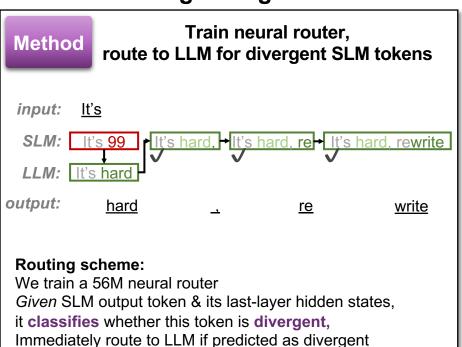


Recent Work: SLM-LLM Mix Inference (R2R)

Label divergent token, then train a neural token-router, utilizing LLMs only for path-divergent tokens during SLM generation



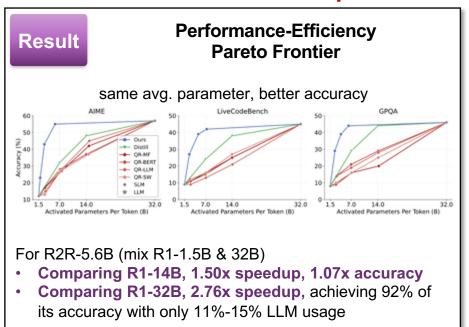
Step3. ask another LLM to verify if difference causes divergence

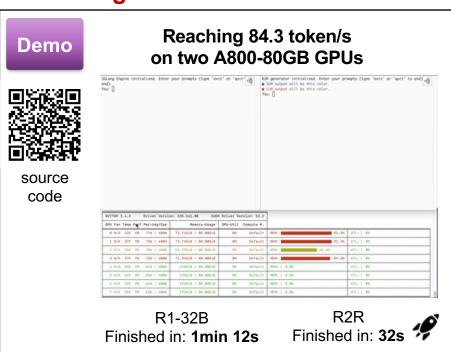


[1] Tianyu, Fu, et al. "Efficiently Navigating Divergent Reasoning Paths with Small-Large Model Token Routing" Submitted to NeurIPS'25. [Under Review]

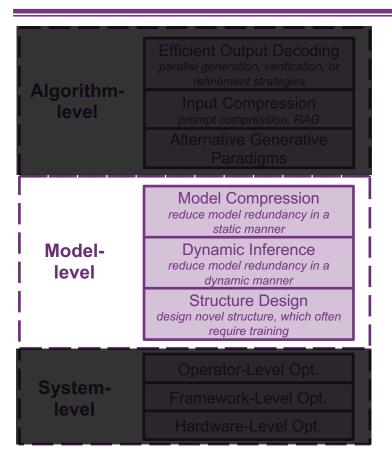
Experimental Results

Mixing R1-1.5B & 32B, uing **R2R** with **5.6B** avg. activated param. per token achieve **performance exceeding R1-14B**





Menu of Techniques



Model Compression

- Quantization
- Sparse Attention
- Weight Pruning
- Sharing
- Knowledge Distillation

Dynamic Inference

- Module-granularity
- Model-granularity

Structure Design

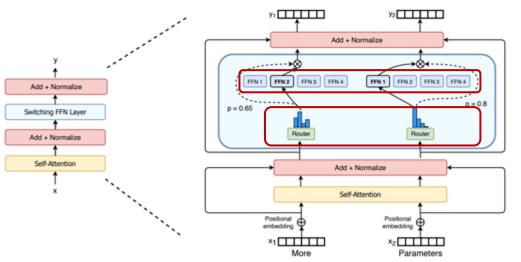
- Mixture-of-Experts (Efficient FFN)
- Efficient Attention

Model-level: Efficient Structure Design

Efficient structure design

Efficient FFN design: Mixture-of-Expert (MoE, one common architecture)

follow the dynamic inference idea)



- Expert: module to process different inputs
 - Shared expert: an expert module that is always active for every input
 - Expert granularity: the FFN intermediate hidden dimension
- Router: direct the input to the appropriate expert networks
 - Basically, each token is routed to a fixed number of experts based on scores produced by the router.
 - In training or multi-batch inference scenarios, load balancing among experts also needs to be considered.^[1]

[1] Mu, Siyuan, and Sen Lin. "A comprehensive survey of mixture-of-experts: Algorithms, theory, and applications." arXiv preprint arXiv:2503.07137 (2025).

MoE

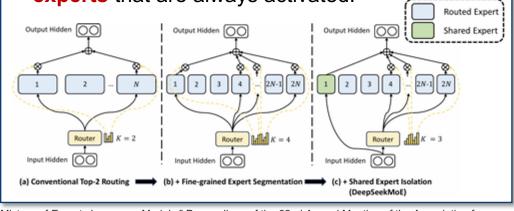
DeepSeekMoE

Motivation

- The designated expert will intend to assemble different types of knowledge in its parameters, which are hard to utilize simultaneously.
- Multiple experts may converge in acquiring shared knowledge in their respective parameters, leading to redundancy in expert parameters.

Methods

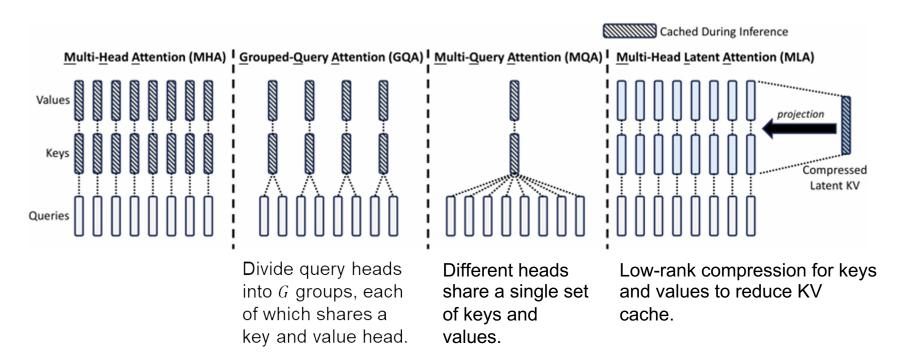
- Segment the experts into a finer grain by splitting the FFN intermediate hidden dimension.
- Isolate certain experts to serve as shared experts that are always activated.



^[1] Dai, Damai, et al. "DeepSeekMoE: Towards Ultimate Expert Specialization in Mixture-of-Experts Language Models." Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2024.

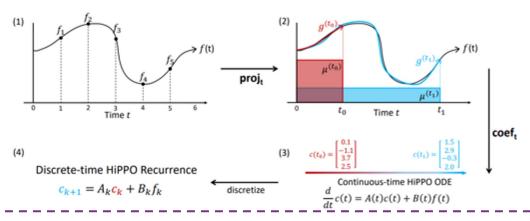
Efficient Attention mechanism

- Efficient structure design
 - Efficient Attention mechanism



Model-level: Efficient Structure Design

- Efficient structure design
 - Non-Transformer architecture design
 - State Space Model (SSM)
 - core idea: compress token information into hidden state

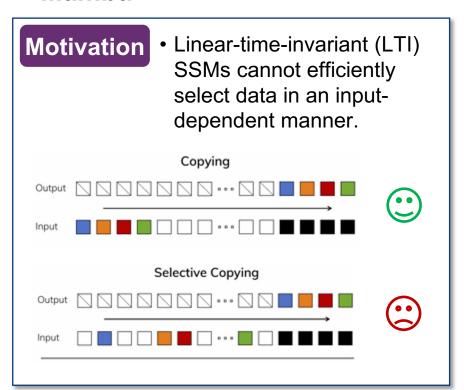


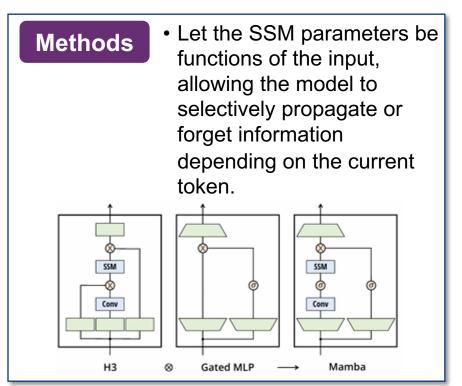
Research directions:

- 1. Design better parametrization or initialization strategy.
- 2. Design better model architecture based on SSM.

Mamba

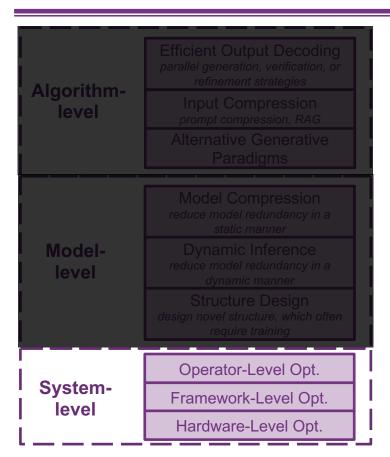
Mamba





[1] The router can direct simpler queries to smaller models and more complex ones to larger models, thereby balancing response quality with cost efficiency.

Menu of Techniques



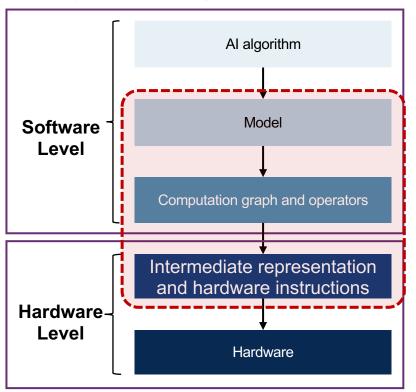
Operator-Level Opt.

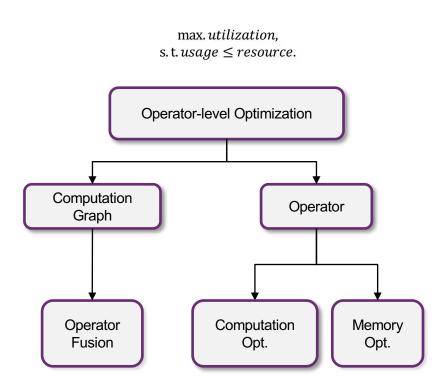
Framework-Level Opt.

Hardware-Level Opt.

Design Space

System Design: Operator-level Optimization

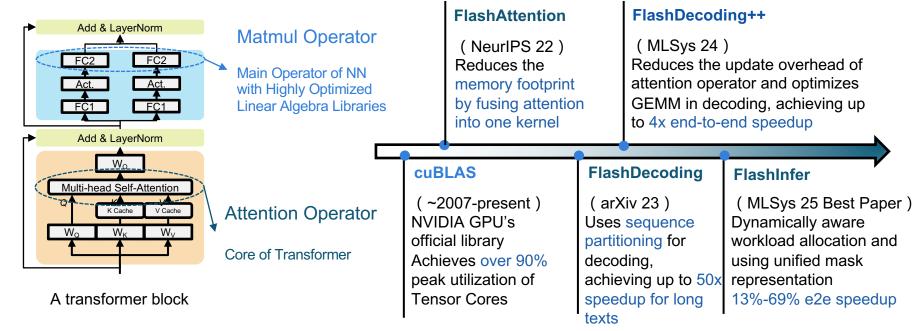




Operator-level Optimization



Optimizes the hardware utilization by tailoring workload mapping to hardware specifications



*Matmul: Matrix Multiplication

Implementation of Matmul on GPU

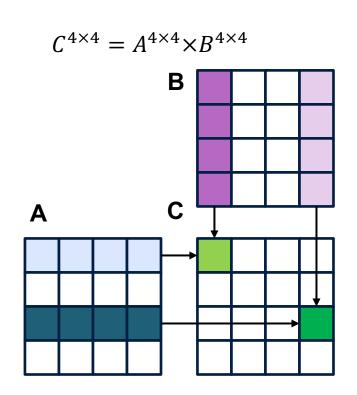
- Observation: No dependencies in C,
- → Parallel computation
 - Method: Each thread loads a row of A and a column of B → computes → writes to C

Computational-to-Memory Ratio

An important metric for measuring GPU utilization Higher value → Higher compute unit utilization

$$Comp-to-Mem Ratio = \frac{computation \ amount}{memory \ access \ amount}$$

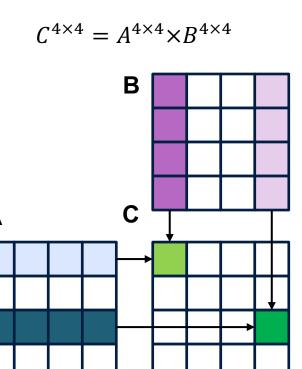
Comp-to-Mem Ratio is only **0.25FLOP/Byte** in this example



*Total computation: 2*4*4*4=128FLOP. Total memory access: 4*(2*4*4*4)=512Byte

Implementation of Matmul on GPU

```
Operator Implementation:
  global void MatMul(int *A, int *B, int *C, int width) {
  //eg. width=4
  int row = threadIdx.y;
  int col = threadIdx.x:
  if (row < width && col < width) {
     int sum = 0:
     for (int k = 0; k < width; k++) {
       sum += A[row * width + k] * B[k * width + col];
     \} //Each thread processed one row of A, one column of B
     C[row * width + col] = sum;
     //Write result to the corresponding position of C
```



Computation Optimization: Tiling

- Motivation: Maximize GPU compute utilization
- Method: Increasing tile size.

Example: 1 thread \rightarrow 2 rows of A & 2 cols of B

→ 4 elements in C

Tiling raises Comp-to-Mem Ratio to **0.5FLOP/Byte**

Oversized tiles reduce parallelism. The key is balance.

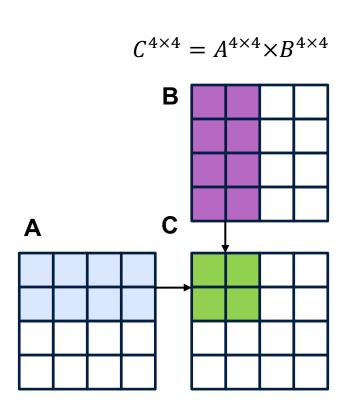
$$C^{4\times4} = A^{4\times4} \times B^{4\times4}$$

$$C$$

^{*}Total compute: 128FLOP, Total access reduces to4*(2*4*4*4)/2=256Byte

Computation Optimization: Tiling

```
Operator Implementation:
  global void MatMulTiling(int *A, int *B, int *C, int width) {
  //eq. TILE WIDTH=2
  int tx = threadIdx.x; int ty = threadIdy.y;
  for (int i = 0; i < TILE WIDTH; i++) {
     for (int j = 0; j < TILE WIDTH; j++) {
        int row = ty * TILE WIDTH + i;
        int col = tx * TILE WIDTH + j;
       int sum = 0:
       for (int k = 0; k < width; k++) {
          sum += A[row * width + k] * B[k * width + col];
        \\Read A and B from global memory and compute
        C[row * width + col] = sum;
  \{\text{//Each thread processes 2 rows of A and 2 columns of B}\}
```

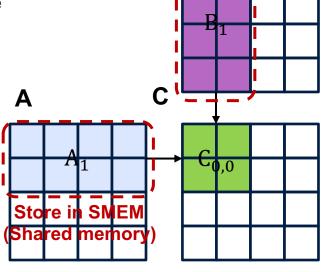


Memory Optimization: Using Shared Memory

- Motivation: Enable data sharing and fast communication between threads
- Method: Load A and B into shared memory for reuse
- → reduce global memory access

Comp-to-Mem Ratio reaches 1FLOP/Byte

$$C^{4\times4} = A^{4\times4} \times B^{4\times4}$$



^{*}store A and B into shared memory, total memory access reduces to 128Byte

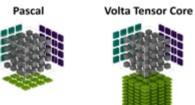
• Memory Optimization: Using Shared Memory Operator Implementation:

```
global void MatMul shared(int *A, int *B, int *C, int width) {
  shared int sharedA[width][width]; // Declare shared memory
  shared int sharedB[width][width];
int tx = threadIdx.x; int ty = threadIdx.y;
sharedA[ty][tx] = A[ty * width + tx]; // Write to SMEM
sharedB[ty][tx] = B[ty * width + tx];
  syncthreads(); // Sync all threads
for (int i = 0; i < TILE WIDTH; i++) {
   for (int j = 0; j < TILE WIDTH; j++) {
       // Read from SMEM & compute (row/col calculation are omitted)
       for (int k = 0; k < width; k++) {
        sum += sharedA[row][k] * sharedB[k][col];
     C[row * width + col] = sum;
```

```
C^{4\times4} = A^{4\times4} \times B^{4\times4}
Α
                               c_{0.0}
 Store in SMEM
Shared memory
```

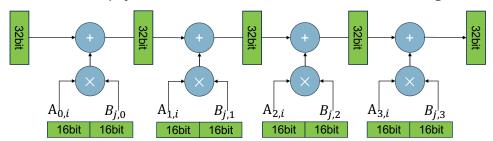
Special Hardware: using Tensor Core

- Why: CUDA Cores bottleneck at large-scale matrix ops in deep learning
- Tensor Core:
 - NVIDIA's dedicated DL cores.
 - First introduced in Volta (2017)



The V100 GPU uses
Tensor Cores achieves 2×
speedup over the P100[1]

Input: 16bit → Multiply-accumulate → Stored in 32-bit registers



[1] https://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf

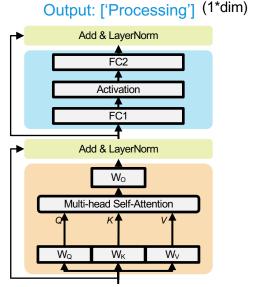
Special Hardware: using Tensor Core

CUDA provides mma.h and wmma API for Tensor Core operations Pesudo code: #include <mma.h> global void MatMul mma(half *a, half *b, float *c, int M, int N, int K) { // Declare fragments (16x16x16 tile size) wmma::fragment<wmma::matrix a, 16, 16, 16, half, wmma::row major> a frag; wmma::fragment<wmma::matrix b, 16, 16, 16, half, wmma::col major> b frag; wmma::fragment<wmma::accumulator, 16, 16, 16, float> c frag; wmma::fill fragment(c frag, 0.0f); // Init accumulator wmma::load matrix sync(a frag, a, K); wmma::load matrix sync(b frag, b, N); wmma::mma sync(c frag, a frag, b frag, c frag); // MM on Tensor Core wmma::store matrix sync(c, c frag, N, wmma::mem row major); // Store result

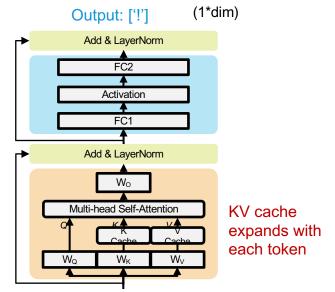
Decoder-only LLM Inference

Decoder-only LLM Inference: Two distinct phases

- Prefill phase: Handling input prompts, saving generated K and V in KV Cache
- Decode phase: Update and use KV cache for computing Attention to generate a new token



Prompt: ['I', 'like', 'natural', 'language'] (4*dim)



Prompt: ['1', 'like', 'natural', 'language', 'Processing] (1*dim)

Prefill Optimization: FlashAttention^[1,2,3]

- FlashAttention optimizes attention computation in prefill phase
- One of the most widely adopted acceleration methods with 16.2k GitHub stars^[4]

FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness

Tri Dao[†], Daniel Y. Fu[†], Stefano Ermon[†], Atri Rudra[‡], and Christopher Ré[†]

†Department of Computer Science, Stanford University

†Department of Computer Science and Engineering, University at Buffalo, SUNY

June 24, 2022



^[1] Dao, Tri, et al. "Flashattention: Fast and memory-efficient exact attention with io-awareness." Advances in neural information processing systems 35 (2022): 16344-16359.

^[2] Dao, Tri. "Flashattention-2: Faster attention with better parallelism and work partitioning." arXiv preprint arXiv:2307.08691 (2023).

^[3] Shah, Jay, et al. "Flashattention-3: Fast and accurate attention with asynchrony and low-precision." Advances in Neural Information Processing Systems 37 (2024): 68658-68685.

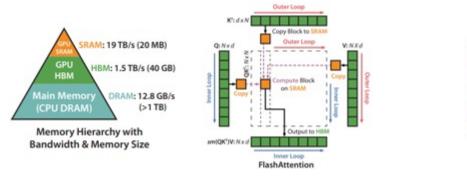
^[4] https://github.com/Dao-AlLab/flash-attention

Prefill Optimization: FlashAttention^[1,2,3]

Why: Complex Attention I/O; Large activation memory

How: Operator fusion, including fwd and bwd

• Results: 2-4x speedup; memory: $O(N^2) \rightarrow O(N)$



Saves I/O & memory via operator fusion



Tiling strategies differ between fwd/bwd passes

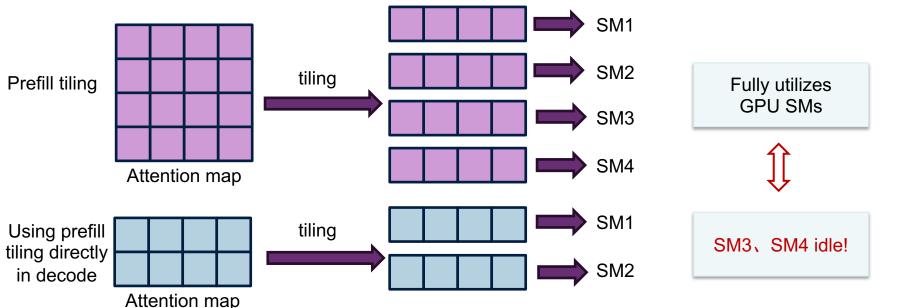
^[1] Dao, Tri, et al. "Flashattention: Fast and memory-efficient exact attention with io-awareness." Advances in neural information processing systems 35 (2022): 16344-16359.

^[2] Dao, Tri. "Flashattention-2: Faster attention with better parallelism and work partitioning." arXiv preprint arXiv:2307.08691 (2023).

^[3] Shah, Jay, et al. "Flashattention-3: Fast and accurate attention with asynchrony and low-precision." Advances in Neural Information Processing Systems 37 (2024): 68658-68685.

Decode Optimization: FlashDecoding^[1]

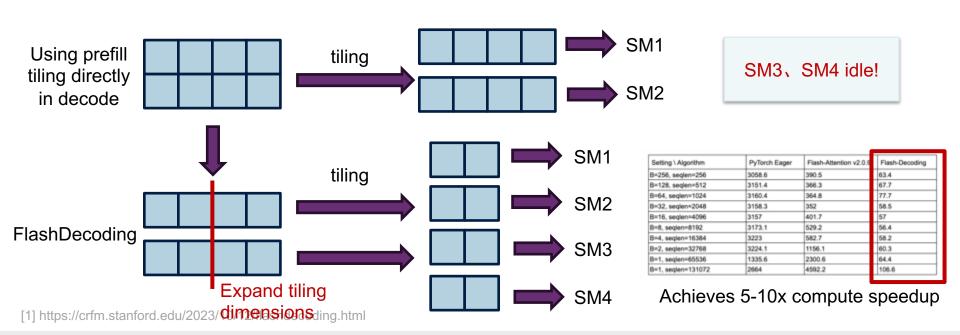
Why: only 1 token/step in decode → using prefill tiling strategy directly causes low GPU utilization



[1] https://crfm.stanford.edu/2023/10/12/flashdecoding.html

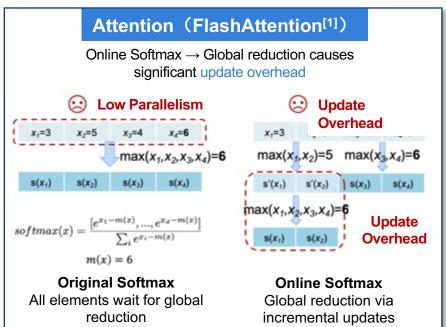
Decode Optimization: FlashDecoding^[1]

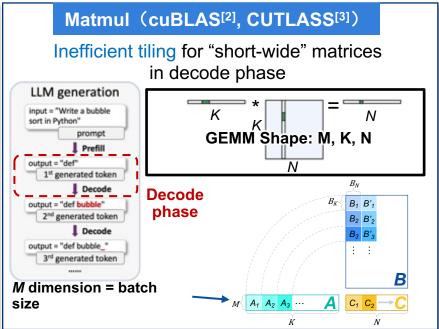
Method: Increases matrix tile count to boost SM utilization, achieving higher GPU efficiency





Attention & Matmul operators still optimizable in LLM inference





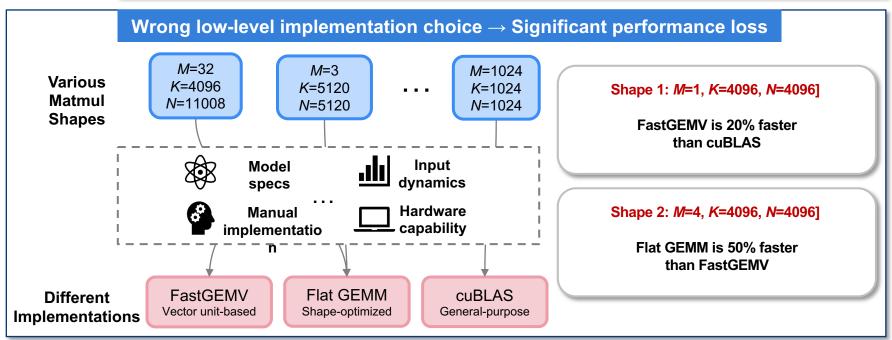
^[1] Dao T, Fu D, Ermon S, et al. Flashattention: Fast and memory-efficient exact attention with io-awareness[J]. Advances in neural information processing systems, 2022, 35: 16344-16359.

^[2] https://developer.nvidia.com/cublas

^[3] https://github.com/NVIDIA/cutlass

Motivation

Suboptimal implementation against various Matmul shapes



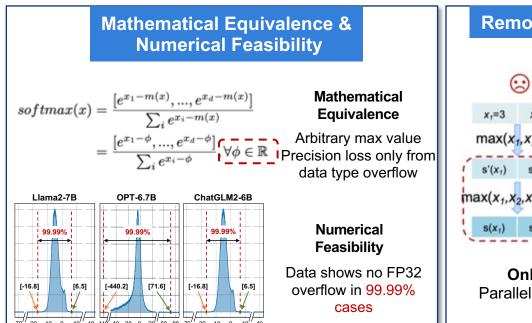
^[1] https://github.com/wangsiping97/FastGEMV

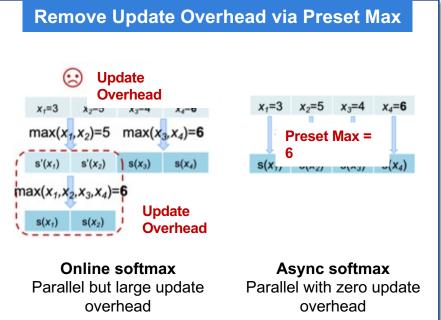
^[2] https://developer.nvidia.com/cublas

^[3] https://github.com/NVIDIA/cutlass

Method

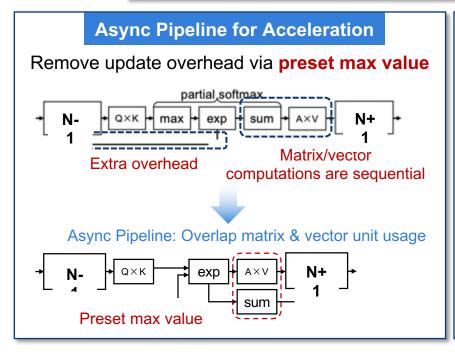
Leverage value distribution in LLM inference to optimize Attention operator^[1]

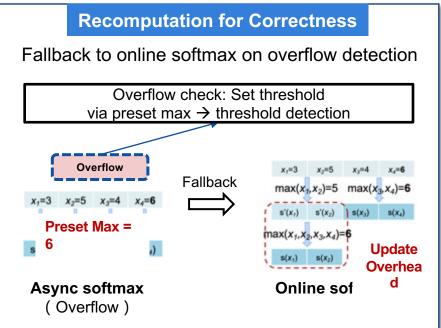




Method

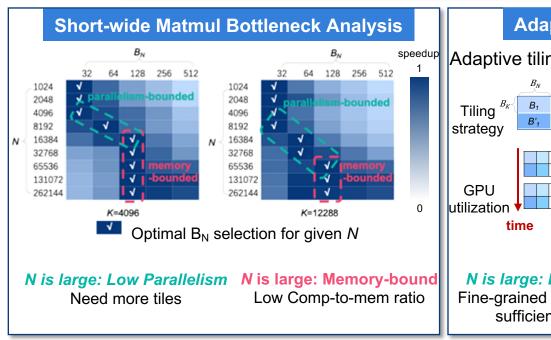
Leverage value distribution in LLM inference to optimize Attention operator^[1]

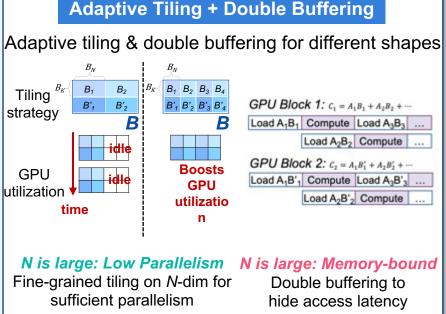




Method

Leverage typical shapes in LLM inference to optimize Matmul operator^[1]

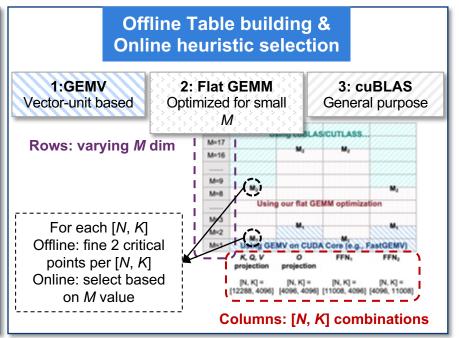




Method

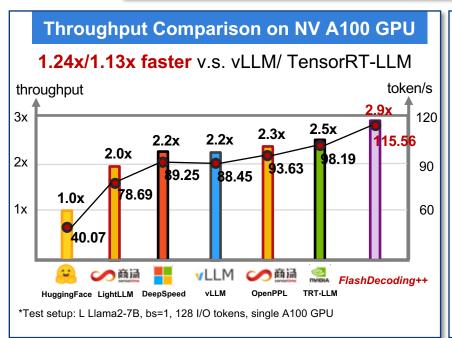
Leverage Matmul shape patterns in LLM inference for dynamic implementation selection^[1]

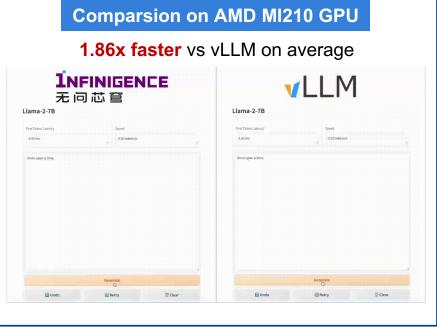
Shape Variation Patterns 4 shape categories in inference with only M dimension varies Observation 1: SegLen*B HD*3 K. Q. V projection Fixed $[N, K] \rightarrow \text{only}$ O projection SegLen*B HD Prefill 4 combinations per phase EEN1 SegLen*B | FD model SeqLen*B FD FFN2 K, Q, V projection Observation 2: O projection Decode Only M dimension phase FFN1 varies with input FFN2 Prefill: M=total token count HD: Hidden dimension size Only 4 shape FD: Dimension size after the first FFN Decode: M=batch size categories! SegLen: Input seguence length



Method

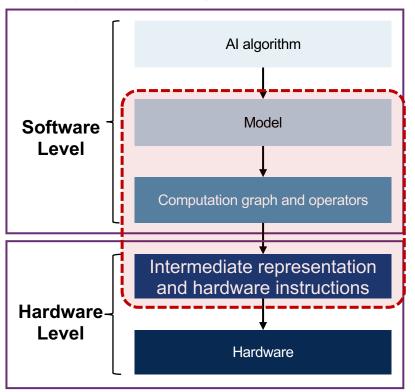
Throughput surpasses SOTA by over 10% 1.88x faster v.s. HuggingFace on average

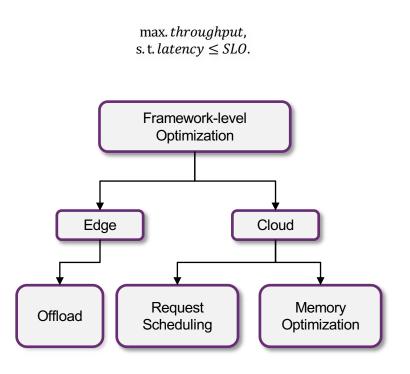




Design Space

System Design: Framework-level Optimization

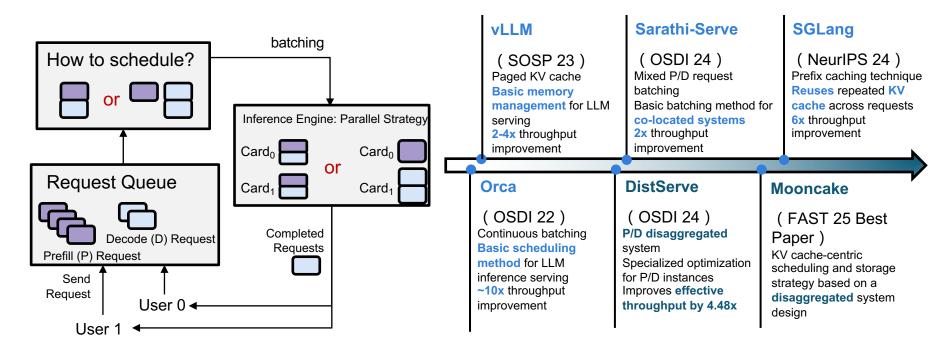




Framework-level Optimization



Optimizing the system throughput adhering to the service-level objective (SLO)



Offloading Techniques

Offloading Technique: Definition and Motivation

 Definition: Offloading parts of the model (model weights, KV cache, etc.) from GPU to other devices (e.g., CPU) for storage or even computation, to save space and improve computational efficiency

Motivation: The large volume of model parameters and KV cache data exceeds the

storage capacity of GPU memory

Table 1: LLaMA2-13B, KV Cache size with context length

| Context length | 10k | 100k | 500k | 1000k |
|----------------|--------|--------|---------|---------|
| KV Cache size | 8.19GB | 81.9GB | 409.6GB | 819.2GE |
| Misc size | 26GB | 26GB | 26GB | 26GB |

LLaMA2-13B: when the context length reaches 100k, the required KV cache reaches nearly 82GB, exceeding the memory capacity of a single GPU

Table 3: Model parameter counting at layer-level (dtype: BF/FP16)

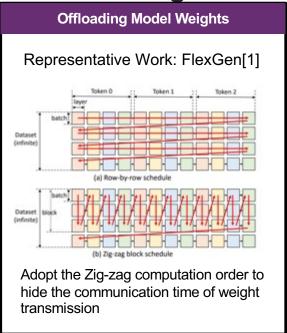
| Layers | Modules | Shapes | No. Parameters | Per Layer | MB | $^{\mathrm{GB}}$ |
|---------------|-----------|------------------------|----------------|-----------|-----------|------------------|
| Layer 0 | Embedding | [129280, 7168] | 926,679,040 | 1.5 B | 2880 | 2.8 |
| | MLA | | 187,107,328 | | | |
| | MLP | 3 * [7168, 18432] | 396,361,728 | | | |
| | LN | 2*7168+1536+512 | 16,384 | | | |
| Layers 1 - 2 | MLA | ********** | 187,107,328 | 0.58 B | 1112 | 1.1 |
| | MLP | 3 * [7168, 18432] | 396,361,728 | | | |
| | LN | 2*7168+1536+512 | 16,384 | | | |
| Layers 3 - 59 | MLA | * | 187,107,328 | 11.5 B | 21950 | 21.44 |
| | Gate | [256, 7168] | 1,835,008 | | | |
| | MoE | 3 * [7168, 2048] * 257 | 11,318,329,344 | | | |
| | LN | 2*7168+1536+512 | 16,384 | | | |
| Layer 60 | MLA | *11 177 | 187,107,328 | 12.4 B | 23712 | 23.16 |
| | Gate | [256, 7168] | 1,835,008 | | | |
| | MoE | 3 * [7168, 2048] * 257 | 11,318,329,344 | | | |
| | LN | 2*7168+1536+512 | 16,384 | | | |
| 11 | Head | [7168, 129280] | 926,679,040 | | | |
| Total | | 2010/10/2010/2010 | | 671 B | 1.280,000 | 1250 |

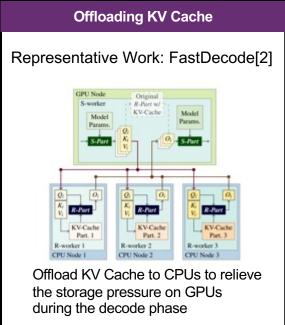
DeepSeek-V3: parameter size reaches 1250GB

[1] Memory Analysis on the Training Course of DeepSeek Models, Zhang et al. Arxiv Preprint 2502.07486.

Offloading Techniques

Offloading Techniques: Categories







- [1] FlexGen: High-Throughput Generative Inference of Large Language Models with a Single GPU, Sheng et al. Arxiv Preprint 2303.06865.
- [2] FastDecode: High-Throughput GPU-Efficient LLM Serving using Heterogeneous Pipelines, He et al. Arxiv Preprint 2403.11421.

[3] https://github.com/kvcache-ai/ktransformers

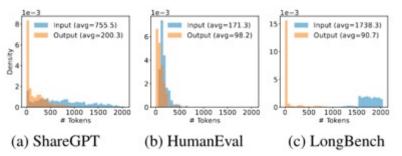
Request Scheduling

Request Scheduling

Definition: Scheduling requests in different phases to GPU instances for batched computation

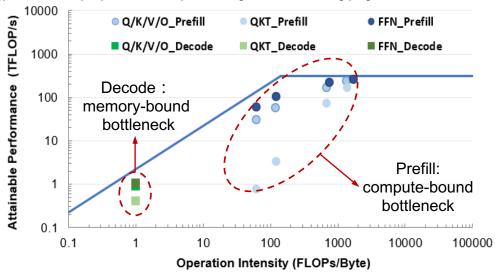
Motivation: Requests for LLM inference have varying lengths and distinct phases (Prefill

and Decode phases), and schedu



Statistical distribution of request lengths across different datasets

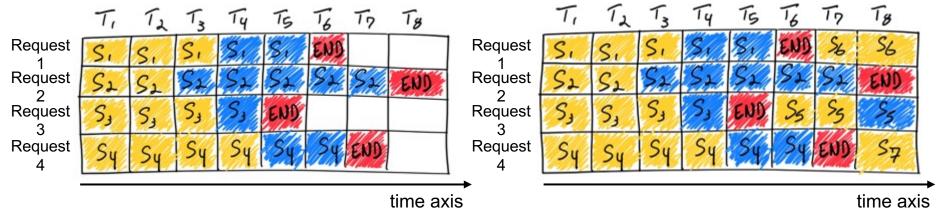
How to perform batch processing for requests with varying lengths?



Request Scheduling

Request Scheduling: Batching

 Orca (G. Yu et al, OSDI'22) proposes the continuous batching technique, which batches requests with varying lengths at the TOKEN granularity. Compared to request-level batching, it improves throughput by 36.9x



Request-level batching
The processing time depends on the longest request, leading to low utilization

Token-level batching (Continuous batching)
Batching requests of different lengths by concatenating
them in the Token dimension

[1] G. Yu, et al. "ORCA: A Distributed Serving System for Transformer-Based Generative Models.", OSDI, 2022.

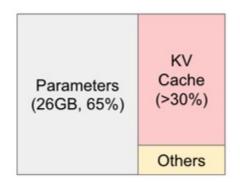
Memory Optimization

Memory Optimization: PagedAttention



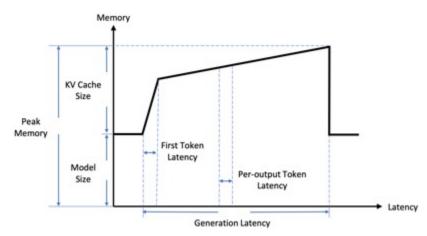
The growth of KV cache and memory fragmentation limit concurrency, resulting suboptimal system throughput.

vLLM (W. Kwon et al, SOSP'23) proposes PagedAttention, which stores KV cache in a paged manner. This approach effectively eliminates memory fragmentation, and improves throughput by 2-4 times.



NVIDIA A100 40GB

KV cache accounts for a large proportion in distributed systems.



KV cache in LLM inference grows with the generation process

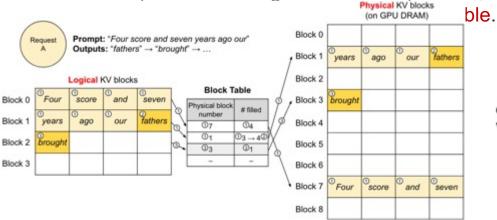
[1] W. Kwon, et al. " Efficient Memory Management for Large Language Model Serving with PagedAttention.", SOSP, 2023.

Memory Optimization

Memory Optimization: PagedAttention

• Memory: KV cache is stored in blocks along the sequence dimension at different physical addresses.

Computation: During attention computation, the physical addresses of the corresponding



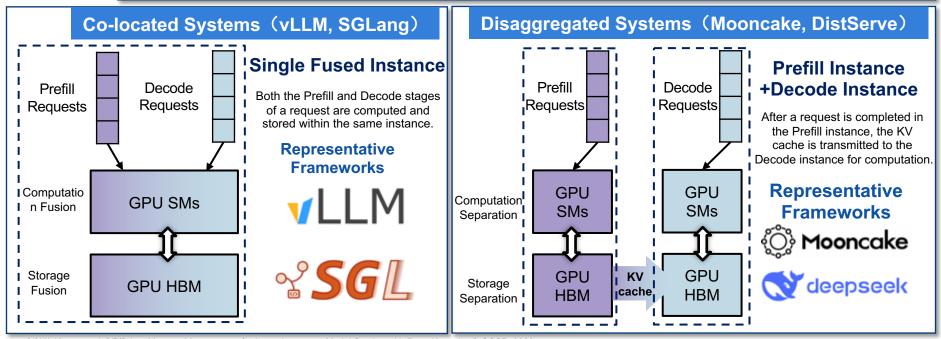
Memory: Logically continuous sequences are stored in blocks at actual physical addresses.

Key and value vectors Block 1 years ago our fathers Block 2 brought forth Query forth vector Block 0 Four and score seven

Computation: Index the KV cache addresses required for computation by looking up a table.

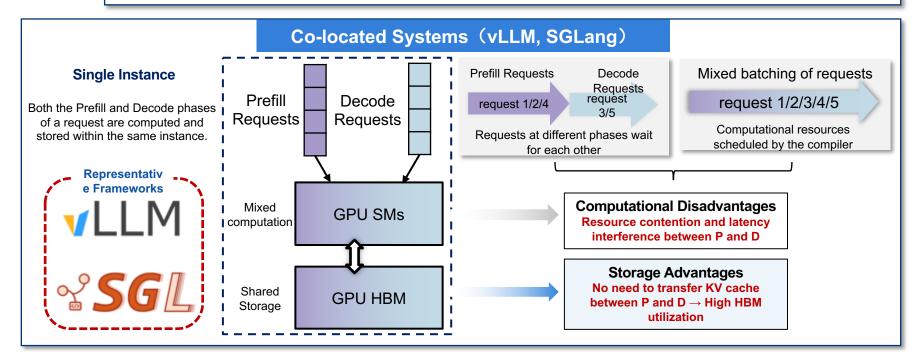
[1] W. Kwon, et al. " Efficient Memory Management for Large Language Model Serving with PagedAttention.", SOSP, 2023.

Key Problem Co-located and disaggregated systems for Prefill/Decode have their own strengths and weaknesses in computation and storage.

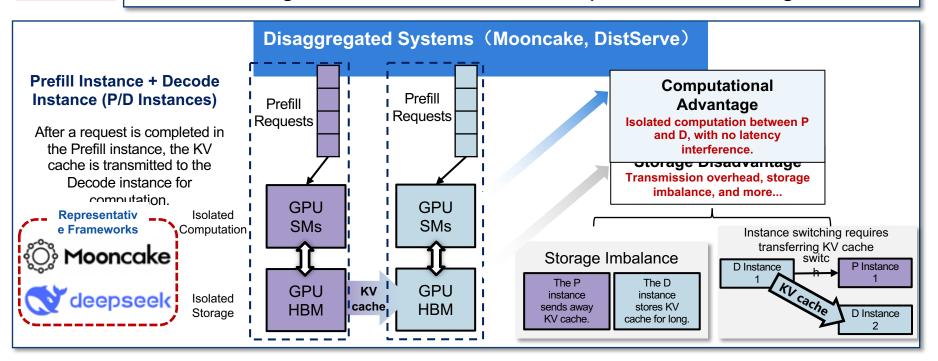


- [1] W. Kwon, et al. " Efficient Memory Management for Large Language Model Serving with PagedAttention.", SOSP, 2023.
- [2] L. Zheng, et al. "SGLang: Efficient Execution of Structured Language Model Programs", NeurlPS, 2024.
- [3] R. Qin, et al. "Mooncake: Trading More Storage for Less Computation", FAST, 2025.

Key Problem Co-located and disaggregated systems for Prefill/Decode have their own strengths and weaknesses in computation and storage.



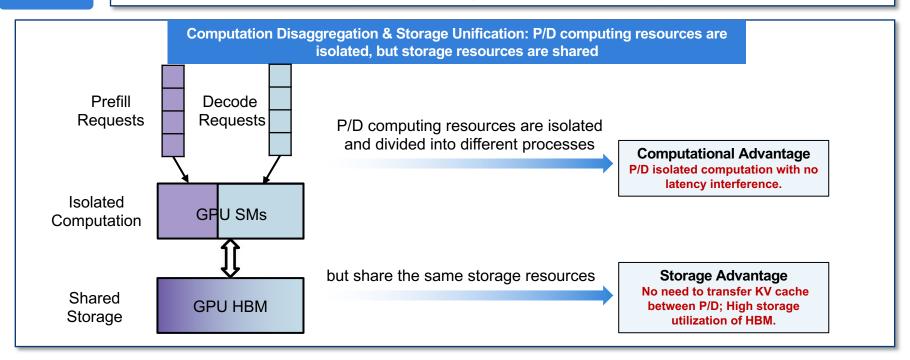
Key Problem Co-located and disaggregated systems for Prefill/Decode have their own strengths and weaknesses in computation and storage.



^[1] R. Qin, et al. "Mooncake: Trading More Storage for Less Computation", FAST, 2025.

Method

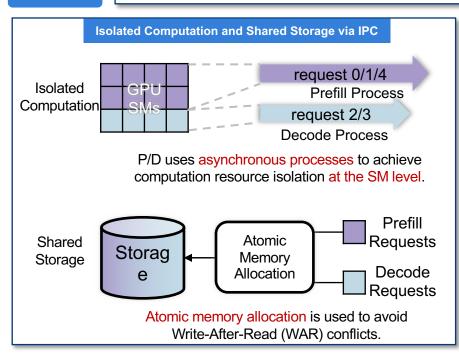
Combine the computational advantage of disaggregated systems and the storage advantage of co-located systems^[1]

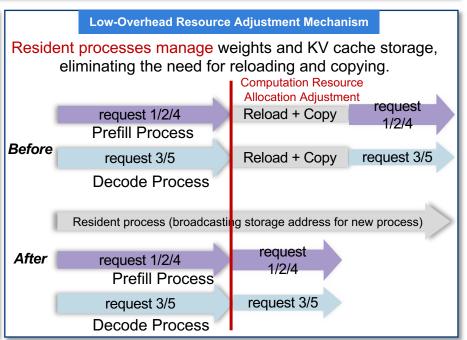


[1] Ke, Hong, et al. "semi-PD: Towards Efficient LLM Serving via Phase-wise Computation Disaggregation and Unified Storage." arXiv preprint arXiv:2504.19867. 2025.

Method

Combine the computational advantage of disaggregated systems and the storage advantage of co-located systems^[1]

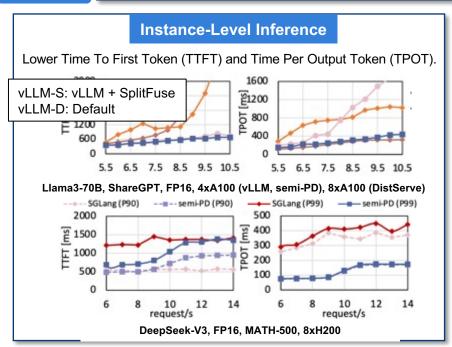


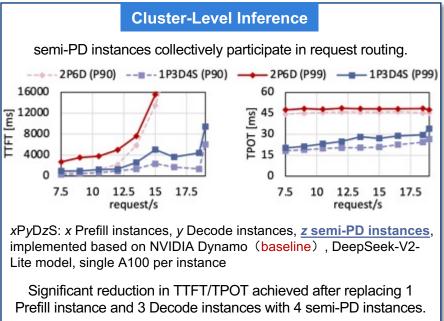


[1] Ke, Hong, et al. "semi-PD: Towards Efficient LLM Serving via Phase-wise Computation Disaggregation and Unified Storage." arXiv preprint arXiv:2504.19867. 2025.

Result

Llama3 series models: 1.55-1.72x improvement in request service rate under given SLOs DeepSeek-V3 model: 1.49-2.58x reduction in latency

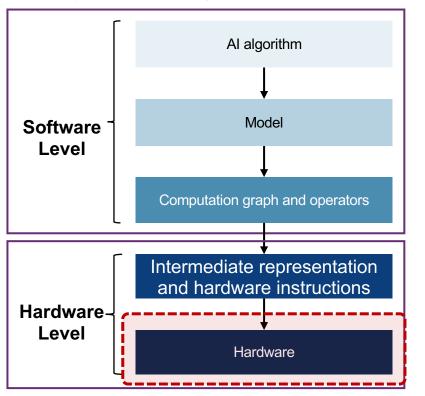


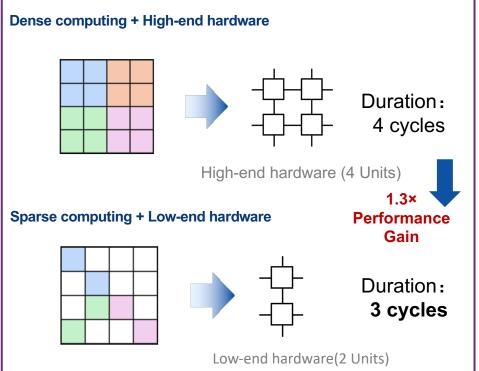


[1] Ke, Hong, et al. "semi-PD: Towards Efficient LLM Serving via Phase-wise Computation Disaggregation and Unified Storage." arXiv preprint arXiv:2504.19867. 2025.

Design Space

System Design: Hardware-level Optimization





Video generation is one of the important modalities, and it is a promising path towards as the physical world's simulators

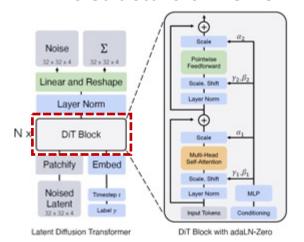
OpenAl's Sora model Videos generated by Sora are almost indistinguishable from real ones [1]



Source: [1] OpenAI, https://openai.com/index/sora/

VGMs are mainly based on diffusion transformers, composed of noising and denoising process

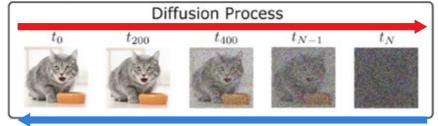
The structure of VGMs



At present, most of the mainstream video generation models adopt the Diffusion Transformer (DiT) architecture^[1]

Diffusion

Forward diffusion (backward propagation): Gradually add Gaussian noise of different amplitudes



Reverse diffusion (inference): Denoise gradually

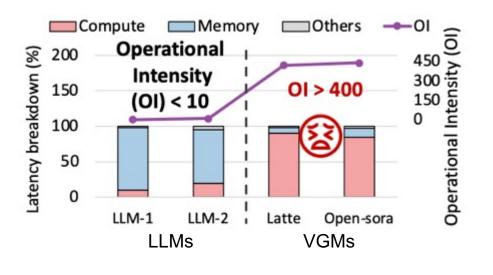
The principle is to endow the model with the ability to generate videos through multiple rounds of noise superposition and denoising training.

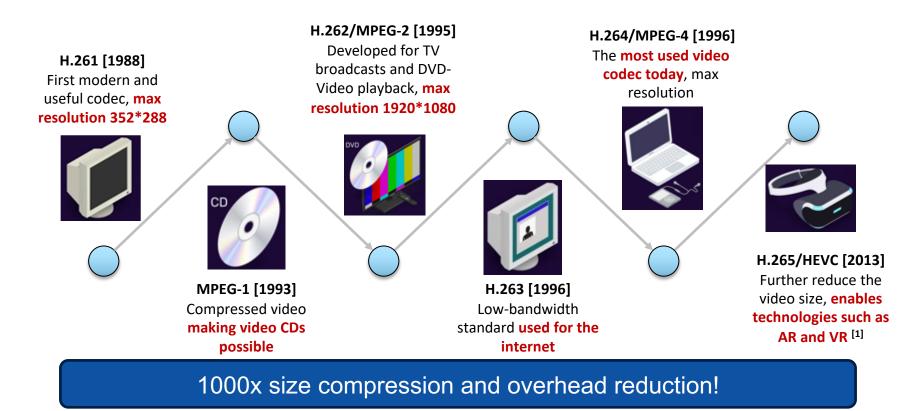
Source: [1] Scalable diffusion models with transformers. ICCV 2023.

Different from LLMs, the inference bottleneck of video generation models (VGMs) has shifted from memory-bound to compute-bound

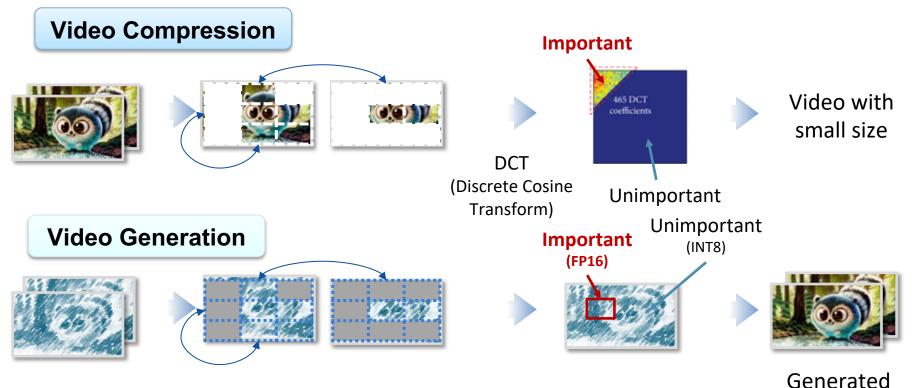
| | LLMs (mainly considering the decode stage) | VGMs (Temporal, spatial, and FFN structure) | |
|----------------------------|---|---|--|
| Computation* | $\approx 12d^2 + 2Nd$ | $\approx 16FNd^2 + 2FN(F+N)d$ | |
| Memory access* | $\approx 12d^2 + 2Nd$ | $\approx 16d^2 + 15FNd$ | |
| Operational intensity (OI) | ≈ 1 | $\approx FN$ | |

F: Video frames. N: Tokens. d: Hidden dimension.
*Calculation of a single block/layer.





 $Source: \ [1] \ https://medium.com/videocoin/the-history-of-video-codecs-infographic-432e5be1154f.$

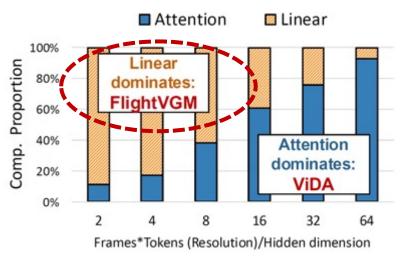


1. Activation sparsification

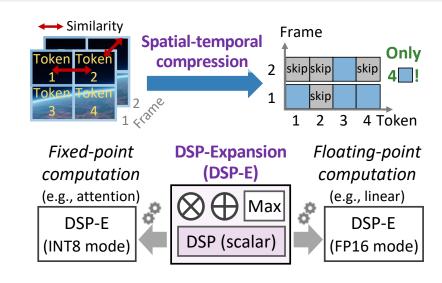
2. Hybrid precision quantization

video

For linear-heavy video generation models, we propose FlightVGM, a HW-SW co-design with temporal-spatial & floating-fixed strategies



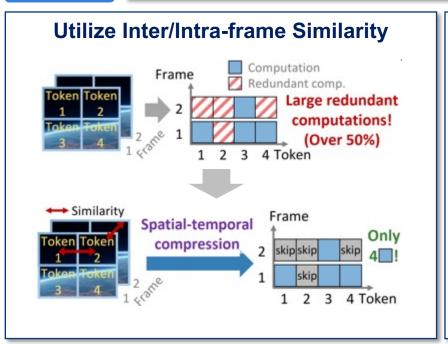
- ➤ Attention dominated: resolution ↑ or dim.↓
- ▶ Linear dominated: resolution↓ or dim.↑

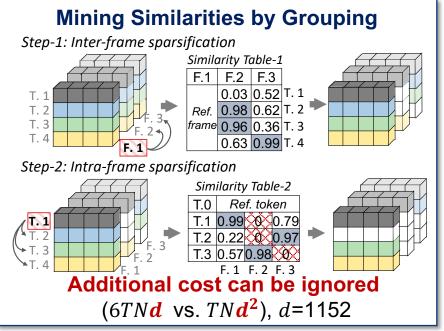


Sparse Computing + Configurable Design

→ Higher Performance

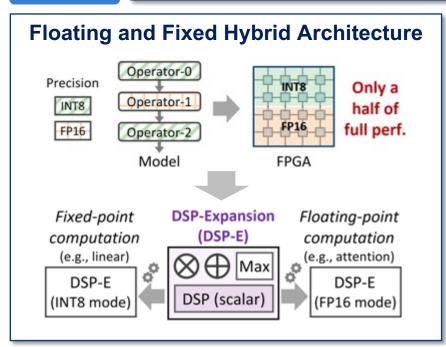
Main Method To address the computational redundancy problem in VGMs, a temporal-spatial joint sparse method is proposed

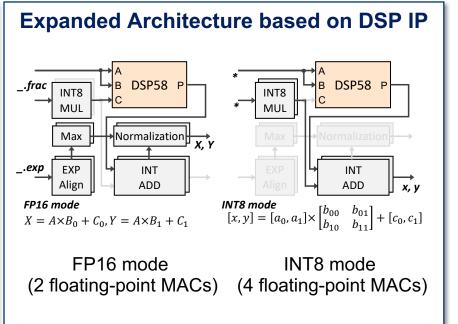




[1] Liu J, Zeng S, Ding L, et al. Flightvgm: Efficient video generation model inference with online sparsification and hybrid precision on fpgas[C]//Proceedings of the 2025 ACM/SIGDA International Symposium on Field Programmable Gate Arrays. 2025: 2-13.

Main Method A DSP58 extension architecture is proposed to enhance utilization under mixed-precision computing





[1] Liu J, Zeng S, Ding L, et al. Flightvgm: Efficient video generation model inference with online sparsification and hybrid precision on fpgas[C]//Proceedings of the 2025 ACM/SIGDA International Symposium on Field Programmable Gate Arrays. 2025: 2-13.

Models and datasets

■ Models: Latte-1 and Open-Sora 1 Table 1: Hardware parameters of GPU and FPGA platforms.

■ Datasets: UCF-101

Metrics

CLIPSIM: Text-Video Alignment

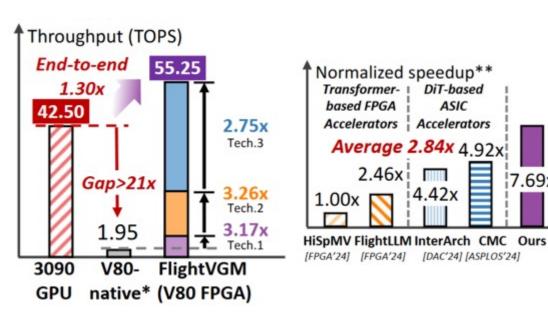
VBench: Video quality

Baseline

| Platform | NVIDIA 3090 GPU | AMD U280 FPGA | AMD V80 FPGA |
|-----------------|--------------------|------------------|-----------------|
| C | 328 | 9024 | 10848 |
| Compute units | Tensor Cores | DSP48s | DSP58s |
| Frequency (MHz) | 1695 | 225 | 300 |
| FP16 PCP (TOPS) | 142 | 5.4 | 6.5 |
| Memory (GB) | 24 | 8 & 32 | 32 & 32 |
| BW (GB/s) | 936 | 38 & 460 | 51 & 819 |

- Generic hardware: NVIDIA 3090 GPU
- FPGA-based accelerator for Transformer: HiSpMV[FPGA'24] and FlightLLM[FPGA'24]
- ASIC-based accelerator for DiT: InterArch[DAC'24] and CMC[ASPLOS'24]

With 21× lower peak compute than the NVIDIA 3090, FlightVGM achieves 1.3× higher speedup and 4.5× better energy efficiency





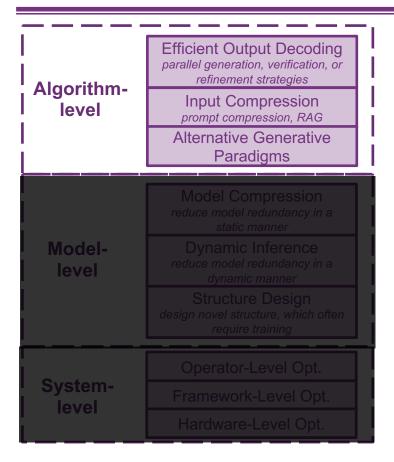
Video generated by original model



Video generated by efficient model
[1] Liu J, Zeng S, Ding L, et al. Flightvgm: Efficient video generation model inference with online sparsification and hybrid precision on fpgas[C]//Proceedings of the 2025 ACM/SIGDA International Symposium on Field Programmable Gate Arrays. 2025: 2-13.

17.69x

Menu of Techniques



Efficient Output Decoding

- Speculative Decoding
- Jacobi Decoding
- Agentic Generation

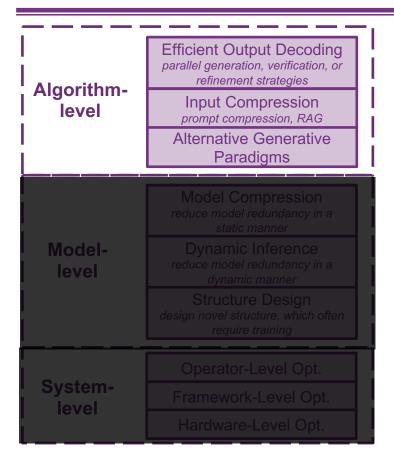
Input Compression

Input Compression

Alternative Generative Paradigms

Diffusion for Text

Menu of Techniques



Efficient Output Decoding

- Speculative Decoding
- Jacobi Decoding
- Agentic Generation

Input Compression

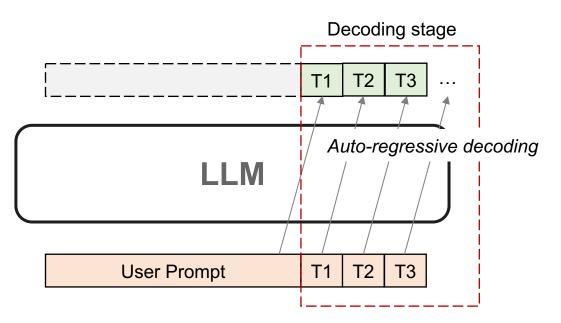
Input Compression

Alternative Generative Paradigms

Diffusion for Text

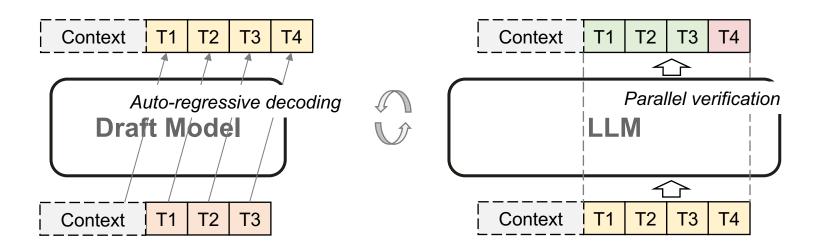
Speculative Decoding: Concepts

- Recall: autoregressive decoding of LLMs
 - Decoding stage: memory bound (data transfer of model weights & KV cache)
 - Redundant computation is left unused!



Speculative Decoding: Concepts

- Core idea of speculative decoding
 - Use a small draft model to generate multiple token for verification
 - The LLM conducts parallel verification (memory bound allows more computation)
 - Key elements: 1) the acceptance rate of generated tokens; 2) the cost of draft model



Speculative Decoding: Demo

Demo from [1]

- Green: accepted tokens
- Red: rejected tokens
- Blue: corrected tokens

```
[START] japan ' s benchmark nikkei 22 75

[START] japan ' s benchmark nikkei 225 index rose 22 76

[START] japan ' s benchmark nikkei 225 index rose 226 69 7 points

[START] japan ' s benchmark nikkei 225 index rose 226 69 points, or 0 1

[START] japan ' s benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 9859

[START] japan ' s benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 7 in

[START] japan ' s benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 7 in

[START] japan ' s benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 in tokyo late

[START] japan ' s benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 in tokyo late

[START] japan ' s benchmark nikkei 225 index rose 226 69 points, or 1 5 percent, to 10, 989 79 in late morning trading [END]
```

Speculative Decoding: Speed-up Estimation

- The speed-up rate of speculative decoding (SD) can be estimated!
 - **S**: the total number of tokens
 - R: the number of SD rounds
 - y: the number of generated tokens in each SD round

$$Speedup = \frac{T_{AR}}{T_{SD}} = \frac{S \cdot T_{T}(B, 1)}{R \cdot \left(\gamma \cdot T_{D}(B, 1) + T_{T}(B, \gamma) + T_{reject}\right)} = \frac{S}{R} \cdot \frac{1}{\gamma \cdot \left[\frac{T_{D}(B, 1)}{T_{T}(B, 1)} + \left[\frac{T_{T}(B, \gamma)}{T_{T}(B, 1)} + \left[\frac{T_{reject}}{T_{T}(B, 1)}\right]\right]}$$

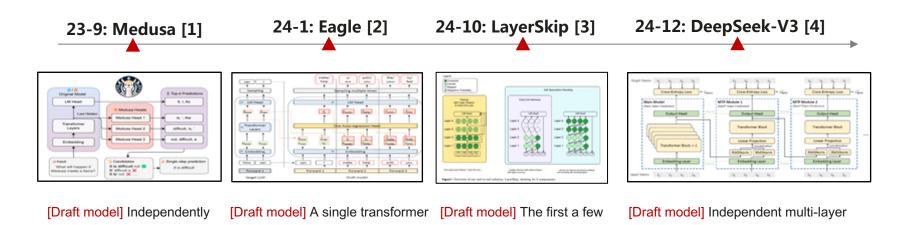
- ① $\frac{T_D(B,1)}{T_T(B,1)}$ The relative latency of draft model to the target model
- ② $\frac{T_T(B,1)}{T_T(B,\gamma)}$ The cost of multi-token verification. A large batch size B is harmful to speed-up rate
- $\frac{T_{reject}}{T_T(B,1)}$ The negligible cost of token sampling for the rejected tokens

^[1] Sadhukhan, Ranajoy, et al. Magicdec: Breaking the latency-throughput tradeoff for long context generation with speculative decoding. ICLR 2025.

^[2] Huang, Zongle, et al. MoESD: Unveil Speculative Decoding's Potential for Accelerating Sparse MoE. arXiv preprint arXiv:2505.19645 (2025).

Speculative Decoding: Representative Works

- How to find an good draft model?
 - consistent with the target model
 - efficient in decoding



layers of the LLM itself

layer taking the output from LLM

trained multi-layer decoders

perception pre-trained together with DS-V3

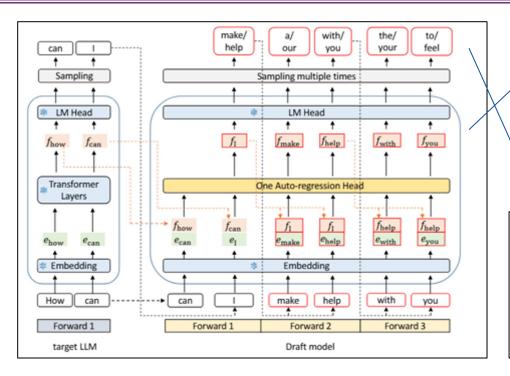
^[1] Tianle Cai, et. al. Medusa: Simple Ilm inference acceleration framework with multiple decoding heads. ICML 2024.

^[2] Yuhui Li, et. al, EAGLE: Speculative Sampling Requires Rethinking Feature Uncertainty, ICML 2024

^[3] Mostafa Elhoushi, et. al. LayerSkip: Enabling Early Exit Inference and Self-Speculative Decoding, arXiv preprint, 2024

^[4] Deepseek Team. DeepSeek-V3 Technical Report.

Speculative Decoding: Eagle



Drat model: a single transformer layer

params: 0.25B ~ 1B

- # training data: 1B tokens

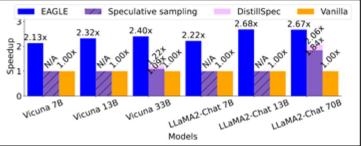
Acceptance rate: 75%

Tree attention: more tokens per pass

Verify more tokens per pass

- E.g., 12 tokens (4 paths)





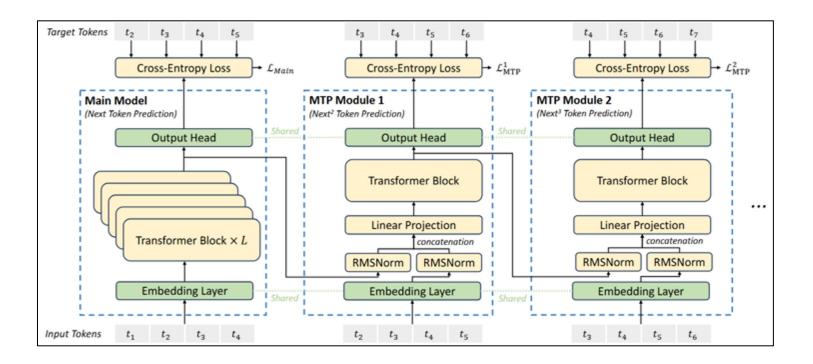
^[1] Yuhui Li, et. al, EAGLE: Speculative Sampling Requires Rethinking Feature Uncertainty, ICML 2024

^[2] Yuhui Li, et, al, Eagle-2: Faster inference of language models with dynamic draft trees. EMNLP 2024

^[3] Yuhui Li, et. al, EAGLE-3: Scaling up Inference Acceleration of Large Language Models via Training-Time Test, arXiv preprint 2503.01840.

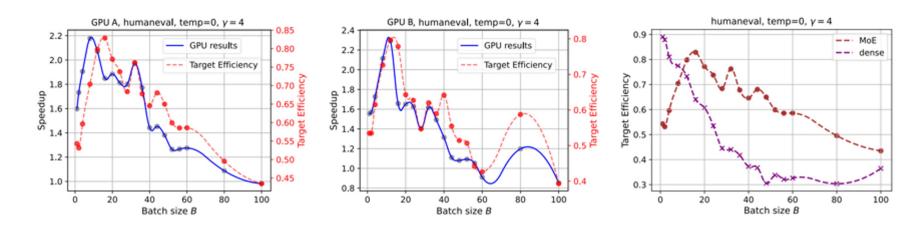
Speculative Decoding: Multi-Token Prediction

Multi-Token Prediction (MTP): trained from scratch with the LLM backbone

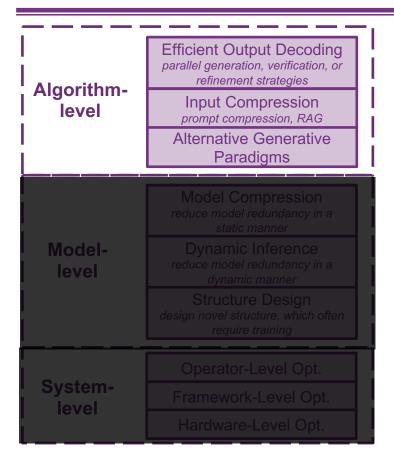


Speculative Decoding for MoE Architectures

- SD is hard to tackle MoE:
 - MoE does not favor small batch: Additional memory costs to load experts
 - SD does not favor large batch: SD is beneficial for memory-bound systems
- SD + MoE is helpful under medium batch size
 - The number of activated experts saturates, but does not reach compute-bound



Menu of Techniques



Efficient Output Decoding

- Speculative Decoding
- Jacobi Decoding
- Agentic Generation

Input Compression

Input Compression

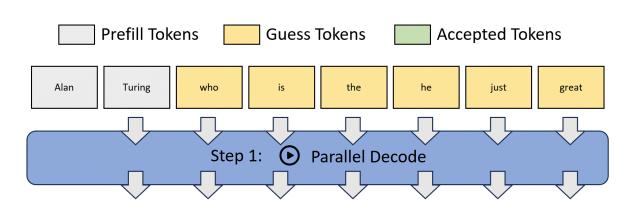
Alternative Generative Paradigms

Diffusion for Text

Jacobian Decoding

Drawbacks of speculative decoding

- Low acceptance rate, decoding time etc.
- A draft model needs to be separately trained and paired with the target LLM



Jacobian decoding:

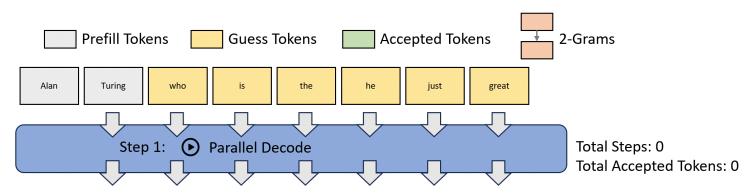
A single LLM without the draft model

Total Steps: 0

Total Accepted Tokens: 0

Lookahead Decoding

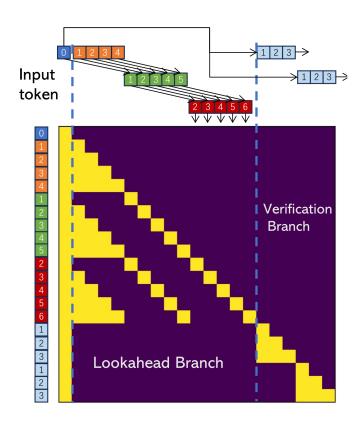
"Lookahead": reuse the promising draft from past N-gram trajectories



2-Gram Pool

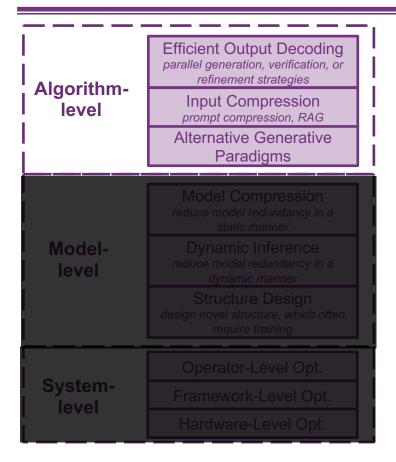
Lookahead Decoding

Lookahead branch maintains a fixed-sized, 2D window to generate n-grams from the Jacobi iteration trajectory.



Verification branch selects and verifies promising n-gram candidates.

Menu of Techniques



Efficient Output Decoding

- Speculative Decoding
- Jacobi Decoding
- Agentic Generation

Input Compression

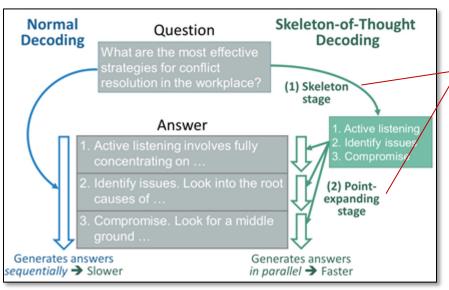
Input Compression

Alternative Generative Paradigms

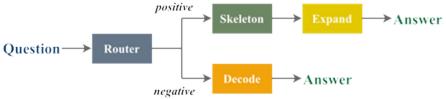
Diffusion for Text

Skeleton-of-Thought (SoT)

SoT: LLM generates the skeleton autoregressively, and then each points in parallel (an attempt in agentic generation for efficiency)

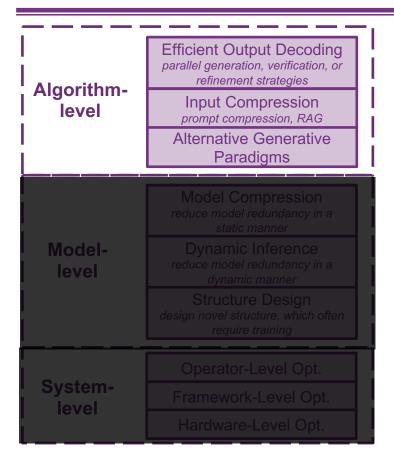


- Skeleton Stage: Guide the LLM to output a concise skeleton of the answer
- Point-expanding Stage: Guide the LLM to expand on each point from the skeleton in parallel
- Achieve up to 2.39x end-to-end speed-up



SoT in practice: A router to classify queries (SoT-R)

Menu of Techniques



Efficient Output Decoding

- Speculative Decoding
- Jacobi Decoding
- Agentic Generation

Input Compression

Input Compression

Alternative Generative Paradigms

Diffusion for Text

Prompt Compression

Prompt compression: eliminate redundant tokens in the prompt

Selective Context: filter out redundant tokens to shorten the input prompt

Original: INTRODUCTION Continual Learning (CL); also known as Lifelong Learning; is a promising learning paradigm to design models that have to learn how to perform multiple-tasks across different environments over their lifetime {To uniform the language and enhance the readability of the paper we adopt the unique term continual learning (CL); I ldeal CL models in the real world should be deal with domain shifts; researchers have recently started to sample tasks from two different datasets. For instance; proposed to train and evaluate a model on Imagenet first and then challenge its performance on the Places 365 dataset; considers more scenarios; starting with Imagenet or Places 365; and then moving on to the VOC/CUB/Scenes datasets. Few works propose more advanced scenarios built on top of more than two datasets.

Filtered: INTRODUCTION Continual Learning (a promising learning paradigm to design models have to how across over To uniform the language and enhance adopt the unique term continual learning Ideal CL models in should deal domain shifts researchers recently started sample tasks two different datasets For instance proposed to train and evaluate on Imagenet first challenge Places 365 considers more scenarios starting Imagenet or Places 365 the VOC/CUB/Scenes datasets Few works propose more advanced scenarios built top more than two datasets

1. Identify the token importance

$$I(x) = -\log_2 P(x_t|x_0, x_1, ..., x_{t-1})$$

2. Group tokens to units

$$I(u) = \sum_{i=t}^{n} I(x_i)$$

3. Sort units in the descending order

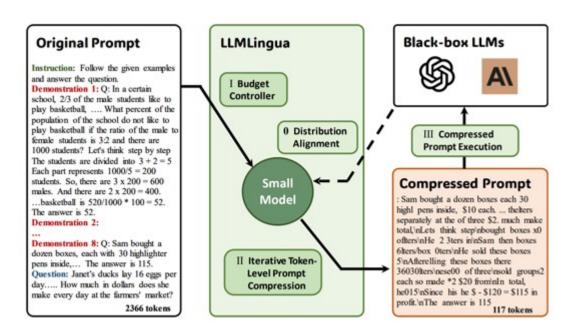
$$I_p = \text{np.percentile}([I(u_0), .., I(u_k)], p)$$

4. Keep units above the threshold

$$C' = U_i \mid I(U_i) \ge I_p, 1 \le i \le n$$

LLMLingua

Compress the prompt with an small language model, with an reduction rate up to 20x.

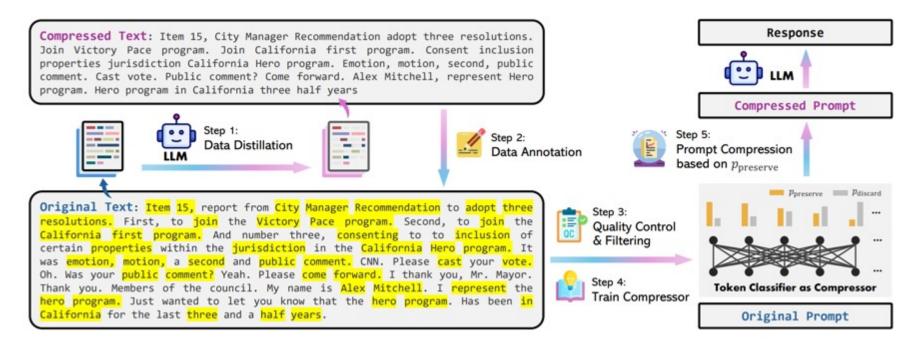


- 0. Distribution alignment
- Instruction tuning of small LLM
- 1. Budget controller
- Calculate the token-wise perplexity and sort in the descending order
- 2. Iterative prompt compression

Group tokens into segments to ensure their dependency, and compute segment-wise perplexities

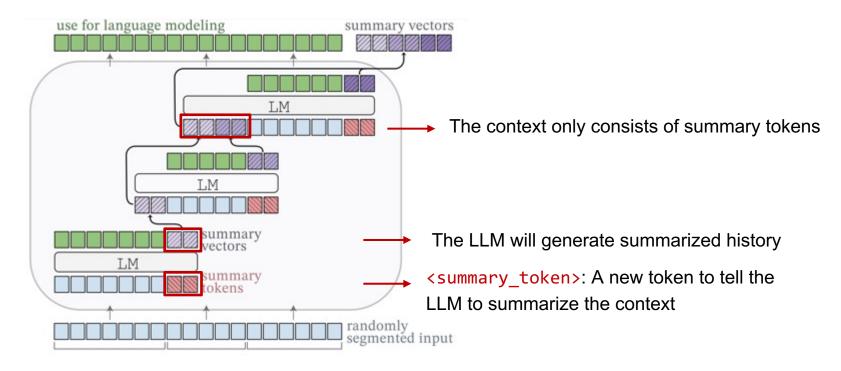
LLMLingua-2

A better way to construct training data for the compressor: instructed by GPT-4



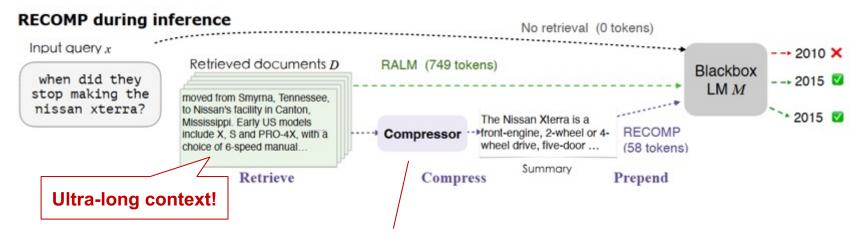
AutoCompressor

The LLM learns to summarize history context given the instruction <summary_token>



RECOMP

- Prompt compression for RAG systems
 - Otherwise, the retrieved documents can be extremely long

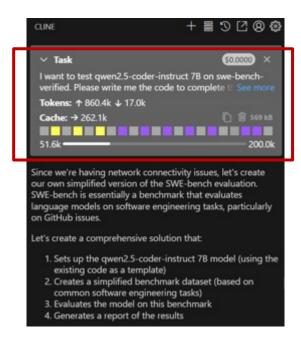


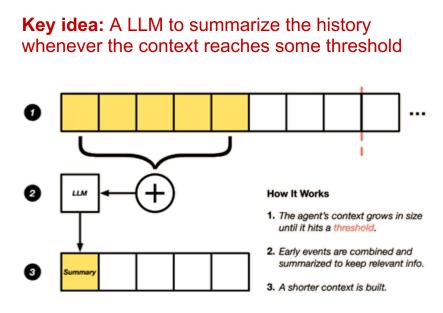
Abstractive Compressor: summarize text from the documents

Extractive Compressor: extract text from the documents

Context Folding for Agentic AI

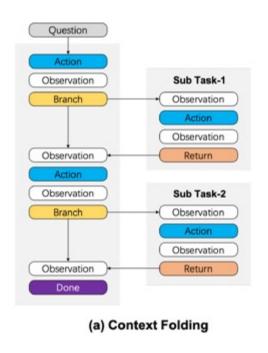
- Context Folding: prompt compression for Agentic Al
 - Deep research, coding agent, etc.

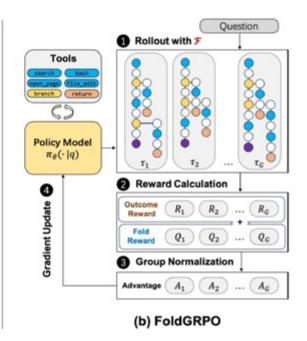




Context Folding for Agentic AI

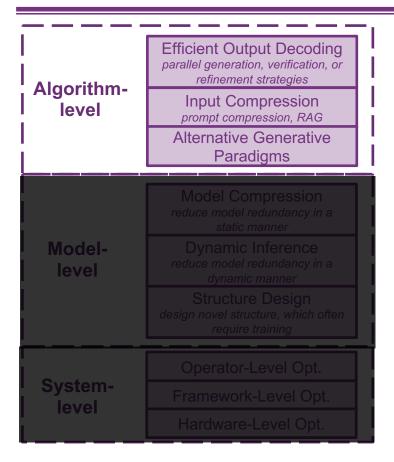
Context Folding: prompt compression for Agentic Al





- Learn to create, solve and summarize sub-tasks
- 2. Reduce the output context from sub-task result
- 3. SFT/RL training to enable the ability to fold context
- Context reduction 10x without accuracy drop on SWE & BrowseComp

Menu of Techniques



Efficient Output Decoding

- Speculative Decoding
- Jacobi Decoding
- Agentic Generation

Input Compression

Input Compression

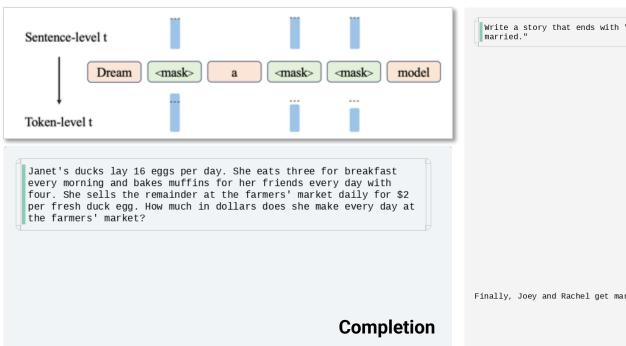
Alternative Generative Paradigms

Diffusion for Text

Diffusion Language Models

- Any other solutions? —— Don't use the autoregressive model?

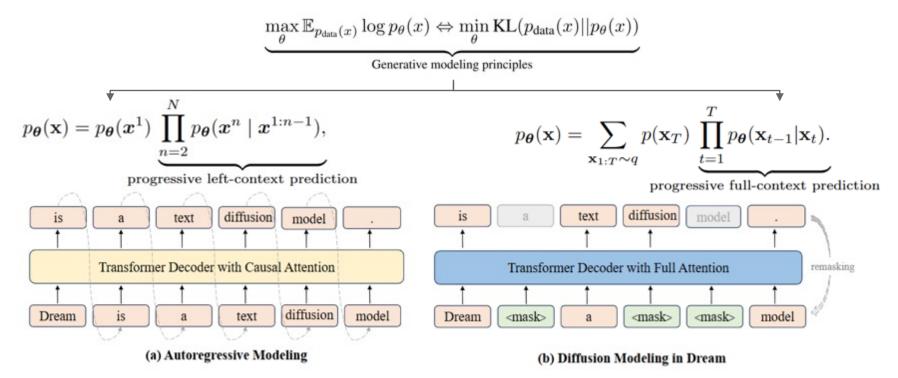
—— To fully parallelized input & output, diffusion language models



Write a story that ends with "Finally, Joey and Rachel get Finally, Joey and Rachel get married. Infiliing

Diffusion Language Models

Autoregressive modeling v.s. Diffusion modeling

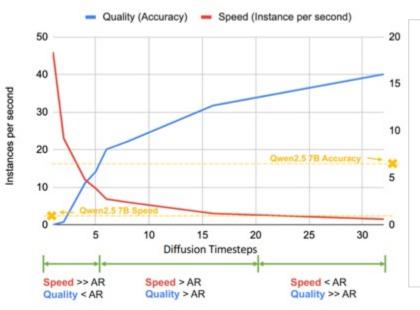


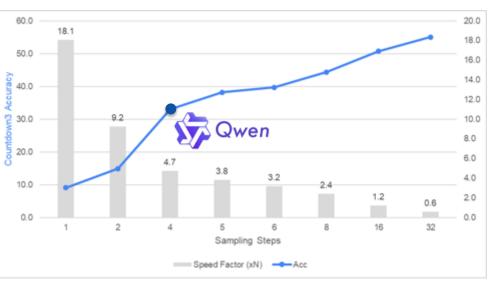
^[1] Jiacheng Ye, et al. Dream 7B: Diffusion Large Language Models. arXiv preprint, 2508.15487.

^[2] Shen Nie, et al. Large Language Diffusion Models. NeurIPS 2025.

Diffusion Language Models

Inference latency of diffusion models





Contents

1 Background

Model-Level Optimization

2 Preliminary

- 6 System-Level Optimization
- 3 Problem Definition & Conceptual Analysis
- Algo-Level Optimization

4 Practical Pipeline

8 Conclusion

Tutorial Review

 Motivation: Scaling up model / data / computation based on Transformer is the mainstream and effective pathway for stronger generative AI till now. The scaling up of the model size and input & output cause efficiency issues.

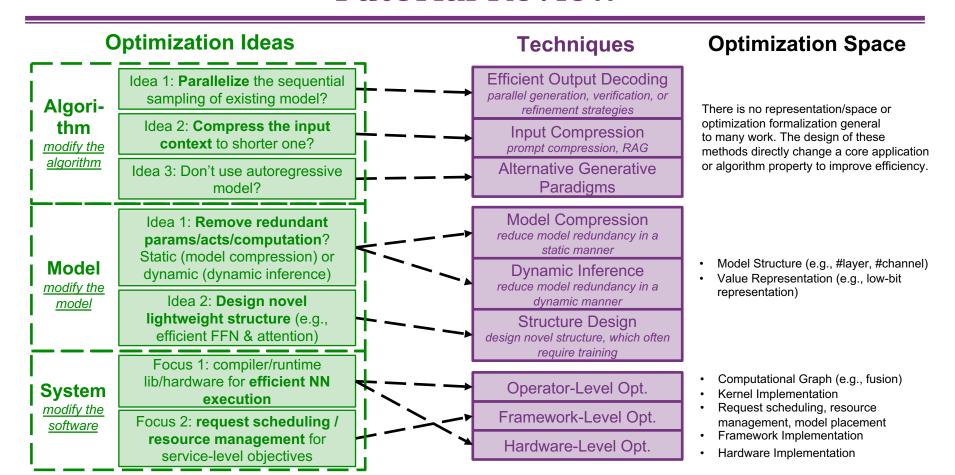
Preliminary:

- Most LLMs use autoregressive model as the generative modeling method, the transformer architecture, in which the attention operation is a core mechanism.
- We introduce basic concepts of software, hardware system, device, chip, microarchitecture, and the interface between software & hardware instruction.
- Al inference is seen as forwarding data on a computational graph, where each node represents a single operator, edge represents dependency. Operators are translated to instructions. Hardware execute instructions.

Tutorial Review

- **Problem Definition**: Usually, latency, memory, energy consumption and throughput will be the ultimate objective or constraint on "efficiency". In the meantime, the intelligence level of AI needs to be retained.
 - **Measured metrics** are actually tested on platform (thus is platform-related), and directly correspond to the objectives / constraints.
 - **Proxy metrics** are estimated with only model specification. In practice, they are useful in diagnosis of the bottleneck and estimation of measured metrics.
- **Practical Pipeline**: We can estimate bottleneck modules and overall objectives, whether each module is compute or memory bound using some simple method (e.g., roofline model), then we can actually profile them (NVIDIA GPU: Nsight system & compute). Finally, we design method accordingly.

Tutorial Review



Future Directions

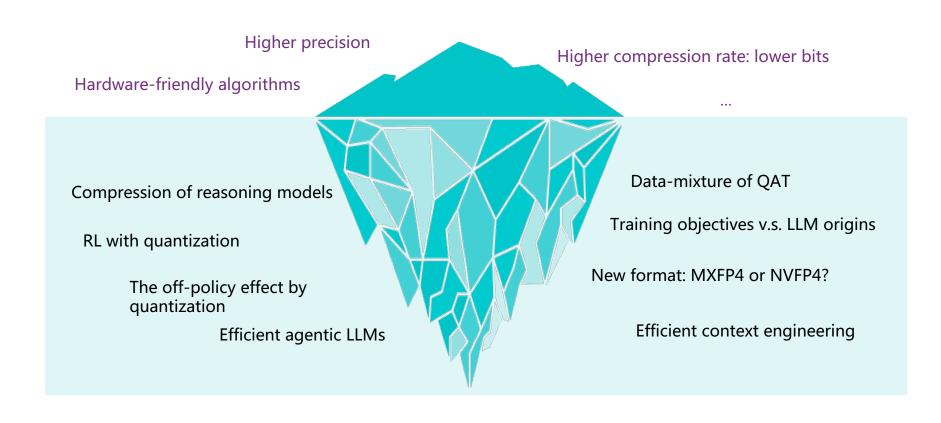
Application Requirements

- Long context (e.g., for complex reasoning)
- Multi-modality input/output
- Multi-model agentic application

Recent Active Directions

- Algorithm-level
 - Agentic generation (multi-agent collaboration)
 - Agent context engineering
- Model-level
 - Latent reasoning
 - Efficient architecture design
- System-level
 - Agent infrastructure

Future Directions



Thank You!

Xuefei Ning¹, Guohao Dai^{2,4}, Haoli Bai³, Lu Hou³, Yu Wang¹, Qun Liu³

¹Tsinghua University ²Shanghai Jiao Tong University ³Huawei ⁴Infinigence-Al

Tutorial Website

https://haolibai.github.io/emnlp-2025-tutorial-efficiency/

