# Mistral 7B

**Name**

박제현

**NLP**

2024/05/28

TRAIN AND TEST
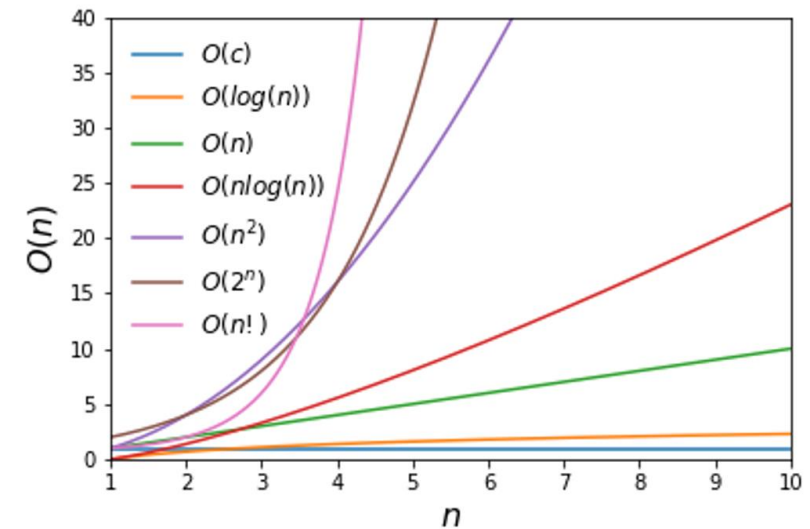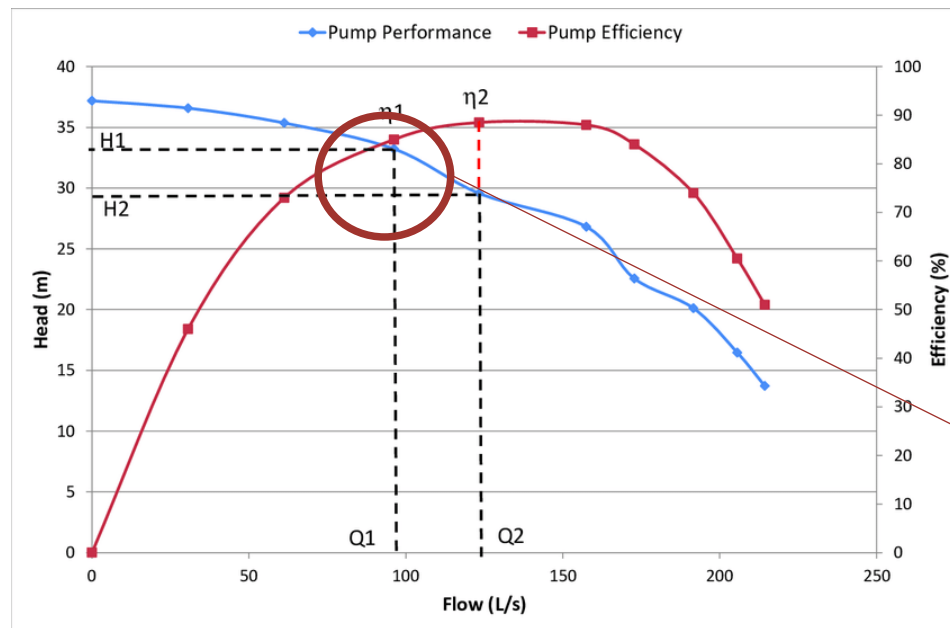
# Contents

- **Introduction**
    - **Balancing Performance and Efficiency in NLP Models**
    - **Recent Appoaches – LLMs Optimization**
    - **Technical Innovations of Mistral 7B**


- **Architectural details**
    - **Sliding Window Attention**
    - **Rolling Buffer Cache**
    - **Pre-fill and Chunking**


- **Main resilts & Analysis**

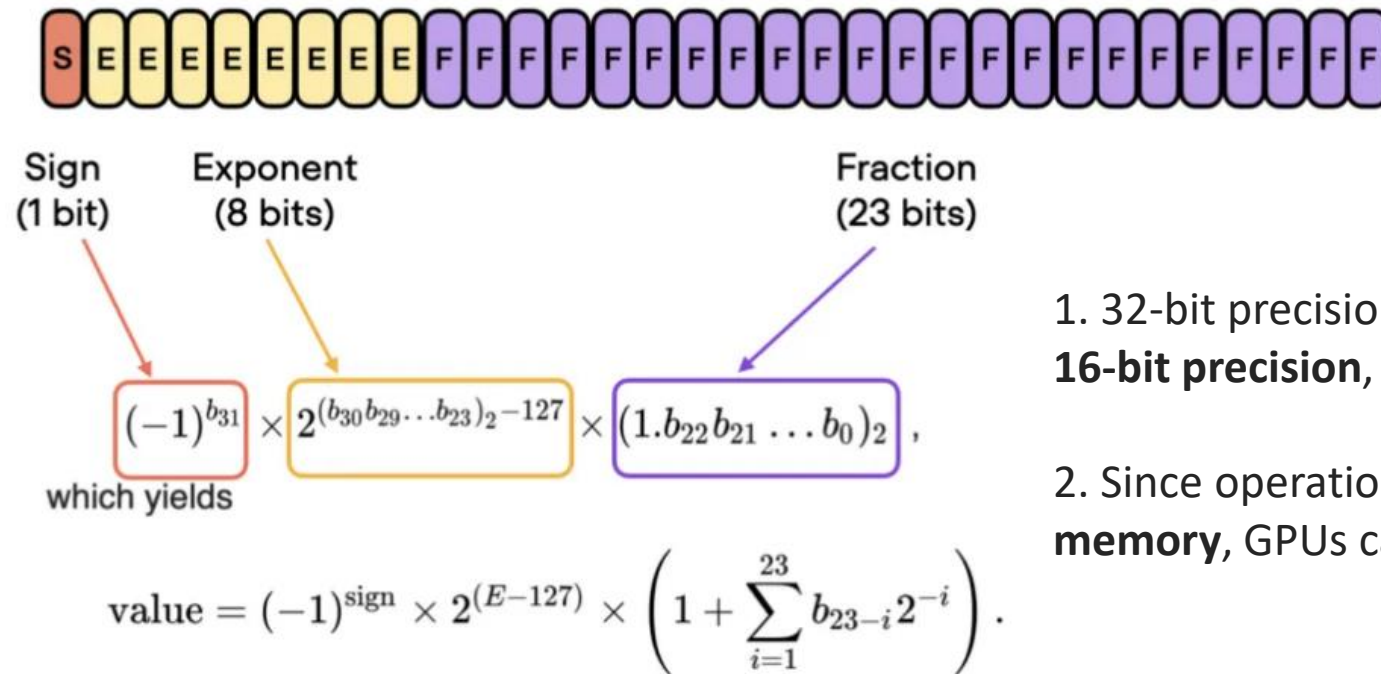# Balancing Performance and Efficiency in NLP Models: Innovations of Mistral 7B

**Performance**



**Model Size & Performance**
**Performance & 1 / Efficiency**

- **Lower Computational costs!**
- **Less Inference latency! -> LoRA**

# Recent Appoaches – LLMs Optimization

**Floating point**

float 32



Sign
(1 bit)

Exponent
(8 bits)

Fraction
(23 bits)

$$(-1)^{b_{31}} \times 2^{(b_{30}b_{29}\ldots b_{23})_2 - 127} \times (1.b_{22}b_{21}\ldots b_0)_2 ,$$

which yields

$$\text{value} = (-1)^{\text{sign}} \times 2^{(E-127)} \times \left(1 + \sum_{i=1}^{23} b_{23-i} 2^{-i}\right).$$

1. 32-bit precision requires twice as much GPU memory as **16-bit precision**, allowing more efficient use of GPU memory.

2. Since operations on **lower precision tensors require less memory**, GPUs can process them more quickly.

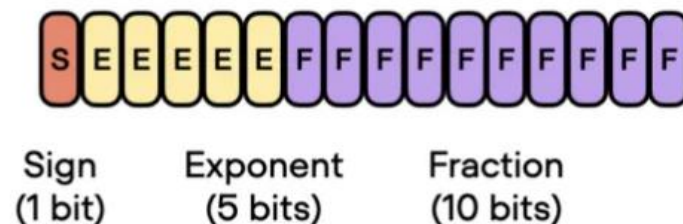# Recent Appoaches – LLMs Optimization

**Don't transfer** all parameters and operations to 16-bit floats. Instead, we **switch between** 32-bit and 16-bit operations during training
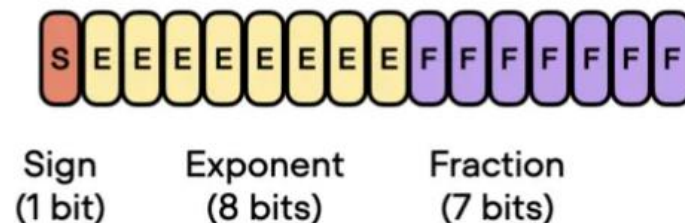
# Recent Appoaches – LLMs Optimization

float 16 ("half" precision)



Sign
(1 bit)

Exponent
(5 bits)

Fraction
(10 bits)

bfloat 16 ("brain" floating point, more "dynamic range" like float 32)


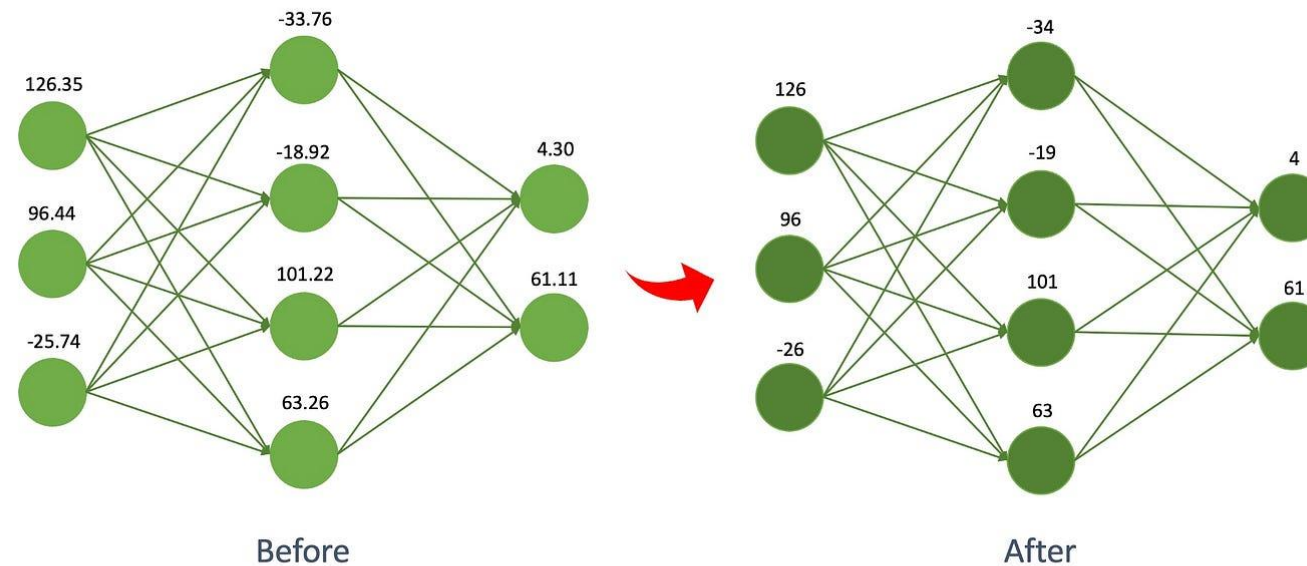
Sign
(1 bit)

Exponent
(8 bits)

Fraction
(7 bits)

**Google developed this format** for machine learning and deep learning applications, particularly in their **Tensor Processing Units** (TPUs). While bfloat16 was originally developed for TPUs, this format is now supported by **several NVIDIA GPUs**.

# Recent Appoaches – LLMs Optimization

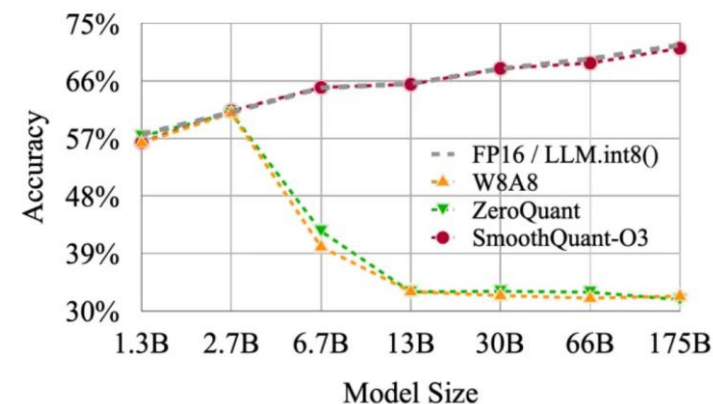Master the Art of Quantization: A Practical Guide
with TensorFlow and PyTorch

Quantization converts the model weights from **float32** to **low-bit integer** representations, for example, **8-bit integers** (and, recently, even **4-bit integers**).

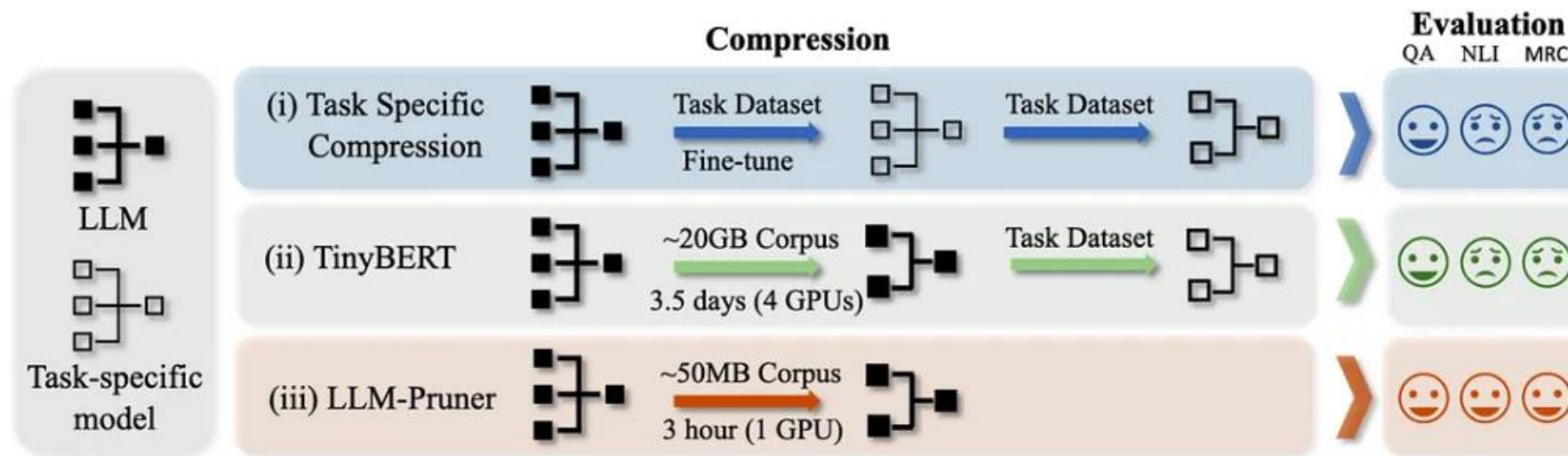# Recent Appoaches – LLMs Optimization

## Quantization

**1. Post-Training Quantization (PTQ)**: A model is **first trained to converge, then we convert its weights** to a lower precision without more training. It is usually quite cheap to implement in comparison to training.

**2. Quantization-Aware Training (QAT)**: Quantization is applied **during pre-training or further fine-tuning**. QAT can perform better but requires extra computation resources and access to representative training data

**- But Both can lower Performance**
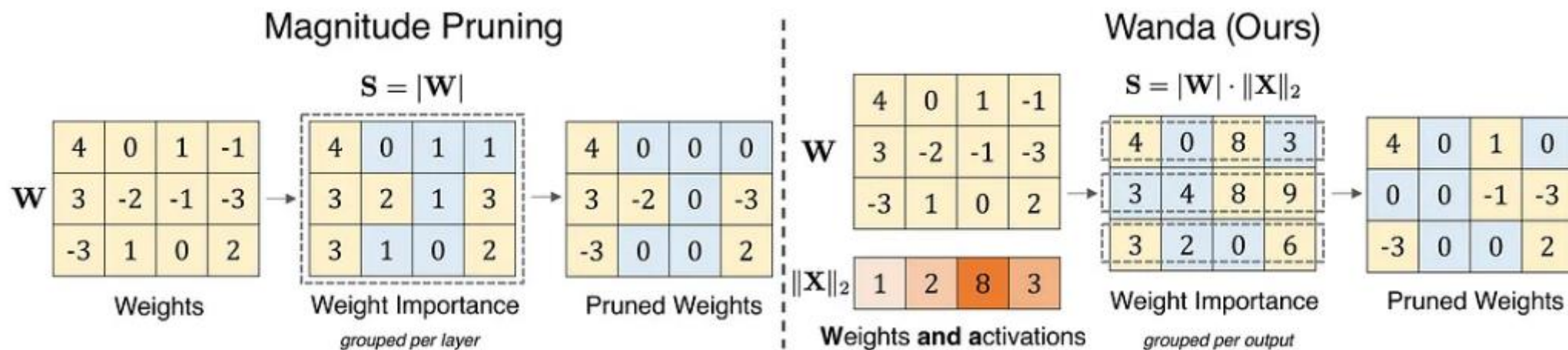
# Recent Appoaches – LLMs Optimization

**Pruning**



Adopts **structural pruning** that selectively removes **non-critical coupled structures** based on **gradient information**, maximally preserving most of the LLM's functionality.
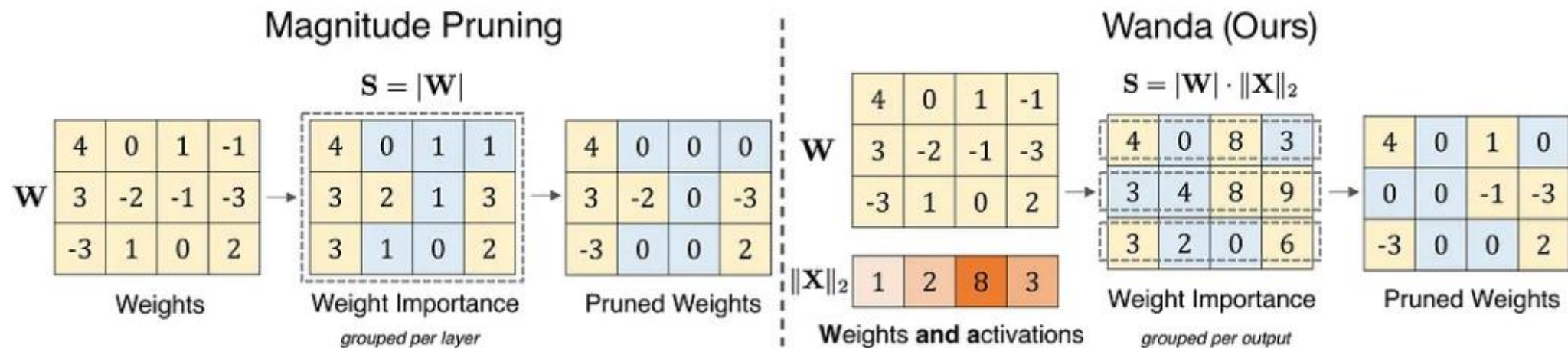
# Recent Appoaches – LLMs Optimization

**Pruning**



Magnitude Pruning — $S = |\mathbf{W}|$

Wanda (Ours) — $S = |\mathbf{W}| \cdot \|\mathbf{X}\|_2$

Compared to **magnitude pruning** which removes weights **solely based on their magnitudes**, Wanda removes weights on a **per-output basis** by the **product of weight magnitudes** and **input activation norms**.

# Balancing Performance and Efficiency in NLP Models: Innovations of Mistral 7B

**Pruning**



Magnitude Pruning — Wanda (Ours)

Compared to **magnitude pruning** which removes weights **solely based on their magnitudes**, Wanda removes weights on a **per-output basis** by the **product of weight magnitudes** and **input activation norms**.

# Technical Innovations of Mistral 7B

**Grouped-Query Attention (GQA):**
- Acceleration of inference speed.
- Reduction of memory requirements and increase in batch size.

**Sliding Window Attention (SWA):**
- Improved efficiency in handling long sequences.
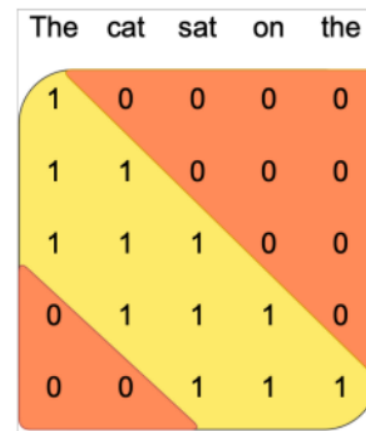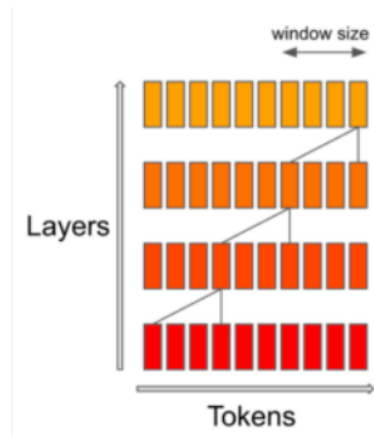- Reduction of computational costs.
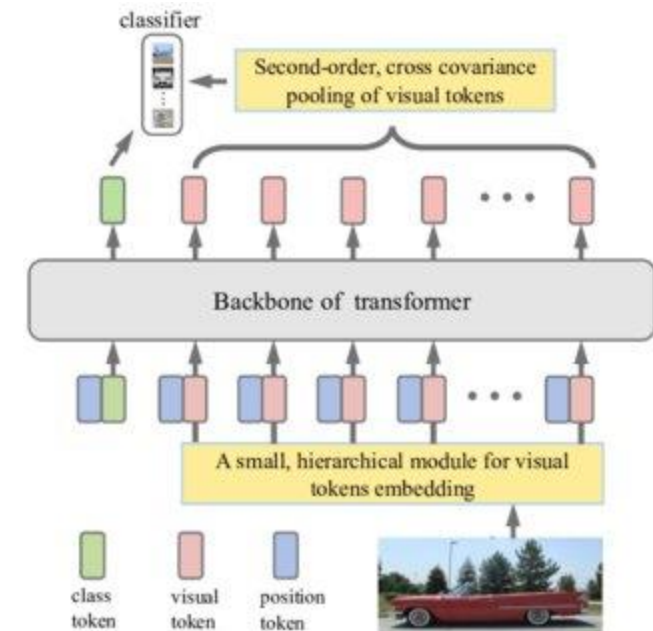
# Sliding Window Attention

Sliding Window Attention

# Sliding Window Attention

**Vanilla Attention:**
Computational complexity: **O(n^2)**, where n is the <u>sequence length</u>.
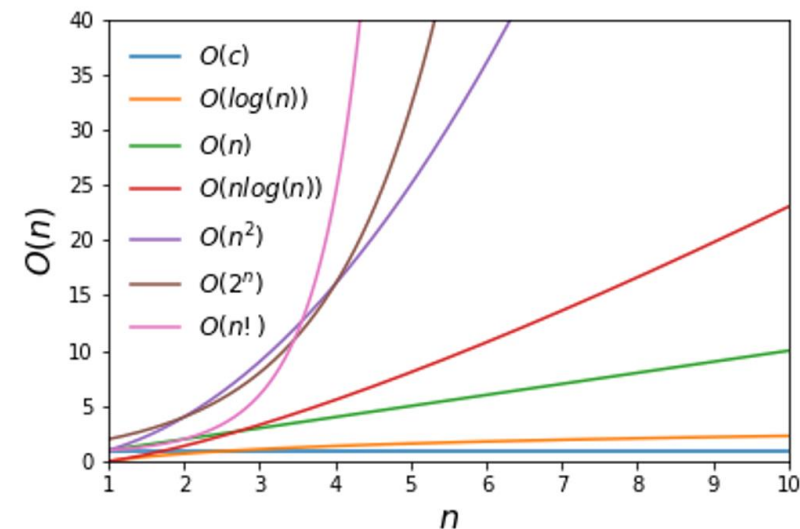Memory usage: Increases **linearly** with the number of tokens.
Doubling the sequence length **quadruples the computation**.

**Sliding Window Attention (SWA):**
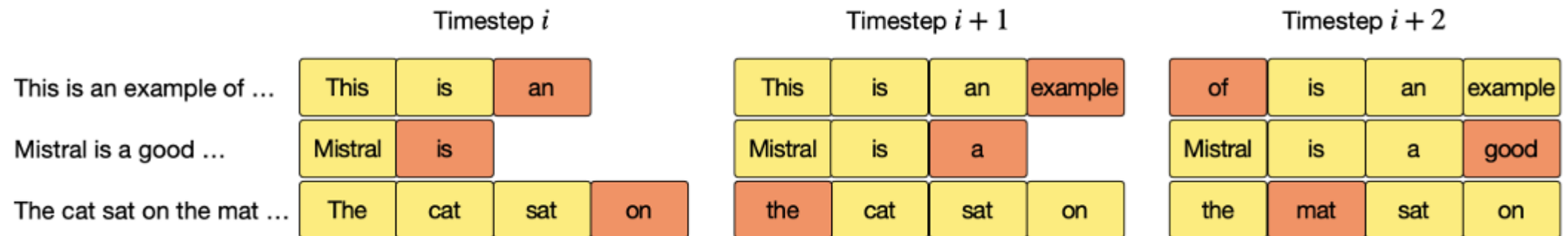Computational complexity: **O(n * W)**, where W is the <u>window size</u>.
Memory usage: Increases **linearly** with W.
The computational load is **constrained by the window size**, even as the sequence length increases.

TRAIN AND TEST

# Rolling Buffer Cache

Rolling Buffer Cache



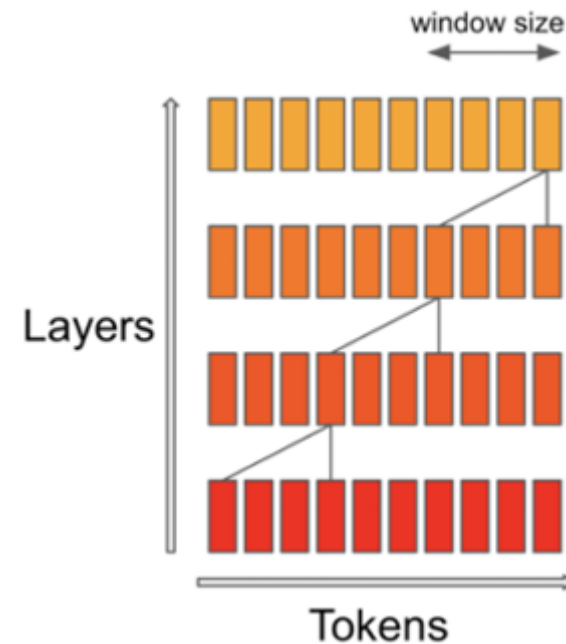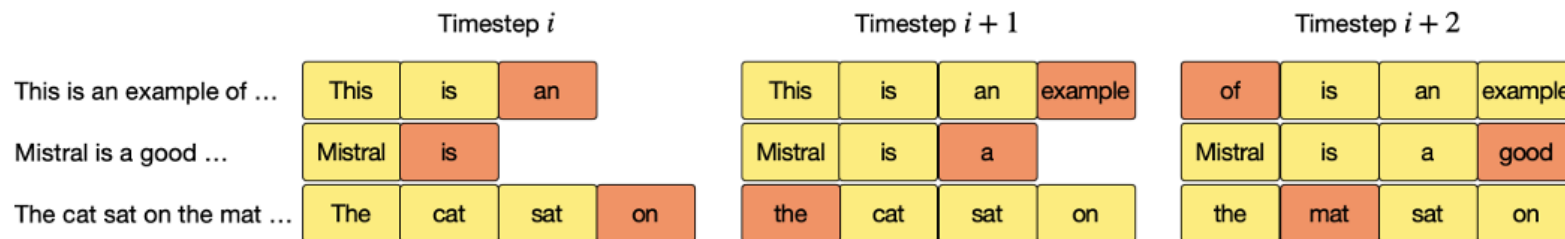The attention span is **limited to a fixed size** $W$.
This means each token will only attend to $W$ **preceding tokens**.

The cache is organized as a rolling buffer with a fixed size $W$.
Keys and values for each timestep $i$ are stored in the cache position $i \bmod W$. (i = 0, 1, 2, 3, 4 . . .)

# Rolling Buffer Cache

# Pre-fill and Chunking

# Main resilts & Analysis
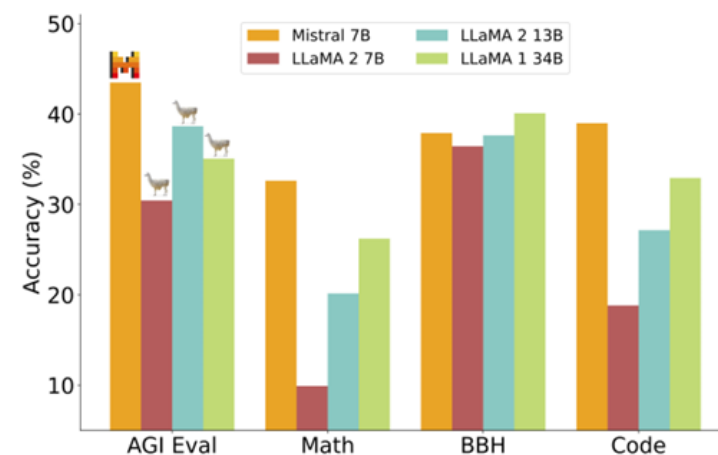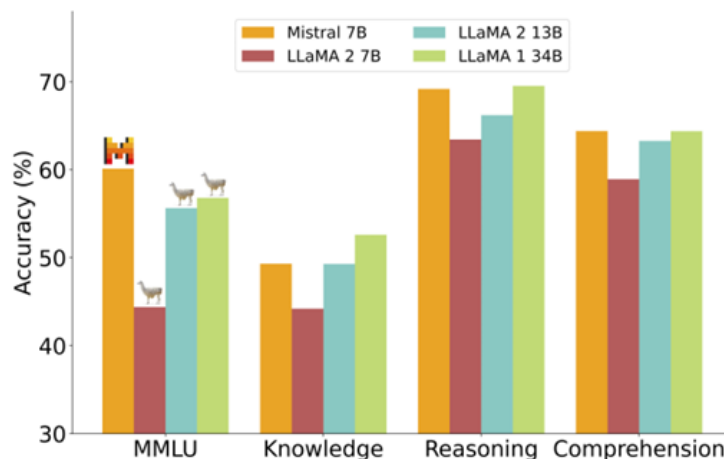
| Model | Modality | MMLU | HellaSwag | WinoG | PIQA | Arc-e | Arc-c | NQ | TriviaQA | HumanEval | MBPP | MATH | GSM8K |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LLaMA 2 7B | Pretrained | 44.4% | 77.1% | 69.5% | 77.9% | 68.7% | 43.2% | 24.7% | 63.8% | 11.6% | 26.1% | 3.9% | 16.0% |
| LLaMA 2 13B | Pretrained | 55.6% | **80.7%** | 72.9% | 80.8% | 75.2% | 48.8% | **29.0%** | **69.6%** | 18.9% | 35.4% | 6.0% | 34.3% |
| Code-Llama 7B | Finetuned | 36.9% | 62.9% | 62.3% | 72.8% | 59.4% | 34.5% | 11.0% | 34.9% | **31.1%** | **52.5%** | 5.2% | 20.8% |
| Mistral 7B | Pretrained | **60.1%** | **81.3%** | **75.3%** | **83.0%** | **80.0%** | **55.5%** | **28.8%** | 69.9% | 30.5% | 47.5% | **13.1%** | **52.2%** |

17

# Main resilts & Analysis

TRAIN AND TEST