# RenderFormer

**Mid-term Project Presentation**

*10th Oct, 2025*

Team 4
Kyaw Ye Thu, Janghyun



*CS580: Computer Graphics*

https://microsoft.github.io/renderformer/

# Outline

- Related Work & their Problems

- Overview of RenderFormer

- Limitations of RenderFormer

- Our Plan to Overcome One Limitation

# Related Work & Their Limitations

# Solving the Rendering Equation

$$L_r(x \to \Theta) = \int_\Psi L(x \leftarrow \Psi) f_r(x, \Psi \to \Theta) \cos \theta_x dw_\Psi$$

**Monte Carlo Estimation of Path Integral**

- P
- M
- Path Guiding
- Multiple Importance Sampling
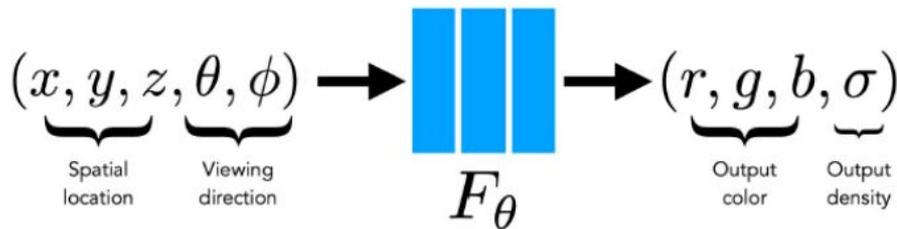- ...

**Recursive**

**Radiosity Techniques (Finite Elements Methods)**

- Pre-computation
- Per-scene Training
- Precomputed Radiance Transfer (PRT)
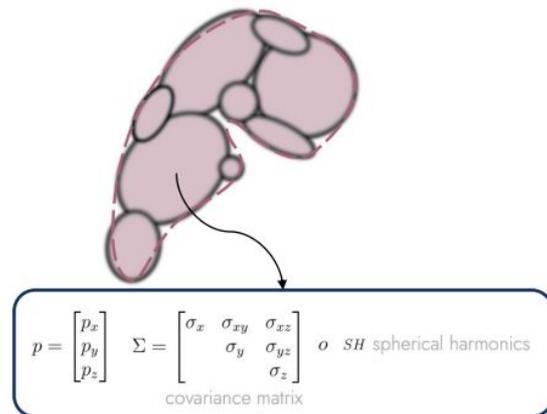- ...

# 3D Inductive Bias (Representation Level)

## NeRF

**Continuous Volumetric Field**: Assumes the scene can be represented as a continuous 5D function.



## 3DGS

**Gaussian Primitives**: Scenes are composed of fuzzy blobs (Gaussians) that can be combined
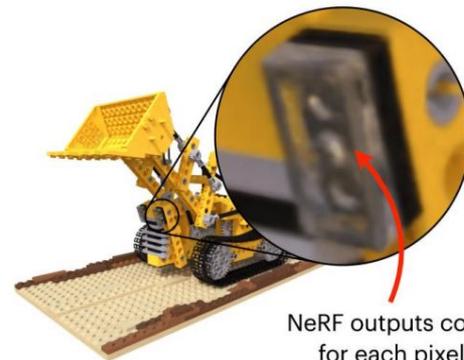
# 3D Inductive Bias

## Why 3D Inductive Biases Are Good

1. **Efficiency**: Reduce the learning space.
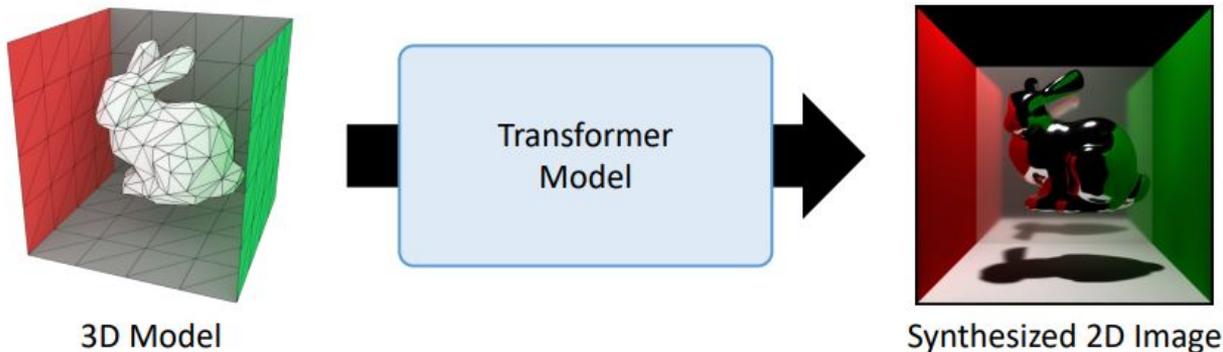2. **Sample efficiency**: With strong priors, models can generalize from fewer examples

## Why 3D Inductive Biases Are Bad

1. **Limited flexibility**: Real-world scenarios may not match the built-in assumptions.
2. **Poor generalization**: A model designed for one type of object/geometry might fail on other types of objects/geometries.



NeRF outputs color for each pixel

# RenderFormer: A Fully Transformer Rendering Pipeline



3D Model

Transformer Model

Synthesized 2D Image
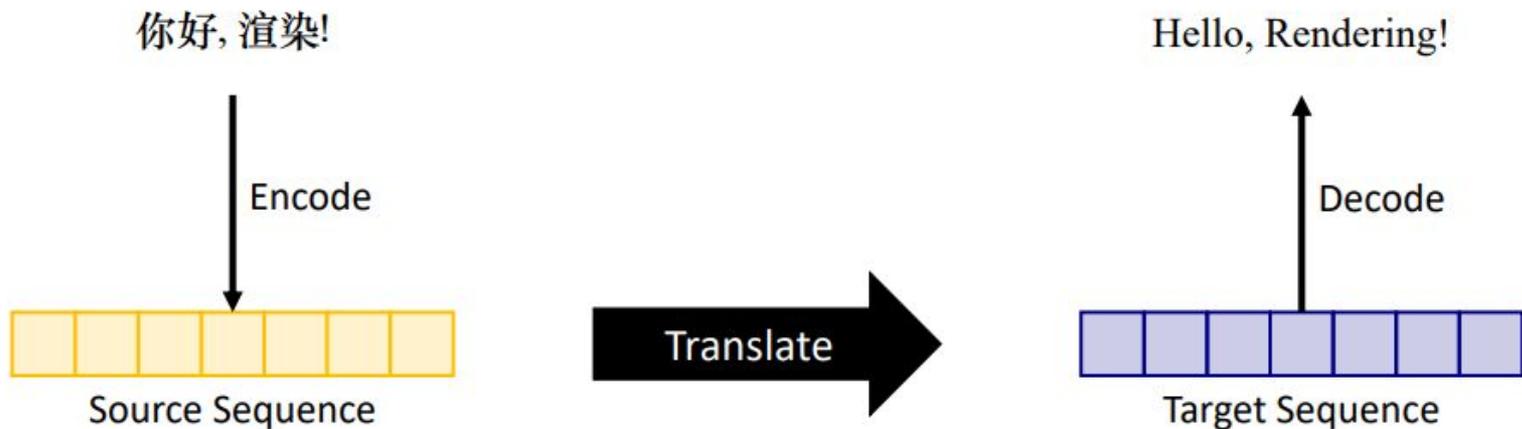
✅ End-to-end with Raw 3D Input

✅ Minimal Prior Constraints

✅ No Per-scene Training

✅ Full Global Illumination Effects

# Overview of RenderFormer

# Sequence-to-Sequence Machine Translation



你好, 渲染!

Encode

Source Sequence

Translate

Hello, Rendering!

Decode

Target Sequence

# Idea: 3D Rendering = Translating 3D to 2D

# Training Data – Template-based Scene Generation



**4,096**
Max #Triangles

**0.01-1.0**
Roughness Range

**8**
Max #Lights

**2.1-2.7**
Light Distance

**30-60**
Camera FOV

**1.5-2.0**
Camera Distance

# Visual Results – Static Scene Rendering

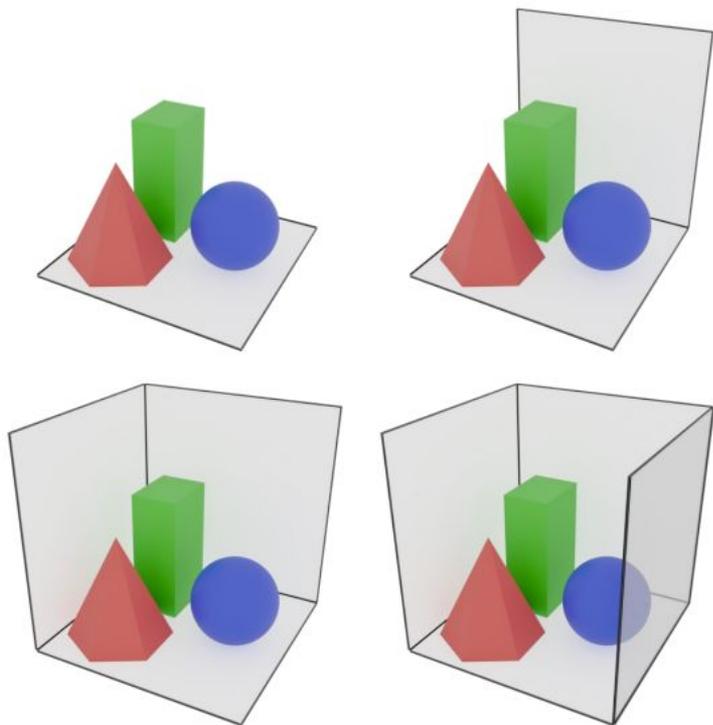# Visual Results – Dynamic Scene Rendering



Cascade Cube Animation
*Tycho Magnetic Anomaly*

Animated Crab
*Bohdan Lvov*

Gyroscope Motion
*reddification*

Animated Character
*mortaleiros*

Marching Cubes Animation
*Tycho Magnetic Anomaly*

Robot Animation
*Gouhadouken*

# Limitations of RenderFormer

# 1. Generalization (# Triangles)



Inside Training Range | Outside Training Range (Extrapolation)

3K | 6K | 11K | 23K | 45K

# 1. Generalization (# light resources)

# 2. Time complexity of Attention



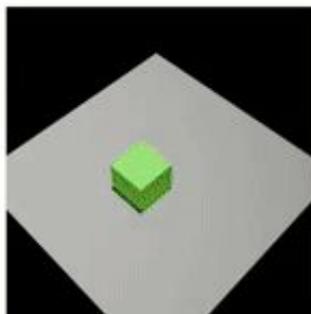| Layer Type | Complexity per Layer | Sequential Operations | Maximum Path Length |
|---|---|---|---|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ | $O(n)$ |
| Convolutional | $O(k \cdot n \cdot d^2)$ | $O(1)$ | $O(log_k(n))$ |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$ | $O(1)$ | $O(n/r)$ |

n = sequence length

d = dimension

k = kernel size of convolutions

A. Vaswani et al., "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, vol. 30

# Time complexity of Renderformer



- The view-dependent stage's time complexity is roughly by O(#Bundle^2)
- The view-independent stage's time complexity is roughly by O($\#tri$^2)
- Because of this, limit the total number of triangles in scenes to 4,096

C. Zeng, Y. Dong, P. Peers, H. Wu, and X. Tong, "RenderFormer: Transformer-based Neural Rendering of Triangle Meshes with Global Illumination," in *Proceedings of the ACM SIGGRAPH 2025 Conference*

# Actually, they used Flash Attention-2



FlashAttention

- They used FA-2 which can help to reduce calculation time and memory usage.
- Flash attention is attention algorithm that tiles softmax-attention to on-chip SRAM.

T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré, "FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness," in *Proceedings of the 36th International Conference on Neural Information Processing Systems* (NeurIPS), 2022

*CS580: Computer Graphics*

# How Flash Attention work?

Standard Attention Implementation

Load Q, K

Write S

$$S = QK^T$$

Memory (HBM) ← Compute

Load S

Write P

P = softmax(s)

Load P, V

$$O = PV$$

Write O

---

**Algorithm 0** Standard Attention Implementation

**Require:** Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.

1: Load $\mathbf{Q}, \mathbf{K}$ by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$, write $\mathbf{S}$ to HBM.
2: Read $\mathbf{S}$ from HBM, compute $\mathbf{P} = \text{softmax}(\mathbf{S})$, write $\mathbf{P}$ to HBM.
3: Load $\mathbf{P}$ and $\mathbf{V}$ by blocks from HBM, compute $\mathbf{O} = \mathbf{PV}$, write $\mathbf{O}$ to HBM.
4: Return $\mathbf{O}$.

T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré, "FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness," in *Proceedings of the 36th International Conference on Neural Information Processing Systems* (NeurIPS), 2022

*CS580: Computer Graphics*

# How Flash Attention work?



Flash Attention

Load $K_j \; V_j$

Load $Q_i \; O_i \; l_i \; m_i$

Kernel operations fused together, reducing reads & writes

Memory (HBM)

Compute

$S_{ij} = Q_i K_j^T$
$\tilde{m} = $ rowmax of S
$P = \exp(s - m)$
$l = $ rowsum of P
$m = \max(\tilde{m}_{ij}, m)$
calculate O from l & m

Write $O_i \; l_i \; v_i$

Initialize O, l and m matrices with zeroes. m and l are used to calculate cumulative softmax. Divide Q, K, V into blocks (due to SRAM's memory limits) and iterate over them, for i is row & j is column.

Outer Loop

$K^T : d \times N$

Copy Block to SRAM

Outer Loop

$Q : N \times d$

$V : N \times d$

$QK^T : N \times N$

Copy

Compute Block on SRAM

Copy

Inner Loop

Inner Loop

Outer Loop

Inner Loop

Output to HBM

$sm(QK^T)V : N \times d$

Inner Loop

FlashAttention

Attention on GPT-2

Matmul

Dropout

Softmax

Mask

Matmul

Fused Kernel

PyTorch

FlashAttention

T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré, "FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness," in *Proceedings of the 36th International Conference on Neural Information Processing Systems* (NeurIPS), 2022

*CS580: Computer Graphics*

# Improvement of Flash Attention-2



(a) FLASHATTENTION
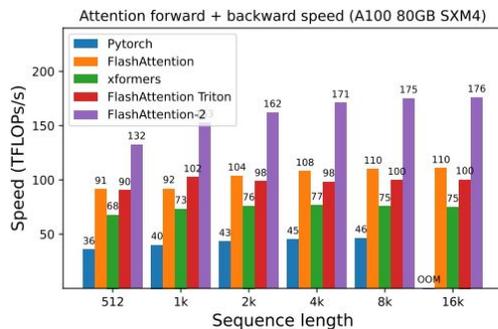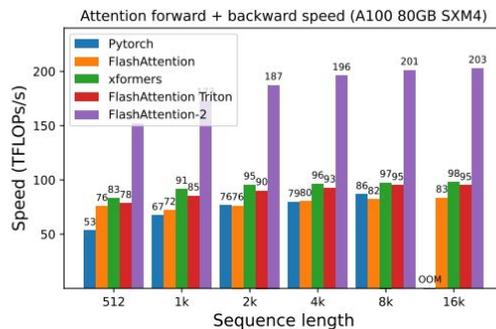
(b) FLASHATTENTION-2

- Flash attention-2 make more fast speed by modifying the process to split Q and share K and V, eliminating the need for continuous synchronization of intermediate calculation results, instead of splitting K and V.
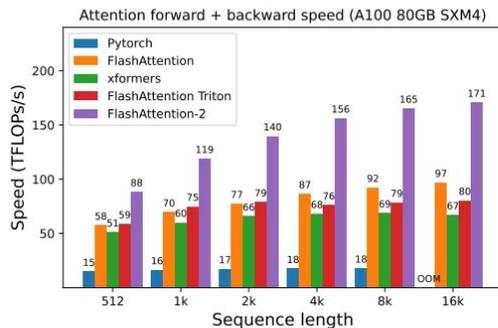
T. Dao, "FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning," in *Proc. Int. Conf. Learn. Representations* (ICLR), 2024

*CS580: Computer Graphics*

# Result of Flash Attention-2



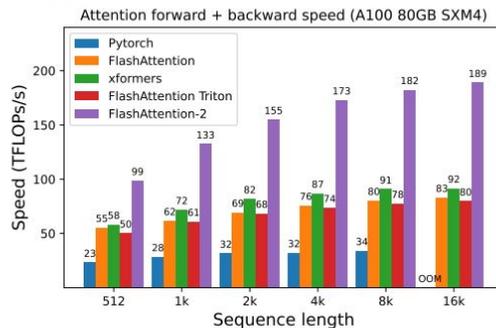Figure 4: Attention forward + backward speed on A100 GPU

(a) Without causal mask, head dimension 64

(b) Without causal mask, head dimension 128

(c) With causal mask, head dimension 64

(d) With causal mask, head dimension 128

T. Dao, "FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning," in *Proc. Int. Conf. Learn. Representations* (ICLR), 2024

*CS580: Computer Graphics*

# Result of Flash Attention-2



FlashAttention Memory Reduction

T. Dao, "FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning," in *Proc. Int. Conf. Learn. Representations* (ICLR), 2024

*CS580: Computer Graphics*

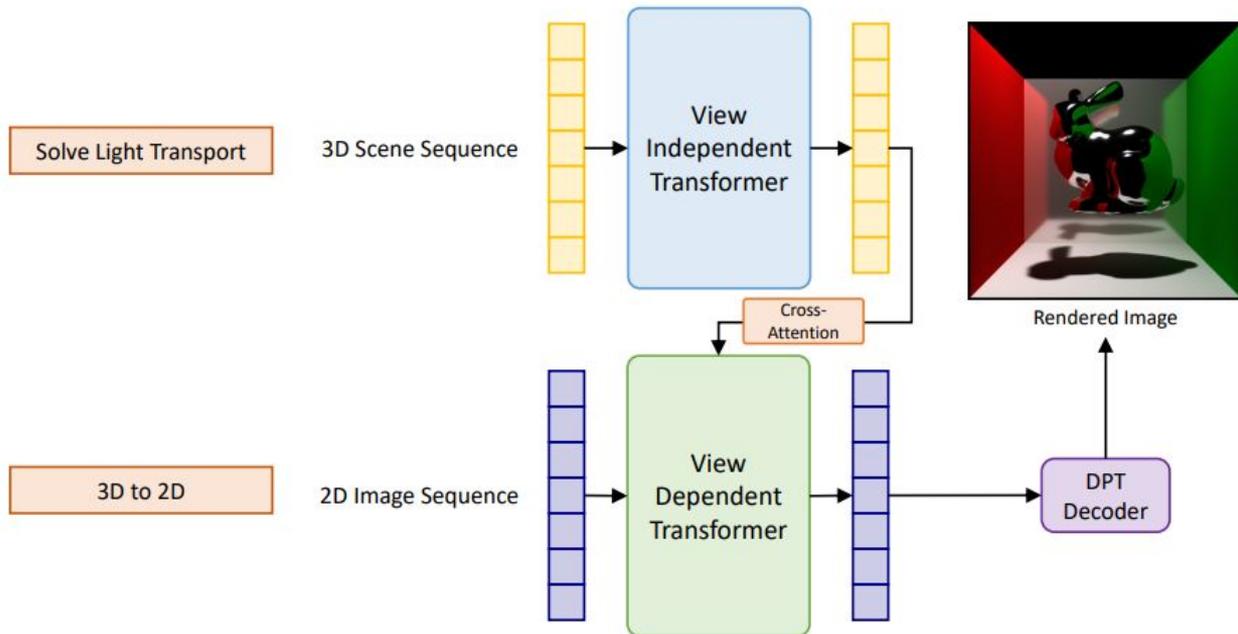# Limitation of Flash Attention-2



- Although they used FA-2 which can help to calculation time, its asymptotic time complexity is same.
- Thus, we need other approach for resolve this issue.

T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré, "FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness," in *Proceedings of the 36th International Conference on Neural Information Processing Systems* (NeurIPS), 2022

*CS580: Computer Graphics*

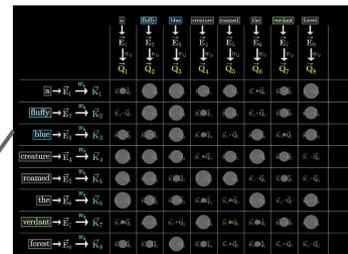# Our Plan

# The Use of Original Transform Architecture

Both stages closely follow the original transformer architecture [Vaswani et al. 2017]



Rendered Image

# How Attention can be Computed Faster

**L x L**

**L x d   d x L     L x d**

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^{\mathsf{T}}}{\sqrt{d_k}}\right) \cdot V$$

L - length of the sequence
d - # dimension of the sequence

# How Attention can be Computed Faster

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^\top}{\sqrt{d_k}}\right) \cdot V$$

**L x L**

**L x d   d x L     L x d**

- The view-independent stage scales roughly by O ($\#tris$^2)
- view-dependent layers scales by O ($\#bundles$^2 + $\#bundles$ × $\#tris$).

# How Attention can be Computed Faster

L x L

L x d   d x L      L x d

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^\mathsf{T}}{\sqrt{d_k}}\right) \cdot V$$

d x L      L x d

d x d
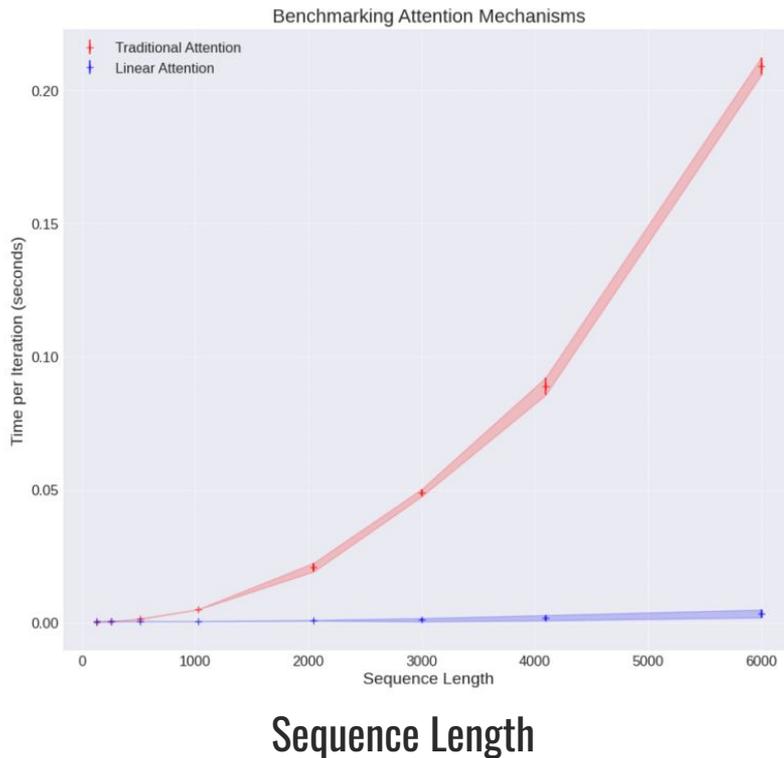
d = 768 (dimension of input sequence)
L = 4,096 (# triangles)

As long as d << L, significantly faster

# Potential Speedup from Linear Attention



Benchmarking Attention Mechanisms

Time per Iteration (seconds)

Sequence Length

1D sequences or 2D images

?

3D

*Figure Source:* https://towardsdatascience.com/linear-attention-is-all-you-need-5fa9c845c1b5/

*CS580: Computer Graphics*

# Our Plan

- Rewrite the **Inference** Part with Linear Attention or Performer

- Compare the Inference Speed and Rendering Result with the original

- Try **training** with Linear Attention or Performer

- Compare the Training Speed and Rendering Result with the original

# Q&A

# Quiz Time