# VR-Pipe

## Streamlining Hardware Graphics Pipeline for Volume Rendering

**Junseo Lee**   Jaisung Kim   Junyong Park   Jaewoong Sim

Seoul National University

# Advance of Graphics Rendering
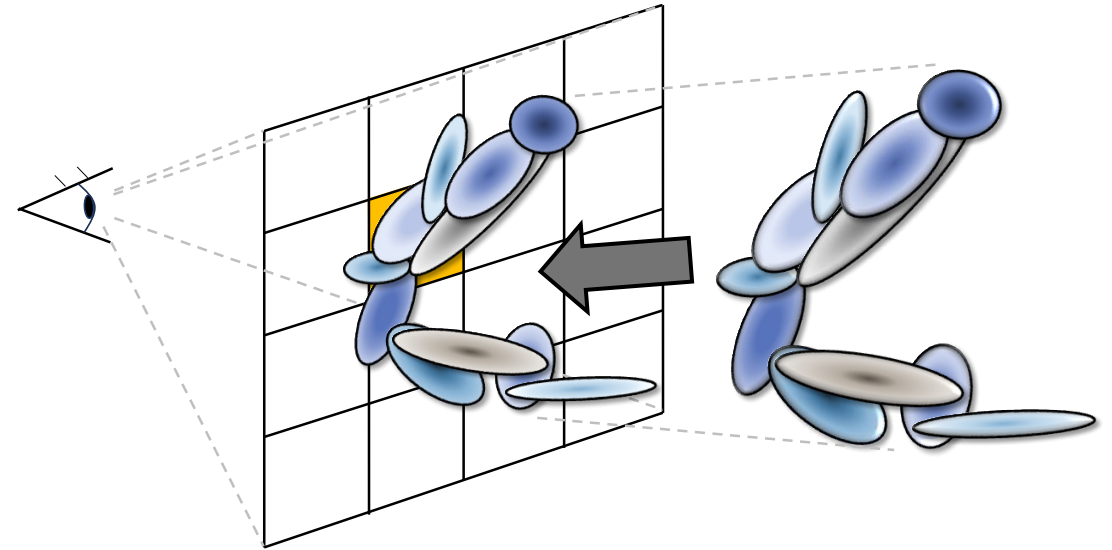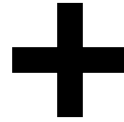
## 3D Gaussian Splatting (3DGS)

# Advance of Graphics Rendering

## 3D Gaussian Splatting (3DGS)

Captured Images



Explicit Representation:
**3D Gaussians**

**Splatting** + Volume Rendering

# Advance of Graphics Rendering

## 3D Gaussian Splatting (3DGS)

Captured Images



Explicit Representation:
**3D Gaussians**

# Advance of Graphics Rendering

## 3D Gaussian Splatting (3DGS)

Captured Images



Explicit Representation:
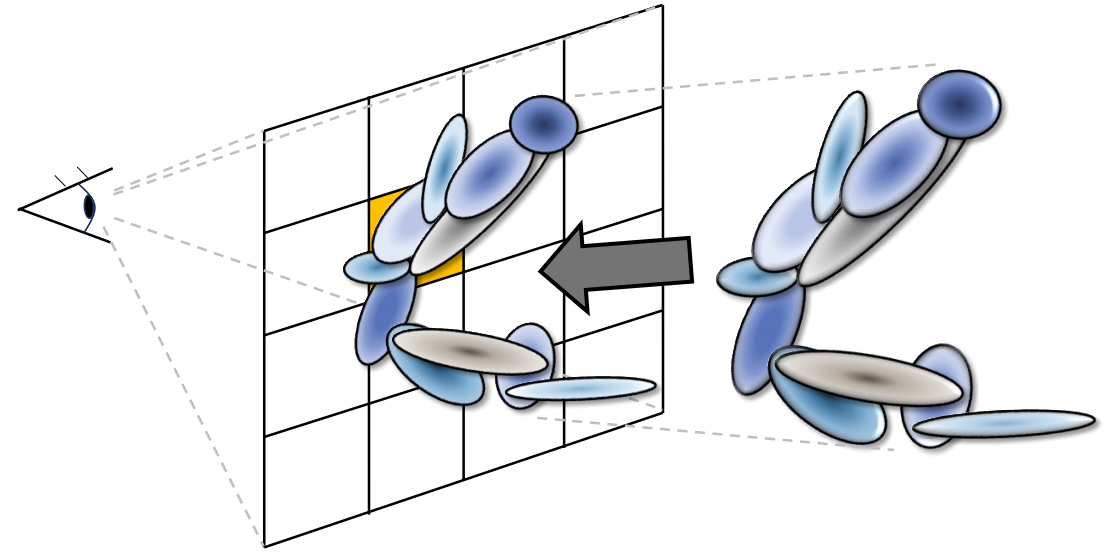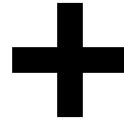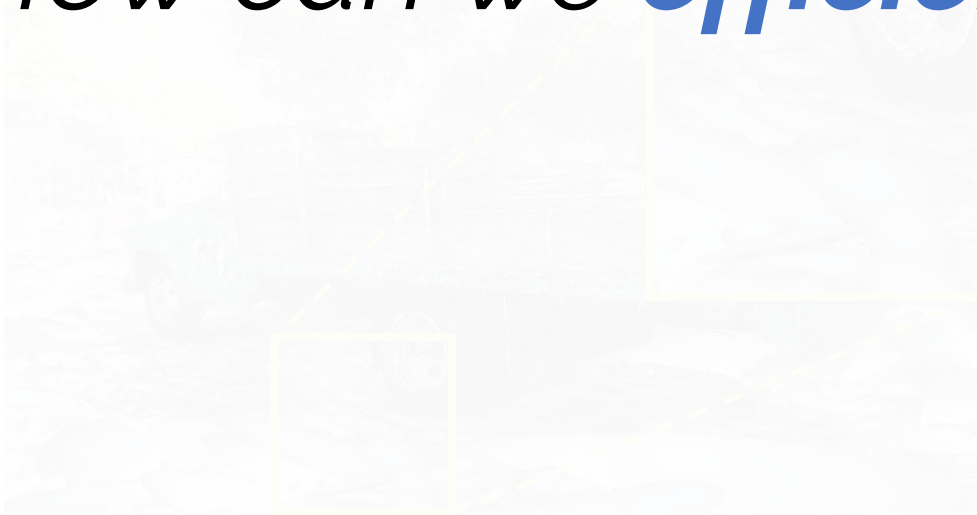**3D Gaussians**

**Splatting** + Volume Rendering

# How can we *efficiently run 3DGS* on a *GPU*?
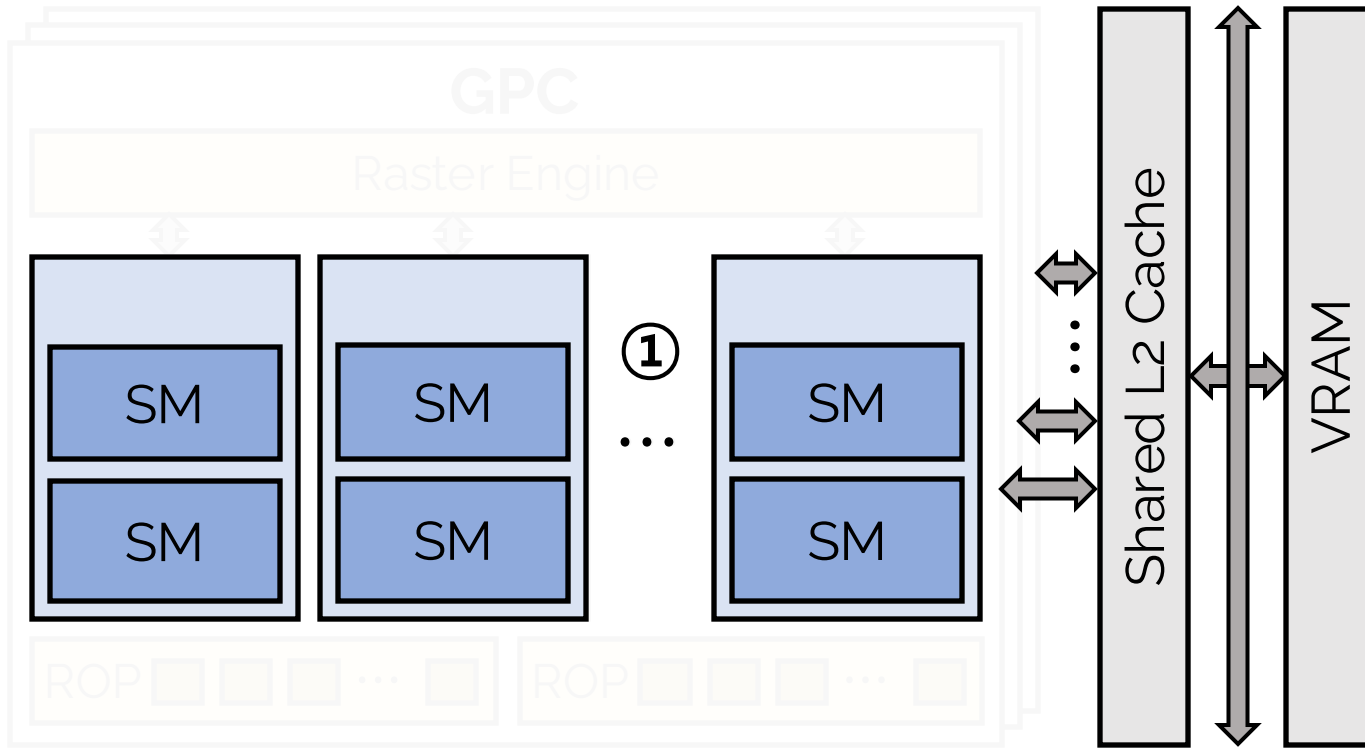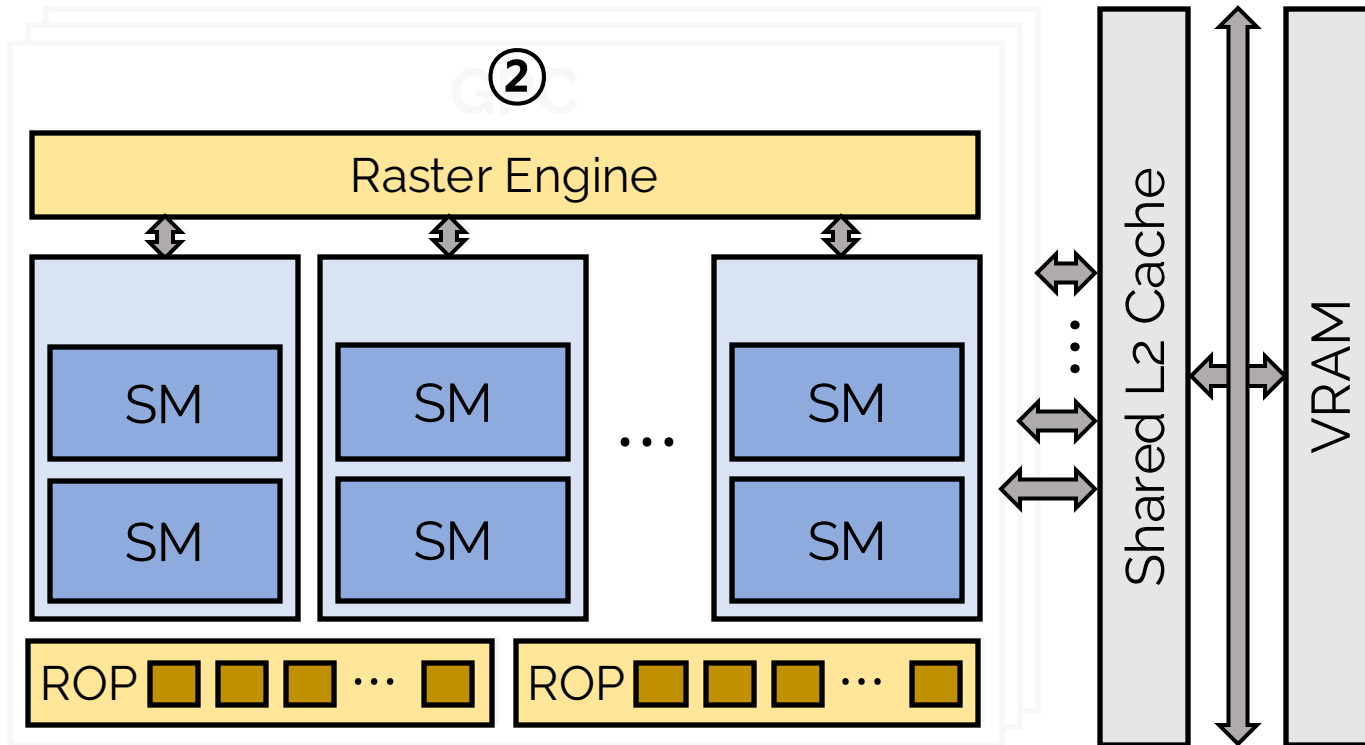
# 3D Gaussian Splatting on a GPU

# 3D Gaussian Splatting on a GPU



## ① SW-based rendering

- Use only **SMs**
- General-purpose computing frameworks (e.g., CUDA, OpenCL)

# 3D Gaussian Splatting on a GPU



① **SW-based rendering**

- Use only **SMs**
- General-purpose computing frameworks (e.g., CUDA, OpenCL)

② **HW-based rendering**

- Use **graphics-specific fixed-function units** w/ **SMs** = **hardware graphics pipeline**
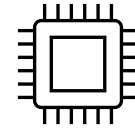- Graphics APIs (e.g., OpenGL, Vulkan)

# Goal of Our Work

## CUDA Optimizations

StopThePop [SIGGRAPH'24]

FlashGS [arXiv'24]

## Specialized Accelerators
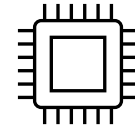
GSCore [ASPLOS'24]

MetaSapiens [ASPLOS'25]

# Goal of Our Work

## CUDA Optimizations

StopThePop [SIGGRAPH'24]
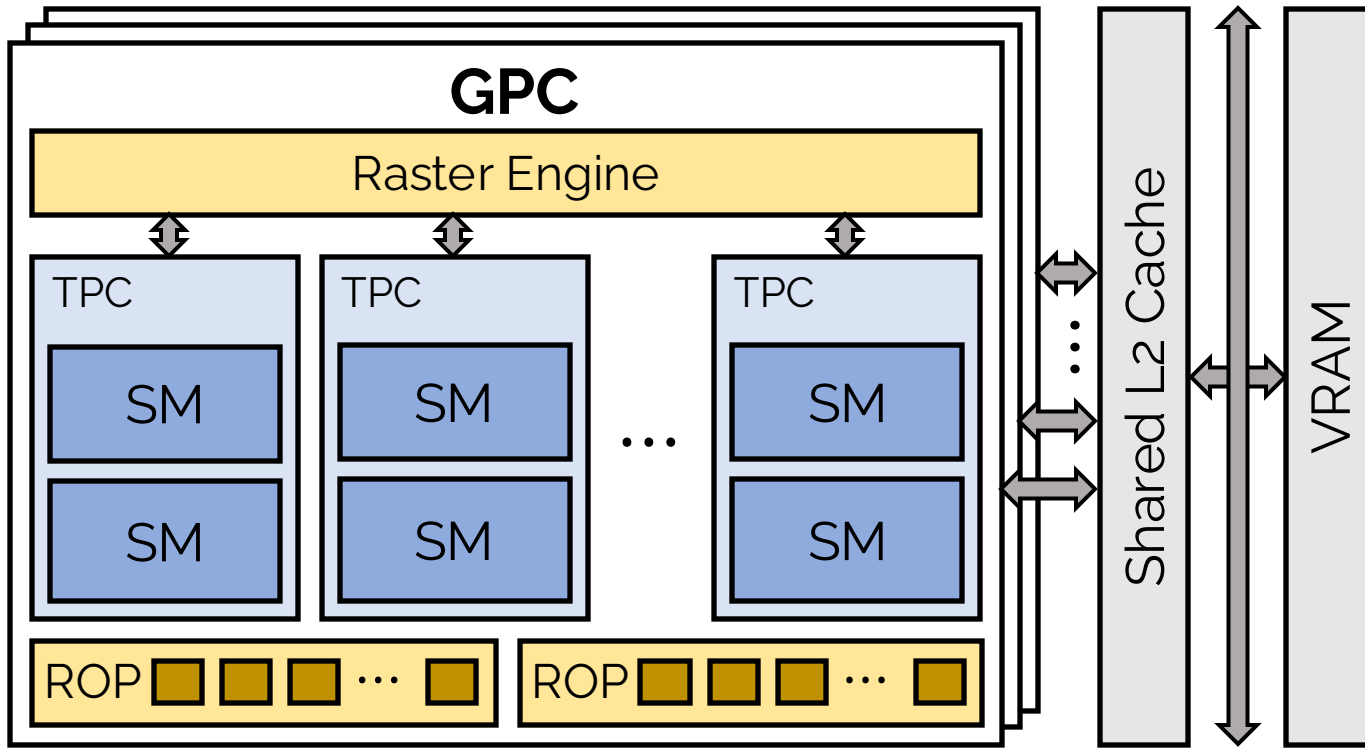
FlashGS [arXiv'24]

## Specialized Accelerators
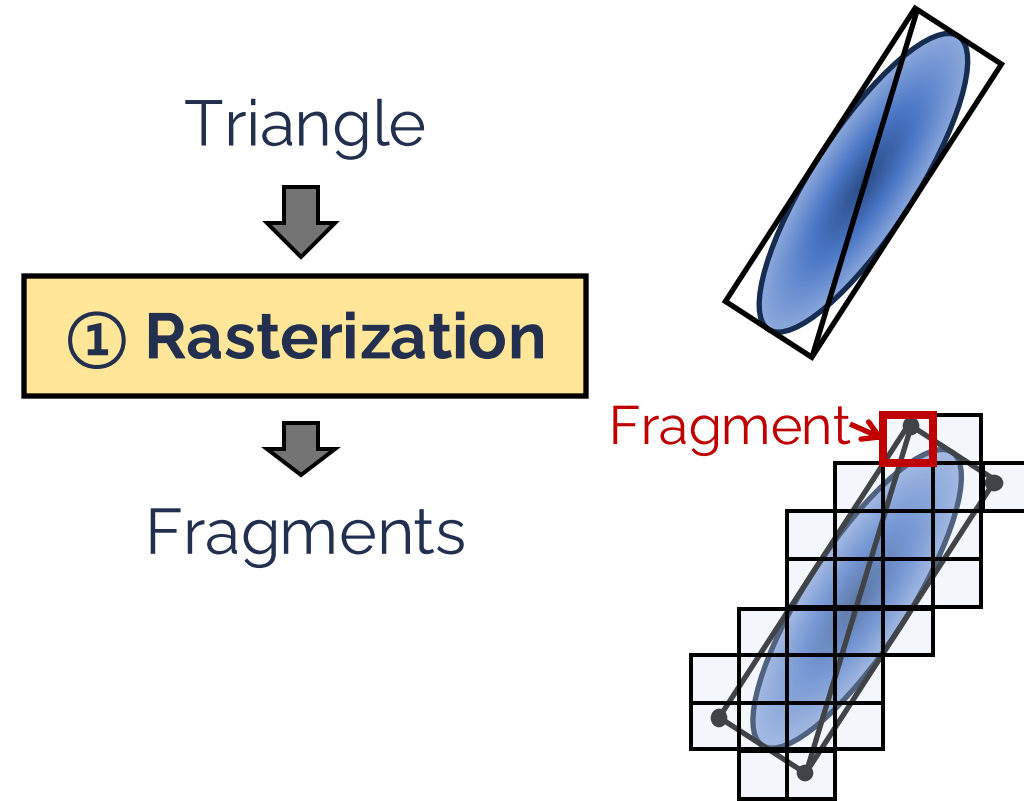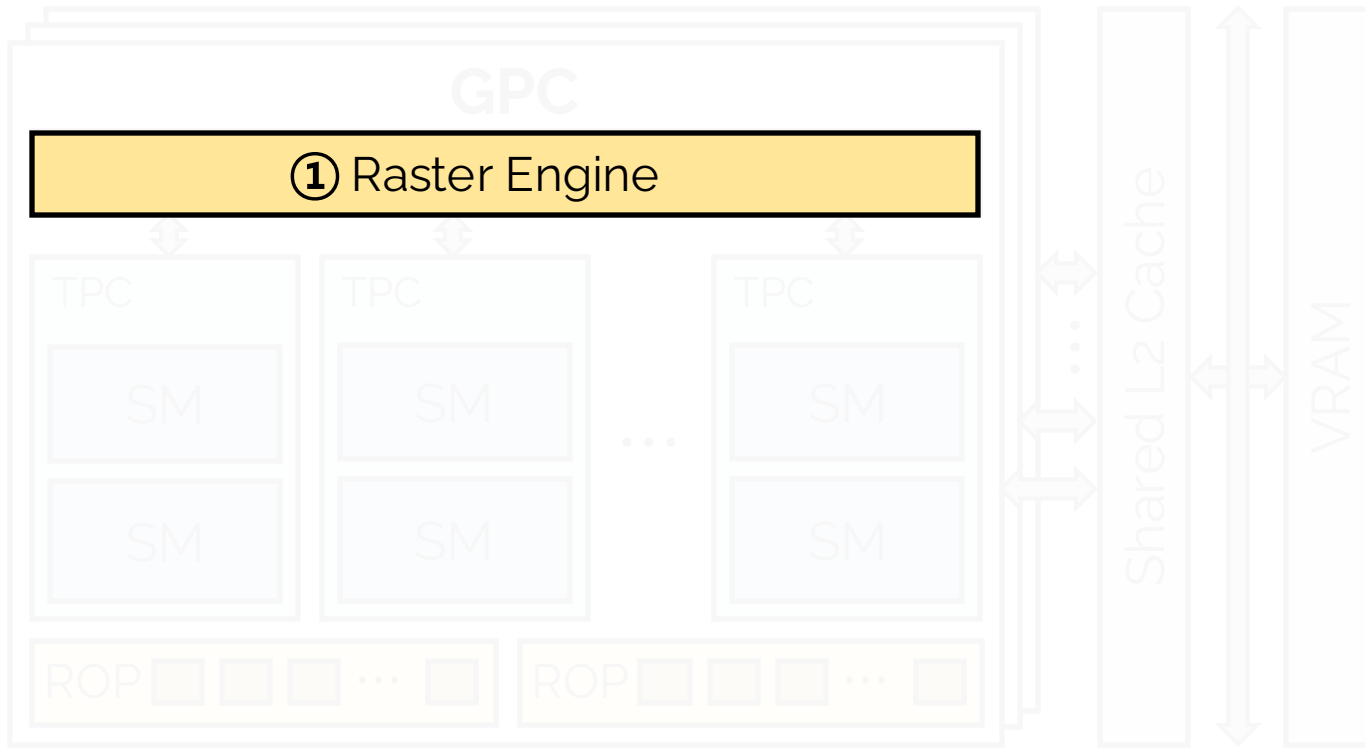
GSCore [ASPLOS'24]

MetaSapiens [ASPLOS'25]

### Our Work

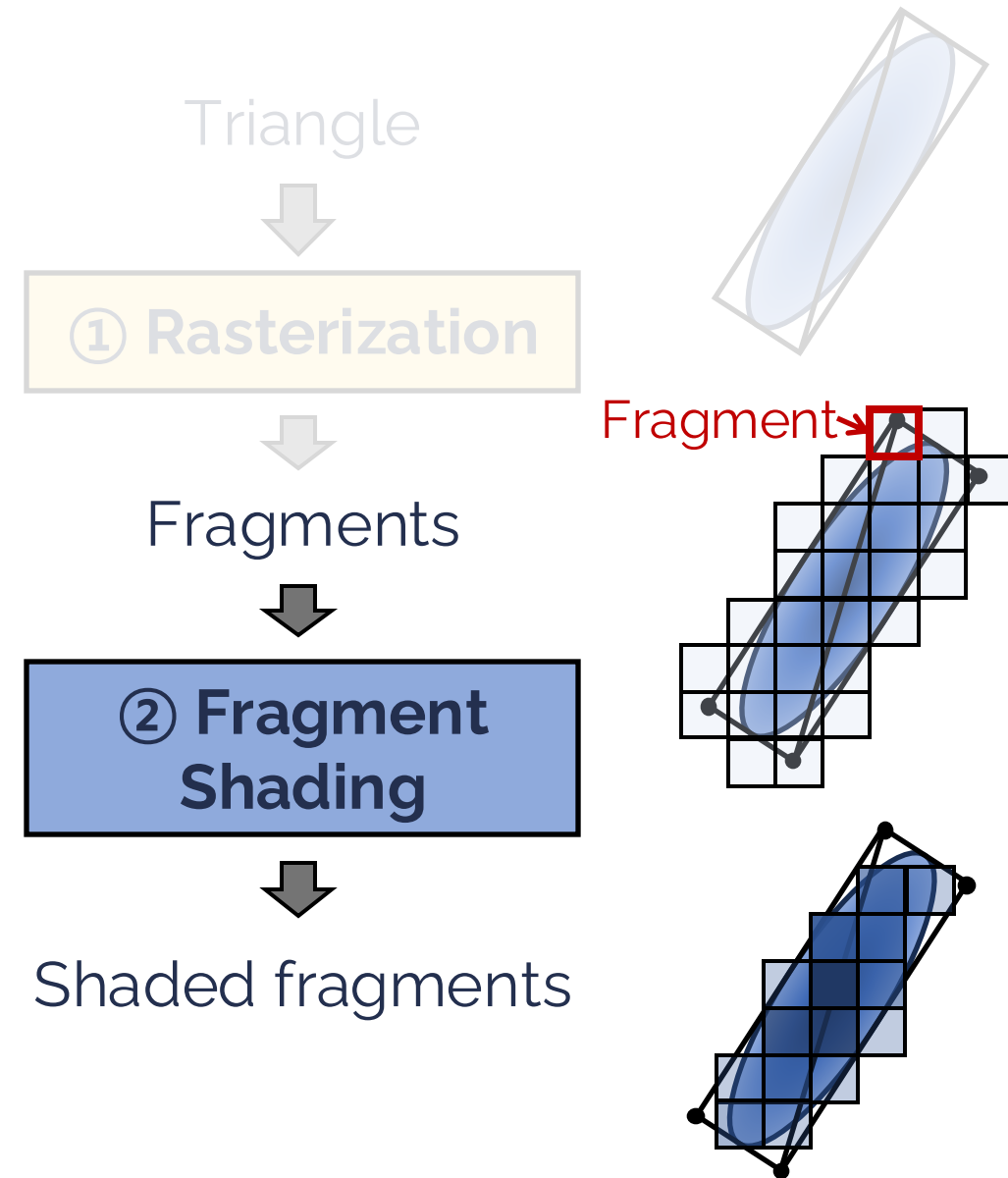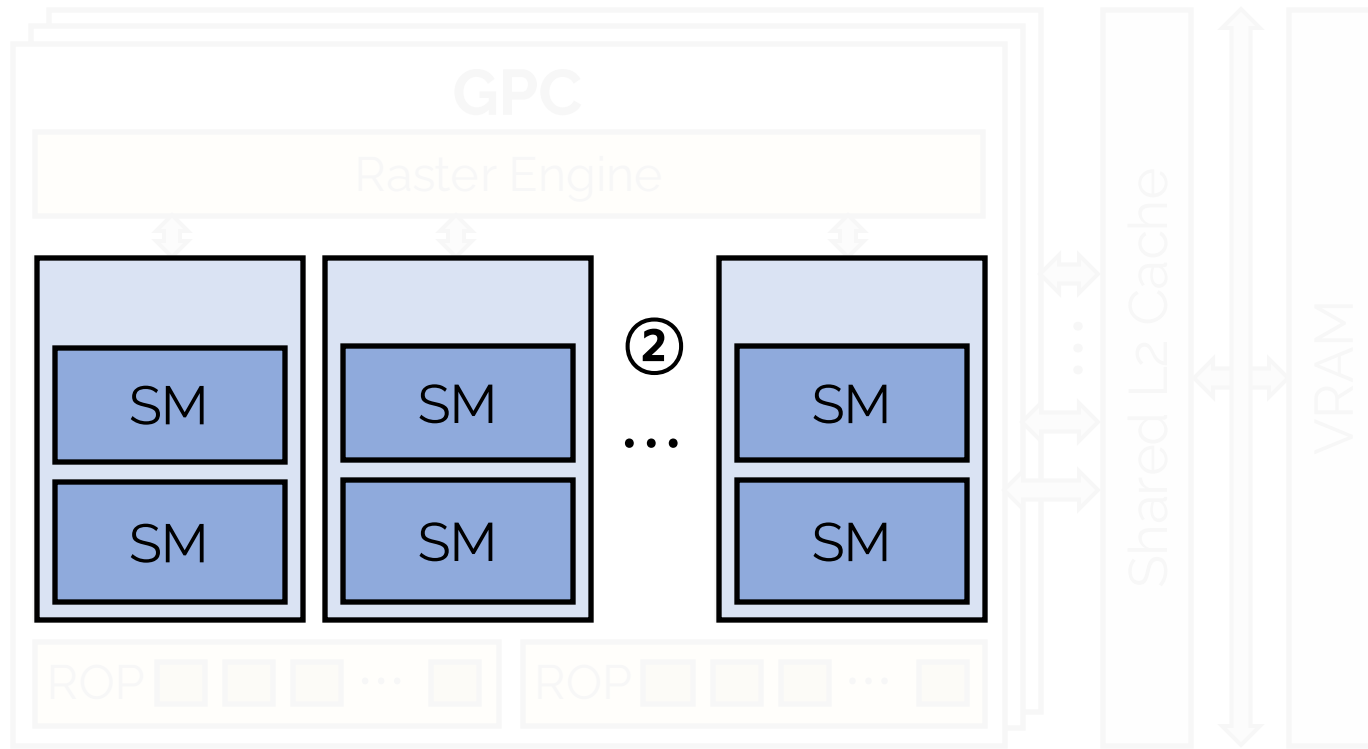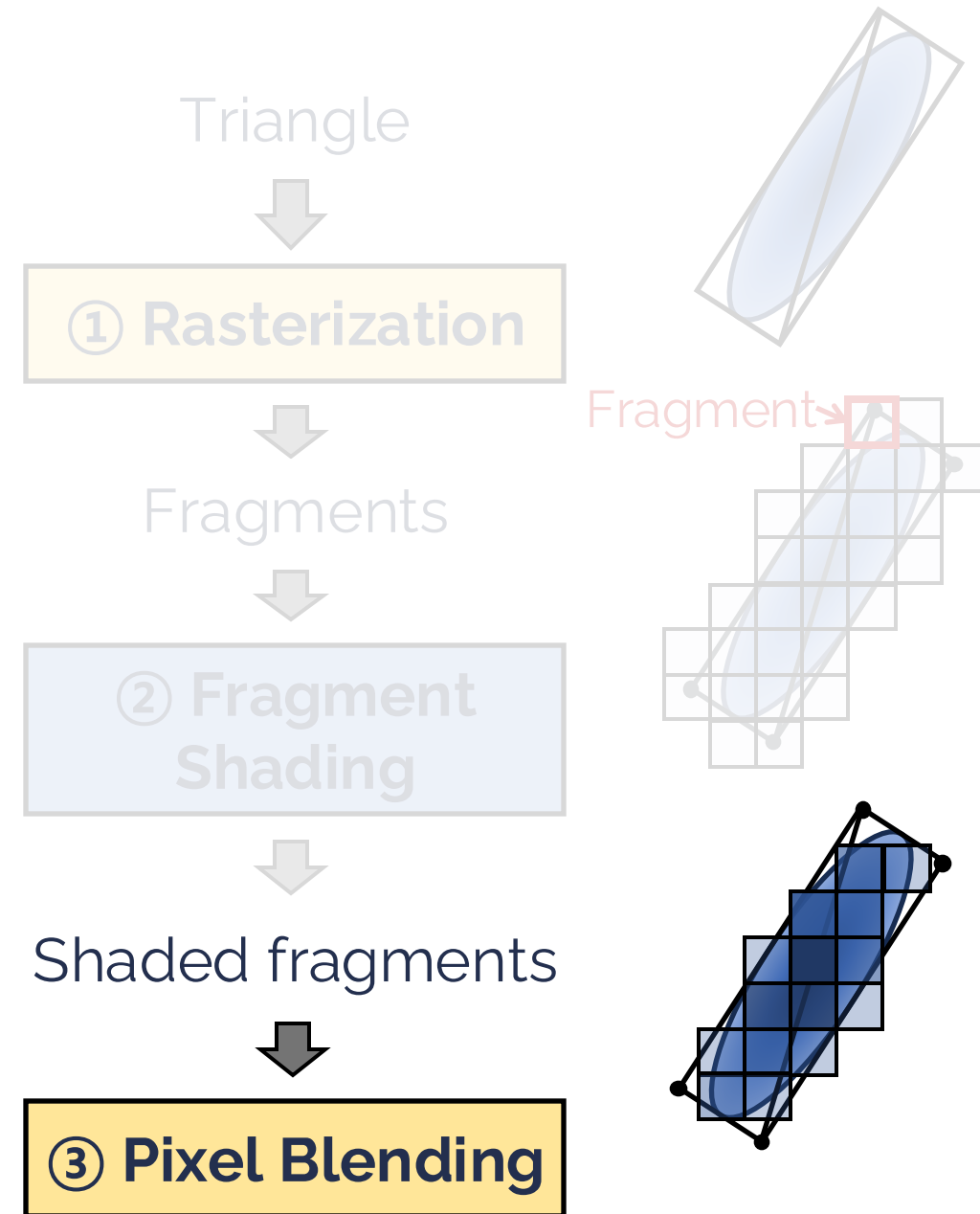Extend the existing **hardware graphics pipeline** for **volume rendering (e.g., 3DGS)**
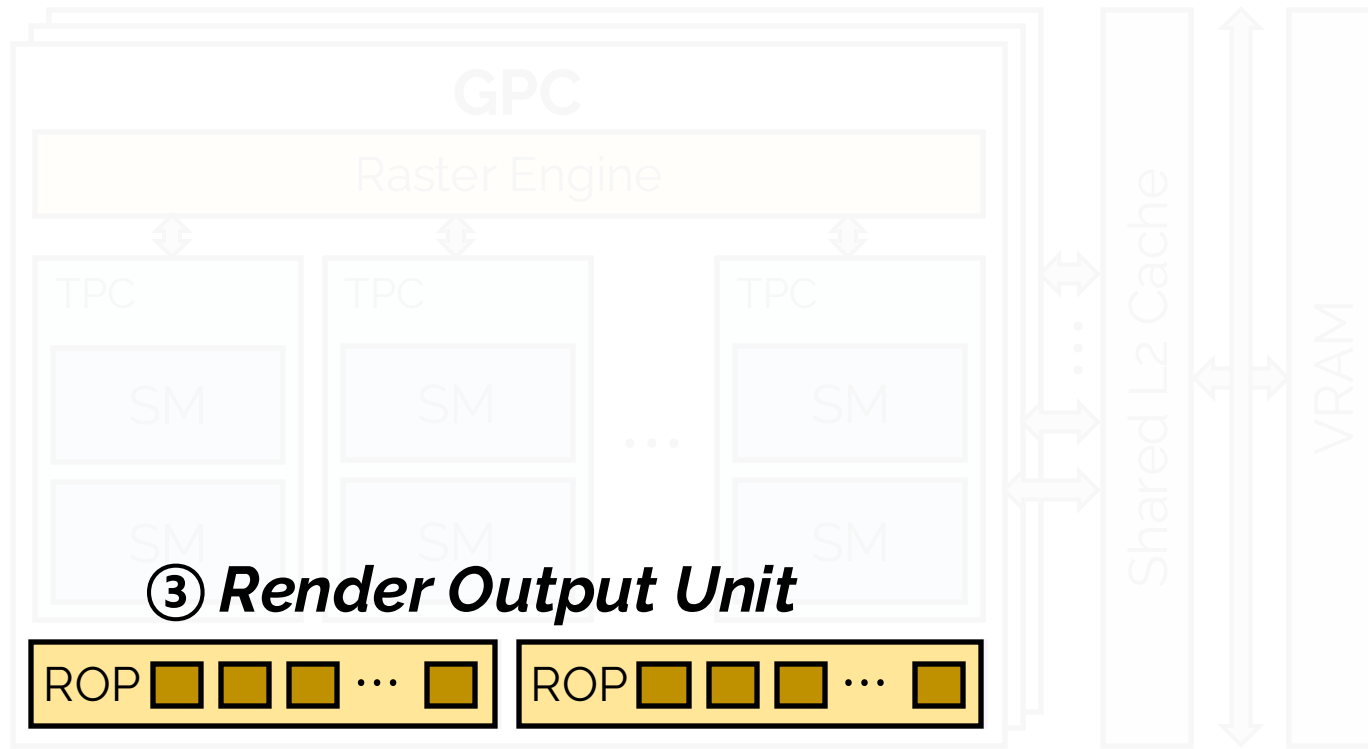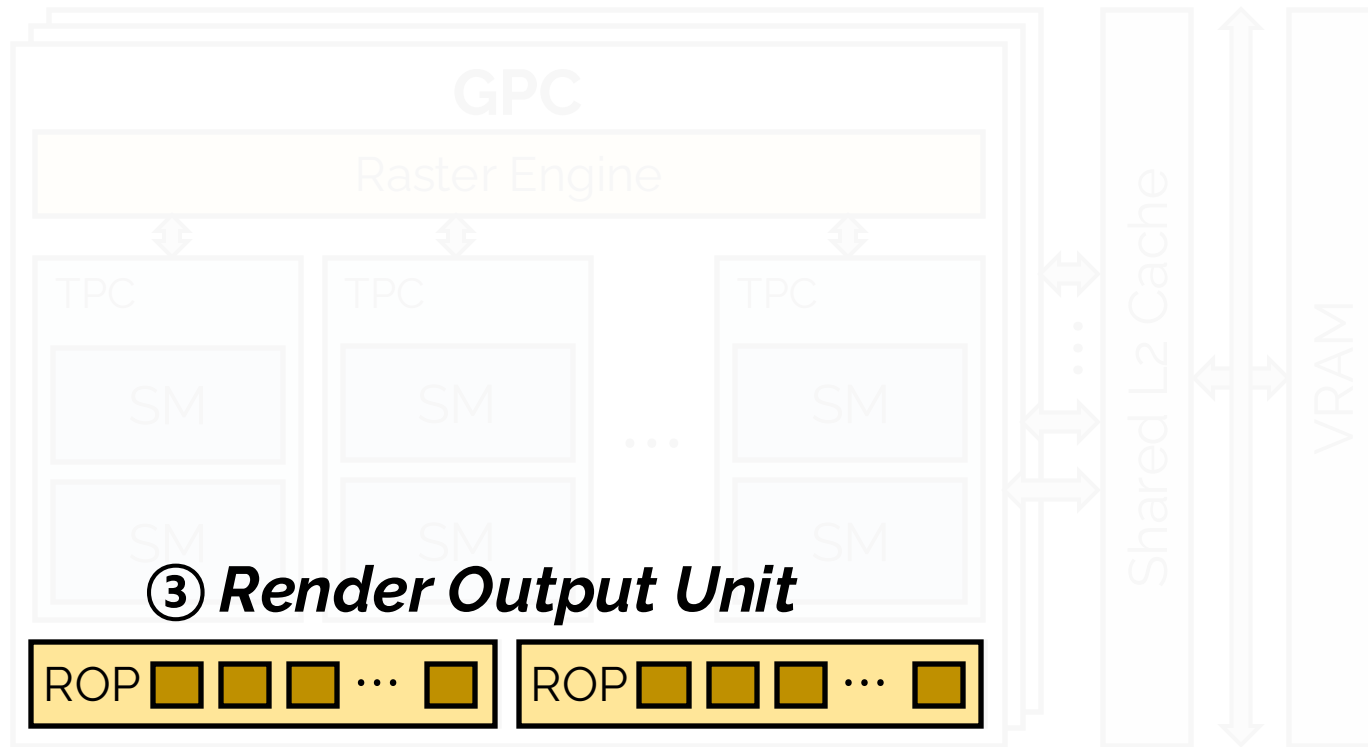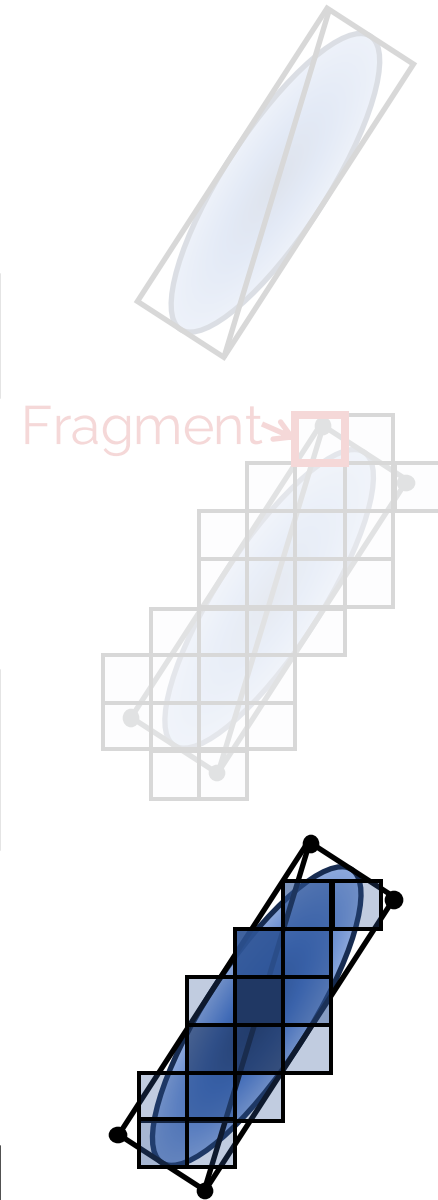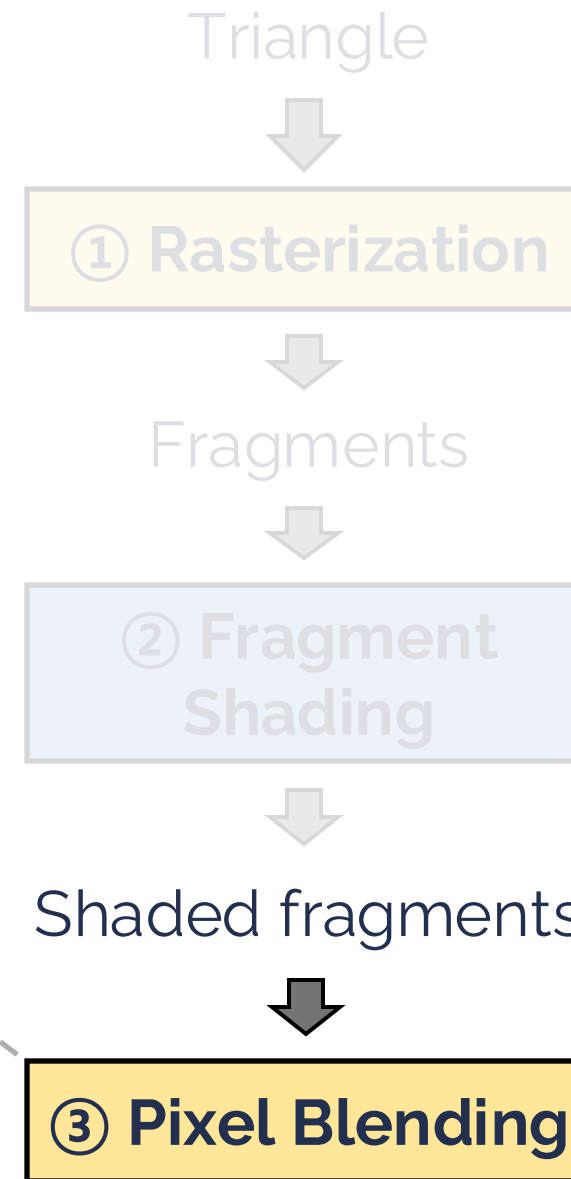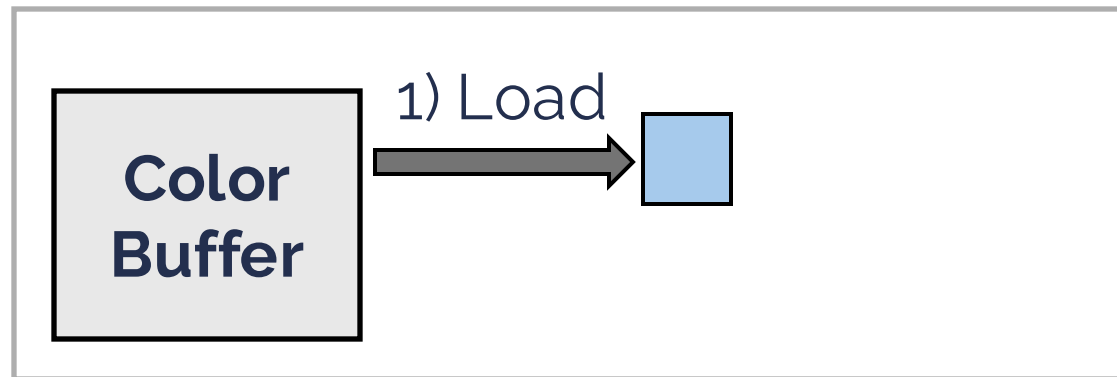
# Hardware Graphics Pipeline

# Hardware Graphics Pipeline

Triangle

① **Rasterization**

Fragments

① Raster Engine

GPC

TPC

TPC

TPC

SM

SM

SM

SM

SM

SM

ROP

ROP

Shared L2 Cache

VRAM

Fragment

# Hardware Graphics Pipeline



Triangle

① **Rasterization**

Fragments

② **Fragment Shading**

Shaded fragments

Fragment

# Hardware Graphics Pipeline

Triangle

① **Rasterization**

Fragment

Fragments

② **Fragment Shading**

③ *Render Output Unit*

ROP ☐ ☐ ☐ ⋯ ☐   ROP ☐ ☐ ☐ ⋯ ☐

Shaded fragments

③ **Pixel Blending**

# Hardware Graphics Pipeline

Triangle

① **Rasterization**

Fragment

Fragments

② **Fragment Shading**

③ *Render Output Unit*

ROP ☐ ☐ ☐ … ☐    ROP ☐ ☐ ☐ … ☐

Shaded fragments

Color Buffer

1) Load

③ **Pixel Blending**

# Hardware Graphics Pipeline

Triangle

① **Rasterization**

Fragment

Fragments

② **Fragment Shading**

③ *Render Output Unit*

ROP ▢ ▢ ▢ ⋯ ▢    ROP ▢ ▢ ▢ ⋯ ▢

Shaded fragments

**Color Buffer**    1) Load    2) Blend

▢ **+** ▢ **=** ▢

③ **Pixel Blending**

# Hardware Graphics Pipeline

Triangle

① **Rasterization**

Fragment

Fragments

② **Fragment Shading**

③ *Render Output Unit*

ROP ⬛ ⬛ ⬛ ⋯ ⬛   ROP ⬛ ⬛ ⬛ ⋯ ⬛

Shaded fragments

**Color Buffer**

1) Load

2) Blend

⬛ **+** ⬛ **=** ⬛

3) Store

③ **Pixel Blending**

5

# Outline

- **Background**
  - 3D Gaussian Splatting (3DGS)
  - Hardware Graphics Pipeline

- **Limitations of Graphics Hardware**

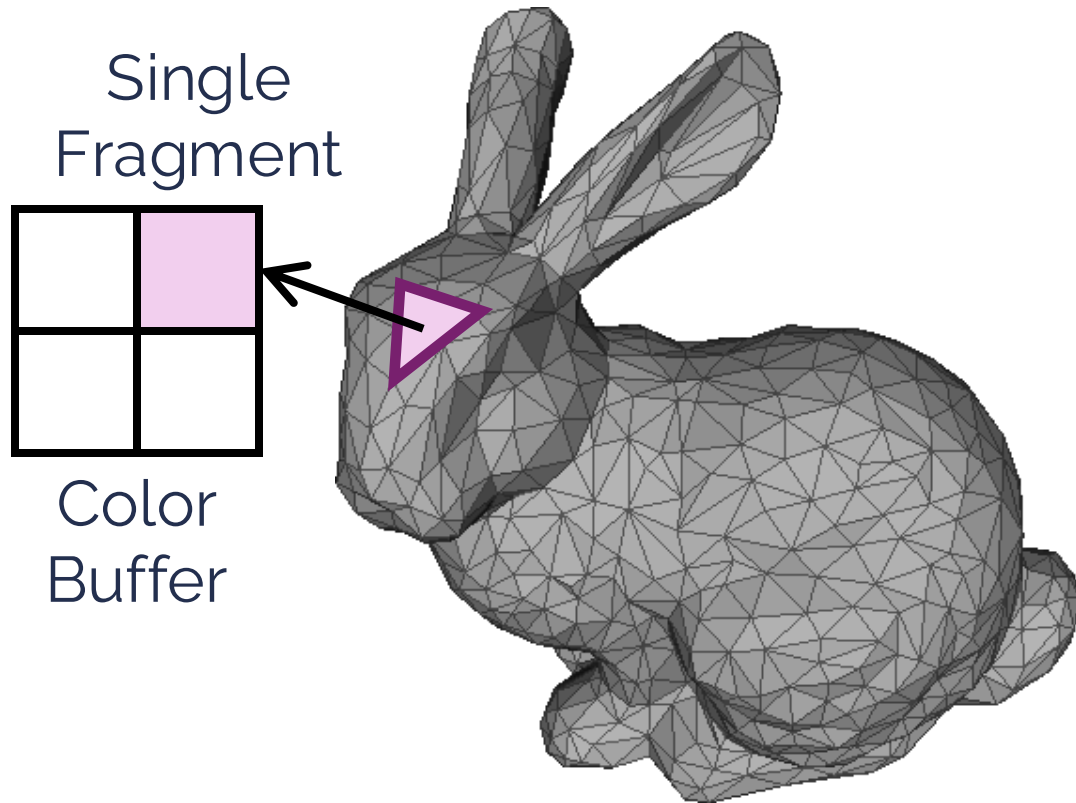- **VR-Pipe: Graphics Hardware Extension for Volume Rendering**
  - Quad Merging with Multi-Granular Tile Binning
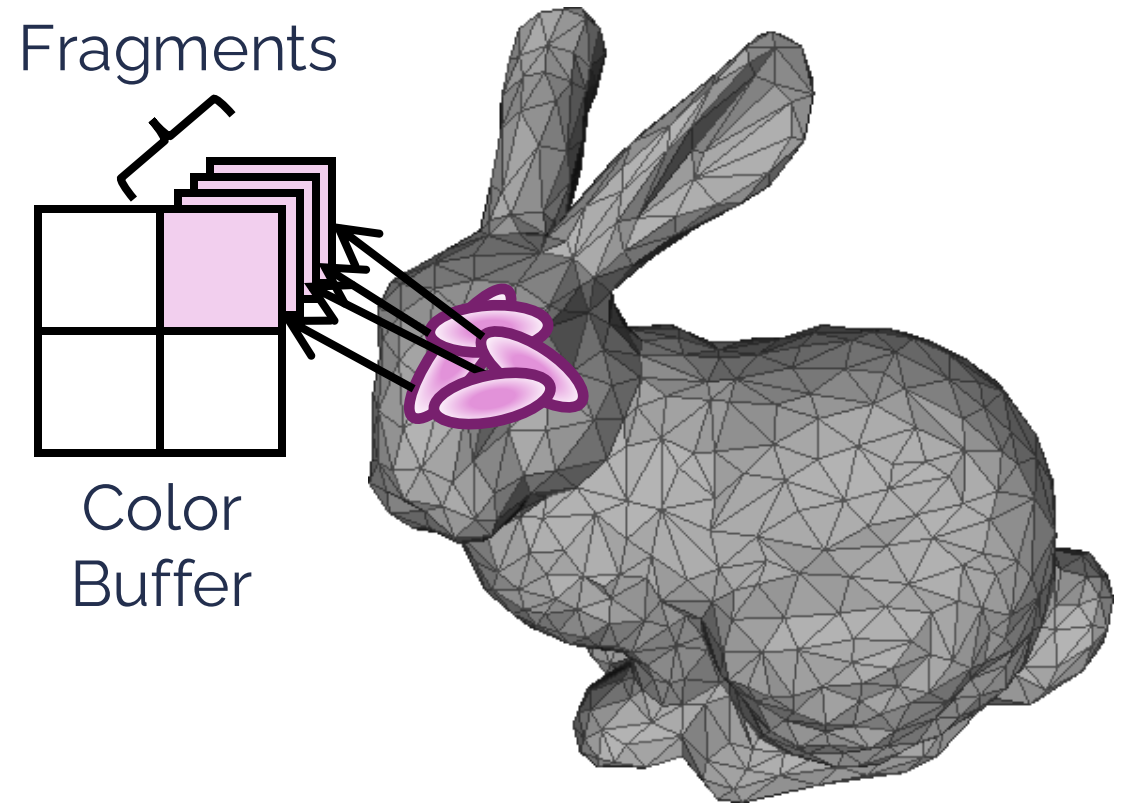  - Hardware Support for Early Termination

- **Evaluation**

- **Conclusion**

# Limitations of Graphics Hardware
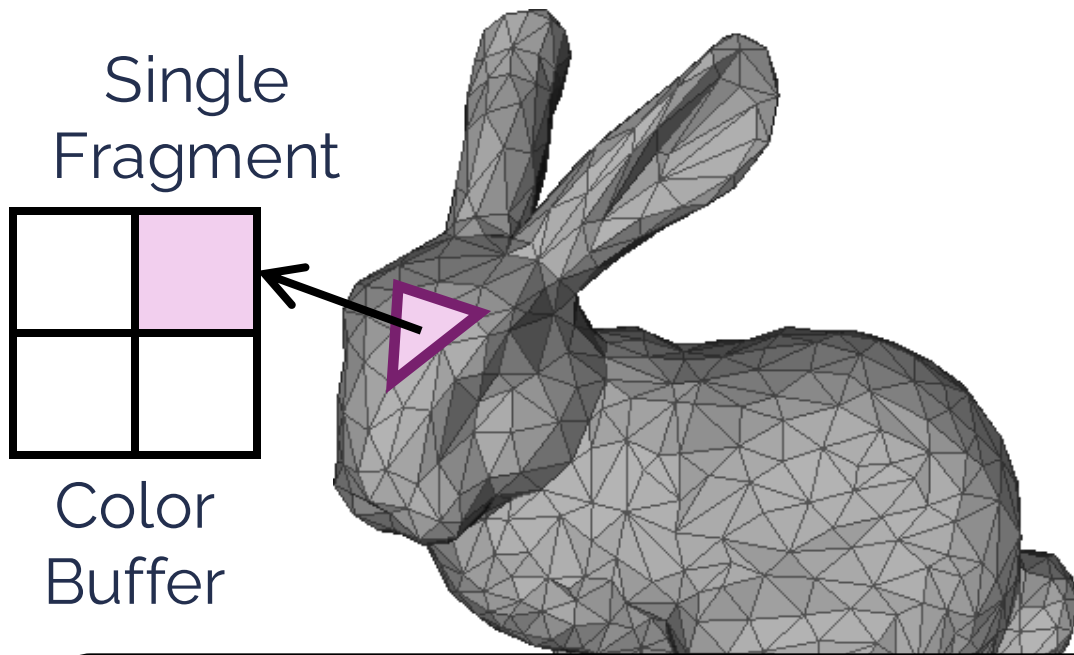
## Mesh-based Rendering

Single Fragment

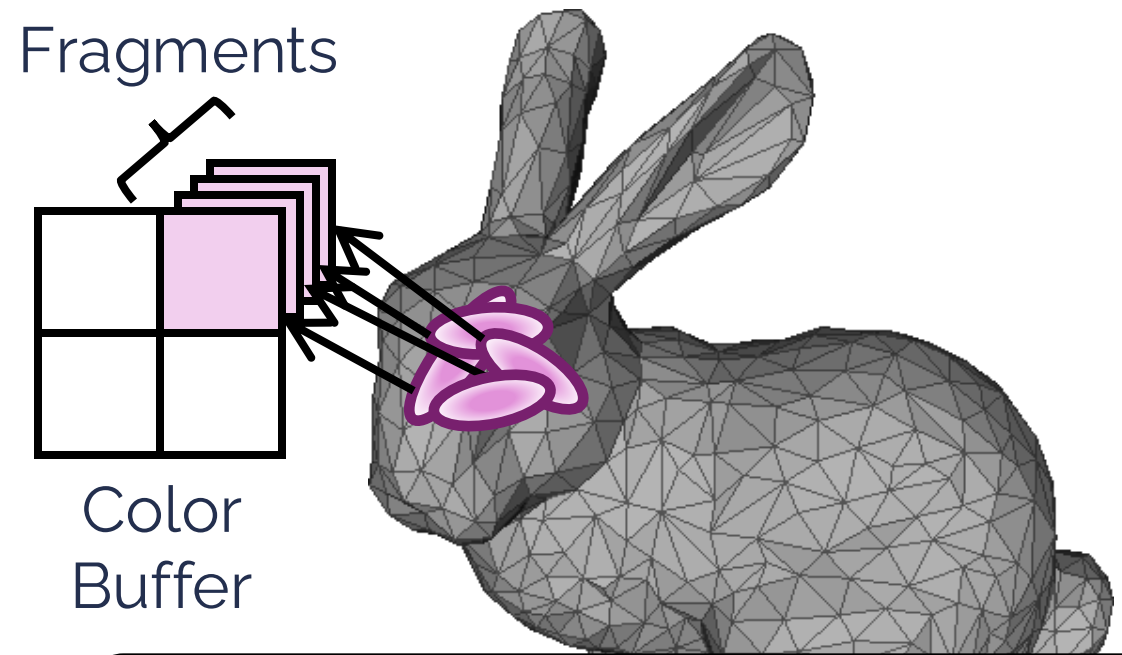Color Buffer

## Volume Rendering (e.g., 3DGS)

Fragments

Color Buffer

# Limitations of Graphics Hardware

## Mesh-based Rendering

Single Fragment

Color Buffer

**One or few *opaque* fragments** per pixel

## Volume Rendering (e.g., 3DGS)

Fragments

Color Buffer

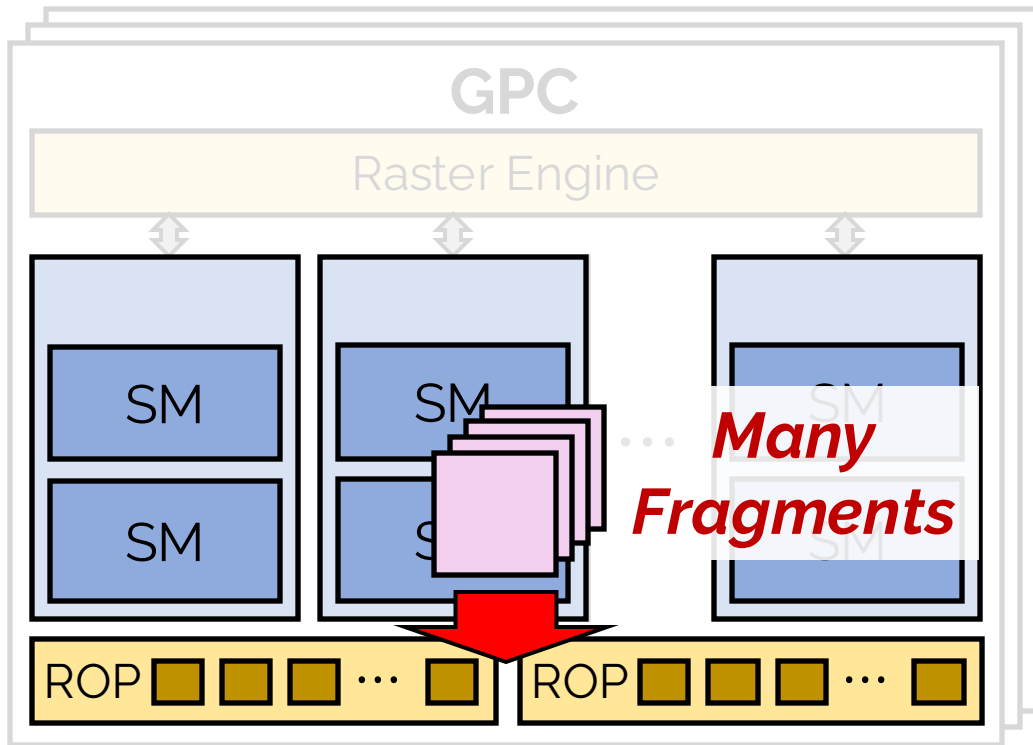**Many *transparent* fragments (i.e., 10-100)** per pixel

# Limitations of Graphics Hardware

# Limitations of Graphics Hardware

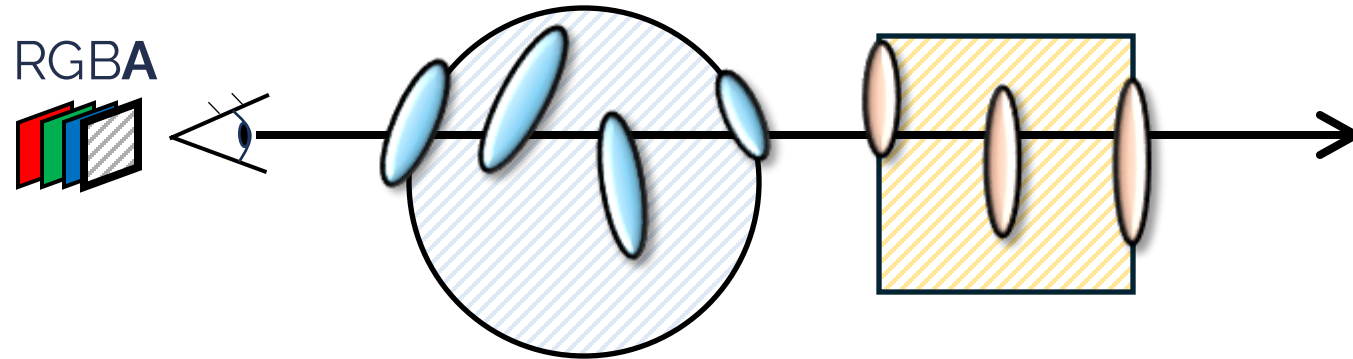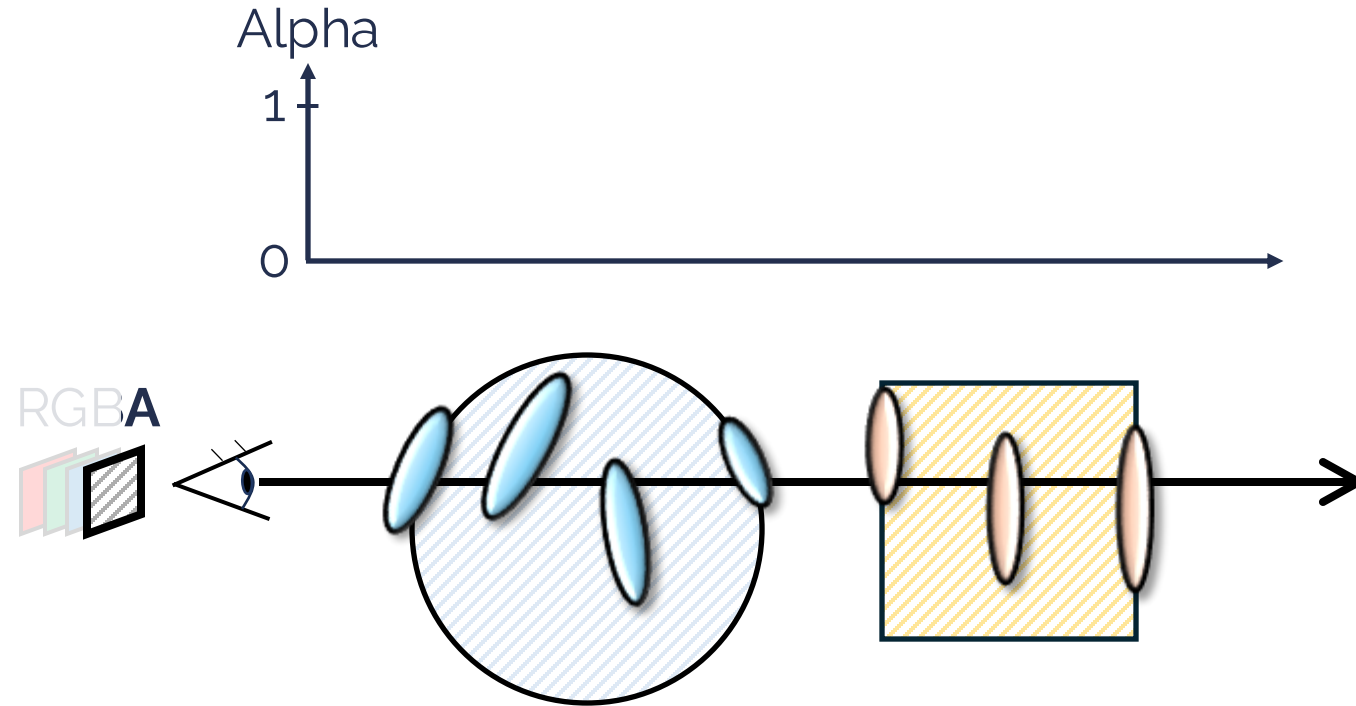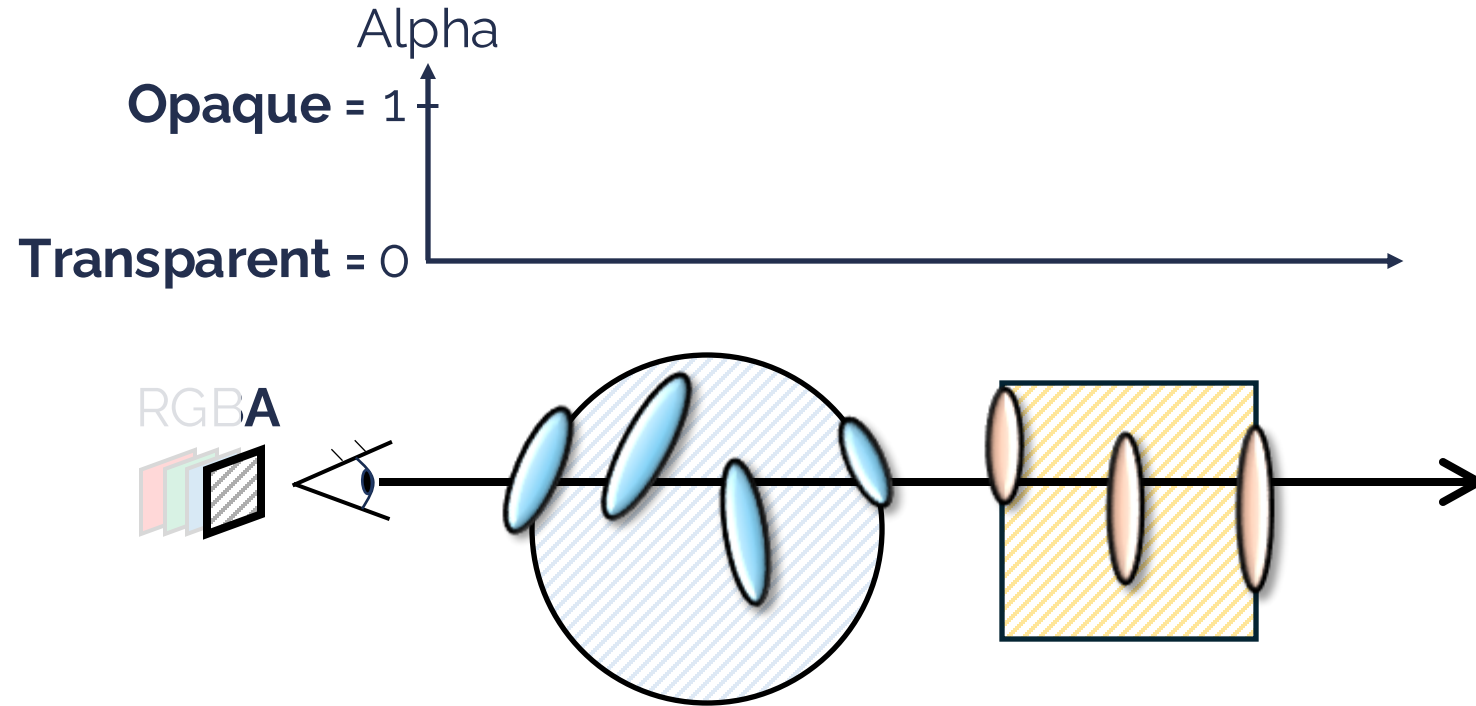# Limitations of Graphics Hardware

## Early Termination

# Limitations of Graphics Hardware

## Early Termination

# Limitations of Graphics Hardware

## Early Termination
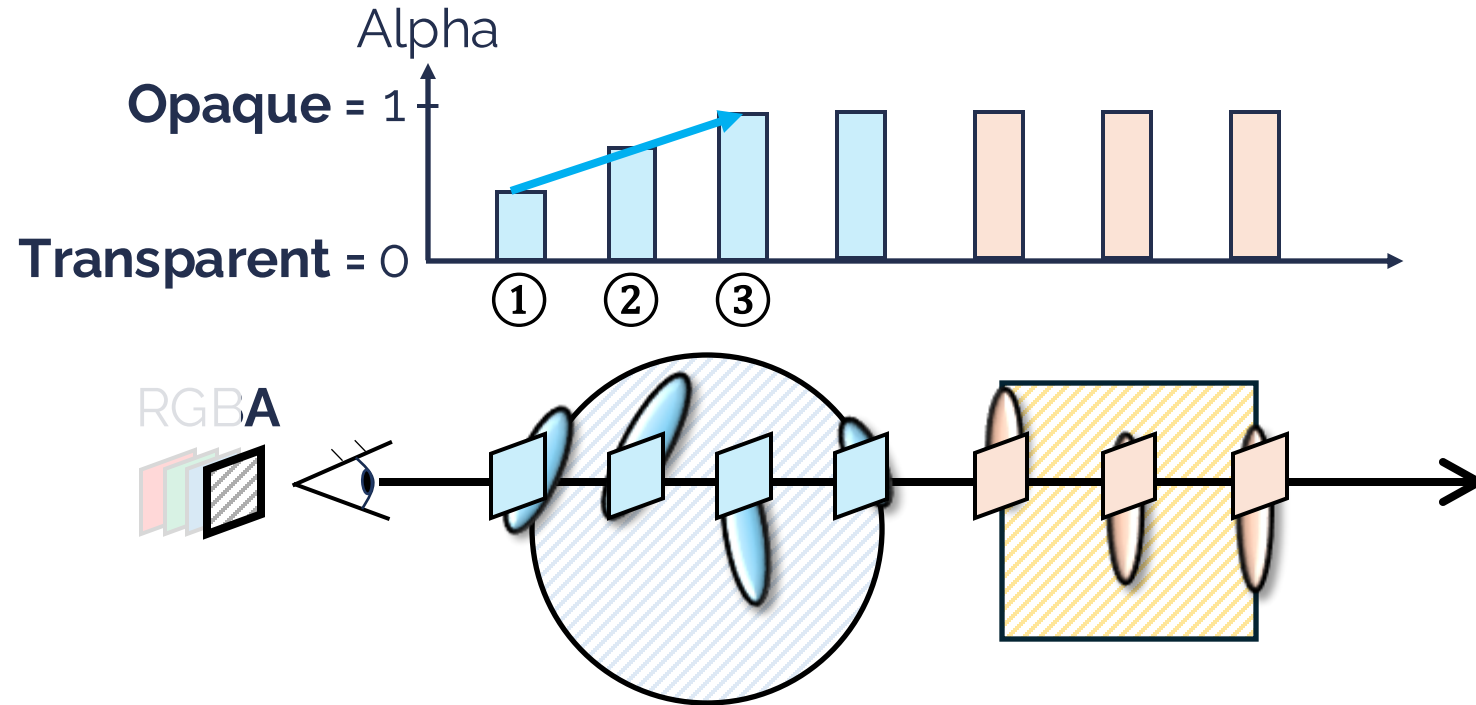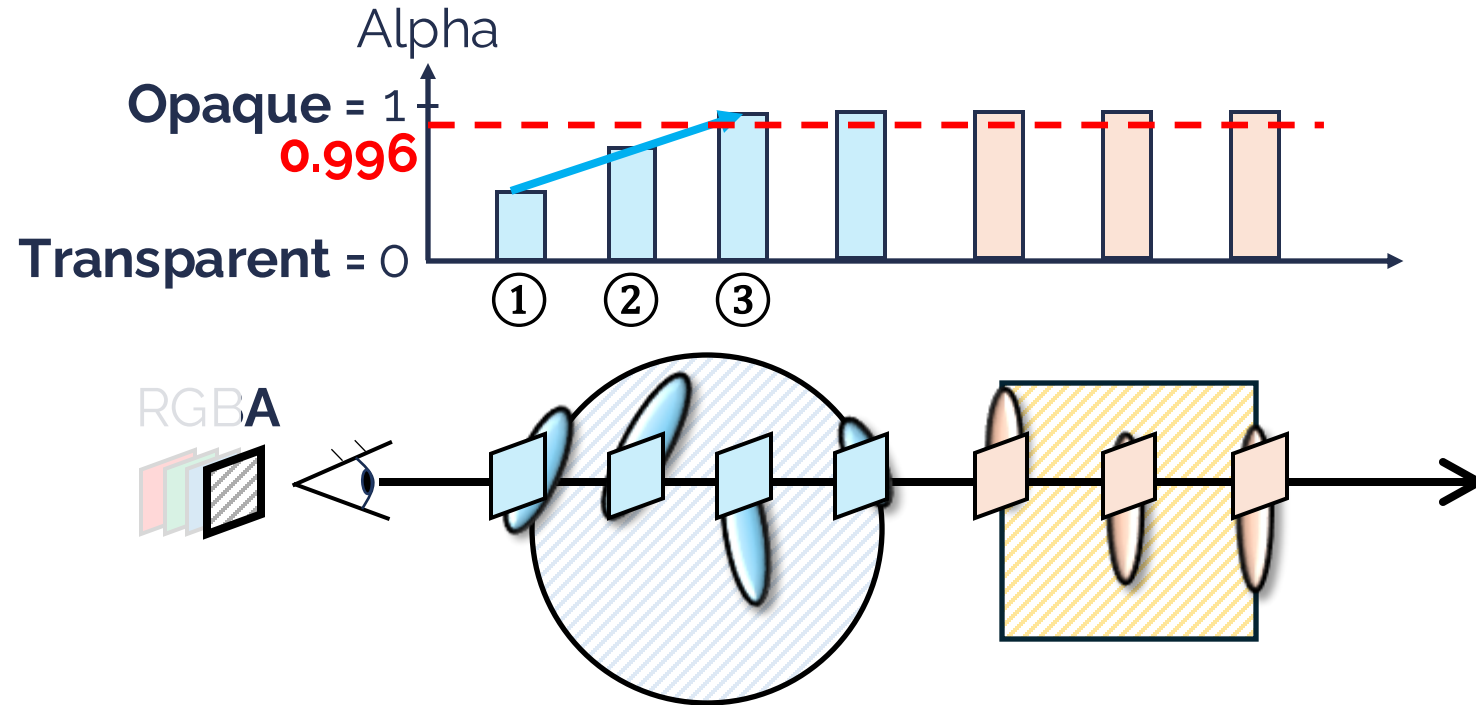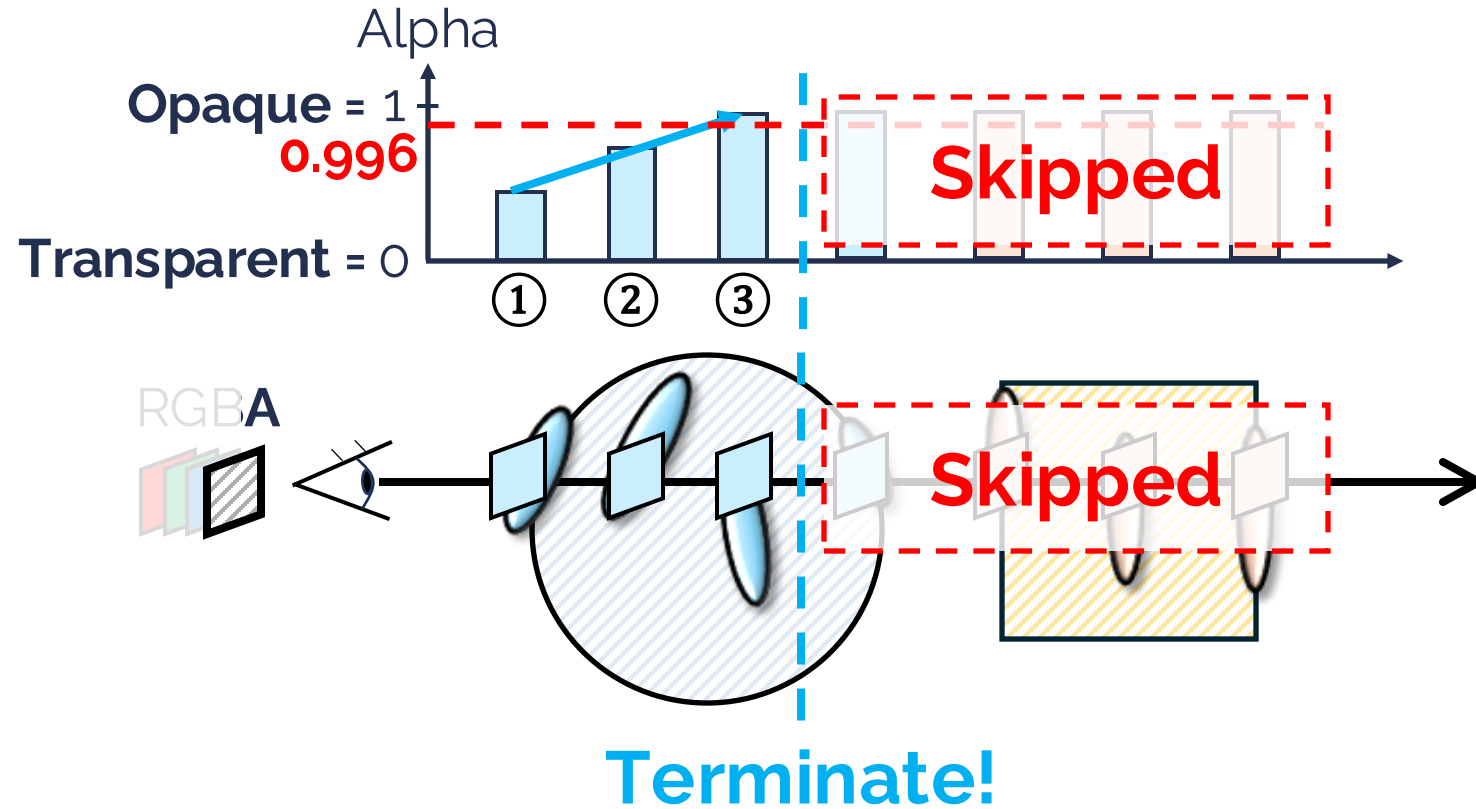
# Limitations of Graphics Hardware

## Early Termination

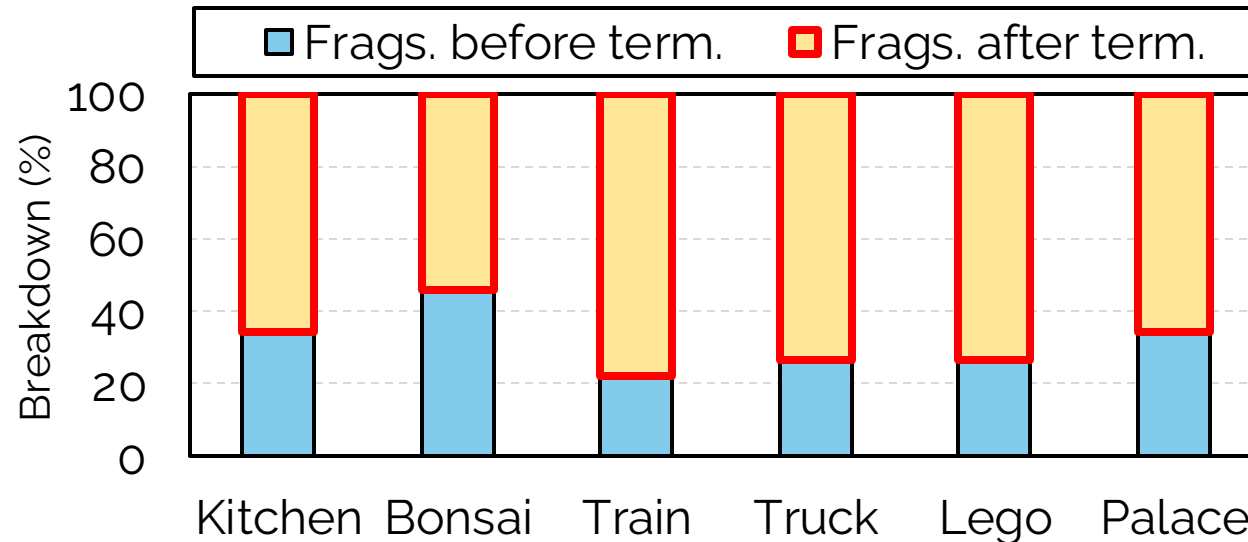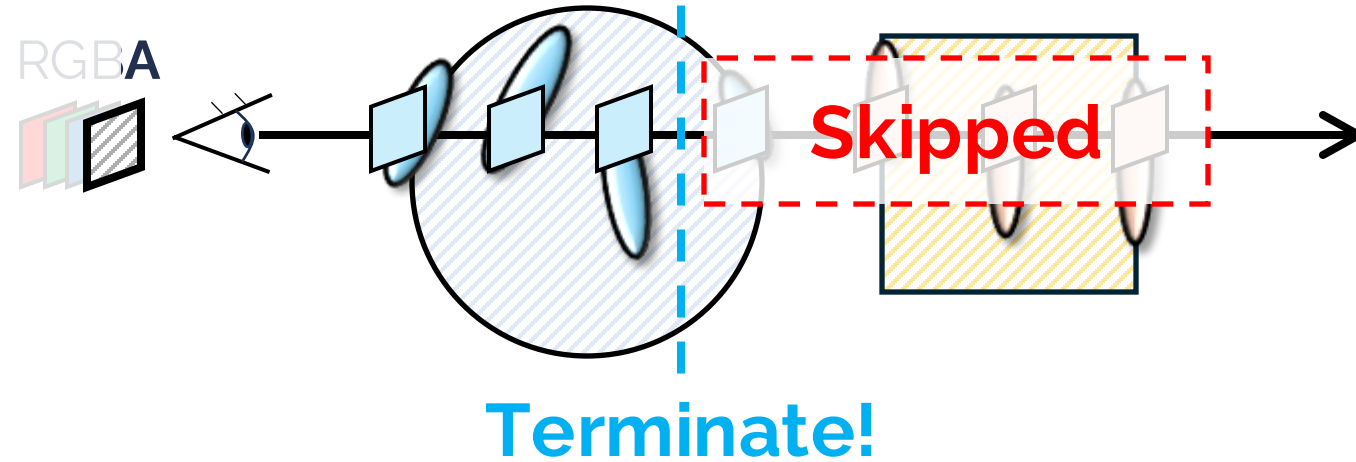# Limitations of Graphics Hardware

## Early Termination

# Limitations of Graphics Hardware

## Early Termination

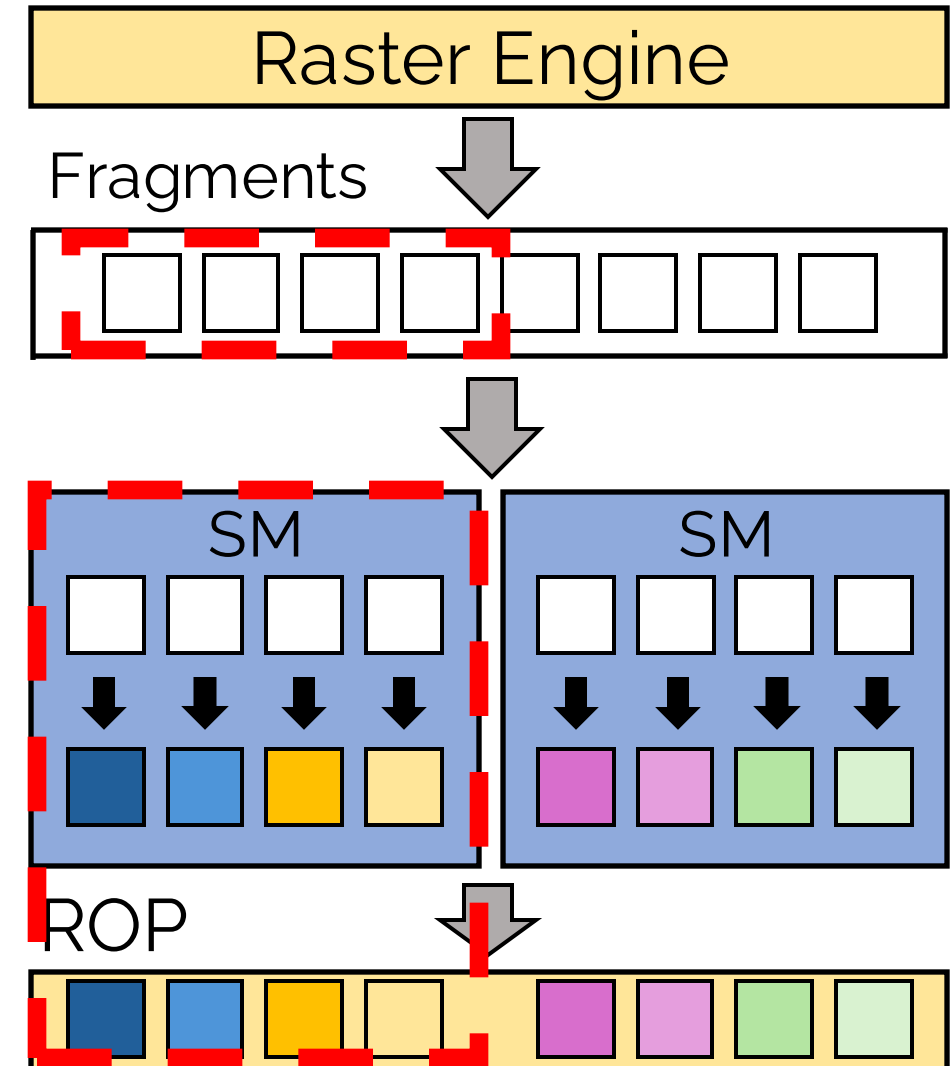# Limitations of Graphics Hardware

## Early Termination

# Limitations of Graphics Hardware
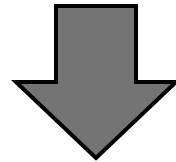
# Limitations of Graphics Hardware

**Observation 1**

Many fragments are **unnecessarily shaded and blended**

Raster Engine

Fragments

SM

SM

ROP

# Limitations of Graphics Hardware
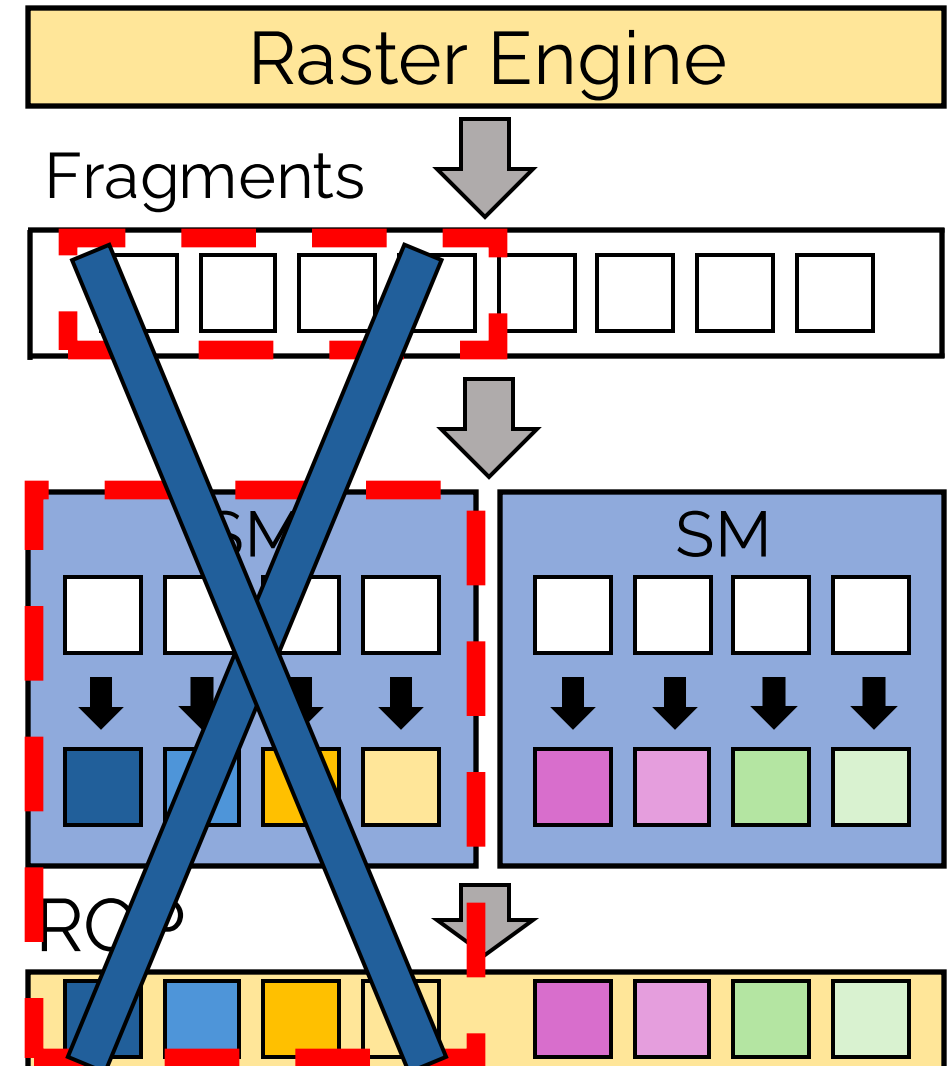
**Observation 1**

Many fragments are **unnecessarily shaded and blended**

↓

**Proposal 1**

Add **hardware support for early termination**
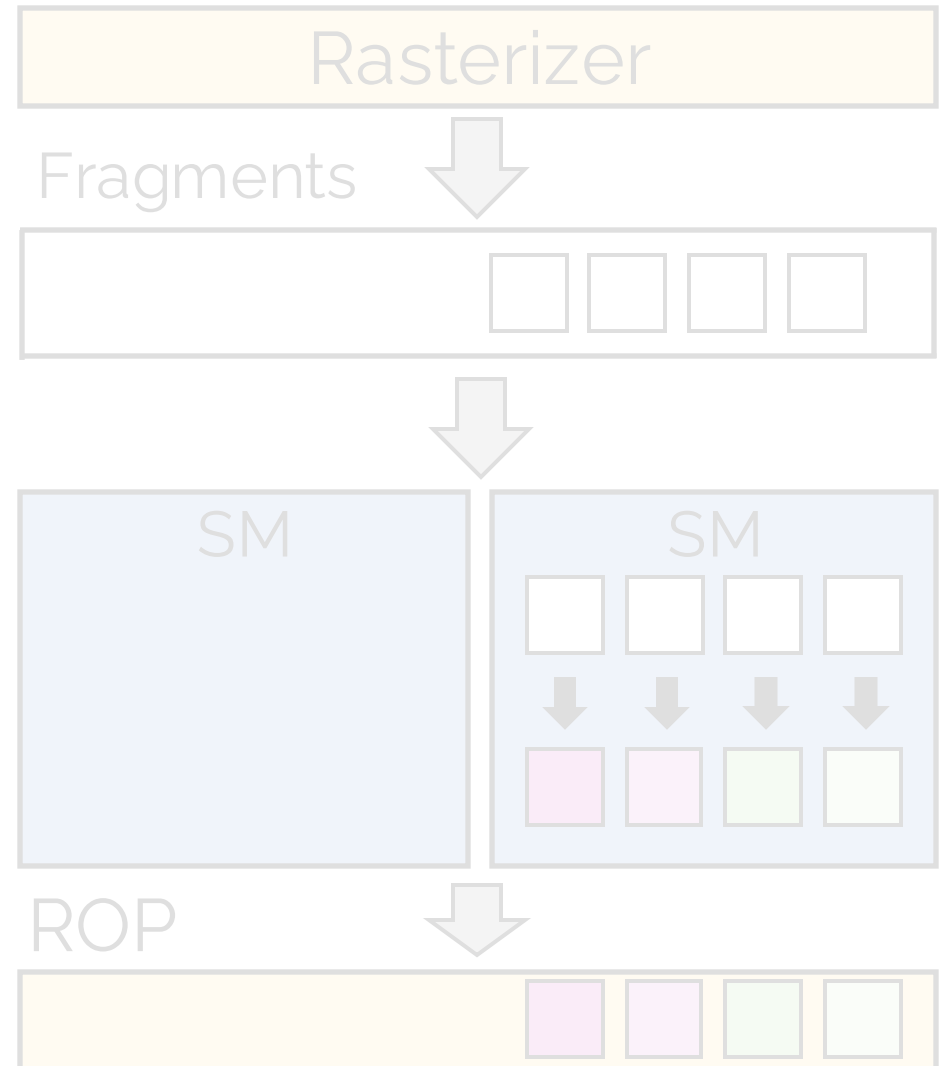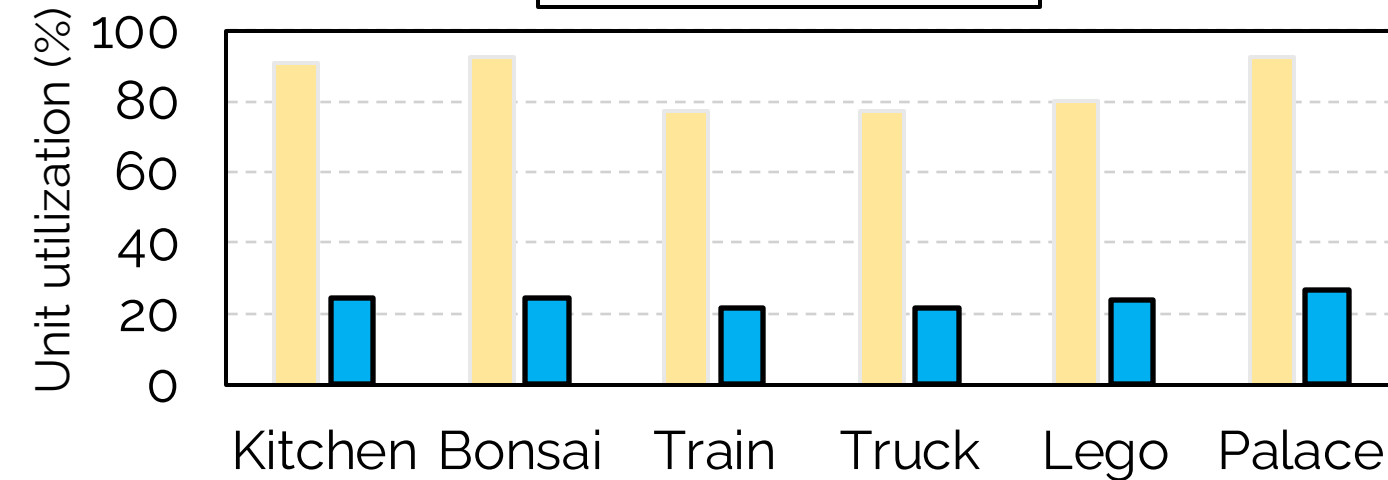
**= Hardware-Based Early Termination (HET)**

Raster Engine

Fragments

SM

SM

ROP

# Limitations of Graphics Hardware
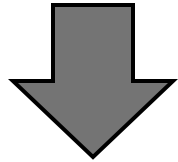
# Limitations of Graphics Hardware

# Limitations of Graphics Hardware

*Observation 2*

**SMs are underutilized**
due to back pressure

*Proposal 2*

**Partially blend fragments**
in SMs

**= Quad Merging (QM)**

Rasterizer

Fragments

SM

SM

ROP

# Limitations of Graphics Hardware



**VR-Pipe improves rendering performance**
by reducing ROP pressure! ☺

# Outline

- **Background**
  - 3D Gaussian Splatting (3DGS)
  - Hardware Graphics Pipeline
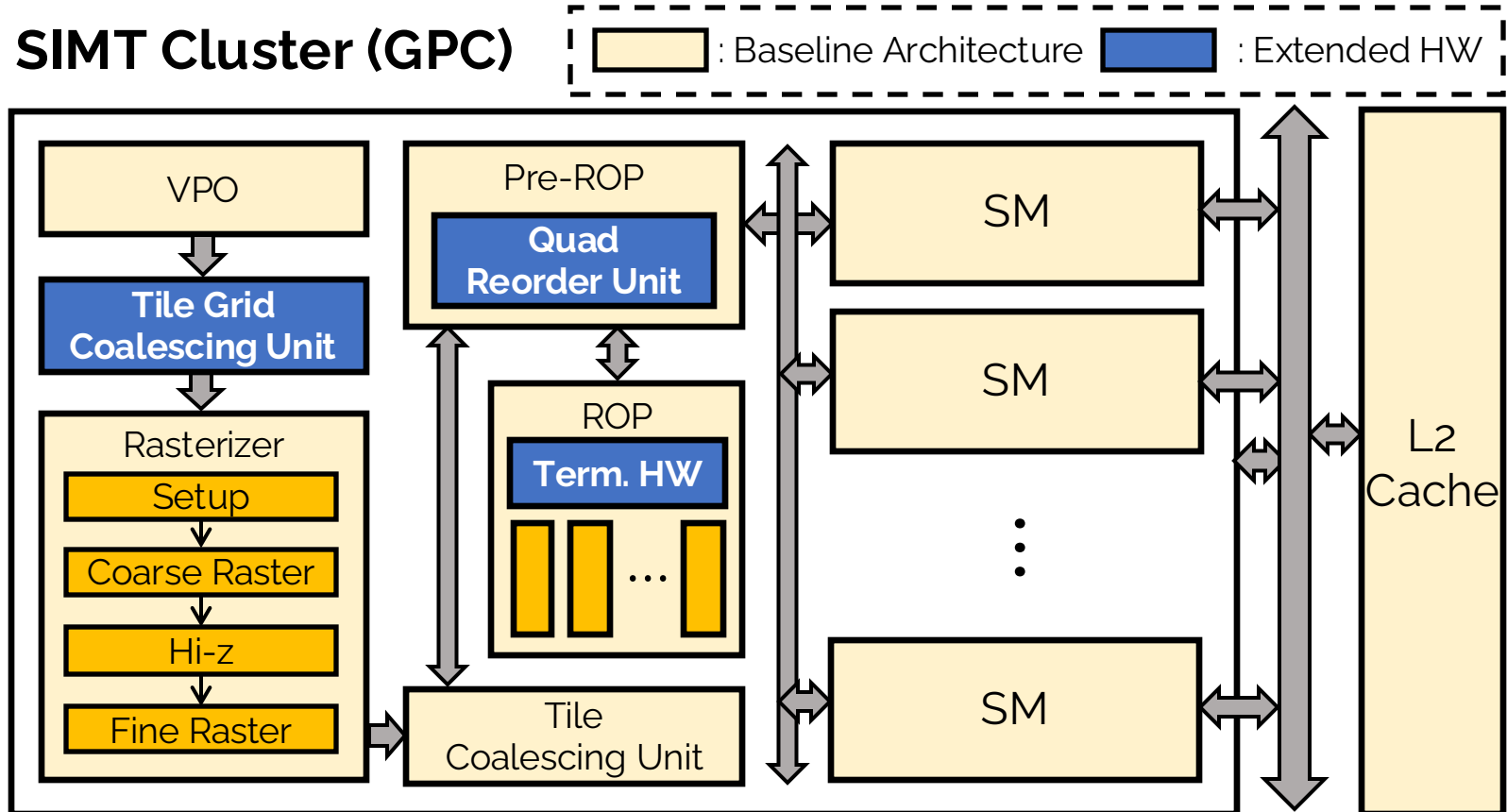
- **Limitations of Graphics Hardware**

- **VR-Pipe: Graphics Hardware Extension for Volume Rendering**
  - Quad Merging with Multi-Granular Tile Binning
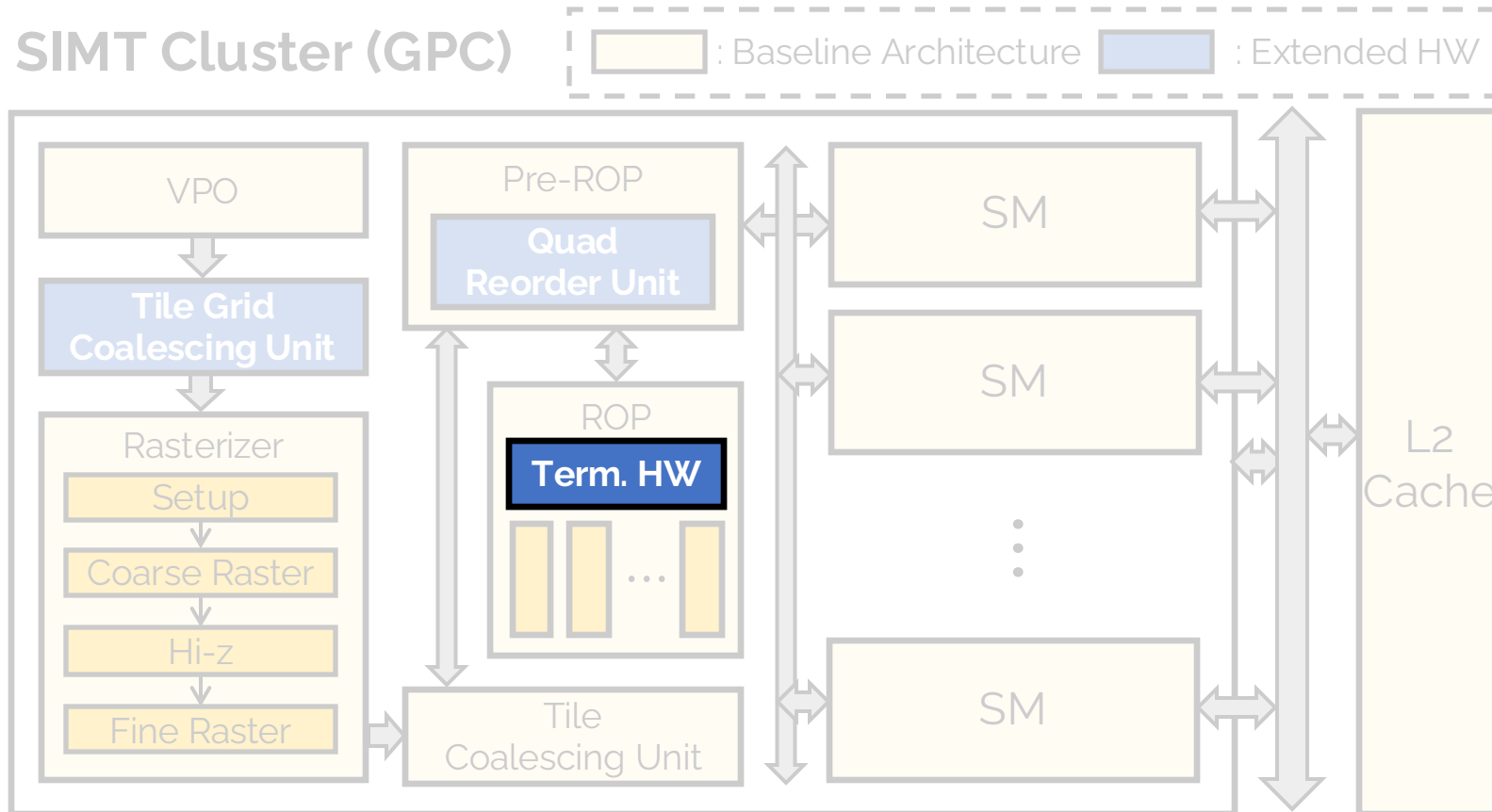  - Hardware Support for Early Termination

- **Evaluation**

- **Conclusion**

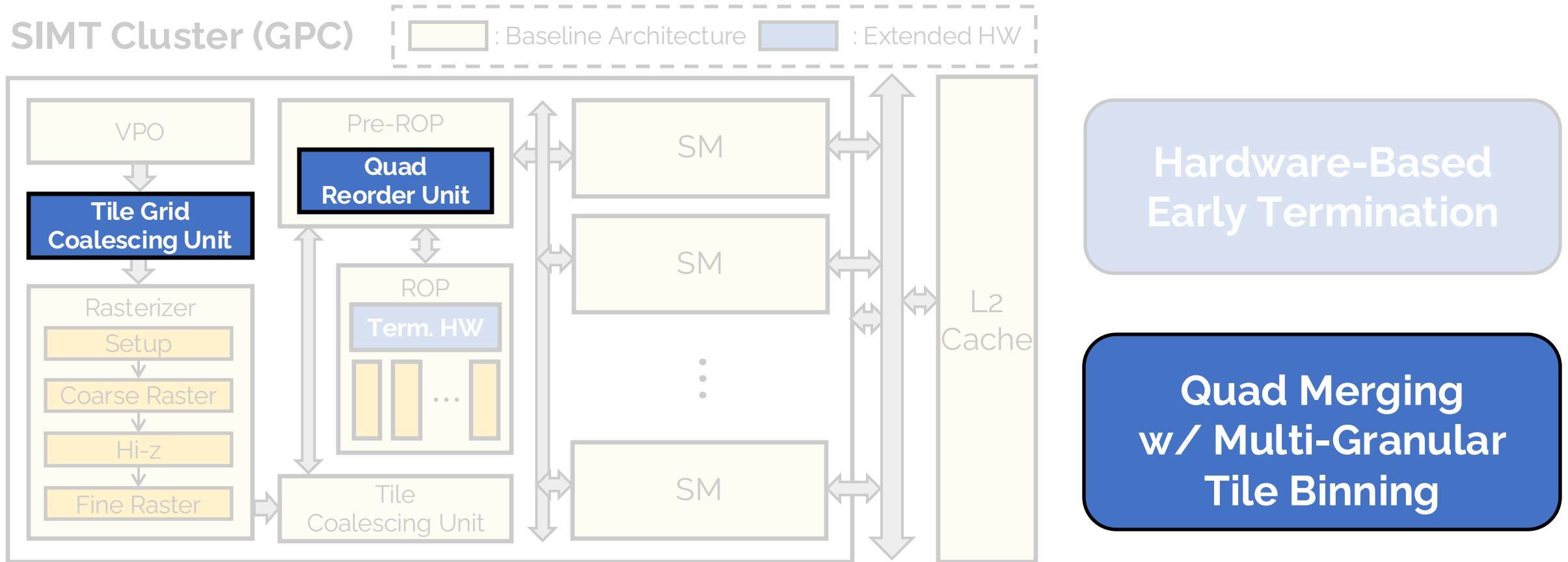# VR-Pipe: GPU Extension for Volume Rendering
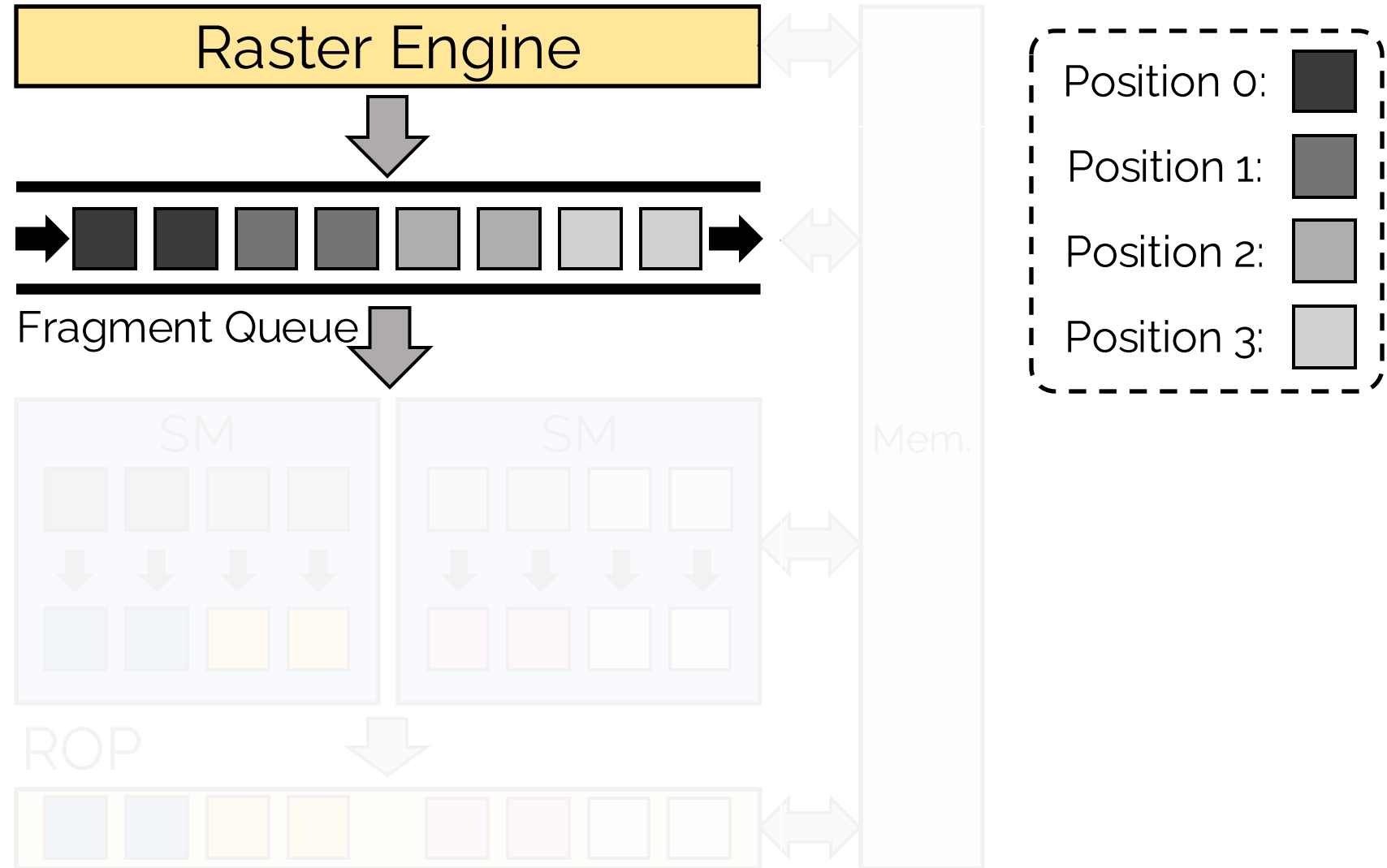
# VR-Pipe: GPU Extension for Volume Rendering

# VR-Pipe: GPU Extension for Volume Rendering

# Hardware-Based Early Termination



Raster Engine
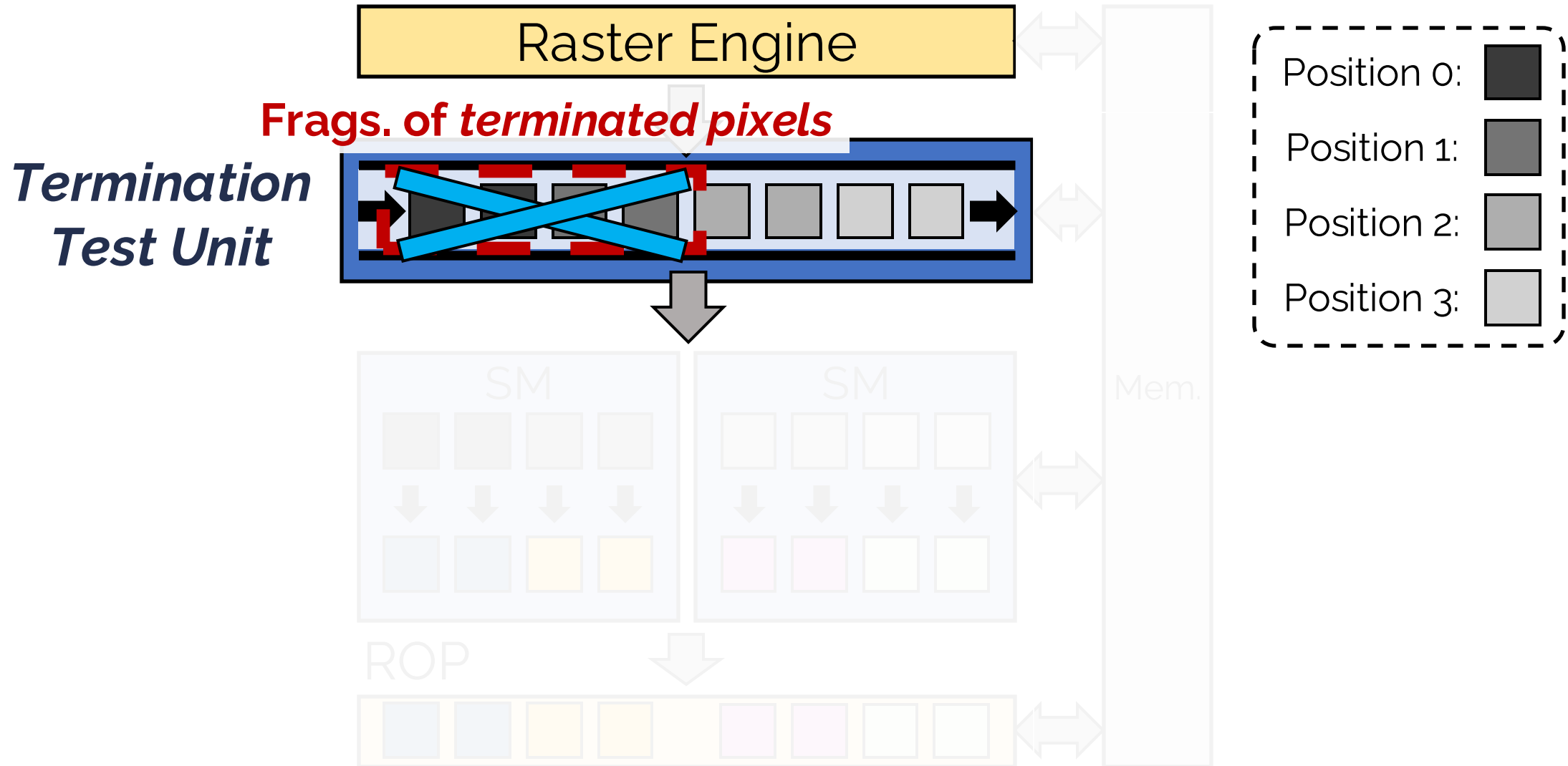
Fragment Queue

SM

SM

Mem.

ROP

Position 0:
Position 1:
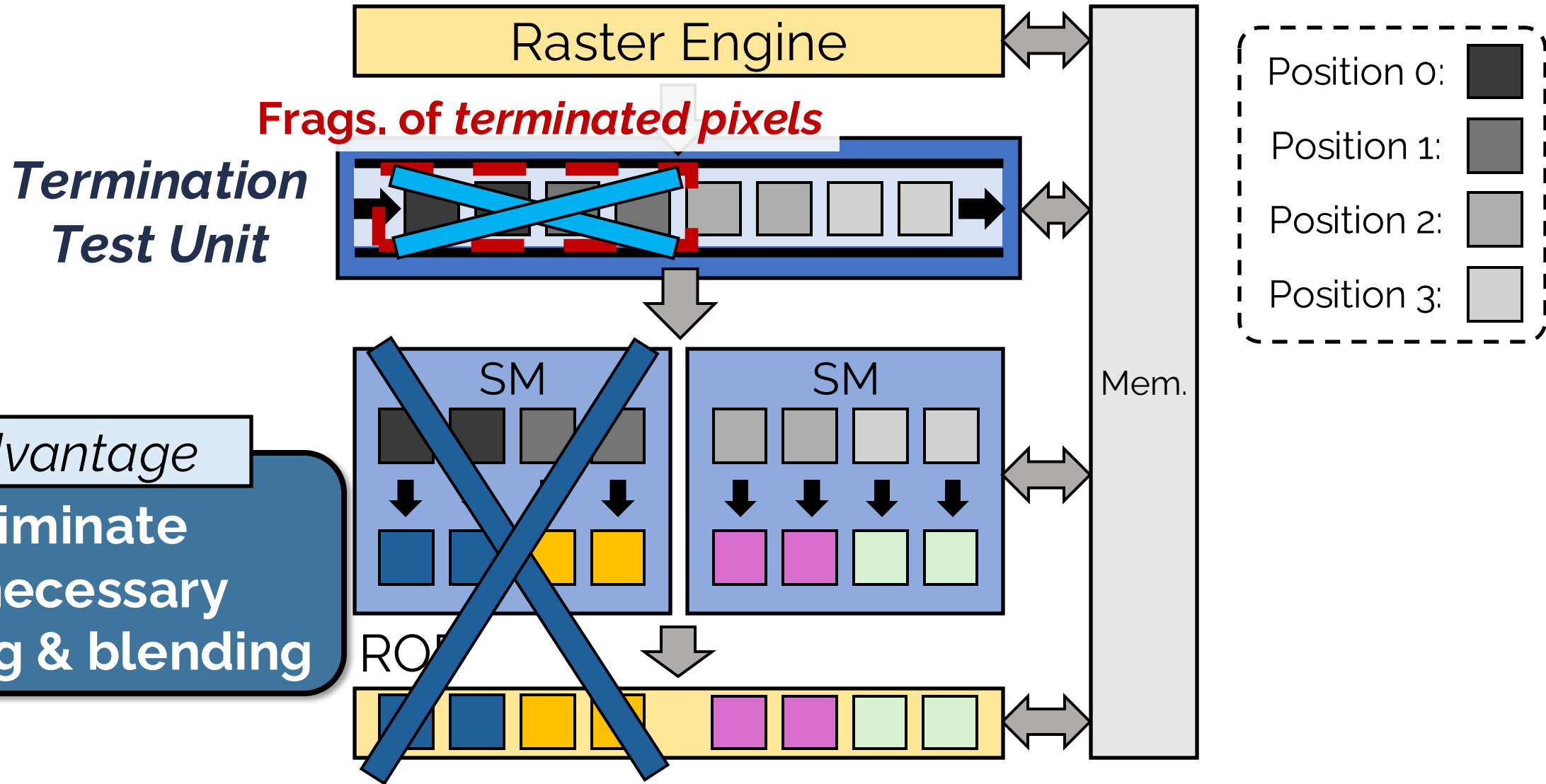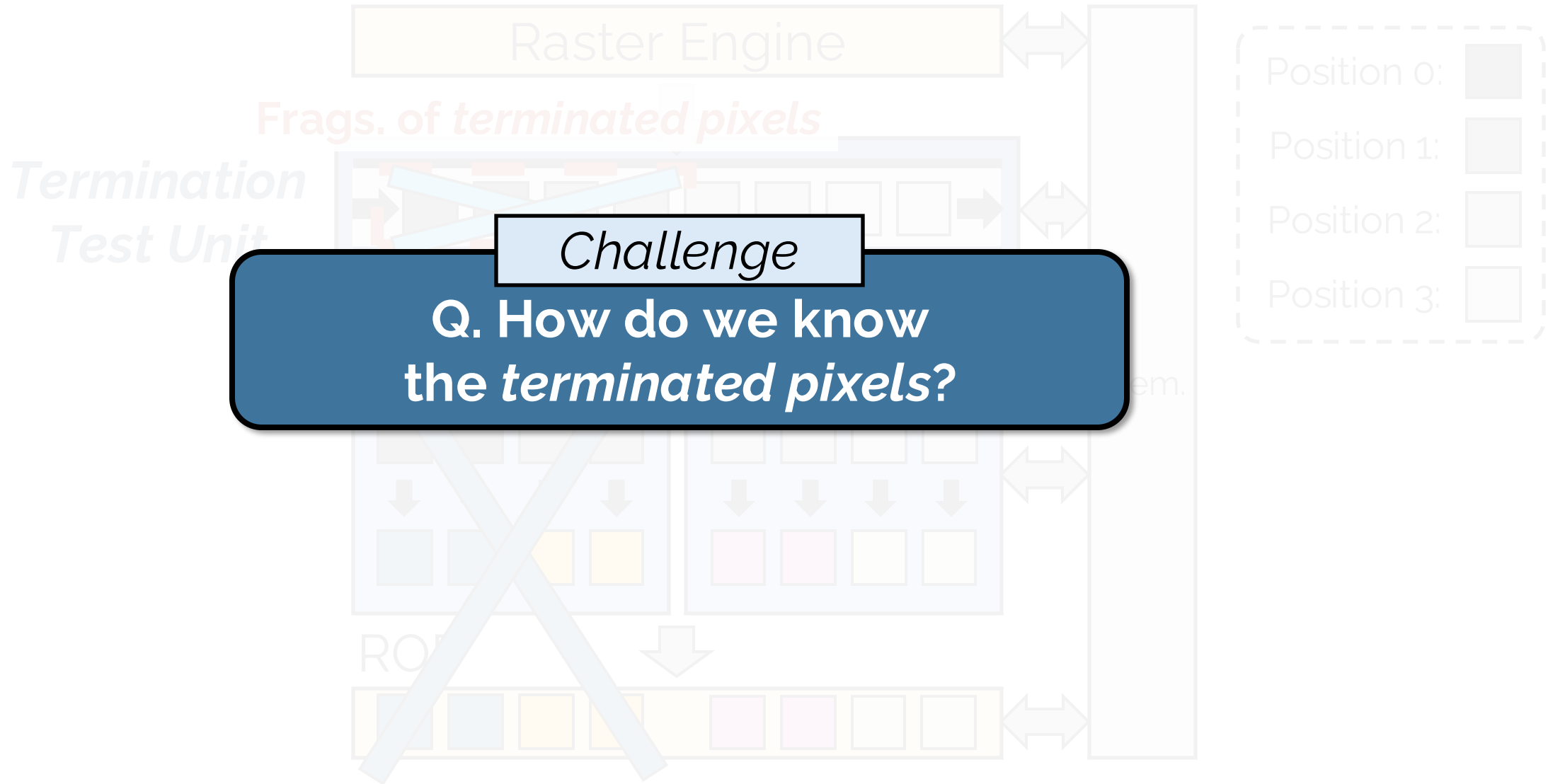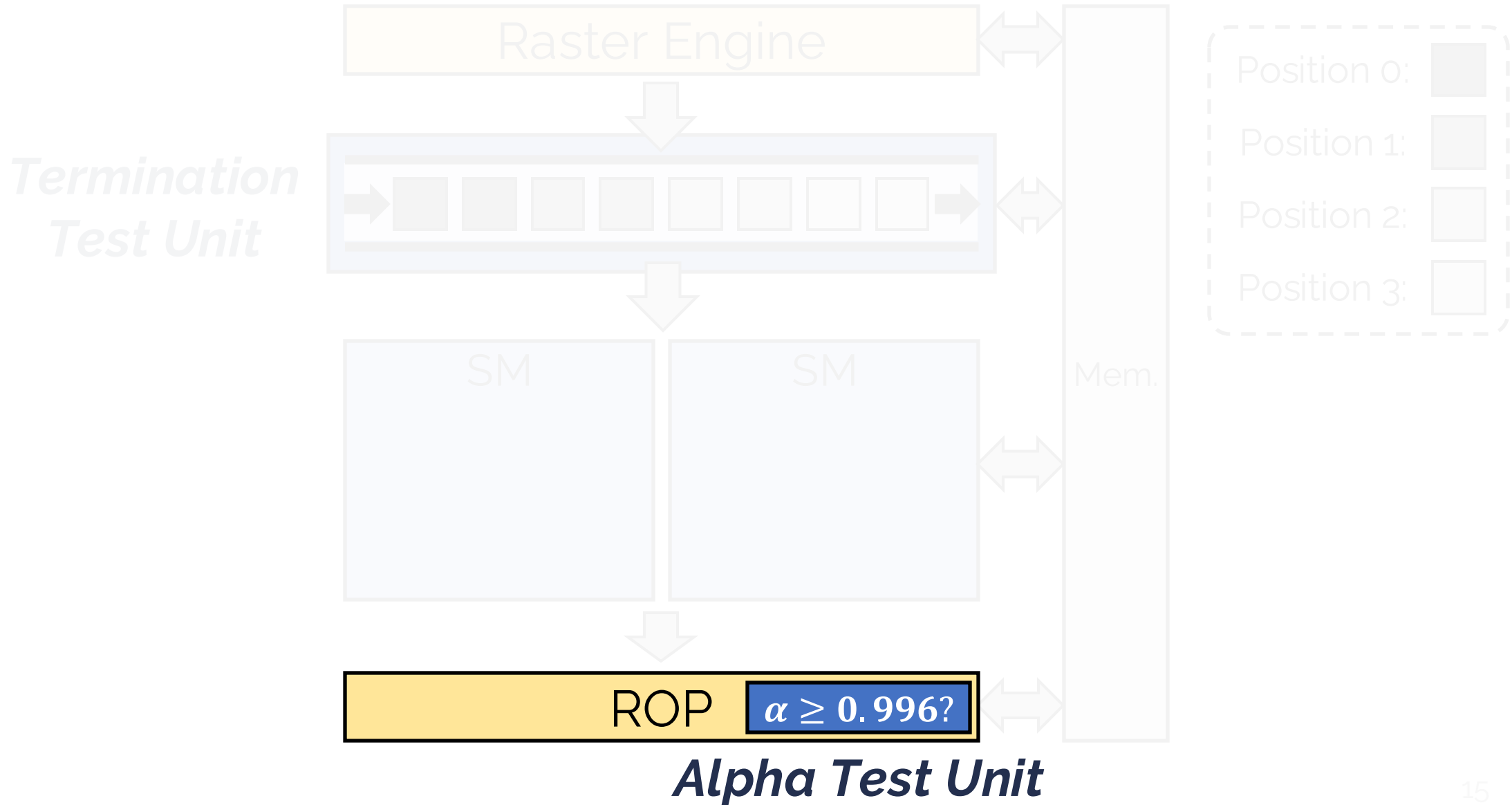Position 2:
Position 3:

# Hardware-Based Early Termination

# Hardware-Based Early Termination

# Hardware-Based Early Termination

Raster Engine

Frags. of *terminated pixels*

*Termination Test Unit*

Position 0:
Position 1:
Position 2:
Position 3:

**Challenge**

**Q. How do we know the *terminated pixels*?**

ROP

14

# Hardware-Based Early Termination



Raster Engine

*Termination Test Unit*

SM

SM

Mem.

ROP $\alpha \geq 0.996$?

*Alpha Test Unit*

Position 0:
Position 1:
Position 2:
Position 3:

# Hardware-Based Early Termination



**Termination Test Unit**

Raster Engine

Mem.

Position 0:
Position 1:
Position 2:
Position 3:

$1$

**Term. flag of Pix. 0**

SM          SM

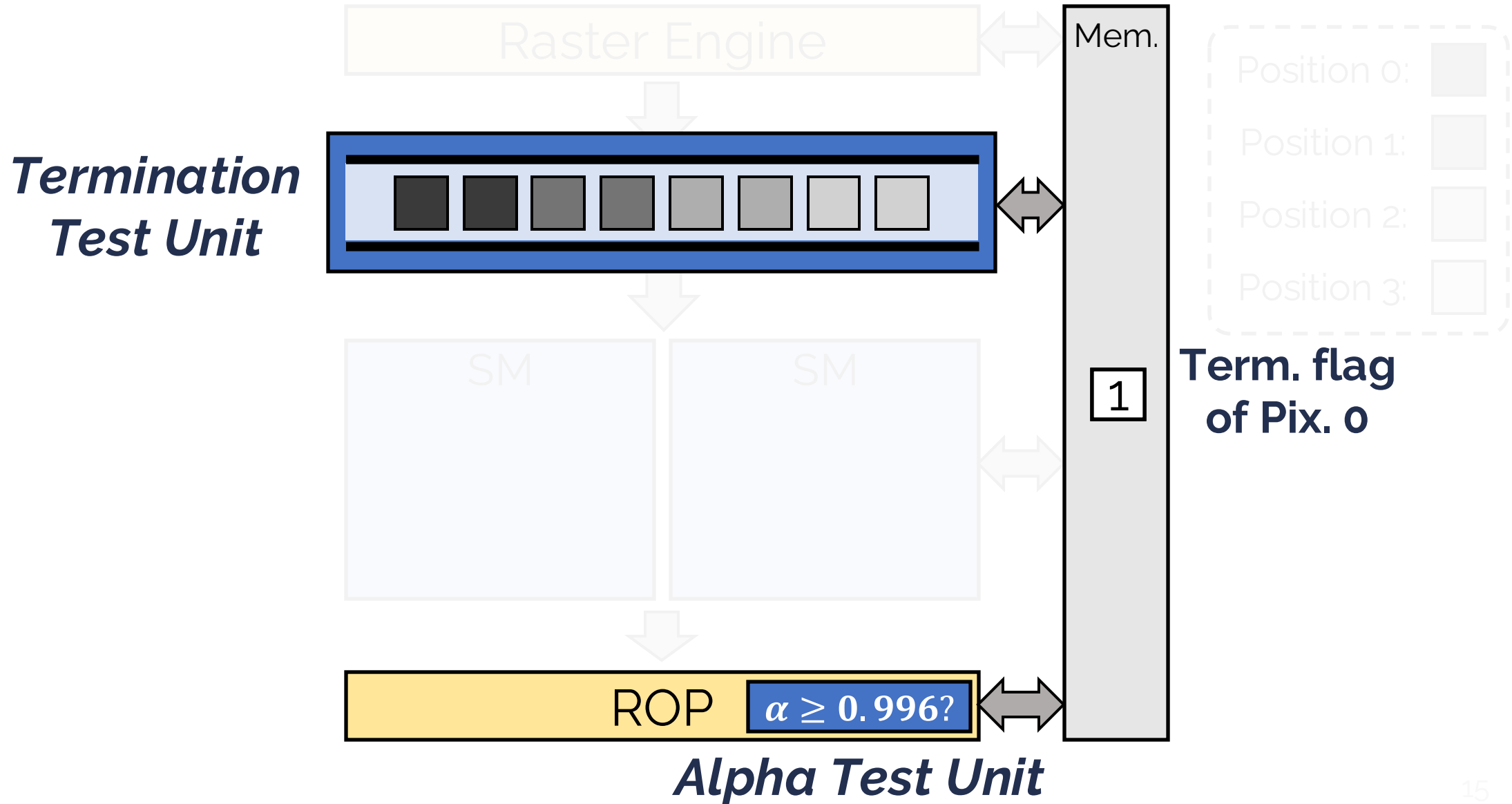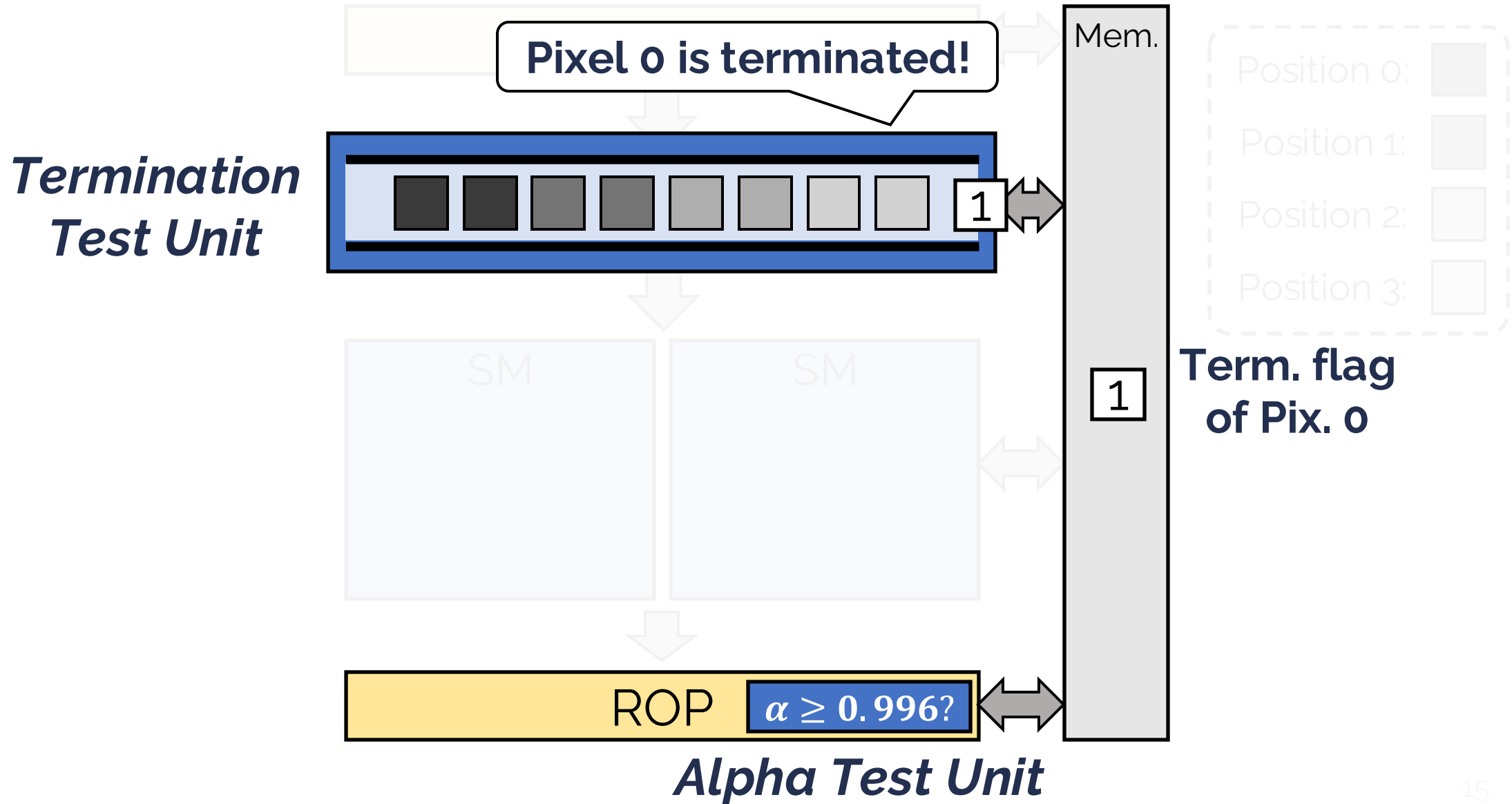ROP   $\alpha \geq 0.996?$
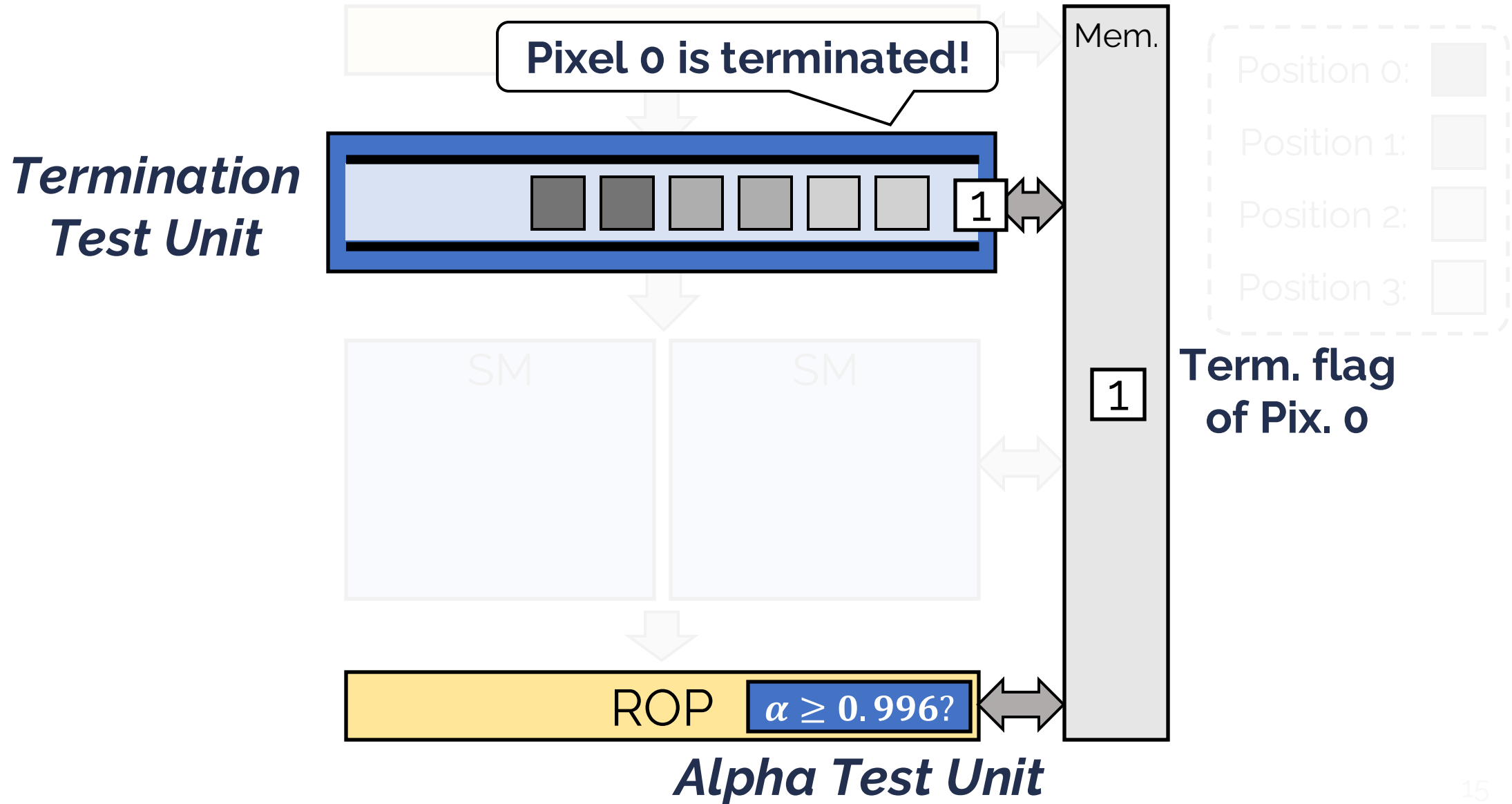
**Alpha Test Unit**

# Hardware-Based Early Termination

# Hardware-Based Early Termination

# Hardware-Based Early Termination



Raster Engine

Mem.

Position 0:

Position 1:

Position 2:

Position 3:

Termination Test Unit

#Pixels

Stencil Buffer

SM

SM

15

# Hardware-Based Early Termination

# Hardware-Based Early Termination



**Repurpose one bit** of a stencil value as a **termination flag**

Mem.

#Pixels

Stencil Buffer

*Unused*

Raster Engine

*Termination Test Unit*

Position 0:
Position 1:
Position 2:
Position 3:

15

# Quad Merging: Key Insight

# Quad Merging: Key Insight

# Quad Merging: Key Insight

# Quad Merging: Key Insight

# Quad Merging: Key Insight



**Partially blend the fragments** using underutilized SMs

# Quad Merging: Challenge

# Quad Merging: Challenge

# Quad Merging: Challenge

# Quad Merging: Challenge

# Quad Merging: Challenge

# Quad Merging: Challenge

# Quad Merging

# Quad Merging

## Quad Reorder Unit
### 1) Reorder the quads

# Quad Merging

**Quad Reorder Unit**
***1) Reorder the quads***

***2) Partially blend using warp shuffling***

Raster Engine

Warp 0        Warp 1

SM        SM

Mem.

ROP

Position 0:
Position 1:
Position 2:
Position 3:

# Quad Merging



**Quad Reorder Unit**

**1) Reorder the quads**

**2) Partially blend using warp shuffling**

Raster Engine

Warp 0    Warp 1

Position 0:
Position 1:
Position 2:
Position 3:

Mem.

*Advantage*

**Reduce the quads**
to be blended in ROPs ☺

ROP

# Outline

- **Background**
  - 3D Gaussian Splatting (3DGS)
  - Hardware Graphics Pipeline

- **Limitations of Graphics Hardware**

- **VR-Pipe: Graphics Hardware Extension for Volume Rendering**
  - Quad Merging with Multi-Granular Tile Binning
  - Hardware Support for Early Termination

- **Evaluation**

- **Conclusion**

# Experimental Setup

## Performance Evaluation

- Emerald (ISCA' 19)
  - Cycle-level simulator w/ graphics hardware modeling based on GPGPU-sim and gem5
  - With extensive modifications based on our analysis

## Workloads

- Mip-NeRF 360: Kitchen, Bonsai
- Tanks & Temples: Train, Truck
- Synthetic-NeRF: Lego
- Synthetic-NSVF: Palace

## Baseline GPU Configuration

| | |
|---|---|
| # GPC | 1 |
| # SMs | 16 (1024 CUDA Cores) |
| Core Frequency | 612 MHz |
| L1D/T | 48KB, 128B line |
| Shared L2 | 4MB, 128B line (sectored) |
| ROP Cache | 16KB, 128B line (sectored) |
| ROP Throughput | 2 quads/cycle (RGBA16F) |
| DRAM | LPDDR3-1600 (16-channel) |

# Performance

QM: Quad Merging
HET: Hardware-based Early Termination

# Performance



QM: Quad Merging
HET: Hardware-based Early Termination

Legend: Baseline, QM, HET, HET+QM

Categories: Kitchen, Bonsai, Train, Truck, Lego, Palace, Geomean

Speedup (y-axis: 0 to 3)

1.35x, 1.80x

# Performance

QM: Quad Merging
HET: Hardware-based Early Termination



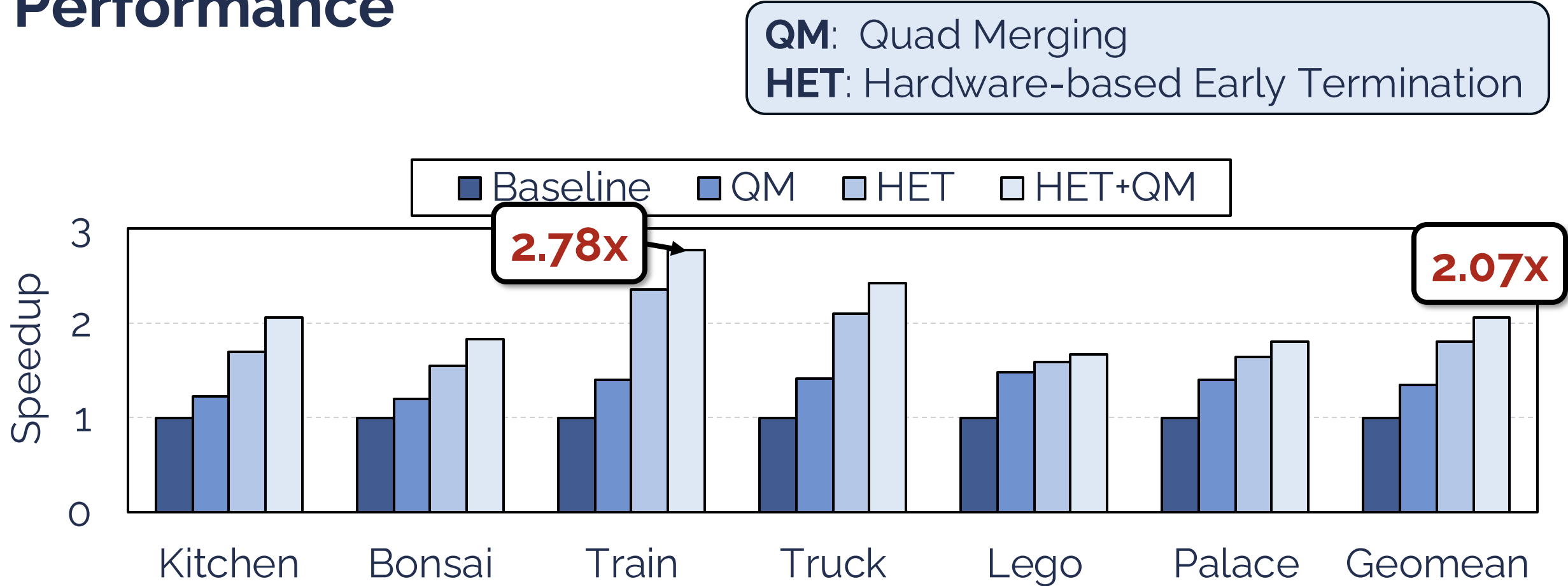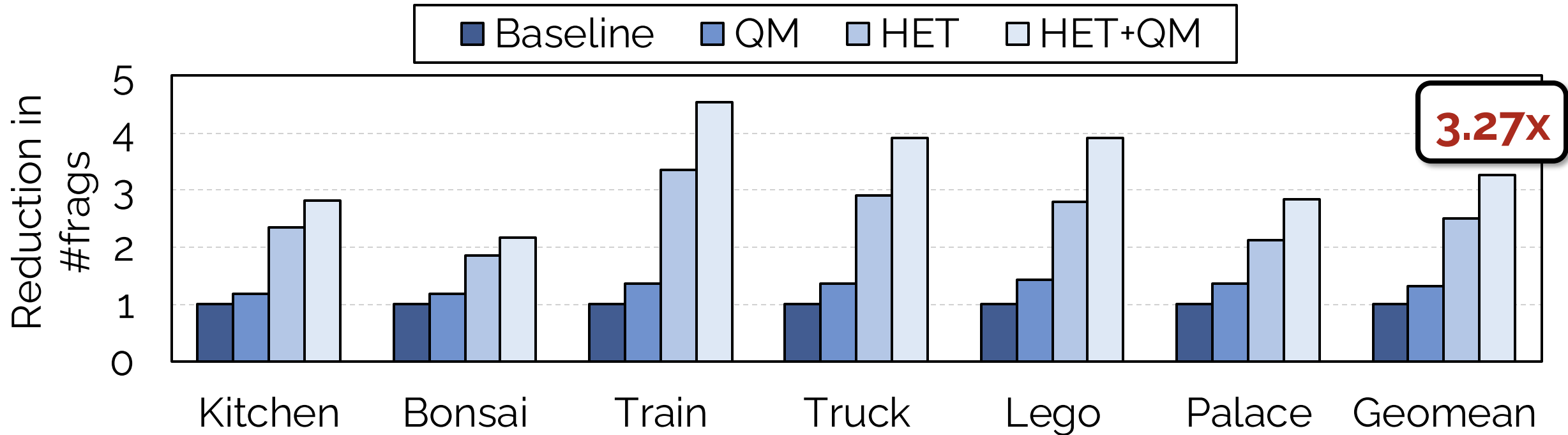**VR-Pipe greatly improves rendering performance**
With minimal hardware overhead in a GPU
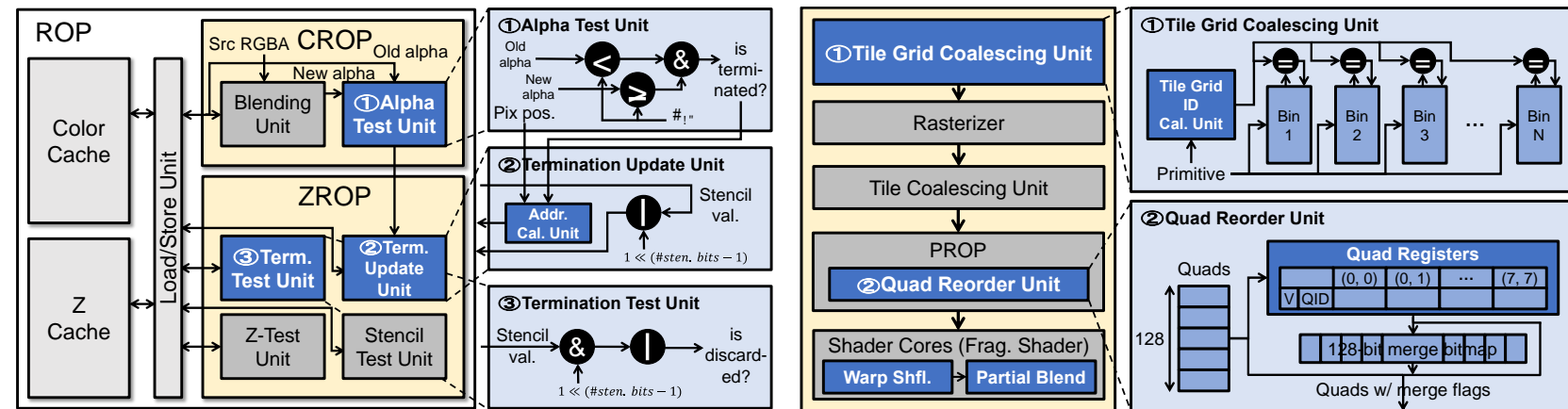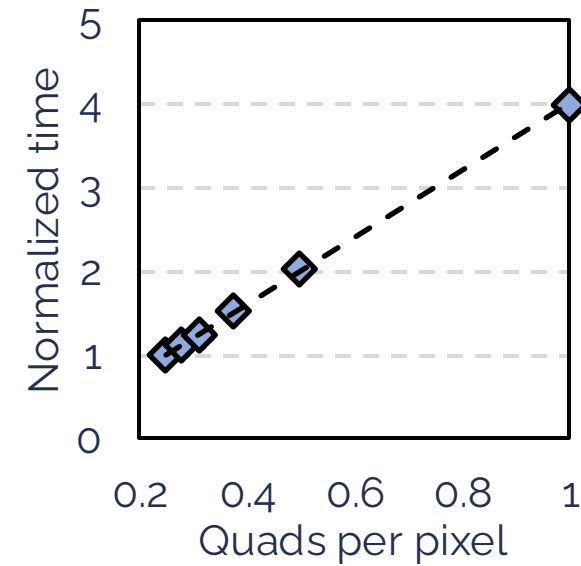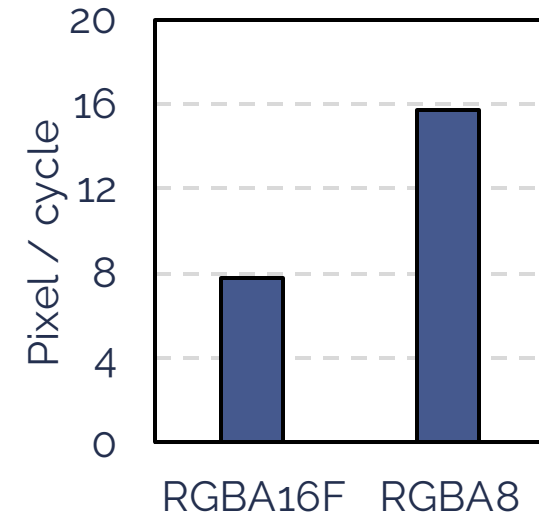
# Source of Performance Gain

Reduction in the Number of Fragments



**VR-Pipe significantly reduces the number of fragments**
blended by ROP

# More Details in Our Paper

- Analysis on Real Graphics Hardware

- Limitations of SW-based Optimizations

- Hardware Implementation Cost

- Details of Proposed Microarchitecture

- Others...

# Conclusion

## Problem

- **High ROP pressure** for blending a number of fragments per pixel
- Lack of native hardware support for early termination

## Solution: VR-Pipe, a GPU hardware extension for volume rendering

- Hardware-based early termination to early-discard the fragments
- Quad merging with multi-granular tile binning to exploit underutilized SMs

## Result

- VR-Pipe achieves up to a **2.78x speedup** over the conventional graphics pipeline with mininal hardware overhead! ☺

# Thank You!

## VR-Pipe

Streamlining Hardware
Graphics Pipeline for
Volume Rendering

**Junseo Lee (junseo.lee@snu.ac.kr)**