# **NeuRex**: A Case for **Neu**ral **Re**ndering **Acc**eleration
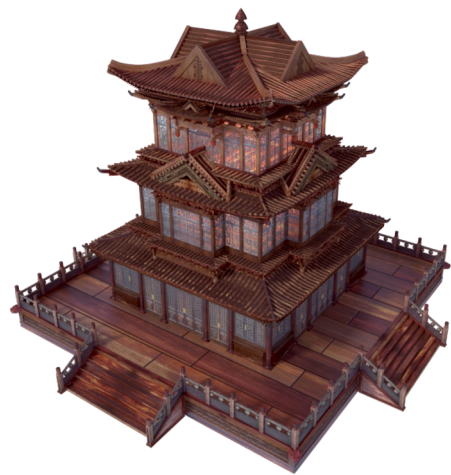
**Junseo Lee**   Kwanseok Choi   Jungi Lee
Seokwon Lee   Joonho Whangbo   Jaewoong Sim
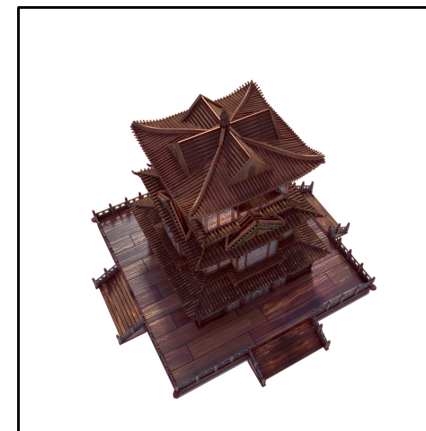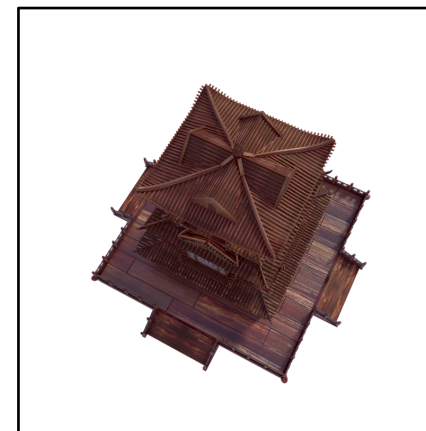
Seoul National University

# Neural Rendering

# Neural Rendering
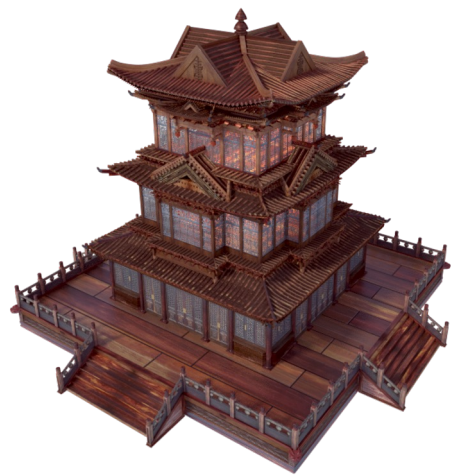
# Neural Rendering

# Neural Rendering

# Neural Rendering

# Neural Rendering

# Neural Rendering



**NeRF Rendering Process**

# NeRF Rendering Process

# NeRF Rendering Process

# NeRF Rendering Process



NeRF Rendering Process

Sampled Points

# NeRF Rendering Process



Sampled Points

Positional Encoding

# NeRF Rendering Process



NeRF Rendering Process

Sampled Points

Positional Encoding

MLP

Color & Density

# NeRF Rendering Process



## Instant-NGP

Sampled Points

Positional Encoding

MLP

Color & Density

# NeRF Rendering Process



NeRF Rendering Process

Instant-NGP

Sampled Points

Multi-resolution Hash Encoding

L0  L1  L2  ...  L14  L15

MLP

Color & Density

# NeRF Rendering Process

## Instant-NGP

Sampled Points

Multi-resolution Hash Encoding

L₀  L₁  L₂  ...  L₁₄  L₁₅

MLP

**1) Still too slow! ☹**

Color & Density

FHD (1920x1080)

Rendering Time (ms)

800

600

400 — **0.5s**  **16.1s**

200

0

RTX 3070    Jetson NX

# NeRF Rendering Process



**NeRF Rendering Process**

## Instant-NGP

Sampled Points

1) Still **too slow**! ☹

2) **Encoding lookup** is bottleneck!

**FHD (1920x1080)**

Rendering Time (ms)

800 · 600 · 400 · 200 · 0

**0.5s** — RTX 3070
**16.1s** — Jetson NX

Percentage (%)

100 · 80 · 60 · 40 · 20 · 0

Others / MLP / ENC — RTX 3070 — **41%**
Others / MLP / ENC — Jetson NX — **61%**

# Execution Flow

# Execution Flow

## *Multi-resolution Hash Encoding*



**Pos.**

$N_{\text{point}}$

$h$

$h$

$h$

$L_0$

$L_1$

...

$L_{15}$

$lerp$

Encoding
Tables

Encoded
Features

Density
MLP

Color
MLP

Output
RGB

$N_{\text{point}}$

# Execution Flow



Pos.

$N_{point}$

$h$

$h$

$h$

$L_0$

$L_1$

$L_{15}$

*lerp*

Encoding
Tables

Encoded
Features

*MLP*

Density
MLP

Color
MLP

Output
RGB

$N_{point}$

# Execution Flow



*Multi-resolution Hash Encoding*

# Multi-resolution Hash Encoding

3D Scene (Bounding Box)

# Multi-resolution Hash Encoding

3D Scene (Bounding Box)

# Multi-resolution Hash Encoding

3D Scene (Bounding Box)

Coarse resolution

Fine resolution

# Multi-resolution Hash Encoding

3D Scene (Bounding Box)

Coarse resolution

Fine resolution

Lev. 0

...

Lev. 15

Hash Tables

# Multi-resolution Hash Encoding



3D Scene (Bounding Box)

Coarse resolution

Fine resolution

Lev. 0

...

Lev. 15

Hash Tables

# Multi-resolution Hash Encoding

3D Scene (Bounding Box)

Coarse resolution

Fine resolution

Lev. 0

Lev. 15

Hash Tables

# Multi-resolution Hash Encoding

3D Scene (Bounding Box)

Coarse resolution

Fine resolution

Lev. 0

...

Lev. 15

Hash Tables

# Multi-resolution Hash Encoding



3D Scene (Bounding Box)

Coarse resolution

Fine resolution

Lev. 0

$h$

Lev. 15

Hash Tables

# Multi-resolution Hash Encoding

Lev. 0

3D Scene (Bounding Box)

Coarse resolution

Fine resolution

$h$

Lev. 15

Hash Tables

# Multi-resolution Hash Encoding

3D Scene (Bounding Box)

Coarse resolution

Fine resolution

$h$

Lev. 0

Interpolate

Encoded Feature

Lev. 15

Hash Tables

# Multi-resolution Hash Encoding



3D Scene (Bounding Box)

Coarse resolution

Fine resolution

Lev. 0

$h$

Lev. 15

$h$

Hash Tables

Interpolate

Encoded Feature

Interpolate

# Outline

- **Background**
  - Neural Rendering
  - Multi-resolution Hash Encoding

- **I-NGP Optimization & Limitations**

- **NeuRex: Efficient Neural Rendering Accelerator**
  - Restricted Hashing
  - Neural Graphics Engine

- **Evaluation**

- **Conclusion**

# I-NGP Optimization



Hash Table Lookup

Pos.

$L_0$

$h$

$lerp$

$L_1$

$h$

$L_{15}$

$h$

Hash Tables

Encoded Features

MLP

# I-NGP Optimization

Hash Table
Lookup

**Pos.**

$h$

$L_0$

$lerp$

Encoded
Features

$h$

$L_1$

$h$

$L_{15}$

Hash
Tables

MLP

# I-NGP Optimization

Hash Table
Lookup

**Pos.**

$L_0$

$lerp$

$h$

$L_1$

$h$

$h$

$L_{15}$

Hash
Tables

Encoded
Features

MLP

# I-NGP Optimization

Hash Table
Lookup

**Pos.**

$L_0$

$lerp$

$h$

$L_1$

$h$

$h$

$L_{15}$

Hash
Tables

Encoded
Features

MLP

**Start!**

# I-NGP Optimization

**High-end GPU**

**L2 Cache (> 4MB)**

$L_{15}$

- L2 cache in high-end GPUs

  > Single hash table

- Implication
  - Minimize off-chip memory access

# Limitations of I-NGP Optimization

Level-wise execution



Pos.

$L_0$

$L_1$

$L_{15}$

Hash Tables

Encoded Features

MLP

# Limitations of I-NGP Optimization

$L_0$

Encoded Features

Pos

**Problem 1. No Performance Portability**
: L2 cache of edge GPU < Single hash table

MLP

$L_{15}$

Hash
Tables

# Limitations of I-NGP Optimization

**Problem 1. No Performance Portability**
: L2 cache of edge GPU < Single hash table

**Problem 2. Resource Underutilization**
: ENC → MLP

# Outline

- **Background**
  - Neural Rendering
  - Multi-resolution Hash Encoding

- **I-NGP Optimization & Limitations**

- **NeuRex: Efficient Neural Rendering Accelerator**
  - Restricted Hashing
  - Neural Graphics Engine

- **Evaluation**

- **Conclusion**

# NeuRex: Hardware-Friendly Hash Encoding

# NeuRex: Hardware-Friendly Hash Encoding

## Original Algorithm

- Access range: **Entire hash table** ☹



Hash Table

# NeuRex: Hardware-Friendly Hash Encoding

## Original Algorithm
- Access range: **Entire hash table** ☹

## Restricted Hashing
- Access range: **Small sub-table** ☺



Hashing

vs.

Hash Table

# NeuRex: Hardware-Friendly Hash Encoding

## Original Algorithm
- Access range: **Entire hash table** ☹

## Restricted Hashing
- Access range: **Small sub-table** ☺



VS.

Hashing

Hash Table

**Subgrid 1**

**Subgrid 6**

| Subtable 0 |
| Subtable 1 |
| Subtable 2 |
| Subtable 3 |
| Subtable 4 |
| Subtable 5 |
| Subtable 6 |
| Subtable 7 |

Hash Table

# NeuRex: Hardware-Friendly Hash Encoding

## Original Algorithm
- Access range: **Entire hash table** ☹

## Restricted Hashing
- Access range: **Small sub-table** ☺



vs.

Hash Table

Subgrid 1

Subgrid 6

| Subtable 0 |
| Subtable 1 |
| Subtable 2 |
| Subtable 3 |
| Subtable 4 |
| Subtable 5 |
| Subtable 6 |
| Subtable 7 |

Hash Table

# NeuRex: Hardware-Friendly Hash Encoding

**Advantage 1. Performance Portability**
: Small on-chip memory is enough



**Off-Chip Memory**

L0  L1  L2  ...  L15
                 Sub

**Accelerator**

Small On-Chip Buffer

# NeuRex: Hardware-Friendly Hash Encoding

**Advantage 1. Performance Portability**
: Small on-chip memory is enough

**Off-Chip Memory**

L0   L1   L2   ...   L15

**Accelerator**

Small
On-Chip
Buffer

Sub

# Advantage 2. Parallelized Execution of ENC and MLP

Subgrid-wise execution

**Subgrid 0**

| Pts. 0 |
| Pts. 1 |
| Pts. 2 |
| Pts. 3 |

| Pts. 4 |
| Pts. 5 |
| Pts. 6 |
| Pts. 7 |

**Subgrid 1**

# Advantage 2. Parallelized Execution of ENC and MLP

Subgrid-wise execution

Subgrid 0

| Pts. 0 |
| Pts. 1 |
| Pts. 2 |
| Pts. 3 |

| Pts. 4 |
| Pts. 5 |
| Pts. 6 |
| Pts. 7 |

Subgrid 1

Lev. 0

| Sub-table 0 |
| Sub-table 1 |

Lev. 1

| Sub-table 0 |
| Sub-table 1 |

⋮

Lev. 15

| Sub-table 0 |
| Sub-table 1 |

Sub
Tables

# Advantage 2. Parallelized Execution of ENC and MLP

Subgrid-wise execution

# Advantage 2. Parallelized Execution of ENC and MLP

Subgrid-wise execution

# Advantage 2. Parallelized Execution of ENC and MLP

Subgrid-wise execution

# Advantage 2. Parallelized Execution of ENC and MLP

Subgrid-wise execution

# NeuRex: Neural Graphics Engine

# NeuRex: Neural Graphics Engine



## Encoding Engine (EE)

- Perform multi-resolution hash encoding

## Tensor Compute Engine (TCE)

- Perform MLP execution

# NeuRex: Neural Graphics Engine



## Encoding Engine (EE)

- Perform multi-resolution hash encoding

# NeuRex: Neural Graphics Engine



## Index Generation Unit (IGU)

- Calculate hash indices
- Compute interpolation weights

# NeuRex: Neural Graphics Engine



## Encoding Lookup Unit (ELU)

- Fetch hash table entries from on-chip buffers

# NeuRex: Neural Graphics Engine



**Interpolation Compute Unit (ICU)**

- Interpolate 8 hash entries

# NeuRex: Neural Graphics Engine



## Encoding Lookup Unit (ELU)

- Fetch hash table entries from on-chip buffers

# Encoding Lookup Unit

# Encoding Lookup Unit

## Grid Cache

- A cache for **coarse levels**' hash encoding (L=0,1,2,…)
- Cache block: Coalesced 8 hash entries

**GC Tag Bank**   **GC Data Bank (8×4B hash table entries)**

| GC Tag Bank | GC Data Bank |
|---|---|
| tag 0 | entry 0 \| entry 1 \| entry 2 \| ⋯ \| entry 7 |
| tag 1 | entry 0 \| entry 1 \| entry 2 \| ⋯ \| entry 7 |
| ⋮ | ⋮ |
| tag $K$-1 | entry 0 \| entry 1 \| entry 2 \| ⋯ \| entry 7 |

# Encoding Lookup Unit

## Grid Cache

- A cache for **coarse levels**' hash encoding (L=0,1,2,…)
- Cache block: Coalesced 8 hash entries

**Coarse Level**

**Fine Level**

**GC Tag Bank**

| tag 0 |
|:-:|
| tag 1 |
| ⋮ |
| tag $K$-1 |

**GC Data Bank (8×4B hash table entries)**

| entry 0 | entry 1 | entry 2 | ⋯ | entry 7 |
|:-:|:-:|:-:|:-:|:-:|
| entry 0 | entry 1 | entry 2 | ⋯ | entry 7 |
| ⋮ | | | | |
| entry 0 | entry 1 | entry 2 | ⋯ | entry 7 |

# Encoding Lookup Unit

## Grid Cache

- A cache for **coarse levels**' hash encoding (L=0,1,2,…)
- Cache block: Coalesced 8 hash entries

| GC Tag Bank | GC Data Bank (8×4B hash table entries) | | | | |
|---|---|---|---|---|---|
| tag 0 | entry 0 | entry 1 | entry 2 | ... | entry 7 |
| tag 1 | entry 0 | entry 1 | entry 2 | ... | entry 7 |
| ⋮ | | | ⋮ | | |
| tag $K$-1 | entry 0 | entry 1 | entry 2 | ... | entry 7 |

**Coarse Level**

**Fine Level**

# Encoding Lookup Unit

Grid Cache
- A cache for **coarse levels**' hash encoding (L=0,1,2,…)
- Cache block: Coalesced 8 hash entries

Subgrid Buffer
- A buffer for **fine levels**' hash encoding (L=15,14,13,…)
- **Restricted hashing** reduces the required on-chip
  memory size for hash encoding! ☺

```
┌──────────────────────────────────────────┐
│                                          │
│   ┌──────────┐          ┌──────────┐     │
│   │          │          │ Subgrid  │     │
│   │ Subtable │   ➡      │ Buffer   │     │
│   │          │          │          │     │
│   └──────────┘          └──────────┘     │
│                                          │
└──────────────────────────────────────────┘
```
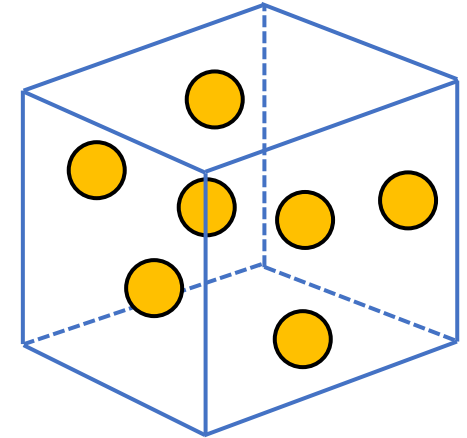
# Outline

- **Background**
  - Neural Rendering
  - Multi-resolution Hash Encoding

- **I-NGP Optimization & Limitations**

- **NeuRex: Efficient Neural Rendering Accelerator**
  - Restricted Hashing
  - Neural Graphics Engine

- **Evaluation**

- **Conclusion**

# Methodology

## RTL Implementation

- Process node: 28nm technology

## Performance Evaluation

- Cycle-level simulator

## Hardware Variants & Baselines

- NeuRex-Server <-> RTX 3070
  - Area: **21.37mm²** <-> 392.5mm²
- NeuRex-Edge <-> Xavier NX
  - Area: **3.14mm²** <-> 350mm²

## Evaluated Workloads

| Dataset | Scene (Resolution) | Type |
|---|---|---|
| Syn-NeRF | Mic (800x800) | Synthetic |
| Syn-NSVF | Palace (800x800) | |
| BlendedMVS | Fountain (768x576) | Real world |
| Tanks& Temples | Family (1920x1080) | |
| Instant-NGP | Fox (1080x1920) | |

# Performance: NeRF



**NeuRex-Server**

**NeuRex-Edge**

# Performance: NeRF



**~3.11x** (NeuRex-Server, Speedup: Mic, Palace, Fountain, Family, Fox)

**~9.88x** (NeuRex-Edge, Speedup: Mic, Palace, Fountain, Family, Fox)

# Performance: NeRF
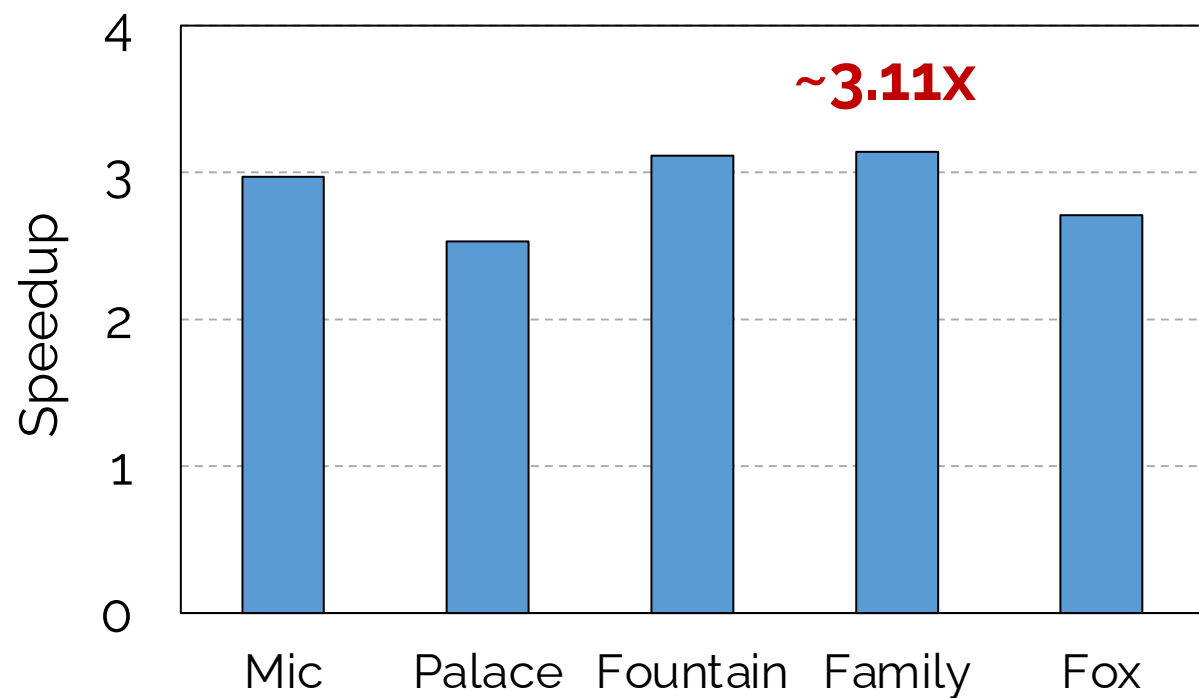


**~3.11x**

**NeuRex-Server**

**~9.88x**

**NeuRex-Edge**

Higher speedup of **NeuRex-Edge**
→ NeuRex enables the performance portability!

# Performance: NeRF

# Rendering Quality

# Rendering Quality

# Rendering Quality



* Higher is better

PSNR (dB) chart comparing Original (green) and Restricted Hashing (orange) across scenes: Mic, Palace, Fountain, Family, Fox.
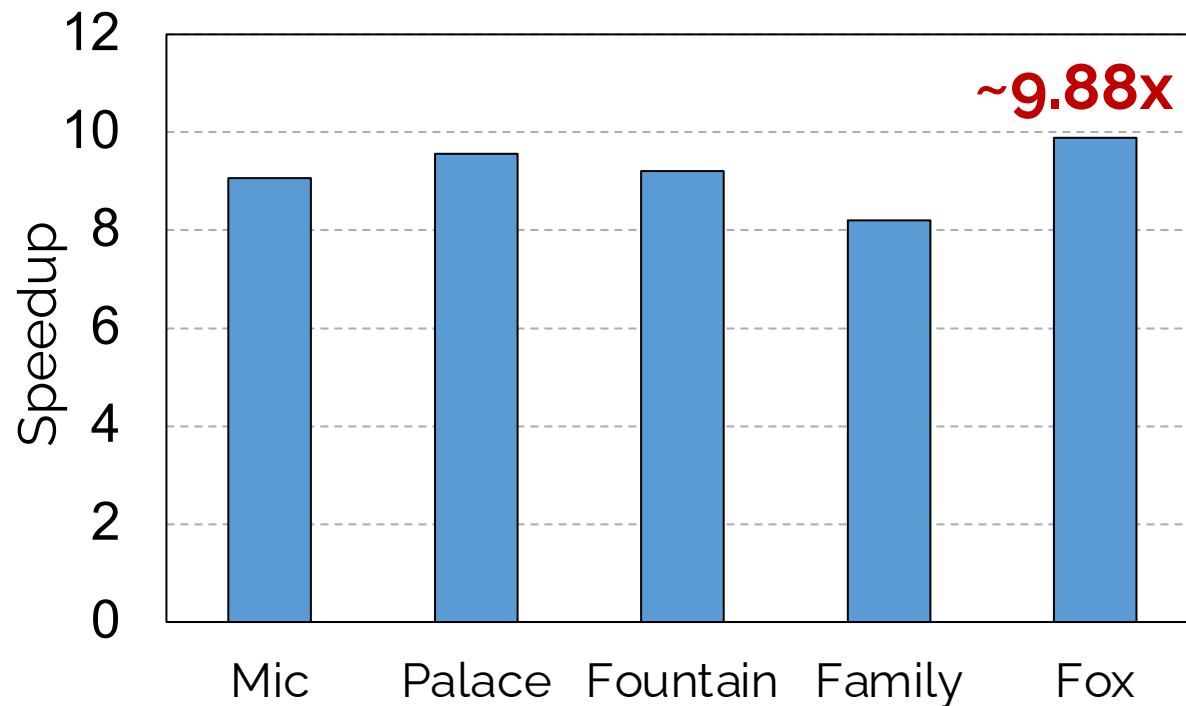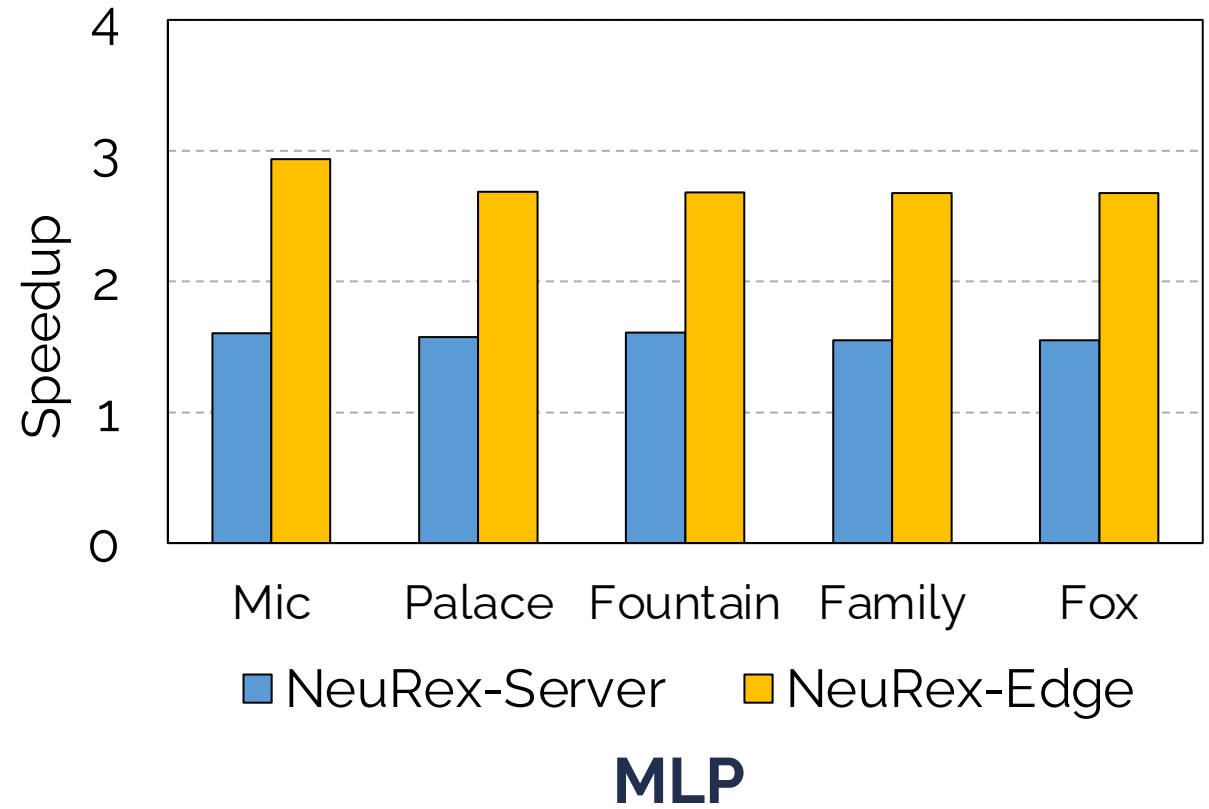
# Rendering Quality



**Reference**

Original Algorithm
(36.63dB)

Restricted Hashing
(36.66dB)

For some scenes, restricted hashing shows
**even better rendering quality** than the original! ☺

# More Details in Our Paper
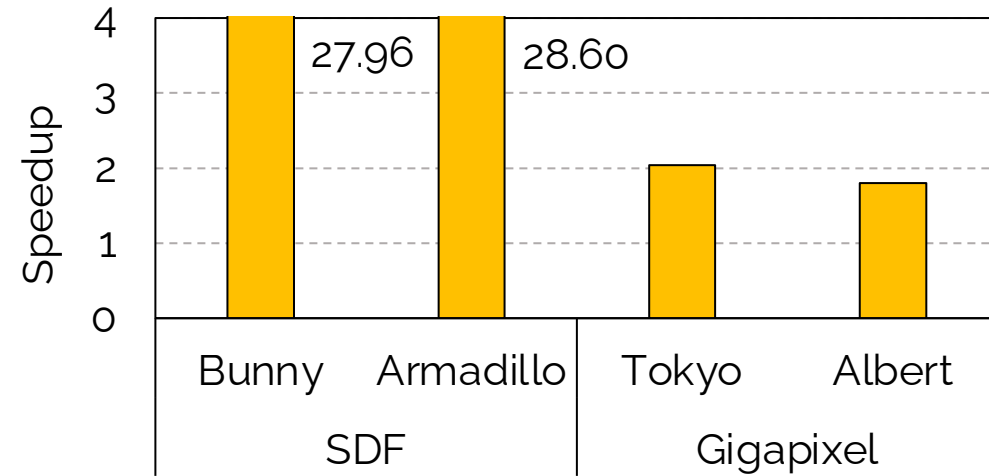
- Neural Rendering Tasks Beyond NeRFs
  - Signed Distance Functions (SDF)
  - 2D Image Approximation (Gigapixel)



- Source of Performance Gains
- Sensitivity Study
- Area and Energy Efficiency
- Others…

# Conclusion

Problem

- Multi-resolution hash encoding is a primary bottleneck in neural rendering with several limitations

Solution: **NeuRex,** an efficient neural rendering accelerator

- *Restricted Hashing* enables performance portability and maximizes resource utilization
- *Neural Graphics Engine* with two specialized on-chip memories

Result

- **NeuRex** achieves up to a 9.88x speedup over the GPU with a substantially small area overhead! ☺

# Thank You!