

# Dynamic 3D Gaussian Fields for Urban Areas

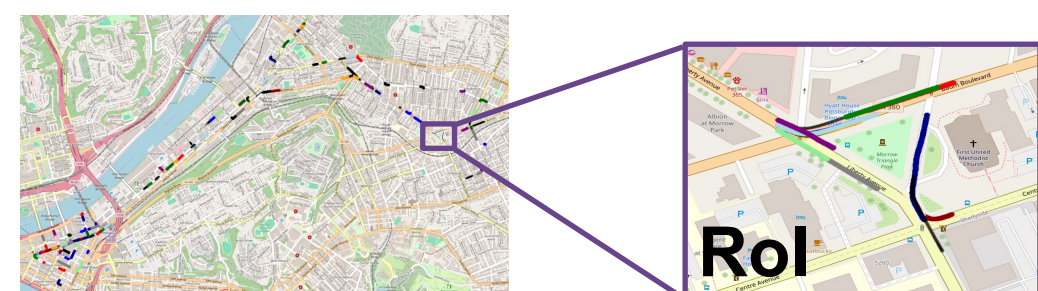


**TL;DR** Given a set of *heterogeneous* sequences of a shared geographic area, we optimize a *single dynamic* scene representation that renders arbitrary viewpoints and configurations at interactive speeds.

## 1. What?

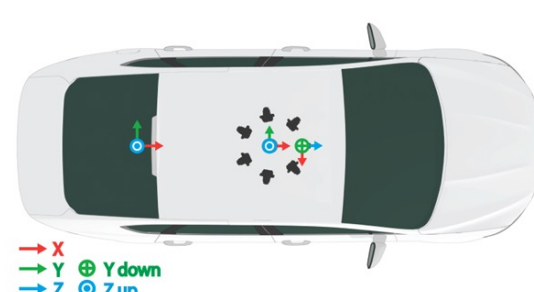
**Given:** Set of sequences  $S$  of a common geographic region of interest (RoI) with:

- Distinct dynamic objects
- Transient objects (construction sites, ...)
- Varying environmental conditions



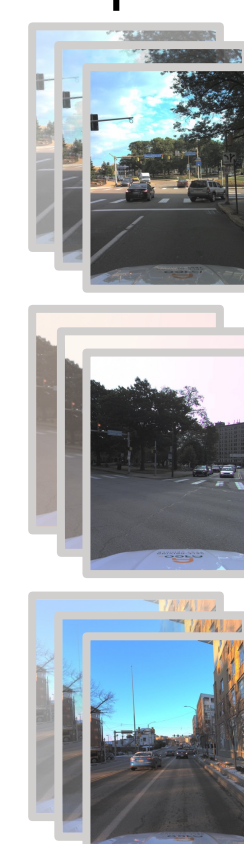
Each sequence  $s \in S$  has...

- Ego-vehicle poses  $\mathbf{P}$
- Camera calibration  $\mathbf{T}, \mathbf{K}$
- Dynamic object poses  $\xi$



**Goal:** Learn function  $f_\theta$  that, for given viewpoint  $(\mathbf{T}, \mathbf{K})$  at sequence  $s$  and time  $t$ , outputs the image  $\mathcal{I} \in [0,1]^{H \times W \times 3}$  with correct appearance & dynamic actors

Inputs



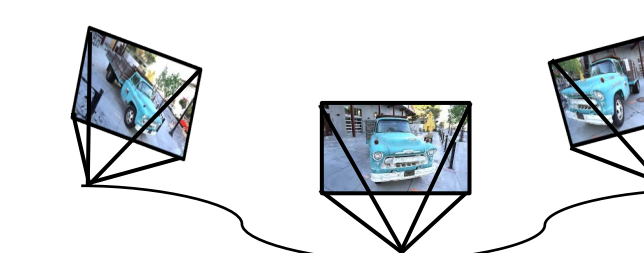
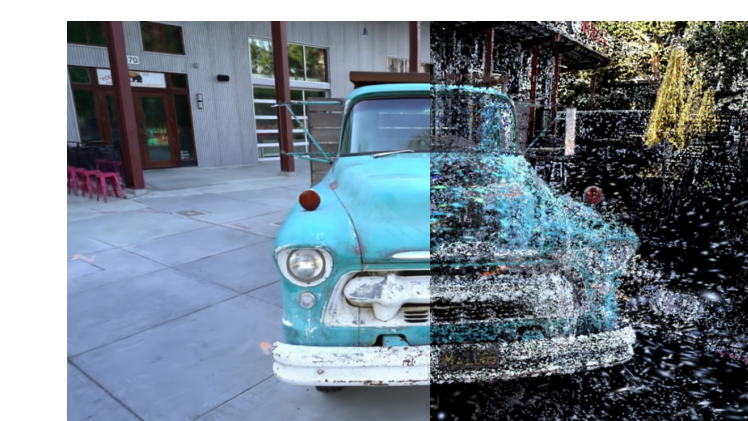
$$f_\theta(\mathbf{T}, \mathbf{K}, t, s) \rightarrow \mathcal{I} \in [0,1]^{H \times W \times 3}$$



## 2. Why?

3D Gaussian Splatting achieved high-quality novel view synthesis at high speed, but:

1. Limited **scalability**: Memory footprint increases linearly with number of primitives
2. Requires **homogeneous input data**: Spherical harmonics cannot model large appearance changes
3. Limited to **static** scenes: 3D primitives have a fixed location



## 3. How?

1. **Gaussians** as efficient **geometry scaffold**
2. **Neural fields** as a compact and flexible **appearance model**
3. **Scene dynamics** via scene **graph** at **global** level and **deformations** at **local** level.

**Scene Graph**  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$   
Stores scene **configurations**

**Nodes**  $\mathcal{V}$ : Latent codes  $\omega$   
Sequence code  $\omega_s^t$   
Object ID  $\omega_o$

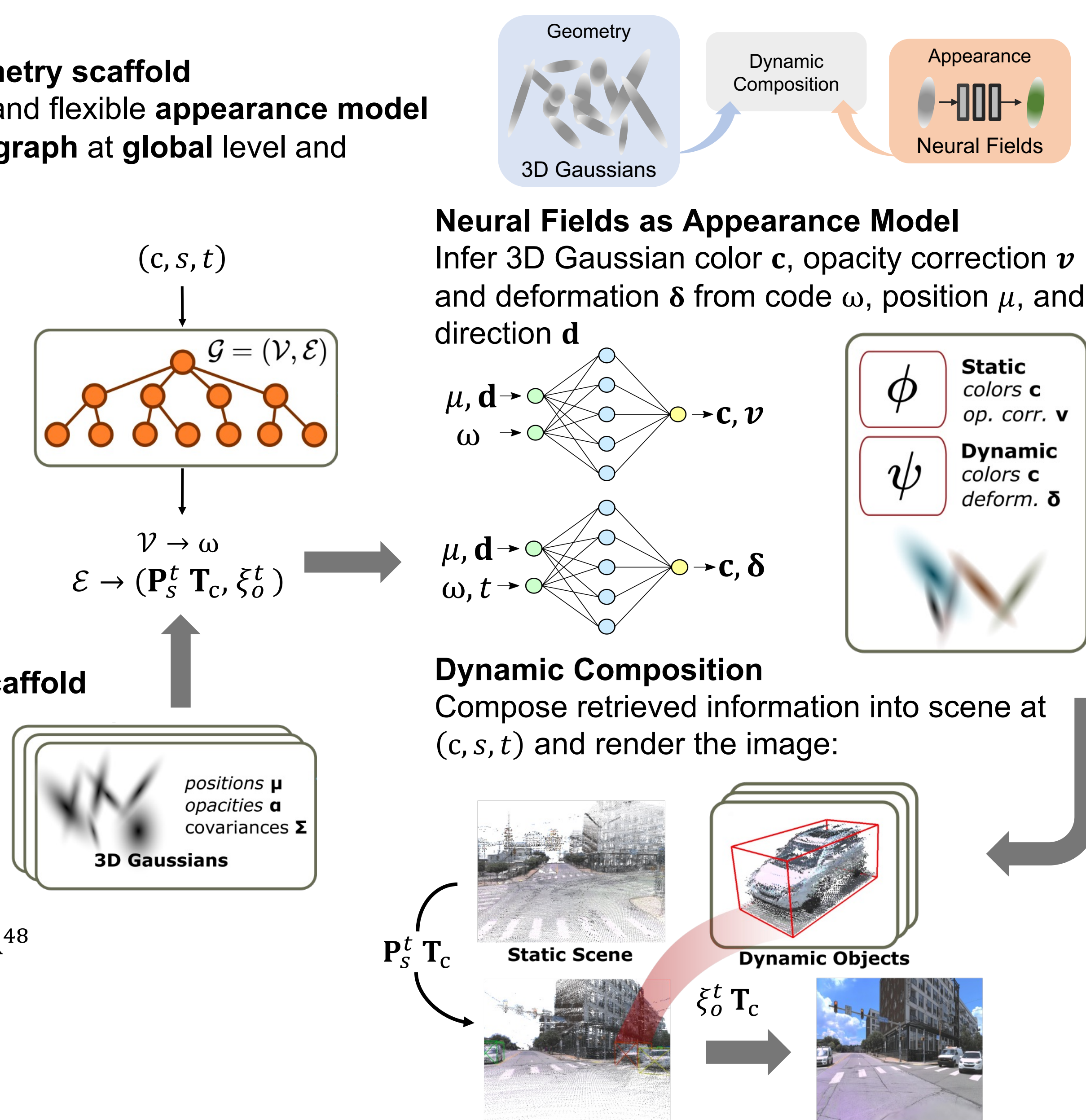
**Edges**  $\mathcal{E}$ : Rigid transformations  
Ego-vehicle poses  $\mathbf{P}$   
Dynamic object poses  $\xi$   
Camera to ego-vehicle  $\mathbf{T}$

**3D Gaussians as Geometry Scaffold**  
Multiple sets of 3D Gaussians

- **Static**  $G_r$   
World space
- **Dynamic**  $\{G_o \in O_s\}$   
Canonical space

Remove SH-Coefficients  $c_{sh} \in \mathbb{R}^{48}$   
→ Reduced **memory footprint**

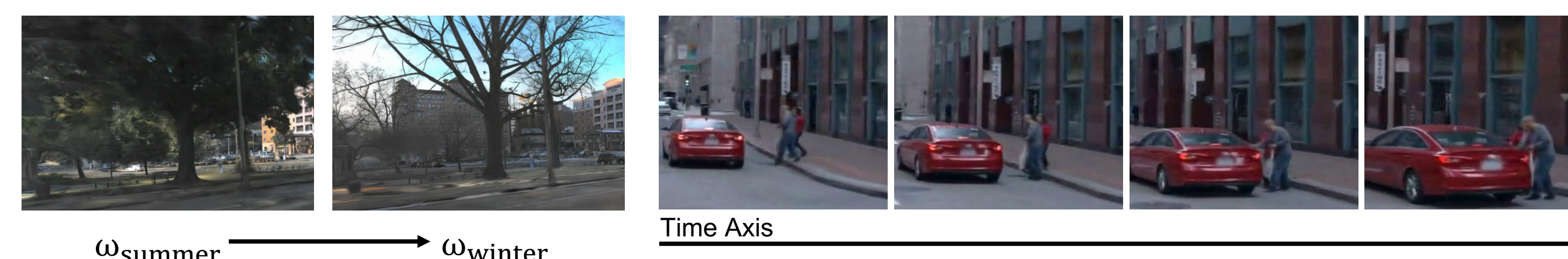
- $N \times 59 \times 4$  bytes = **2.25 GiB**
- $N \times 11 \times 4$  bytes = **0.42 GiB**
- ... for 10M Gaussians



## 4. And?

We can now **model multiple complex traffic scenarios**

- Change scene appearance *and* geometry by exchanging latent code  $\omega$
- Model articulated motion: walking, holding a shopping bag, opening car door, ...



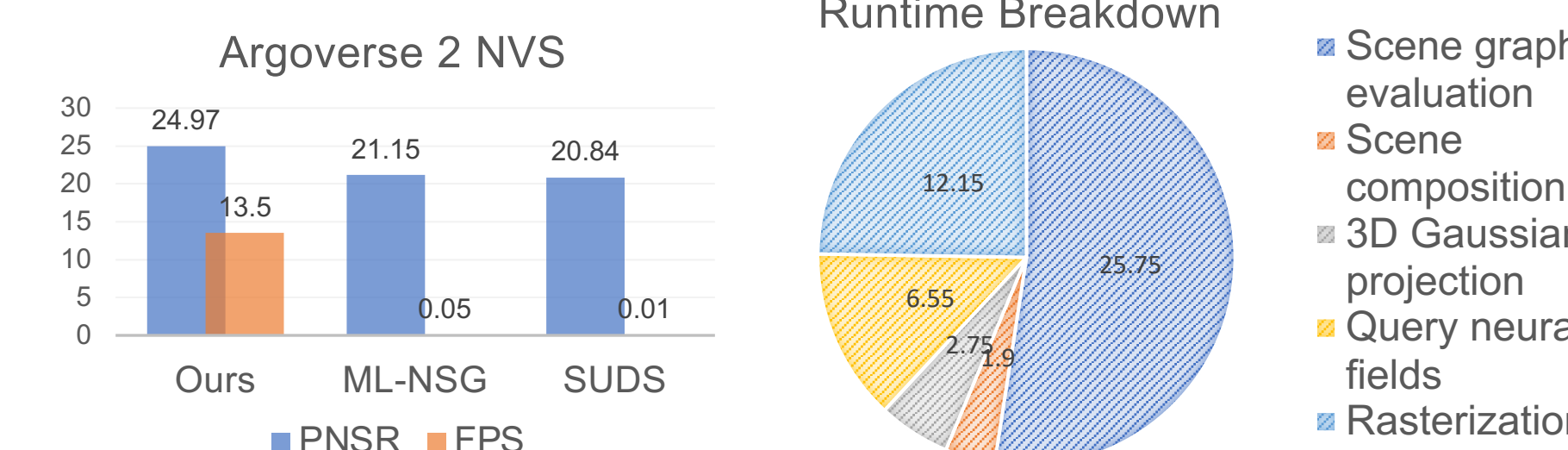
## Key insights

- Neural fields**
- Achieve **comparable accuracy** to spherical harmonics
  - Scale better
  - Are *much* more **flexible**
  - Do **not hurt speed** critically

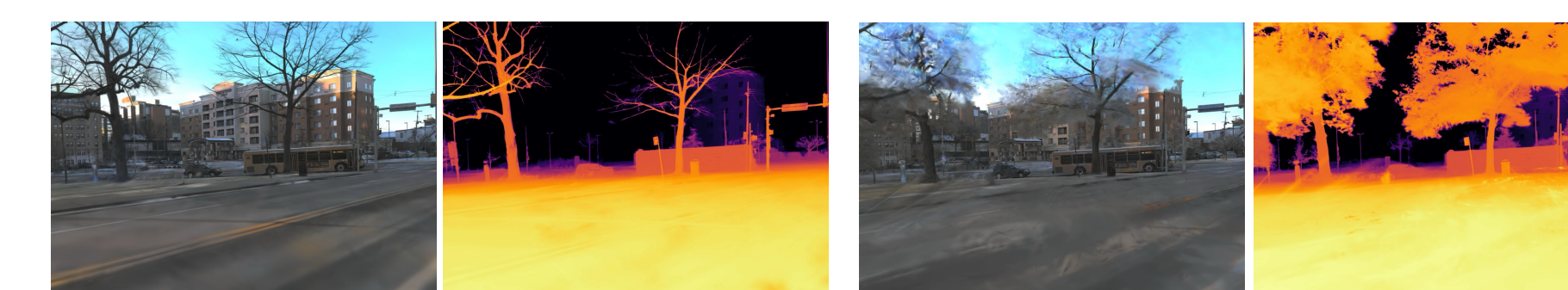
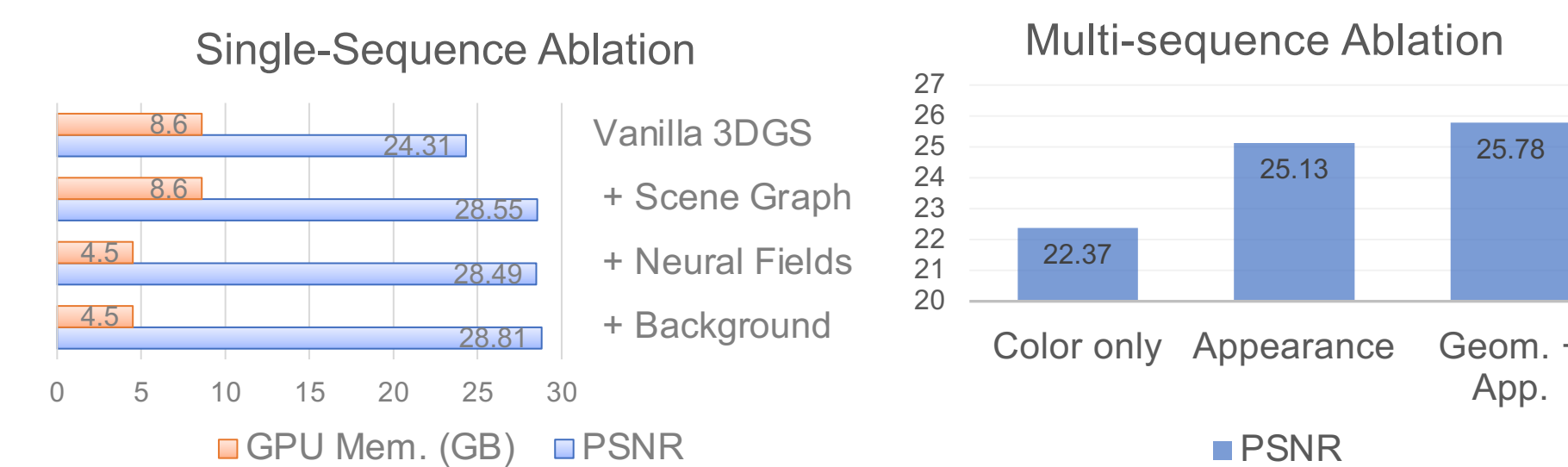
**Opacity correction term  $\mathbf{v}$**

- Essential for *realistic* rendering of heterogeneous captures (Geom. + App.)

- Helpful tricks**
- Multi-GPU training & ADC
  - Camera optimization
  - Background modeling
  - 4D Hash grids



We beat SOTA by over 3 dB while being over 200x faster!



Qualitative example: Geom. + App. vs Appearance