

Embodied Vision

4D Modeling

Tsung-Wei Ke

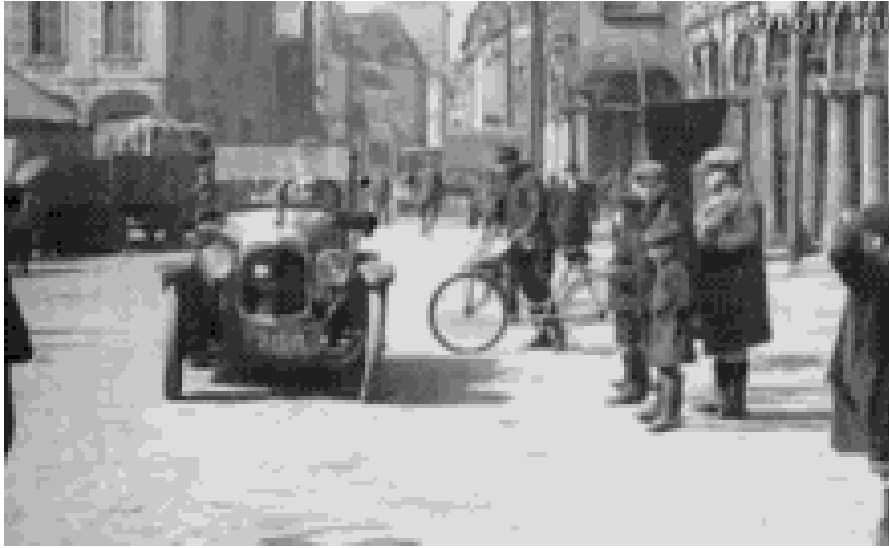
Spring 2026



Disclaimer

- This lecture borrows contents partly from
 - [Computer Vision](#) by Svetlana Lazebnik at UIUC
 - [Computer Vision](#) by Kris Kitani at CMU

We Live in a Dynamic World



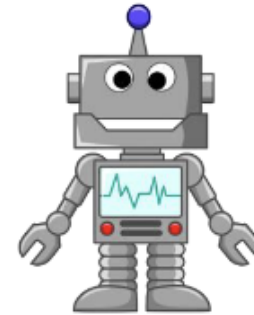
<https://i.gifer.com/fxxv.gif>

<https://ca.sports.yahoo.com/news/mvp-frontrunners-vladimir-guerrero-jr-shohei-ohtani-head-to-head-142258898.html>

We Live in a Dynamic World

Action:

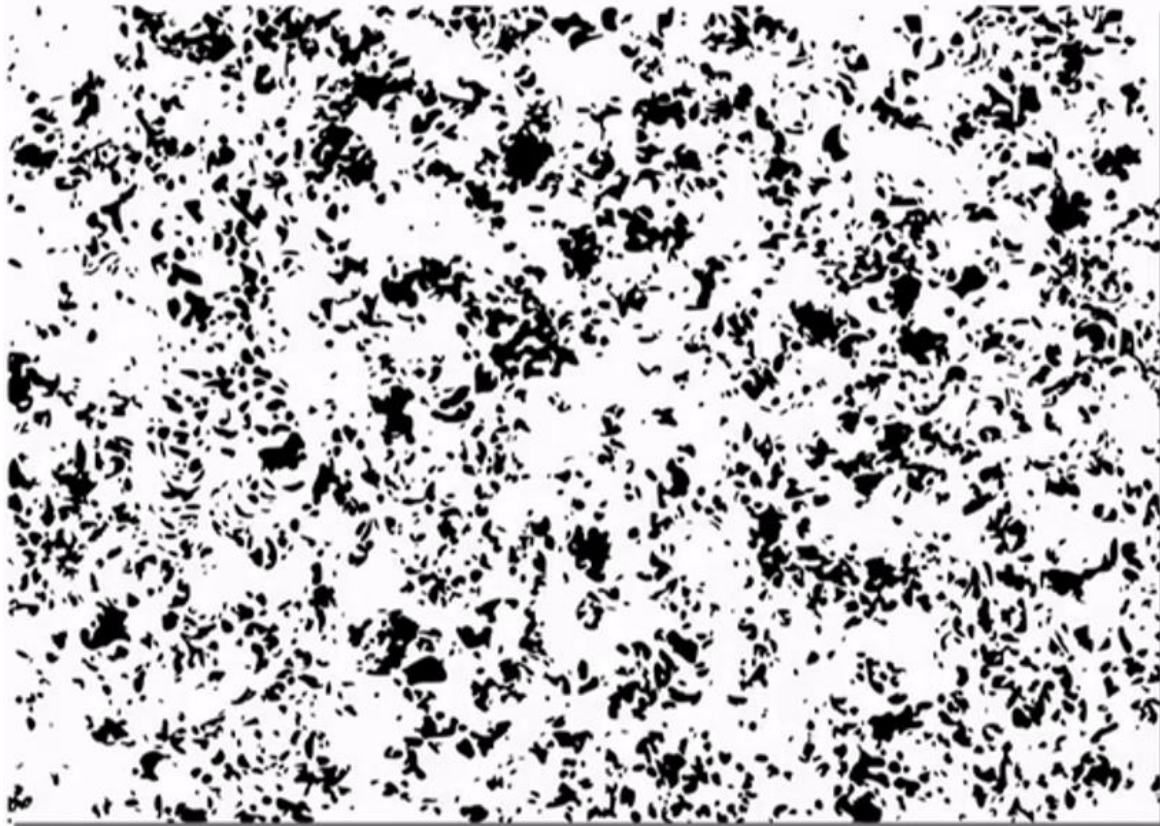
- What is the right action



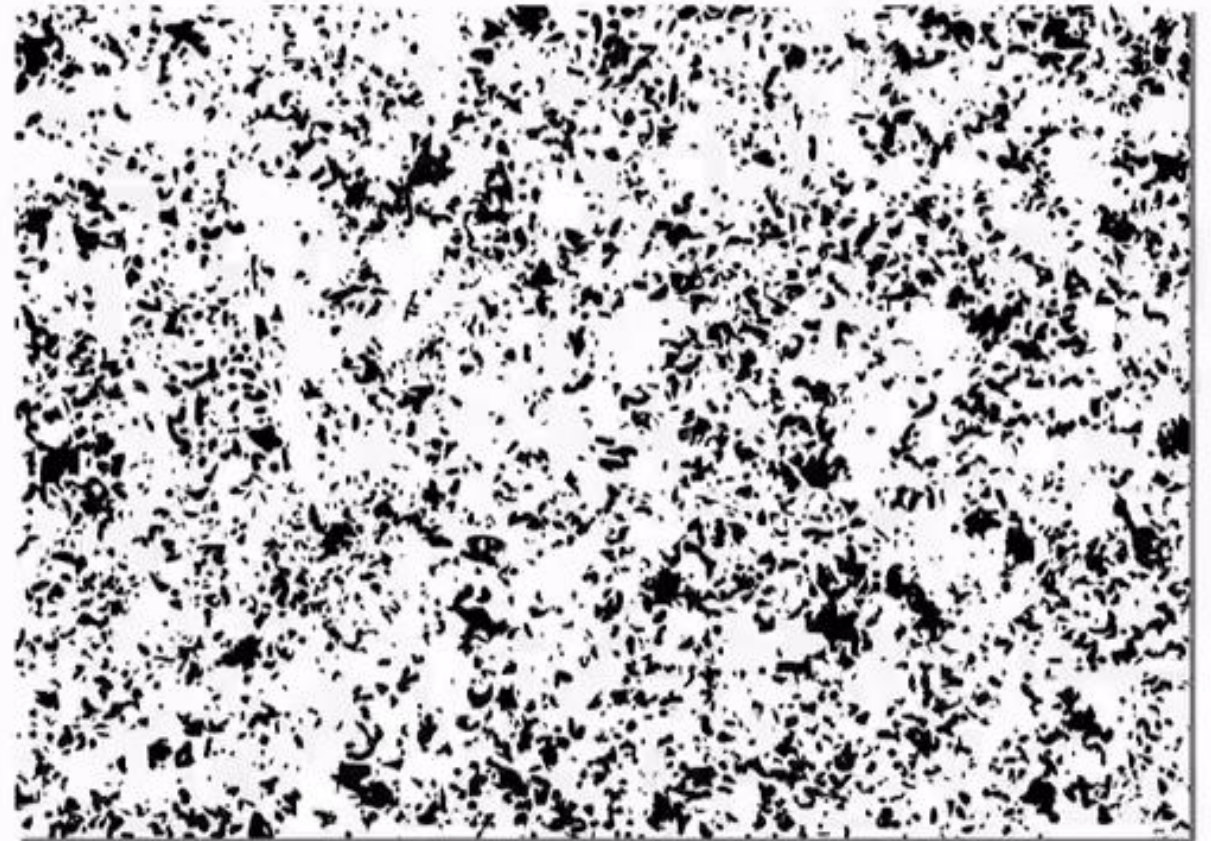
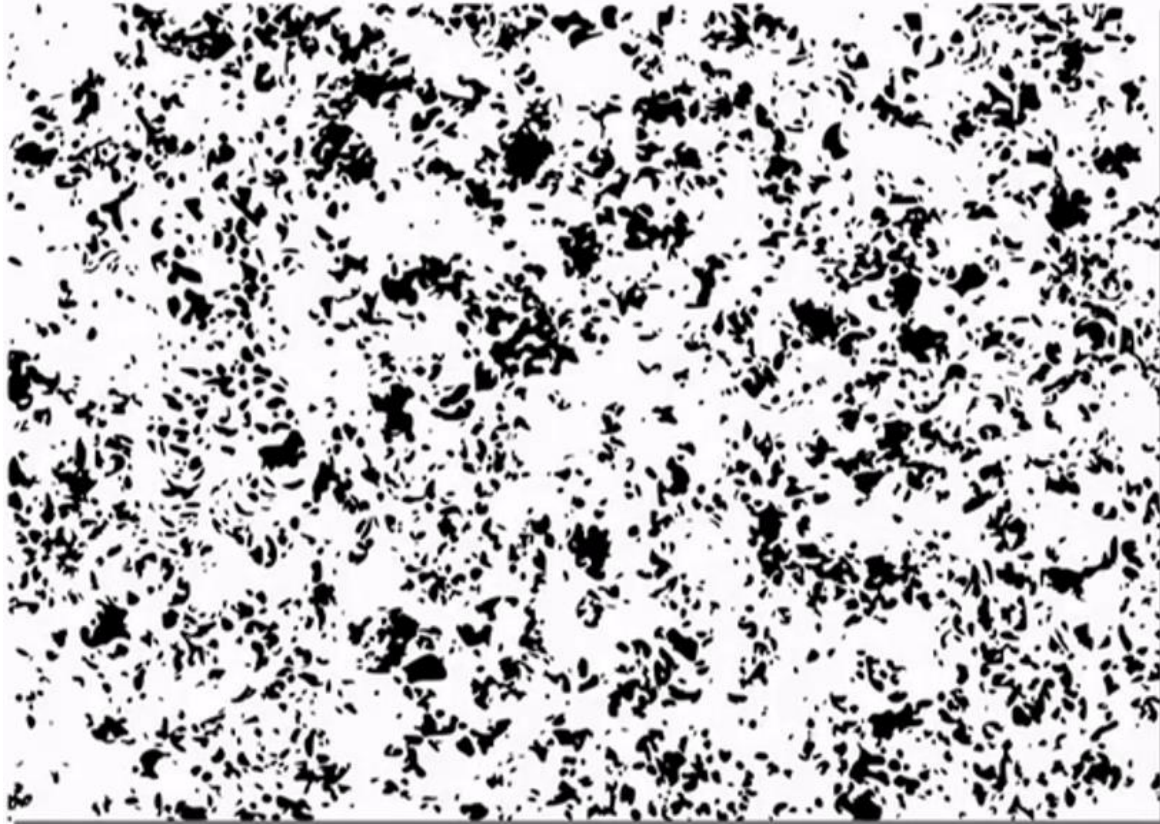
Perception:

- What and where are the objects
- How do objects move

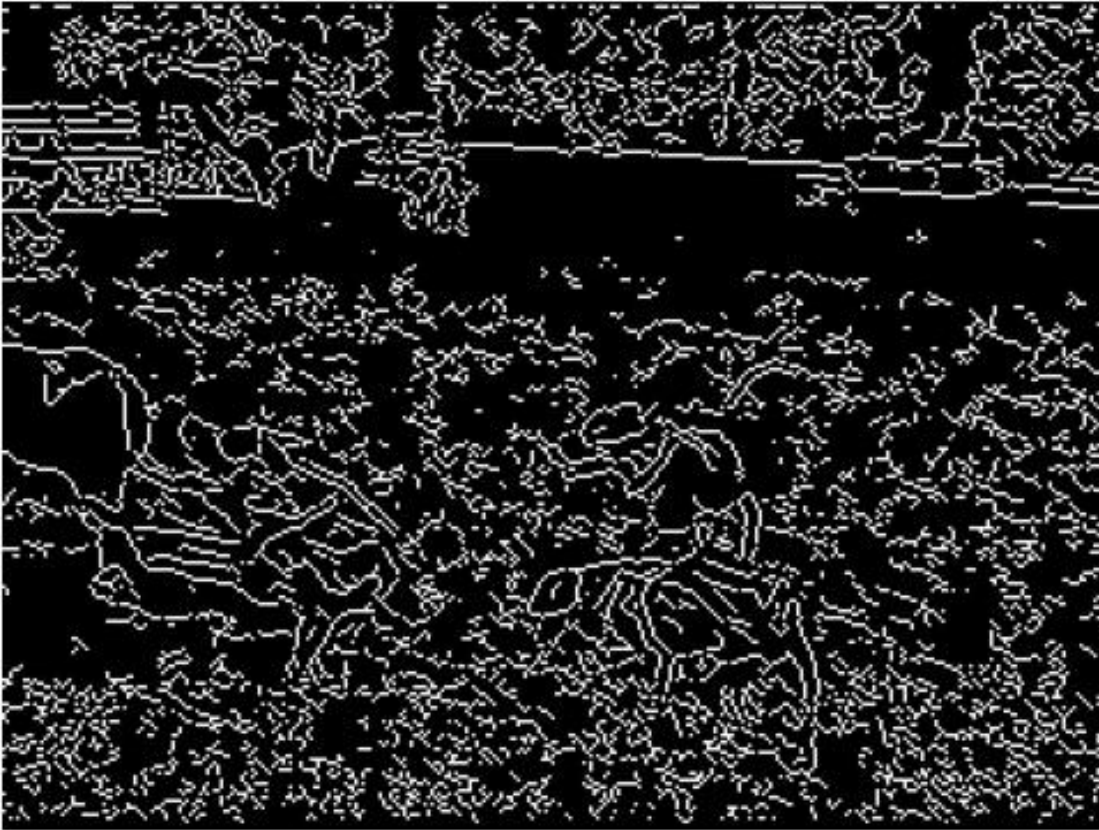
Motion Perception Does Not Rely on Object Recognition



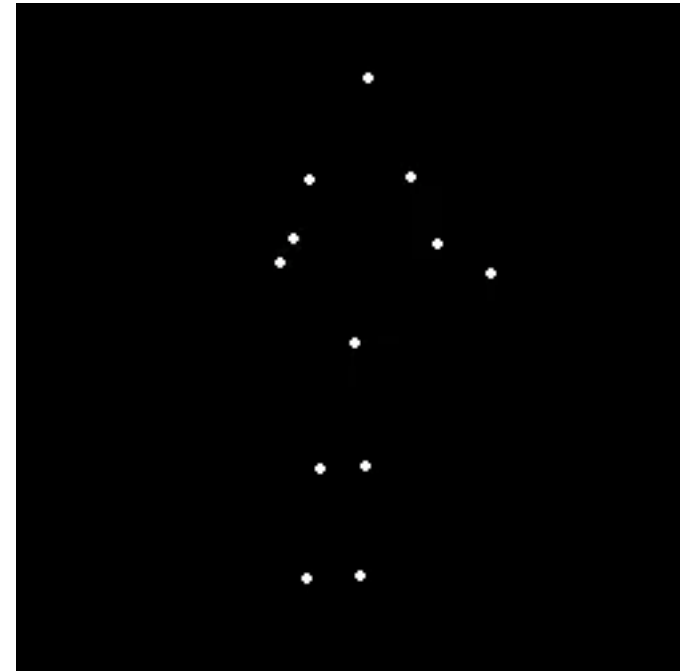
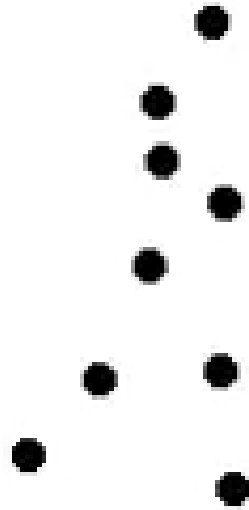
Motion Perception Does Not Rely on Object Recognition



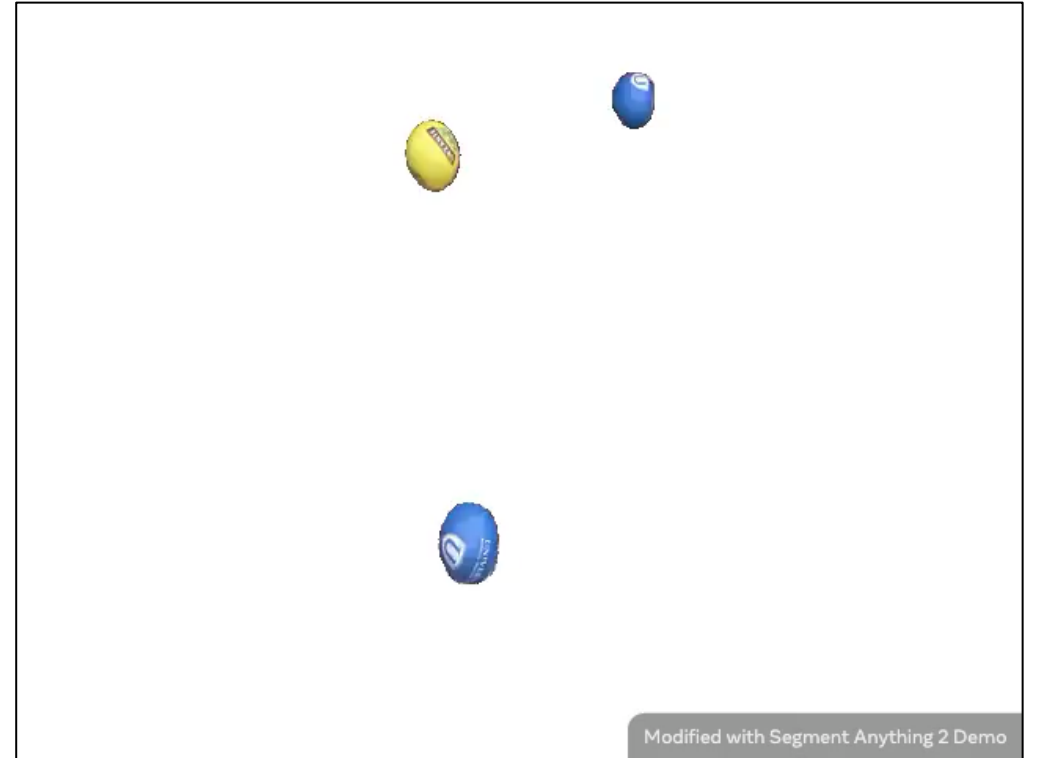
Motion Perception Does Not Rely on Object Recognition



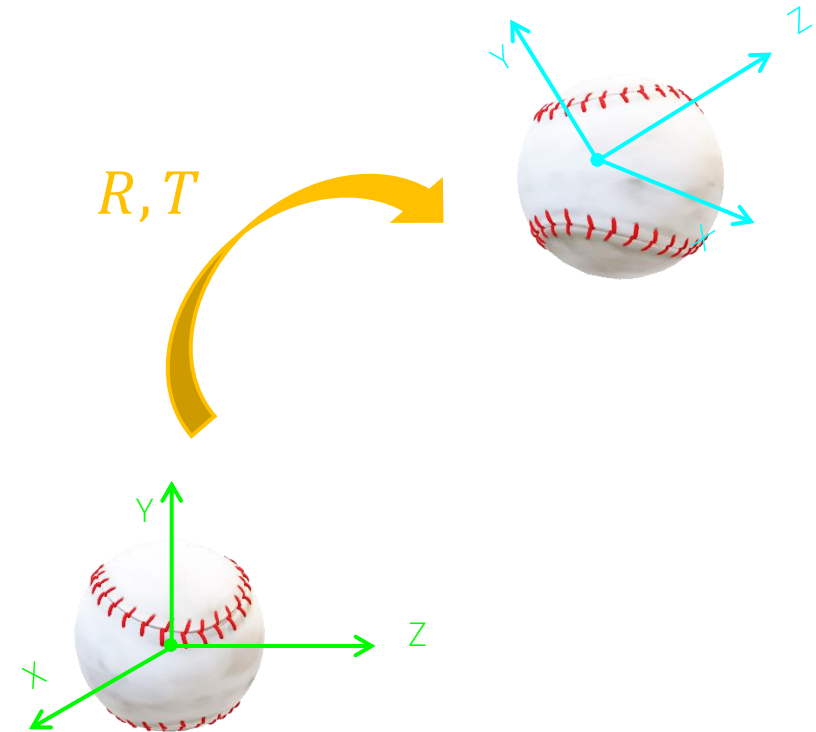
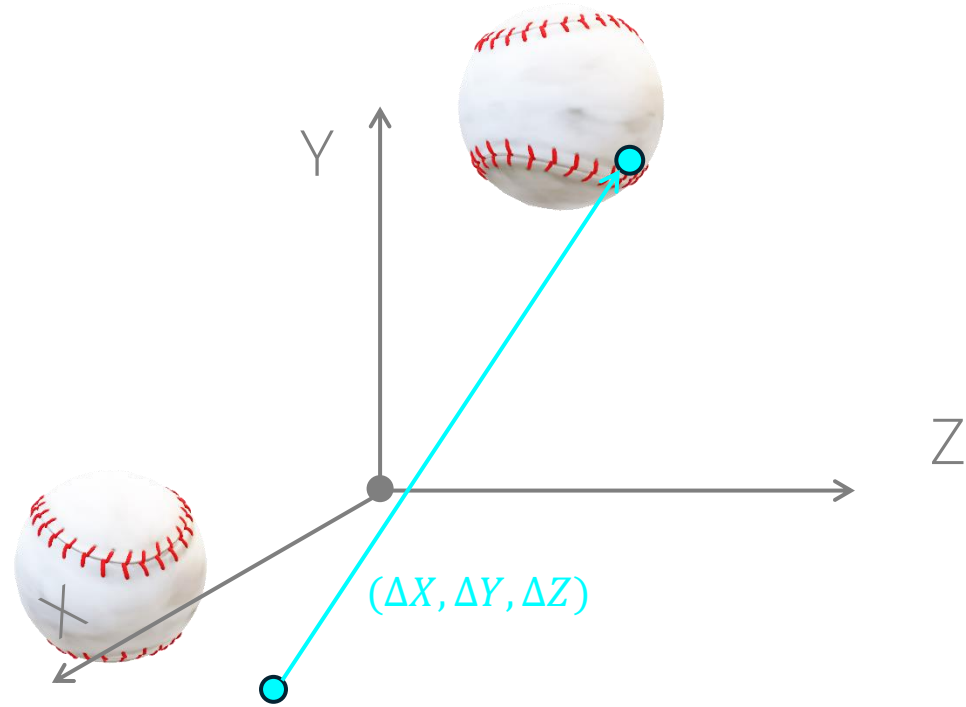
Motion Perception Does Not Rely on Object Recognition



How to Represent Motion?



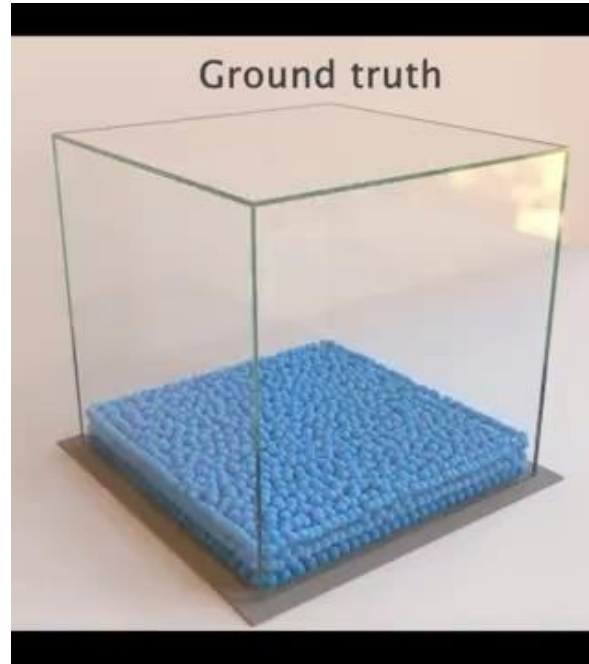
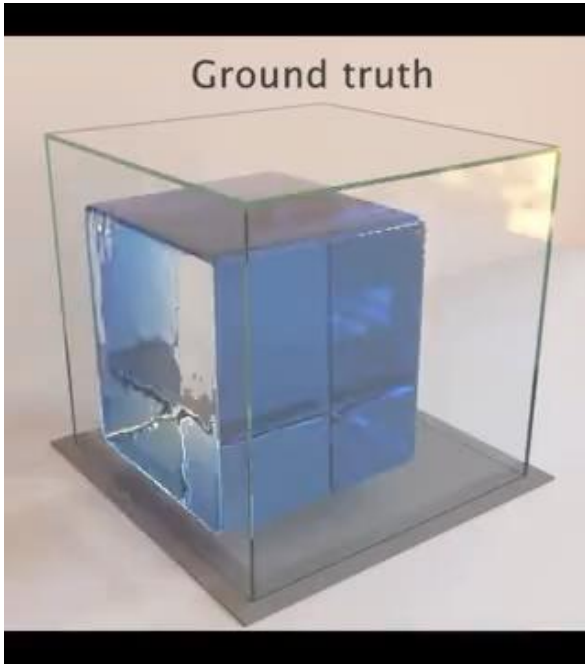
Rigid-Body Motion as a Motion Representation



How to Represent Motion?



Particle Motion as a Motion Representation



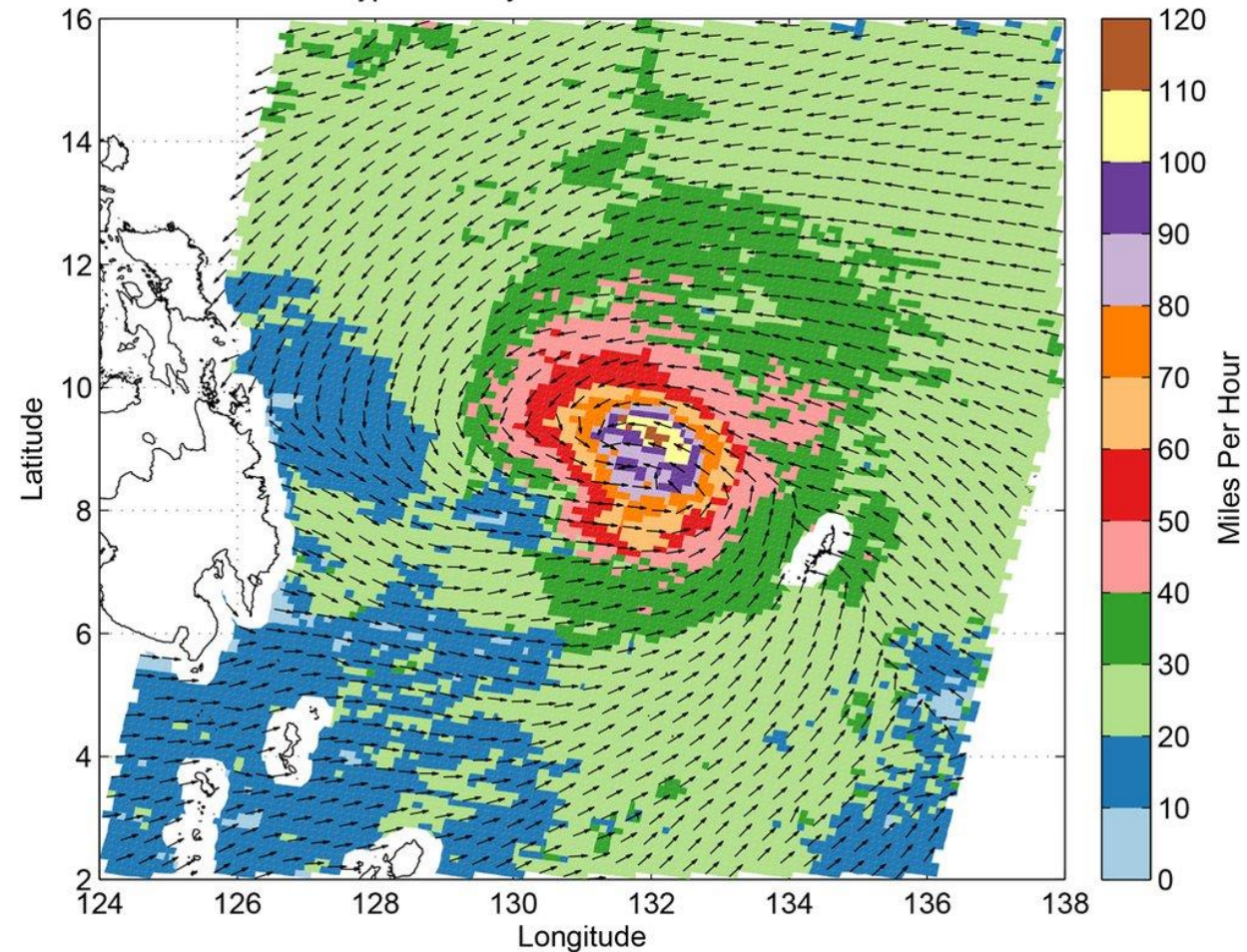
Learning to simulate complex physics with graph networks. Google Deepmind.



Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis. Luiten et al. 3DV 2024.

Vector Field as a Motion Representation

Typhoon Haiyan at Nov 7th 1:30 UTC



<https://youtu.be/mdN800kx2ko?si=DPRdwEUYnObugwKI>

<https://www.jpl.nasa.gov/news/nasa-peers-into-one-of-earths-strongest-storms-ever>

Particle Motion vs. Vector Field Describe the Same Motion from Different Coordinate Systems

Particle motion
(Lagrangian representations)



Vector field
(Eulerian representations)



- We will revisit for physical modeling non-rigid bodies

Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

Content

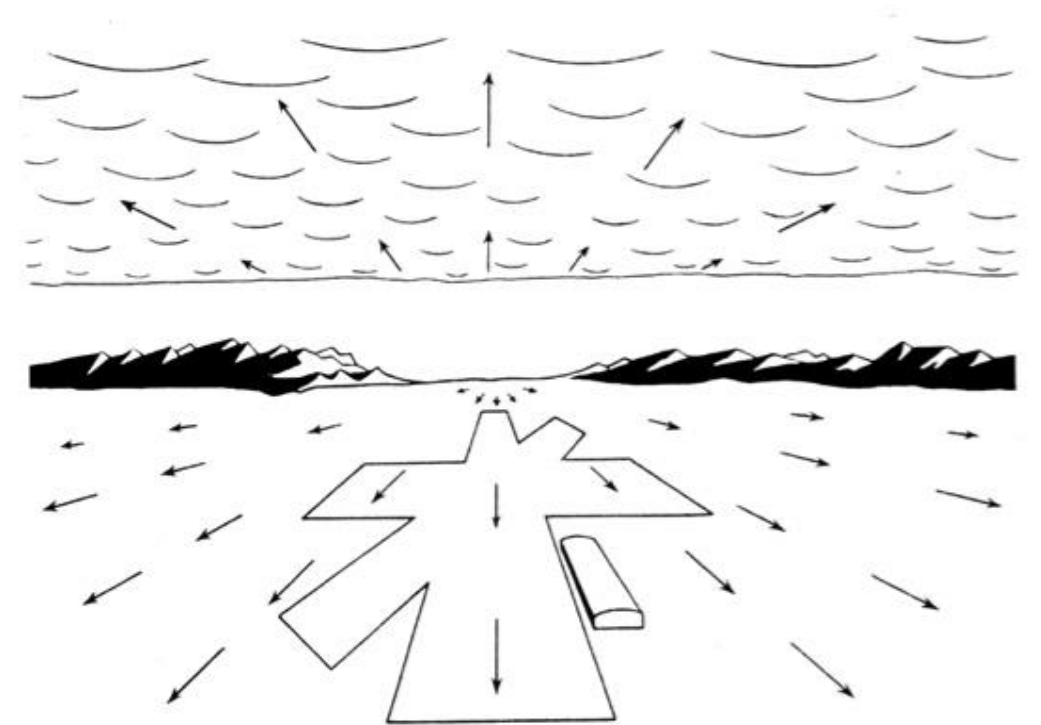
- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

Measure Motion from 2D Images

For 2D images, we don't know the exact 3D motions. We can only measure the pixel apparent motion

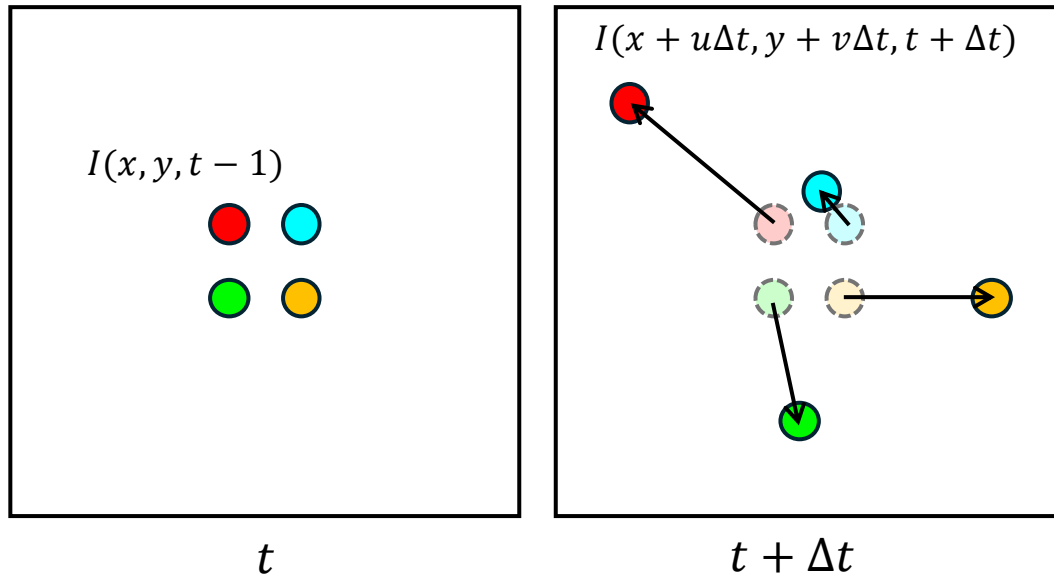


<https://in.pinterest.com/pin/594686325771980172/>



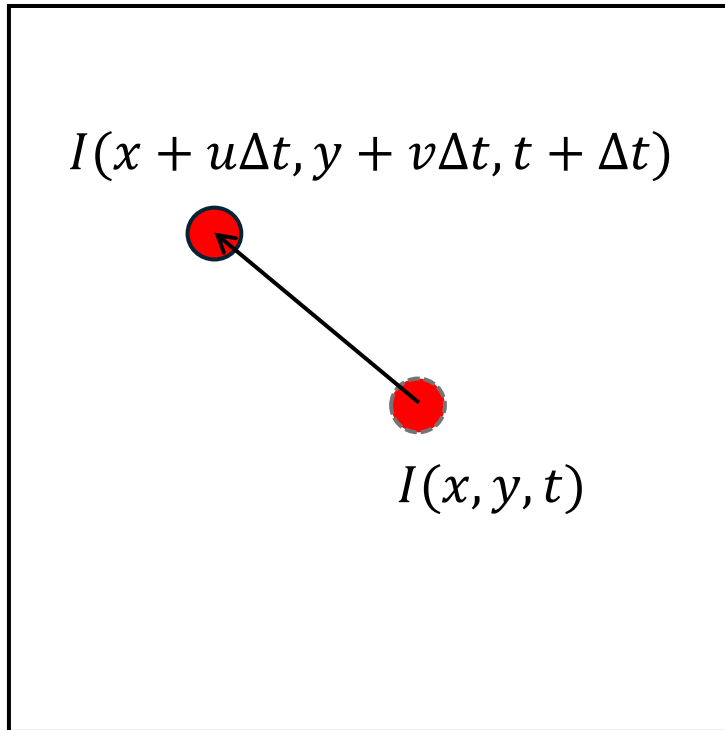
<https://www.cns.nyu.edu/~david/courses/perception/lecturenotes/motion/motion.html>

Optical Flow: Pixel Apparent Motion



- Given frames at times t and $t + \Delta t$, estimate the apparent motion field that describes the velocity $\mathbf{u}(x, y, t)$ and $\mathbf{v}(x, y, t)$ of a pixel at location x, y and time $t + \Delta t$.
- Problem: how can we estimate the velocity field $\mathbf{u}(x, y, t)$ and $\mathbf{v}(x, y, t)$ by observing a sequence of images?

Brightness Constancy Constraints



- **Brightness Constancy Constraints:** the brightness of a particular point remains constant across time

$$I(x + u\Delta t, y + v\Delta t, t + \Delta t) \approx I(x, y, t)$$

- **Brightness Constancy Equation:**

$$\boxed{\frac{DI}{Dt}} = 0 \Rightarrow \boxed{\frac{\partial I}{\partial t}} + \boxed{\frac{\partial I}{\partial x}} \frac{\partial x}{\partial t} + \boxed{\frac{\partial I}{\partial y}} \frac{\partial y}{\partial t} = 0$$

Total derivative

Partial derivative

Derivation of Brightness Constancy Equation

- Multi-variable Taylor Expansion:

$$f(x, y) \approx f(a, b) + \frac{\partial f(a, b)}{\partial x} (x - a) + \frac{\partial f(a, b)}{\partial y} (y - b)$$

- Brightness Constancy Equation:

$$I(x + u\Delta t, y + v\Delta t, t + \Delta t) \approx I(x, y, t) + \frac{\partial I(x, y, t)}{\partial x} u\Delta t + \frac{\partial I(x, y, t)}{\partial y} v\Delta t + \frac{\partial I(x, y, t)}{\partial t} \Delta t$$

$$\Rightarrow \frac{\partial I(x, y, t)}{\partial x} u + \frac{\partial I(x, y, t)}{\partial y} v + \frac{\partial I(x, y, t)}{\partial t} = 0 \Rightarrow \frac{\partial I}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial I}{\partial y} \frac{\partial y}{\partial t} + \frac{\partial I}{\partial t} = 0$$

Take a Closer Look at the Equation

$$\frac{\partial I}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial I}{\partial y} \frac{\partial y}{\partial t} + \frac{\partial I}{\partial t} = 0$$

- Image gradients $\frac{\partial I}{\partial x}$ and $\frac{\partial I}{\partial y}$: spatial derivatives of image intensity, which can be estimated by forward difference / Sober filter / Scharr filter ...
- Temporal gradients $\frac{\partial I}{\partial t}$: temporal derivatives of image intensity, which can be estimated by subtracting images at different timestep
- Optical flow $\frac{\partial x}{\partial t}$ and $\frac{\partial y}{\partial t}$: the unknown variables we want to estimate

Image Gradients and Temporal Gradients

Image gradients $\frac{\partial I}{\partial x}$

t

1	1	1	1	1
1	1	1	1	1
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10

t

1	1	1	1	1
1	1	1	1	1
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10

−

t

1	1	1	1	1
1	1	1	1	1
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10

=

-	0	0	0	-
-	0	0	0	-
-	9	0	0	-
-	9	0	0	-
-	0	0	0	-

$t + 1$

1	1	1	1	1
1	1	1	1	1
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10

Temporal gradients $\frac{\partial I}{\partial t}$

t

1	1	1	1	1
1	1	1	1	1
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10

−

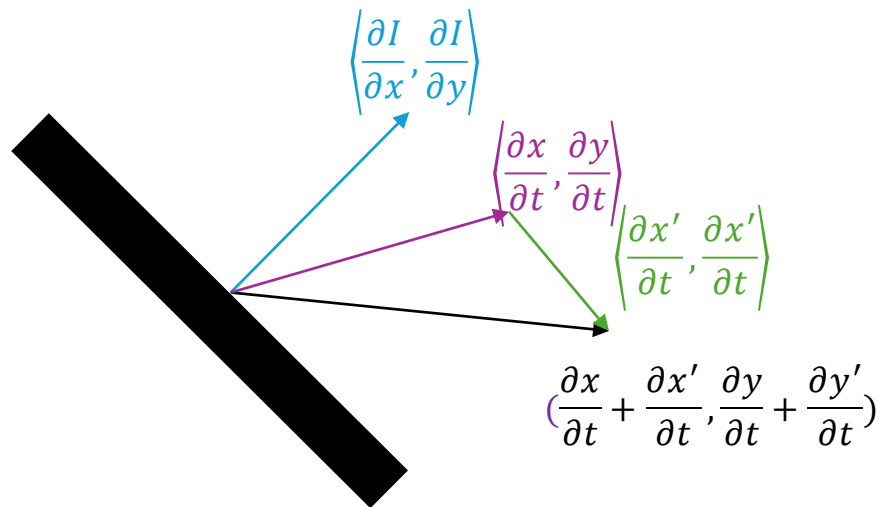
$t + 1$

1	1	1	1	1
1	1	1	1	1
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10

=

0	0	0	0	0
0	0	0	0	0
0	9	9	9	9
0	0	0	0	0
0	0	0	0	0

Given Image and Temporal Gradient, Can We Recover Optical Flow Uniquely?

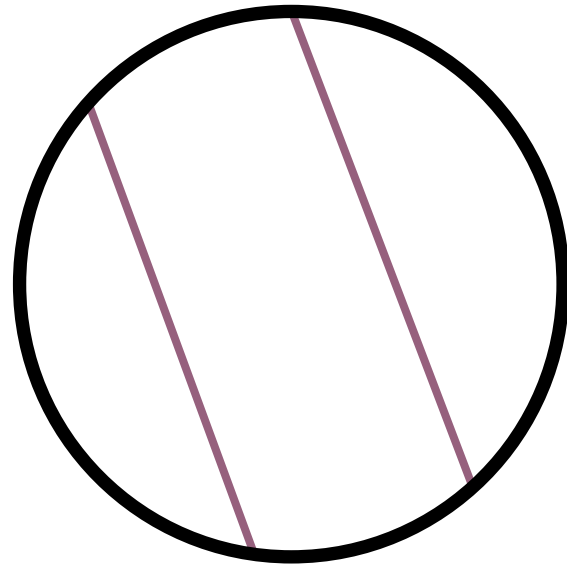


- Let's first write the brightness constancy constraint in vector form:

$$\left\langle \frac{\partial I}{\partial x'}, \frac{\partial I}{\partial y} \right\rangle \cdot \left\langle \frac{\partial x}{\partial t'}, \frac{\partial y}{\partial t} \right\rangle + \frac{\partial I}{\partial t} = 0$$

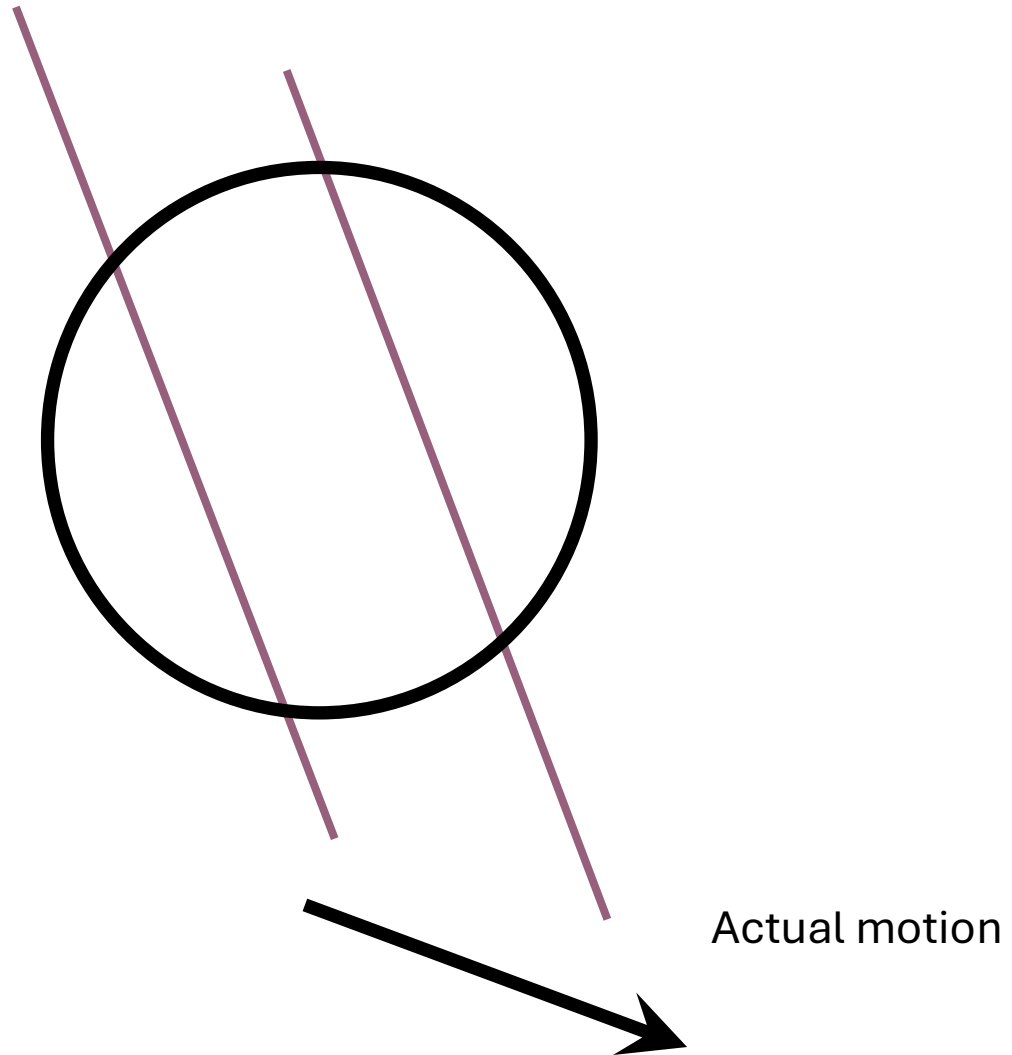
- We have another solution with $\frac{\partial x}{\partial t} + \frac{\partial x'}{\partial t}$ and $\frac{\partial y}{\partial t} + \frac{\partial y'}{\partial t}$ when $\left\langle \frac{\partial I}{\partial x'}, \frac{\partial I}{\partial y} \right\rangle \cdot \left\langle \frac{\partial x'}{\partial t}, \frac{\partial x'}{\partial t} \right\rangle = 0$
- In other words, the component perpendicular to the gradient (i.e., parallel to the edge) cannot be recovered!

The Aperture Problem

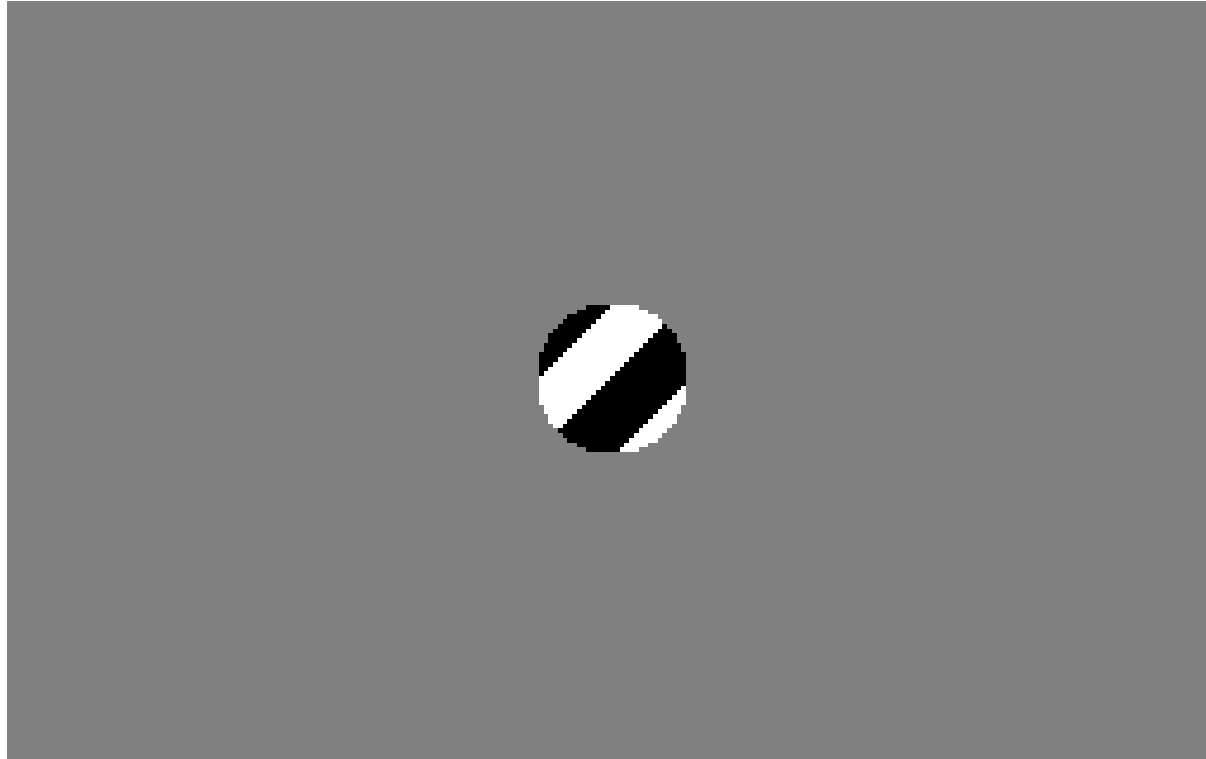


Perceived motion

The Aperture Problem

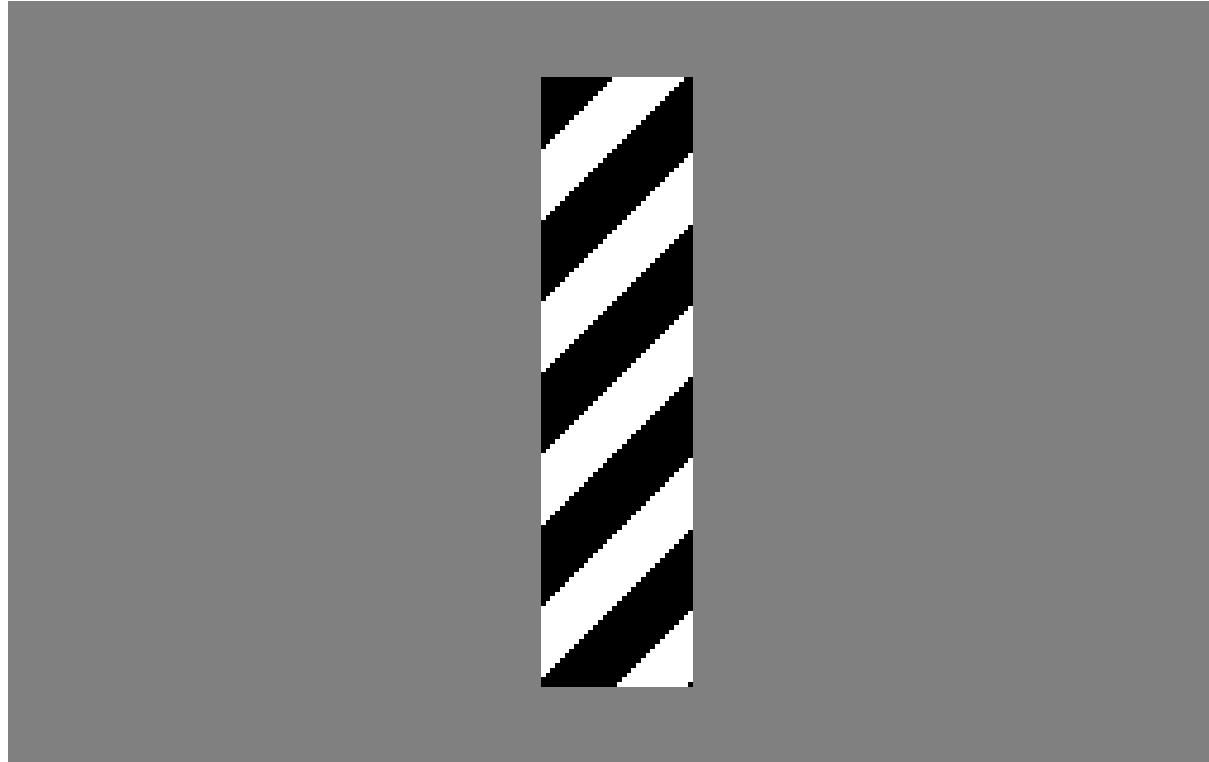


The Barber Pole Illusion



http://en.wikipedia.org/wiki/Barberpole_illusion

The Barber Pole Illusion



http://en.wikipedia.org/wiki/Barberpole_illusion

How to Address the Ambiguity?

- Brightness Constancy Equation: $\frac{\partial I}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial I}{\partial y} \frac{\partial y}{\partial t} + \frac{\partial I}{\partial t} = 0$ has two unknowns $(\frac{\partial x}{\partial t}, \frac{\partial y}{\partial t})$ and 1 equation for 1 pixel
- From the previous examples, we know that looking at more pixels of the same motion helps to resolve the ambiguity
- Idea: solve optical flow with a group of pixels moving in the same velocity!

Lukas-Kanade Method: Address the Ambiguity with Spatial Consistency Constraint

- **Spatial Consistency Constraint:** assume neighboring pixels have the same motion
- If we have n pixels in the neighborhood, then we can set up a linear least squares system:

$$\begin{bmatrix} \frac{\partial I(x_1, y_1)}{\partial x} & \frac{\partial I(x_1, y_1)}{\partial y} \\ \vdots & \vdots \\ \frac{\partial I(x_n, y_n)}{\partial x} & \frac{\partial I(x_n, y_n)}{\partial y} \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial t} \end{bmatrix} = - \begin{bmatrix} \frac{\partial I(x_1, y_1)}{\partial t} \\ \vdots \\ \frac{\partial I(x_n, y_n)}{\partial t} \end{bmatrix}$$

Estimate Optical Flow

$$\begin{bmatrix} \frac{\partial I(x_1, y_1)}{\partial x} & \frac{\partial I(x_1, y_1)}{\partial y} \\ \vdots & \vdots \\ \frac{\partial I(x_n, y_n)}{\partial x} & \frac{\partial I(x_n, y_n)}{\partial y} \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial t} \end{bmatrix} = - \begin{bmatrix} \frac{\partial I(x_1, y_1)}{\partial t} \\ \vdots \\ \frac{\partial I(x_n, y_n)}{\partial t} \end{bmatrix}$$

A
 d
 b

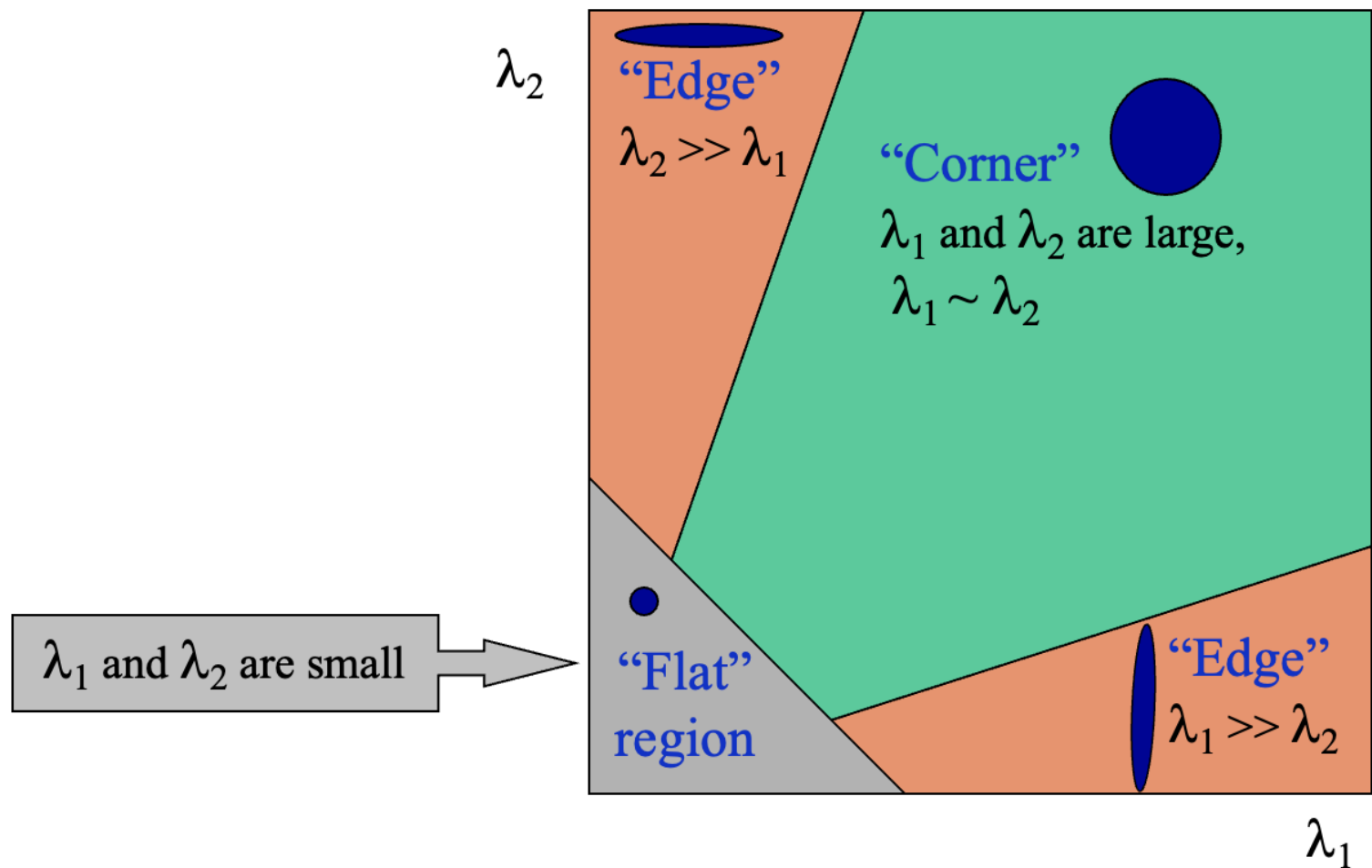
$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial t} \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A$
 d
 $A^T b$

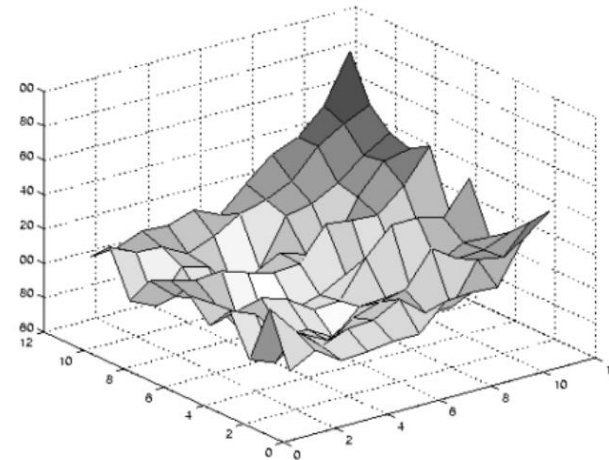
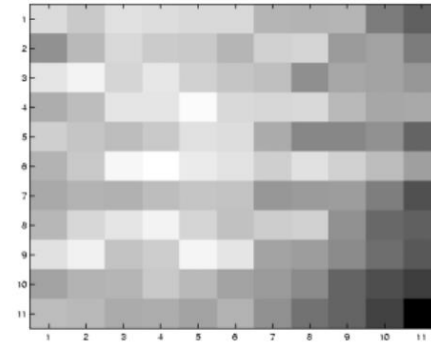
- This is a least square problem: the solutions of $Ad = b$ are the solutions of $A^T Ad = A^T b$
- Solution: $d = (A^T A)^{-1} A^T b$
- When is the system solvable?

Recall: second moment matrix

- Estimation of optical flow is well-conditioned precisely for regions with high “cornerness”:



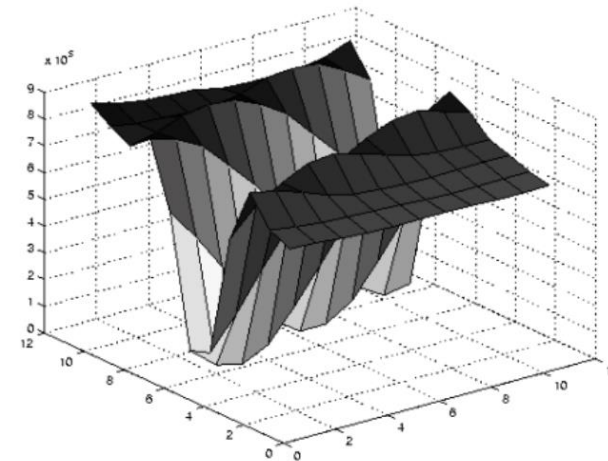
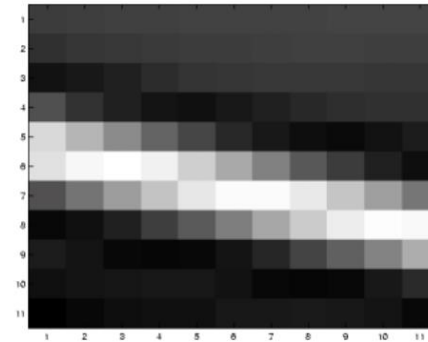
Low texture region



$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \nabla I (\nabla I)^T$$

- gradients have small magnitude
- small λ_1 , small λ_2

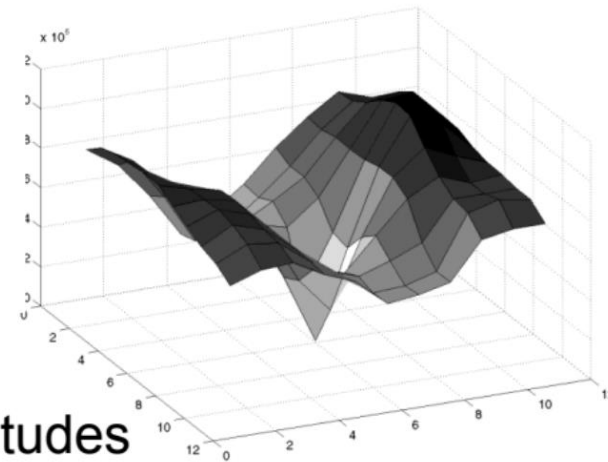
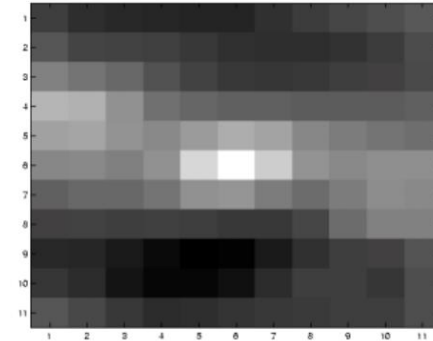
Edge



$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \nabla I (\nabla I)^T$$

- large gradients, all the same
- large λ_1 , small λ_2

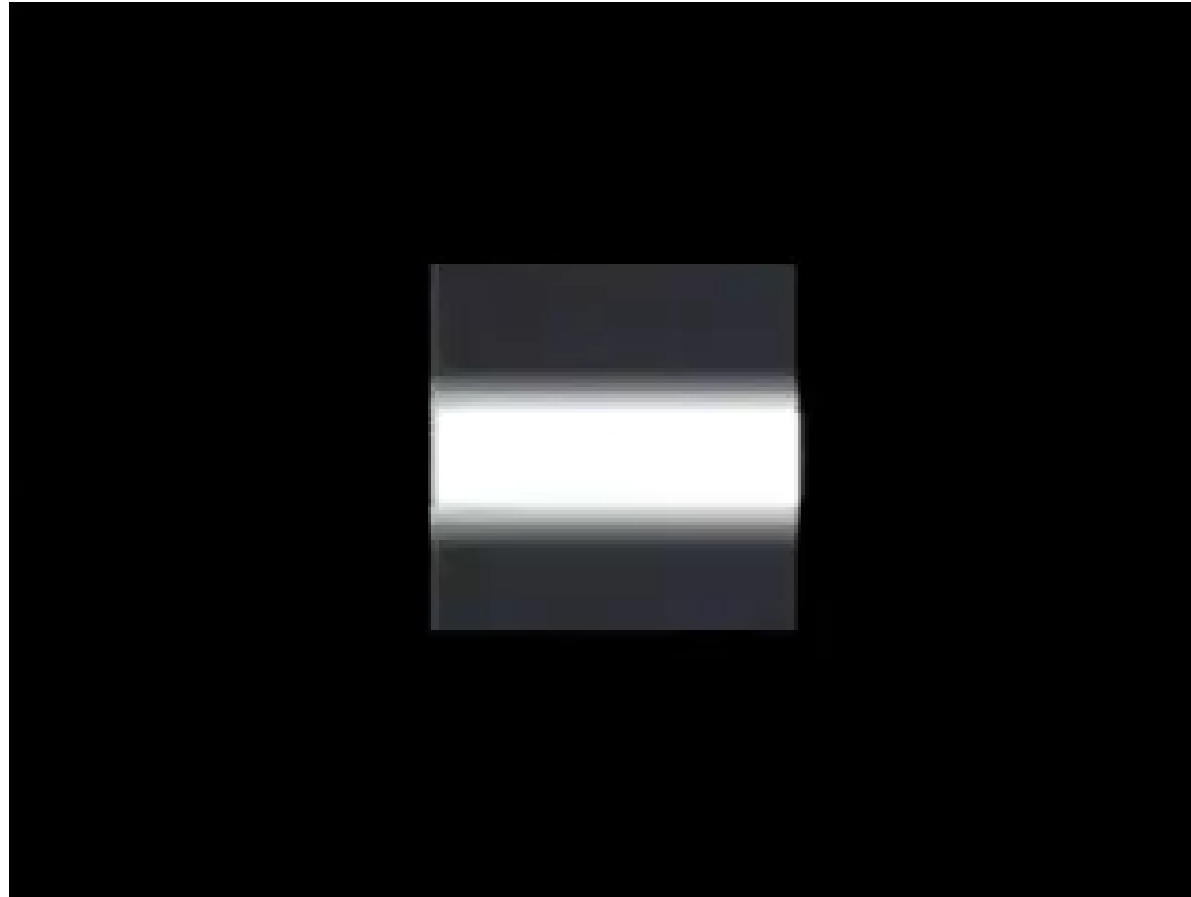
High texture region



$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \nabla I (\nabla I)^T$$

- gradients are different, large magnitudes
- large λ_1 , large λ_2

Conditions for solvability: A Bad Case of Single Straight Edge



Conditions for solvability: A Good Case

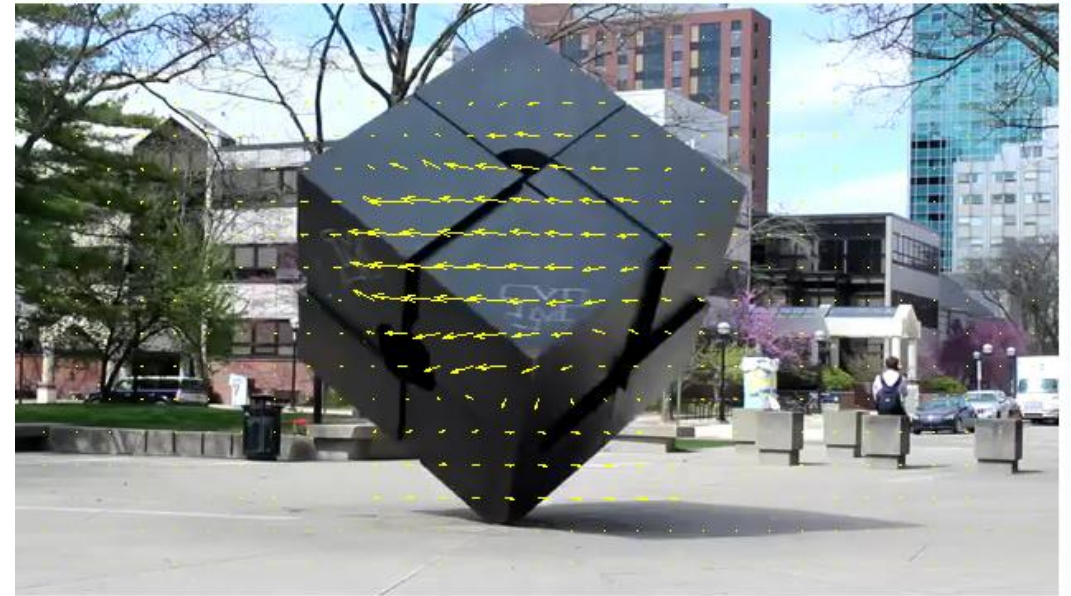


Lukas-Kanade Flow Example

Input frames



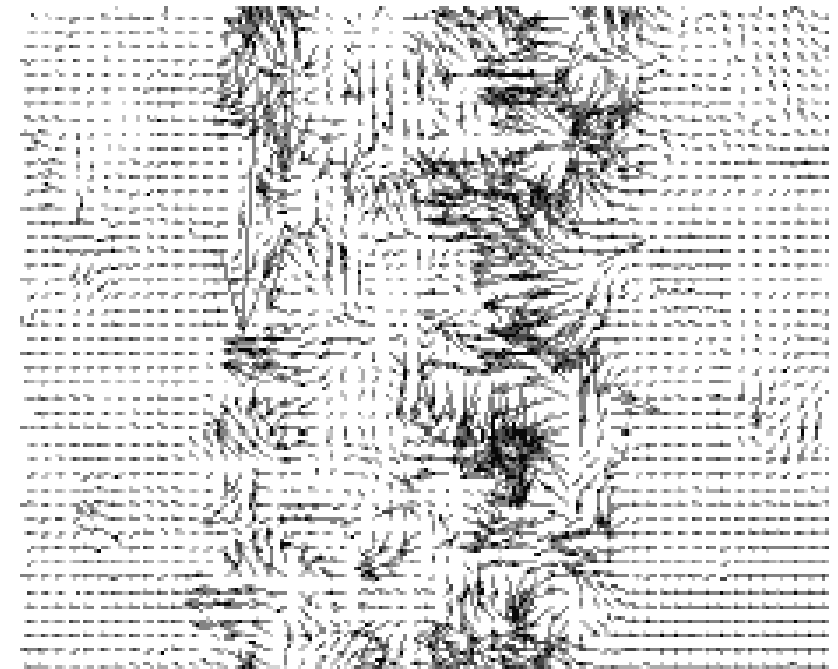
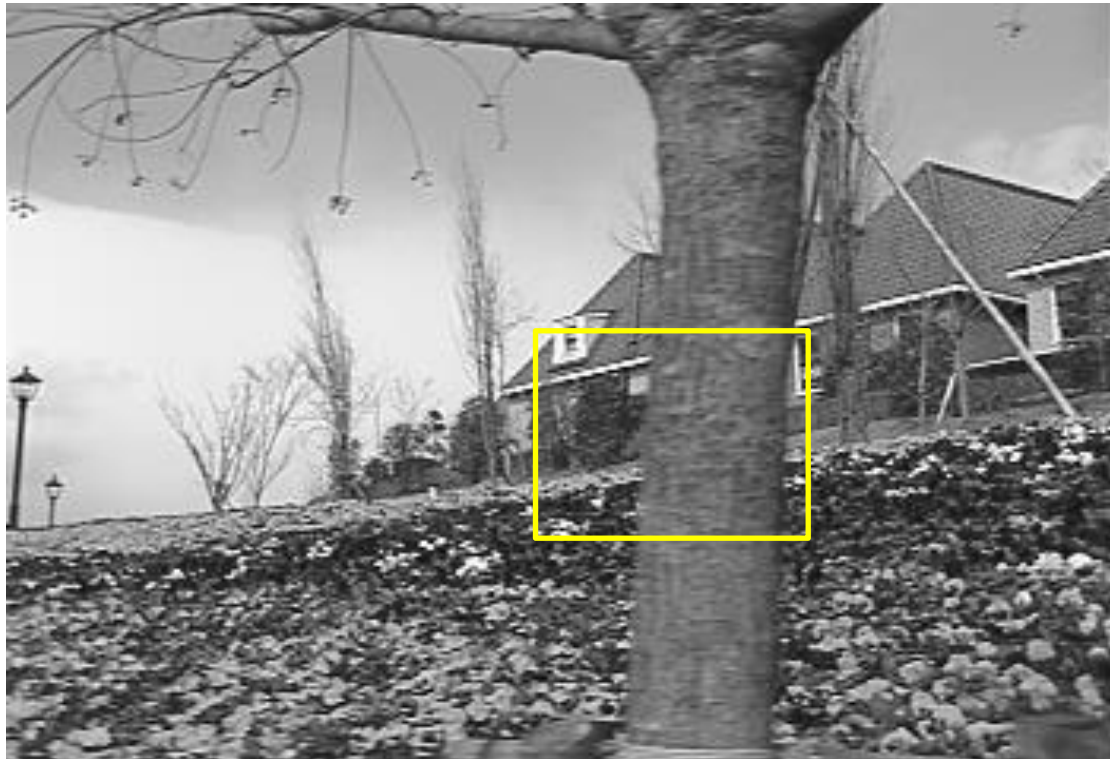
Output



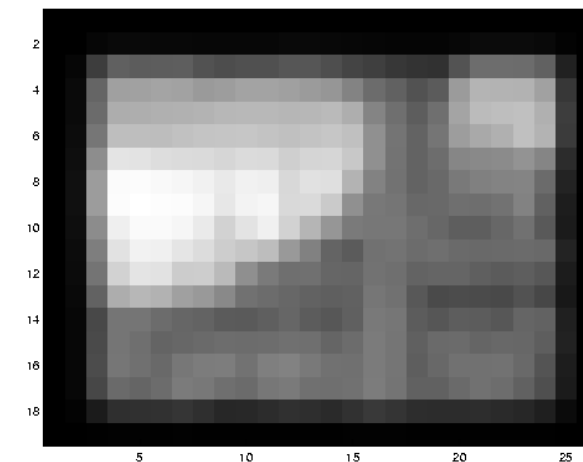
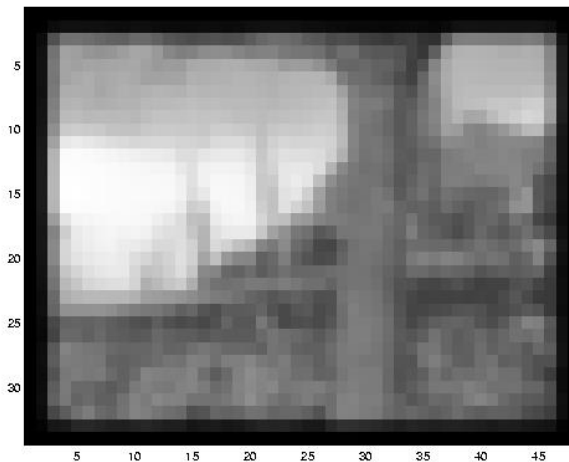
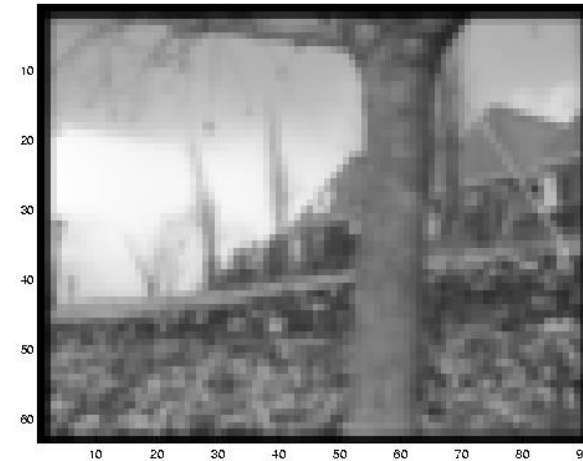
When Does Lukas-Kanade Method Fail?

- Up to now, we have assumed:
 1. 1st order Taylor expansion assumes small motions
 2. Spatial consistency constraint
 3. Brightness constancy constraint

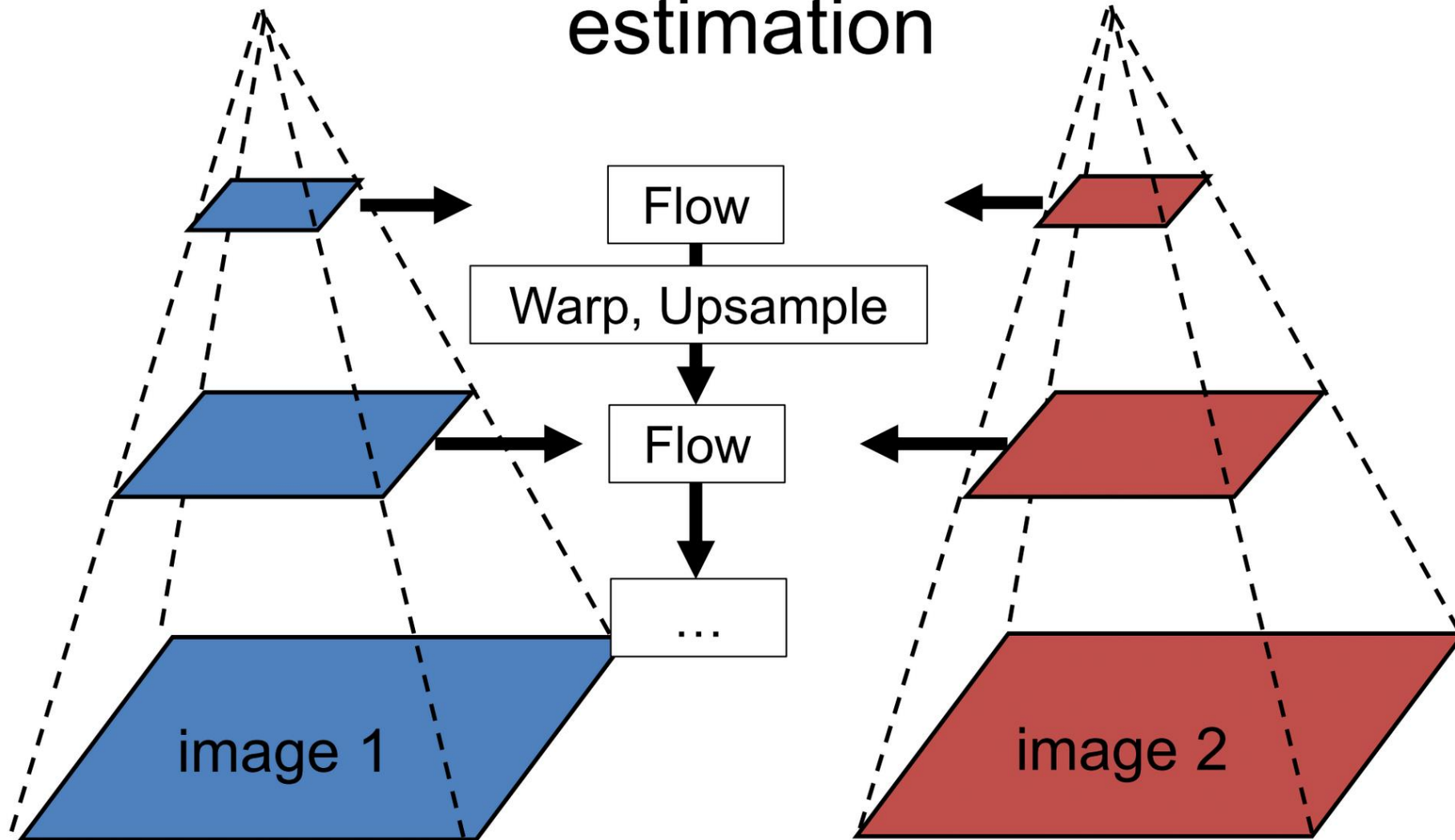
Lukas-Kanade Method Fails when Motion is Large



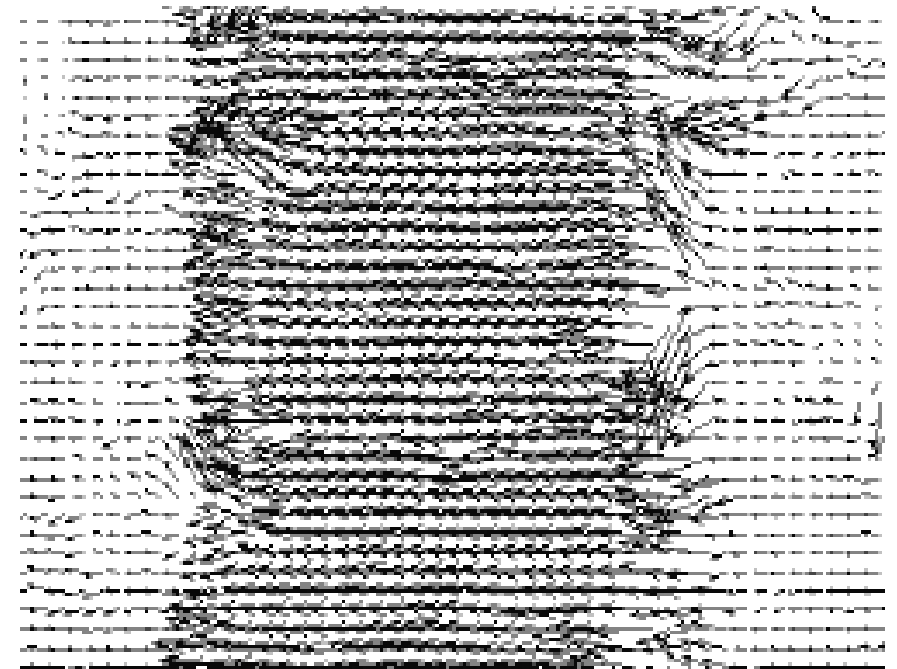
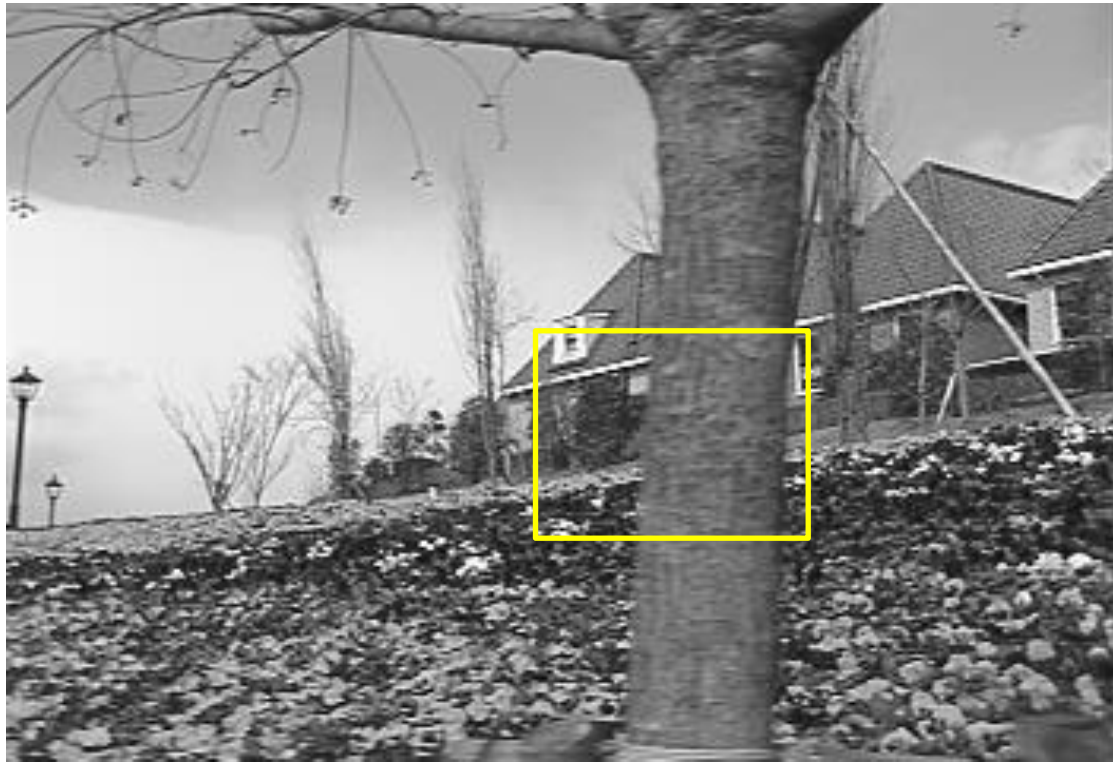
Solution: Reduce the Resolution!



Coarse-to-fine optical flow estimation



After Fixing



Lukas-Kanade Method Fails when Neighboring Pixels Don't Move in the Same Way

- Solution: Figure out which pixels move together then come back and fix

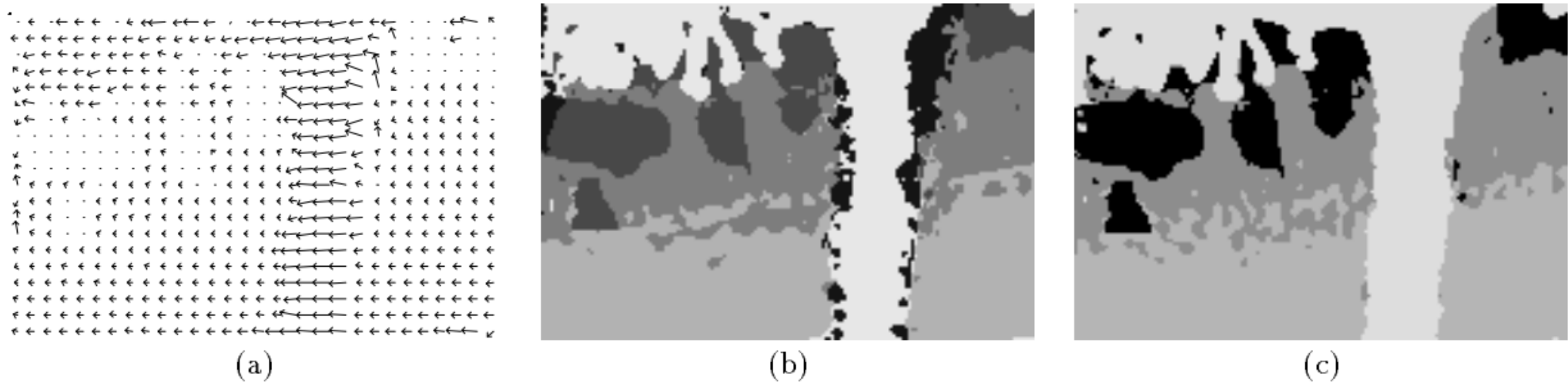
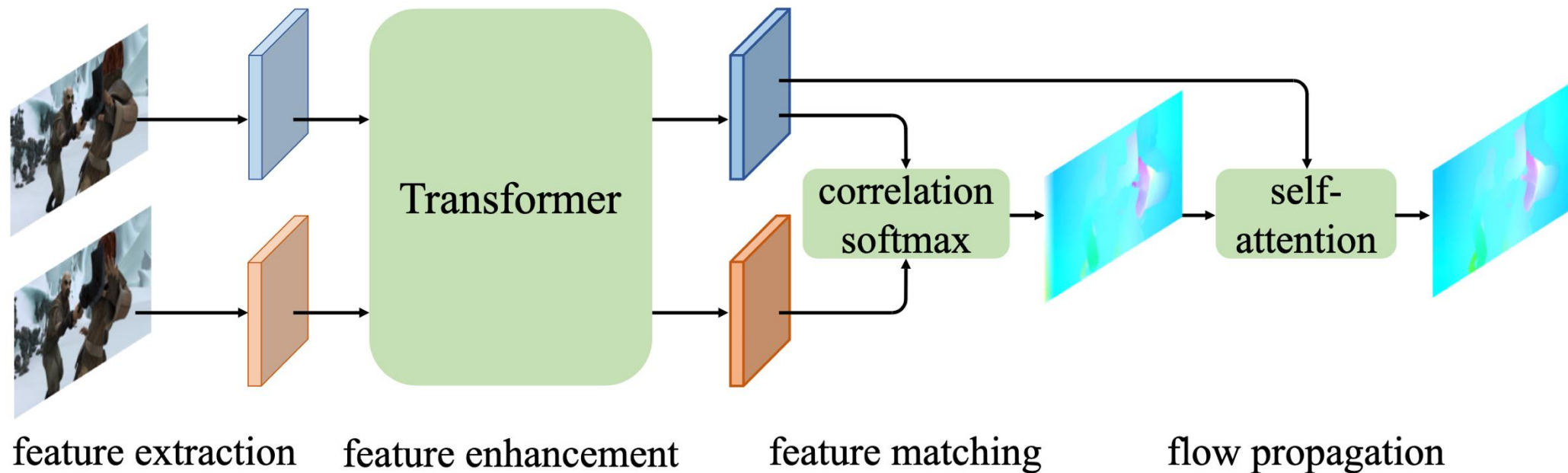


Figure 11: (a) The optic flow from multi-scale gradient method. (b) Segmentation obtained by clustering optic flow into affine motion regions. (c) Segmentation from consistency checking by image warping. Representing moving images with layers.

Lukas-Kanade Method Fails when Brightness Constancy Constraint Doesn't Hold

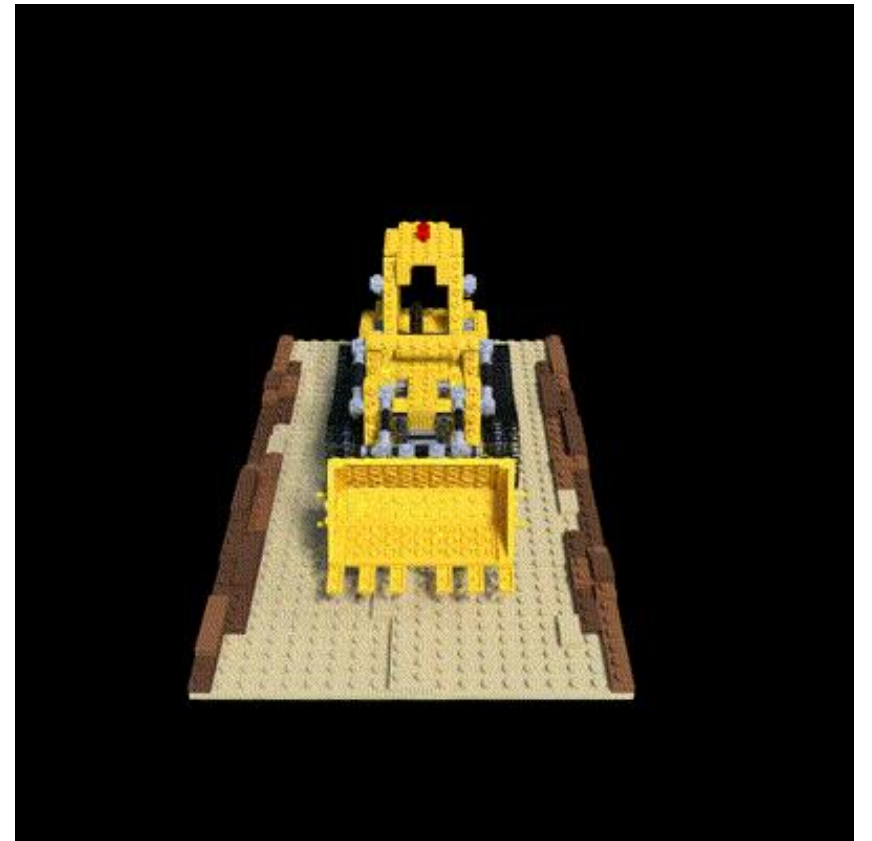
- Solution: Use other forms of matching (e.g. interest points) or learn a neural network for matching!



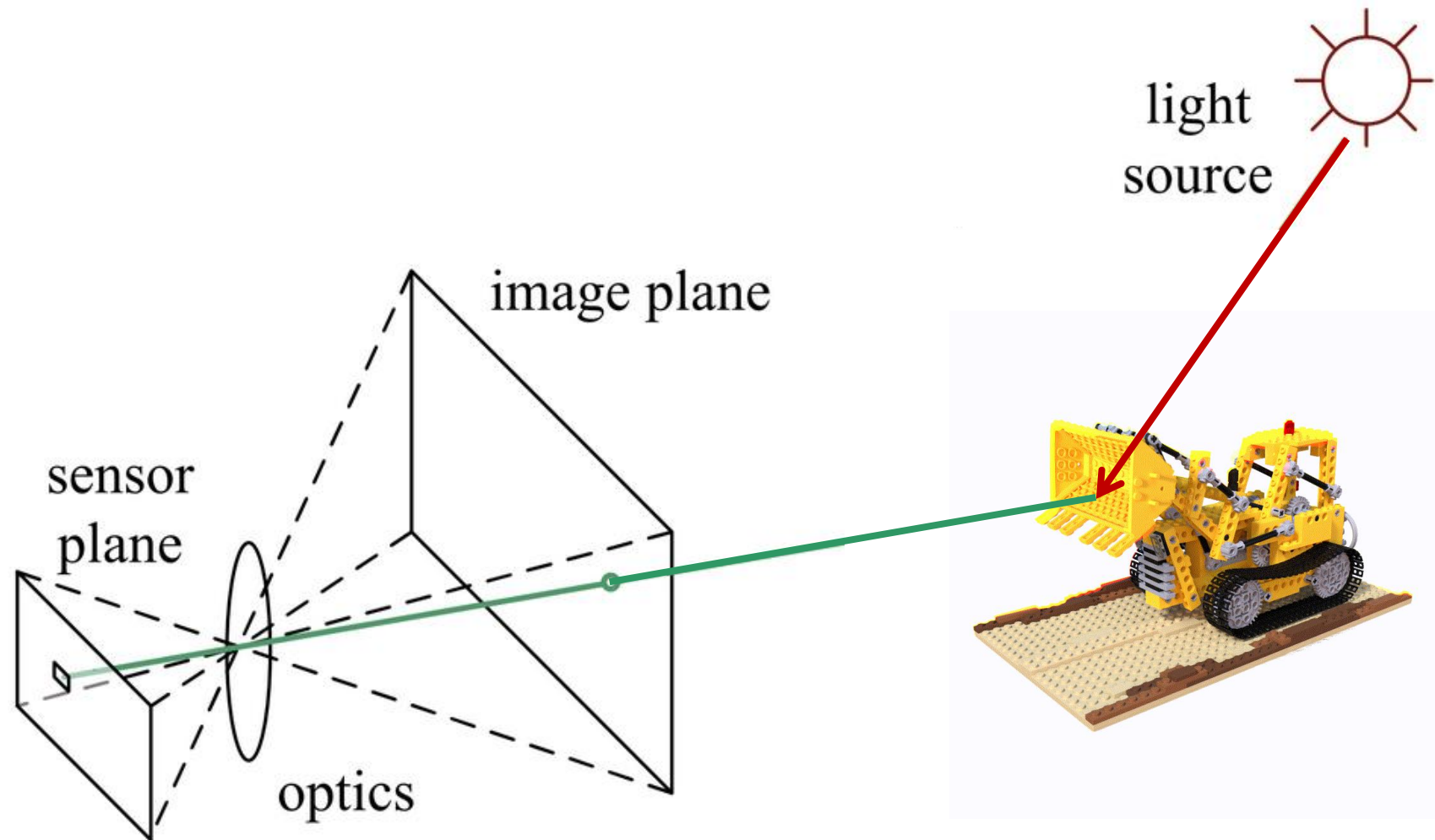
Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

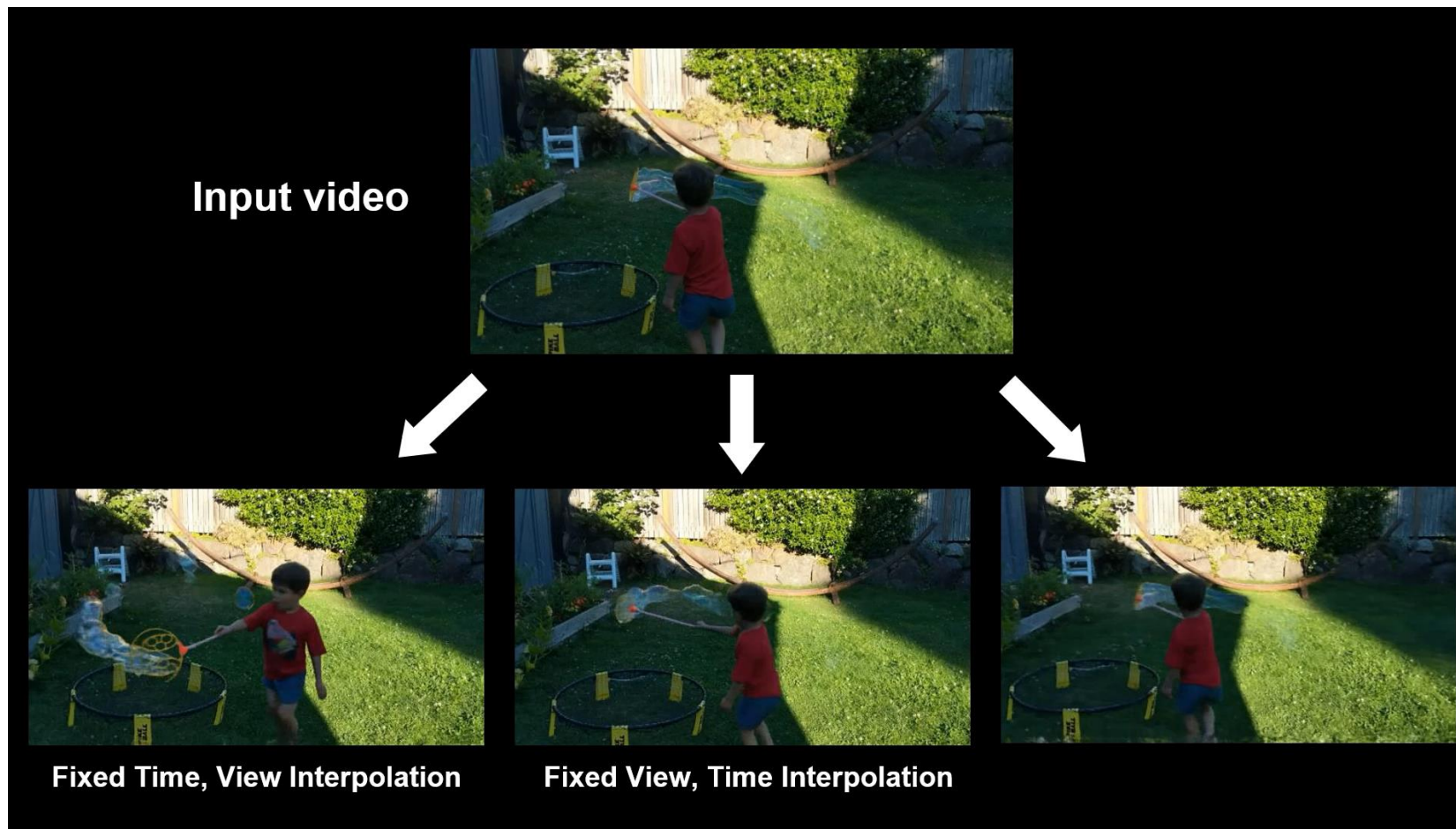
Optical Flow Only Describes 2D Motion. However, We Live in a 3D World, and 3D Motion Matters



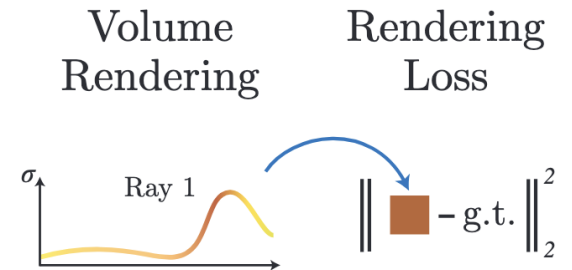
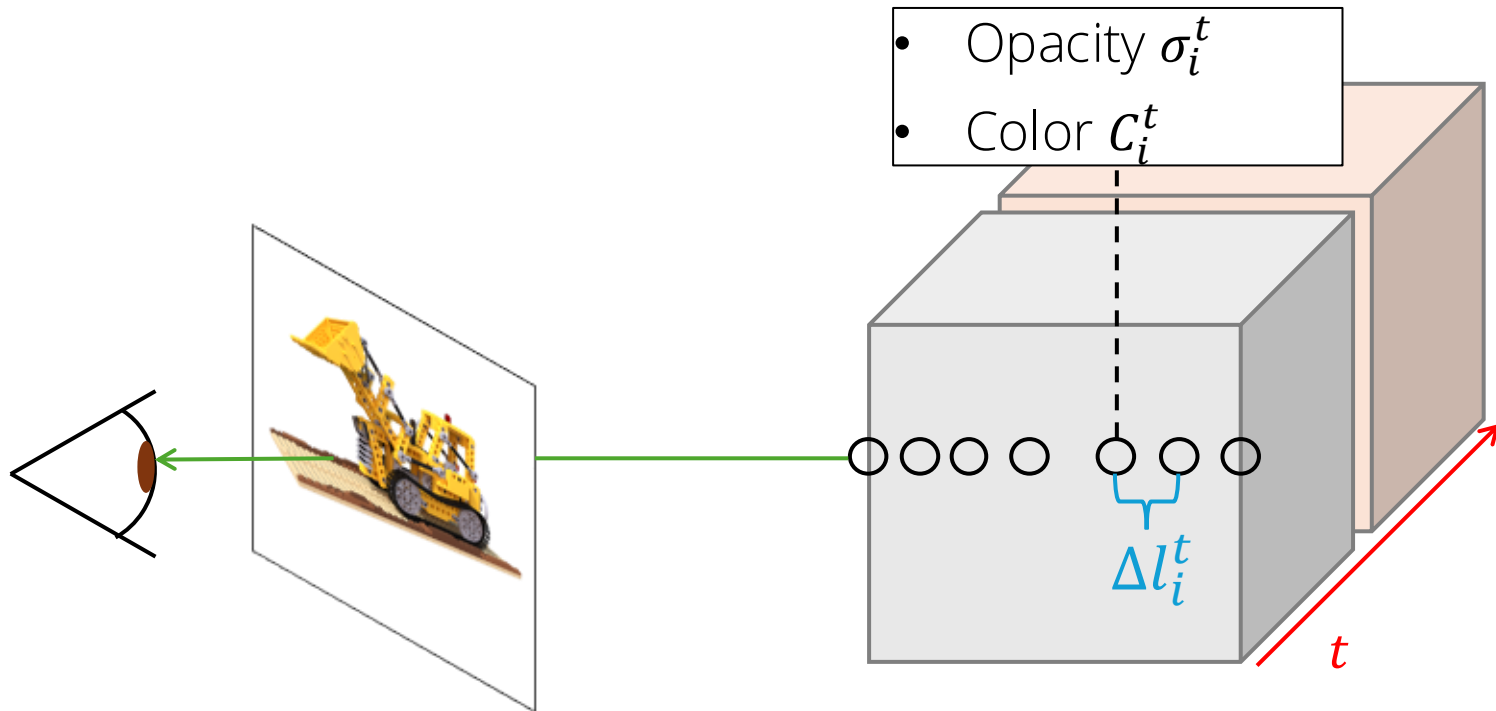
How Can We Model 3D World in Time?



How Can We Model 3D World in Time by Observing Monocular Videos?



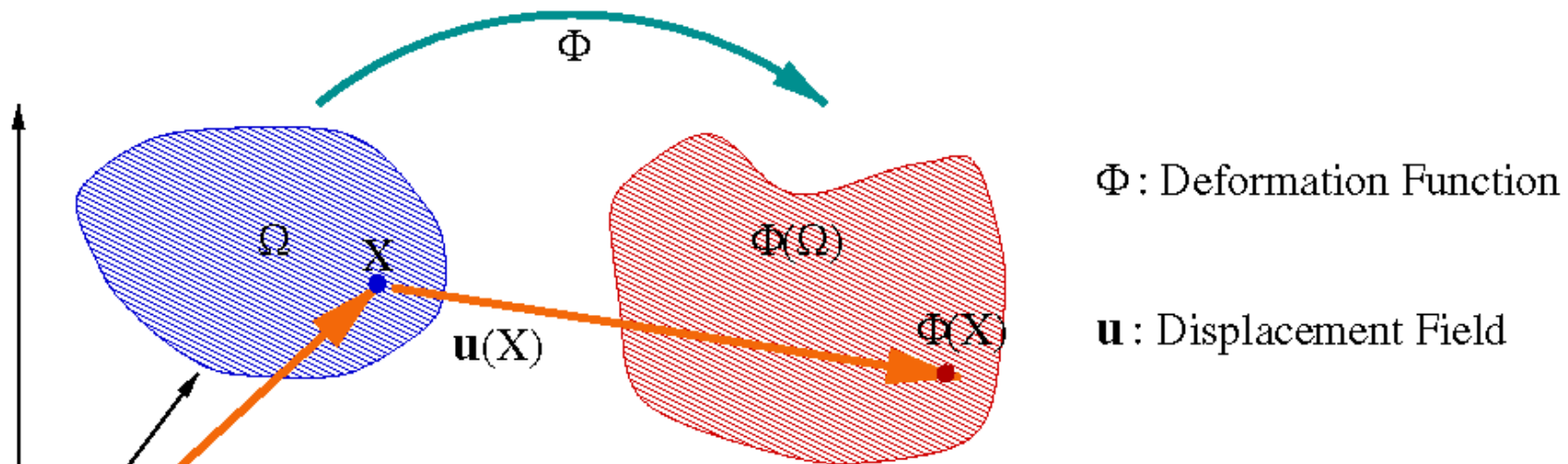
Idea 1: Extend 3D Volume Rendering to 4D Volume Rendering



$$C(x, \omega, t) \approx \sum_{i=1}^N \left(\prod_{j=1}^{i-1} e^{-\sigma_j^t \Delta l_j^t} \right) (1 - e^{-\sigma_i^t \Delta l_i^t}) C_i^t$$

Learning objective: $\min_{\sigma_i^t, C_i^t} \|C(x, \omega, t) - \text{g.t.}\|$

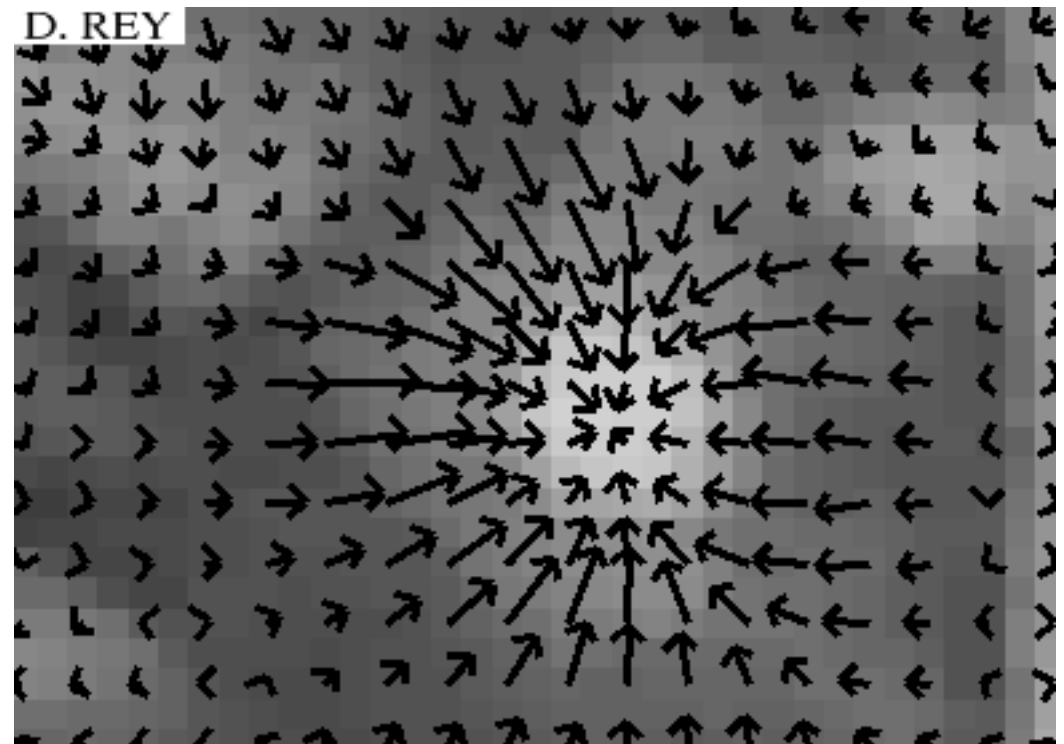
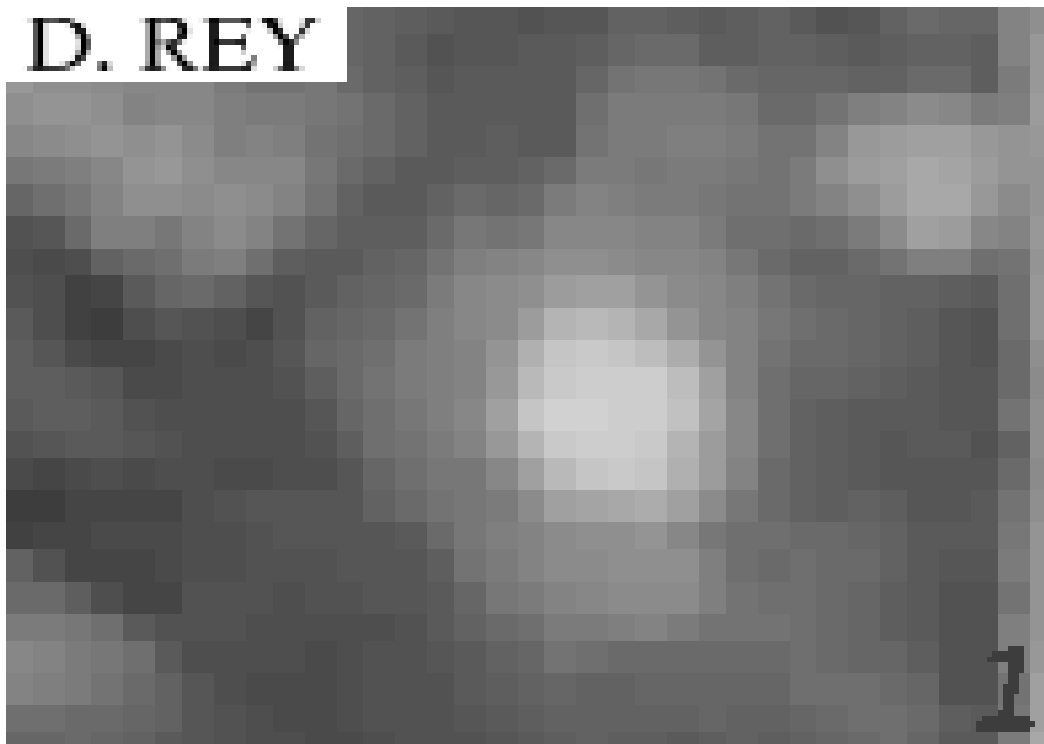
Idea 2: Model the Motion Explicitly



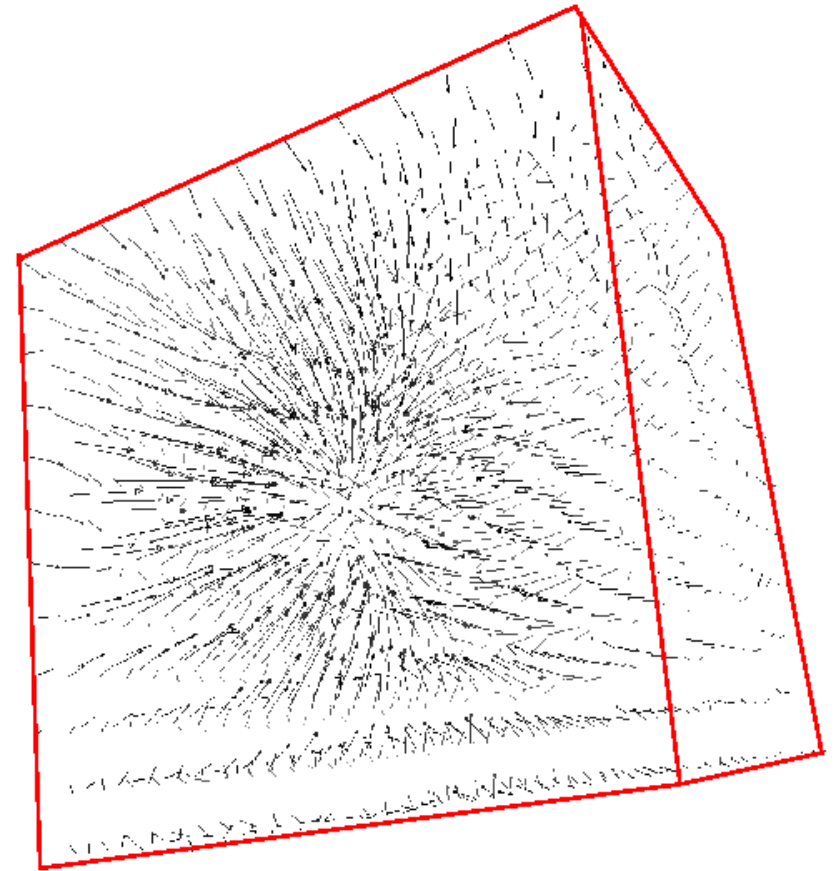
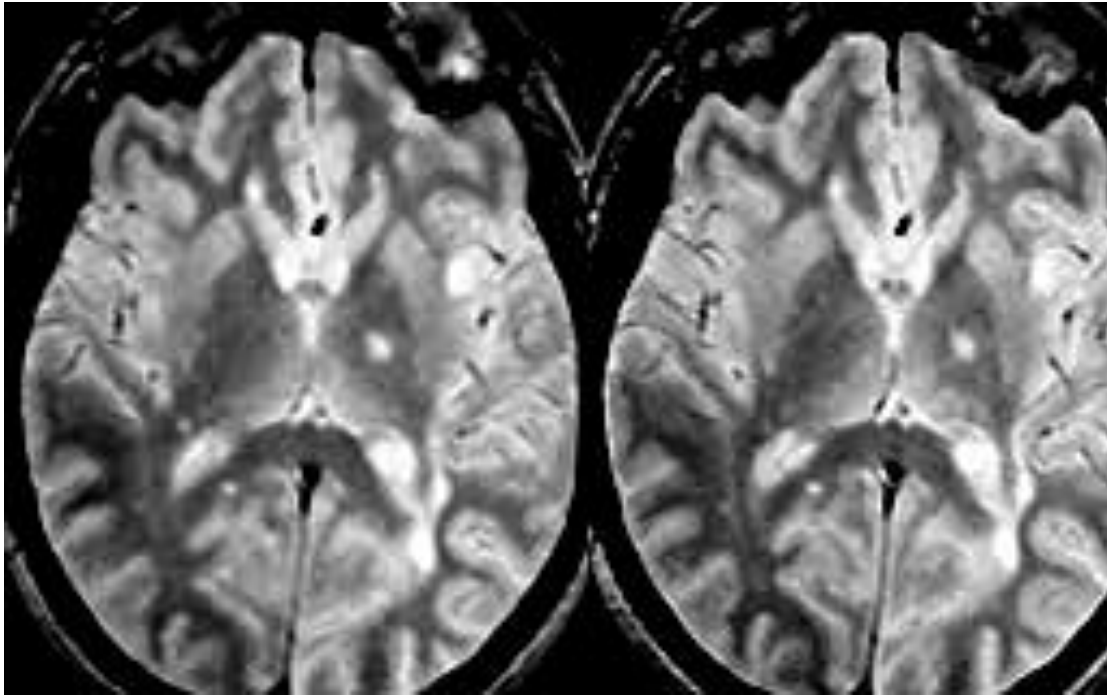
$$\phi(x) = x + \mathbf{u}(x) \implies \nabla\phi = Id + \nabla\mathbf{u}$$

Idea 2: Model the Motion Explicitly

- Optical flow: Pixel Apparent Motion
- 2D Motion flow: the exact particle 3D motion, projected/sliced at certain surface



Idea 2: Model the Motion Explicitly



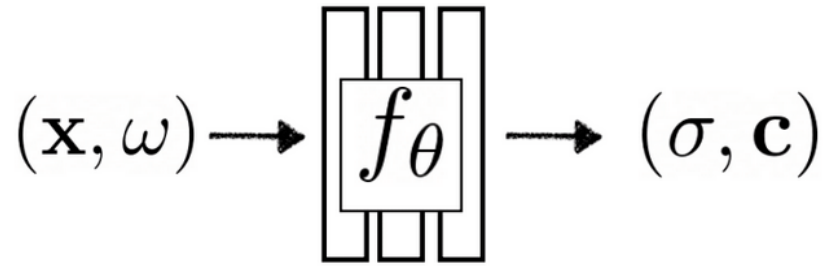
Idea 3: Model Multi-view Videos



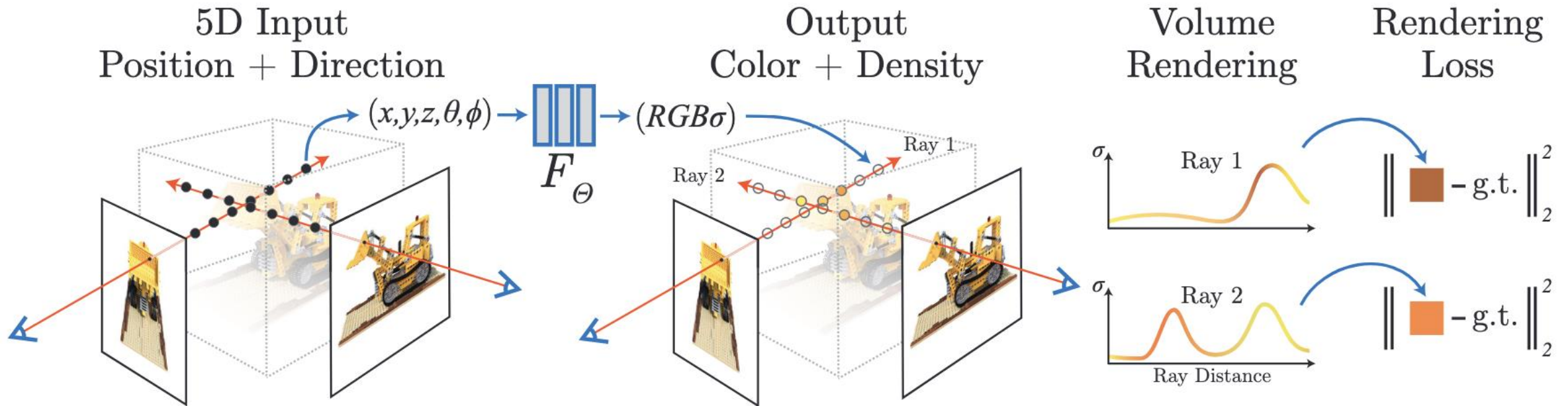
Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

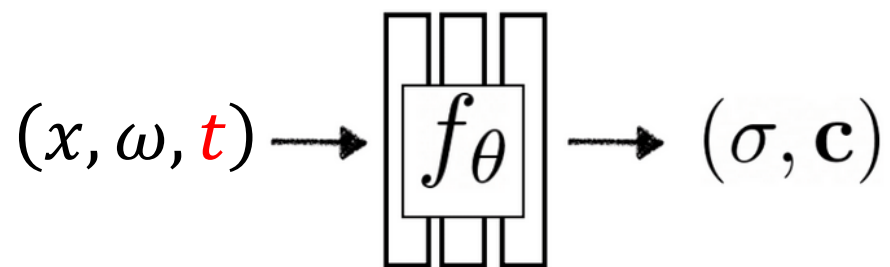
NEural Radiance Field (NERF)



- Use a neural network f_{θ} to predict at each point \mathbf{x} and view direction ω :
 - Opacity σ
 - Color \mathbf{c}

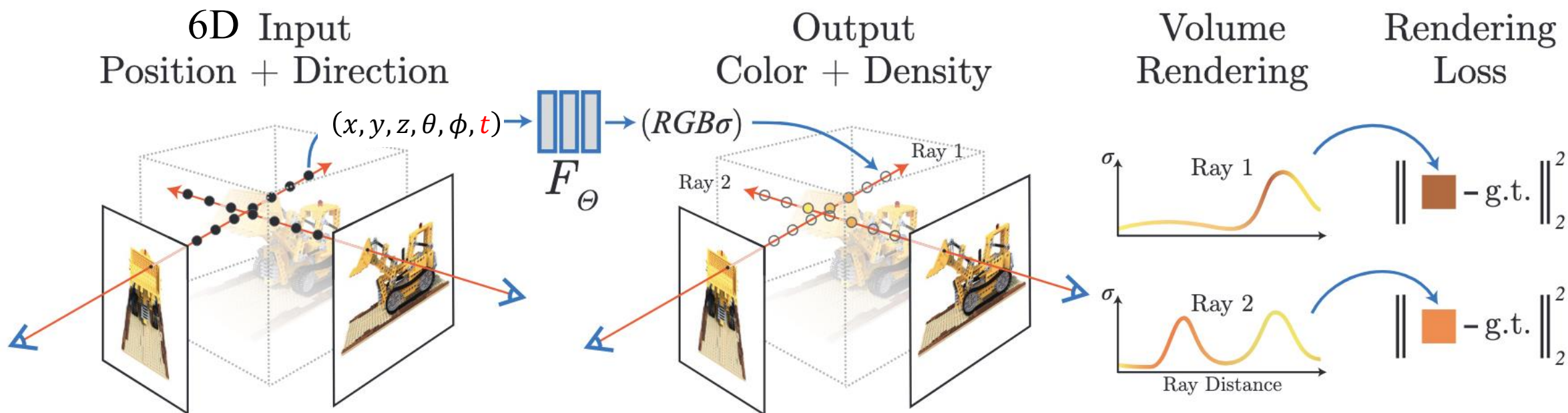


Idea 1: Extend NERF to SpatioTemporal NERF



- Use a neural network f_{θ} to predict at each point x , view direction ω , and time step t :

- Opacity σ
- Color \mathbf{c}



SpatioTemporal NERF Doesn't Work Out-of-Box

- **Problem:** A single video contains only one observation of the scene at any point in time. 3D geometry is under-constrained, e.g. “flat TV” solution



SpatioTemporal NERF Doesn't Work Out-of-Box

- Solution:

- Depth prediction loss: Enforce geometry by reconstructing depth

$$\mathcal{L}_{\text{depth}} = \sum_{(\mathbf{r}, t) \in \mathcal{R}} \left\| \frac{1}{\hat{D}(\mathbf{r}, t)} - \frac{1}{D(\mathbf{r}, t)} \right\|_2^2, \quad \hat{D}(\mathbf{r}, t) = \int_{s_n}^{s_f} T(s, t) \sigma(\mathbf{r}(s), t) s \, ds,$$

- Empty-space loss: Depth loss is not enough, since the predicted depth is a *weighted sum* of depth values along the ray, resulting in haze-like artifacts. Penalize non-zero volume densities away from the depth surface

$$\mathcal{L}_{\text{empty}} = \sum_{(\mathbf{r}, t) \in \mathcal{R}} \int_{s_n}^{d_t(\mathbf{u}) - \varepsilon} \sigma(\mathbf{r}(s), t) \, ds,$$

SpatioTemporal NERF Doesn't Work Out-of-Box

- **Problem:** A large portion of spaces that is hidden from the input frame's viewpoint at any given time is still not constrained



Occluded regions in the input video

SpatioTemporal NERF Doesn't Work Out-of-Box

- Solution:

- Static scene loss: Assume most parts of the scene are static, we can see the same regions, occluded at the current frame, from other *different viewpoints* at other *time*

$$\mathcal{L}_{\text{static}} = \sum_{(\mathbf{x}, t) \in \mathcal{X}} \|F(\mathbf{x}, t) - F(\mathbf{x}, t')\|_2^2,$$

$$F: (x, t) \rightarrow (c, \sigma)$$



(a) Input frame

(b) w/o static loss

(c) w/ static loss

Results



(a) Input

(b) Mesh

(c) Inpainted

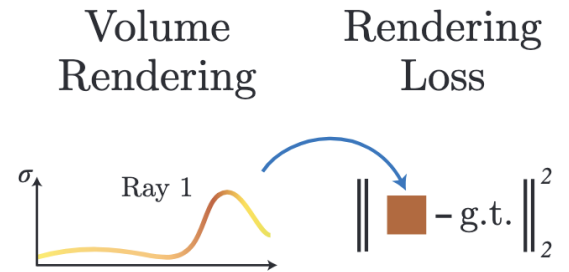
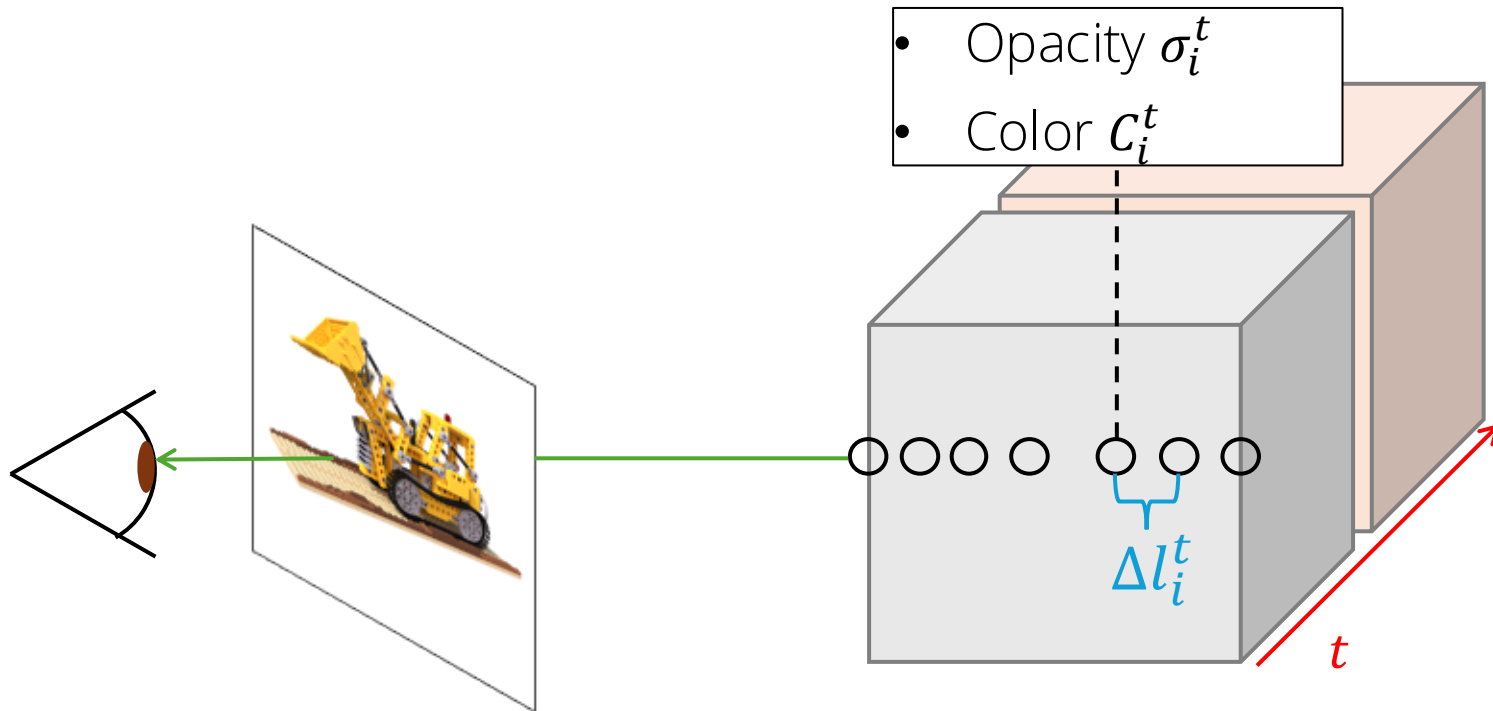
(d) NeRF-T

(e) Ours

Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

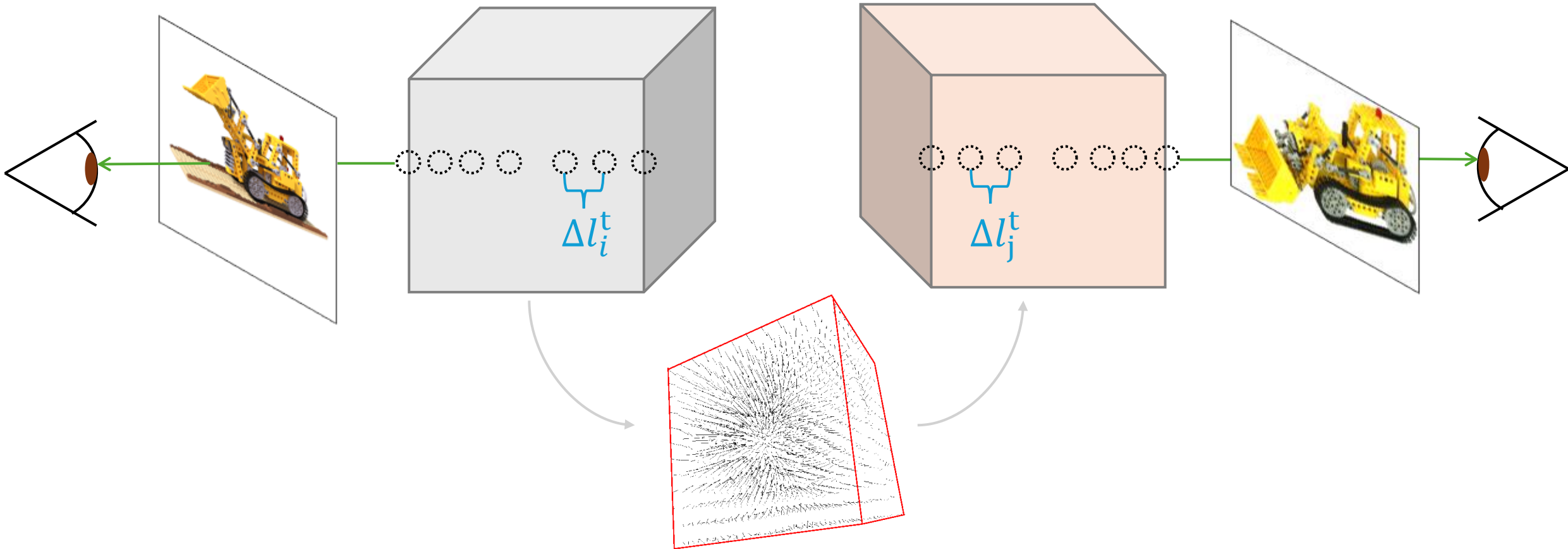
However, 4D NERF Doesn't Model Motion Explicitly. We don't Know the Position of the Same Point Across Time.



$$C(x, \omega, t) \approx \sum_{i=1}^N \left(\prod_{j=1}^{i-1} e^{-\sigma_j^t \Delta l_j^t} \right) \left(1 - e^{-\sigma_i^t \Delta l_i^t} \right) C_i^t$$

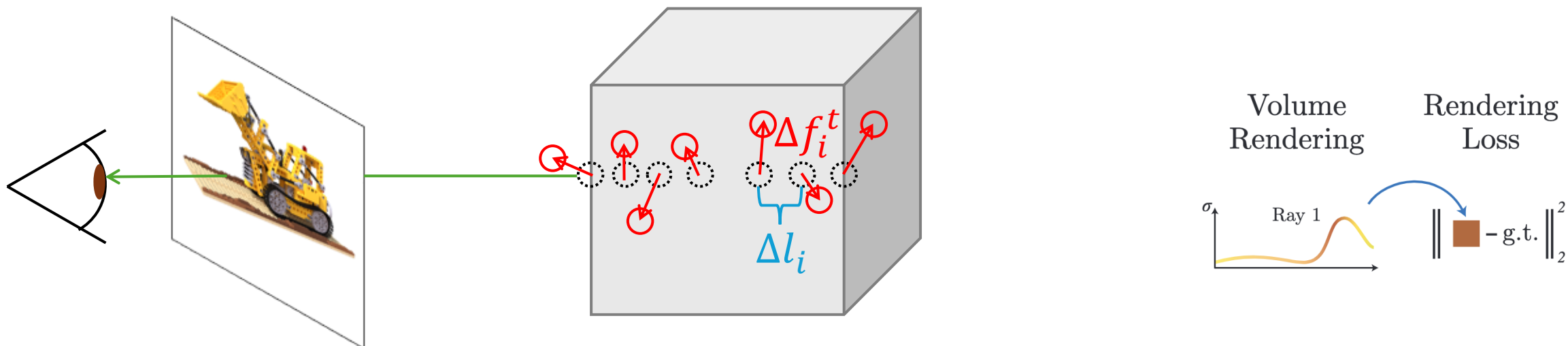
Learning objective: $\min_{\sigma_i^t, C_i^t} \|C(x, \omega, t) - \text{g.t.}\|$

Idea 2: Model the Motion Explicitly



Idea 2: Model the Motion Explicitly

Δf_i^t denotes the motion field

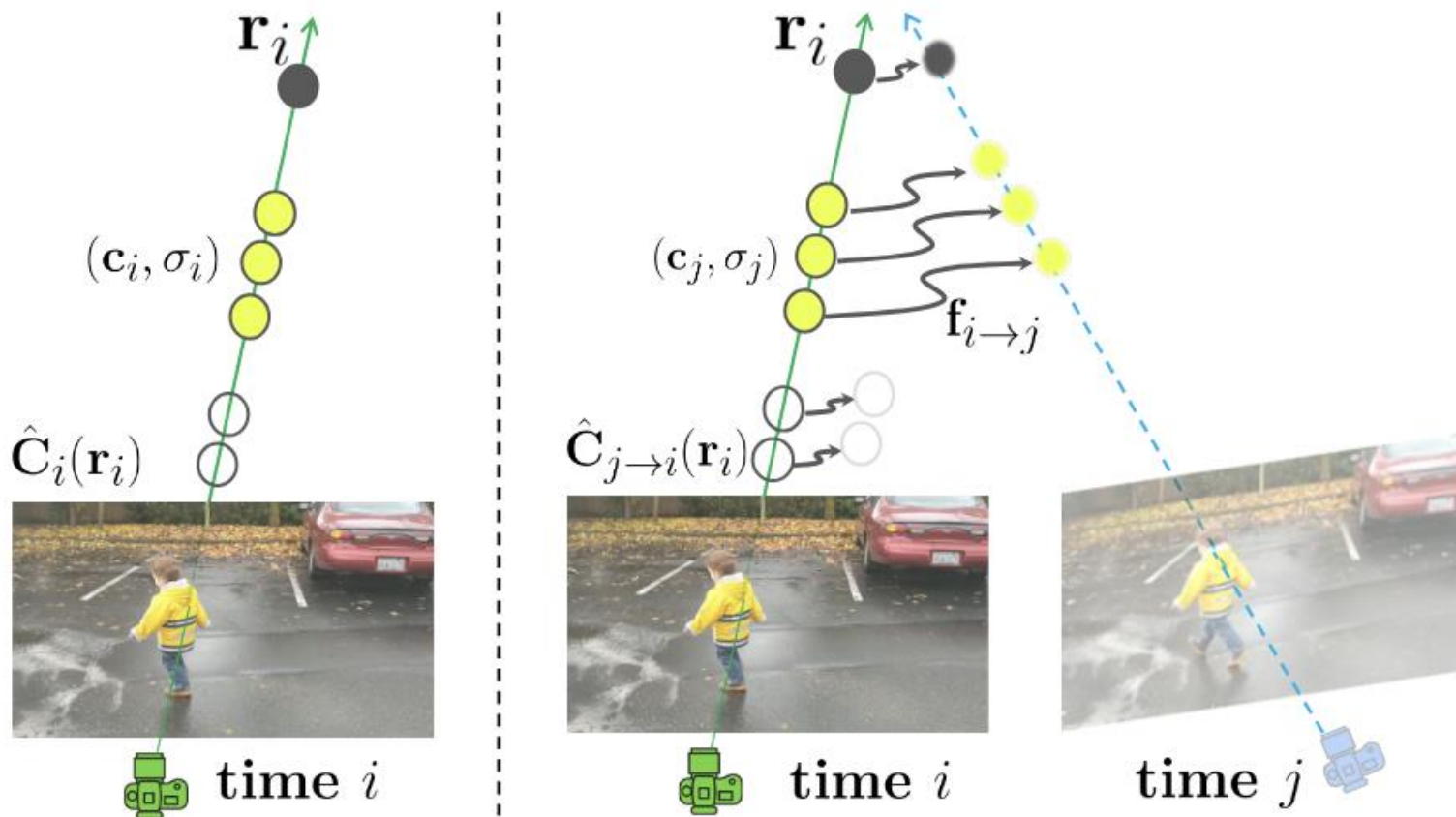


$$C(x, \omega, t) \approx \sum_{i=1}^N \left(\prod_{j=1}^{i-1} e^{-\sigma_{j+\Delta f_j^t} \Delta l_j} \right) (1 - e^{-\sigma_{i+\Delta f_i^t} \Delta l_i}) C_{i+\Delta f_i^t}$$

Learning objective: $\min_{\sigma_i, C_i} \|C(x, \omega, t) - \text{g.t.}\|$

Neural Scene Flow Fields for Space-Time View Synthesis of Dynamic Scenes

- To render a dynamic scene, we need to model the appearance (\mathbf{c}_i, σ_i) and the motion $(\mathbf{f}_{i \rightarrow j})$ of the scene
- Define $\mathbf{f}_{i \rightarrow j}(\mathbf{x})$ as the 3D offset vector moving position \mathbf{x} from time i to time j



Scene Flow Field Warping

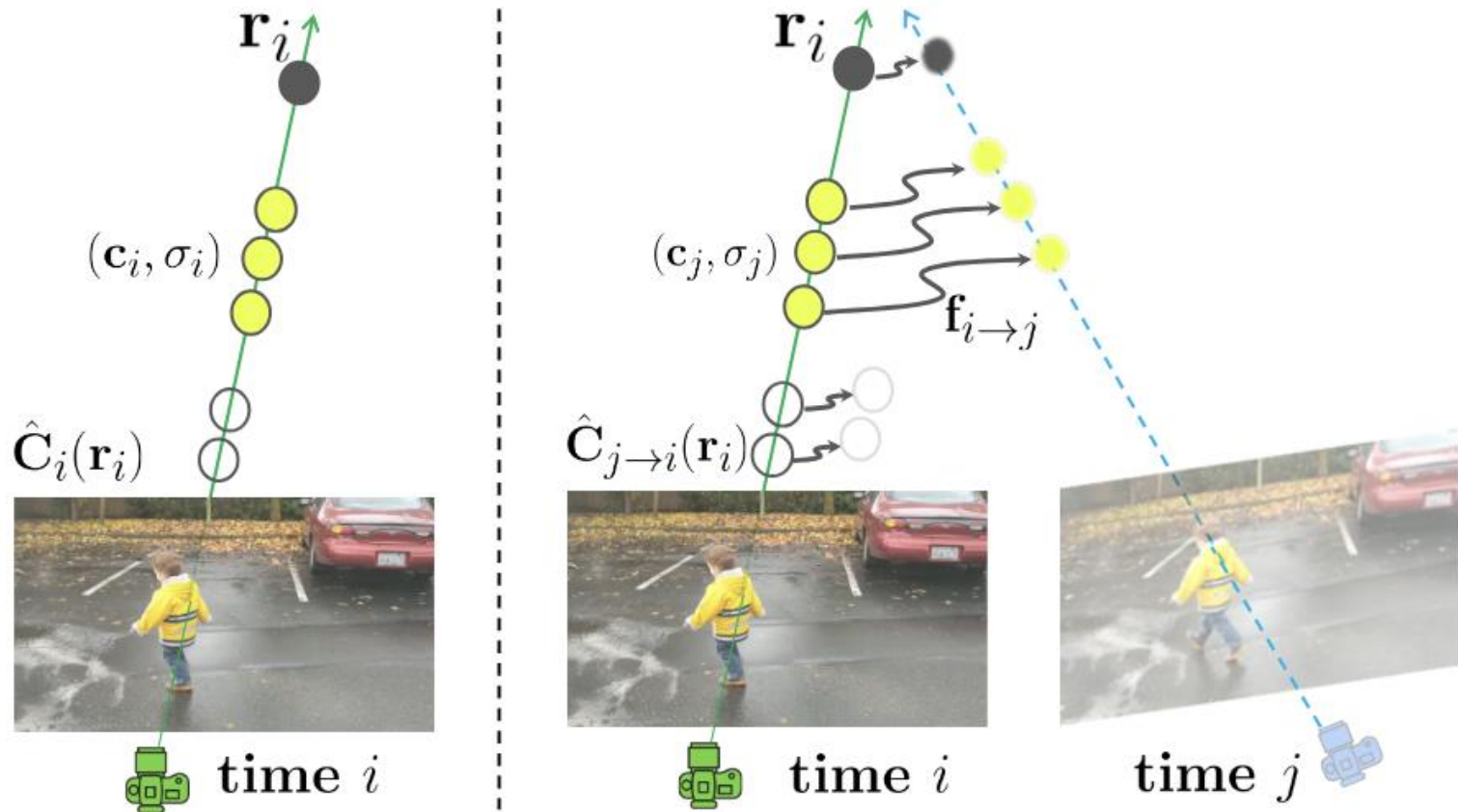
- Idea: We can render an object at the current time i , if we look up the scene created at time j

$$\hat{\mathbf{C}}_{j \rightarrow i}(\mathbf{r}_i) =$$

$$\int_{t_n}^{t_f} T_j(t) \sigma_j(\mathbf{r}_{i \rightarrow j}(t)) \mathbf{c}_j(\mathbf{r}_{i \rightarrow j}(t), \mathbf{d}_i) dt$$

where $\mathbf{r}_{i \rightarrow j}(t) = \mathbf{r}_i(t) + \mathbf{f}_{i \rightarrow j}(\mathbf{r}_i(t))$.

The new position at time j for the position \mathbf{x} from time i



Scene Flow Field Warping

- Idea: We can render an object at the current time i , if we look up the scene created at time j

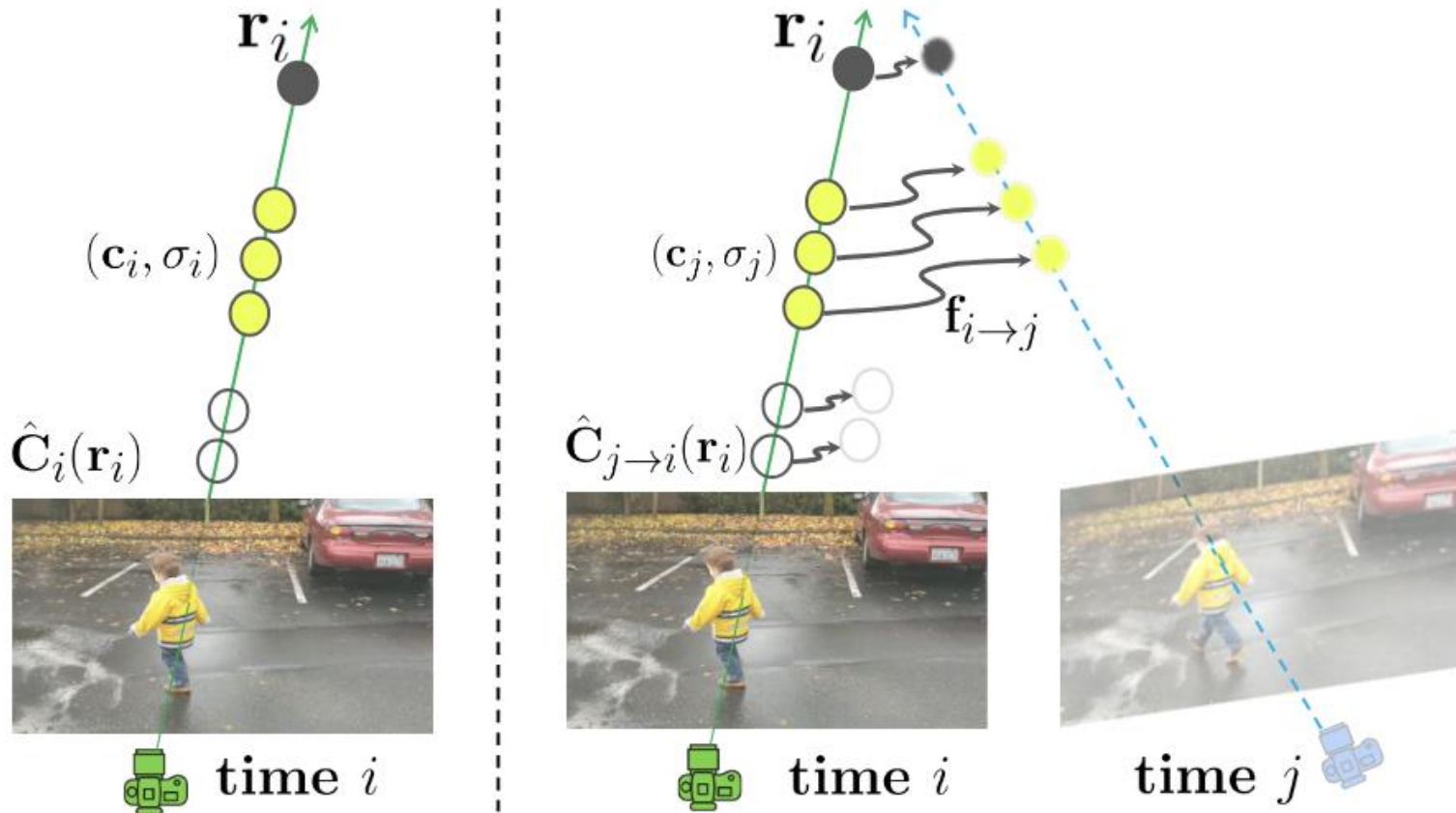
$$\hat{\mathbf{C}}_{j \rightarrow i}(\mathbf{r}_i) =$$

$$\int_{t_n}^{t_f} T_j(t) \sigma_j(\mathbf{r}_{i \rightarrow j}(t)) \mathbf{c}_j(\mathbf{r}_{i \rightarrow j}(t), \mathbf{d}_i) dt$$

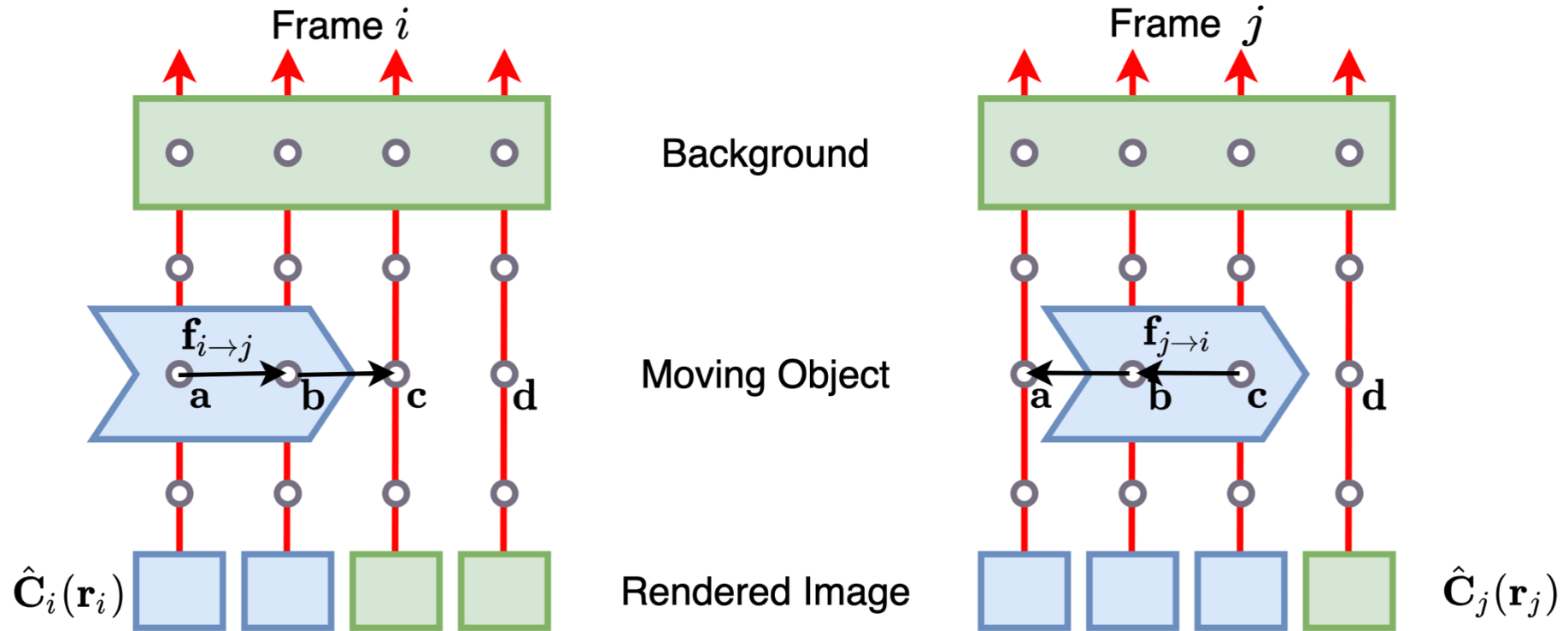
where $\mathbf{r}_{i \rightarrow j}(t) = \mathbf{r}_i(t) + \mathbf{f}_{i \rightarrow j}(\mathbf{r}_i(t))$.

- Optimize appearance and motion by rendering:

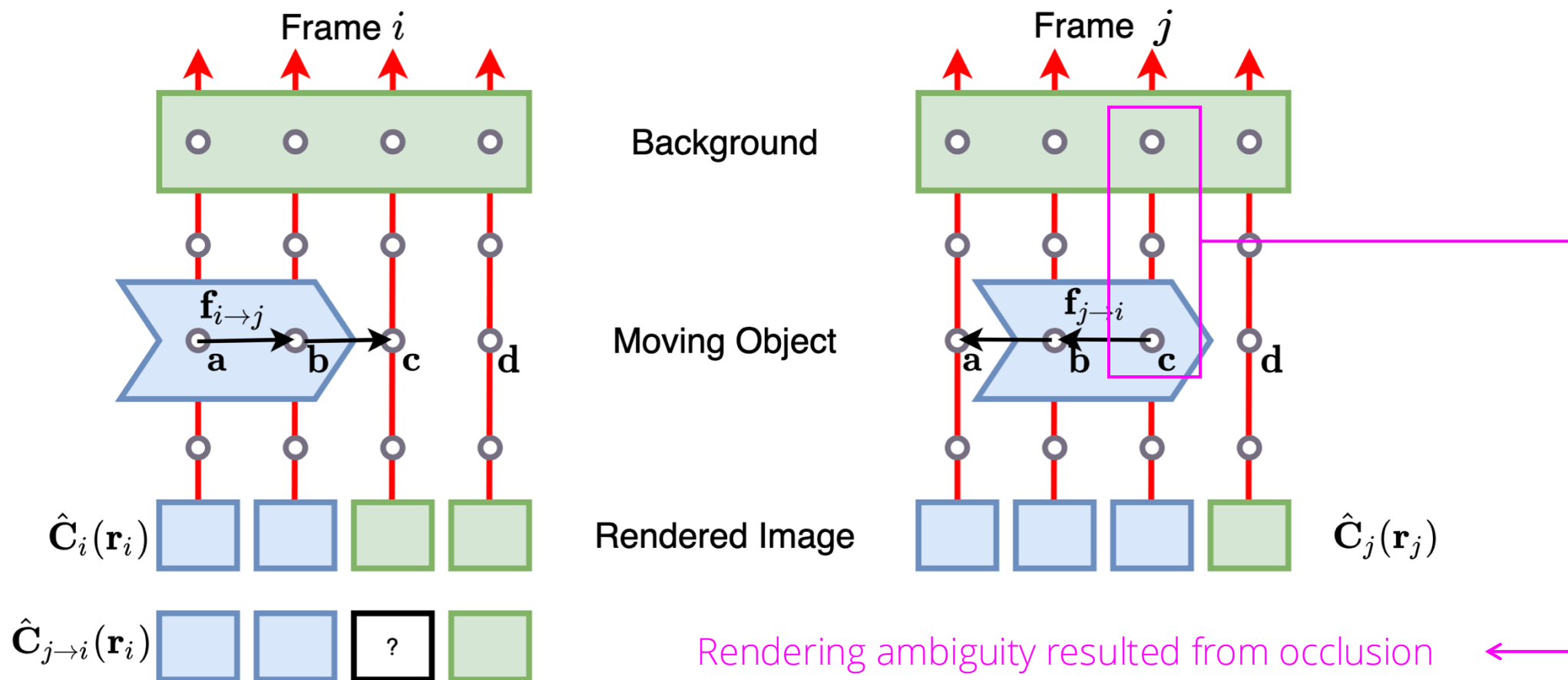
$$\mathcal{L}_{\text{pho}} = \sum_{\mathbf{r}_i} \sum_{j \in \mathcal{N}(i)} \|\hat{\mathbf{C}}_{j \rightarrow i}(\mathbf{r}_i) - \mathbf{C}_i(\mathbf{r}_i)\|_2^2$$



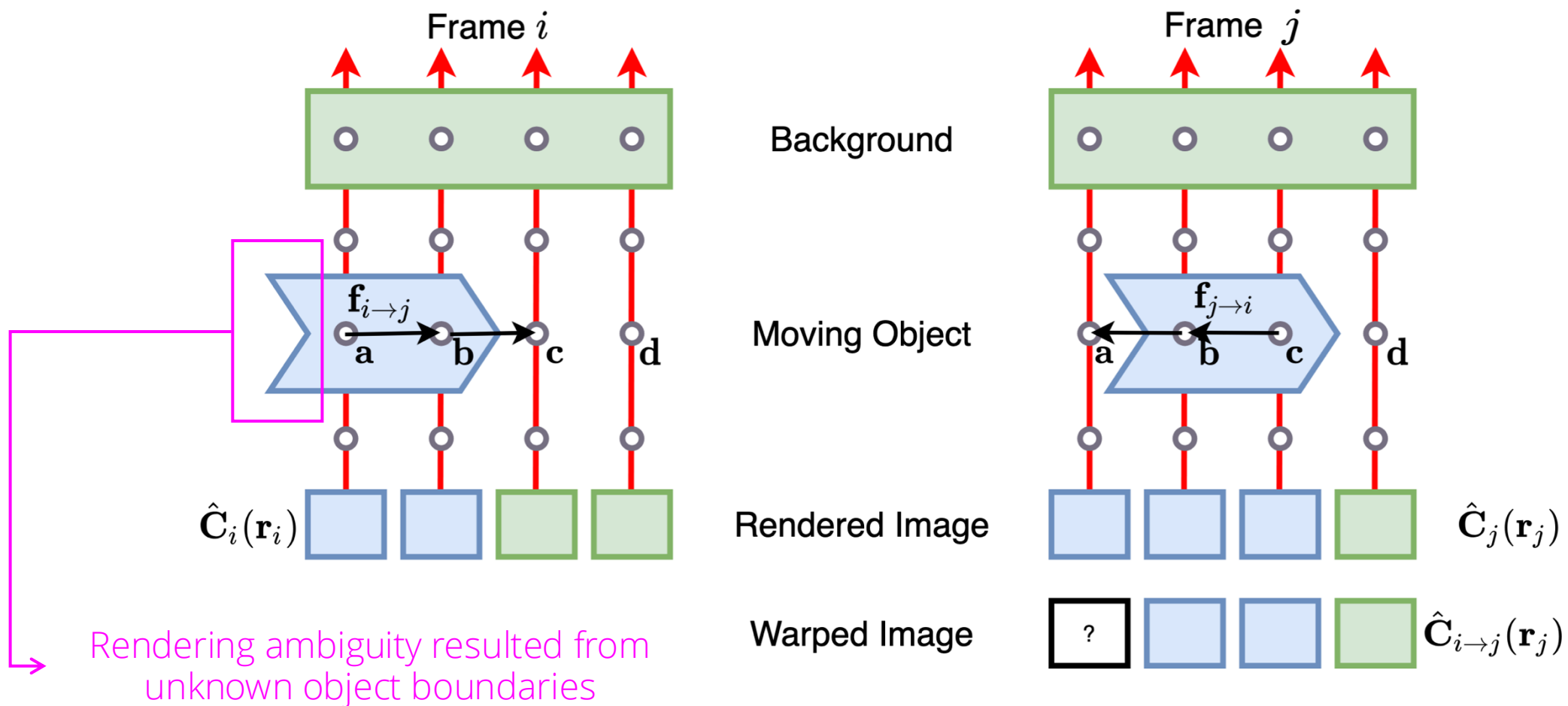
Rendering Loss is Ambiguous with Occlusion



Rendering Loss is Ambiguous with Occlusion



Rendering Loss is Ambiguous with Occlusion



Rendering Loss is Ambiguous with Occlusion

- Solution:

- Predict if the occlusion of points from time i to j : occlusion weight field $w_{i \rightarrow j} \in [0, 1]$
- The occlusion weight indicates where and how much strength the temporal photo-consistency loss should be applied
- Aggregate occlusion weights along a ray:

$$\hat{W}_{j \rightarrow i}(\mathbf{r}_i) = \int_{t_n}^{t_f} T_j(t) \sigma_j(\mathbf{r}_{i \rightarrow j}(t)) w_{i \rightarrow j}(\mathbf{r}_i(t)) dt$$

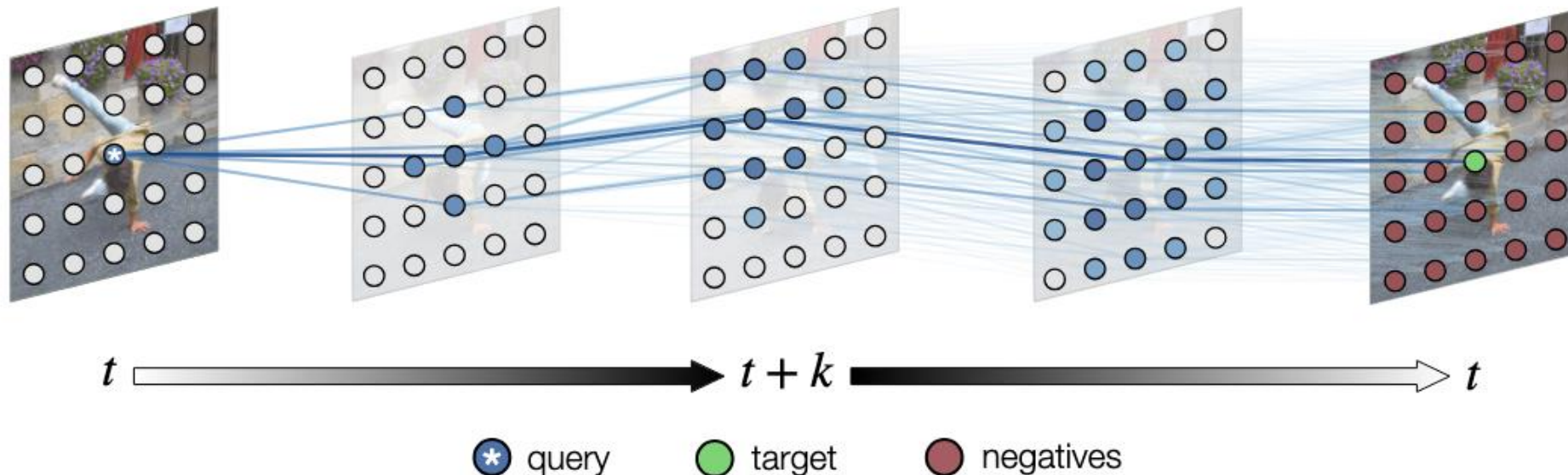
$$\mathcal{L}_{\text{pho}} = \sum_{\mathbf{r}_i} \sum_{j \in \mathcal{N}(i)} \hat{W}_{j \rightarrow i}(\mathbf{r}_i) \|\hat{\mathbf{C}}_{j \rightarrow i}(\mathbf{r}_i) - \mathbf{C}_i(\mathbf{r}_i)\|_2^2$$

$$+ \beta_w \sum_{\mathbf{x}_i} \|w_{i \rightarrow j}(\mathbf{x}_i) - 1\|_1,$$

Accurate Flow is Difficult to Learn

- Solution:
 - Cycle consistency constraint: a point should stay at the same location when moving forward and backward

$$\mathcal{L}_{\text{cyc}} = \sum_{\mathbf{x}_i} \sum_{j \in i \pm 1} w_{i \rightarrow j} \|\mathbf{f}_{i \rightarrow j}(\mathbf{x}_i) + \mathbf{f}_{j \rightarrow i}(\mathbf{x}_{i \rightarrow j})\|_1$$



Accurate Flow is Difficult to Learn

- **Solution:** Data-driven priors
 - Guidance with predicted 2D optical flow

$$\mathcal{L}_{\text{geo}} = \sum_{\mathbf{r}_i} \sum_{j \in \{i \pm 1\}} \|\hat{\mathbf{p}}_{i \rightarrow j}(\mathbf{r}_i) - \mathbf{p}_{i \rightarrow j}(\mathbf{r}_i)\|_1.$$

$\hat{\mathbf{p}}_{i \rightarrow j}$ denotes the perspective projection of the expected 3D point location displaced by the scene flow

- Guidance with predicted depth

$$\mathcal{L}_z = \sum_{\mathbf{r}_i} \|\hat{Z}_i^*(\mathbf{r}_i) - Z_i^*(\mathbf{r}_i)\|_1$$

\hat{Z}_i denotes the expected termination depth computed along a ray (check previous slides)

Accurate Flow is Difficult to Learn

- **Solution:** Decouple static and dynamic scene rendering
 - Create a static and dynamic radiance field
 - Blend both fields by factor v

$$\hat{\mathbf{C}}_i^{\text{cb}}(\mathbf{r}_i) = \int_{t_n}^{t_f} T_i^{\text{cb}}(t) \sigma_i^{\text{cb}}(t) \mathbf{c}_i^{\text{cb}}(t) dt,$$

$$\sigma_i^{\text{cb}}(t) \mathbf{c}_i^{\text{cb}}(t) =$$

$$v(t) \mathbf{c}(t) \sigma(t) + (1-v(t)) \mathbf{c}_i(t) \sigma_i(t)$$

$$\mathcal{L}_{\text{cb}} = \sum_{\mathbf{r}_i} \|\hat{\mathbf{C}}_i^{\text{cb}}(\mathbf{r}_i) - \mathbf{C}_i(\mathbf{r}_i)\|_2^2.$$



Combined render

Static only

Dynamic only

Results



3D Photo [64]

NeRF [46]

Yoon *et al.* [81]

Ours

GT

Input frames



Fixed Time, View Interpolation



Fixed View, Time Interpolation



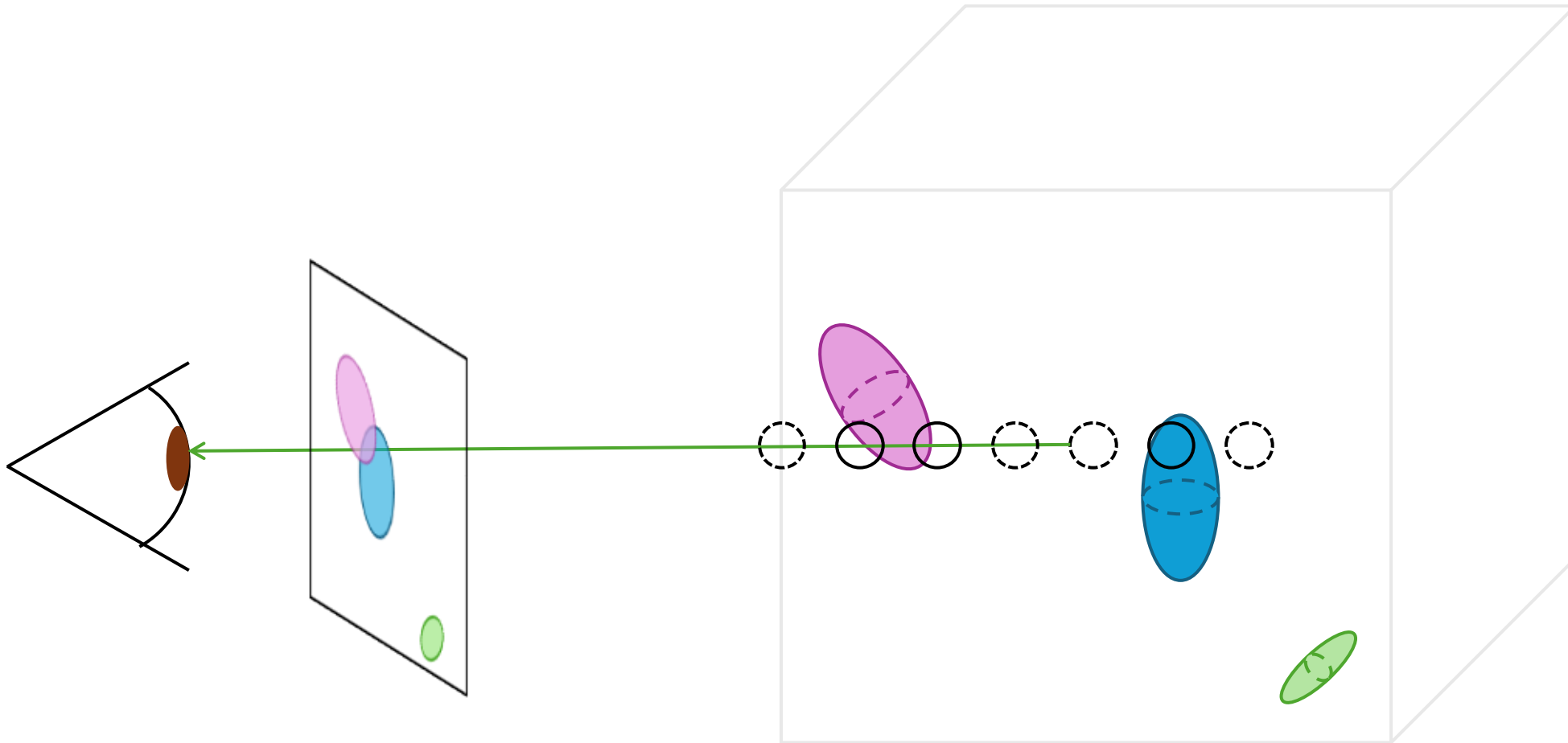
Space-Time Interpolation



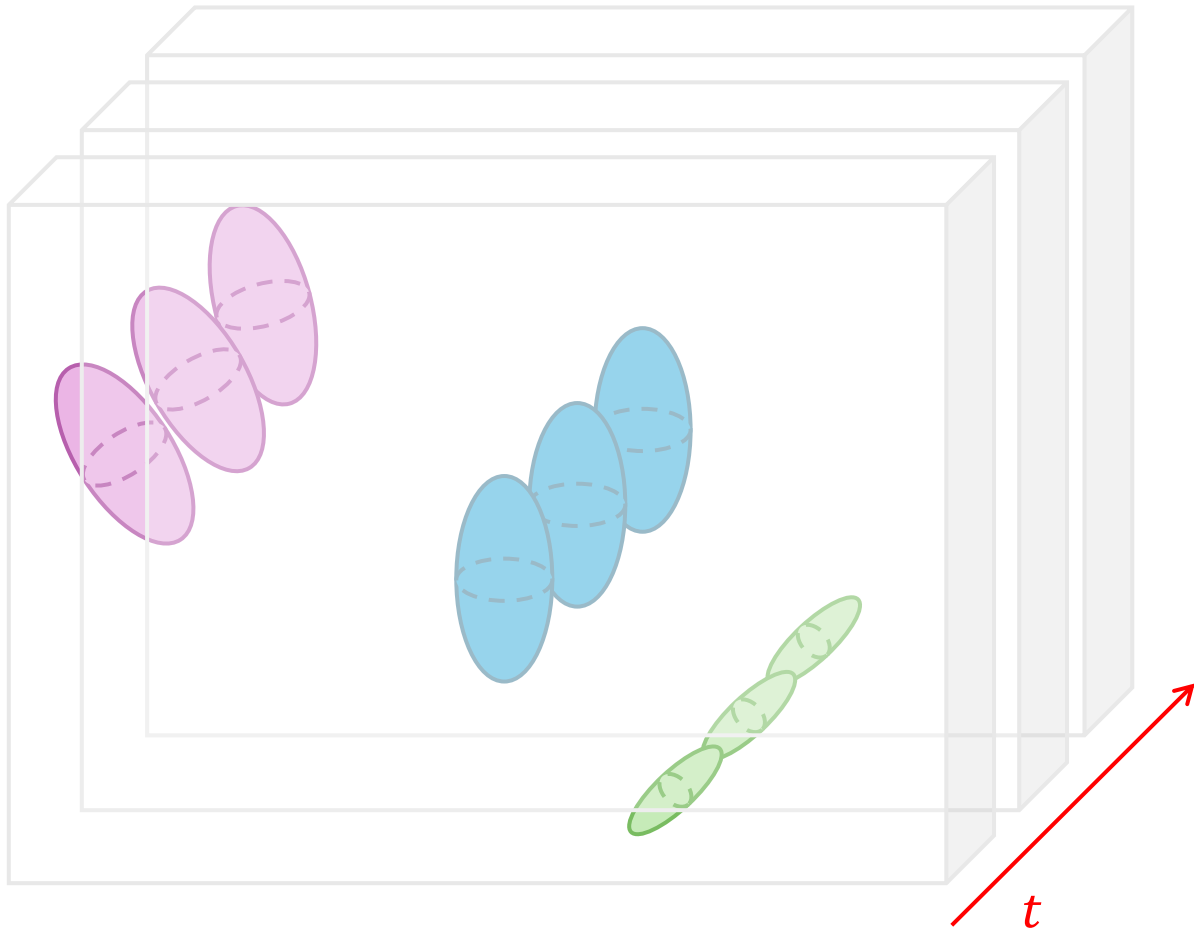
Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

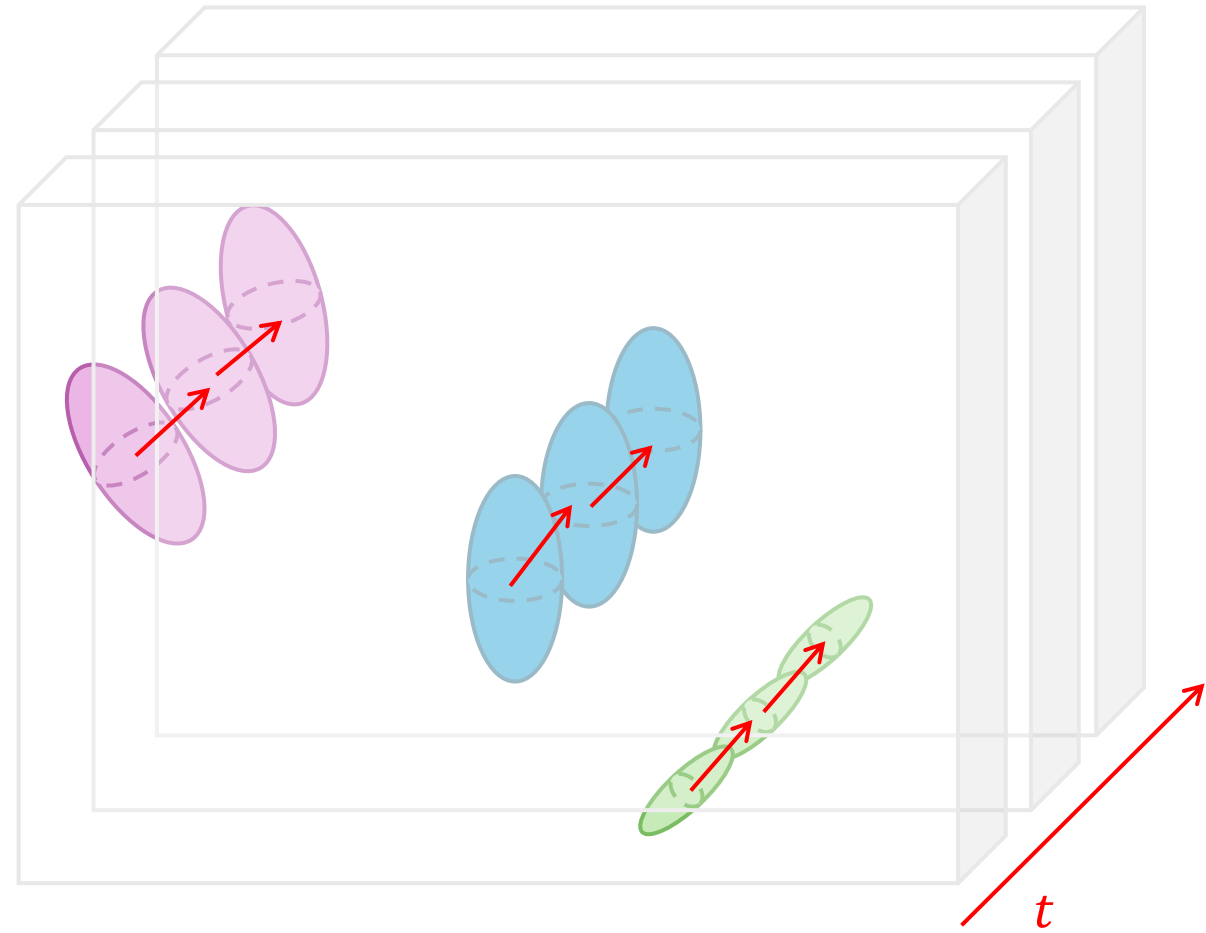
Volume Rendering with 3D Gaussian Splatting



Two Ways of 4D Volume Rendering



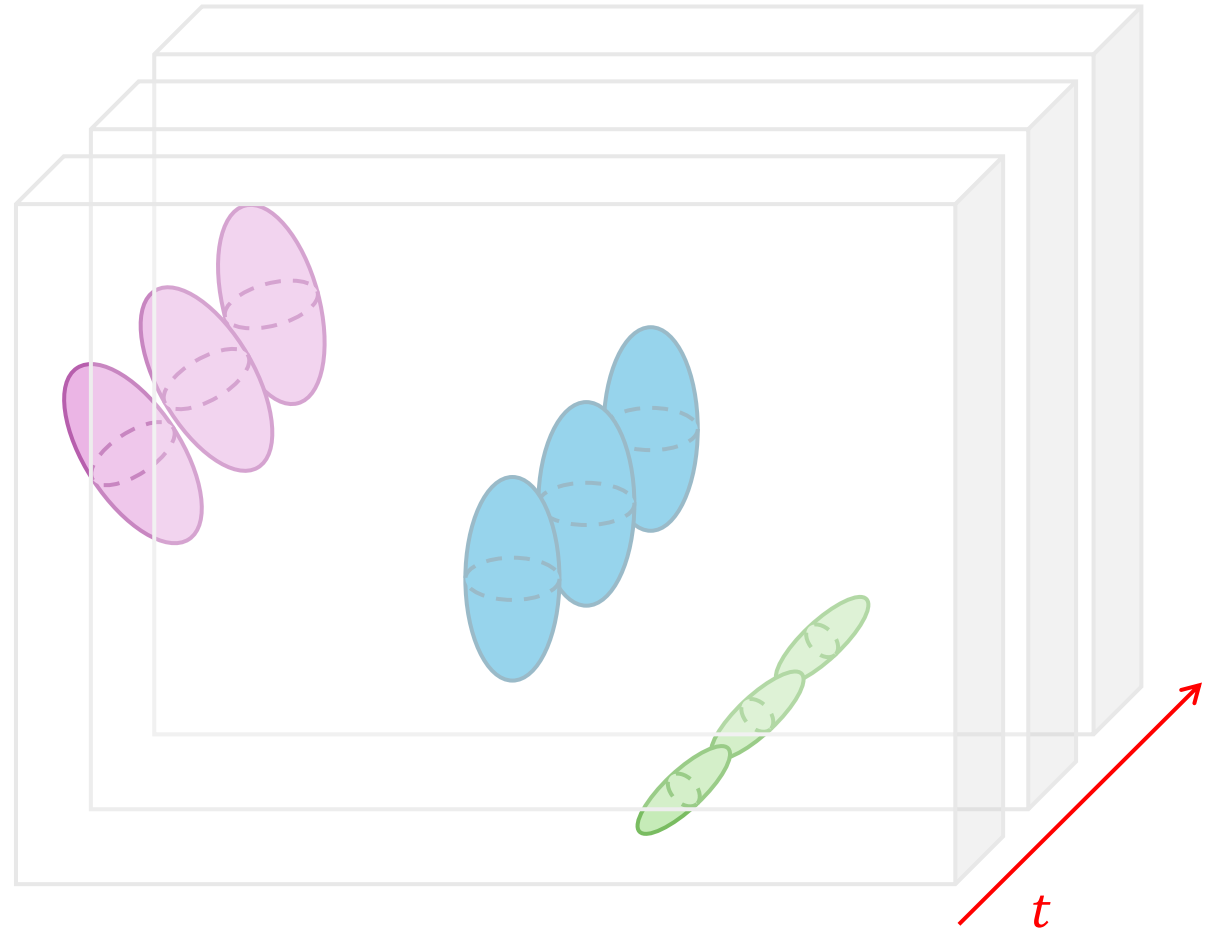
4DGS: Represent a dynamic scene with a set of 4D GS



Dynamic 3DGS: Represent a dynamic scene with a set of 3D GS and their motions

Method 1: Extend 3DGS to 4DGS

- **Idea:** Generalize 3D GS to a 4D formulation.
- **How to render?** Slice a 4D GS into a 3D GS at time t , and render it with 2D rasterization



A Simplified Illustration for Slicing a 4D GS into a 3D GS

- **Idea:** Generalize 3D GS to a 4D formulation.
- **How to render?** Slice a 4D GS into a 3D GS at time t , and render it with 2D rasterization
- **What does it mean?** In a small interval, all motions can be approximated by linear motions, and more complex non-linear cases can be represented by combination of multiple Gaussians

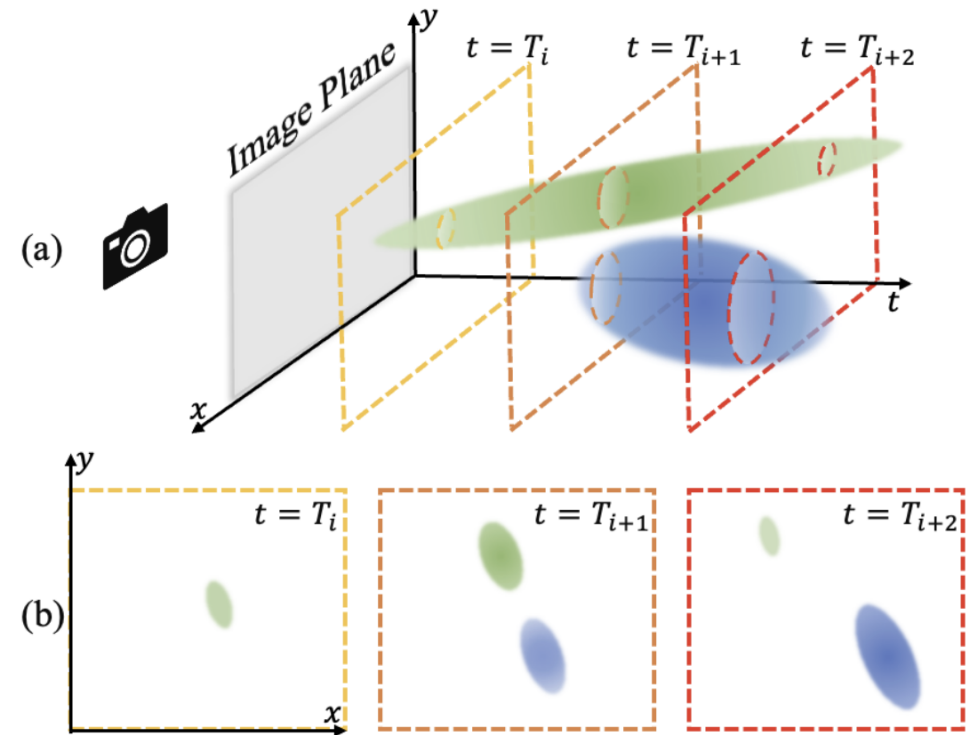


Figure 2: A Simplified 2D Illustration of the Proposed Temporal Slicing. (a) We model 2D dynamics with 3D XYT ellipsoids and slice them with different time queries. (b) The sliced 3D ellipsoids form 2D dynamic ellipses at each timestamp.

Temporally Sliced 4DGS is a 3DGS

- Formulation:

$$G_{4D}(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{4D})^T \boldsymbol{\Sigma}_{4D}^{-1}(\mathbf{x}-\boldsymbol{\mu}_{4D})}$$

- A 4D center position: $\boldsymbol{\mu}_{4D} = (\mu_x, \mu_y, \mu_z, \mu_t)$
- A 4D covariance matrix: $\boldsymbol{\Sigma}_{4D} = \mathbf{R}_{4D} \mathbf{S}_{4D} \mathbf{S}_{4D}^T \mathbf{R}_{4D}^T$
- \mathbf{S}_{4D} denotes a diagonal 4D matrix for scaling
- \mathbf{R}_{4D} denotes a 4D rotator matrix (a 4D GS reduces to a 3D GS if the last four components of \mathbf{R}_{4D} are set to 0)

Generalize to both static and dynamic scene!

Temporally Sliced 4DGS is a 3DGS

- We can write the 4D covariance matrix as:

$$\Sigma_{4D} = \begin{pmatrix} \mathbf{U} & \mathbf{V} \\ \mathbf{V}^T & \mathbf{W} \end{pmatrix} \text{ and } \Sigma_{4D}^{-1} = \begin{pmatrix} \mathbf{A} & \mathbf{M} \\ \mathbf{M}^T & \mathbf{Z} \end{pmatrix}$$

- A temporally sliced 4DGS is a 3DGS

$$G_{3D}(\mathbf{x}, t) = e^{-\frac{1}{2}\lambda(t-\mu_t)^2} e^{-\frac{1}{2}[\mathbf{x}-\boldsymbol{\mu}(t)]^T \Sigma_{3D}^{-1} [\mathbf{x}-\boldsymbol{\mu}(t)]},$$

$$\lambda = \mathbf{W}^{-1}, \quad \Sigma_{3D} = \mathbf{A}^{-1} = \mathbf{U} - \frac{\mathbf{V}\mathbf{V}^T}{\mathbf{W}}, \quad \boldsymbol{\mu}(t) = (\mu_x, \mu_y, \mu_z)^T + (t - \mu_t) \frac{\mathbf{V}}{\mathbf{W}}.$$

Let's Take a Closer Look

- A temporally sliced 4DGS is a 3DGS

$$G_{3D}(\mathbf{x}, t) = \boxed{e^{-\frac{1}{2}\lambda(t-\mu_t)^2}} e^{-\frac{1}{2}[\mathbf{x}-\boldsymbol{\mu}(t)]^T \boldsymbol{\Sigma}_{3D}^{-1}[\mathbf{x}-\boldsymbol{\mu}(t)]},$$

$$\lambda = \mathbf{W}^{-1},$$

$$\boldsymbol{\mu}(t) = (\mu_x, \mu_y, \mu_z)^T + (t - \mu_t) \frac{\mathbf{V}}{\mathbf{W}}.$$

- As t goes by, a Gaussian point first appears once t is sufficiently close to its temporal position μ_t and starts to grow. It reaches the peak opacity when $t = \mu_t$
- By controlling the temporal position and scaling factor, 4D Gaussian can represent challenging dynamics, e.g., motions that suddenly appear or disappear

Let's Take a Closer Look

- A temporally sliced 4DGS is a 3DGS

$$G_{3D}(\mathbf{x}, t) = e^{-\frac{1}{2}\lambda(t-\mu_t)^2} e^{-\frac{1}{2}[\mathbf{x}-\boldsymbol{\mu}(t)]^T \boldsymbol{\Sigma}_{3D}^{-1} [\mathbf{x}-\boldsymbol{\mu}(t)]},$$

$$\boldsymbol{\mu}(t) = (\mu_x, \mu_y, \mu_z)^T + (t - \mu_t) \frac{\mathbf{V}}{W}.$$

- the sliced 3D Gaussian exhibits a new motion term of $(t - \mu_t) \frac{V}{W}$ added to the center position $(\mu_x, \mu_y, \mu_z)^T$
- $\frac{V}{W}$ indicates the motion speed in the current timestamp
- By modeling a scene with 4D-Rotator GS, we can acquire speed field for free

Emerging Optical Flow from 4D Rendering



w/o 4D Consistency L. w/ 4D Consistency L.

Ground Truth

4D-Rotator Gaussian Splatting is Difficult to Learn

- Solution:

- 4D Consistency Loss: Neighboring Gaussian points should have similar speed s

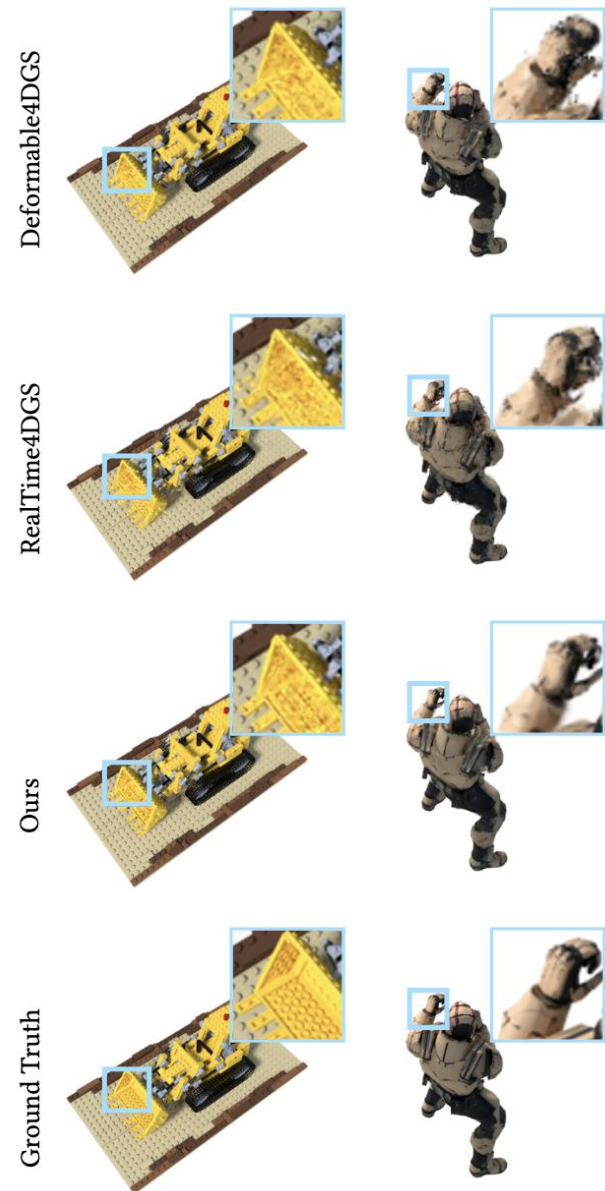
$$L_{\text{consistent4D}} = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{s}_i - \frac{1}{K} \sum_{j \in \Omega_i} \mathbf{s}_j \right\|_1.$$

- Entropy Loss: Remove non-surface Gaussian points by encouraging the opacity to be close to 0 or 1

$$L_{\text{entropy}} = \frac{1}{N} \sum_{i=1}^N -o_i \log(o_i).$$

Results

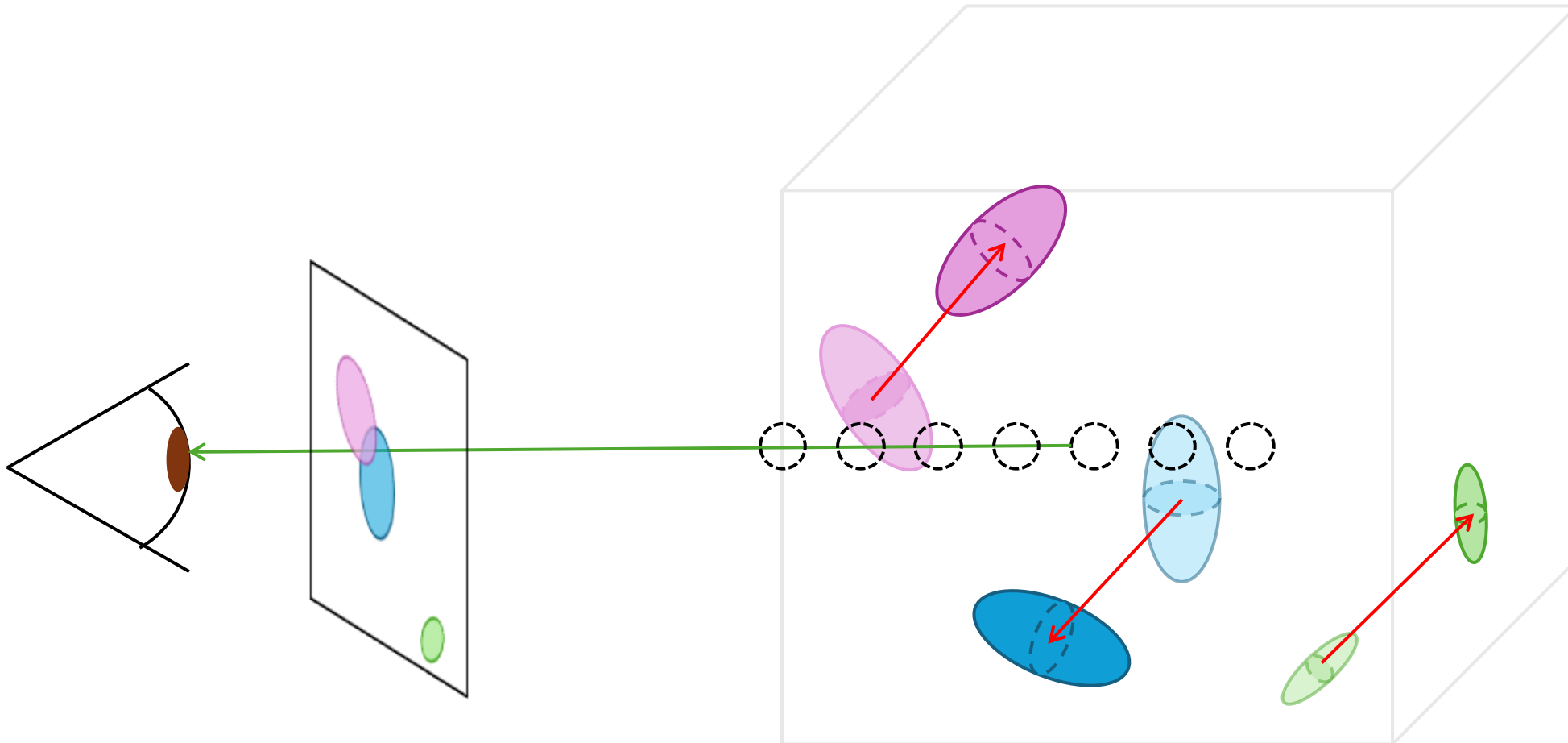
ID	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Train \downarrow	FPS \uparrow
<i>a</i>	DyNeRF [Li et al. 2022b]* \dagger	29.58	-	0.08	1344 h**	0.015
<i>b</i>	StreamRF [Li et al. 2022a]	28.16	0.85	0.31	79 min	8.50
<i>c</i>	HyperReel [Attal et al. 2023]	30.36	0.92	0.17	9 h	2.00
<i>d</i>	NeRFPlayer [Song et al. 2023a] \dagger	30.69	-	0.11	6 h	0.05
<i>e</i>	K-Planes [Fridovich-Keil et al. 2023]	30.73	0.93	0.07	190 min	0.10
<i>f</i>	MixVoxels [Wang et al. 2023b]	30.85	0.96	0.21	91 min	16.70
<i>g</i>	Deformable4DGS [Wu et al. 2023]	28.42	0.92	0.17	72 min	39.93
<i>h</i>	RealTime4DGS [Yang et al. 2024]	29.95	0.92	0.16	8 h	72.80
<i>i</i>	Ours	31.62	0.94	0.14	60 min	277.47



Results



Method 2: 3D GS and their Motions



Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis

- **Idea:** Represent a dynamic scene with the same set of 3D Gaussians and their motions
- **Particle-based dynamic perspective:** Each Gaussian can be thought of as softly representing an area of 3D physical space which is occupied by solid matter



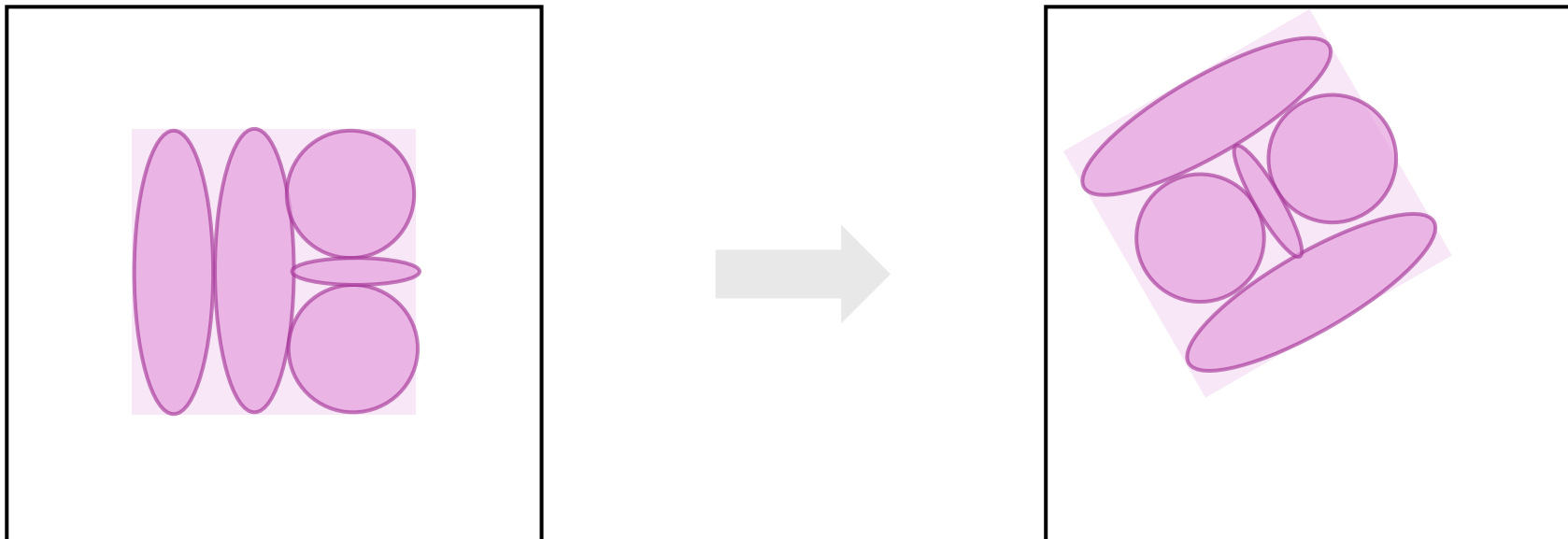
Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis

- Dynamic 3D GS:
 - A 3D center for each timestep (x_t, y_t, z_t)
 - A 3D rotation for each timestep (qw_t, qx_t, qy_t, qz_t)
 - A fixed 3D size across time (s_x, s_y, s_z)
 - A fixed color (r, g, b)
 - A fixed opacity o
 - A background indicator bg
- $7t + 8$ parameters in total



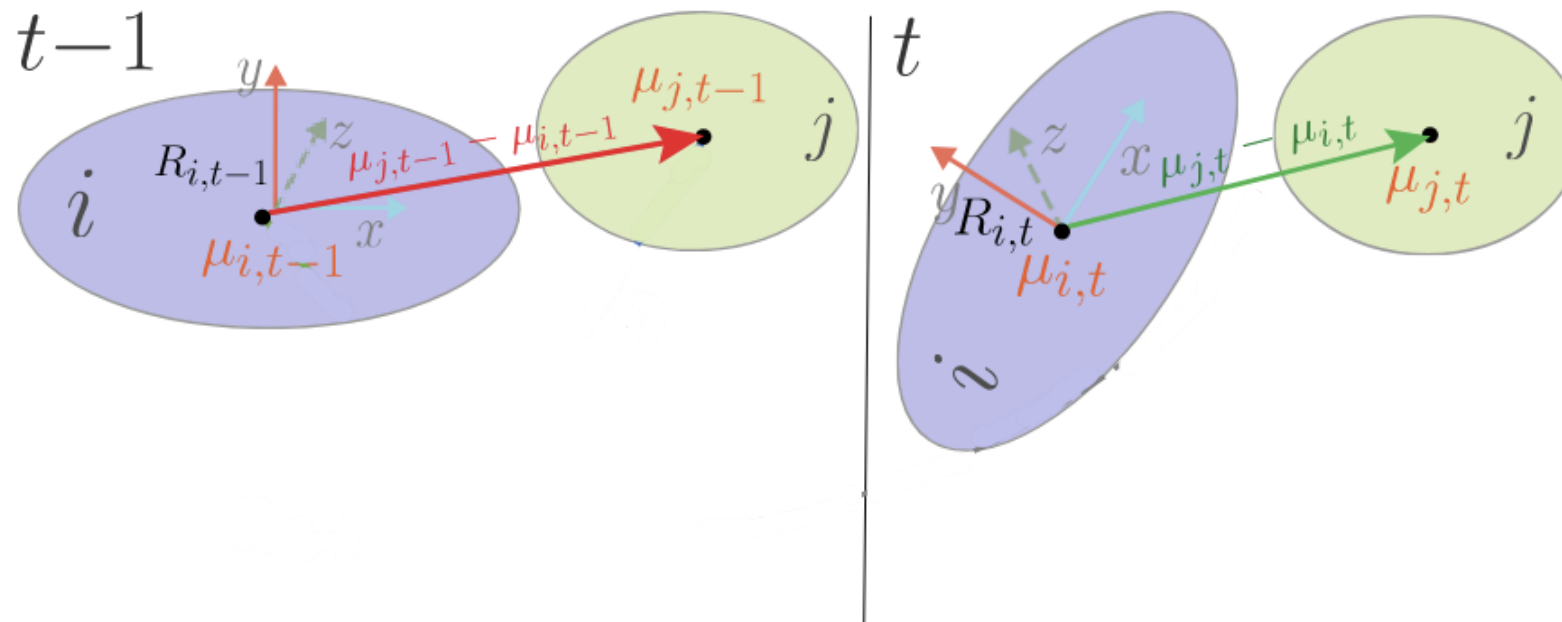
Dynamic 3D Gaussian Doesn't Work Out-of-Box

- **Problem:** Fixing color, opacity and size of Gaussians doesn't work on long videos, especially across areas of the scene where there is a large area of near uniform color. 3D Gaussians move freely in the area of uniform color.



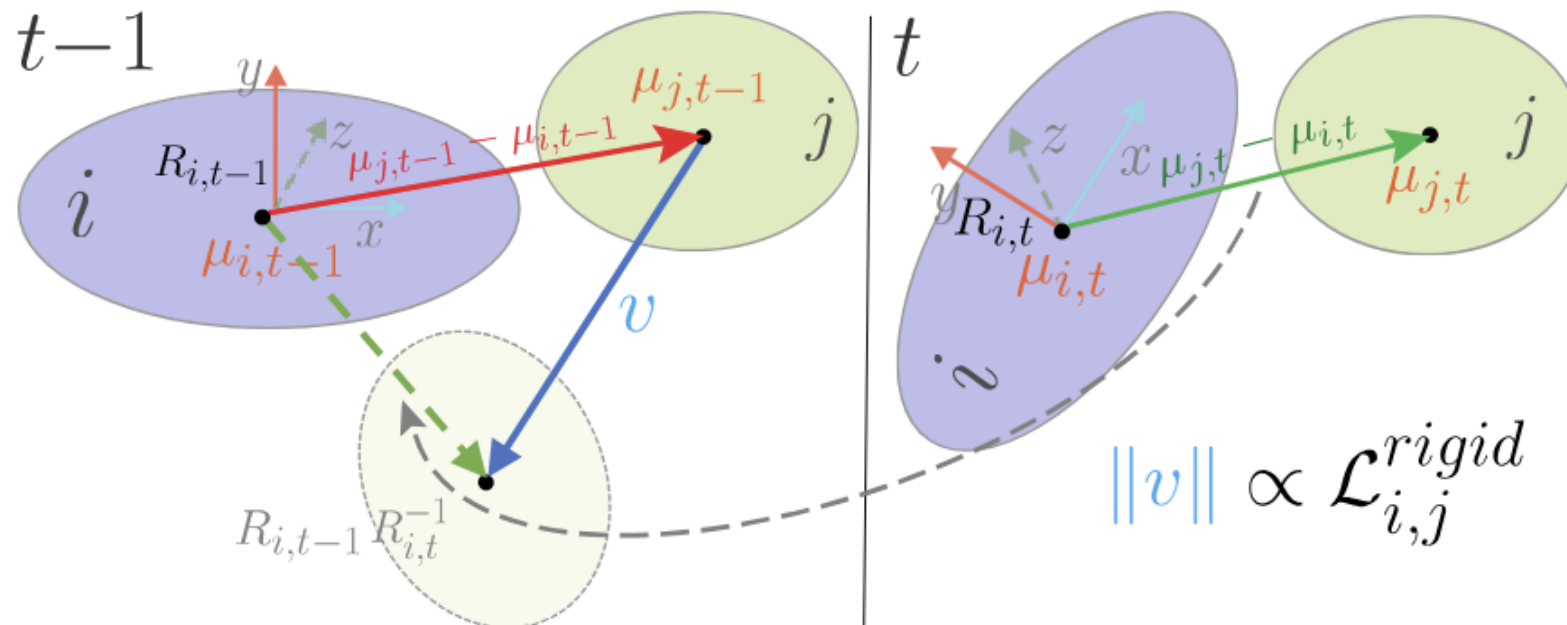
Dynamic 3D Gaussian Doesn't Work Out-of-Box

- **Solution:** Physically-Based Priors
 - Local-rigidity constraint: locally neighboring Gaussians should move like a rigid body



Dynamic 3D Gaussian Doesn't Work Out-of-Box

- **Solution:** Physically-Based Priors
 - Local-rigidity constraint: locally neighboring Gaussians should move like a rigid body



Dynamic 3D Gaussian Doesn't Work Out-of-Box

- **Solution:** Physically-Based Priors

- Local-rigidity constraint: locally neighboring Gaussians should move like a rigid body

$$\mathcal{L}_{i,j}^{\text{rigid}} = w_{i,j} \left\| (\mu_{j,t-1} - \mu_{i,t-1}) - R_{i,t-1} R_{i,t}^{-1} (\mu_{j,t} - \mu_{i,t}) \right\|_2$$

$$\mathcal{L}^{\text{rigid}} = \frac{1}{k|\mathcal{S}|} \sum_{i \in \mathcal{S}} \sum_{j \in \text{knn}_{i,k}} \mathcal{L}_{i,j}^{\text{rigid}}$$

- This loss predicted $\mu_{i,t}$ and $R_{i,t}$ to be geometrically consistent, facilitating more precise tracking

Learning Dynamic 3D Gaussian is Slow

- Solution:
 - Local-orientation loss: Enforce neighboring Gaussians to have the same rotation

$$\mathcal{L}^{\text{rot}} = \frac{1}{k|\mathcal{S}|} \sum_{i \in \mathcal{S}} \sum_{j \in \text{knn}_{i;k}} w_{i,j} \left\| \hat{q}_{j,t} \hat{q}_{j,t-1}^{-1} - \hat{q}_{i,t} \hat{q}_{i,t-1}^{-1} \right\|_2 \quad w_{i,j} = \exp \left(-\lambda_w \left\| \mu_{j,0} - \mu_{i,0} \right\|_2^2 \right)$$

- Local-translation loss: Enforce neighboring Gaussians to have the same translation

$$\mathcal{L}^{\text{iso}} = \frac{1}{k|\mathcal{S}|} \sum_{i \in \mathcal{S}} \sum_{j \in \text{knn}_{i;k}} w_{i,j} \left| \left\| \mu_{j,0} - \mu_{i,0} \right\|_2 - \left\| \mu_{j,t} - \mu_{i,t} \right\|_2 \right|$$

Results

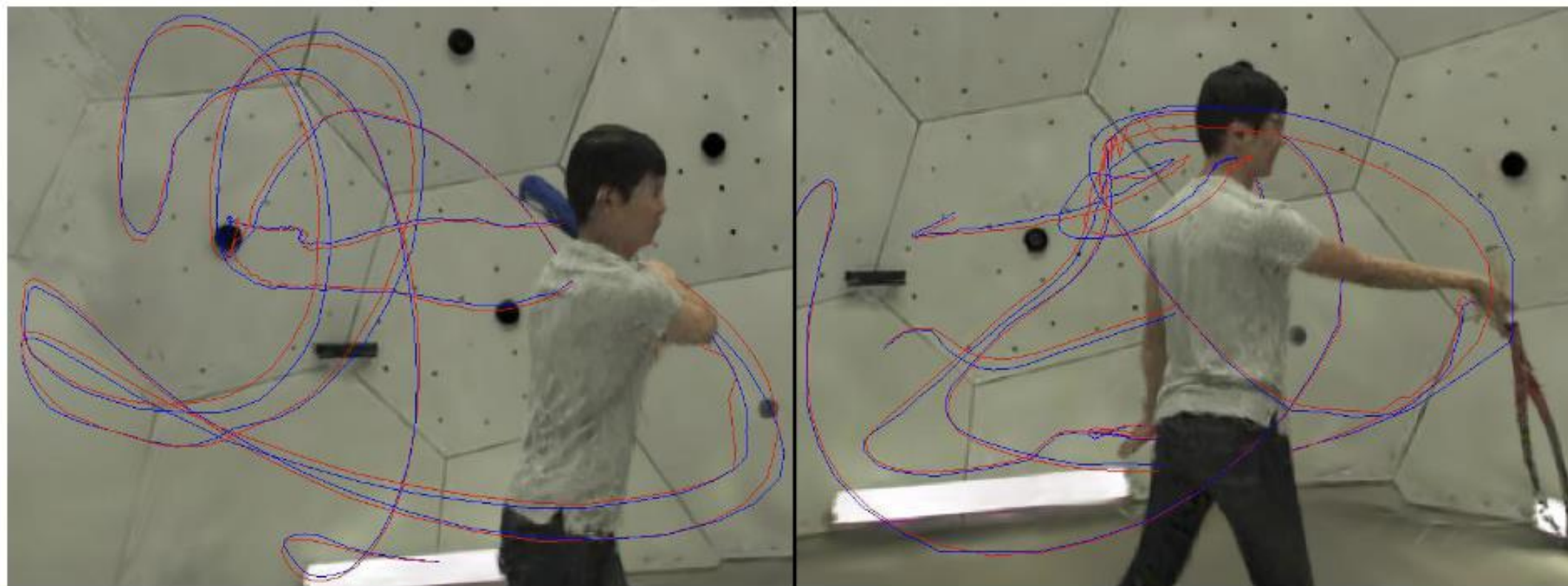


Results



Emerging Tracking from Image Rendering

Red curve: ground-truth tracks; Blue curve: predicted tracks



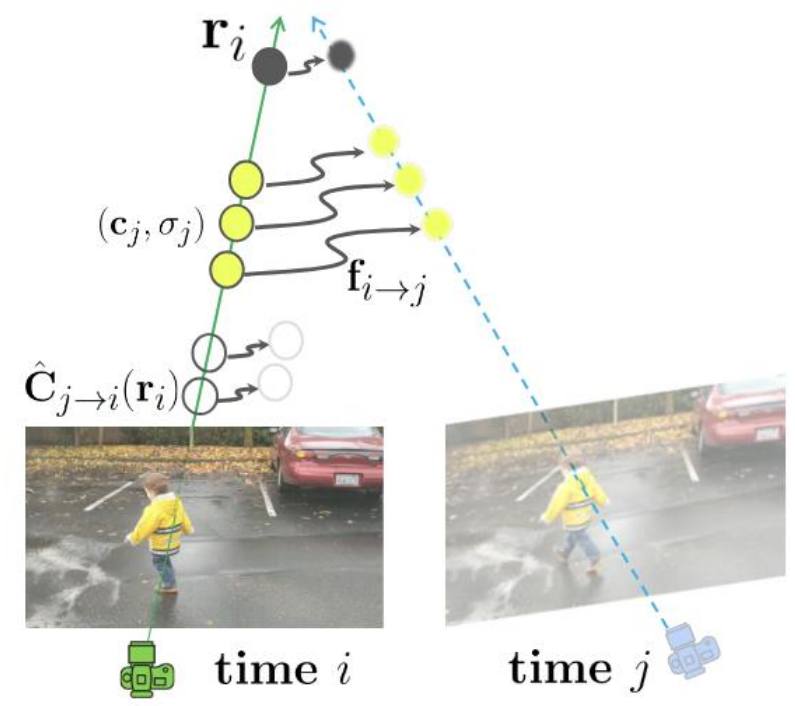
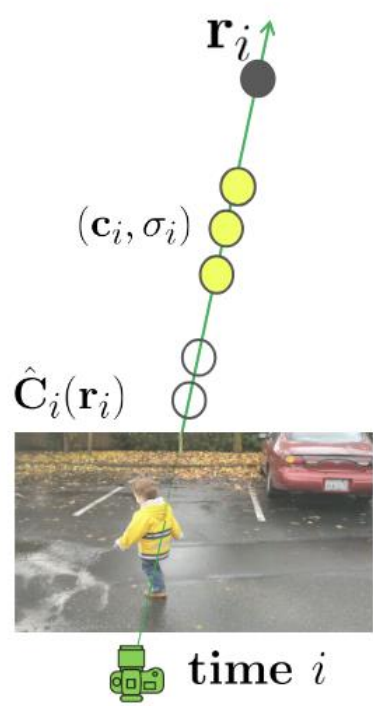
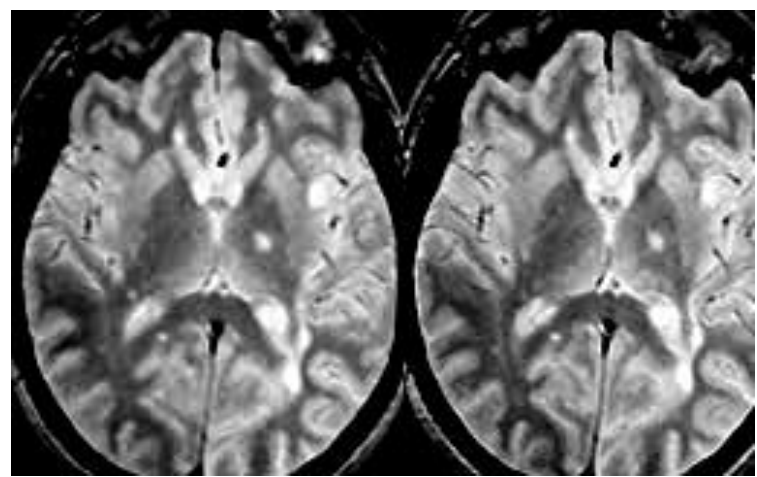
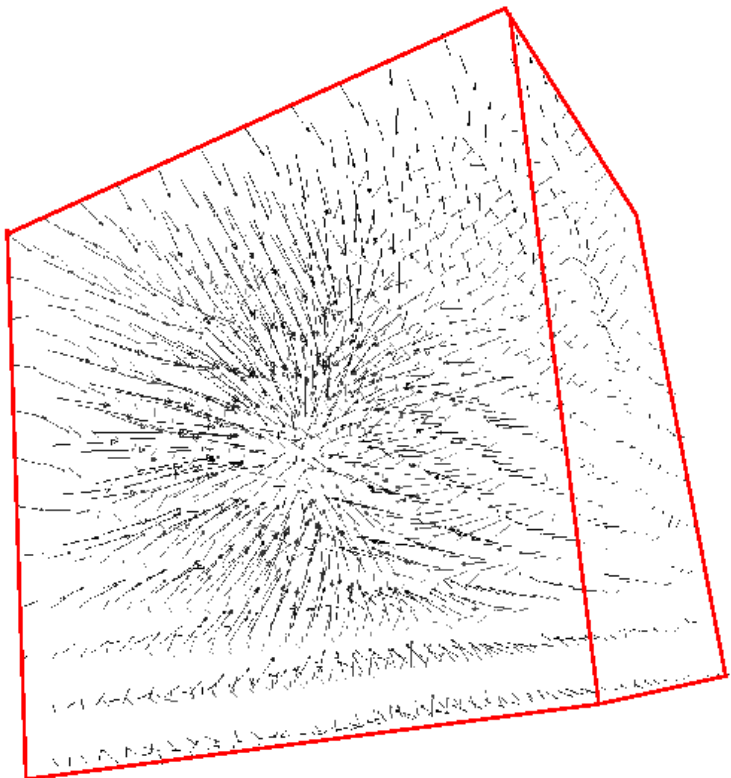
Applications: Augmented Reality



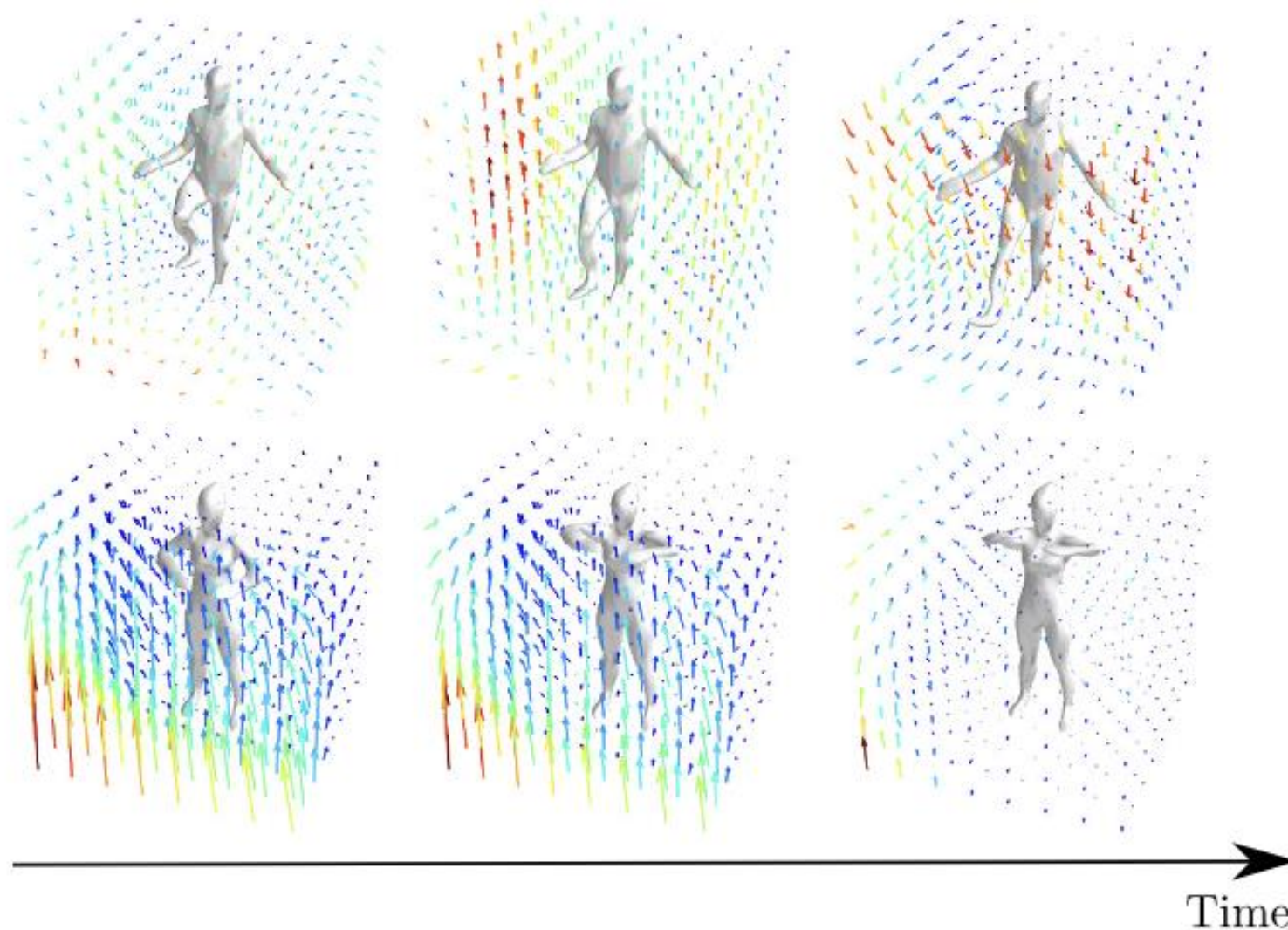
Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

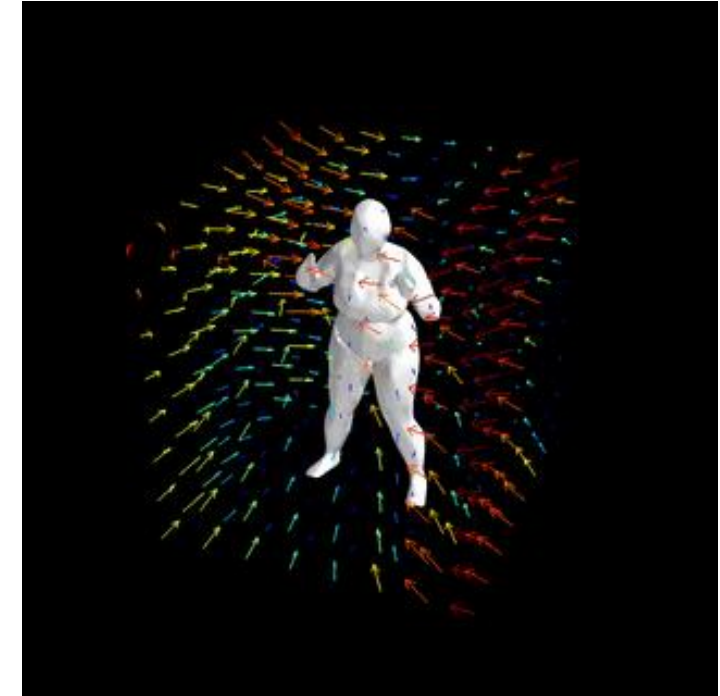
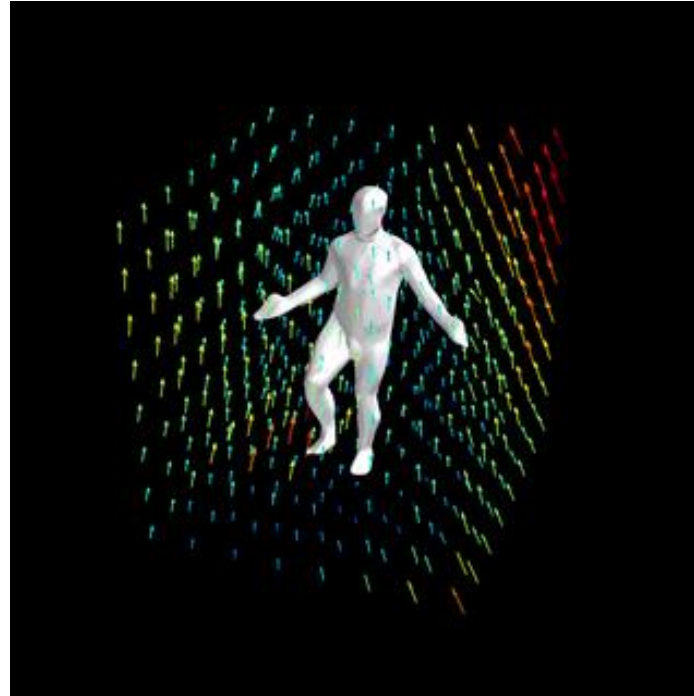
Model a Vector Field



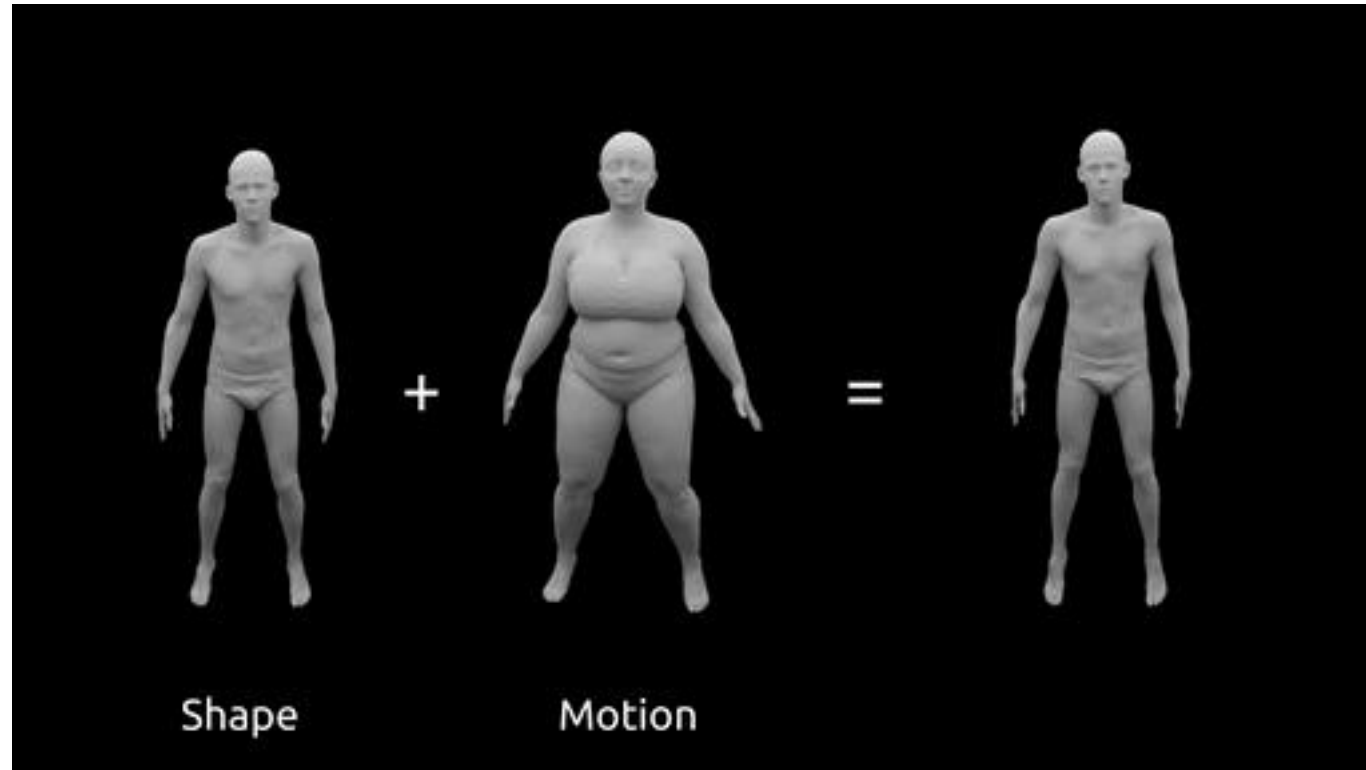
Decoupling Shape and Motion Representations



Modeling High-Resolution Motions with a Vector Field

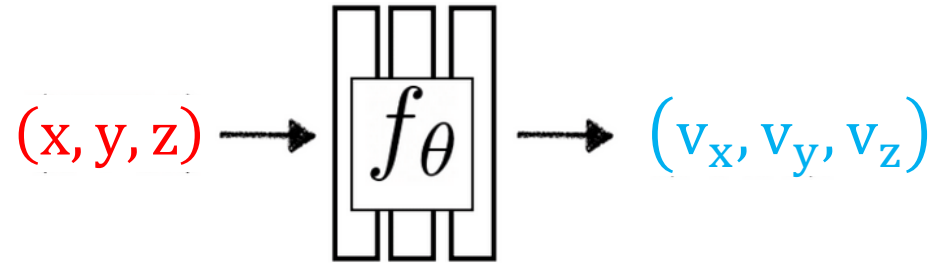


Transfer Motions to Other Instances



How to Model a Dense Motion Field?

Prediction of a neural network



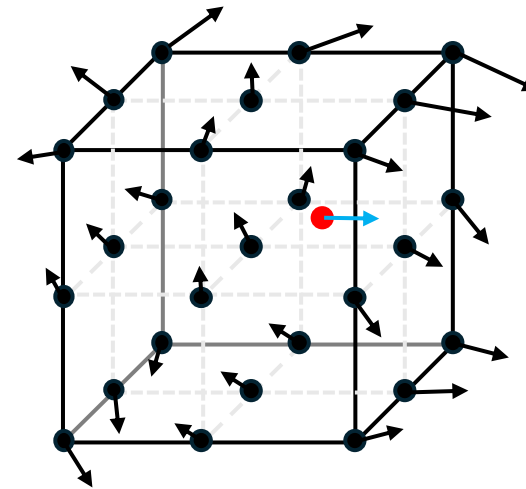
Pros:

- Capture high-frequency motions

Cons:

- Difficult to train

Interpolation of a grid



$$v_{xyz} = \sum_{i \in \mathcal{N}_{xyz}} w_i v_i$$

Pros:

- Share motions across particles (low rank approximation of motions)

Cons:

- Incapable of capturing high-frequency motions

Common RGB Videos Contain only Low-frequency Motions

Low-frequency motions

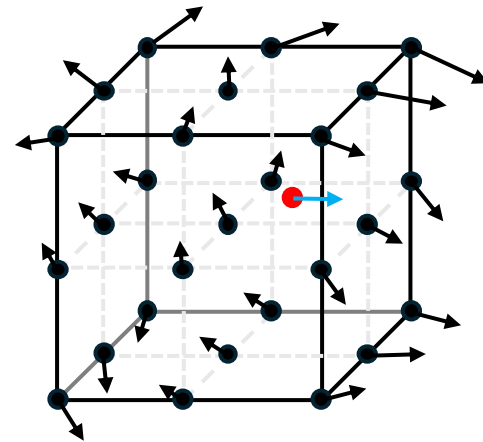


High-frequency motions



Interpolating motions within a grid is a good choice on common RGB videos

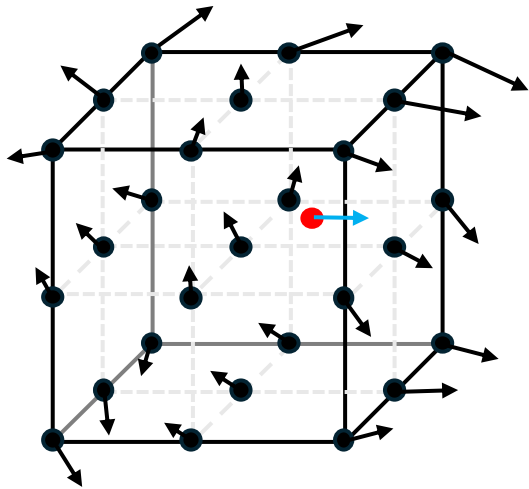
Goal: Infer 3D GS and the Motion Field from RGB Videos for 4D Reconstruction



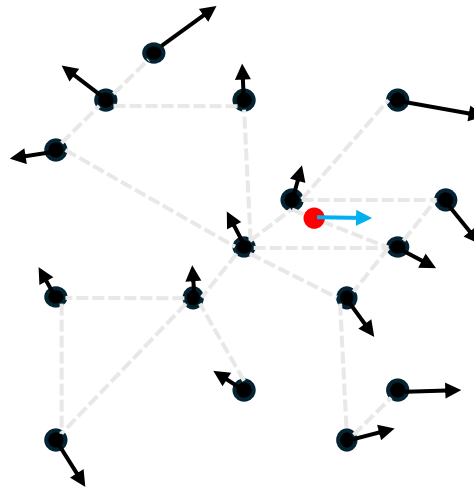
Video source: Shape of Motion: 4D Reconstruction from a Single Video

We Can Extend Grid-based Interpolation to Graph-based Interpolation

Interpolation of a grid



Interpolation of a graph



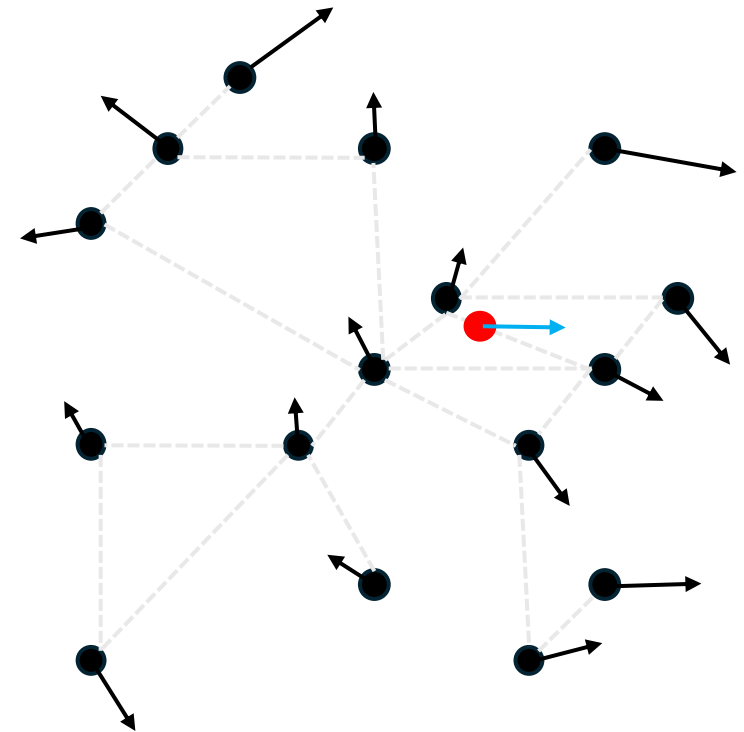
$$V_{xyz} = \sum_{i \in \mathcal{N}_{xyz}} w_i v_i$$

Graph-based Interpolation

- **Rigid transformation of 3D GS:** Transform a 3D GS with parameters (μ, Σ) using a rigid transformation $Q = [R; t]$

$$\mu' = R\mu + t; \quad \Sigma' = R\Sigma$$

- **Motivation:** Real-world motion typically behaves rigidly, smoothly, and compactly, meaning that the actual solution is low-rank and driven by a few key “eigen” motions.
- **4D Motion Scaffold:** Encodes these local “eigen” motions and interpolates the dense deformation field with a structured graph $(\mathcal{V}, \mathcal{E})$



Graph-based Interpolation

- Graph nodes $\mathcal{V} = \{v^{(m)}\}$: each node is a 6-DoF rigid motion trajectory

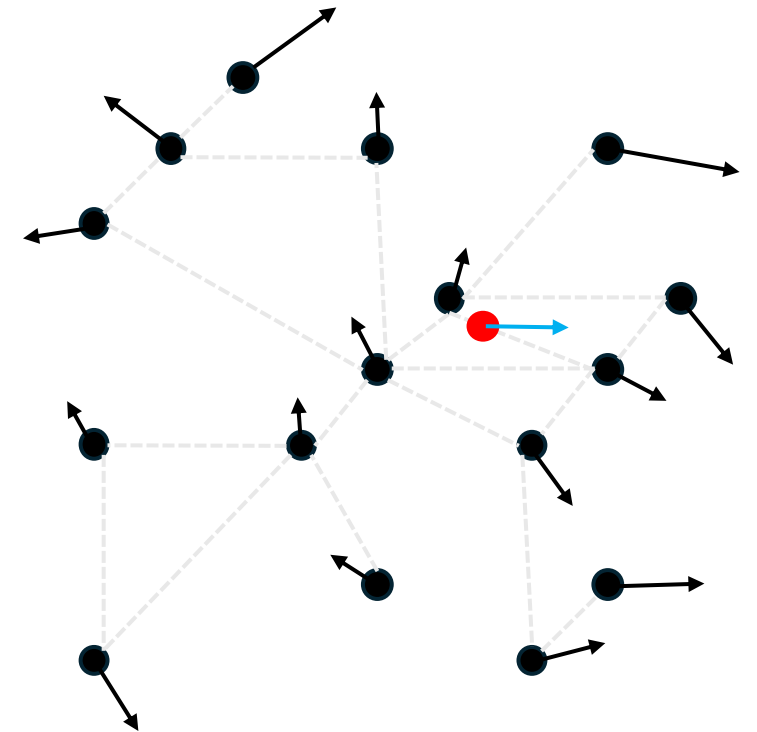
$$v^{(m)} = \left(\left[Q_1^{(m)}, Q_2^{(m)}, \dots, Q_T^{(m)} \right], r^{(m)} \right)$$

where $Q = [R; \mathbf{t}]$ and $r^{(m)}$ denotes control radius

- Graph edge:

$$\mathcal{E}(m) = \text{KNN}_{n \in \{1, \dots, M\}} [D_{\text{curve}}(m, n)],$$

$$D_{\text{curve}}(m, n) = \max_{t=1, 2, \dots, T} \|\mathbf{t}_t^{(m)} - \mathbf{t}_t^{(n)}\|,$$



Graph-based Interpolation

- Graph nodes $\mathcal{V} = \{v^{(m)}\}$: each node is a 6-DoF rigid motion trajectory

$$v^{(m)} = \left(\left[Q_1^{(m)}, Q_2^{(m)}, \dots, Q_T^{(m)} \right], r^{(m)} \right)$$

where $Q = [R; t]$ and $r^{(m)}$ denotes control radius

- Graph edge:

$$\mathcal{E}(m) = \text{KNN}_{n \in \{1, \dots, M\}} [D_{\text{curve}}(m, n)],$$

$$D_{\text{curve}}(m, n) = \max_{t=1, 2, \dots, T} \|\mathbf{t}_t^{(m)} - \mathbf{t}_t^{(n)}\|,$$

- Deformation field \mathcal{W} :

$$\mathcal{W}(\mathbf{x}, \mathbf{w}; t_{\text{src}}, t_{\text{dst}}) = \text{DQB} \left(\{w_i, \Delta \mathbf{Q}^{(i)}\}_{i \in \mathcal{E}(m^*)} \right)$$

$$\text{DQB}(\{(w_i, \mathbf{Q}_i)\}_{i=1}^L) = \frac{\sum_{i=1}^L w_i \hat{\mathbf{q}}_i}{\|\sum_{i=1}^L w_i \hat{\mathbf{q}}_i\|_{DQ}} \in SE(3),$$

$$\Delta \mathbf{Q}^{(i)} = \mathbf{Q}_{t_{\text{dst}}}^{(i)} (\mathbf{Q}_{t_{\text{src}}}^{(i)})^{-1}$$

$$w_i(\mathbf{x}, t_{\text{src}}) = \exp(-\|\mathbf{x} - \mathbf{t}_{t_{\text{src}}}^{(i)}\|_2^2 / 2r^{(i)}) \in \mathbb{R}^+.$$

MoSca: Dynamic Gaussian Fusion from Casual Videos via 4D Motion Scaffolds

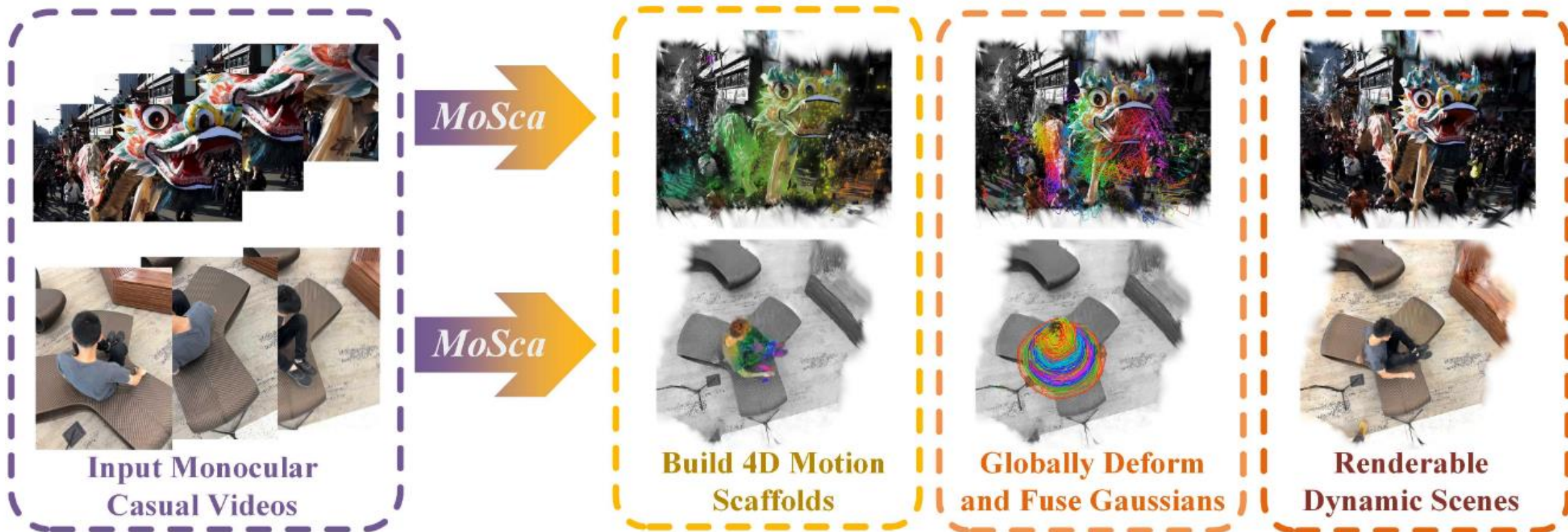
Jiahui Lei¹ Yijia Weng² Adam W. Harley² Leonidas Guibas² Kostas Daniilidis^{1,3}

¹ University of Pennsylvania

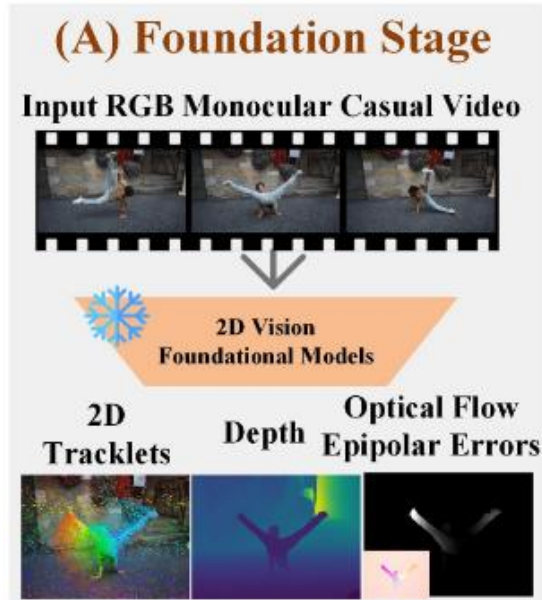
² Stanford University

³ Archimedes, Athena RC

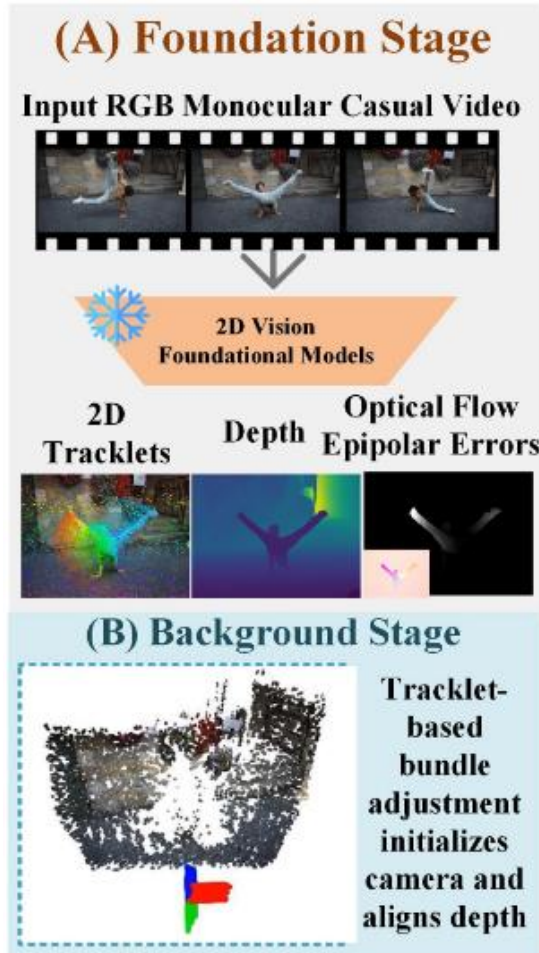
{leijh, kostas}@cis.upenn.edu, {yijiaw, aharley, guibas}@cs.stanford.edu



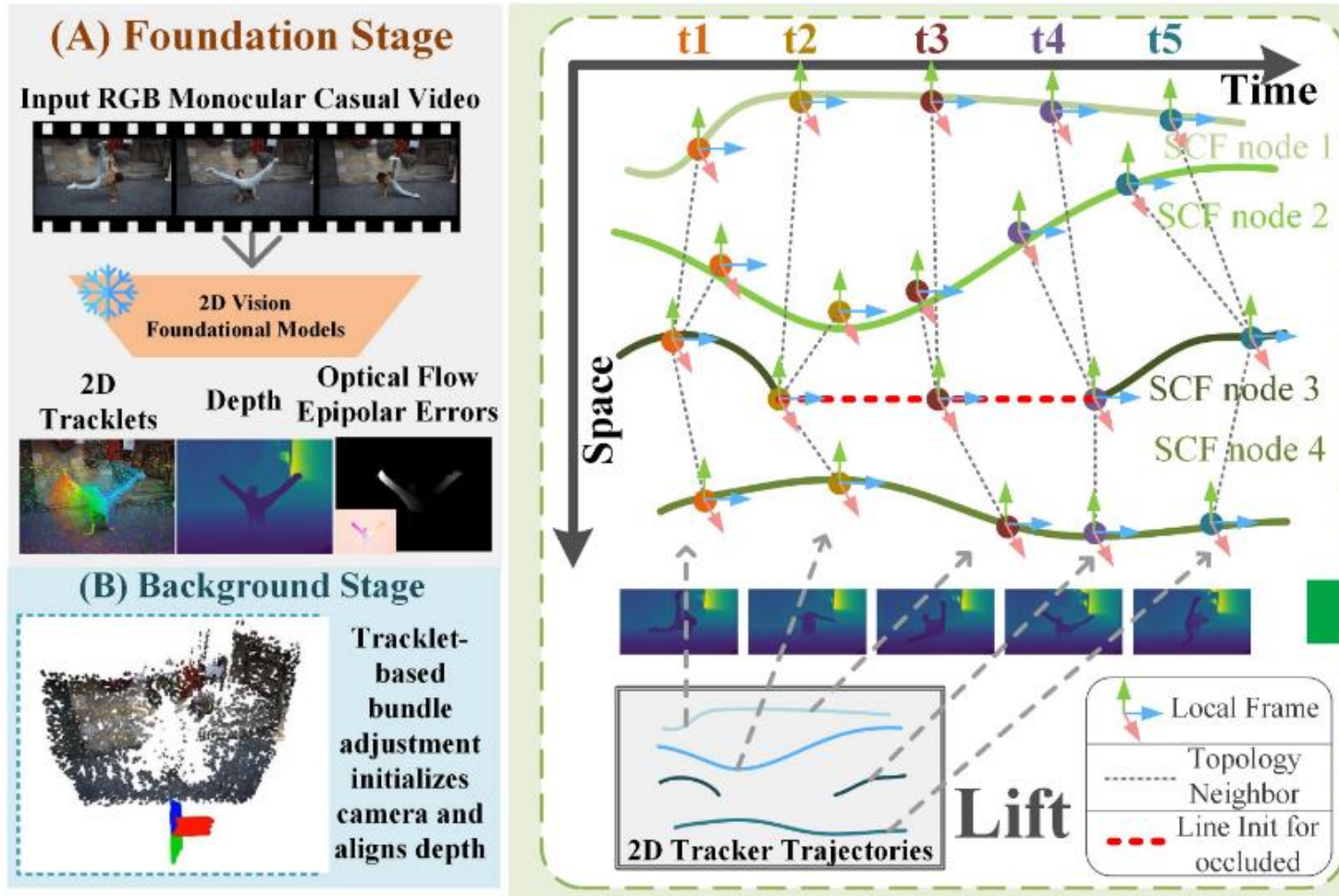
Stage 1: Extract 2D Tracklets, Depth and Optical Flow Epipolar Errors from 2D Foundation Models



Stage 2: Estimate Camera Motions by Separating Foreground and Background



Stage 3: Initiate Graph Nodes \mathcal{V} by Lifting 2D Tracklets to 3D



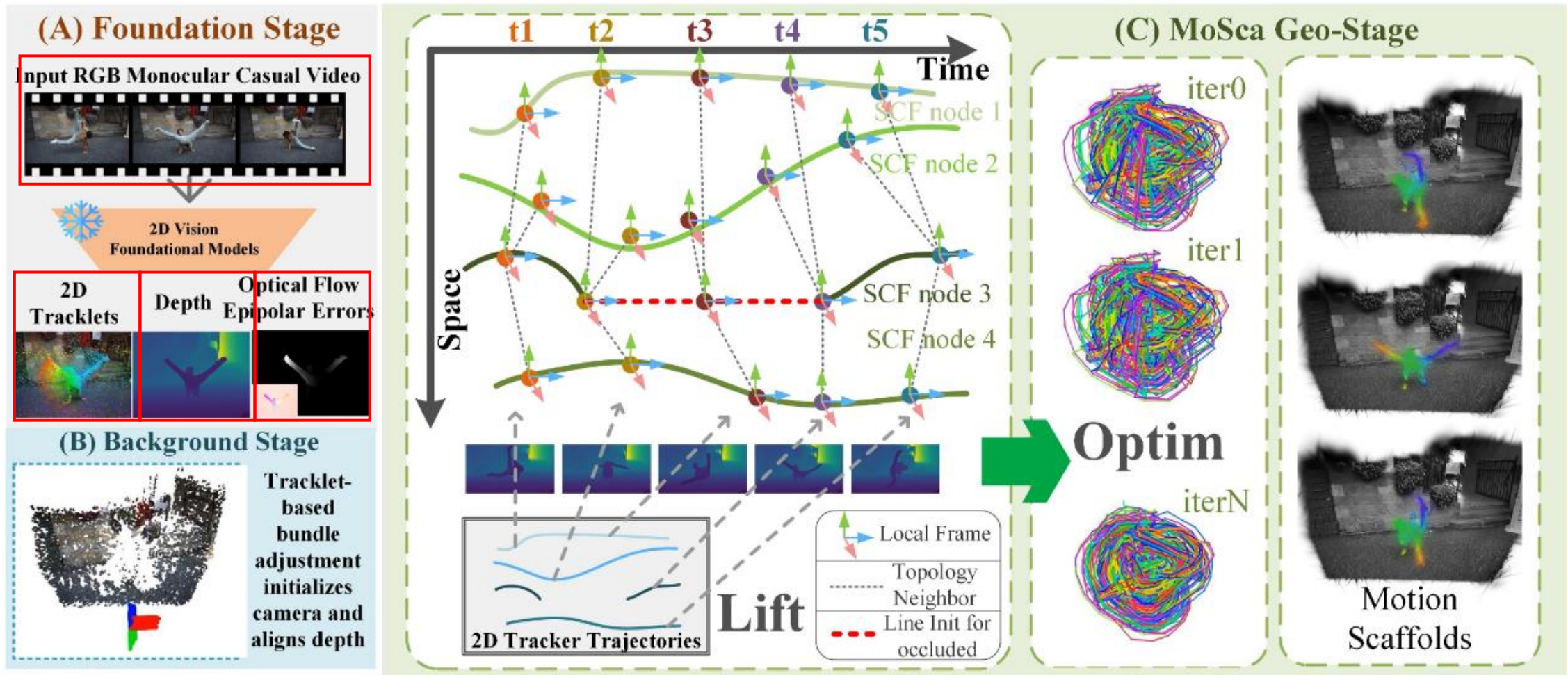
- Graph node $\mathbf{v}^{(m)} = \left(\left[Q_1^{(m)}, Q_2^{(m)}, \dots, Q_T^{(m)} \right], r^{(m)} \right)$

$$Q_t^{(i)} = \left[I, \mathbf{h}_t^{(i)} \right]$$

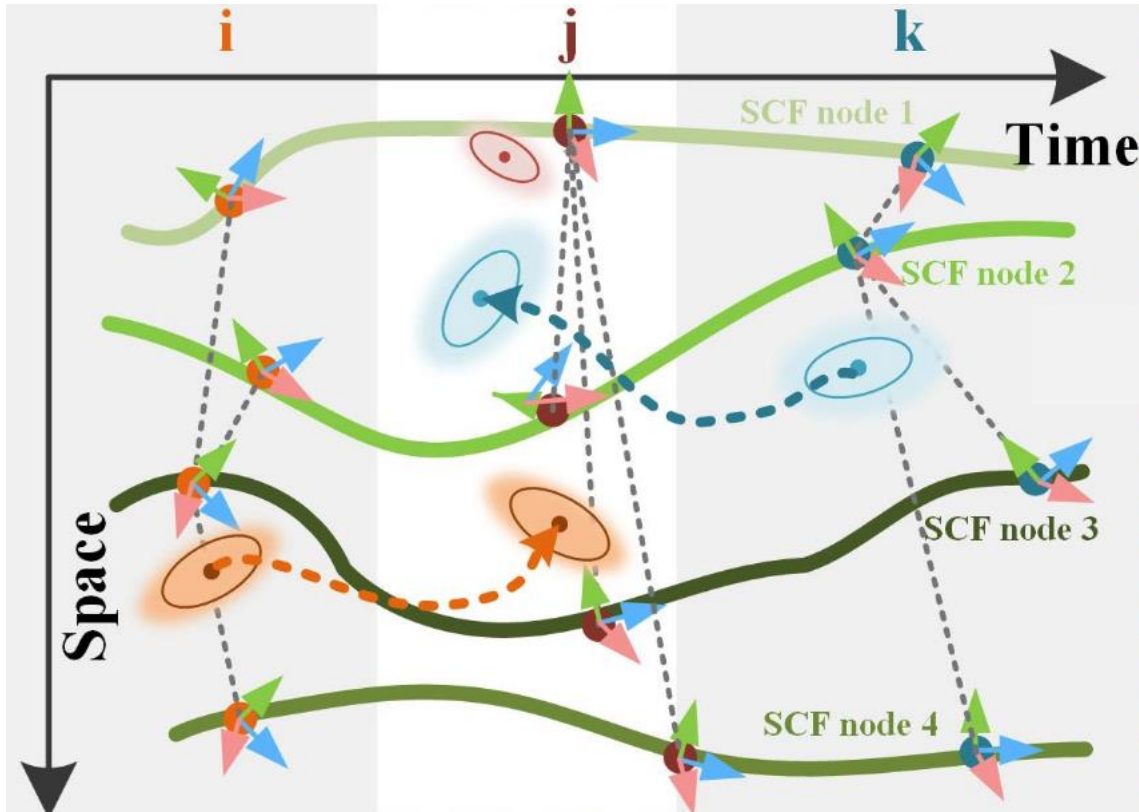
And $\mathbf{h}_t^{(i)}$ denotes the unprojected depth of 2D tracklet i at time t

Stage 4: Optimize the 4D Motion Scaffolds

Supervise with RGB, depth, optical flow, tracking, smoothness and rigid-body losses.



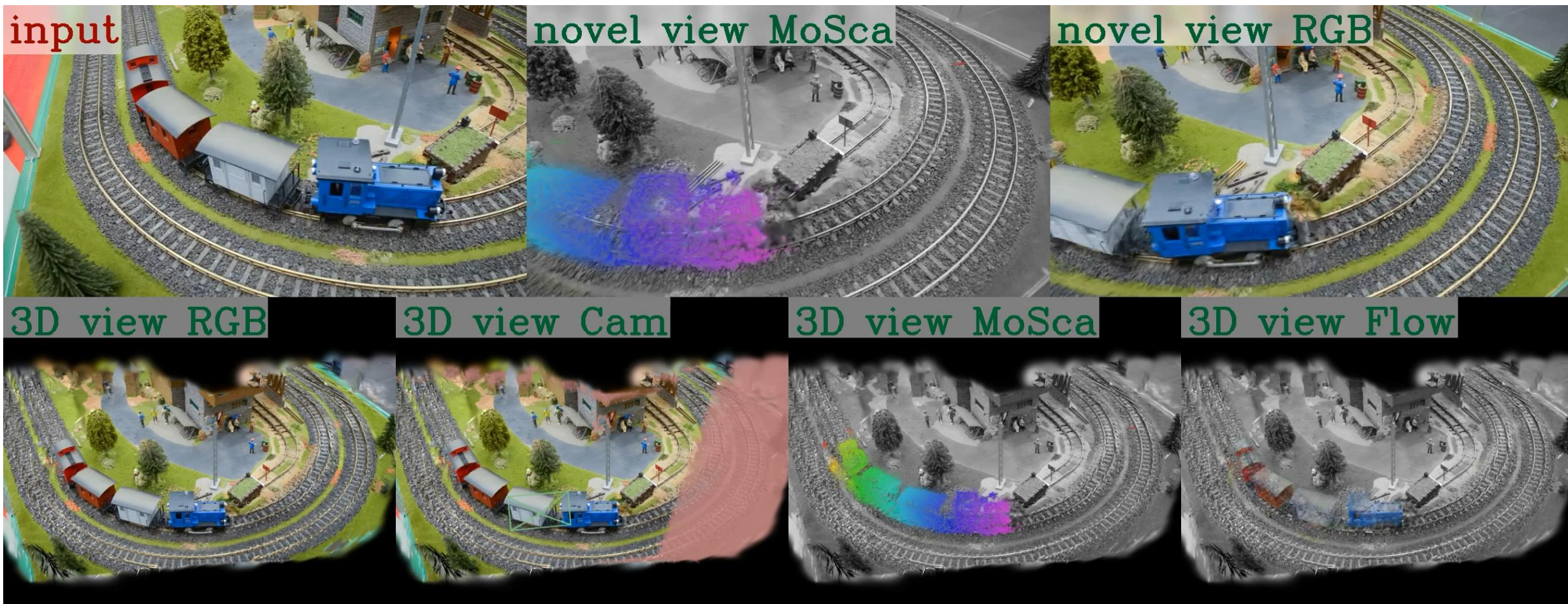
Why Motion Scaffolds are Good Motion Representations?



- The global deformation field \mathbf{w} can transform points at any time globally, enabling the fusion of all observed video frames into a single coherent representation.
- Let t_j^{ref} denote the time step when the j -th GS is initiated and $\Delta\mathbf{w}_j$ be the per-Gaussian learnable skinning weight correction:

$$\mathcal{G}(t) = \{(\mathbf{T}_j(t)\mu_j, \mathbf{T}_j(t)R_j, s_j, o_j, c_j) \mid \mathbf{T}_j(t) = \mathcal{W}(\mu_j, \mathbf{w}(\mu_j, t_j^{\text{ref}}) + \Delta\mathbf{w}_j; t_j^{\text{ref}}, t)\}_{j=1}^N$$

Results



Results



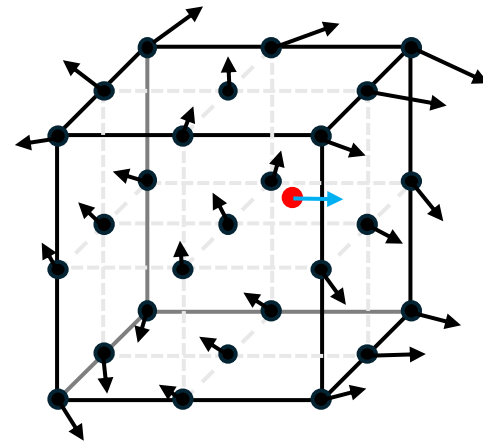
Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

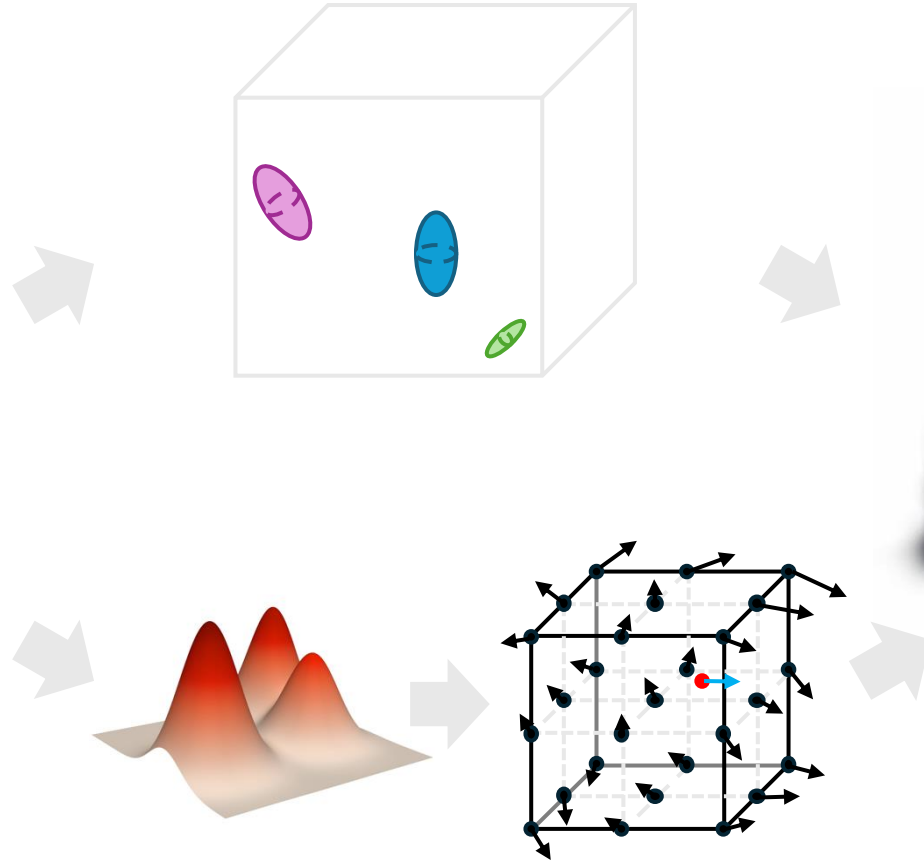
The Motion Field of the Input Video May Be High-Dimensional, which is difficult to infer...



Video source: Shape of Motion: 4D Reconstruction from a Single Video

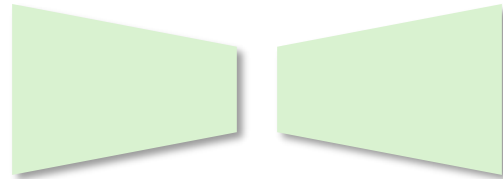
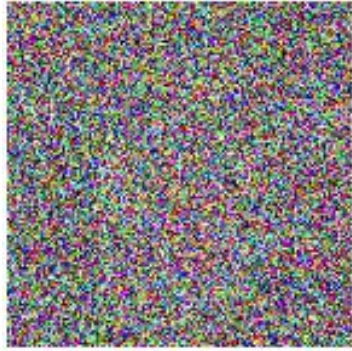


Idea: First Compress Motions into Compact Latent Features, then Decode Back to High-Dimensional Fields



Video source: Shape of Motion: 4D Reconstruction from a Single Video

We Have the Same Problem in Visual Generation



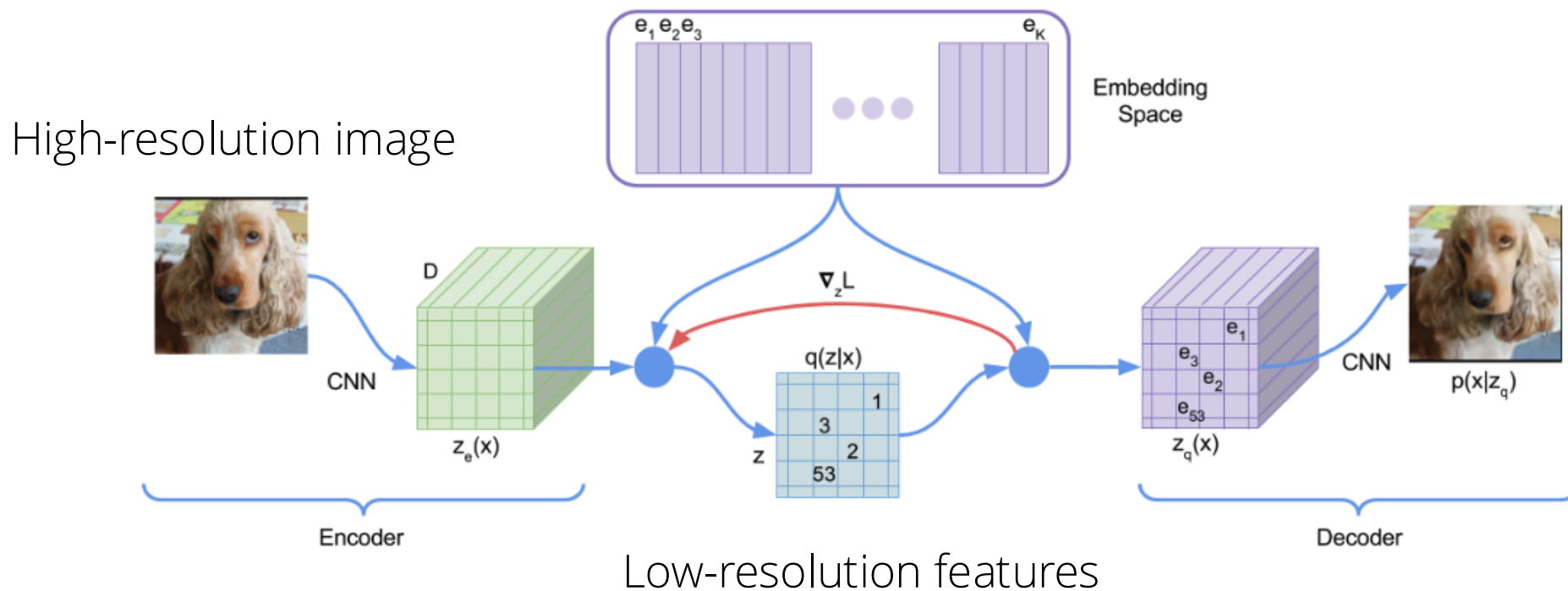
Generative model



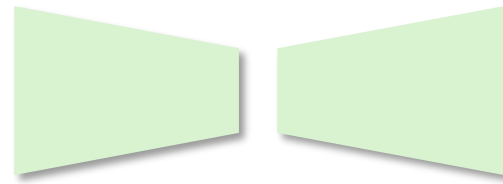
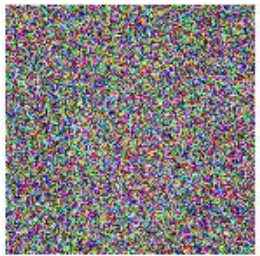
- Generating rgb pixels directly is super expensive!
 - No structural prior: a 16x16 red patch can be denoted by much simpler representations
 - High computational cost

Idea: Encode Images in a Compact Feature Space

- Learn to encode images by auto-encoding
 - information bottleneck + reconstruction loss

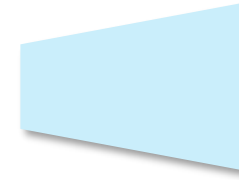
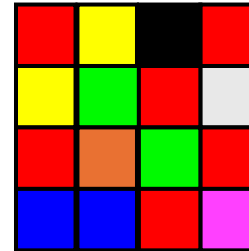


Idea: Generate Images in a Compact Feature Space



Generative model

Low-resolution features

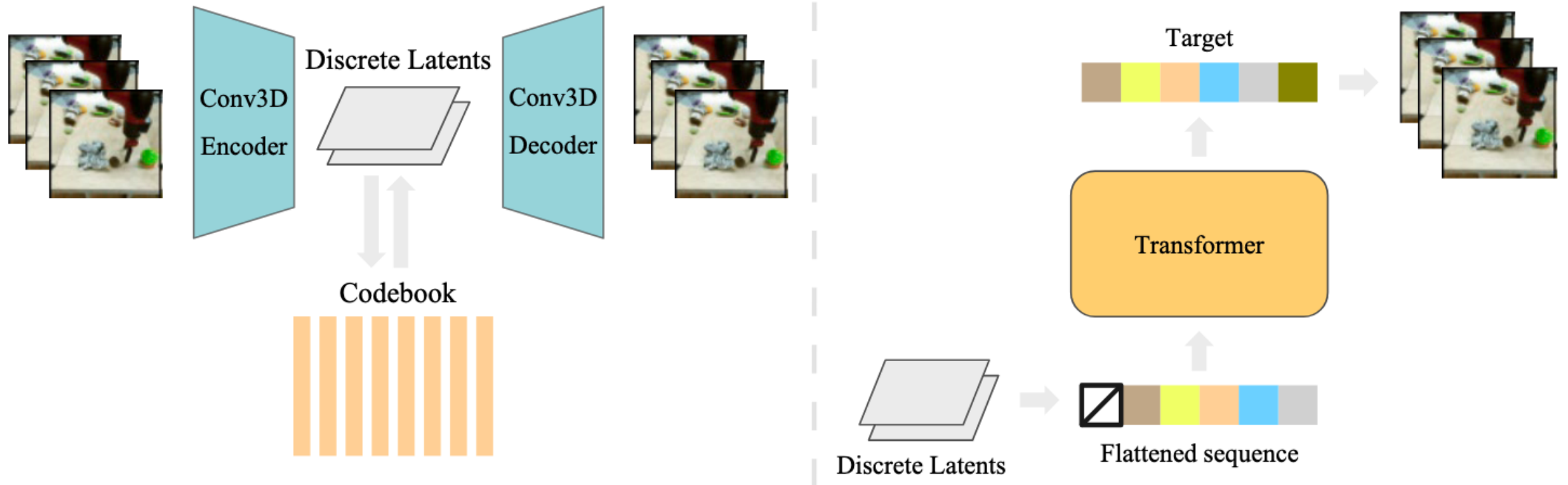


Decoder

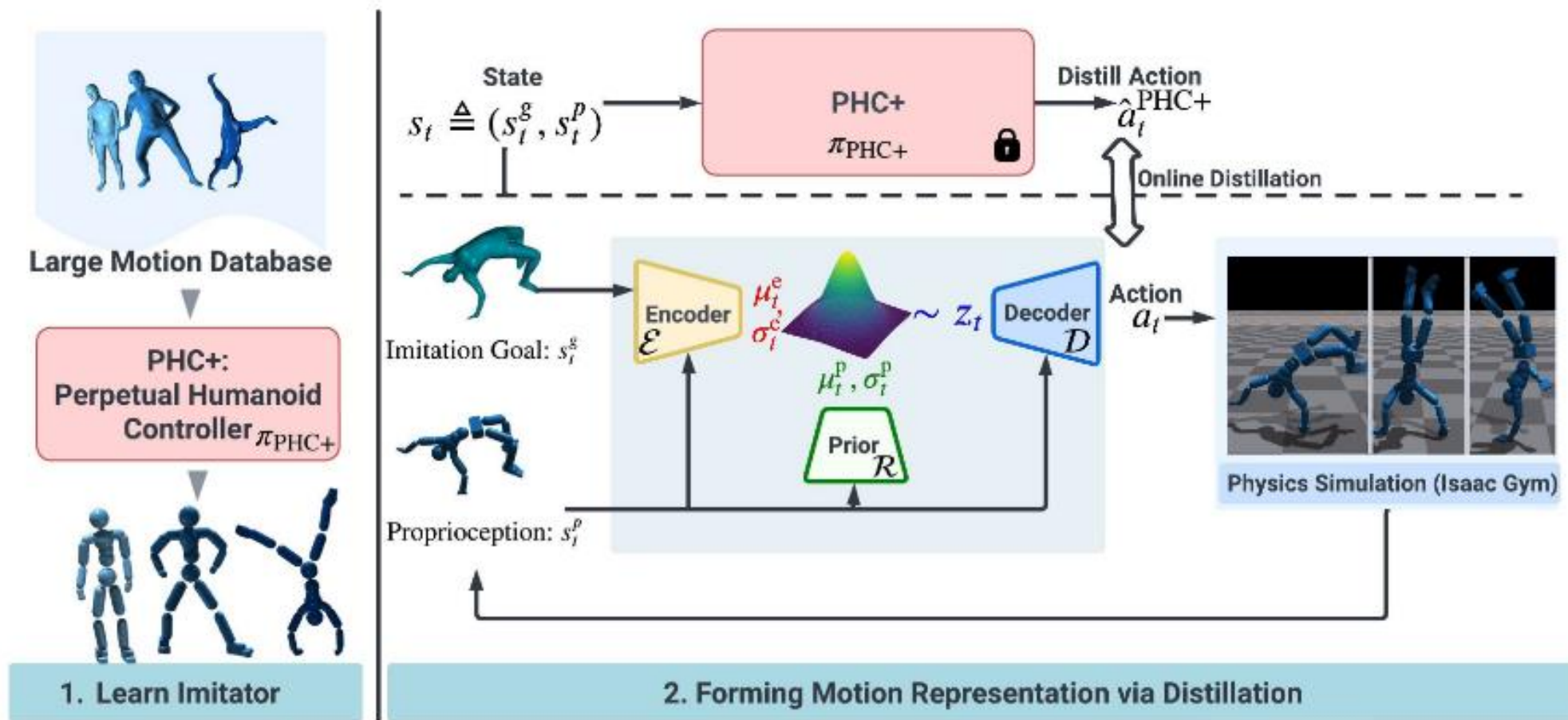
High-resolution image



Same Idea for Video Generation



Same Idea for Human Motion Generation

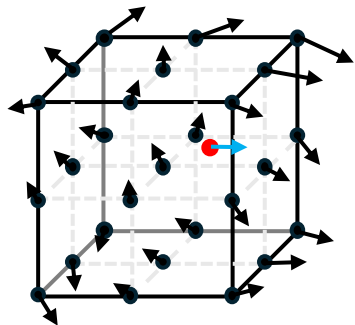


We Need An AutoEncoder that Compresses a Motion Field

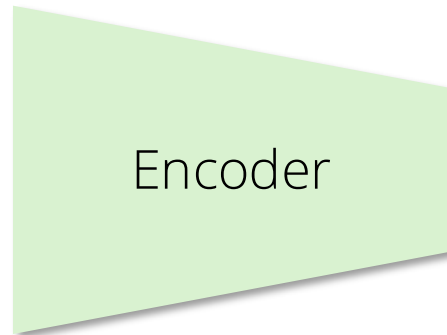
Static Object



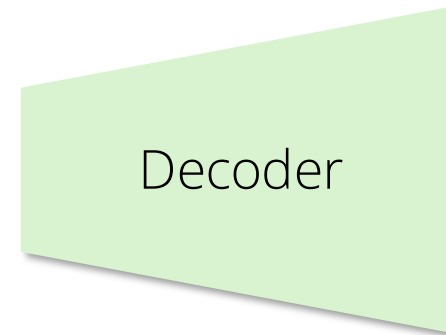
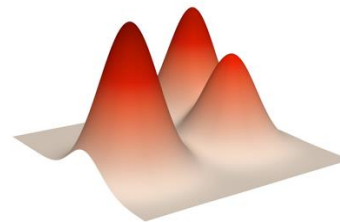
Animated Object



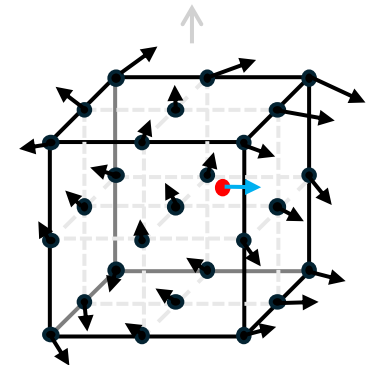
Motion Field



Encoder

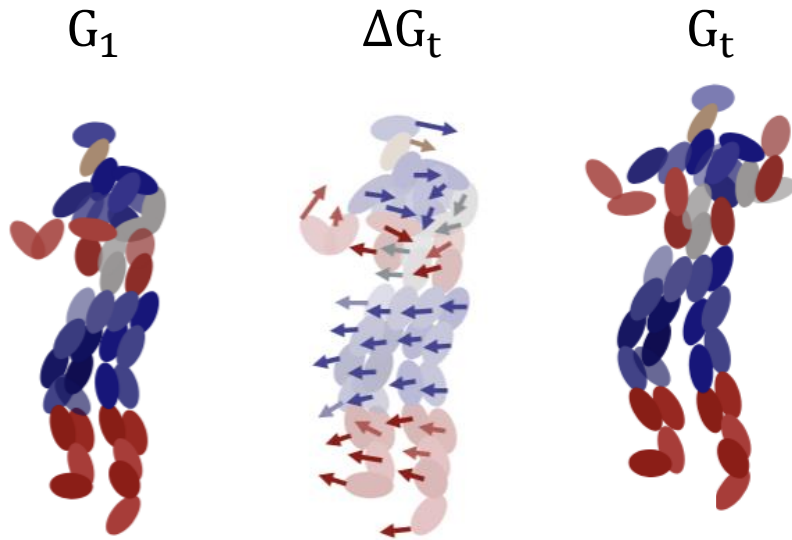


Decoder



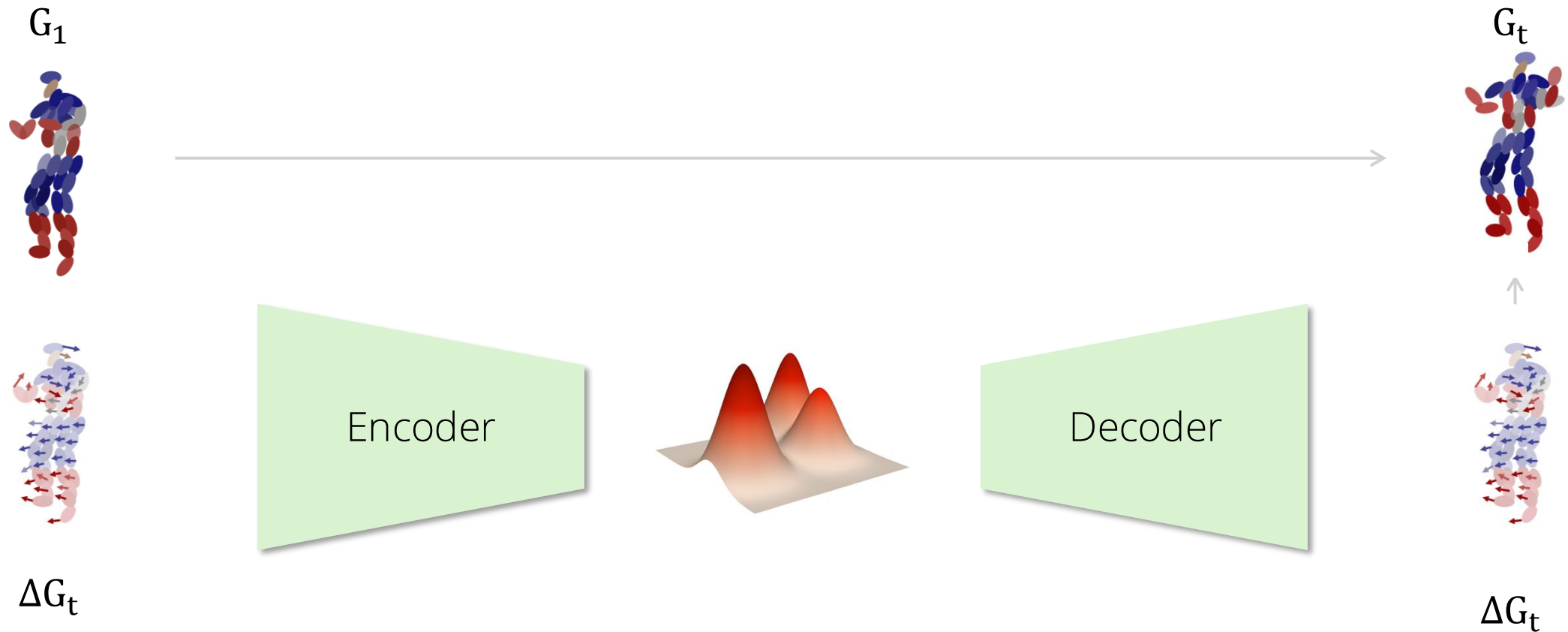
Motion Field

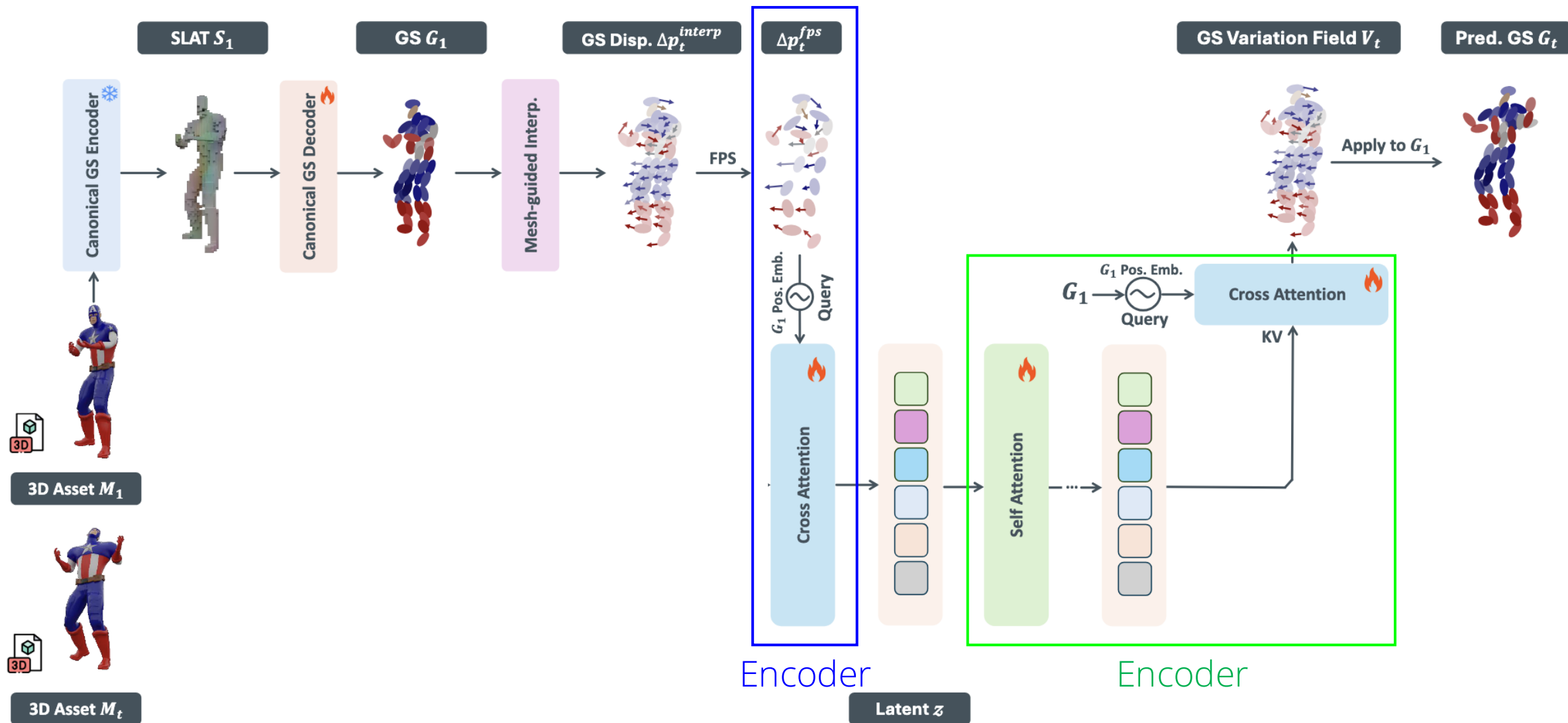
Gaussian Variation Field



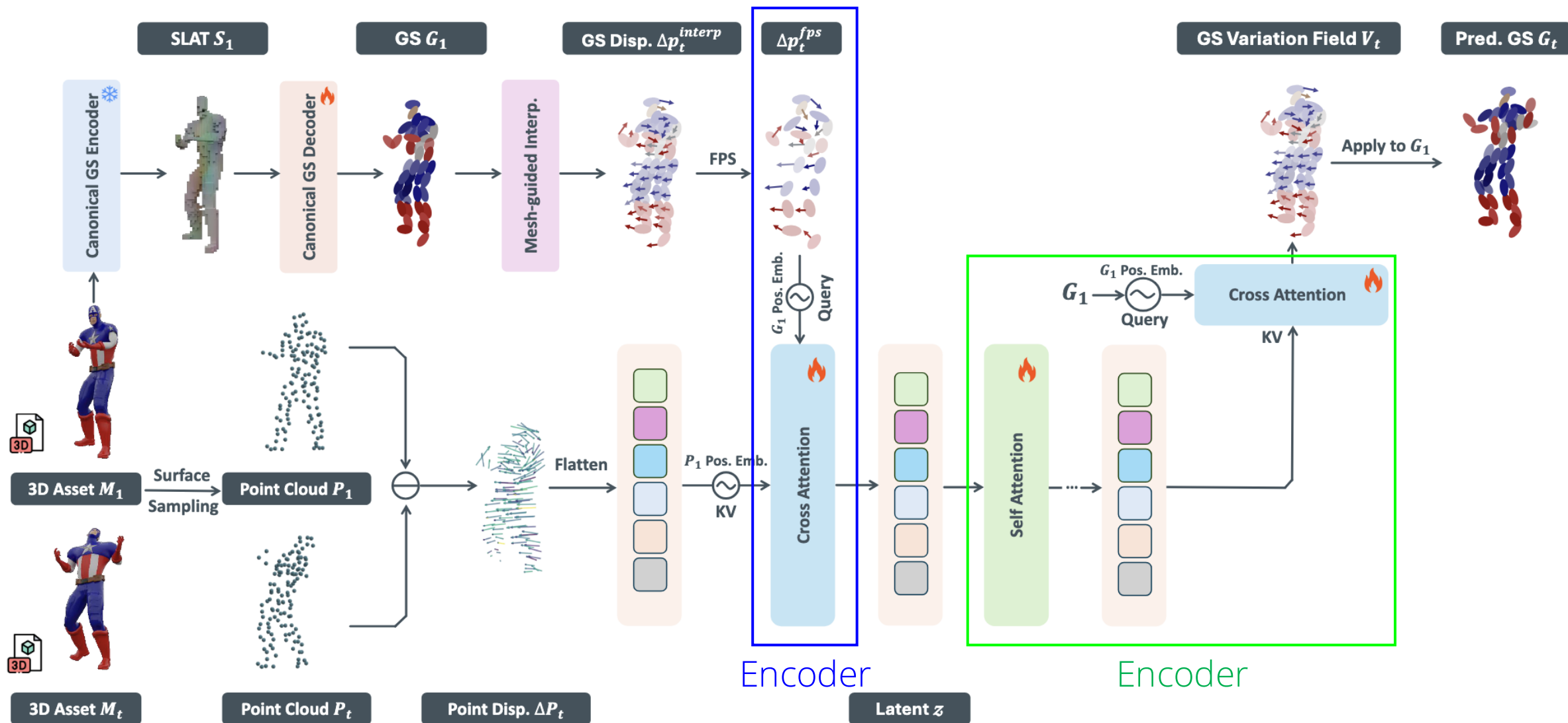
- Goal: generate a sequence of 3DGS models $\mathcal{G} = \{G_t\}_{t=1}^T$ that captures both the shape, appearance, and motion of the object.
- Gaussian Variation Field $\{\Delta G_t\}_{t=1}^T$ describes each Gaussian's attribute variations relative to G_1 over time.

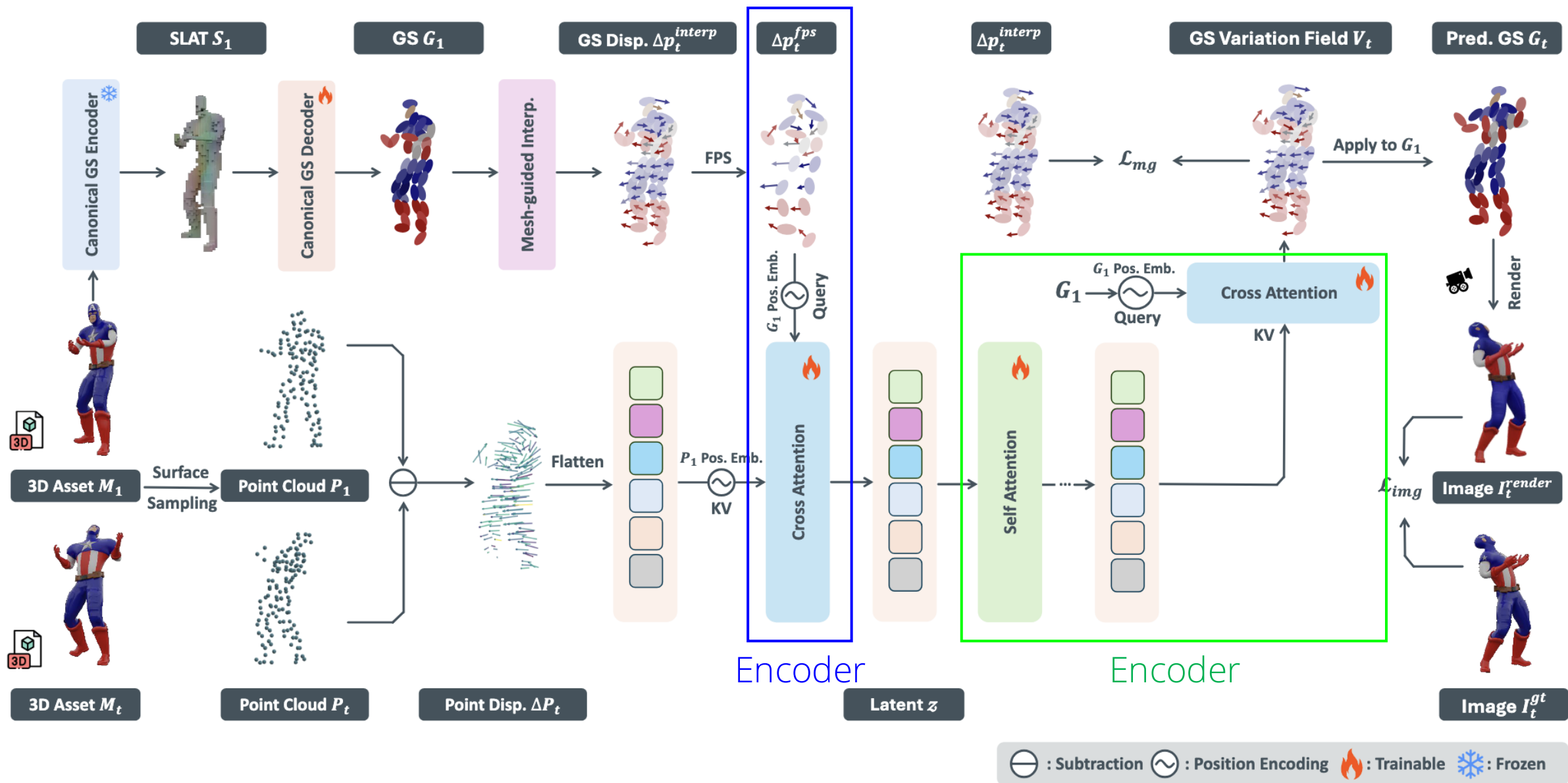
We Need An AutoEncoder that Compresses a Gaussian Variation Field





⊖ : Subtraction ⊗ : Position Encoding 🔥 : Trainable ❄️ : Frozen





Results



Input video



Consistent4D



STAG4D



L4GM

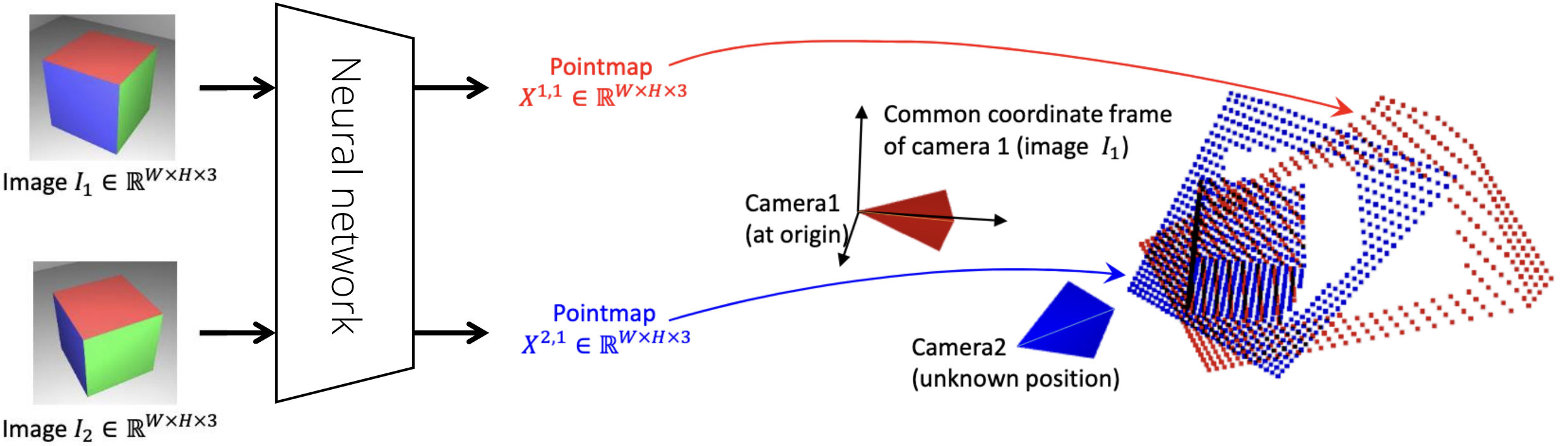


Ours

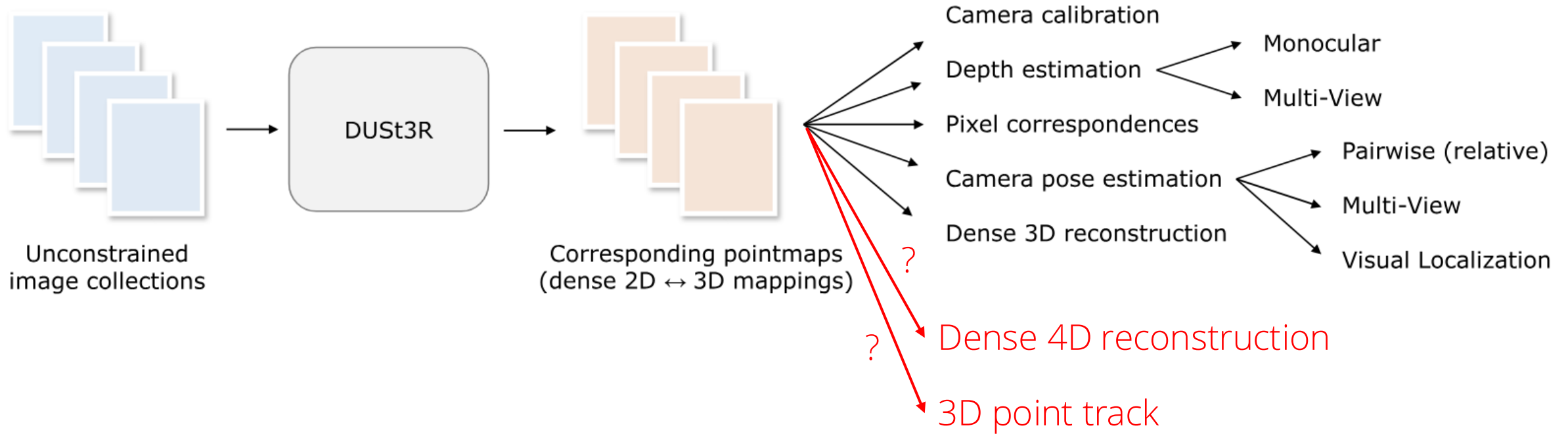
Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

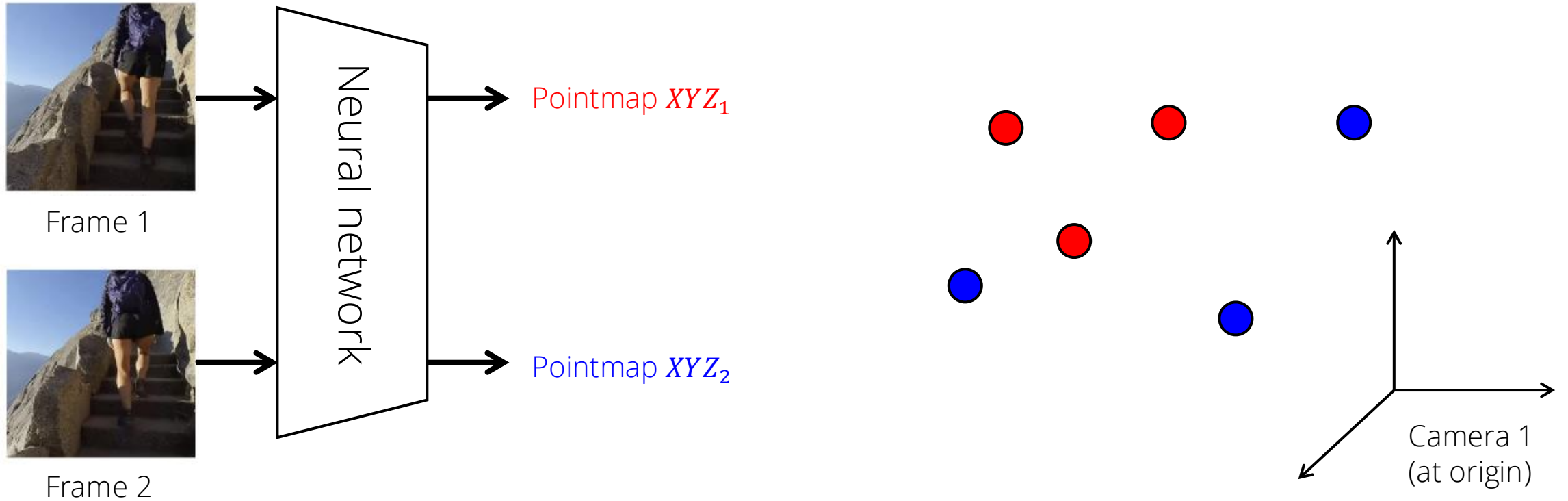
DUSt3R: Dense and Unconstrained Stereo 3D Reconstruction



A Model that Solves 3D Reconstruction How About 4D Reconstruction?

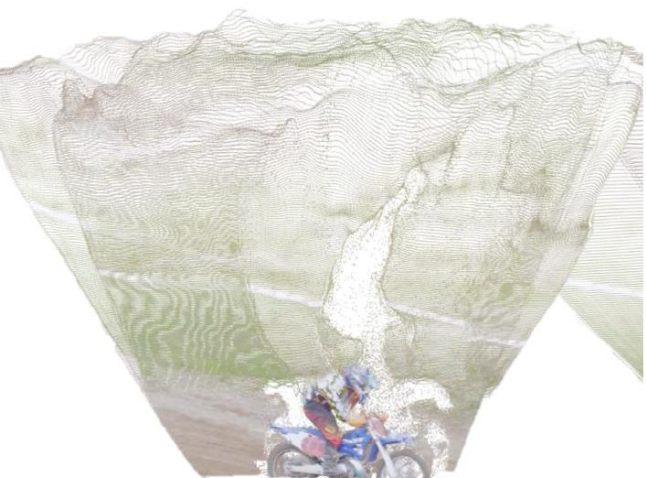


A Naïve Solution: Video Frames as Inputs



This Naïve Approach Doesn't Work Out-of-Box

Issue 1: Alignment based on dynamic object



DUST3R ✘



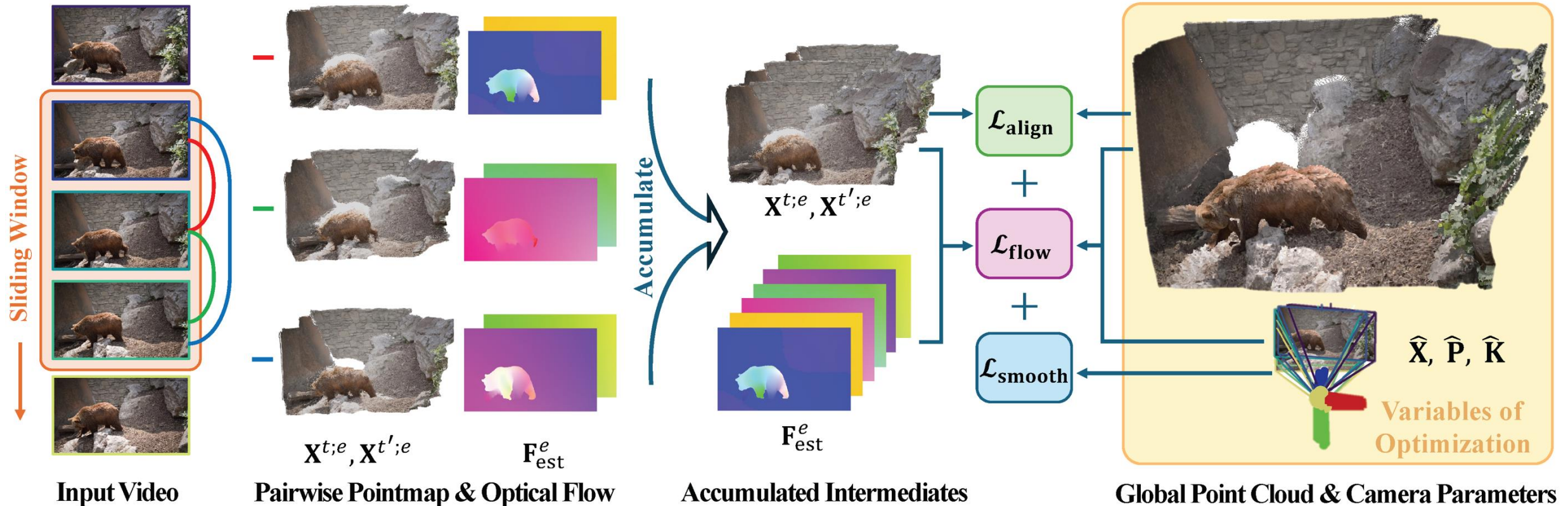
Issue 2: Foreground depth not estimated correctly



DUST3R ✘

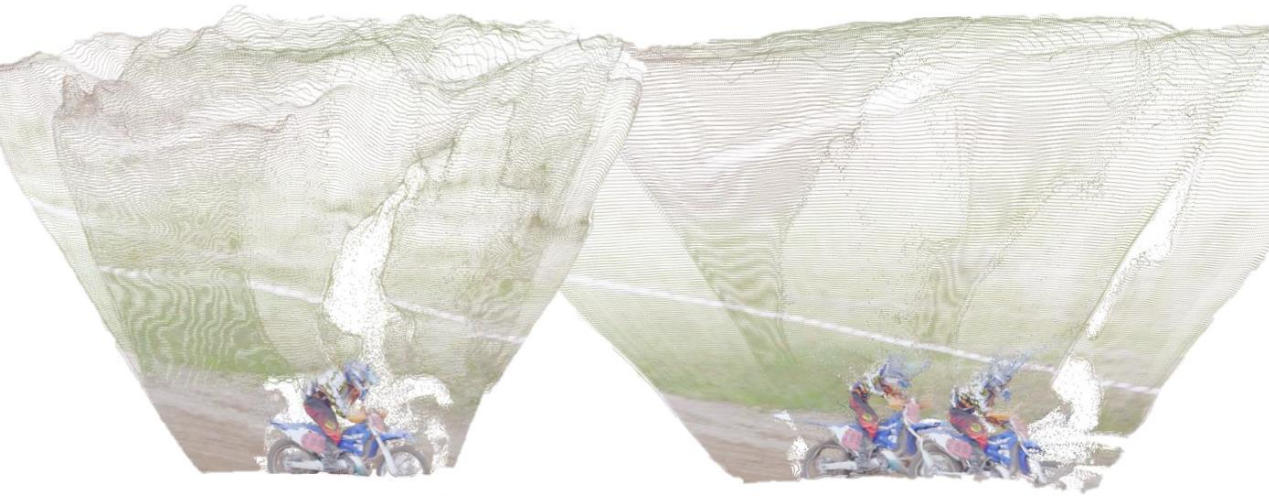


Idea: Globally Align Point Maps Based on Static Points



Results

Issue 1: Alignment based on dynamic object



DUST3R ❌

MonST3R

Issue 2: Foreground depth not estimated correctly



DUST3R ❌

MonST3R



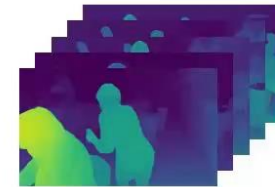
Results



Video Input



Dynamic Point Cloud & Camera Pose



Video Depth



Camera Intrinsics



Dynamic / Static Mask

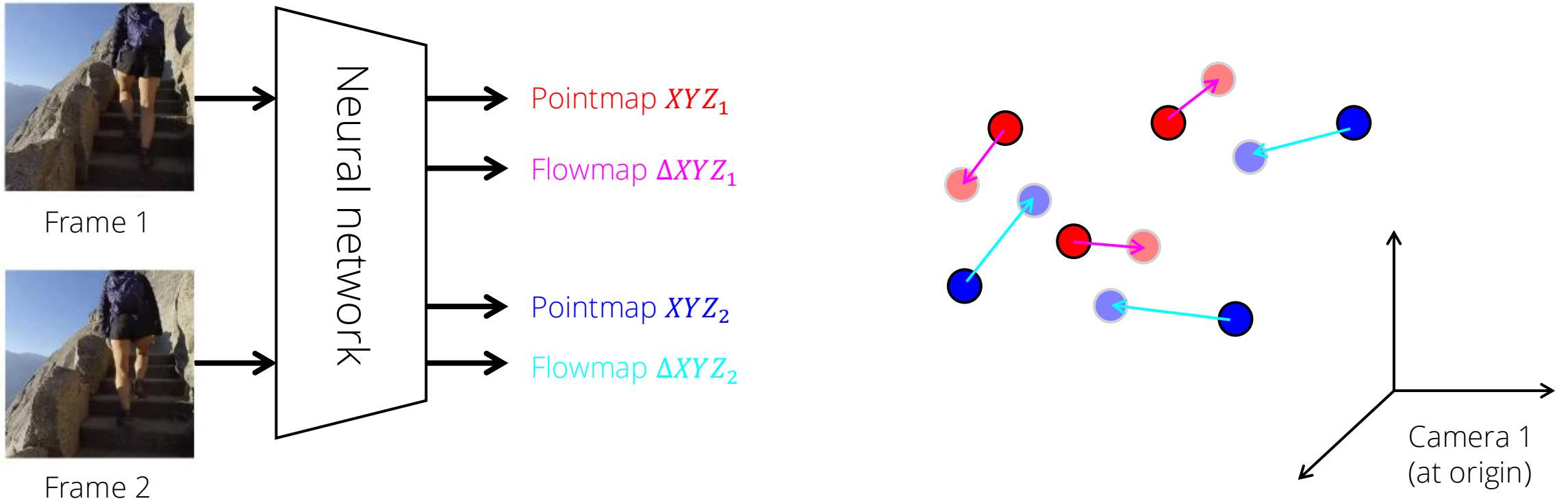
MONSt3R Has Two Problems

- Reconstructed 4D scenes lack correspondence
- Reconstructing a pair of video frames is inefficient

MONSt3R Has Two Problems

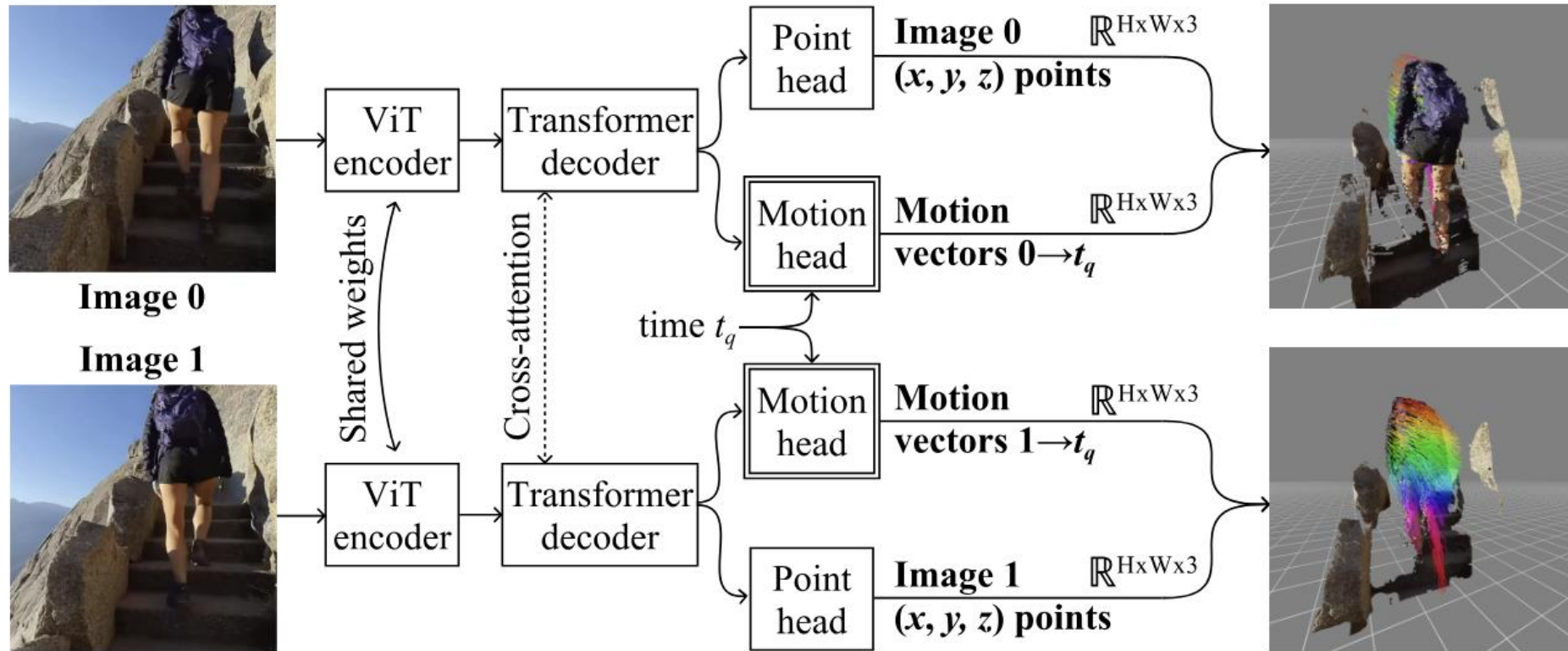
- Reconstructed 4D scenes lack correspondence
- Reconstructing a pair of video frames is inefficient

4D Reconstruction with Correspondence

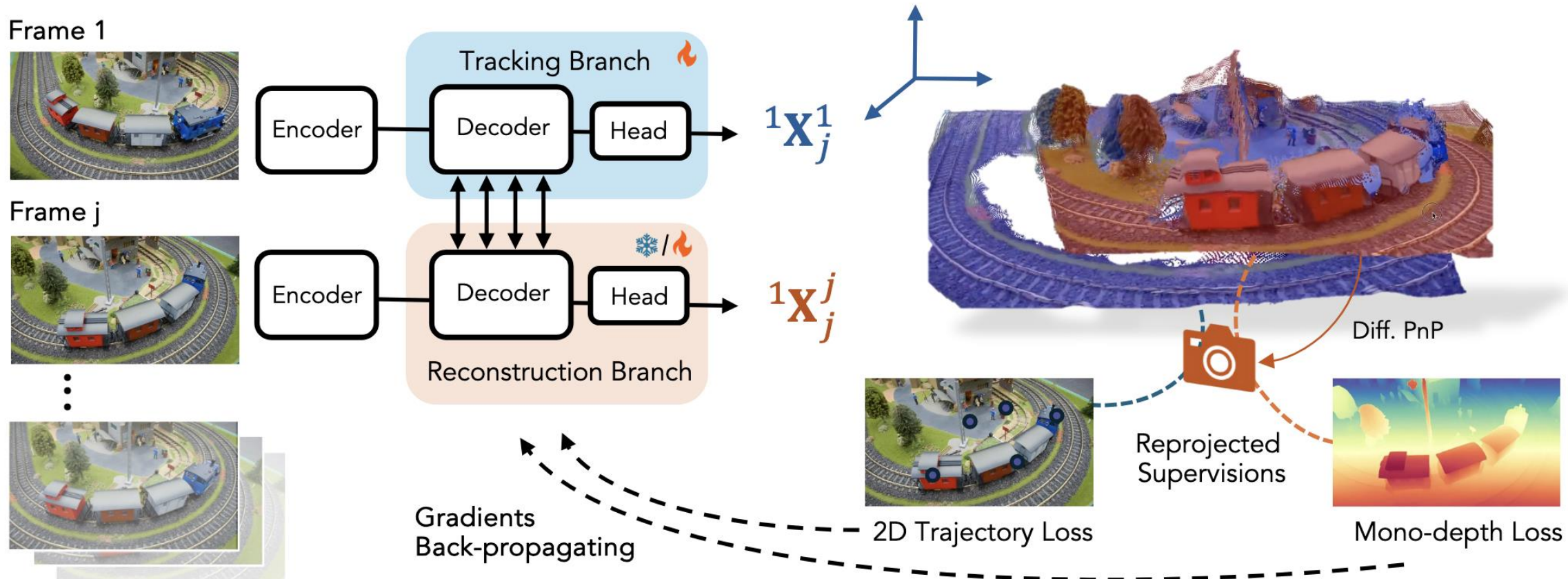


Everyone work on this direction...

For Example, Stereo4D



For Example, St4RTrack



For Example, Dynamic Point Map

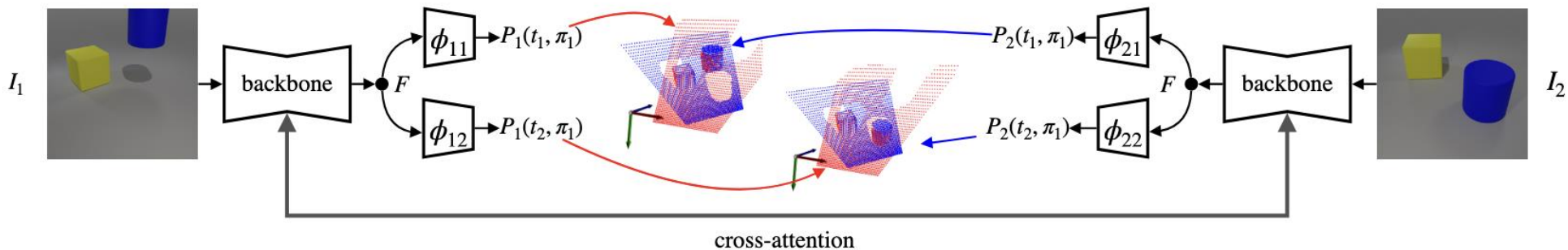


Figure 2. We propose to extend DUST3R to predict **Dynamic Point Maps (DPM)**. Each image in the pair is mapped to two point maps that correspond to the timestamps of the two images (pairs share the same colour in the figure). All points map are defined in the reference frame of image I_1 , undoing the effect of viewpoint change. Scene flow and space-time correspondences can be inferred immediately.

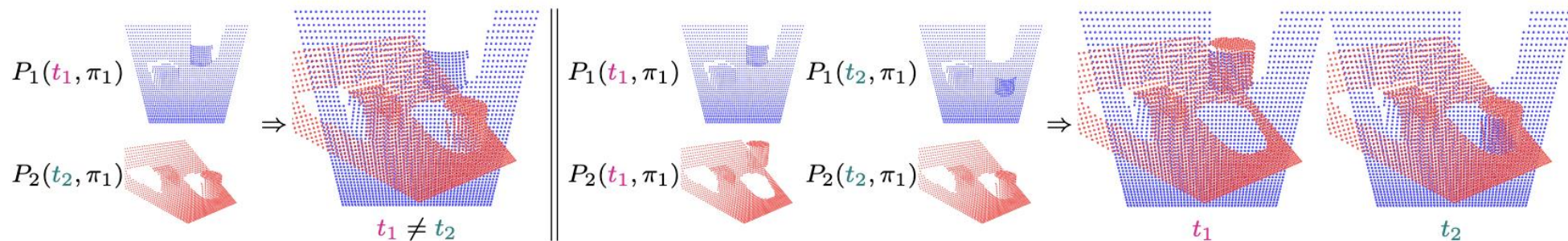
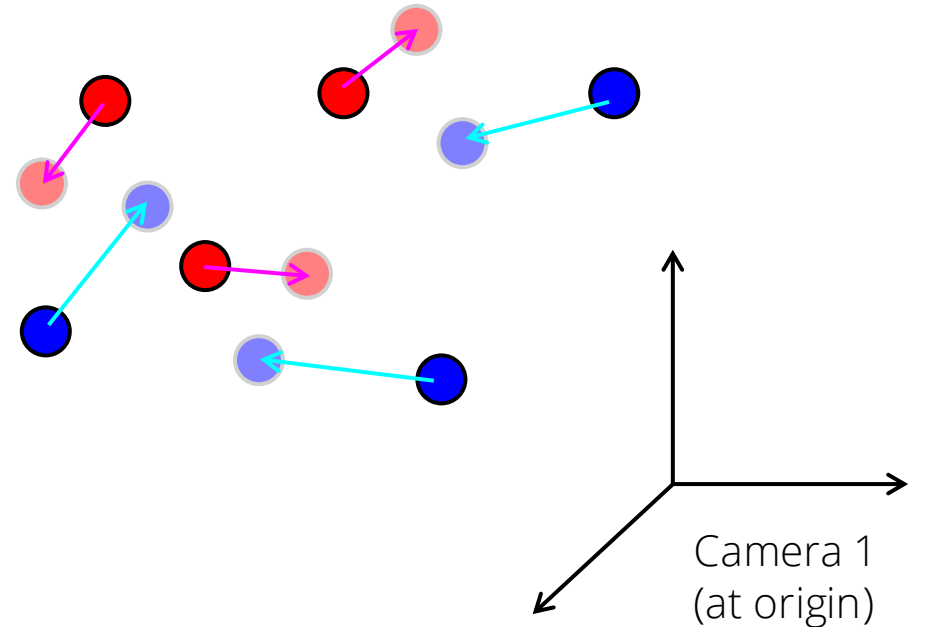
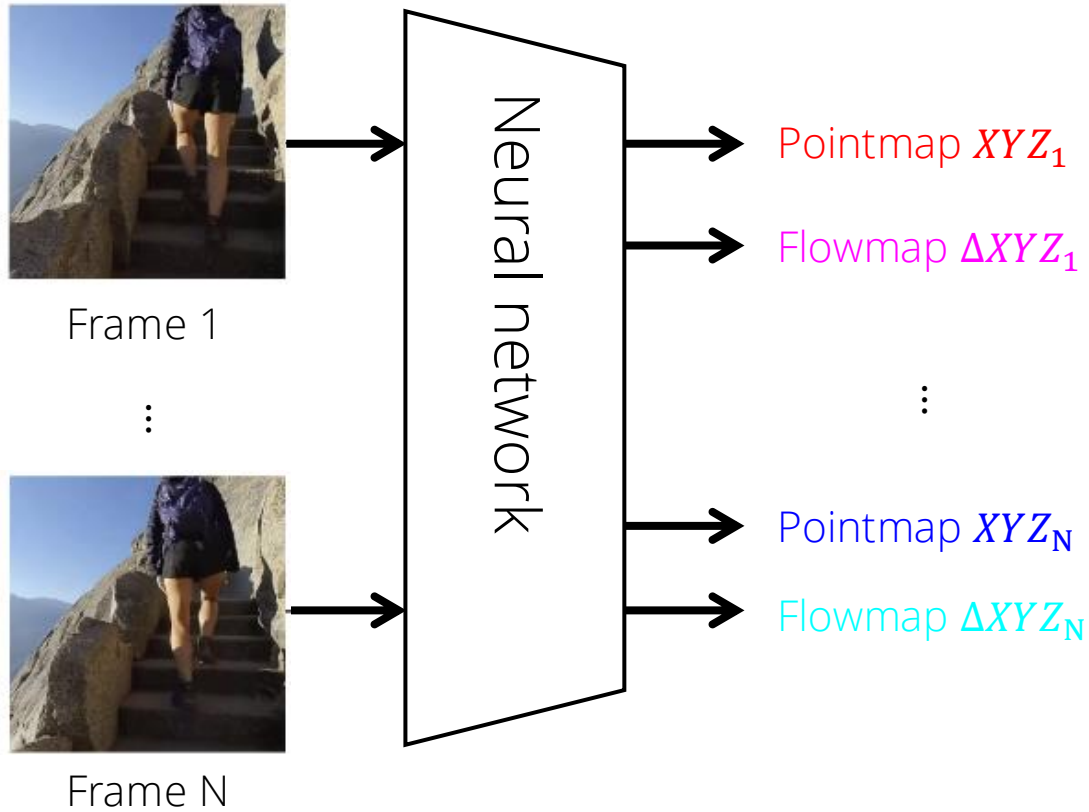


Figure 3. Left: Standard point maps applied to dynamic scenes as in MonST3R [80] fail to represent dynamics. The cylinder, which is moving downwards, breaks invariance when the point maps are overlaid. Right: Our Dynamic Point Maps correctly represent dynamics by also controlling time in addition to viewpoint. They allow to restore invariance while still representing the motion of the cylinder.

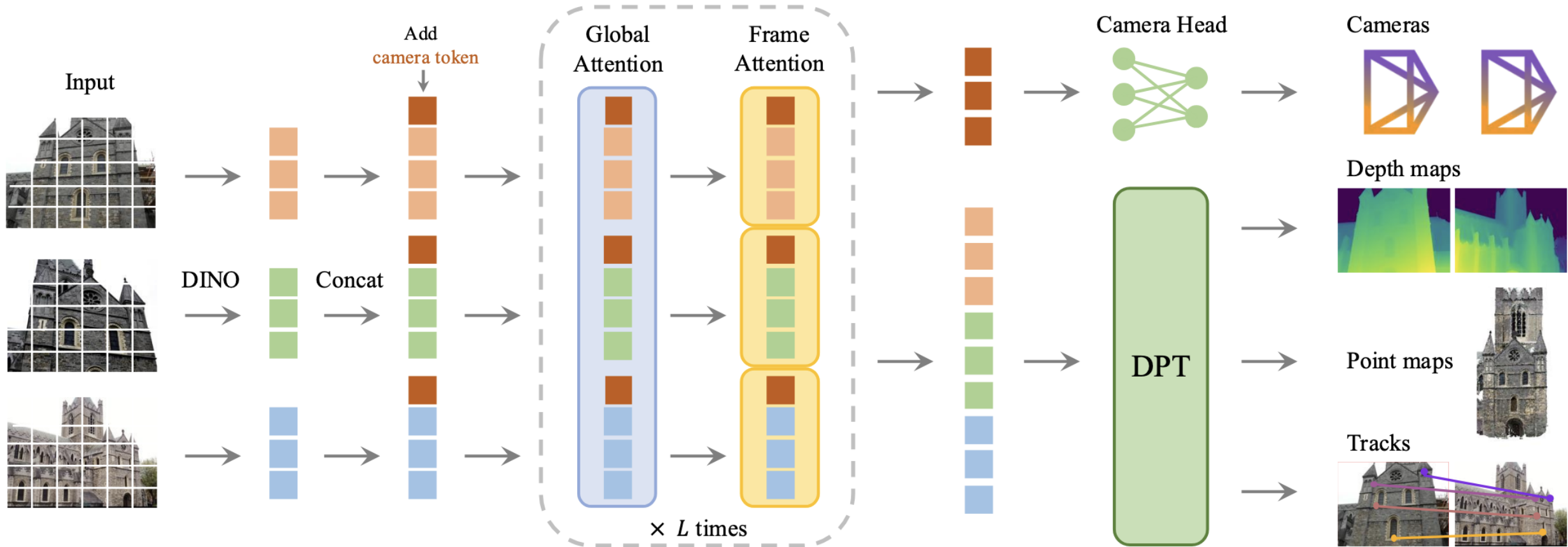
MONSt3R Has Two Problems

- Reconstructed 4D scenes lack correspondence
- Reconstructing a pair of video frames is inefficient

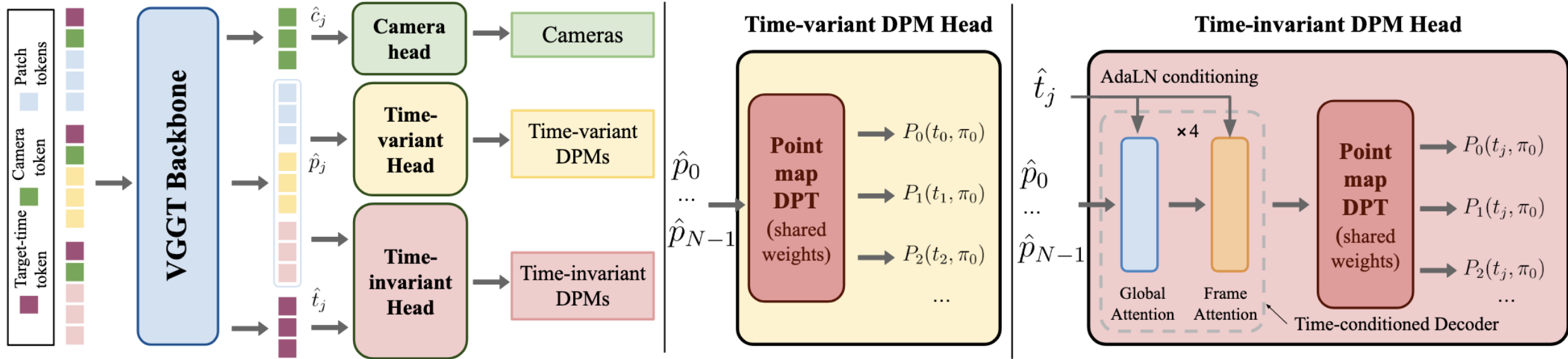
Idea: Use Large Visual Transformer



VGGT Goes Beyond Image Pairs



Append a Motion Head to VGGT



Results

Tennis

Content

- 2D Motion Representations
- 3D Motion Representations
 - 4D Volume Rendering with NERF
 - Explicit Motion Field with NERF
 - 4D Volume Rendering with 3DGS
 - Explicit Motion Field with 3DGS
 - Implicit Motion Field with 3DGS
 - 4D Foundation Models
 - Multi-view Video Generation (we'll cover in the next lecture)

What We Will Cover Next Week

- Generative Modeling