

Neural Networks as **PRIOR-BASED** **INFERENCE ALGORITHMS**

6.S980 | 2022-05-19

Project Proposal Details (due Oct 18th)

- 1-2 letter pages in CVPR 2022 layout (find on Canvas).
- Sections:
 - **Abstract / Problem definition:** What are you trying to do?
 - **Brief related work:** How is it done today, and what are the limits of current practice?
 - **Your proposed method:** What is new in your approach and why do you think it will be successful?
 - **Why is it interesting:** Who cares? If you are successful, what difference will it make?
 - Description of the first, simplest experiment that will show you if your idea has merit
 - Description of final experiment - best case, what will you be able to show?

Project Details

- Project proposals due in ~2 weeks (Oct 18th).
- **Vincent's project office hours:**
 - Wednesday, October 5th, 10 - 11 pm.
 - Friday, October 7th, 2 - 3 pm (**weekly**)
 - 10-minute time slots per team.
 - Sign up [here](#).
- My advice: Try preparing the answers to the questions on the previous slide - that's the fastest, most effective pitch you can do!

Example Paper Discussion

Decomposing NeRF for Editing via Feature Field Distillation

Sosuke Kobayashi
Preferred Networks, Inc.
`sosk@preferred.jp`

Eiichi Matsumoto
Preferred Networks, Inc.
`matsumoto@preferred.jp`

Vincent Sitzmann
Massachusetts Institute of Technology
`sitzmann@mit.edu`

pfnet-research.github.io/distilled-feature-fields/

Abstract

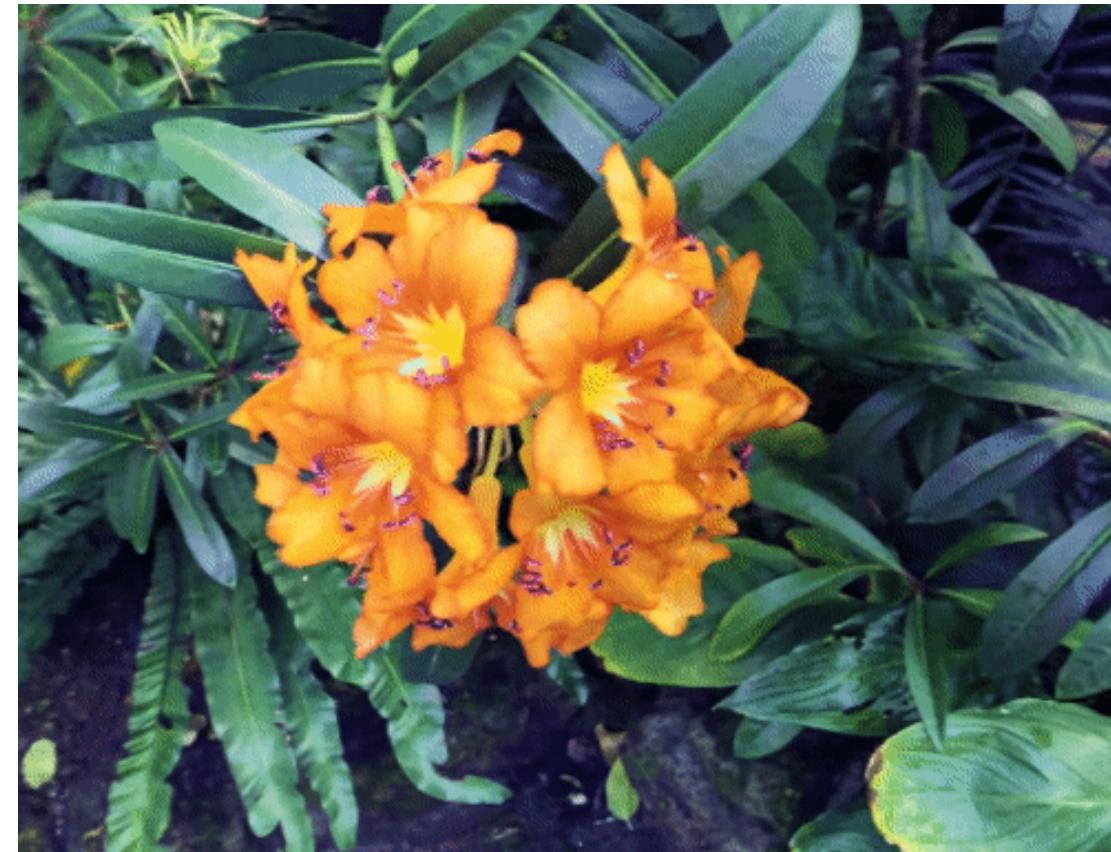
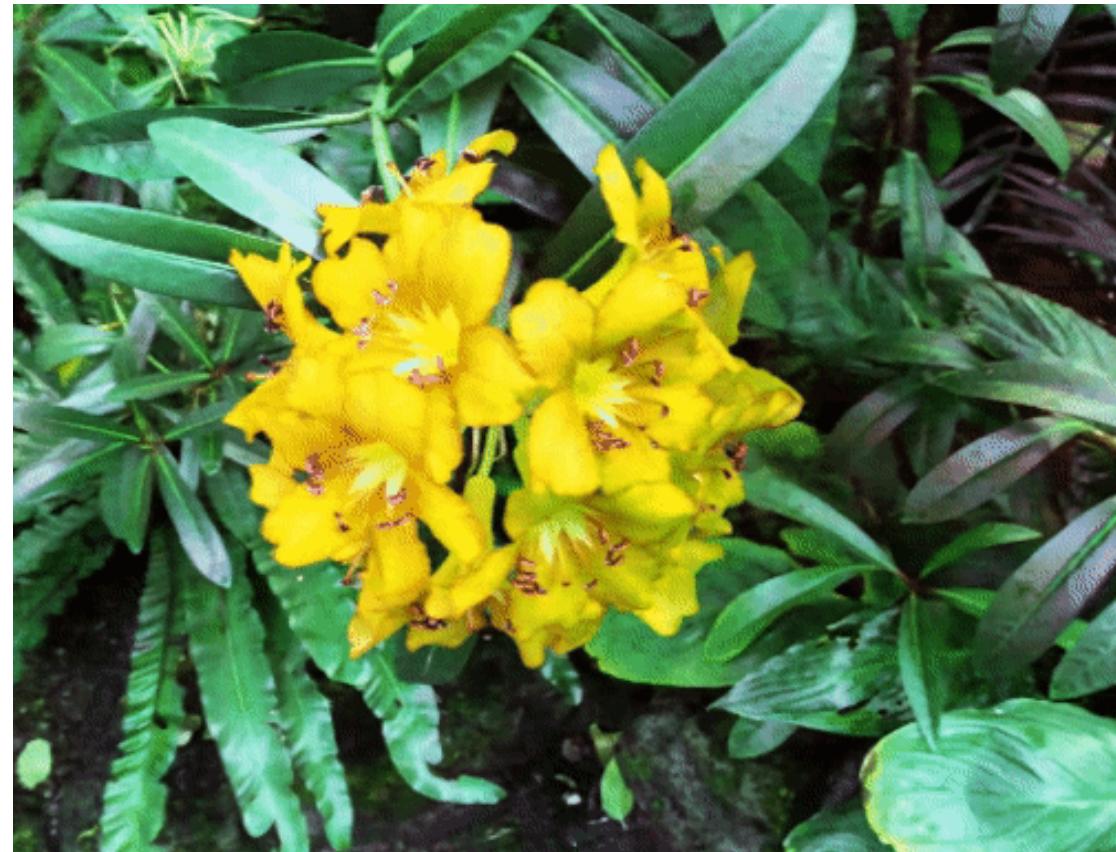
Emerging neural radiance fields (NeRF) are a promising scene representation for computer graphics, enabling high-quality 3D reconstruction and novel view synthesis from image observations. However, editing a scene represented by a NeRF is challenging, as the underlying connectionist representations such as MLPs or voxel grids are not object-centric or compositional. In particular, it has been difficult to selectively edit specific regions or objects. In this work, we tackle the problem of semantic scene decomposition of NeRFs to enable query-based local editing of the represented 3D scenes. We propose to distill the knowledge of off-the-shelf, self-supervised 2D image feature extractors such as CLIP-LSeg or DINO into a 3D feature field optimized in parallel to the radiance field. Given a user-specified query of various modalities such as text, an image patch, or a point-and-click selection, 3D feature fields semantically decompose 3D space

Problem Statement: What is it trying to do?

- Given: **Many** 2D images + cameras of **single** 3D scene
- Today: can reconstruct 3D scenes via diff. rendering, but not easily editable!
- Goal: Reconstruct, then **edit** 3D scene representation
- Edits are color, deletion, copy-paste, moving of scene parts, ...



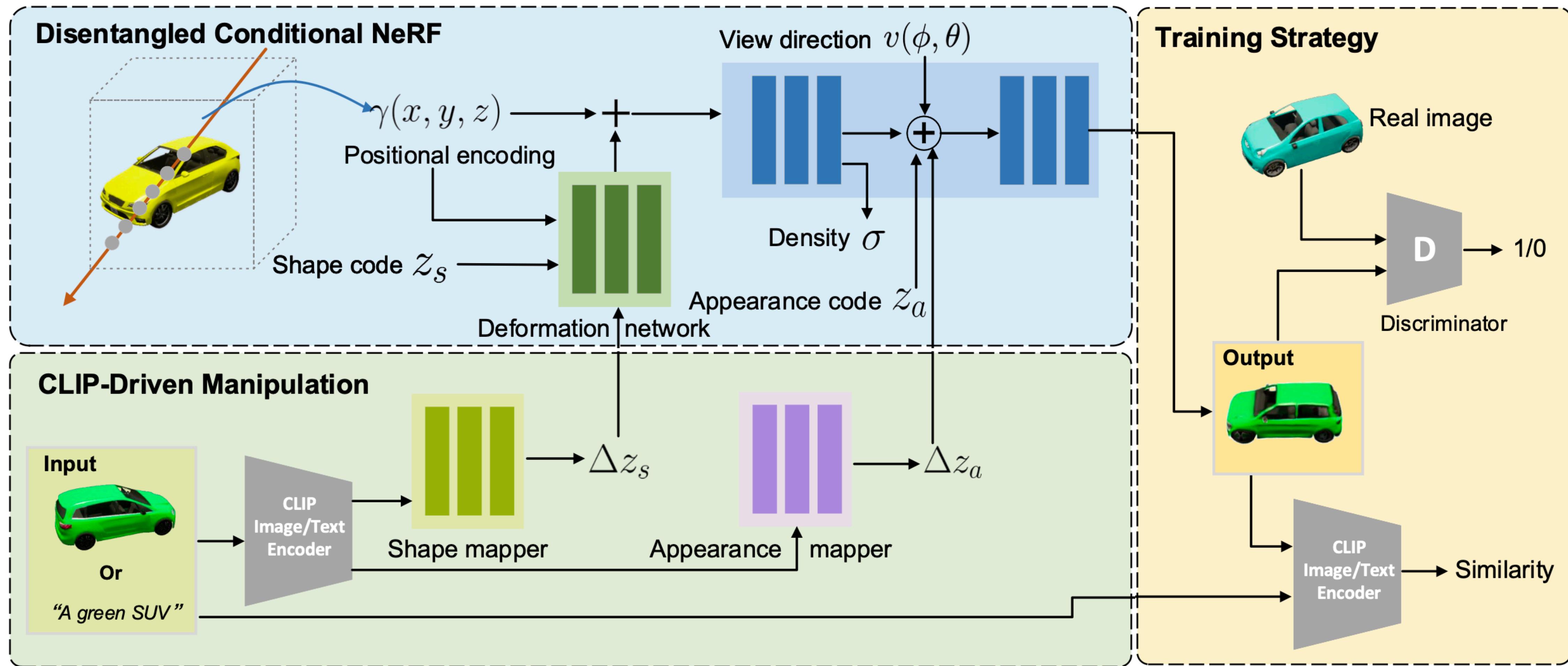
Related Work: CLIPNeRF



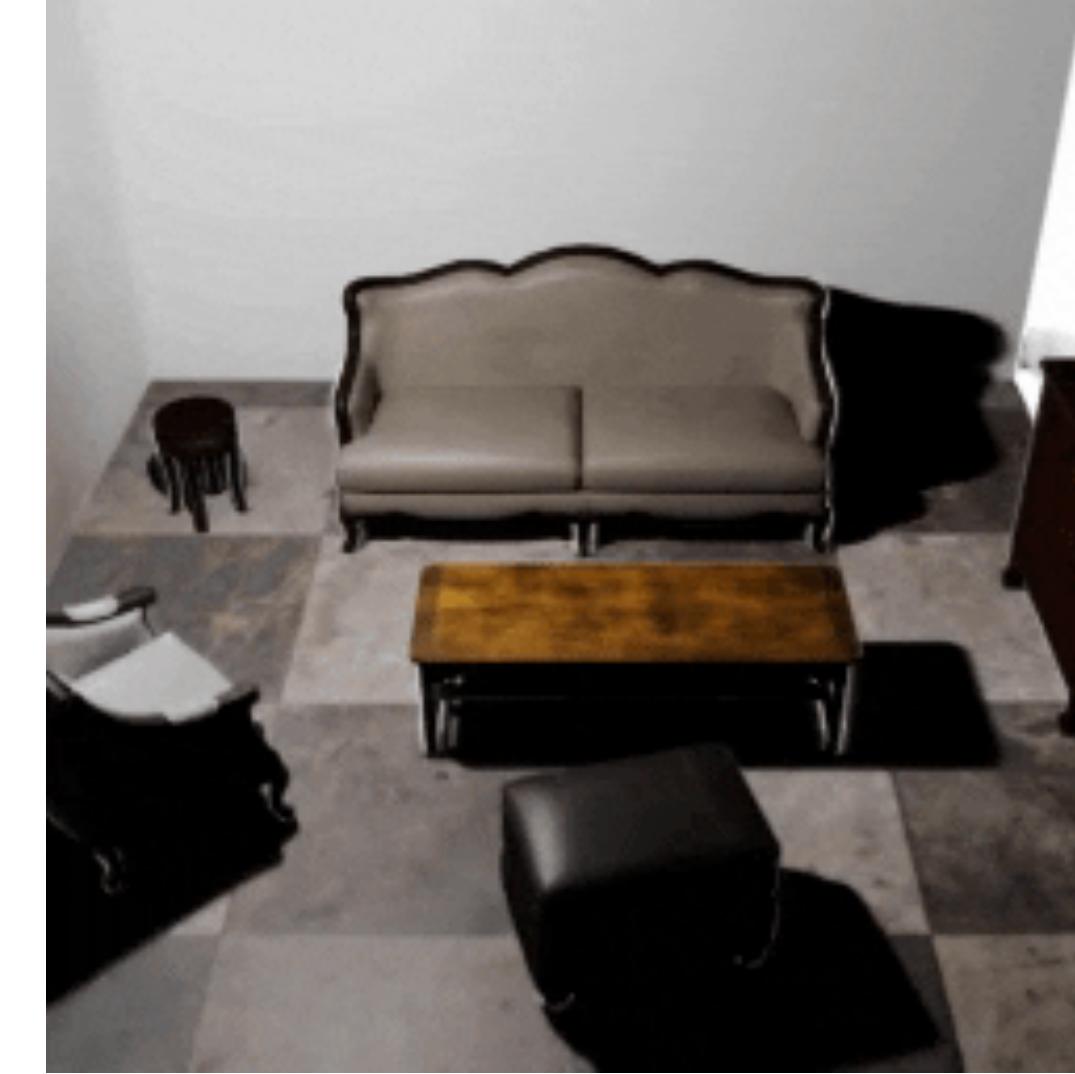
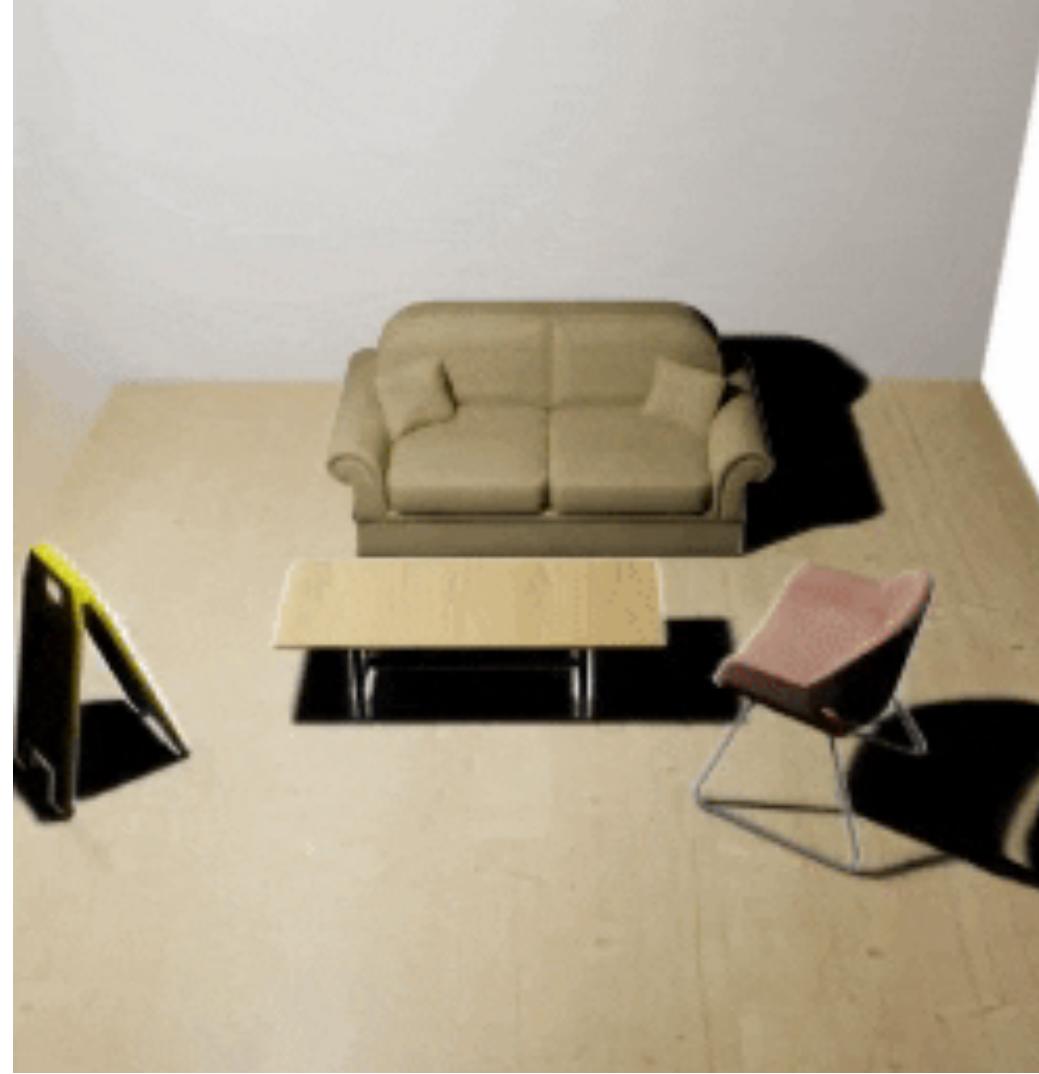
, ,

Wang et al. [84] and Jain et al. [32] use CLIP to edit or generate a single-object NeRF with a text prompt query by optimizing the NeRF parameters to generate images matched with the text. While such methods are promising, they do not enable selective editing of specific scene regions. For example, the prompt “yellow flowers” may affect unintended scene regions, such as the leaves of a plant.

Related Work: CLIPNeRF



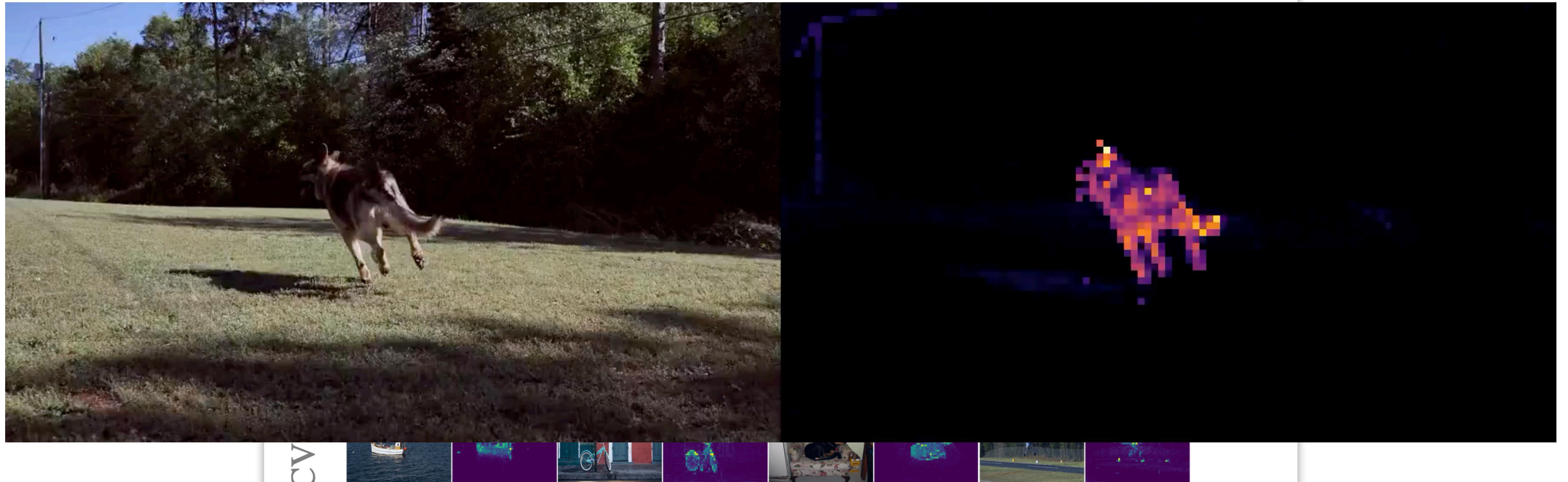
Related Work: Object-Centric Representations



”

Other studies also explored reconstruction with more structured hybrid representations via pipelines specialized to a domain (e.g., traffic scene) [58, 22, 36] or situation (e.g., each object data is independently accessible) [25, 24, 92]. Note that this line of work defines and constrains domains or the types of segmentation during or before training, and thus limits the degrees of freedom for editable scenes and objects.

Background: Dino Feature Extractor



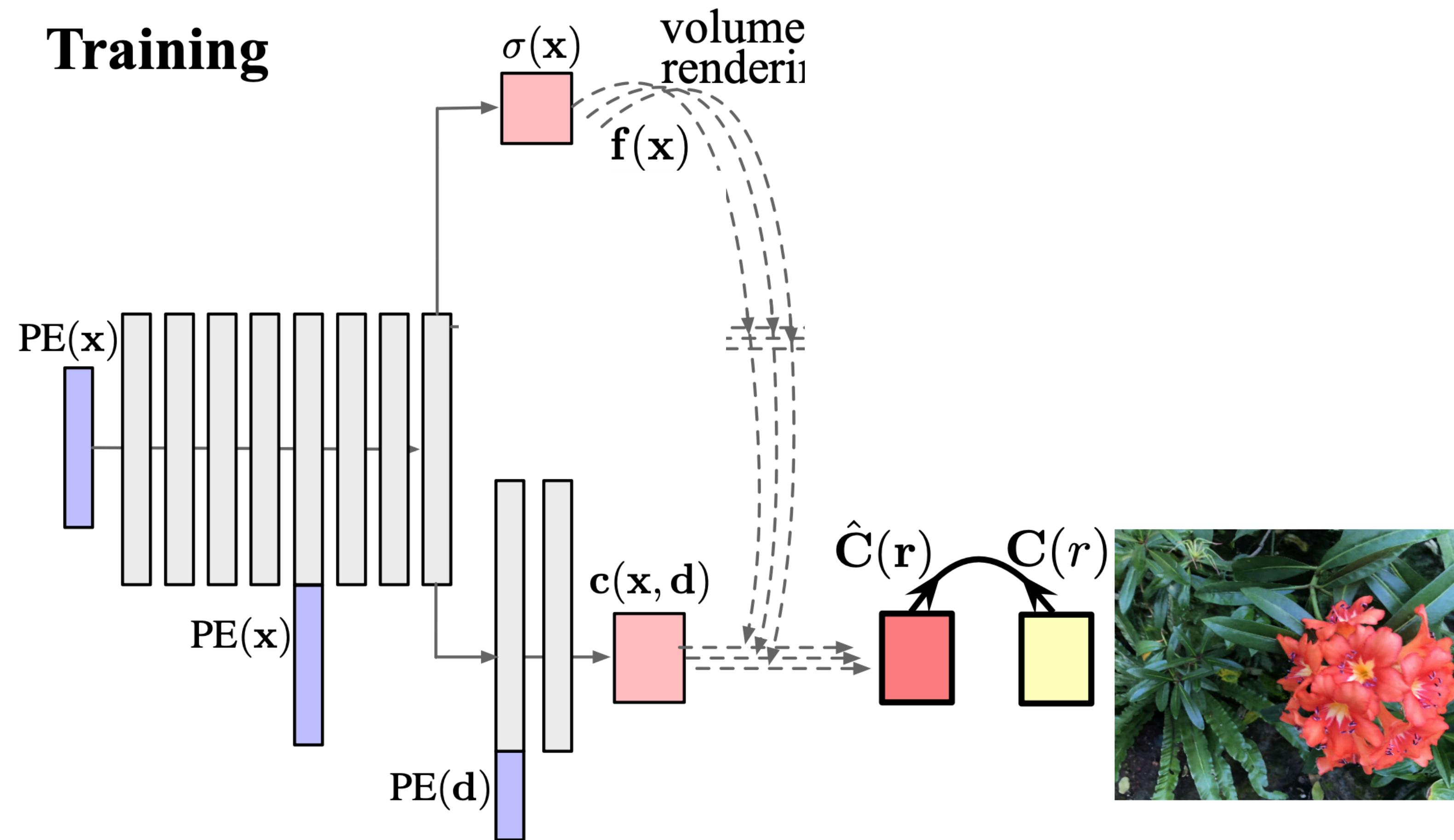
[CS.CV]

Figure 1: **Self-attention from a Vision Transformer with 8×8 patches trained with no supervision.** We look at the self-attention of the [CLS] token on the heads of the last layer. This token is not attached to any label nor supervision. These maps show that the model automatically learns class-specific features leading to unsupervised object segmentations.

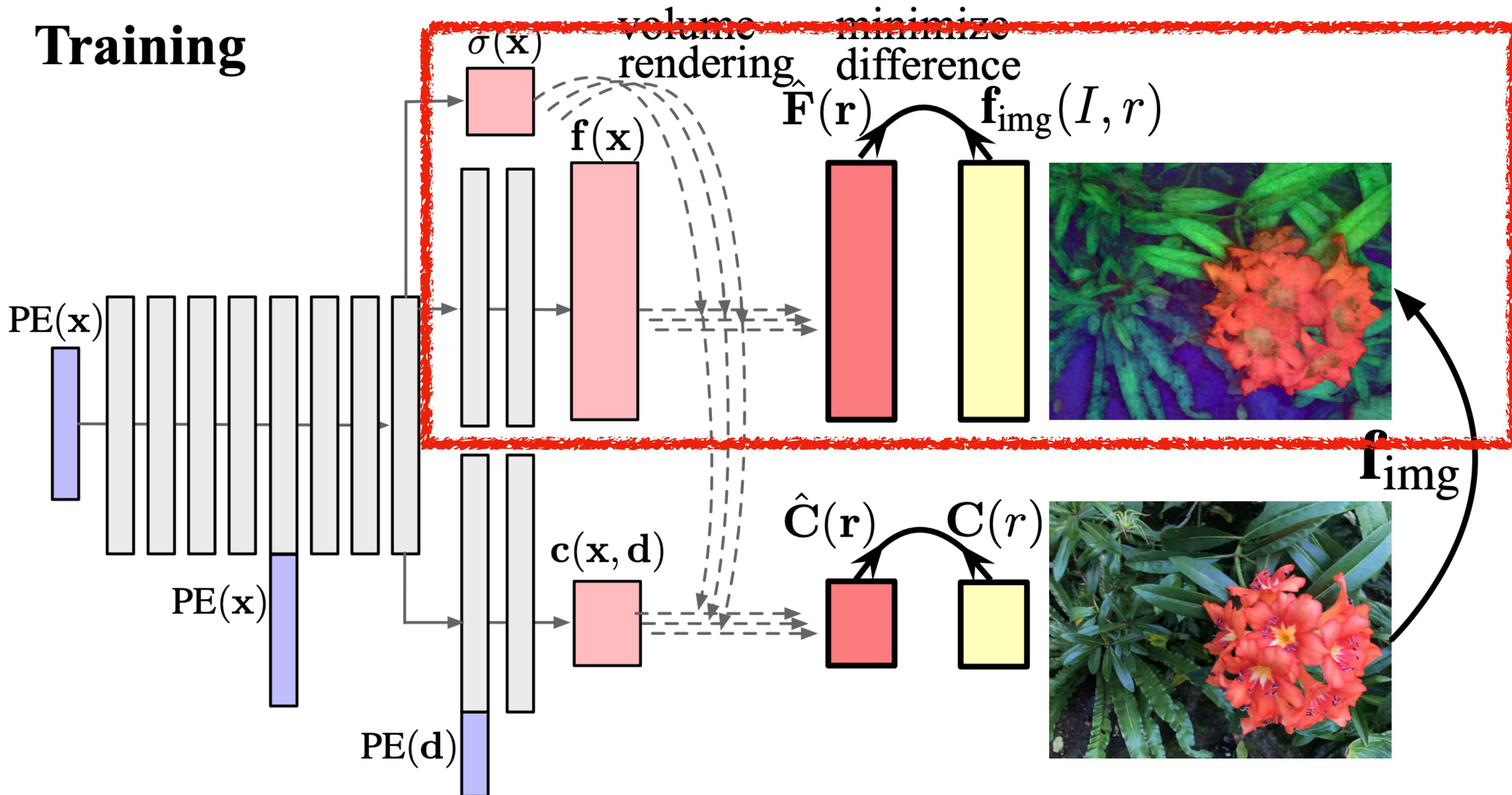
Background: Dino Feature Extractor



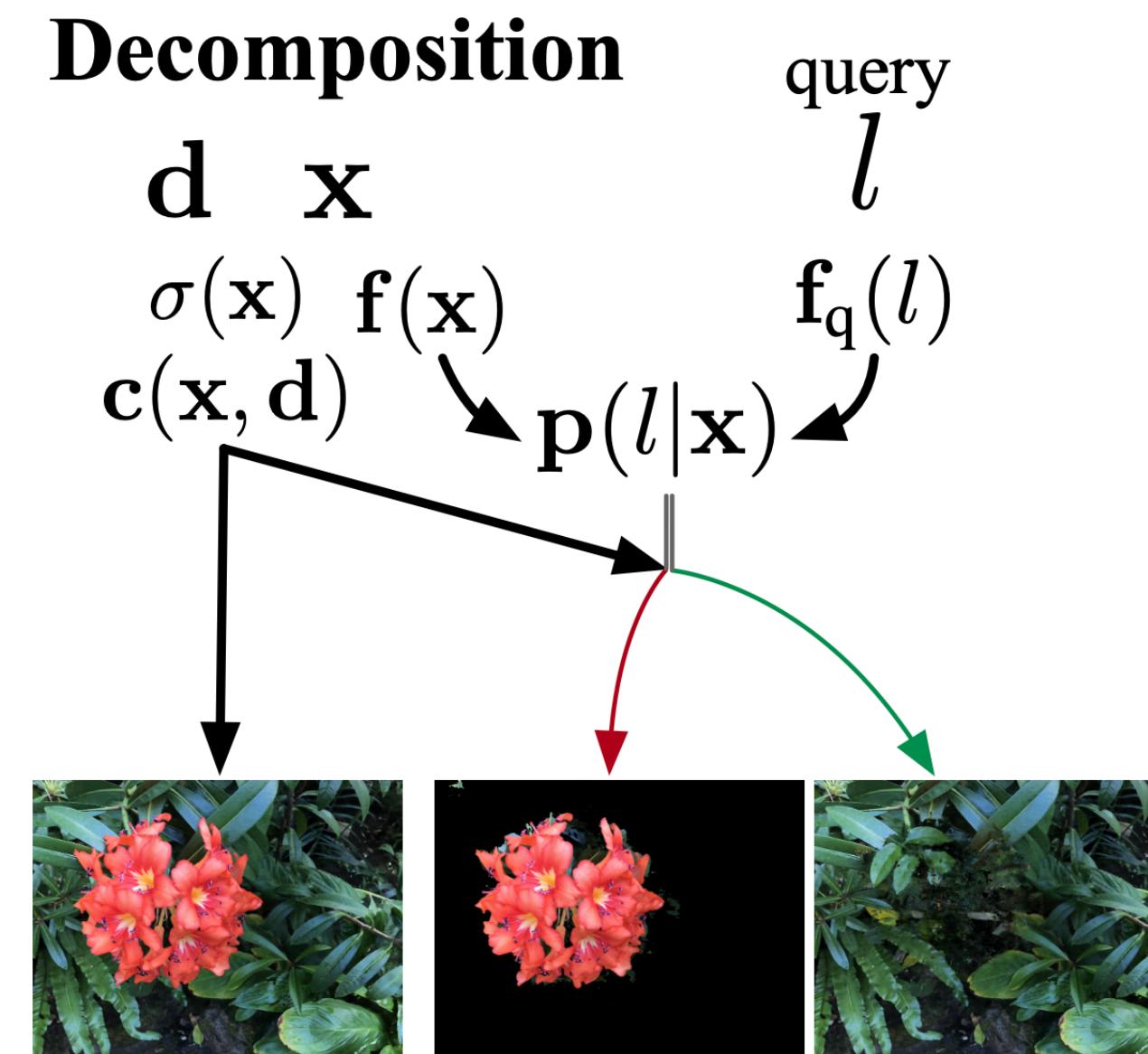
Method: Training Feature Field



Method: Training Feature Field



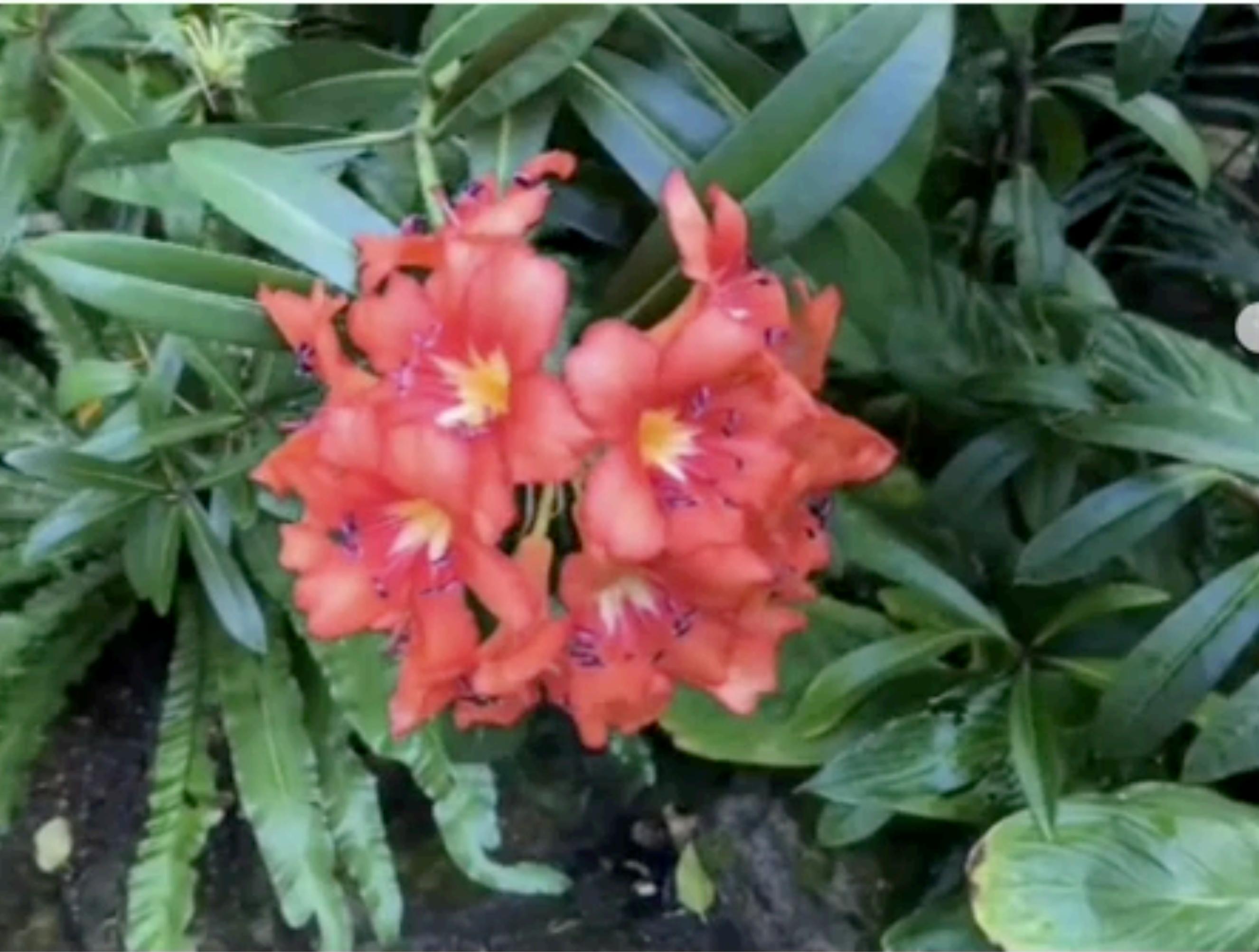
Method: Selecting Regions



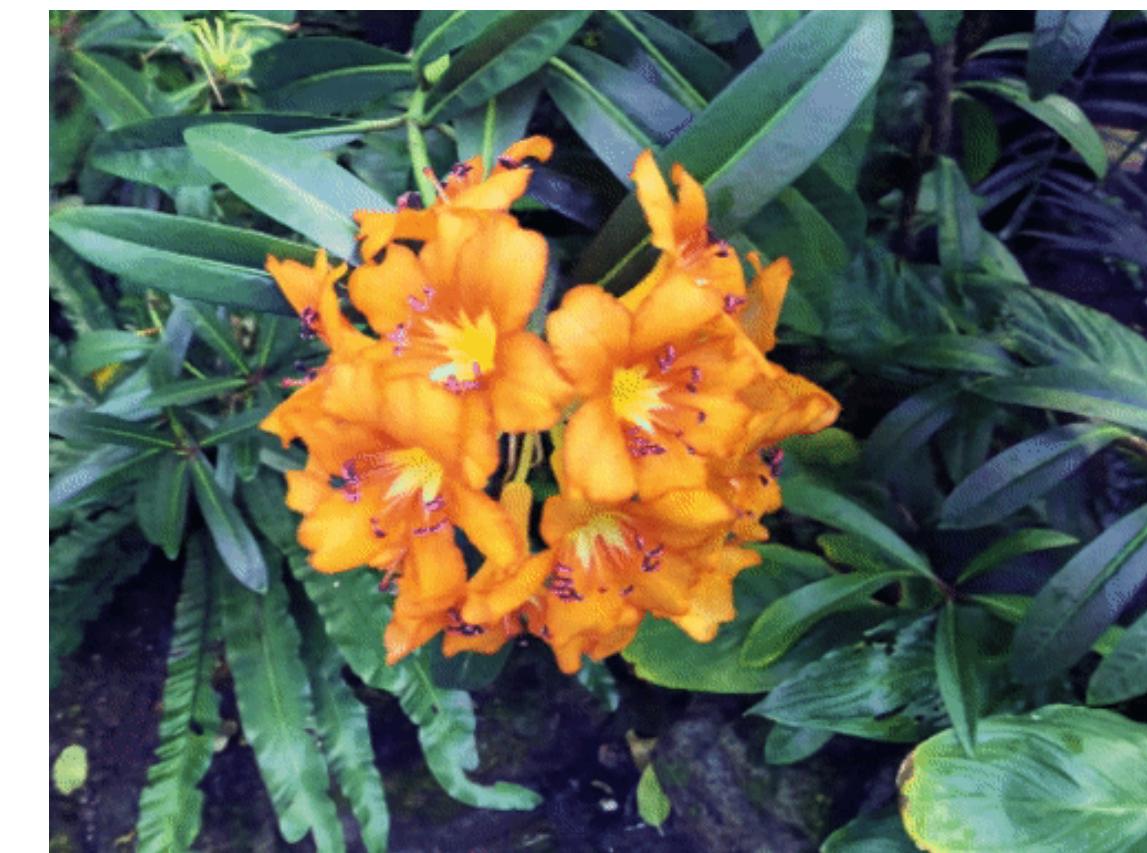
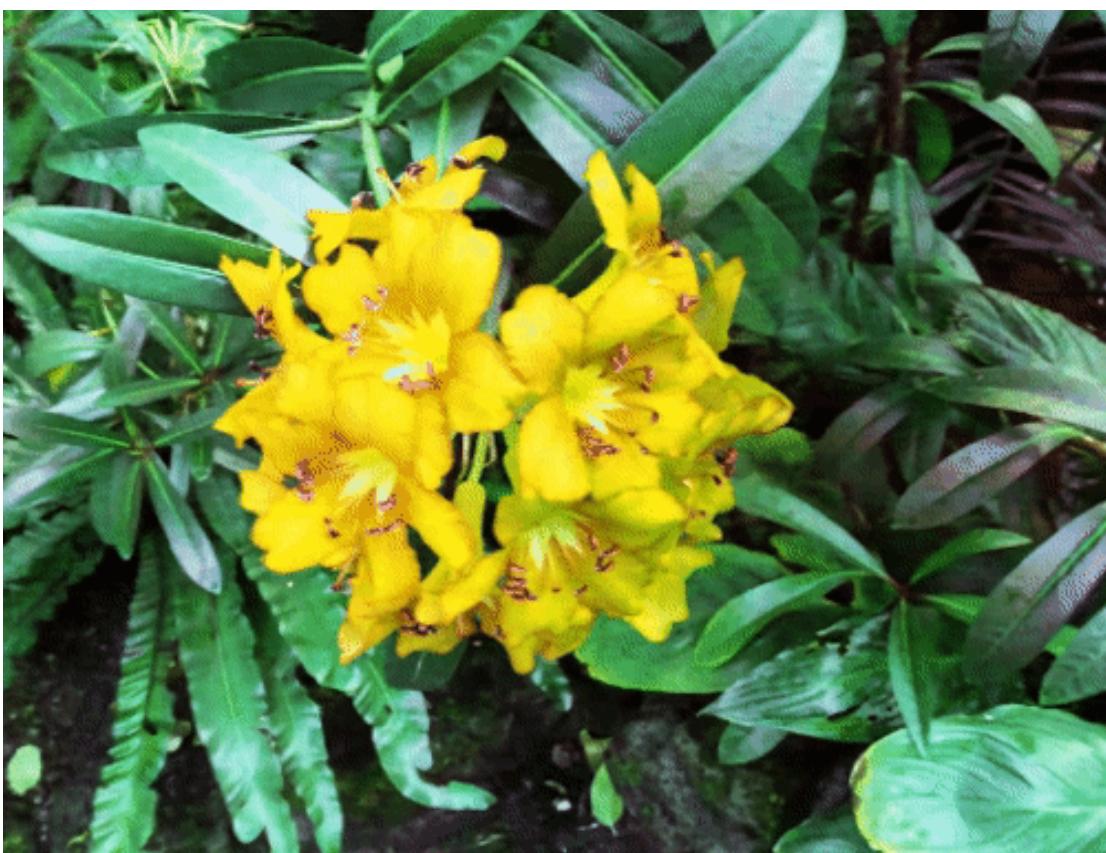
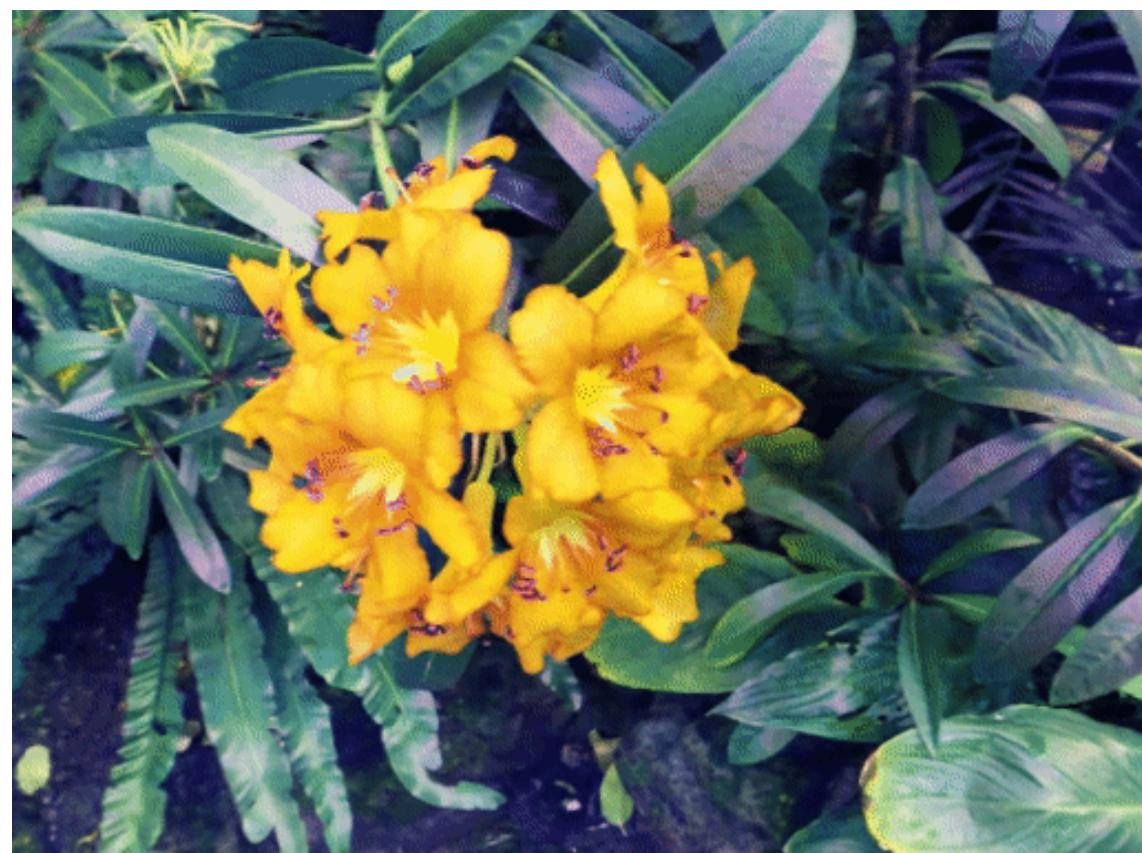
A trained DFF model can perform 3D zero-shot segmentation by its feature field \mathbf{f} and a query encoder \mathbf{f}_q . Probability of a label l of a point \mathbf{x} in the 3D space, $\mathbf{p}(l|\mathbf{x})$, is calculated by dot product of the 3D feature $\mathbf{f}(\mathbf{x})$ and text label feature $\mathbf{f}_q(l)$ followed by a softmax:

$$\mathbf{p}(l|\mathbf{x}) = \frac{\exp(\mathbf{f}(\mathbf{x})\mathbf{f}_q(l)^T)}{\sum_{l' \in \mathcal{L}} \exp(\mathbf{f}(\mathbf{x})\mathbf{f}_q(l')^T)} . \quad (5)$$

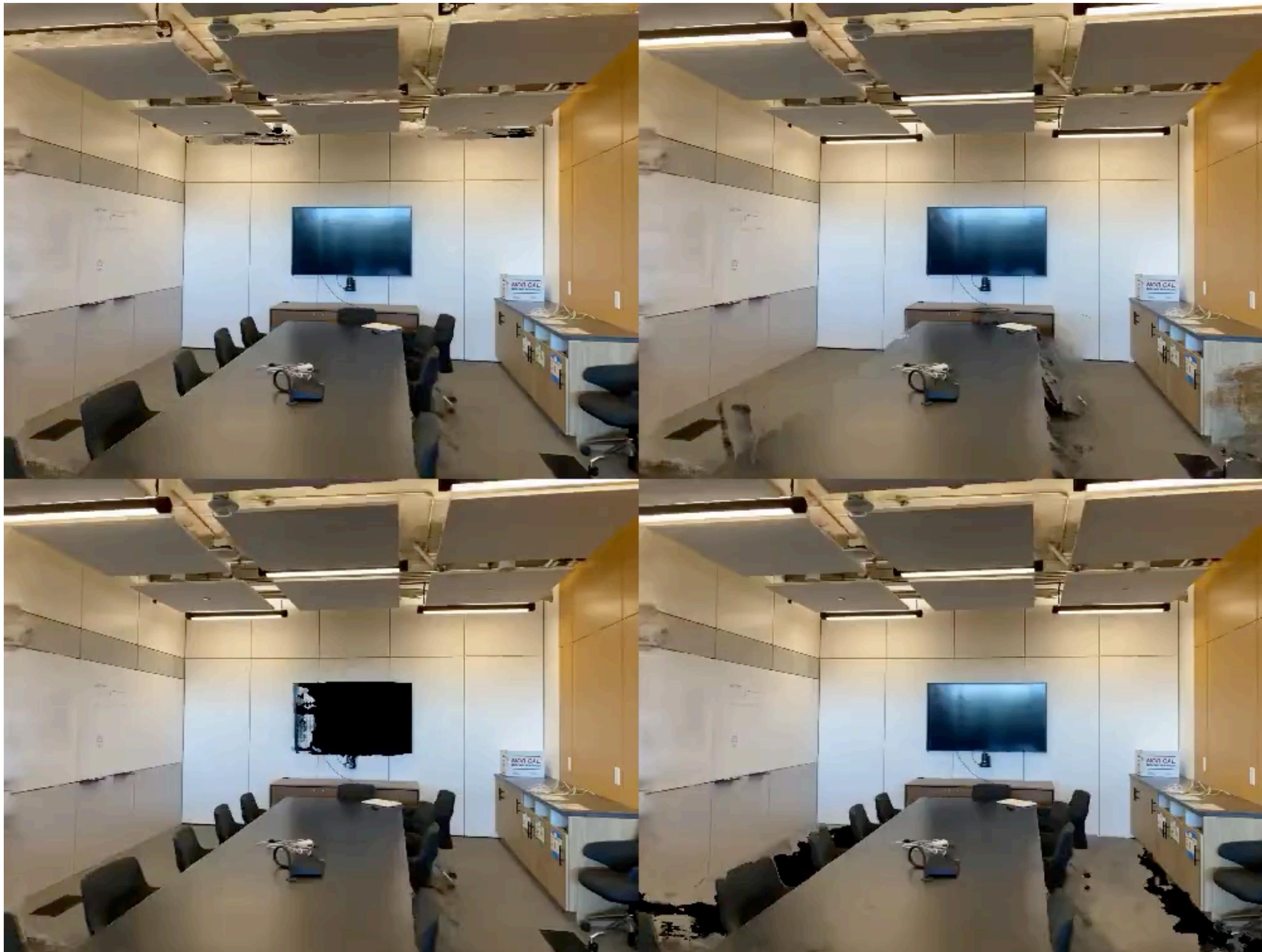
Results



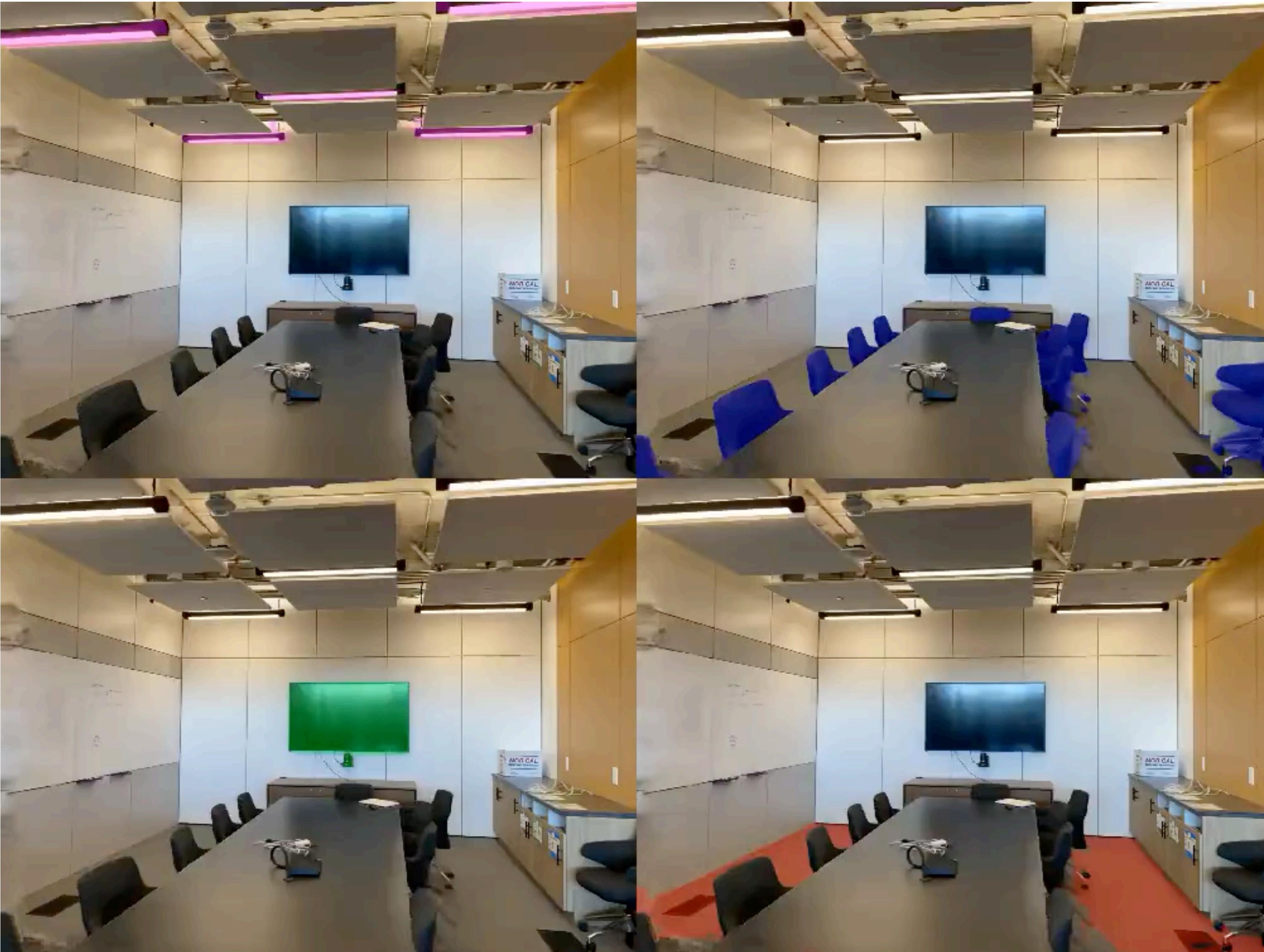
Baseline



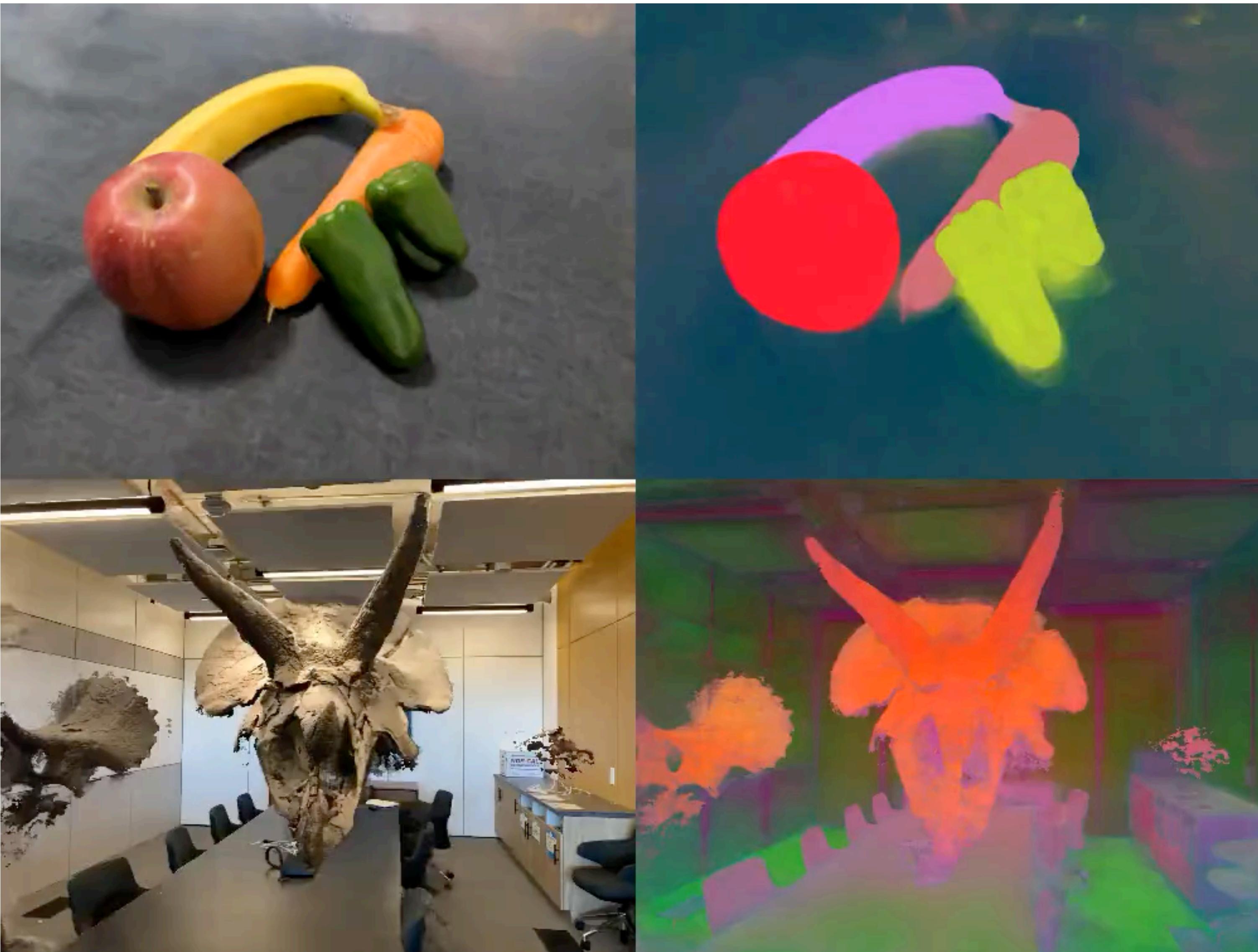
Deleting objects!



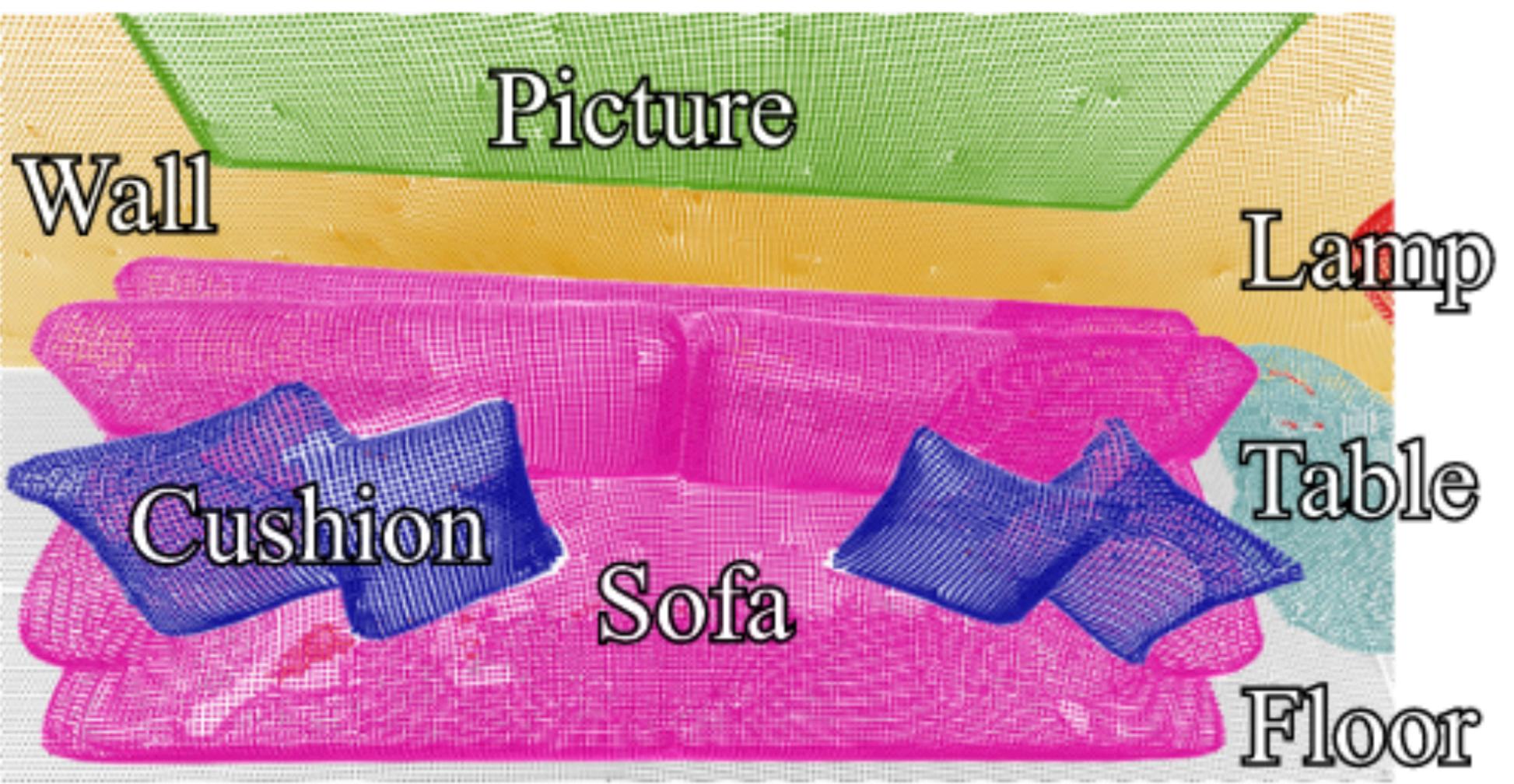
Re-coloring objects!



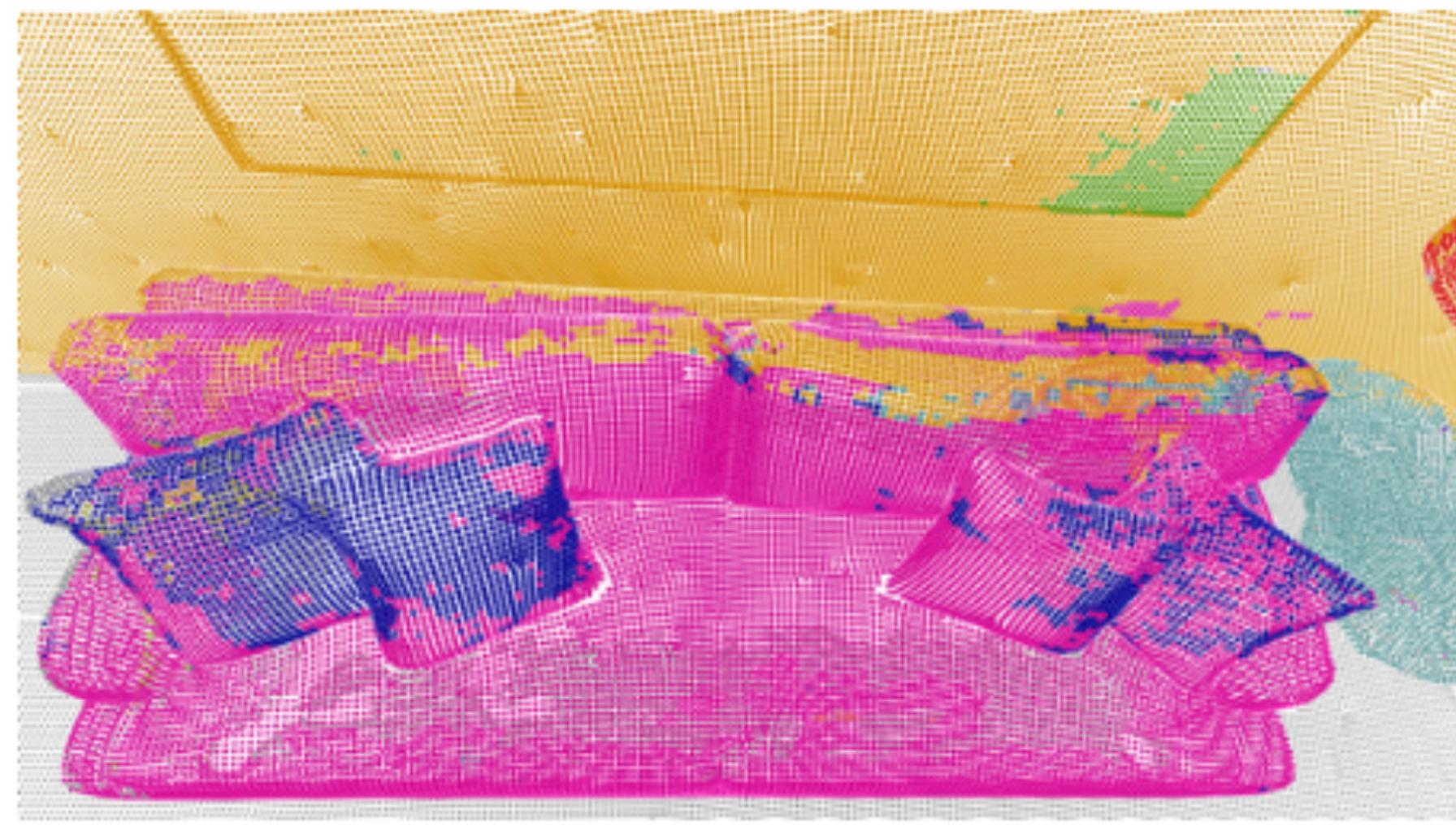
Re-sizing and copy-pasting



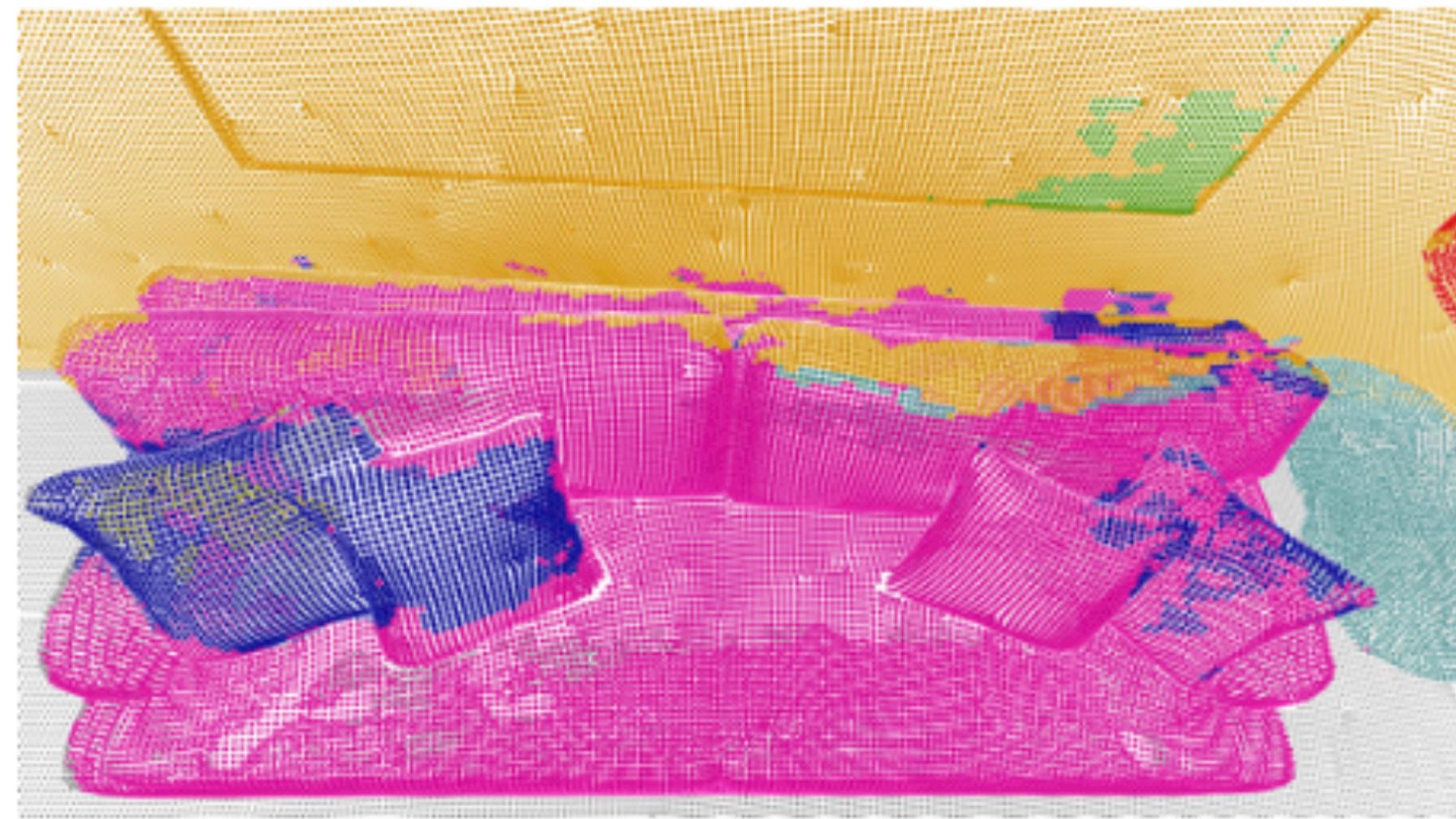
Ground truth



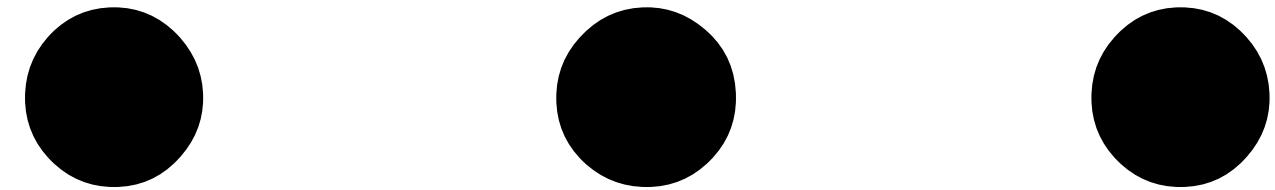
fine MLP



coarse MLP



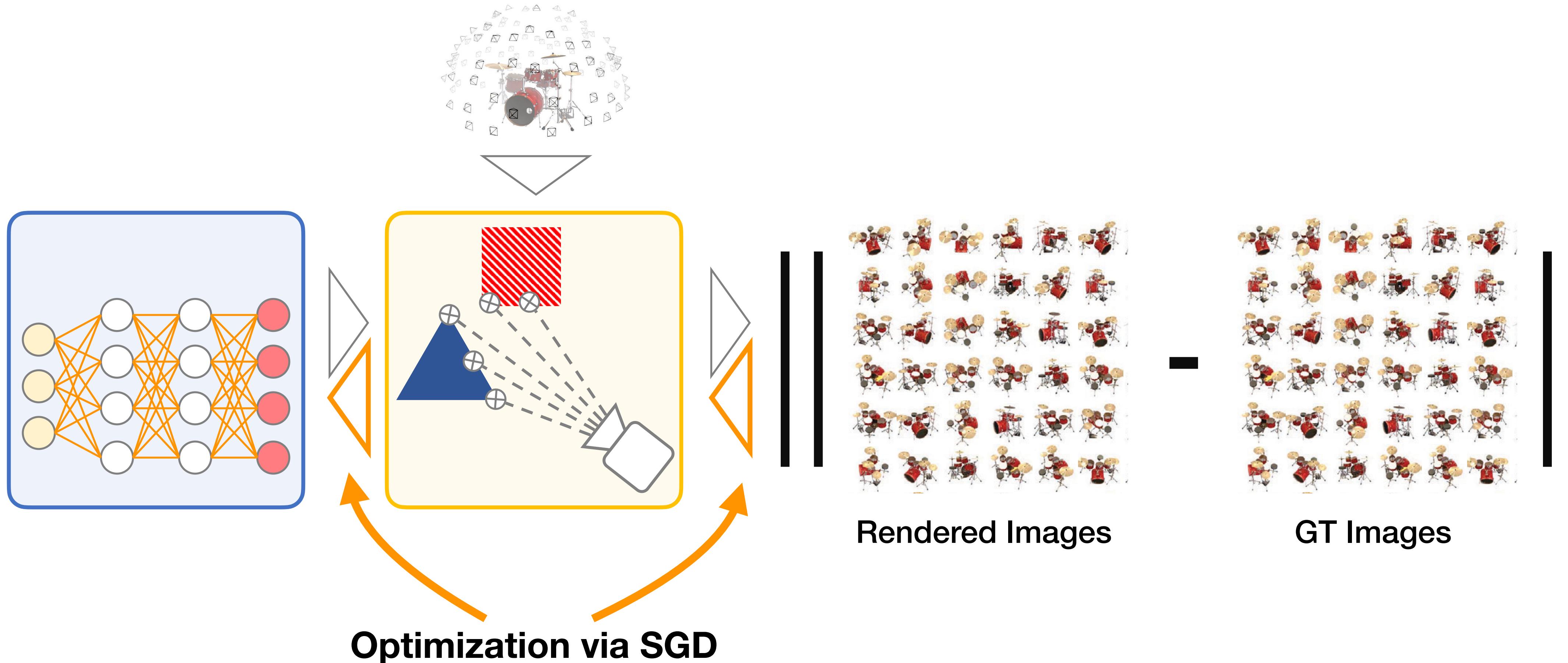
Limitations and Future Work



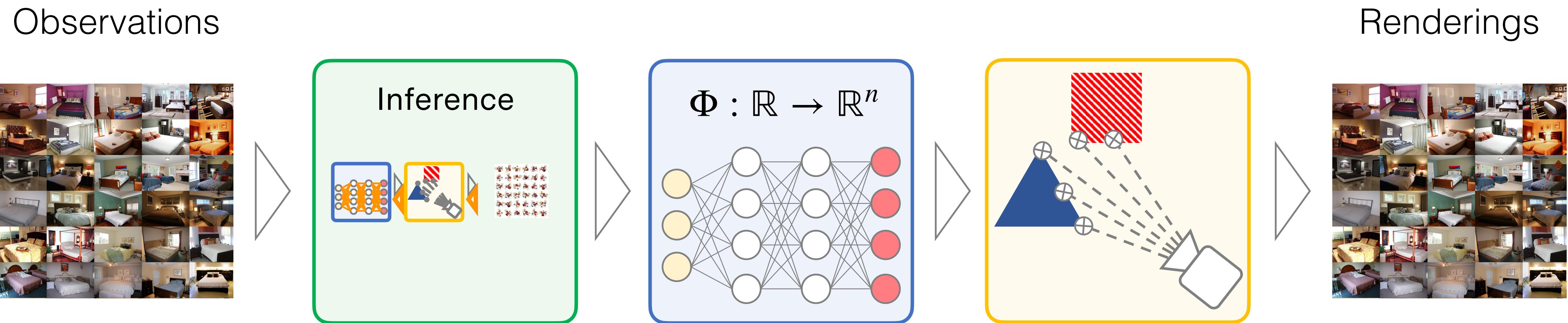
Example Paper Discussion

—END—

Recap: Diff. Rendering

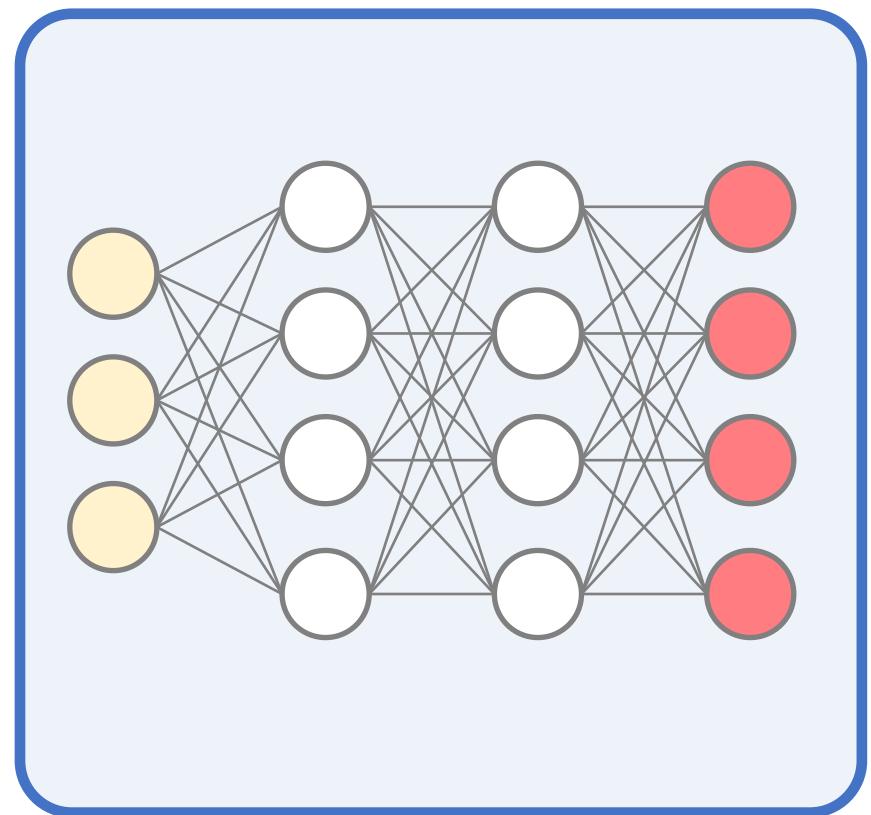


Recap: Diff. Rendering is inference

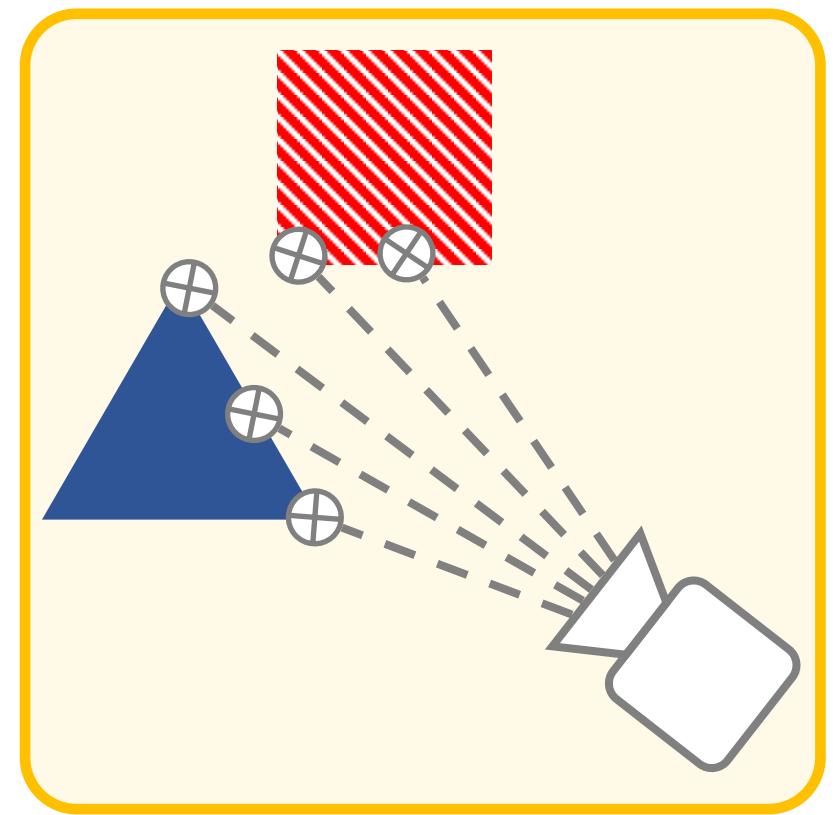


We inferred “hidden/latent” variables 3D appearance and geometry, given “observable” variables RGB images and camera poses.

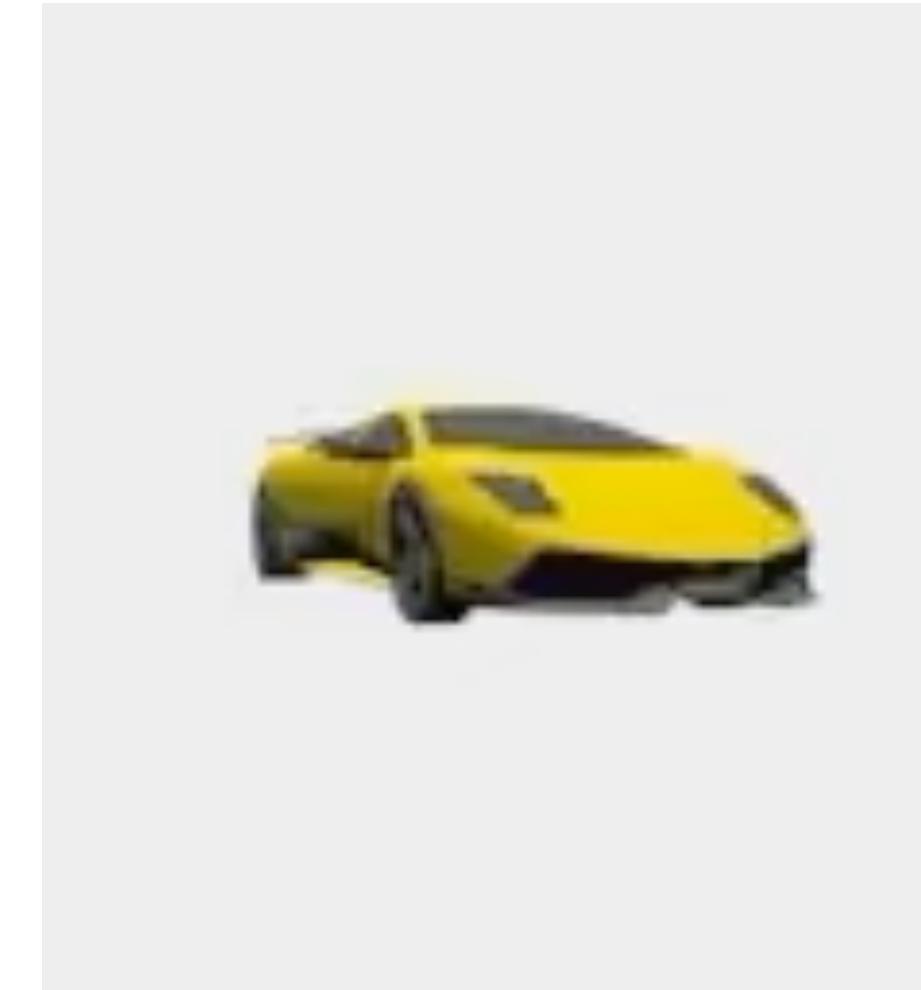
Recap: What if observations don't constrain scene representation?



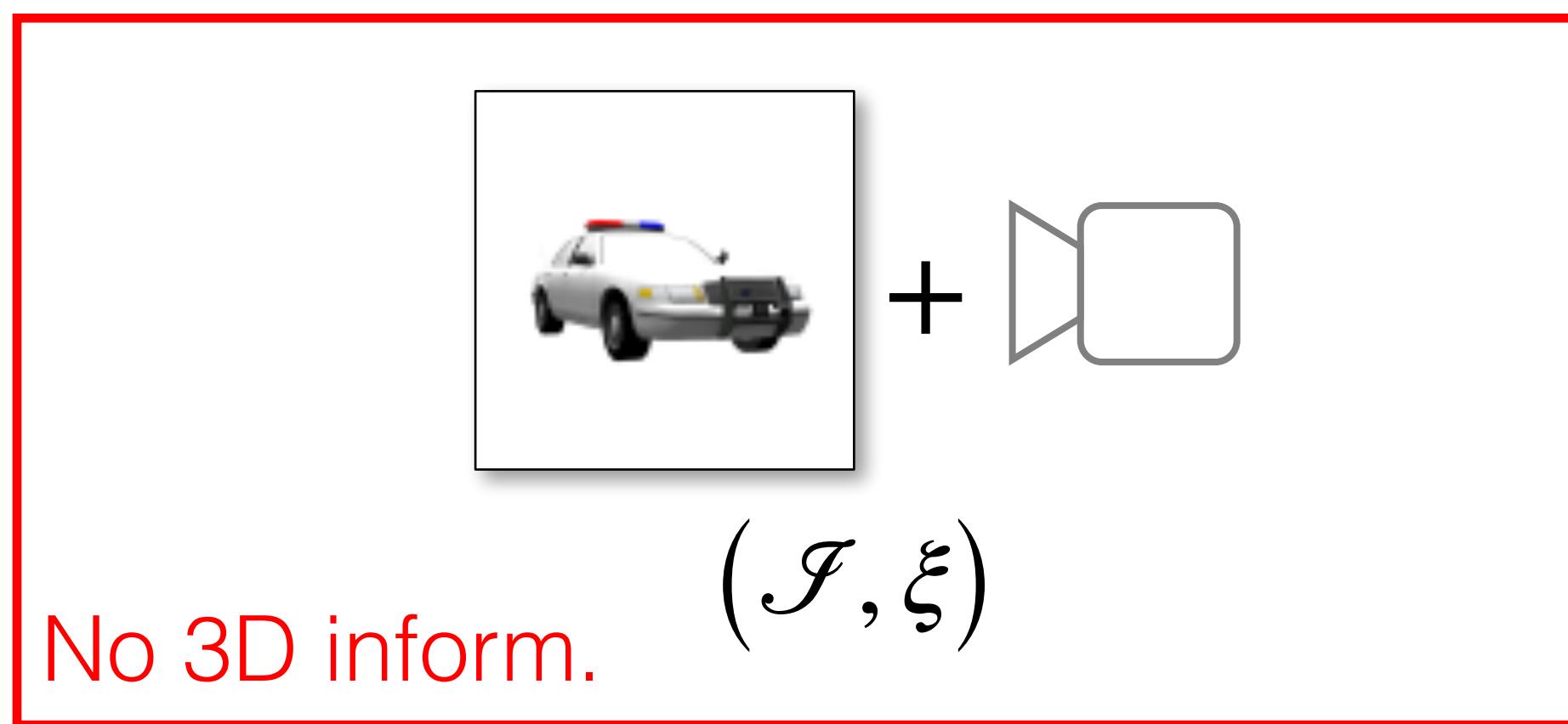
Srn_ϕ



Render_θ



Input



$$\underset{\phi, \theta}{\operatorname{argmin}} \| \text{Render}_\theta(\text{Srn}_\phi, \xi) - \mathcal{I} \|$$



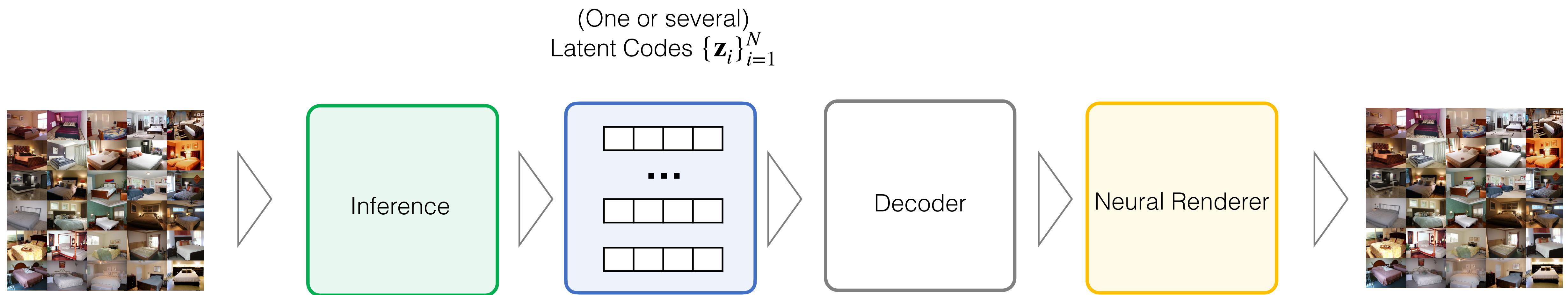
Normal map

RGB

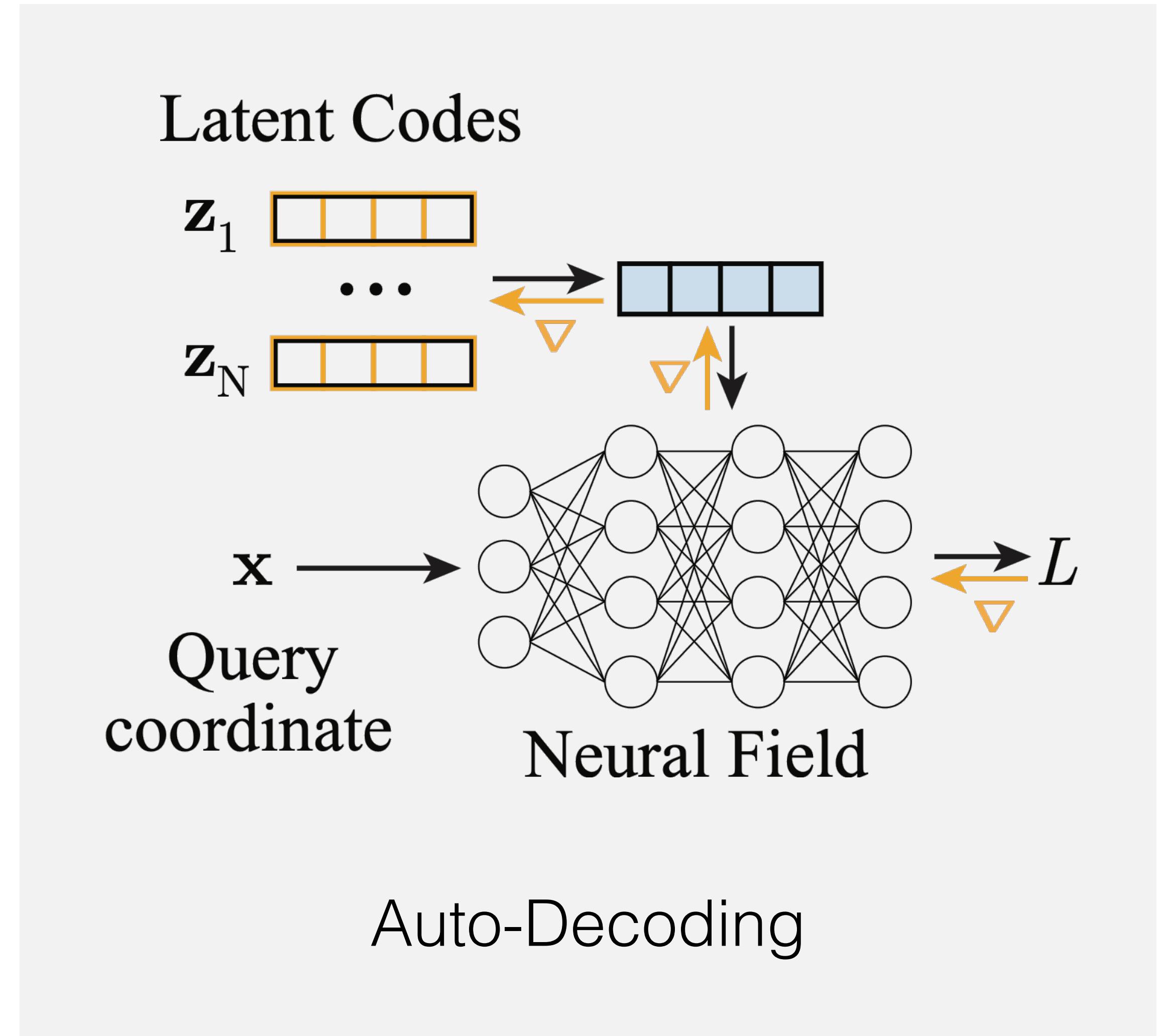
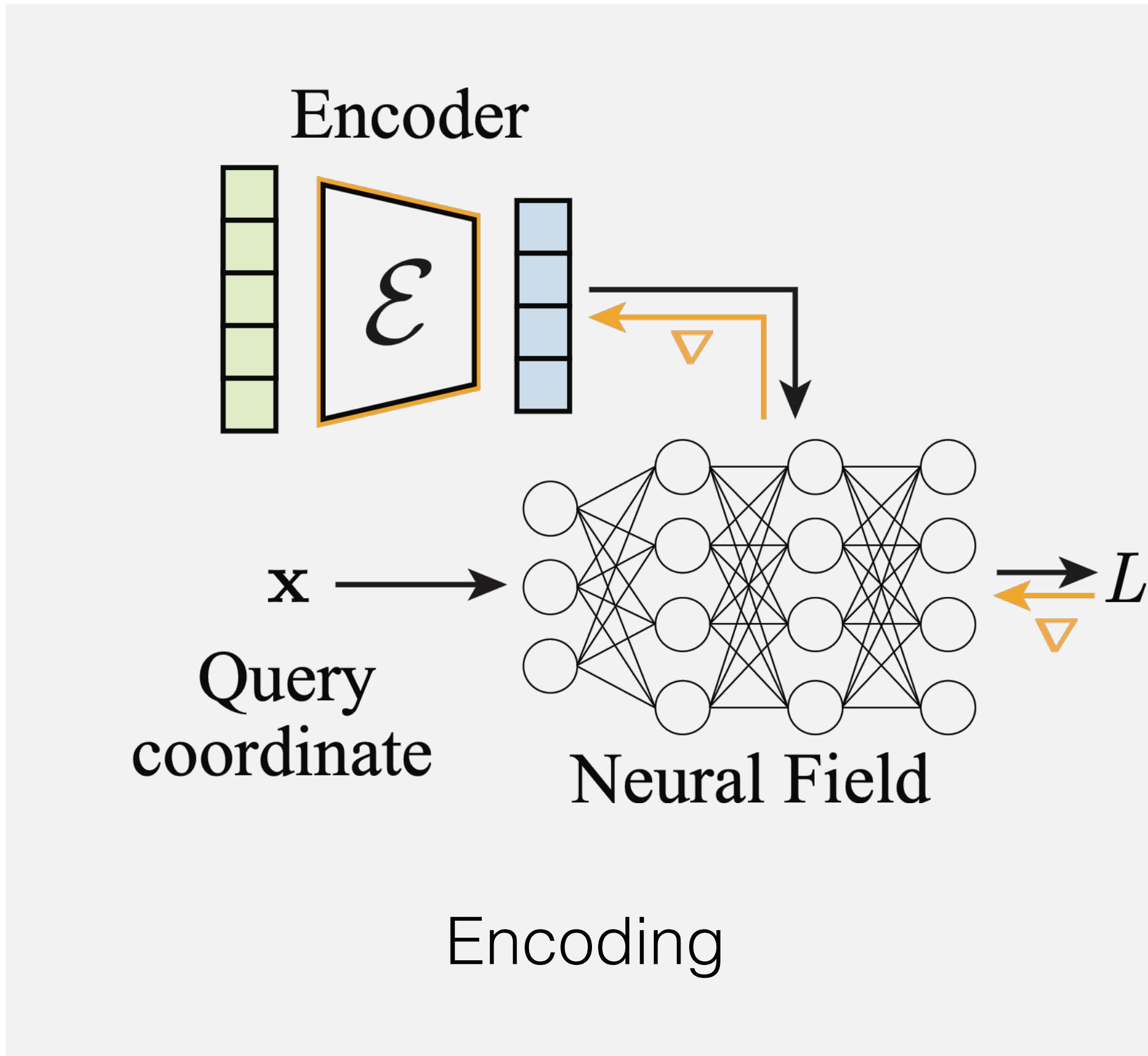


GT

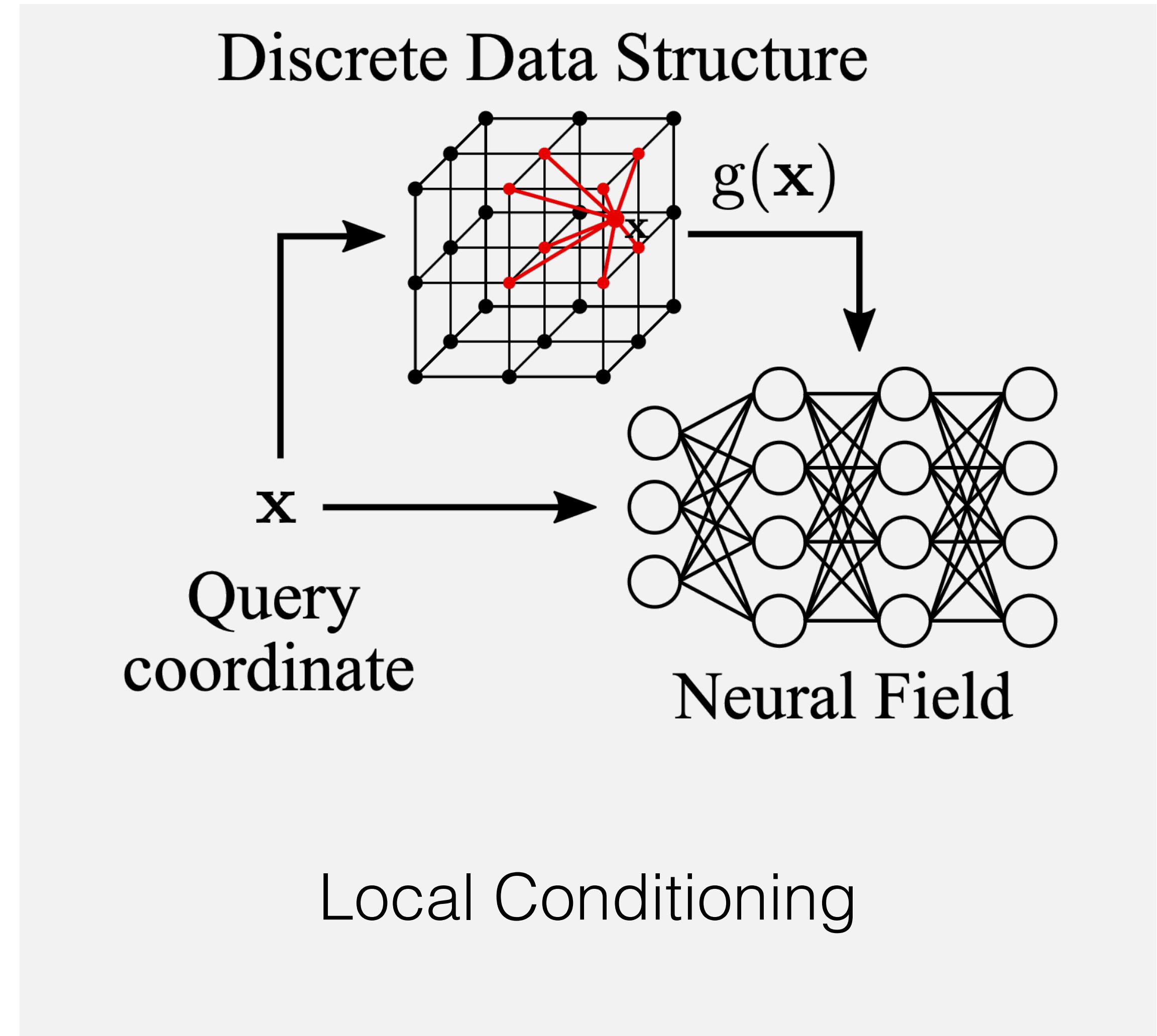
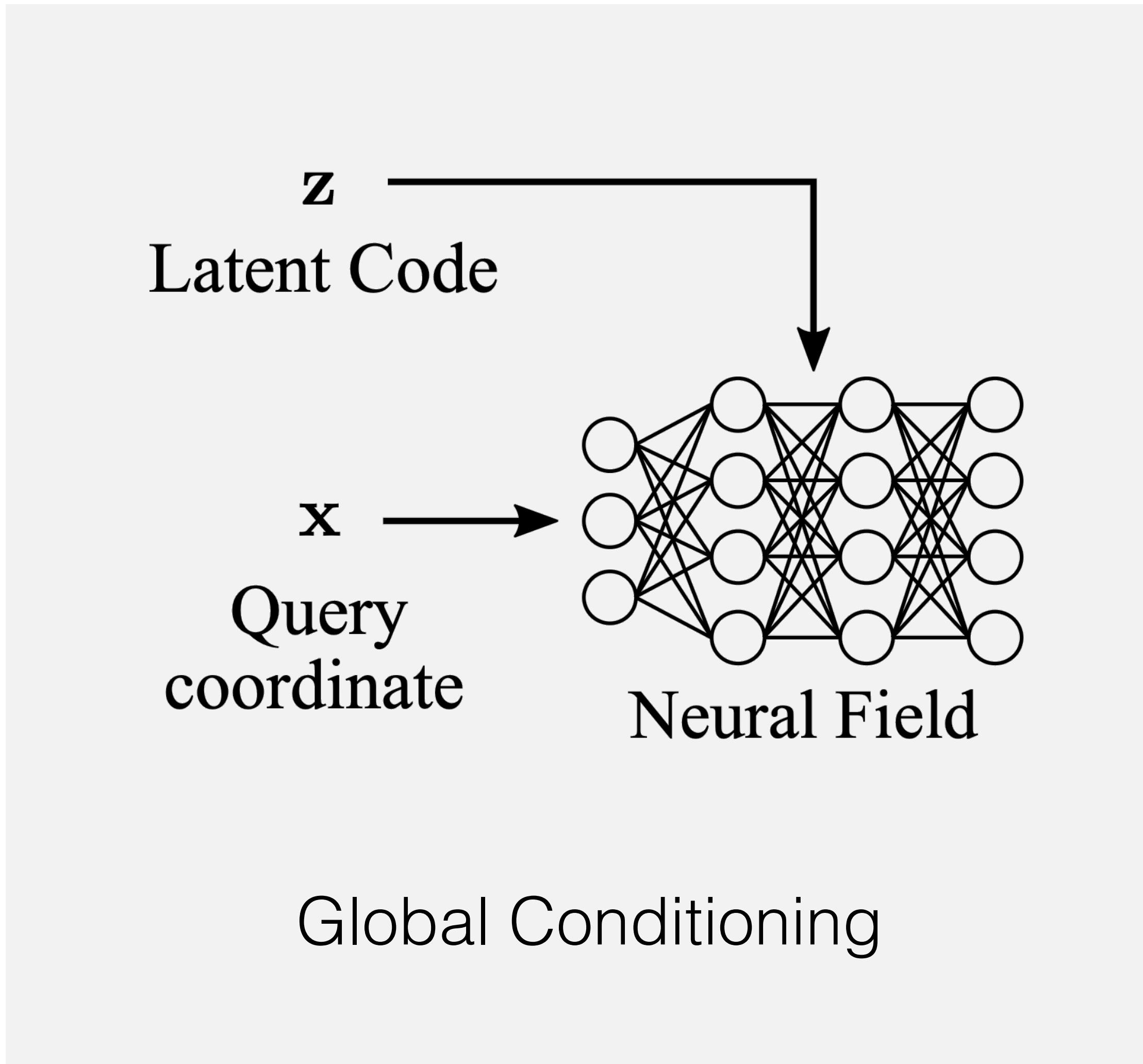
Recap: General Framework for auto-encoding based Scene Representation Learning



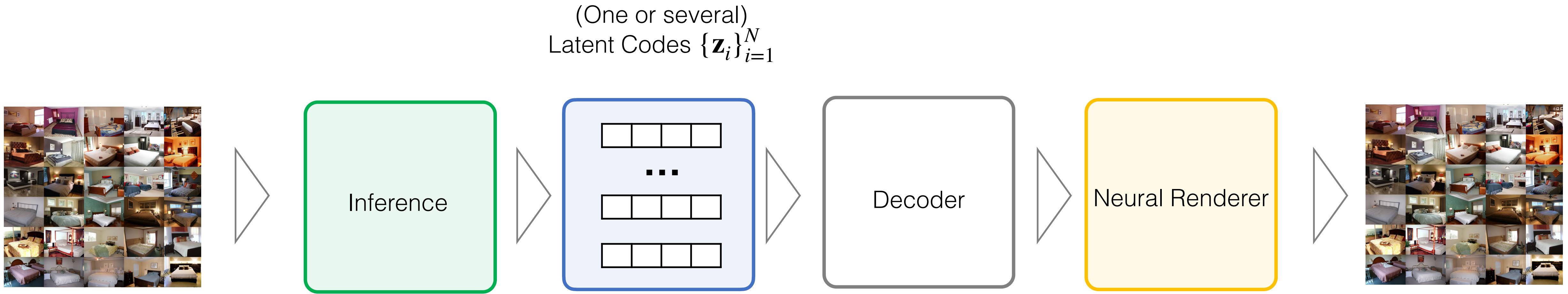
Recap: Encoding & Auto-Decoding



Recap: Global Conditioning and Local Conditioning



Today: Advanced Inference Topics



Amortized Neural Rendering (=Light Fields)

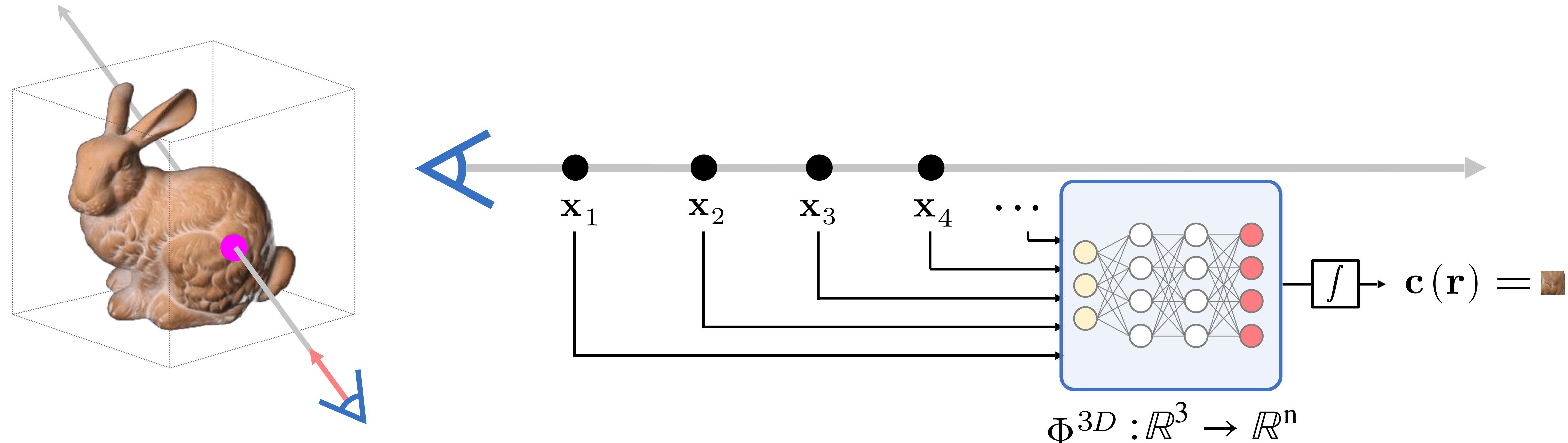
Conditioning with Transformers

Operating on Neural Scene Representations

Geometric Deep Learning

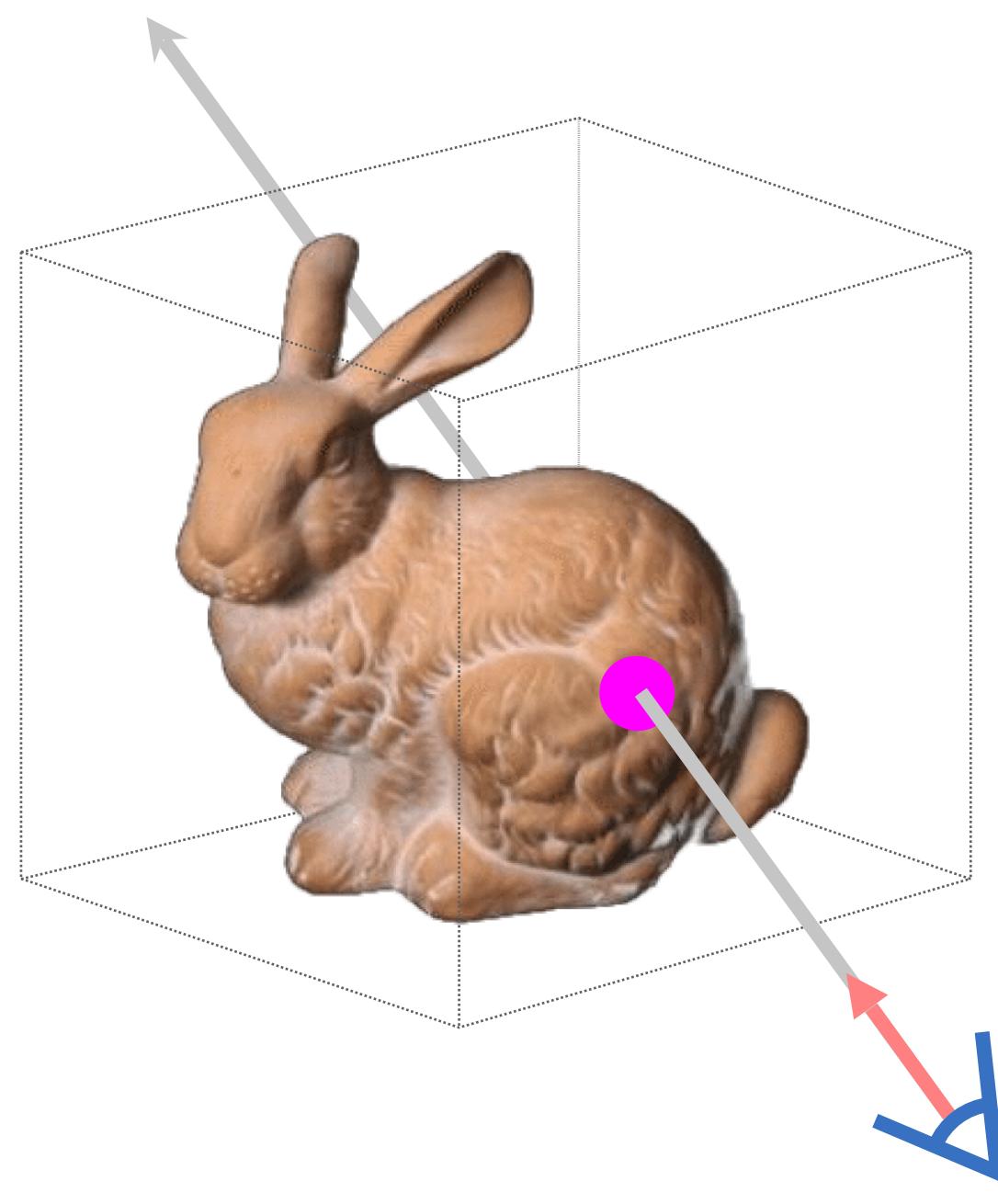
...

General structure of Neural Renderers for 3D-structured Representations



Hundreds of samples per ray.

Time- and memory-intensive.



A

x_1

x_2

x_3

x_4

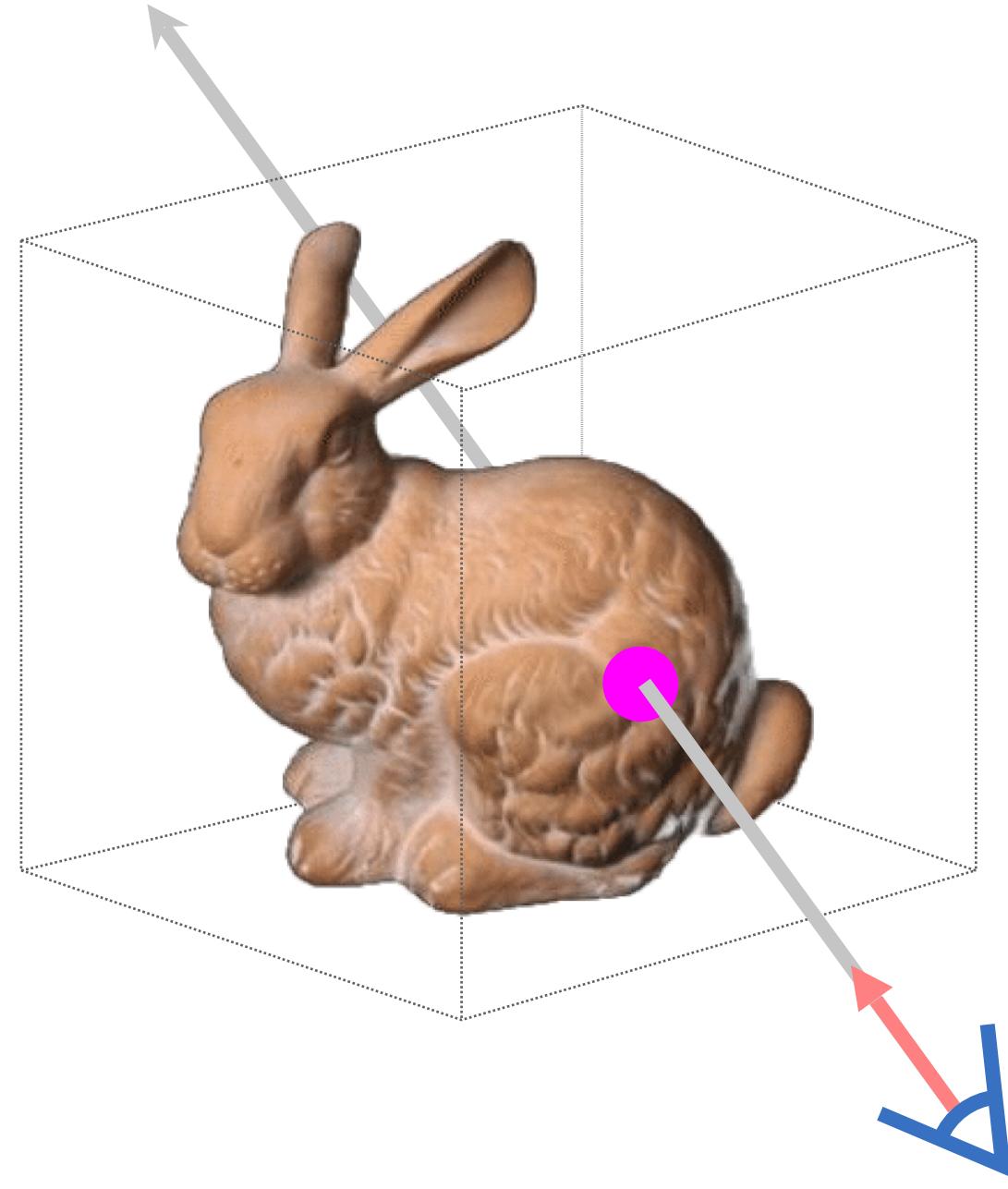
...

Light Field

$$\Phi^{3D} : \mathbb{R}^3 \rightarrow \mathbb{R}^n$$



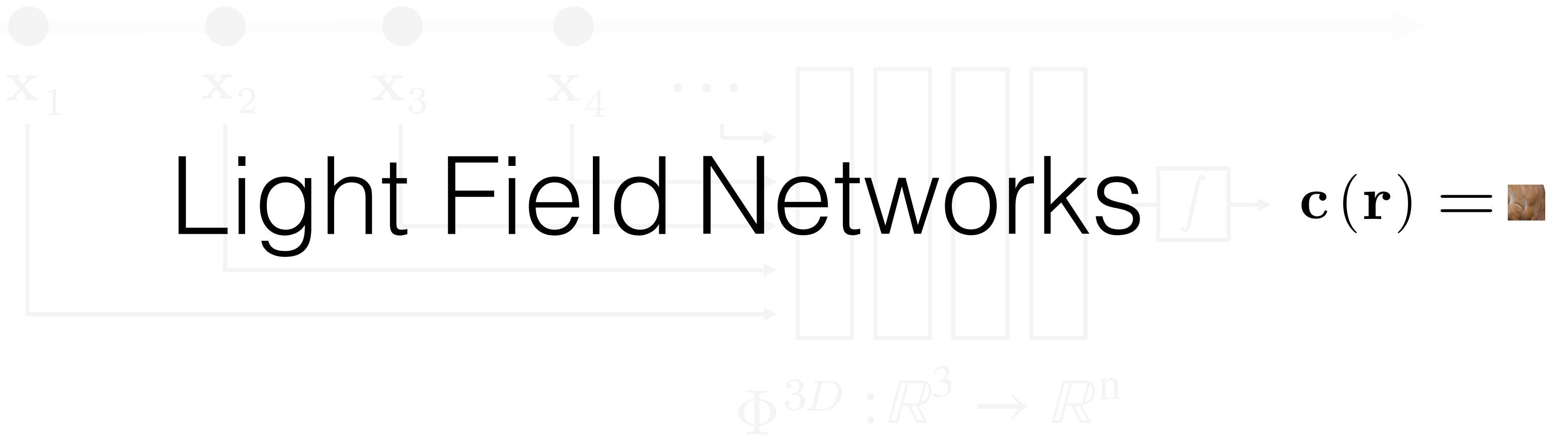
$$c(\mathbf{r}) =$$



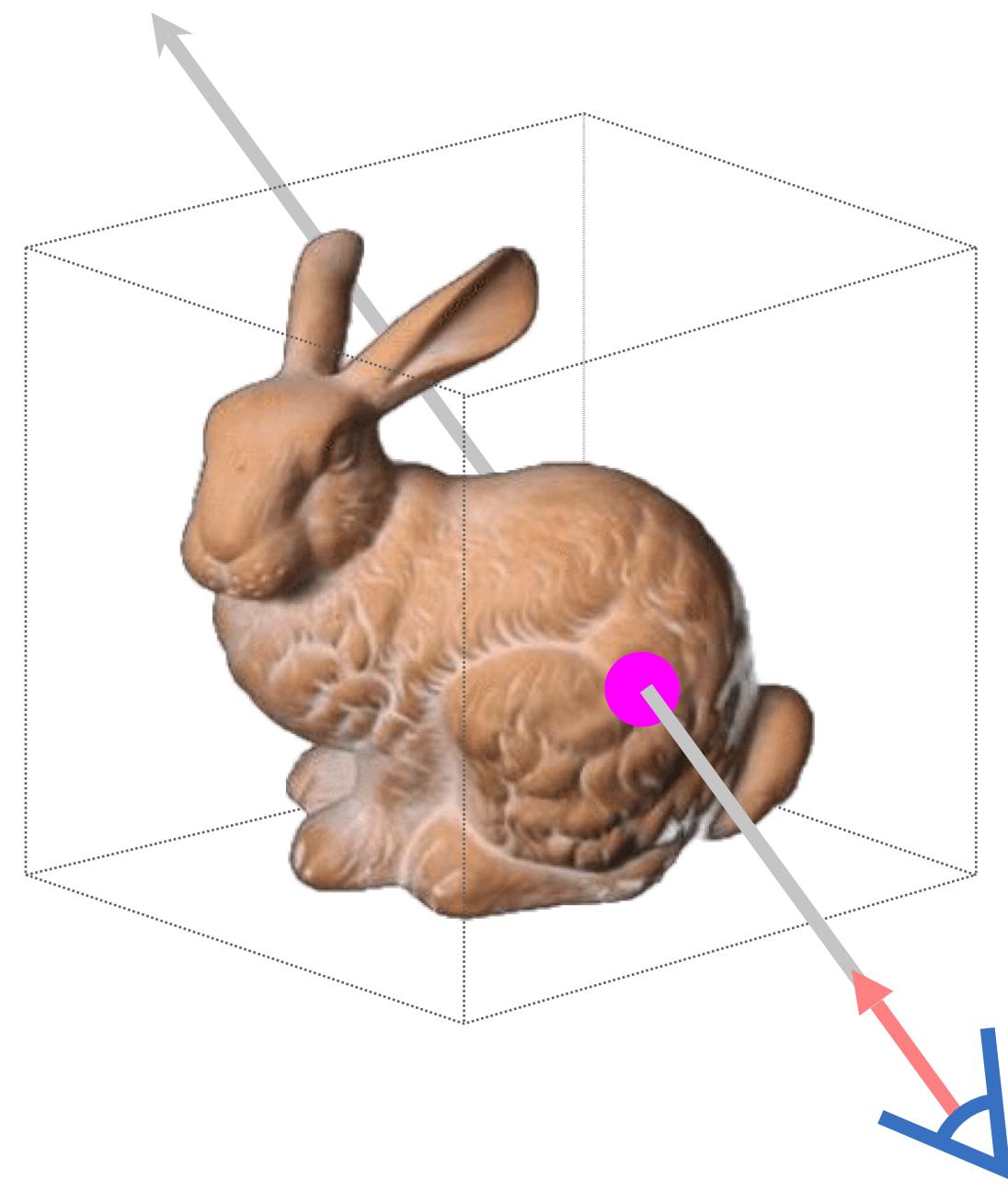
A

x_1

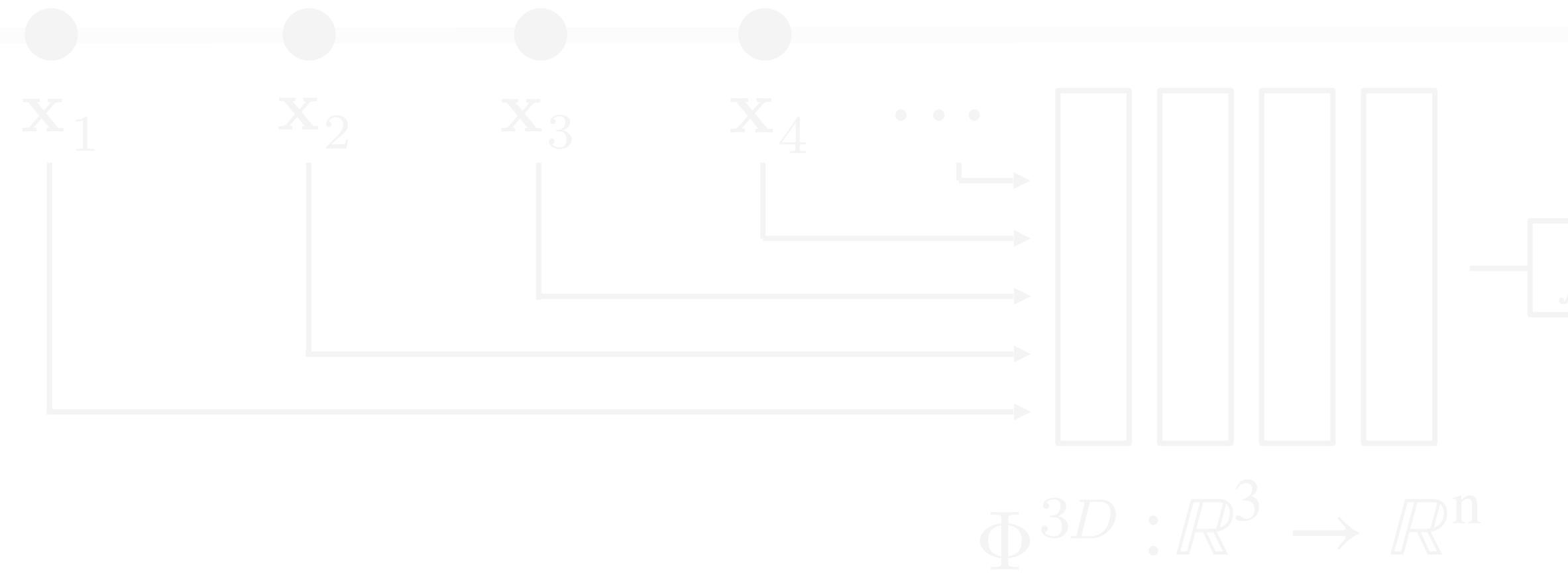
Light Field Networks



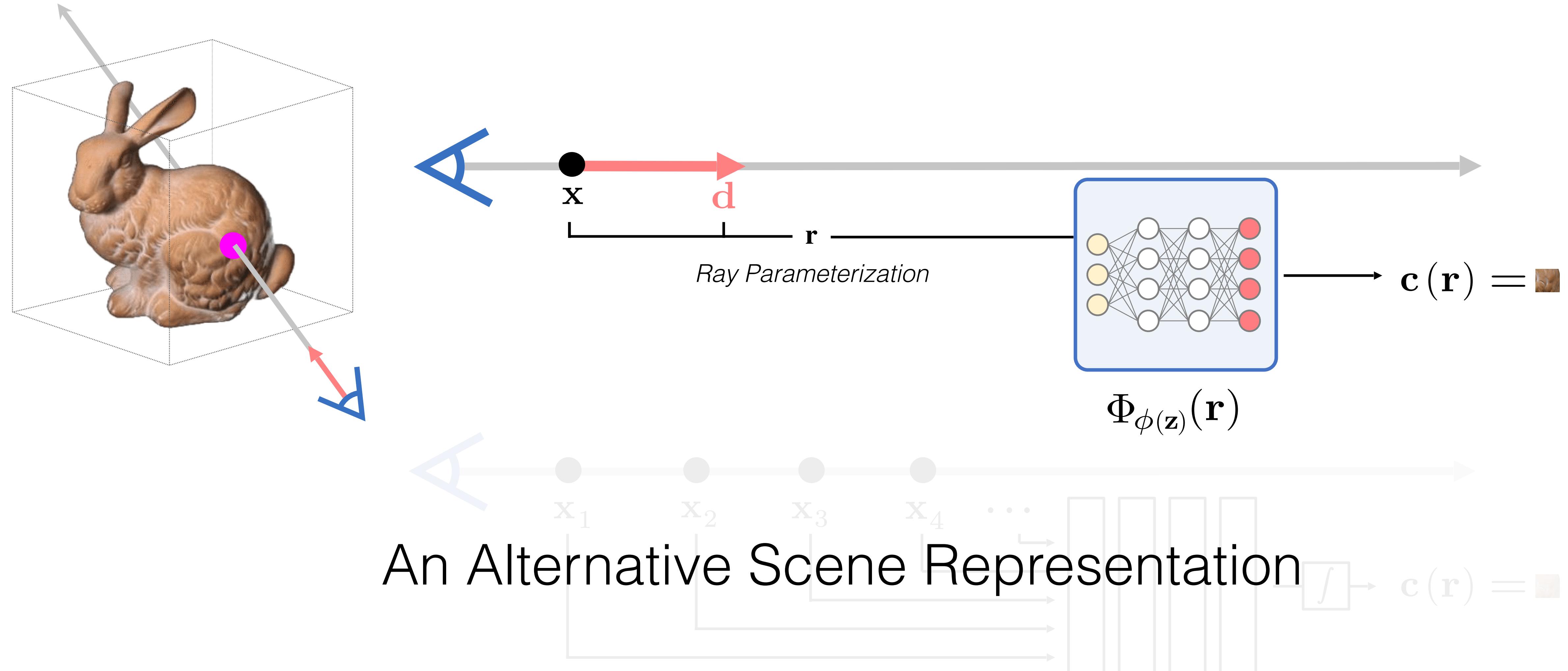
Light Field Networks



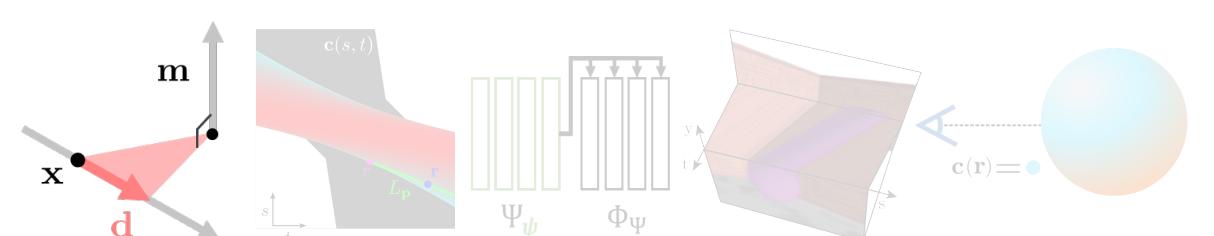
A



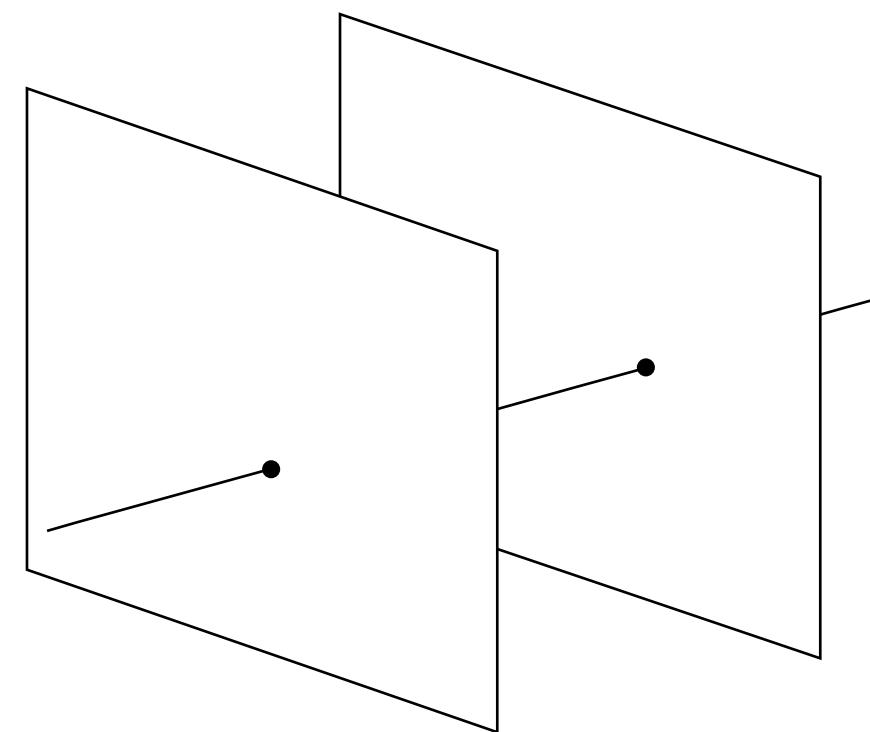
Light Field Networks



Conventional Light Field Parameterizations

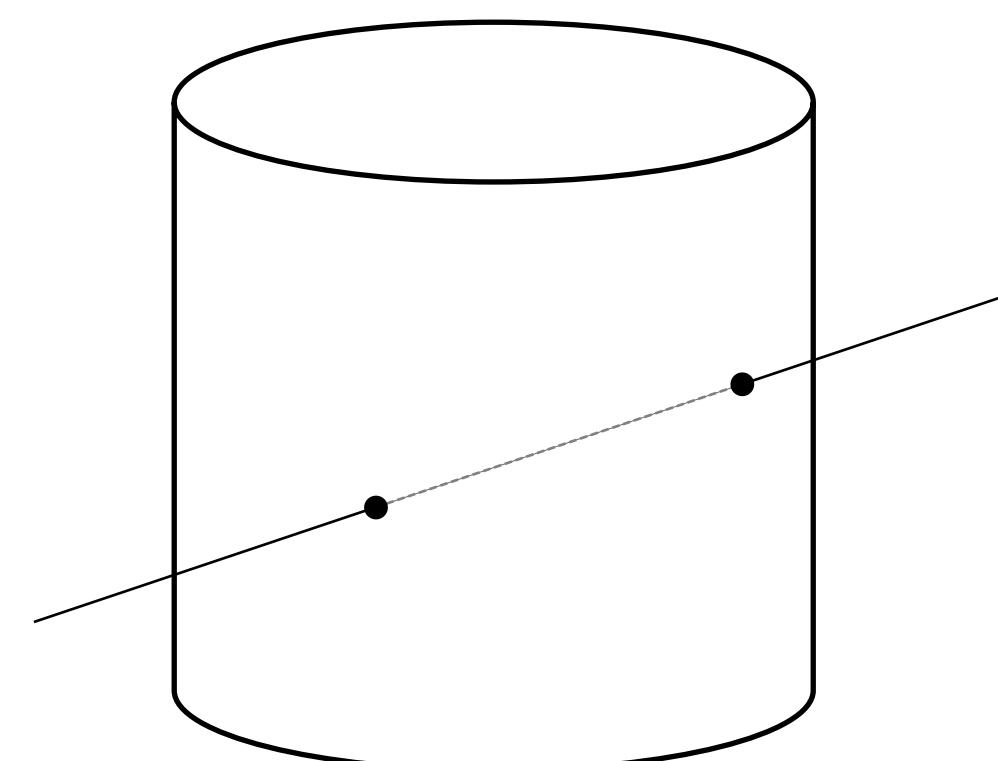


Two-Plane



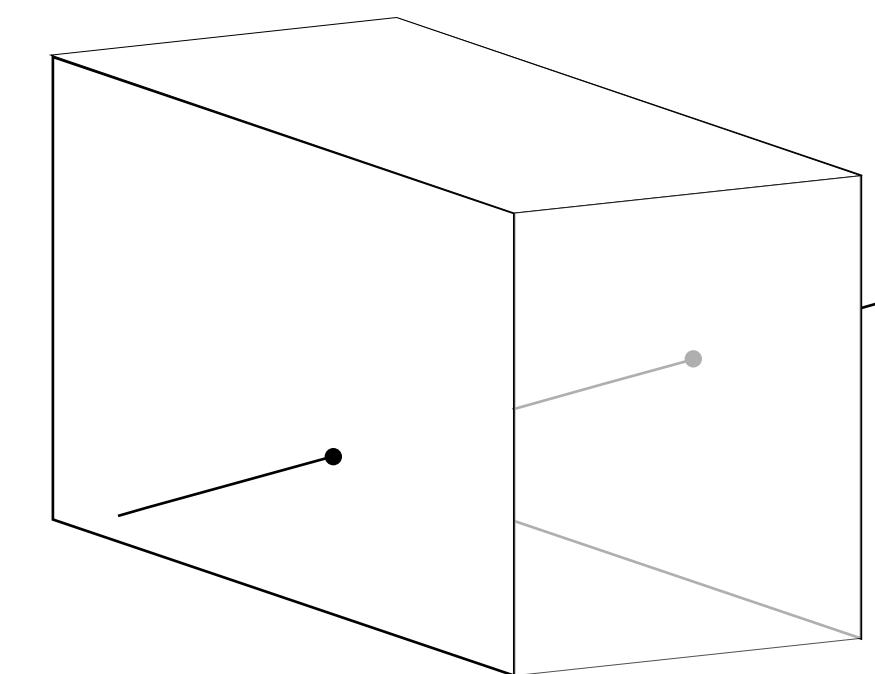
Not 360°

Cylindrical



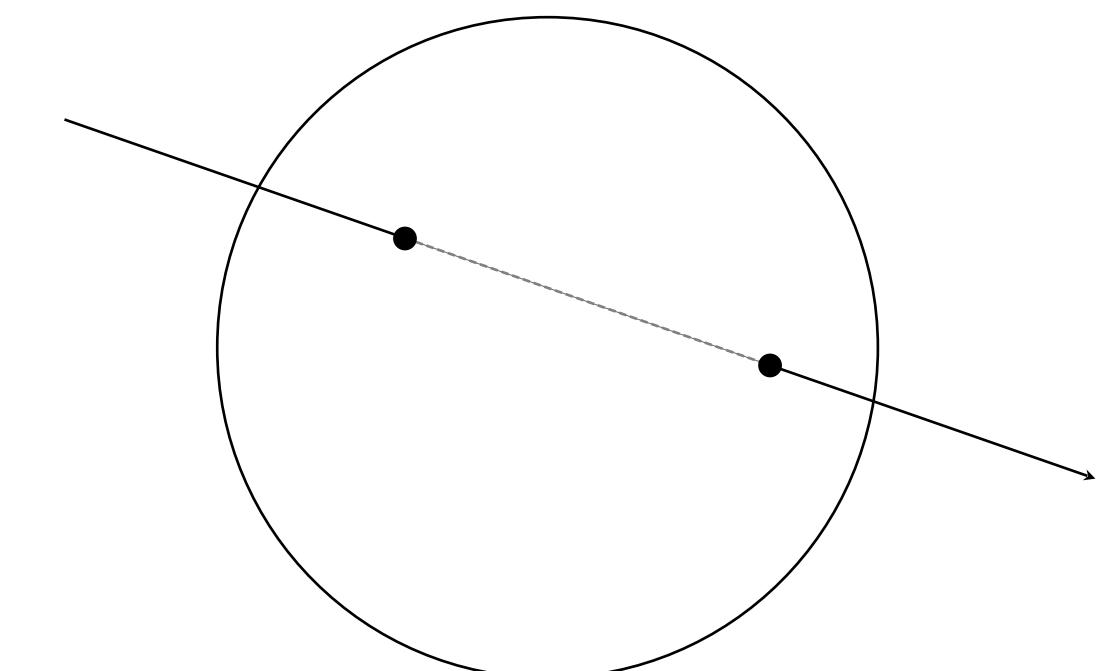
Not 360°

Lumigraph



Not Continuous

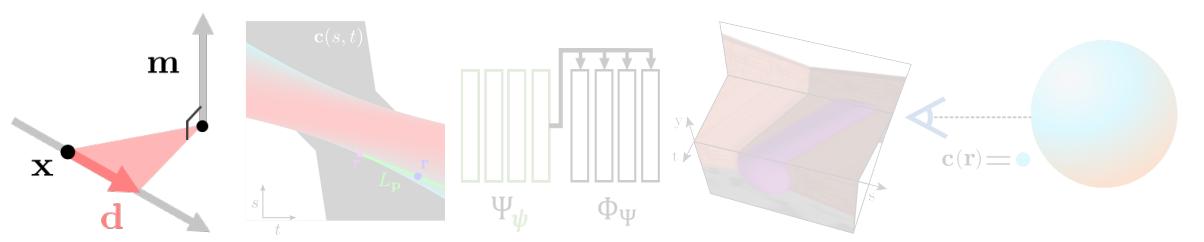
Two-Sphere

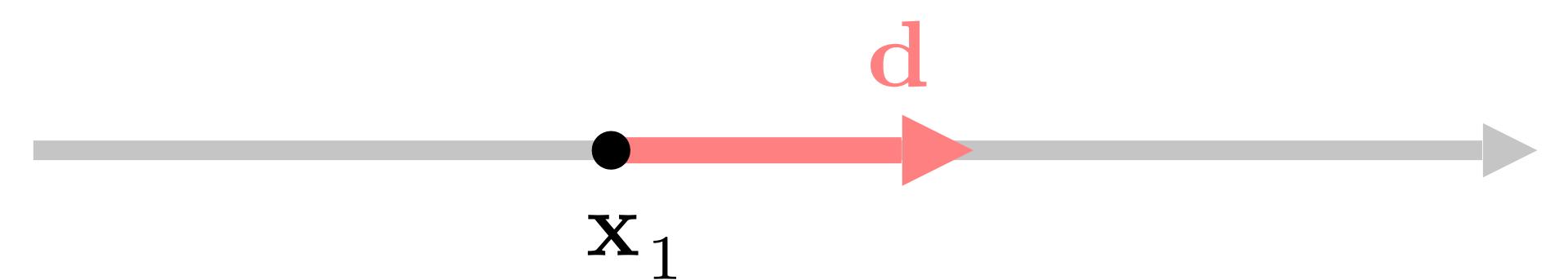


Bounded Scenes

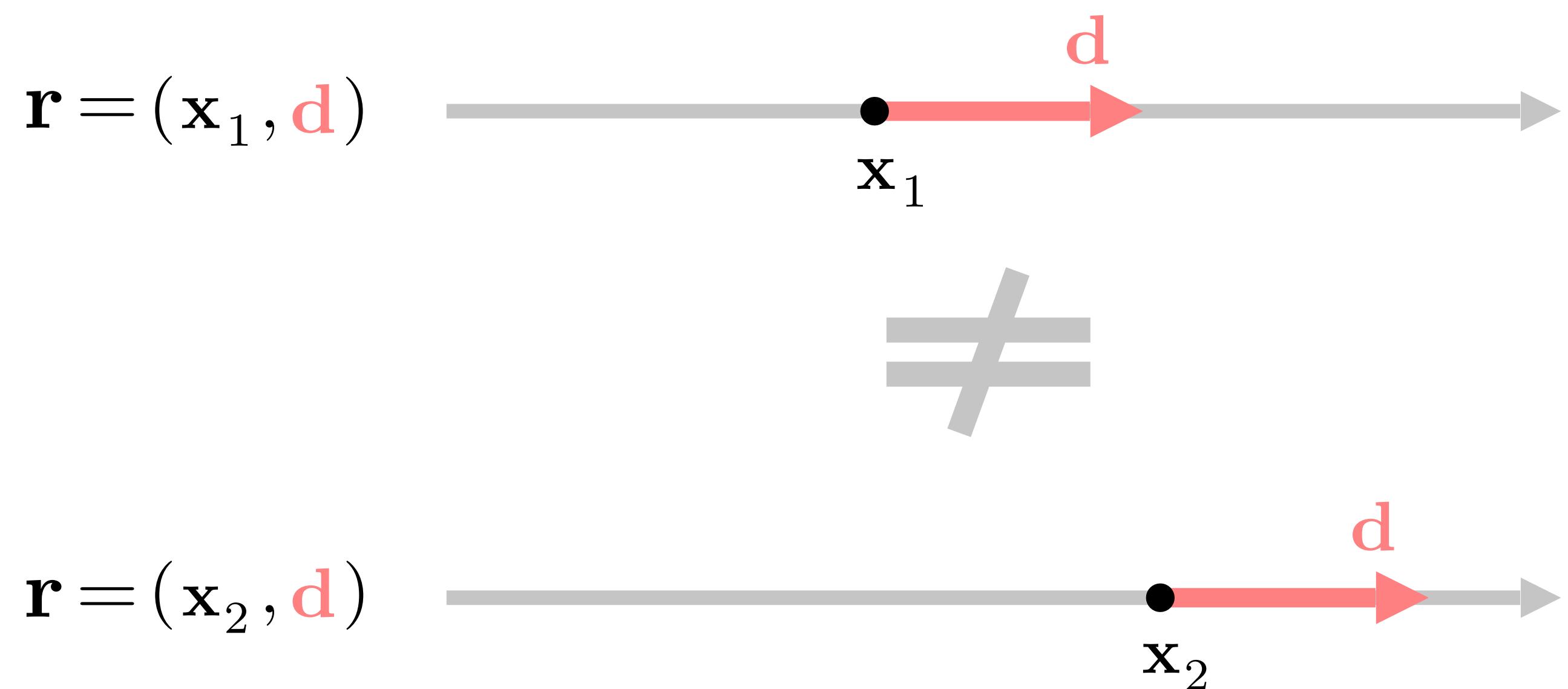
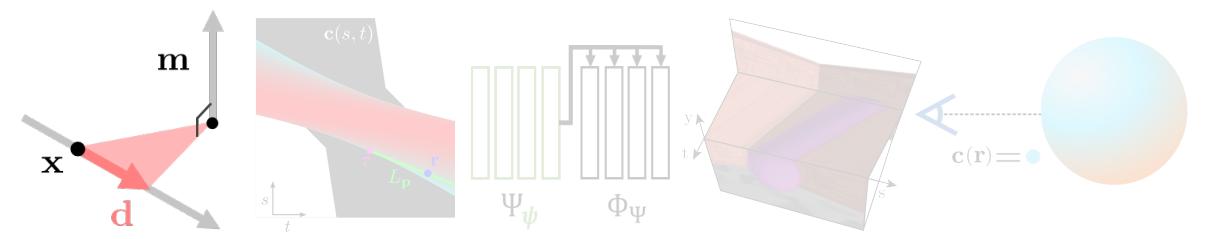
Difficult to use as a complete scene representation

“Point-direction” coordinates



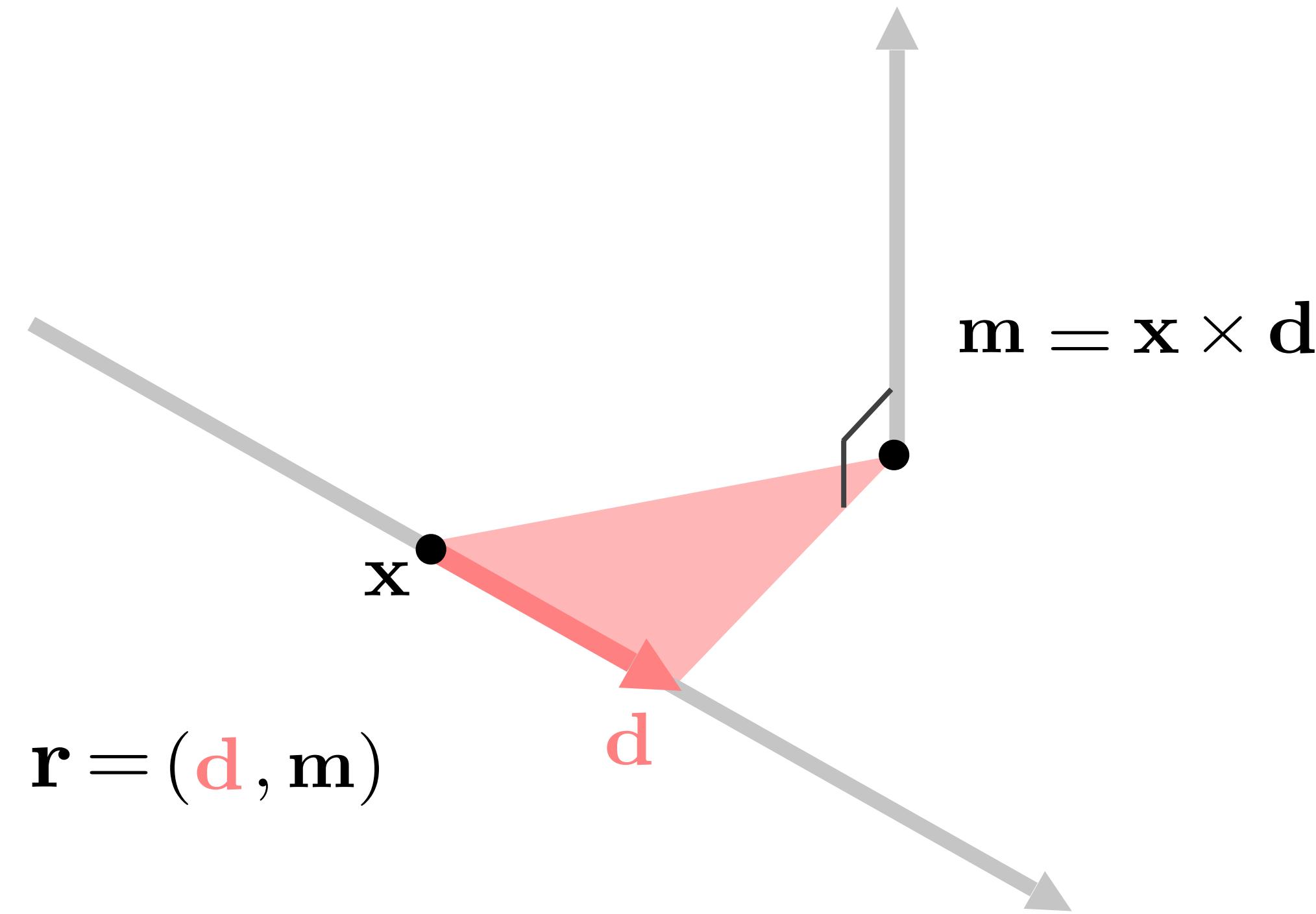
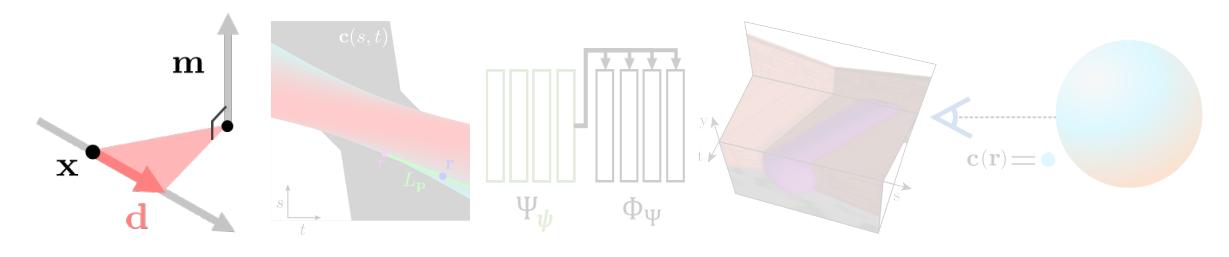
$$\mathbf{r} = (\mathbf{x}_1, \mathbf{d})$$


“Point-direction” coordinates



Not unique: Same ray, two different coordinates.

Plücker coordinates



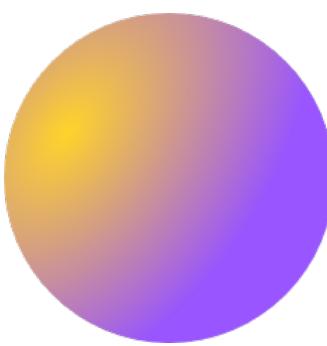
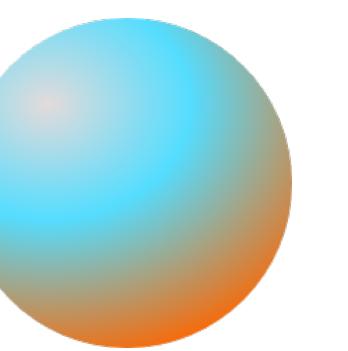
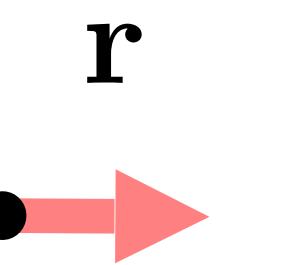
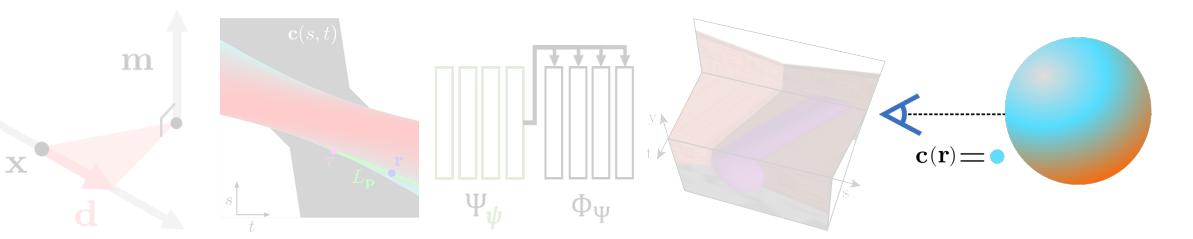
Invariant to choice of x .

Parameterize all rays without special cases.

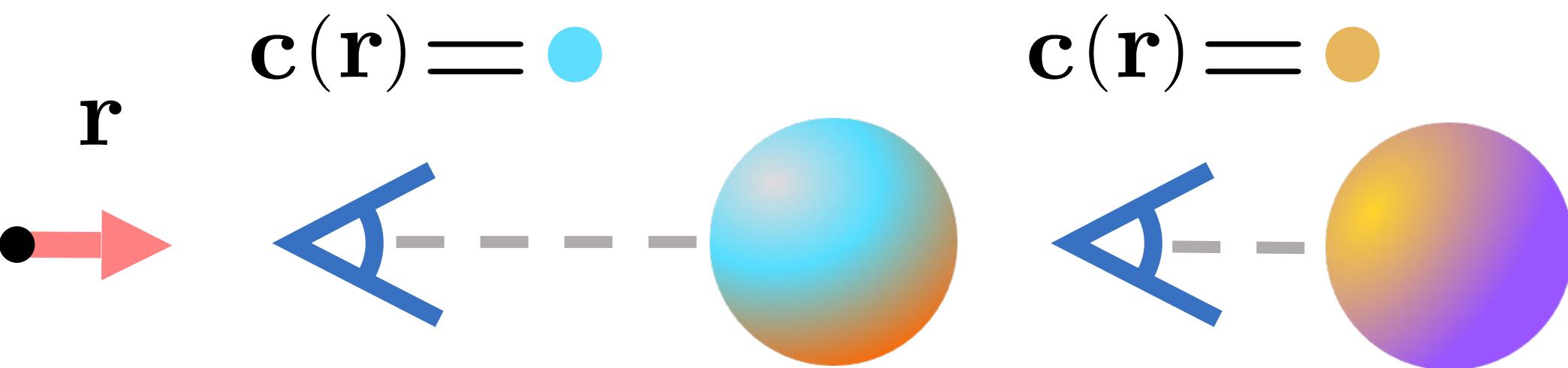
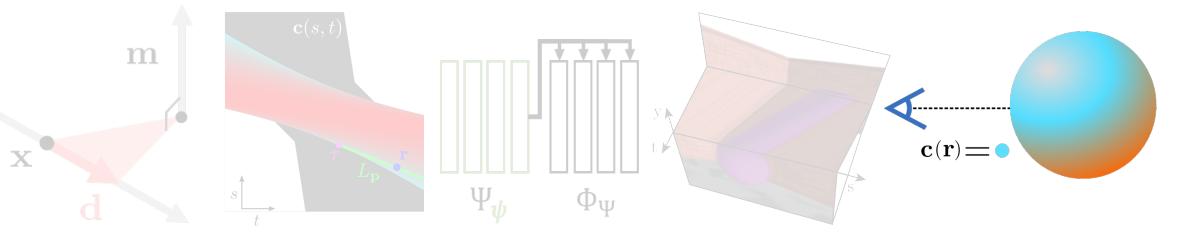
Impractical for discrete representations, since $\mathbf{r} \in \mathbb{R}^6$.

Invariant to choice of x .

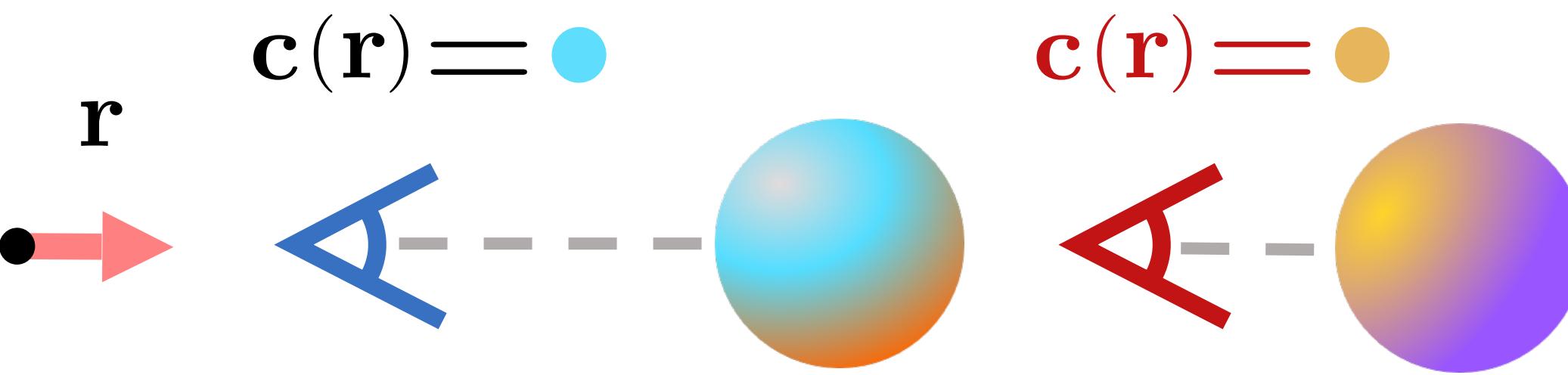
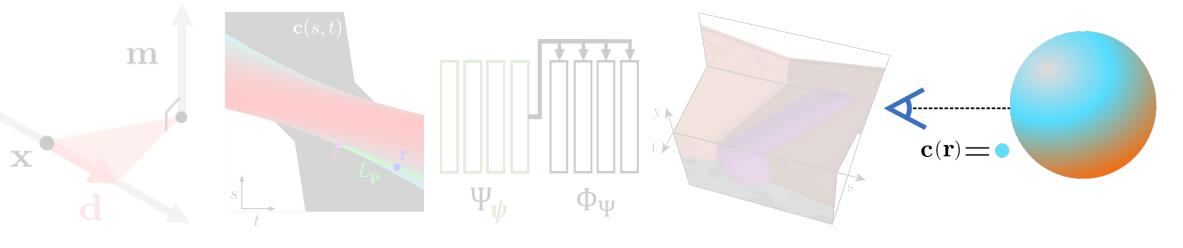
Impact of ray parameterization



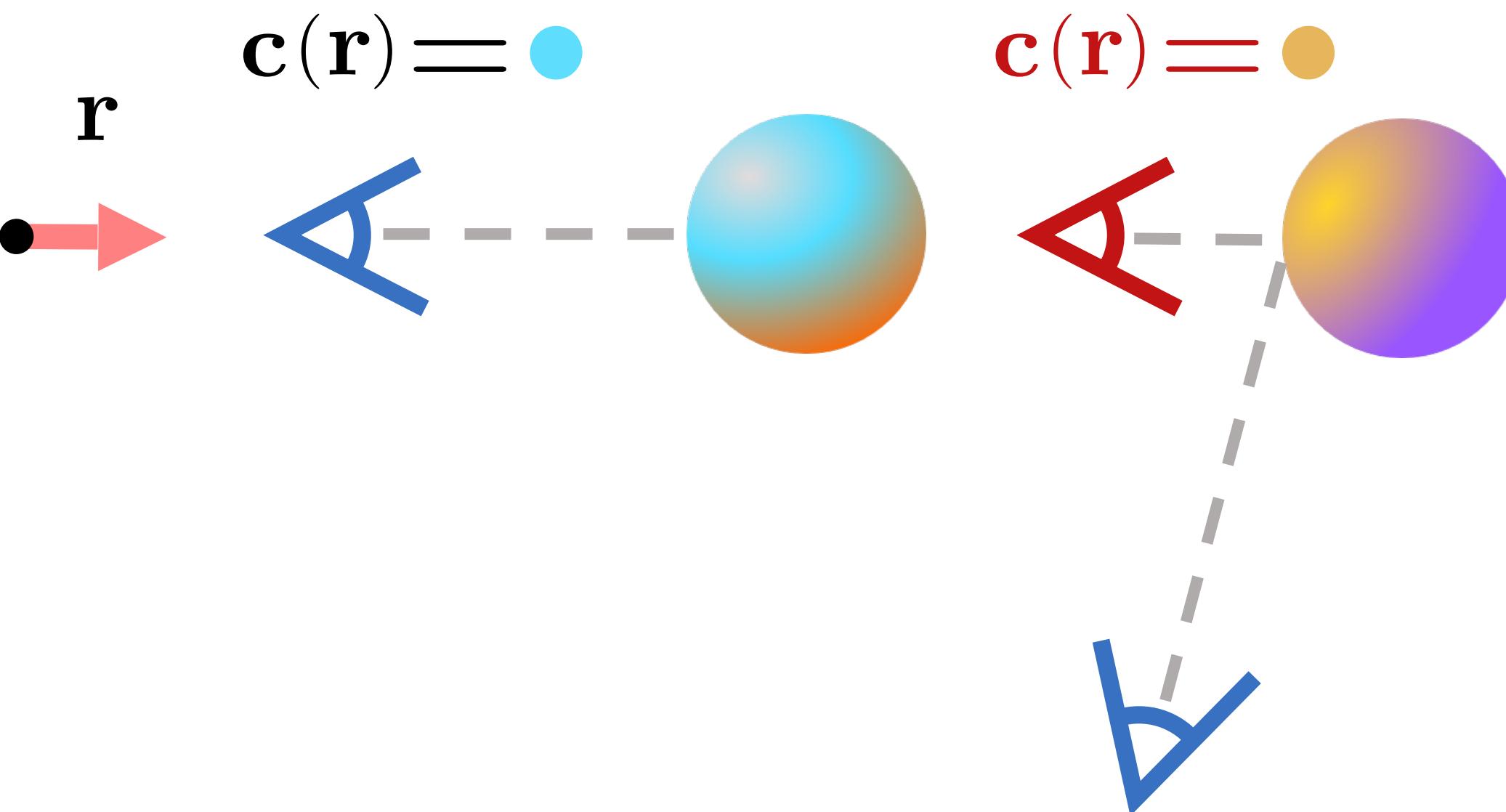
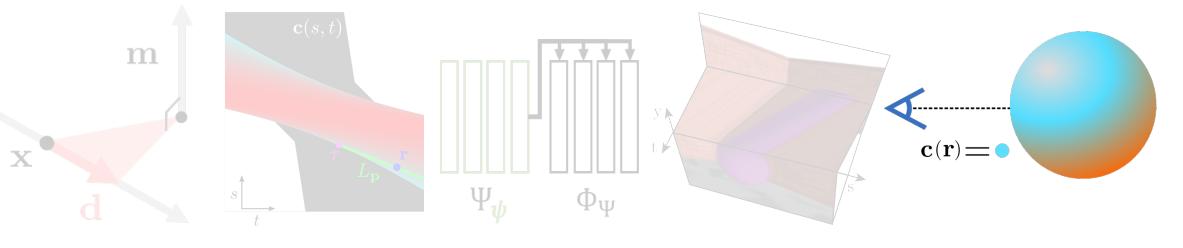
Limitations



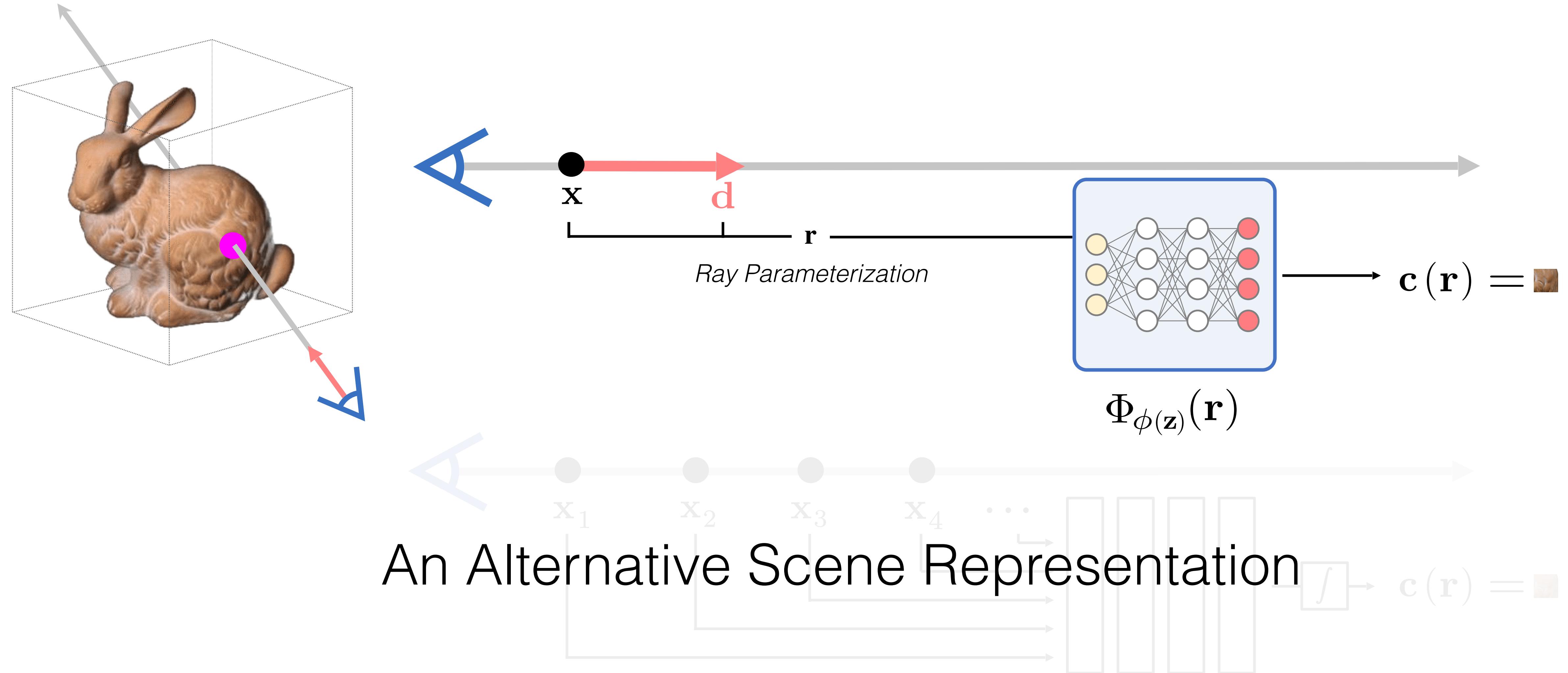
Limitations



Limitations



Light Field Networks



Overfitting doesn't work out-of-the-box: No built-in multi-view consistency!

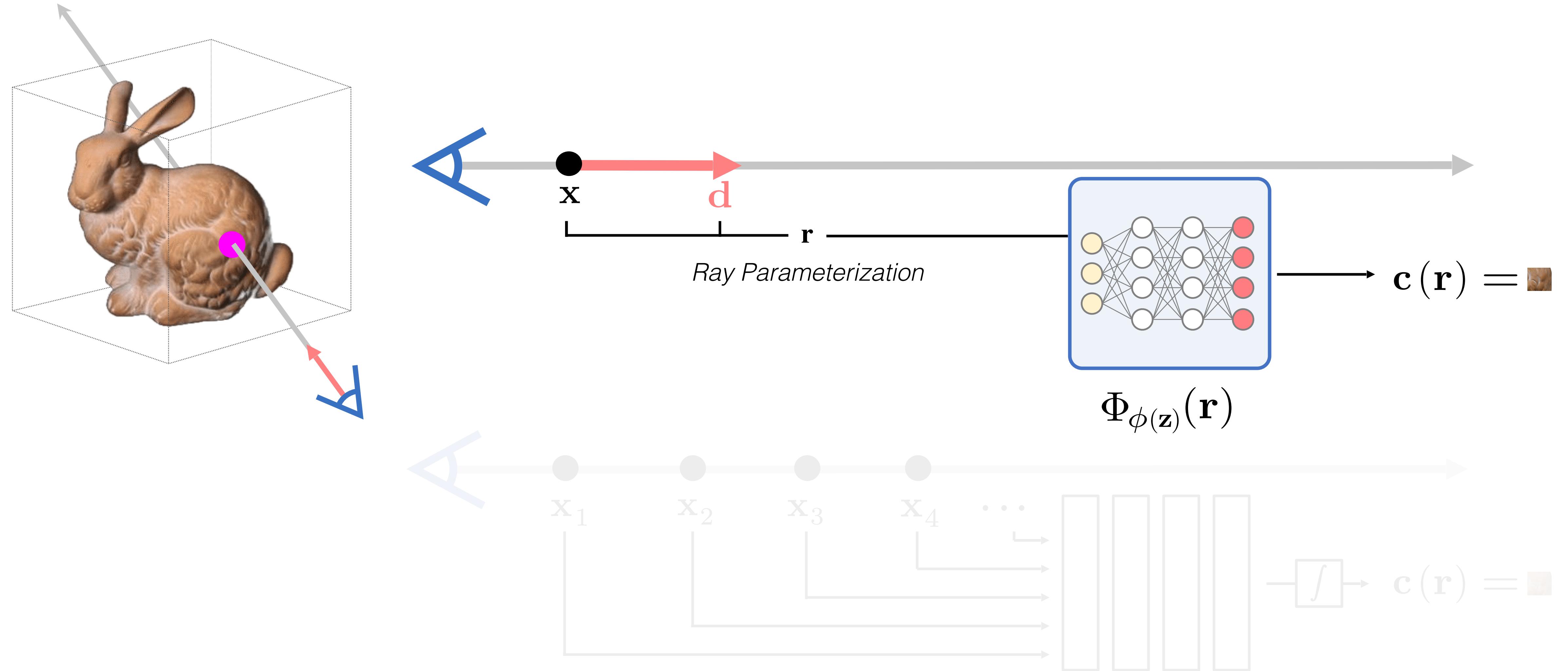


Context Views

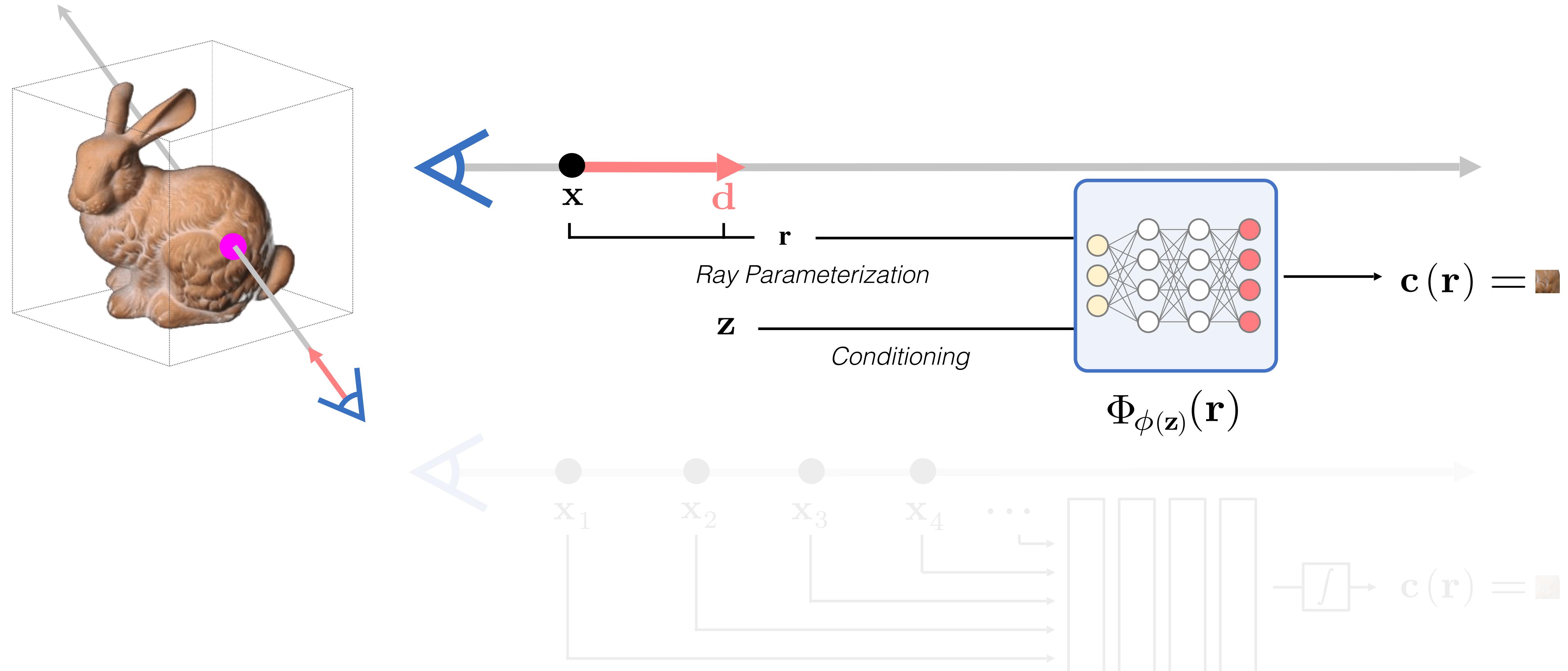


Intermediate Views

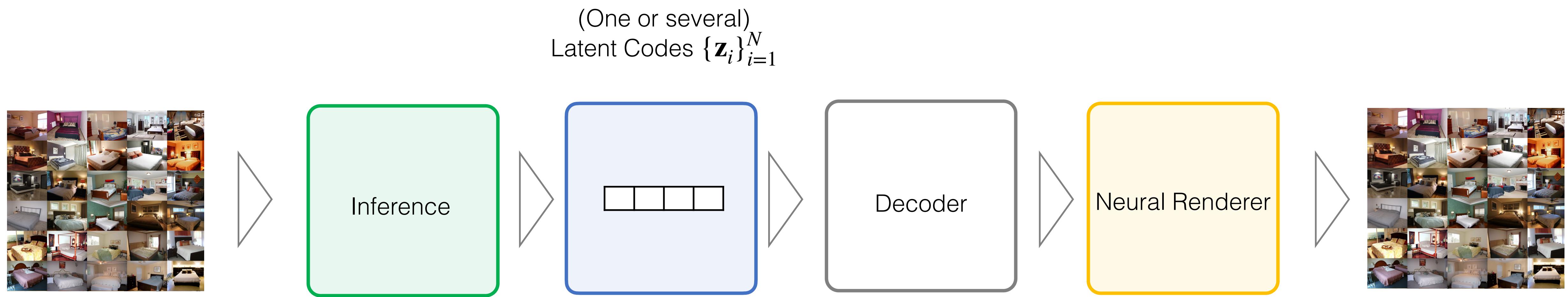
Light Field Networks



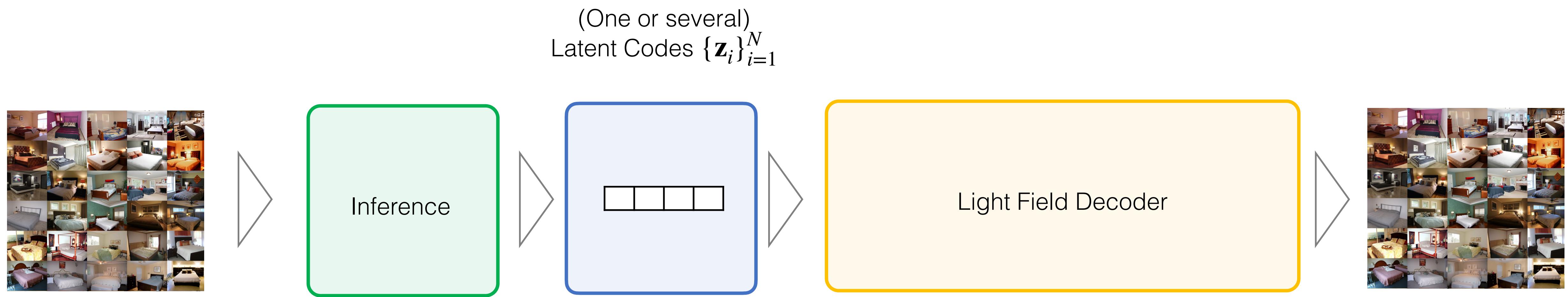
Now: Learn Prior over Light Fields!



Light Field Networks



Light Field Networks



More difficult inference problem, but cheaper & maximally general renderer.

Light Field Networks
500 FPS
1 evaluation per ray



Volumetric Rendering (pixelNeRF)
0.033 FPS
196 evaluations per ray



100x speed



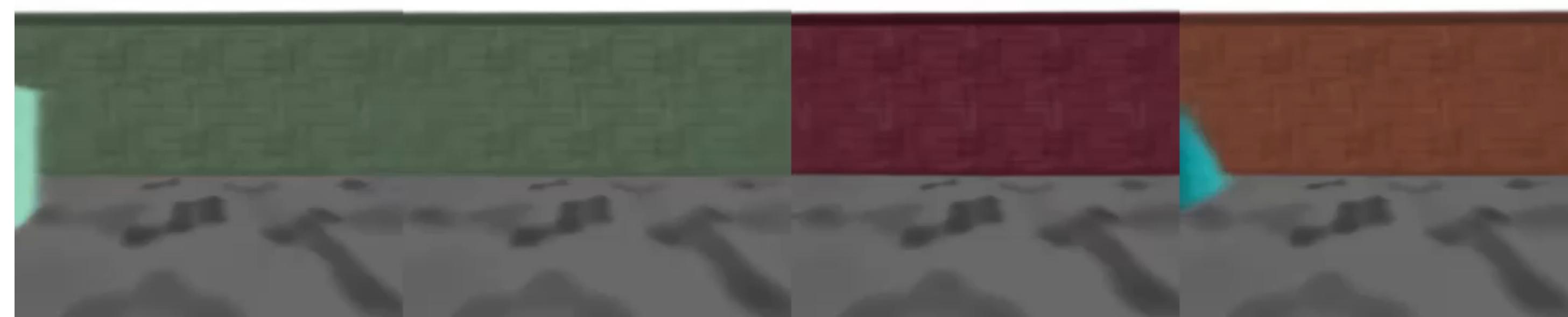
Real-time.

>100x reduction in memory: Can be trained on small GPUs!

Light Field Networks

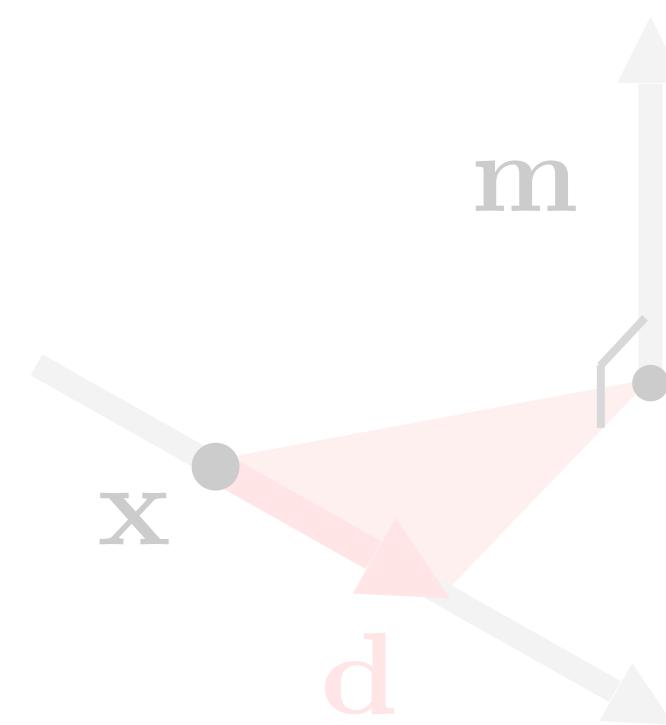
500 FPS

1 evaluation per ray

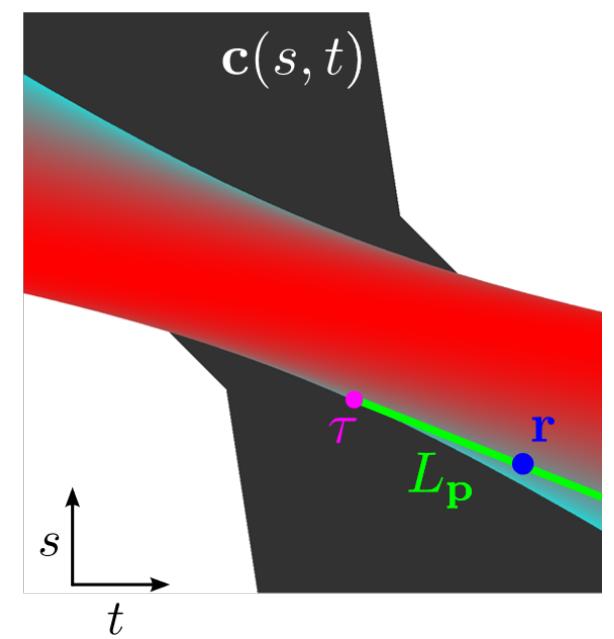


How LFNs encode 3D geometry

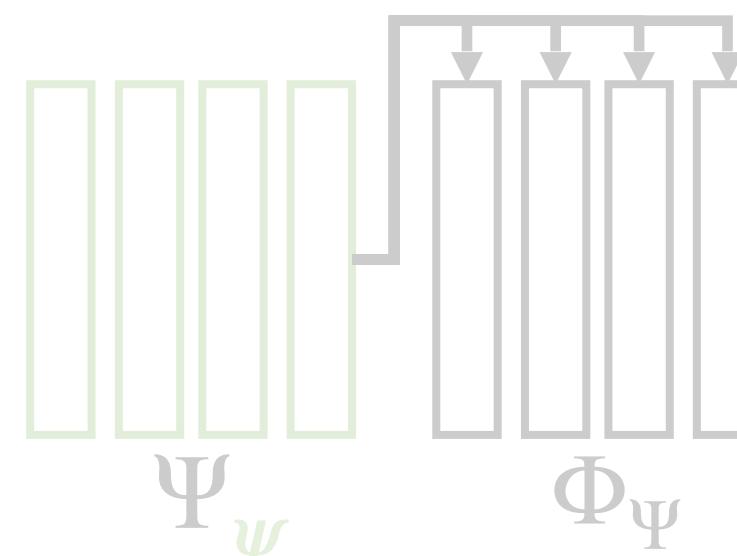
Parameterization



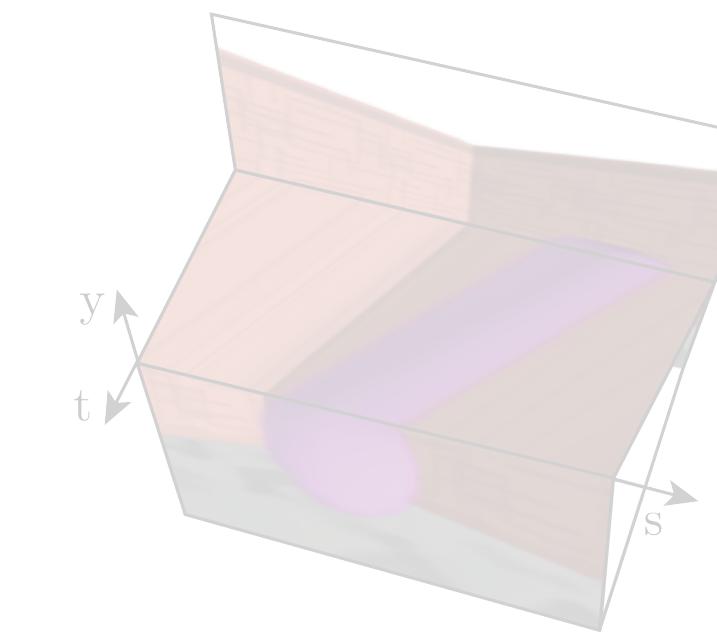
LFN Geometry



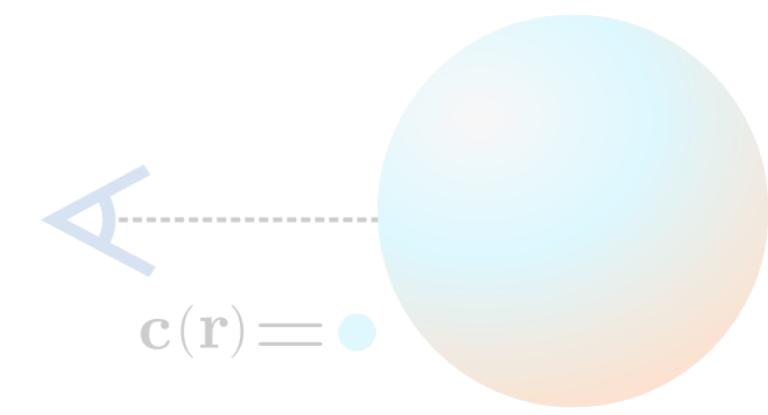
Meta-Learning



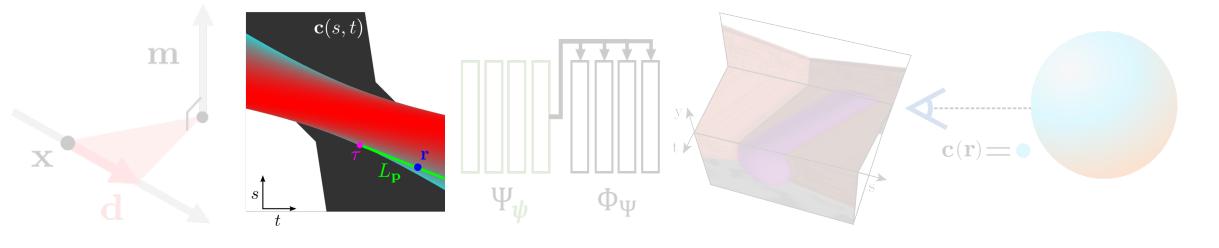
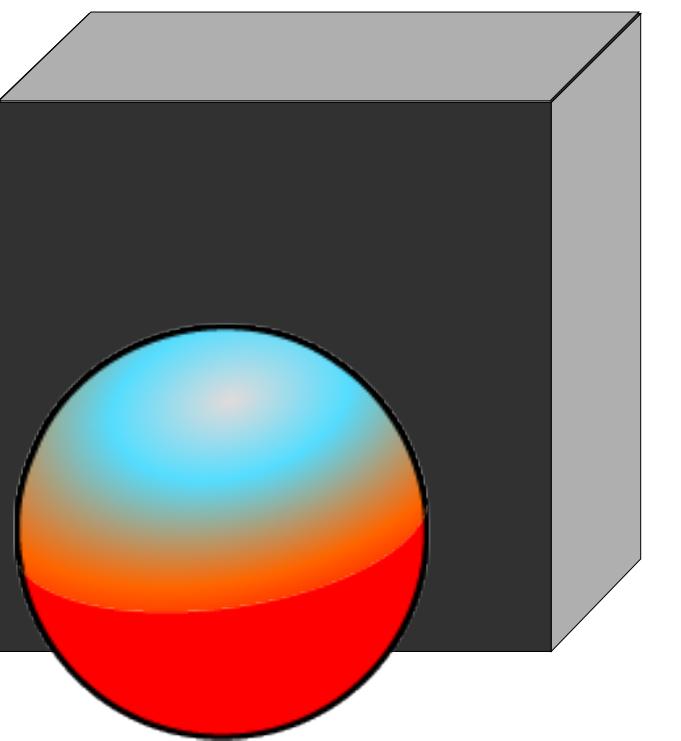
Results



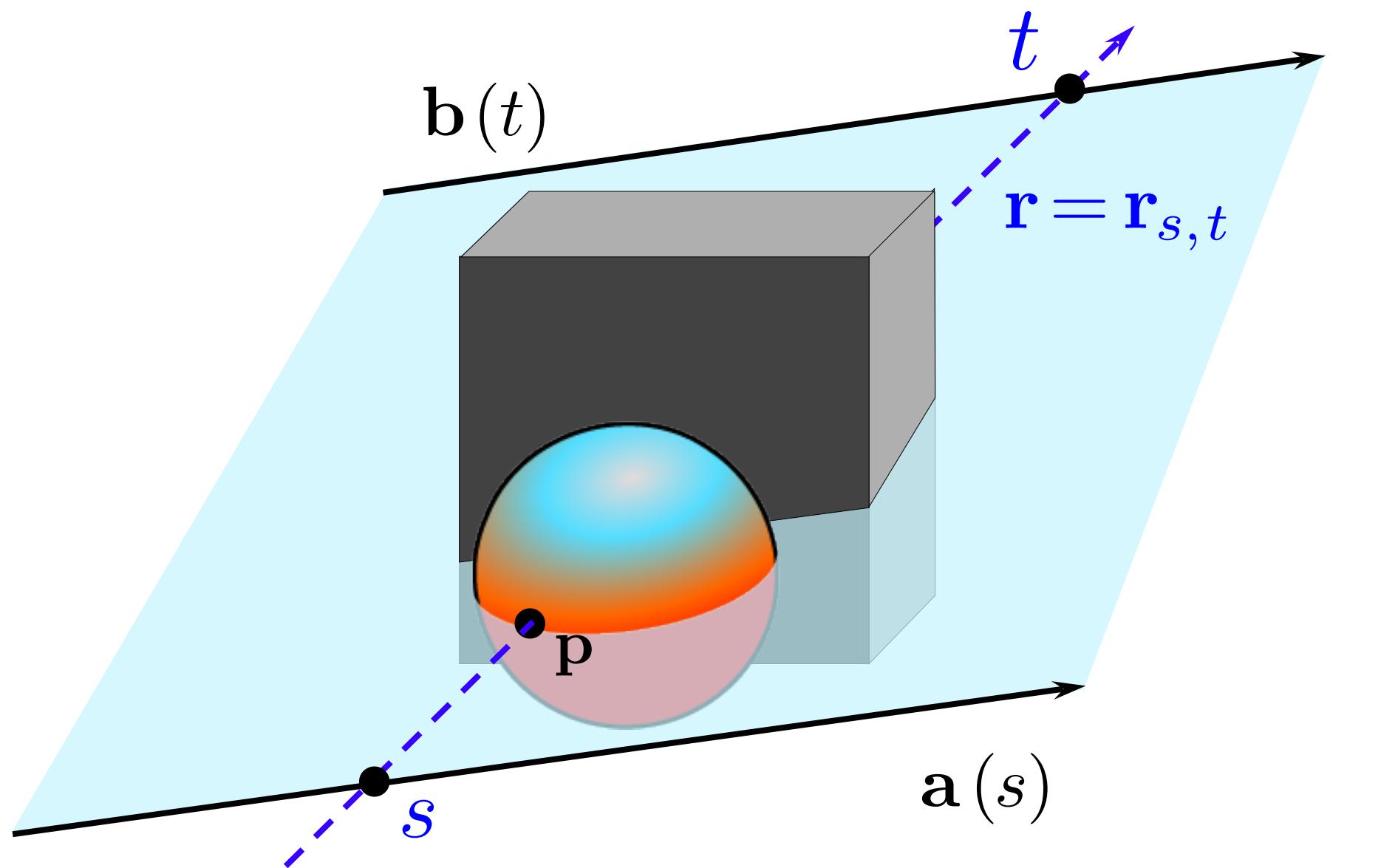
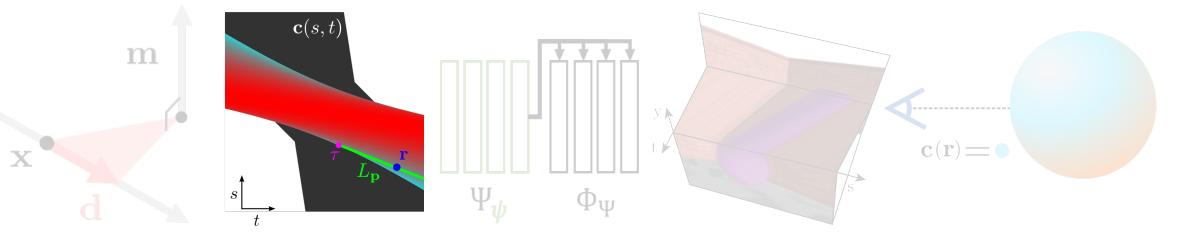
Limitations



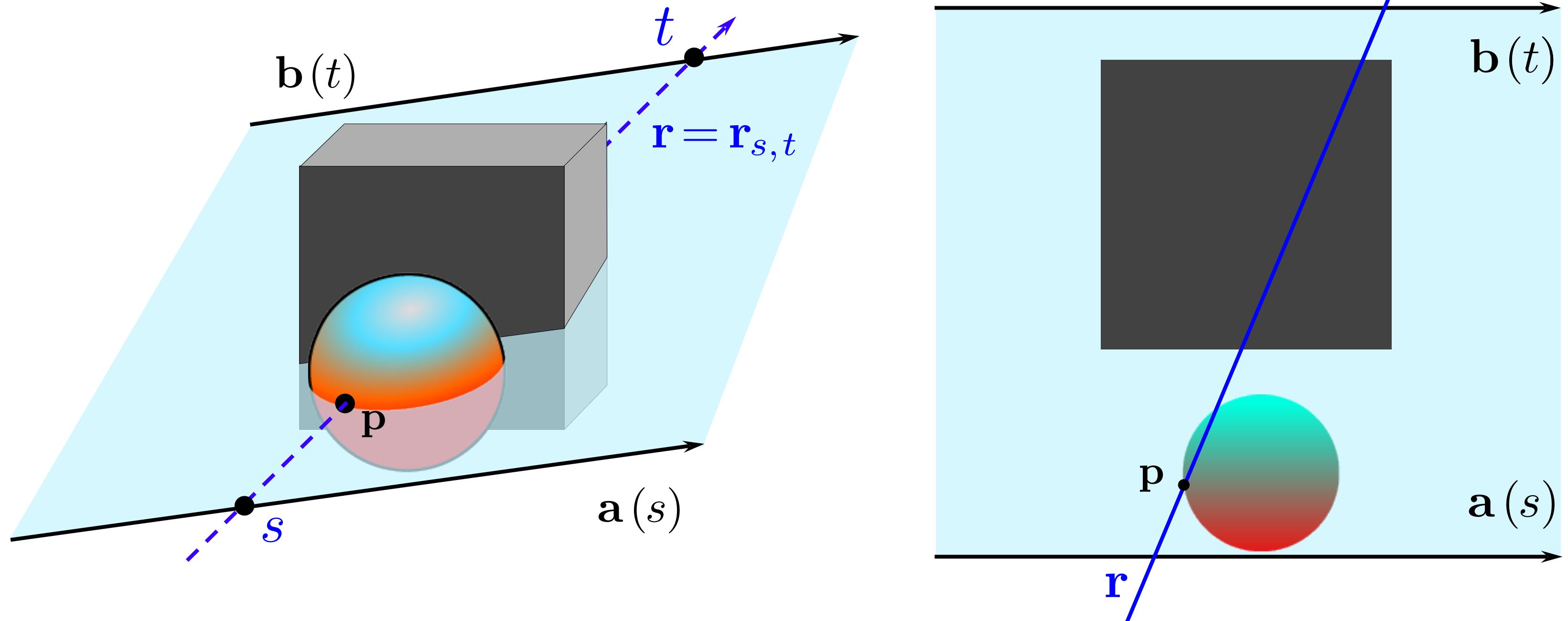
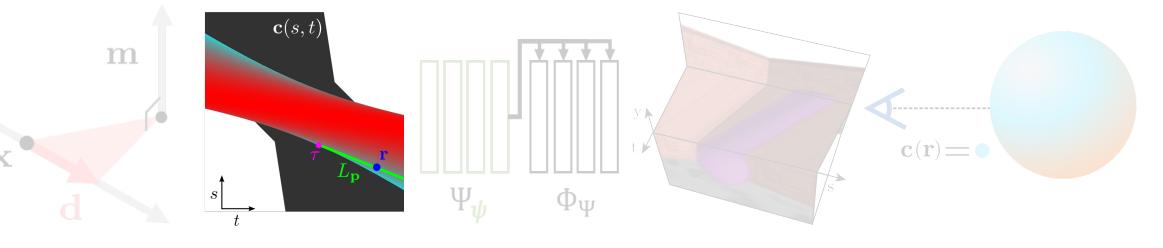
The geometry of LFNs



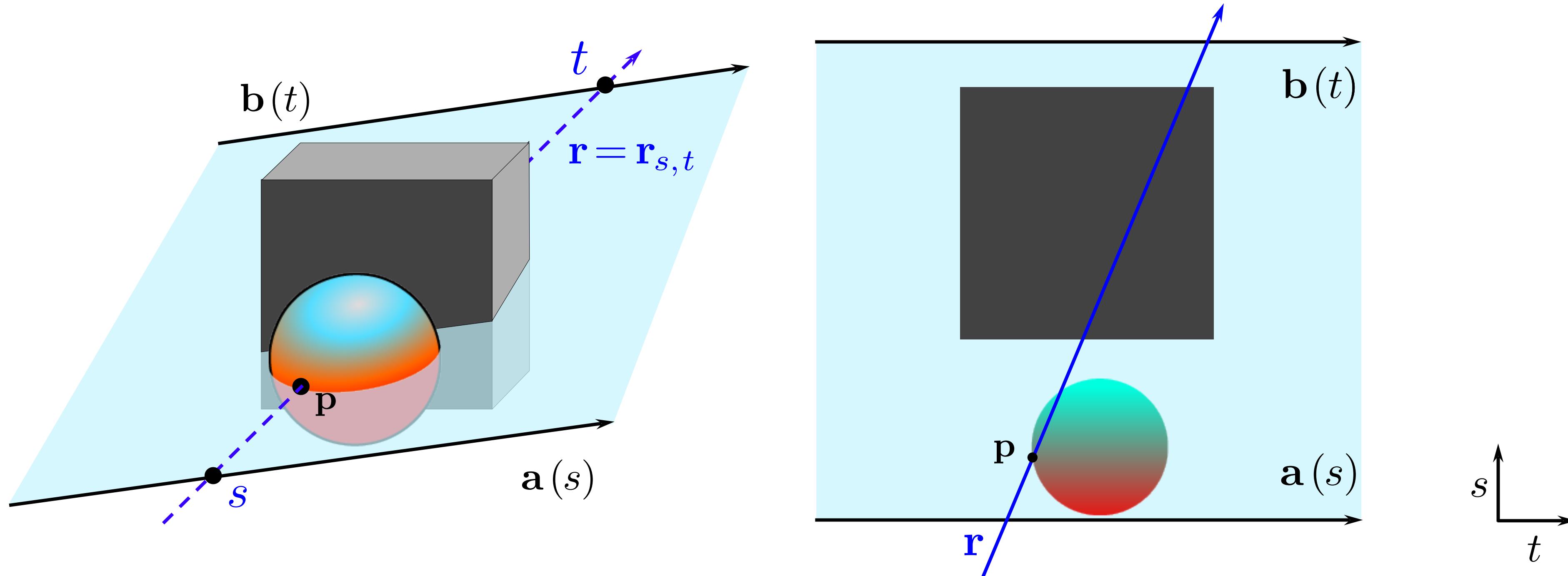
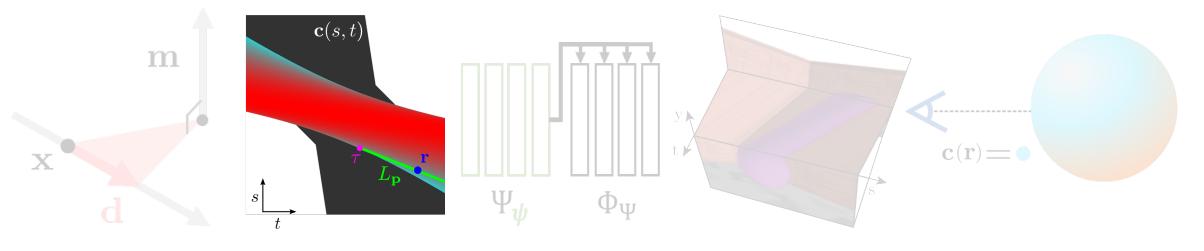
The geometry of LFNs



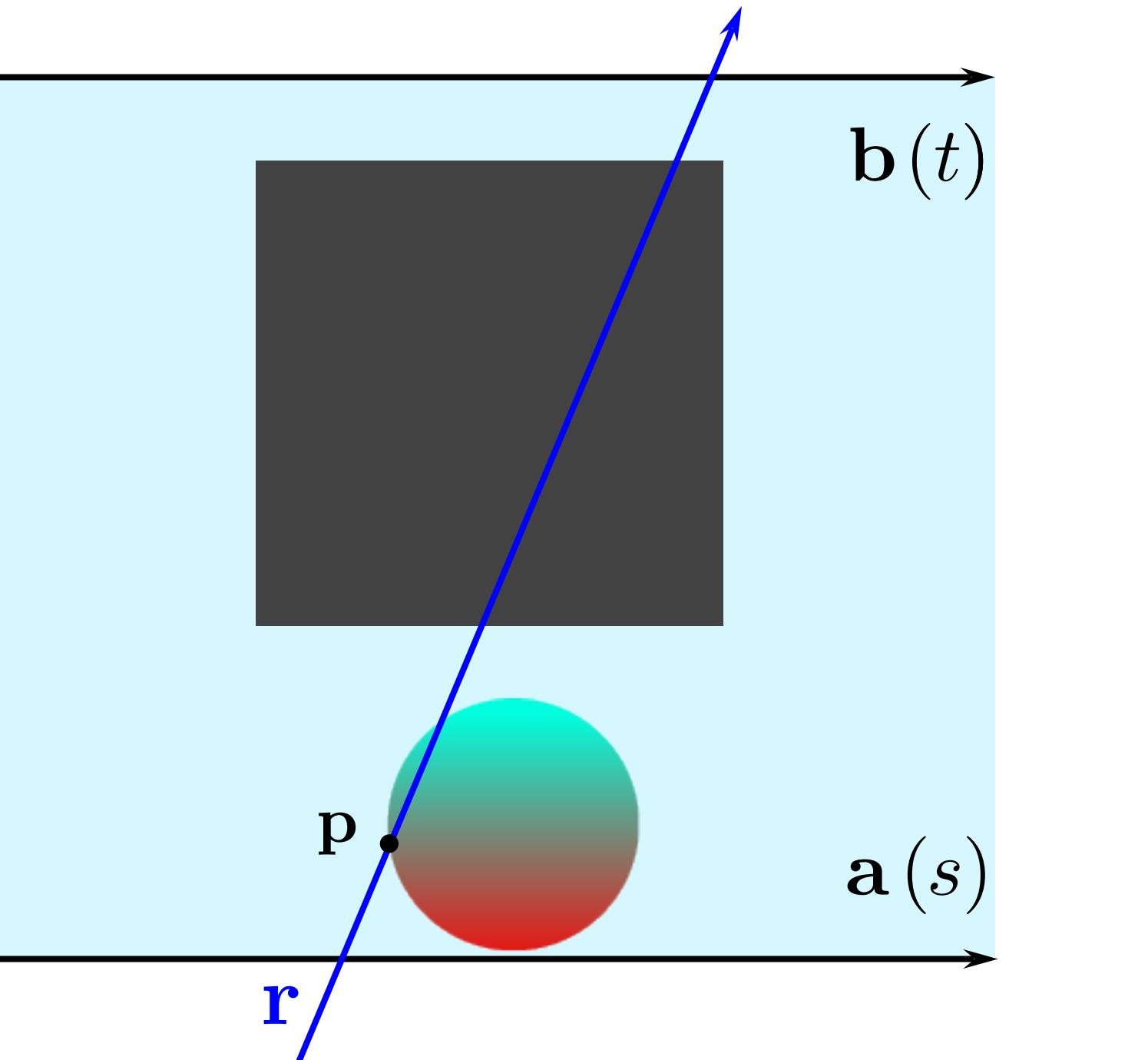
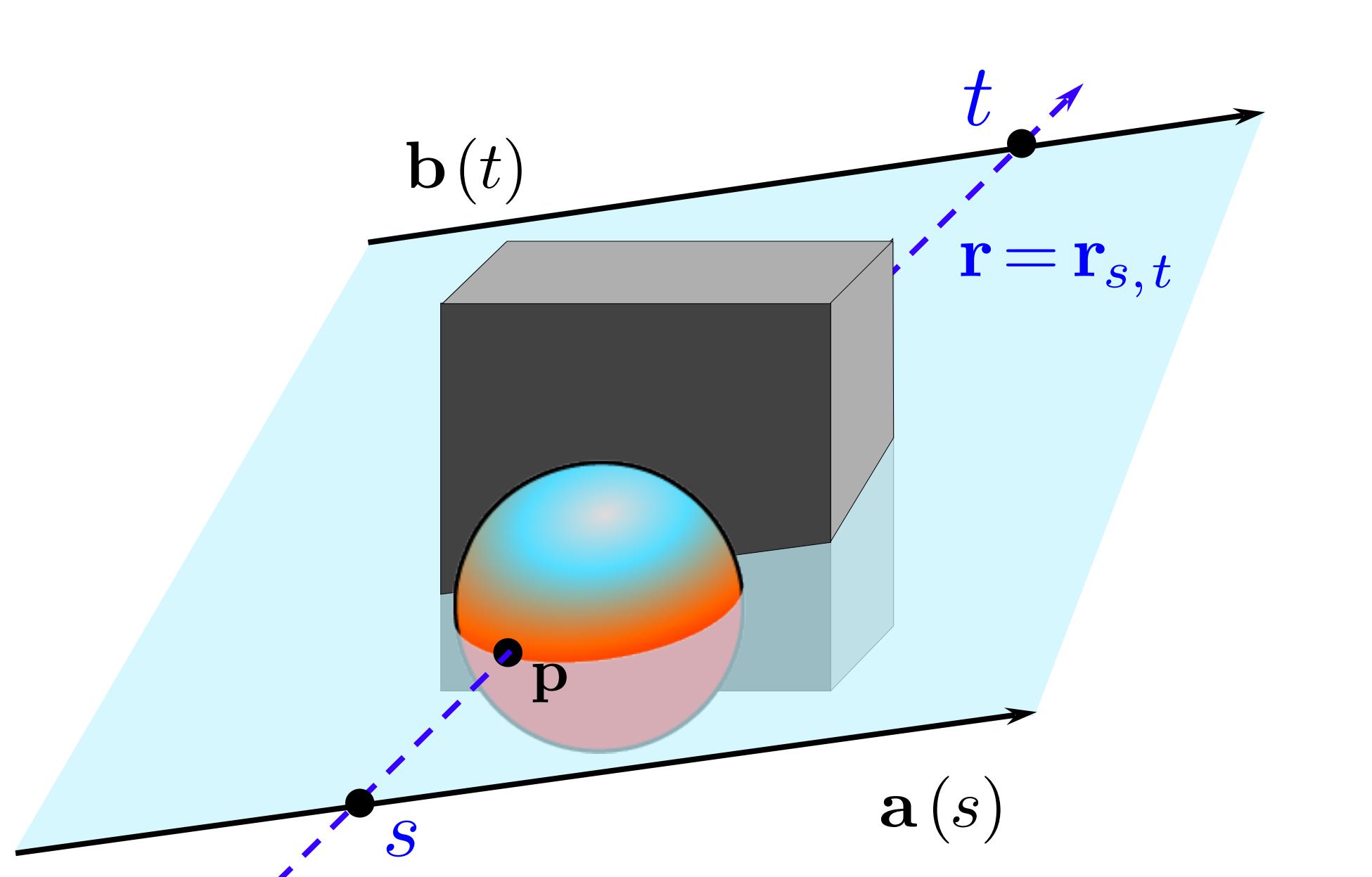
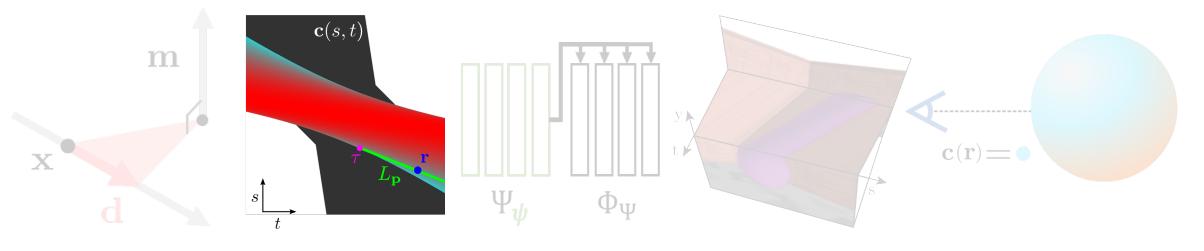
The geometry of LFNs



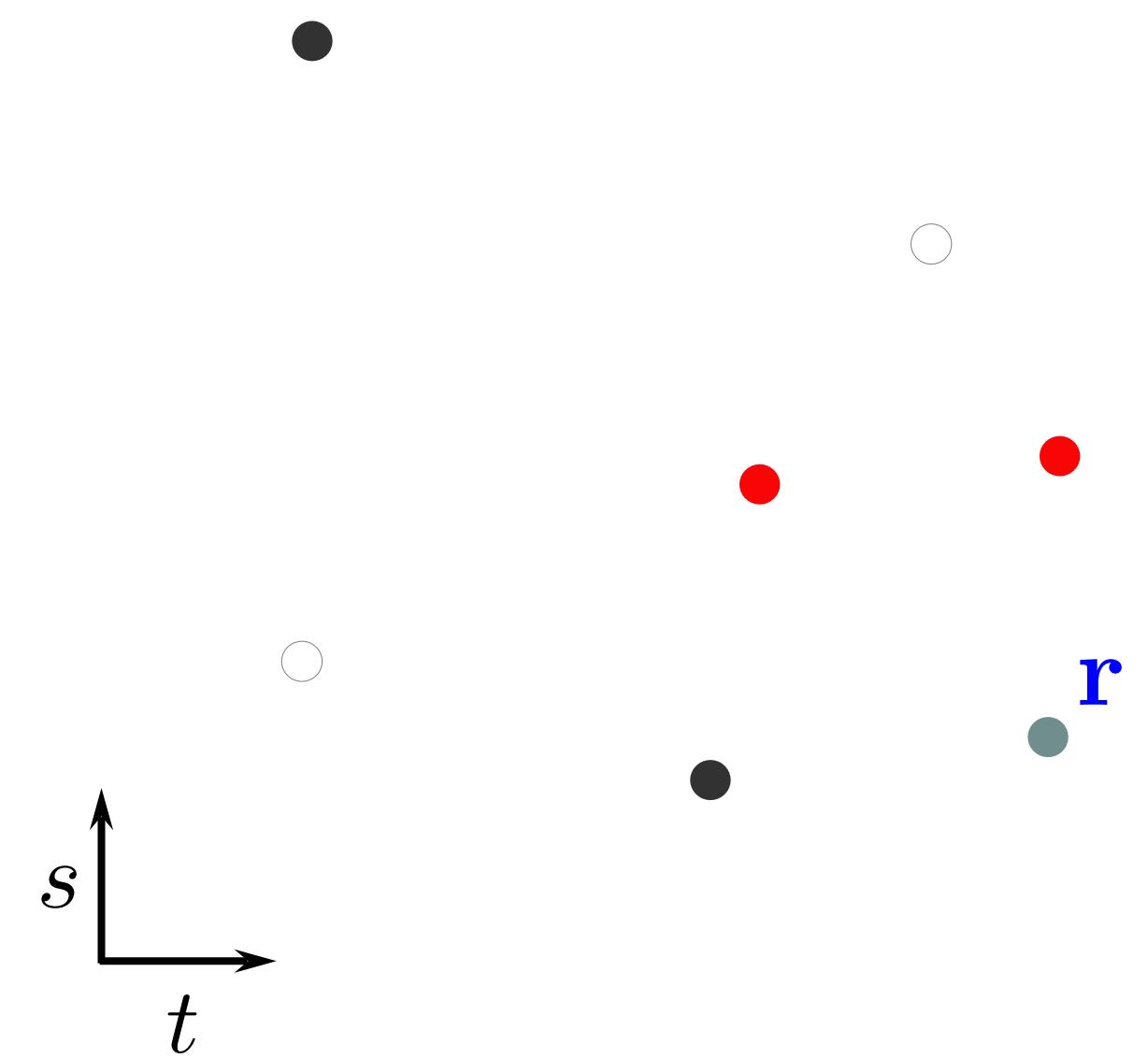
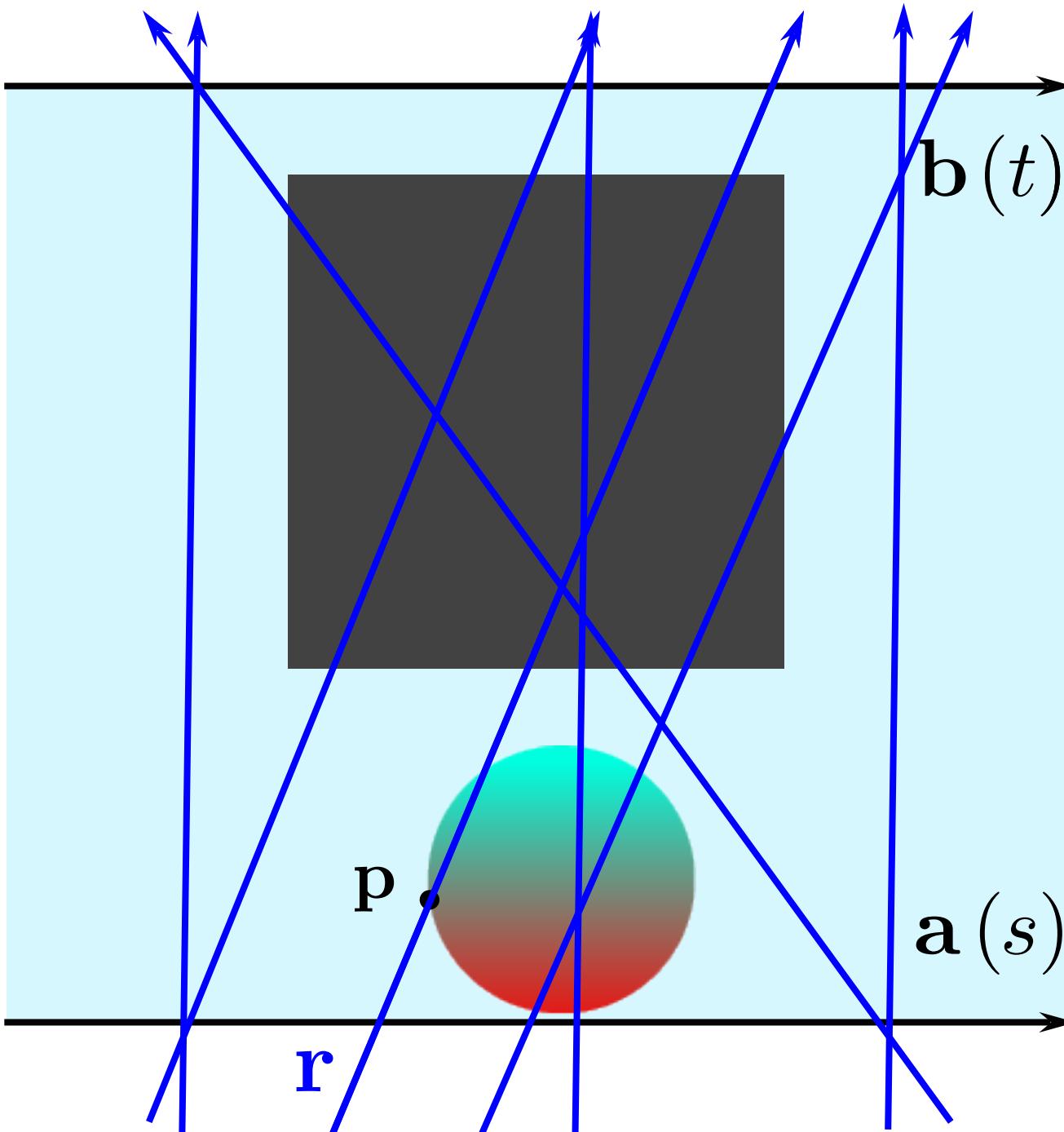
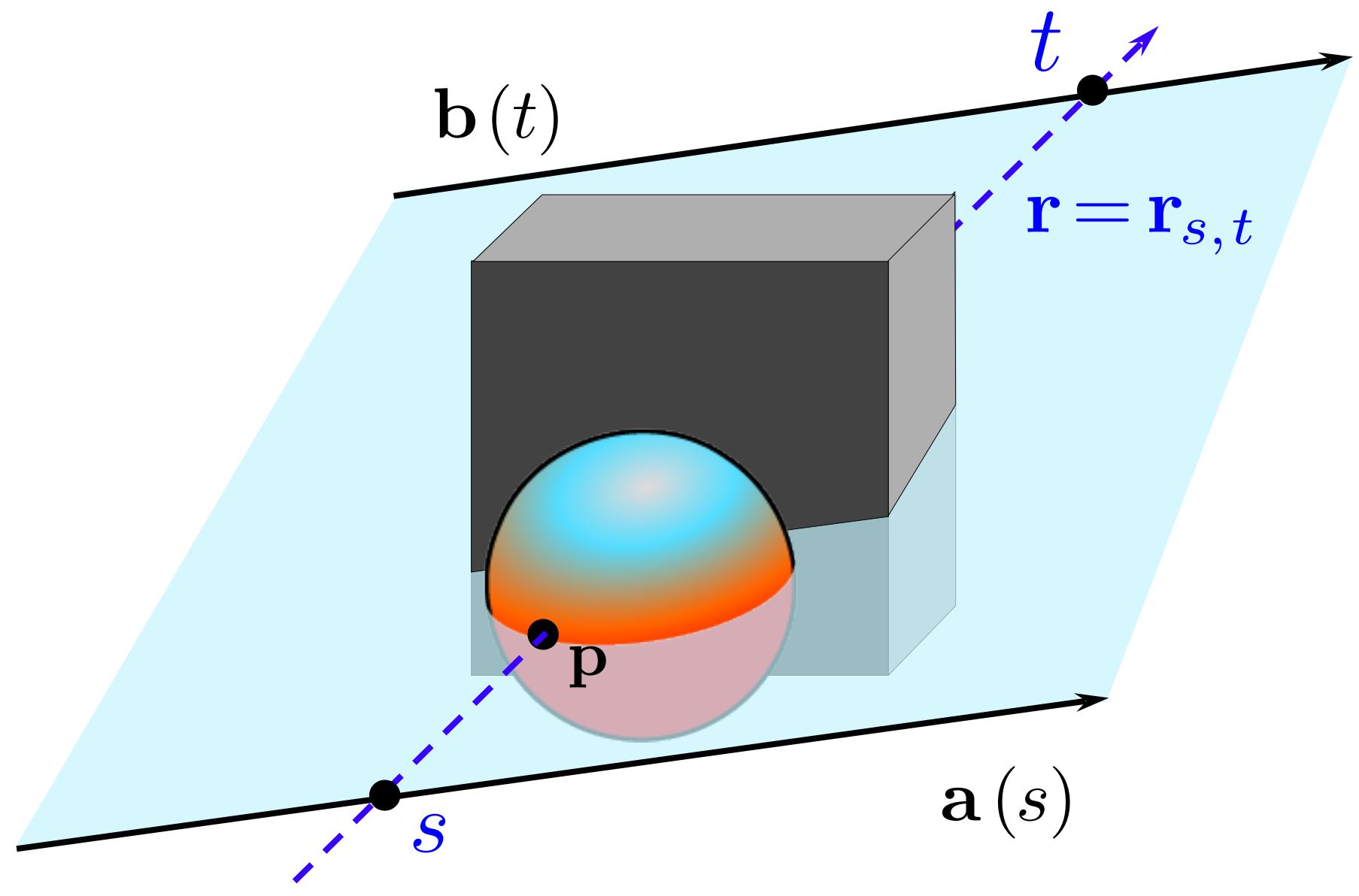
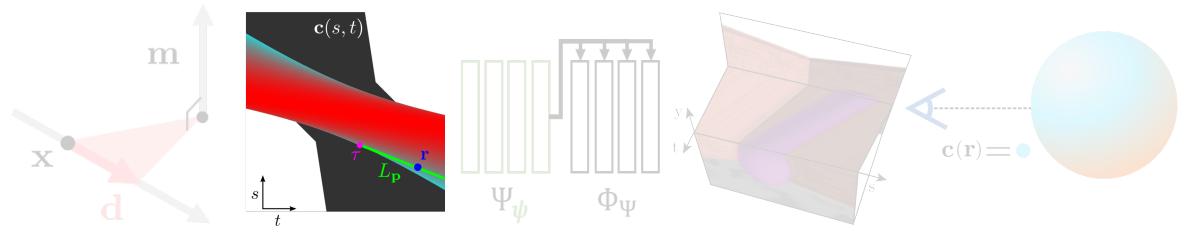
The geometry of LFNs



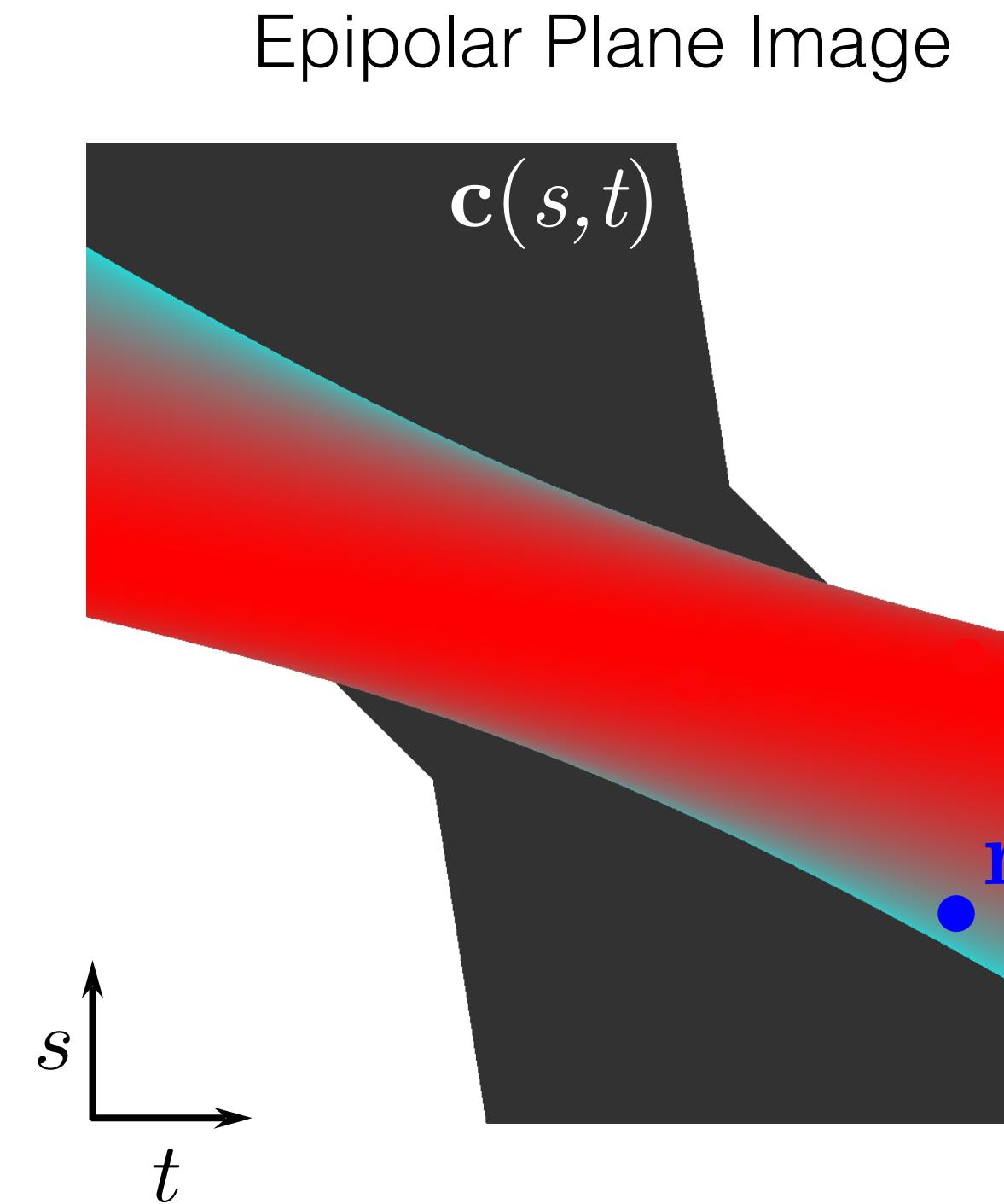
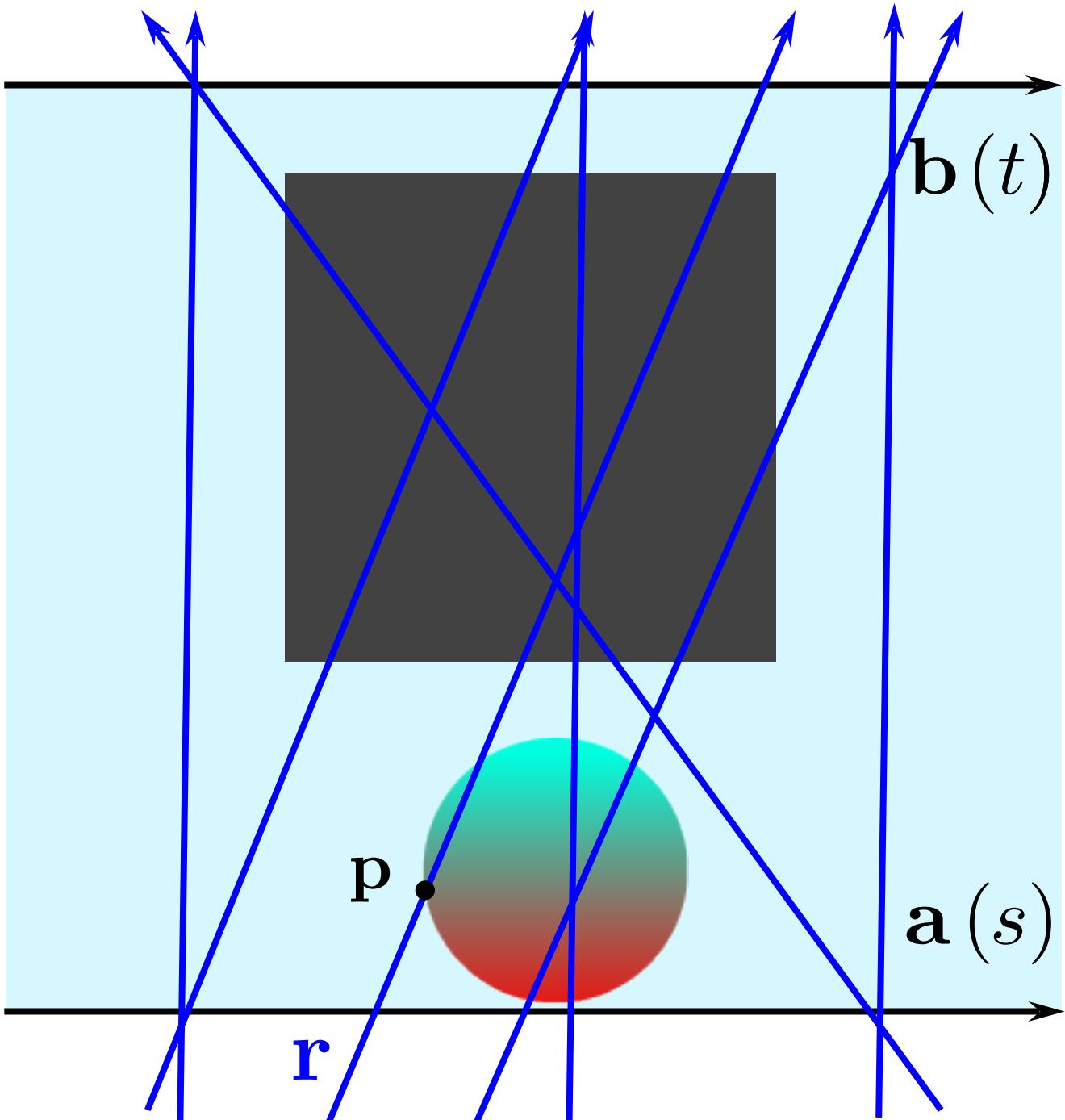
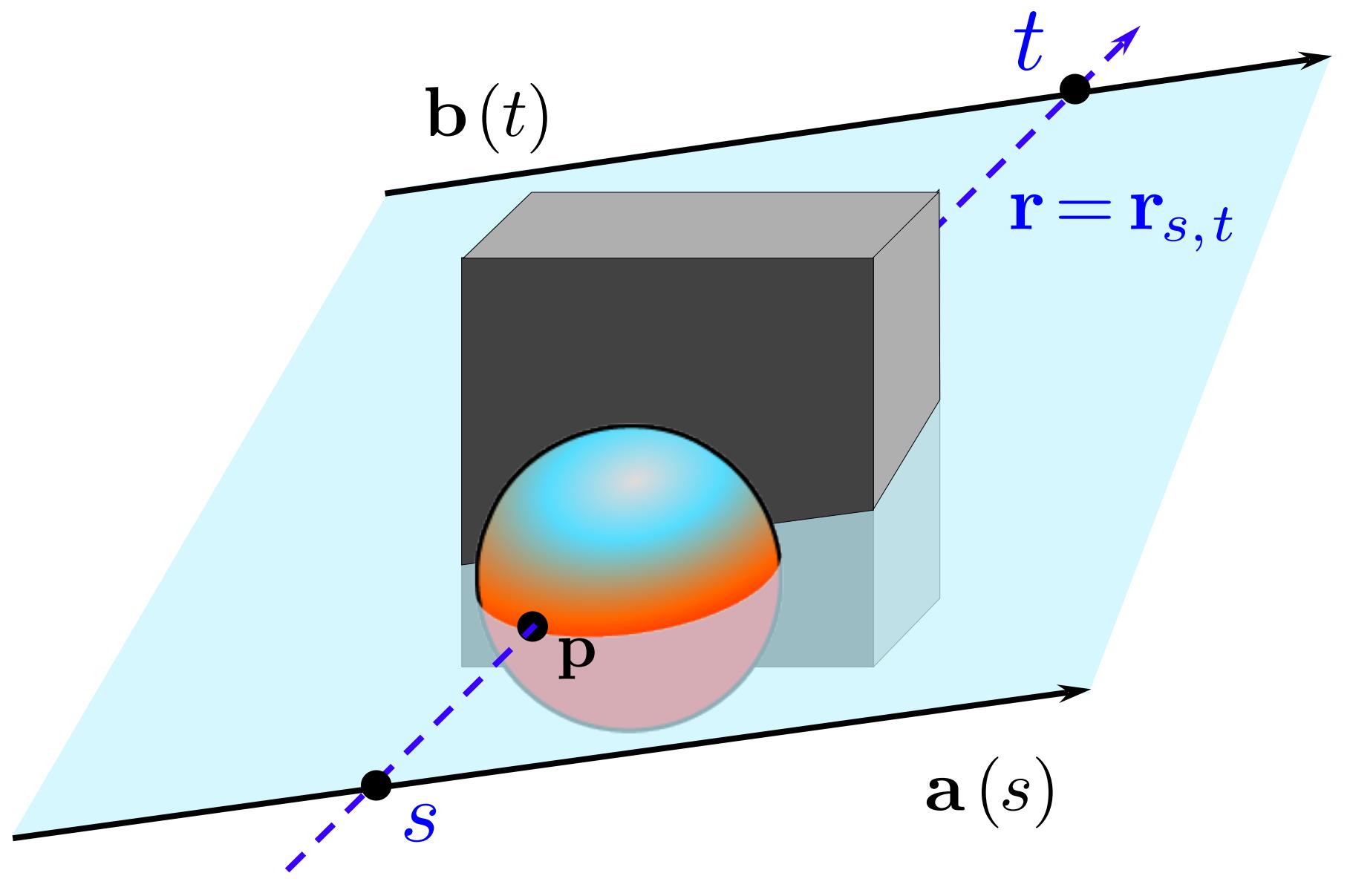
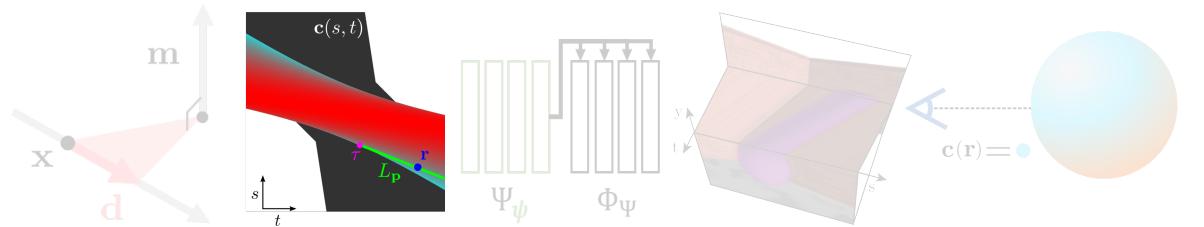
The geometry of LFNs



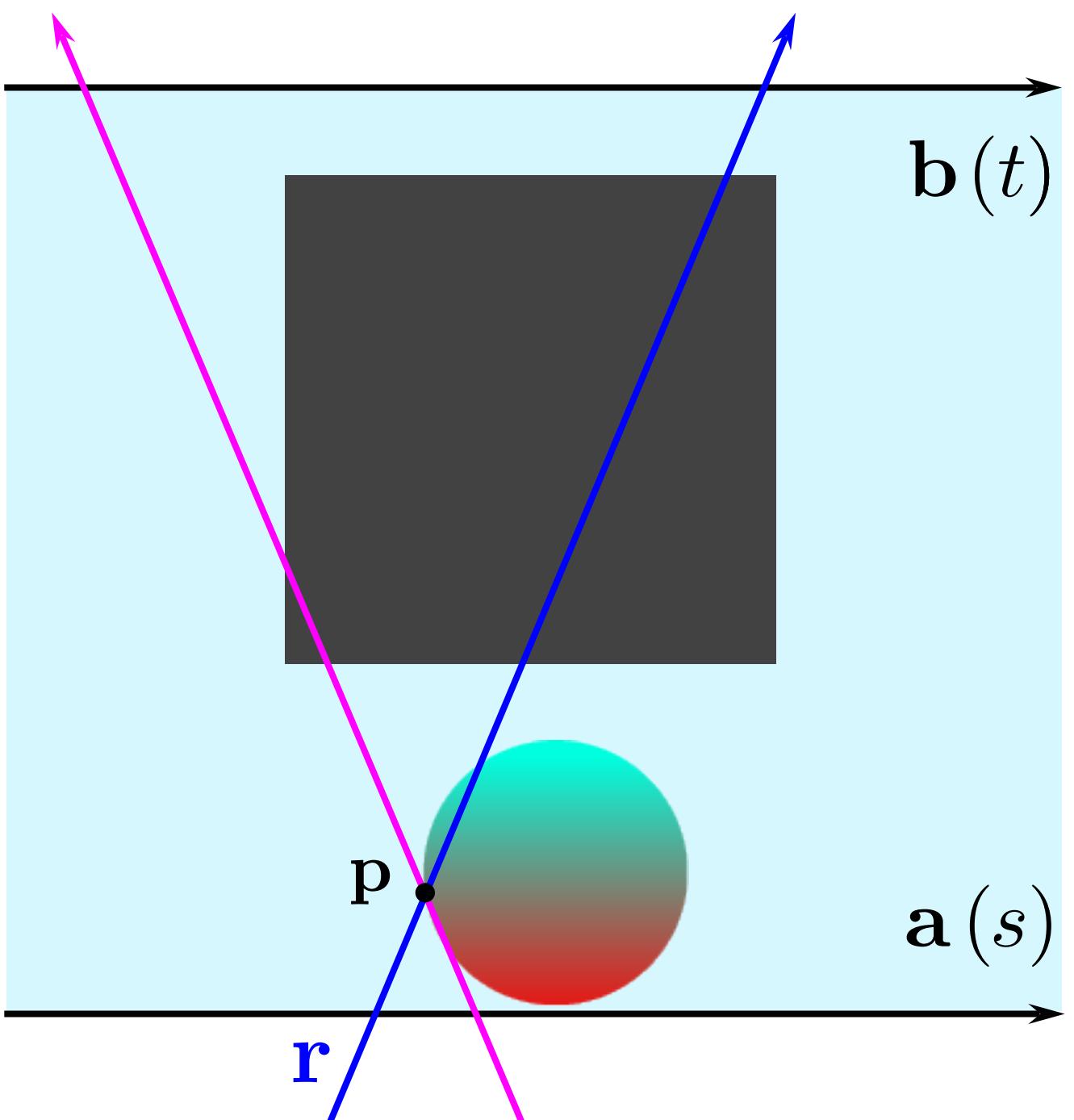
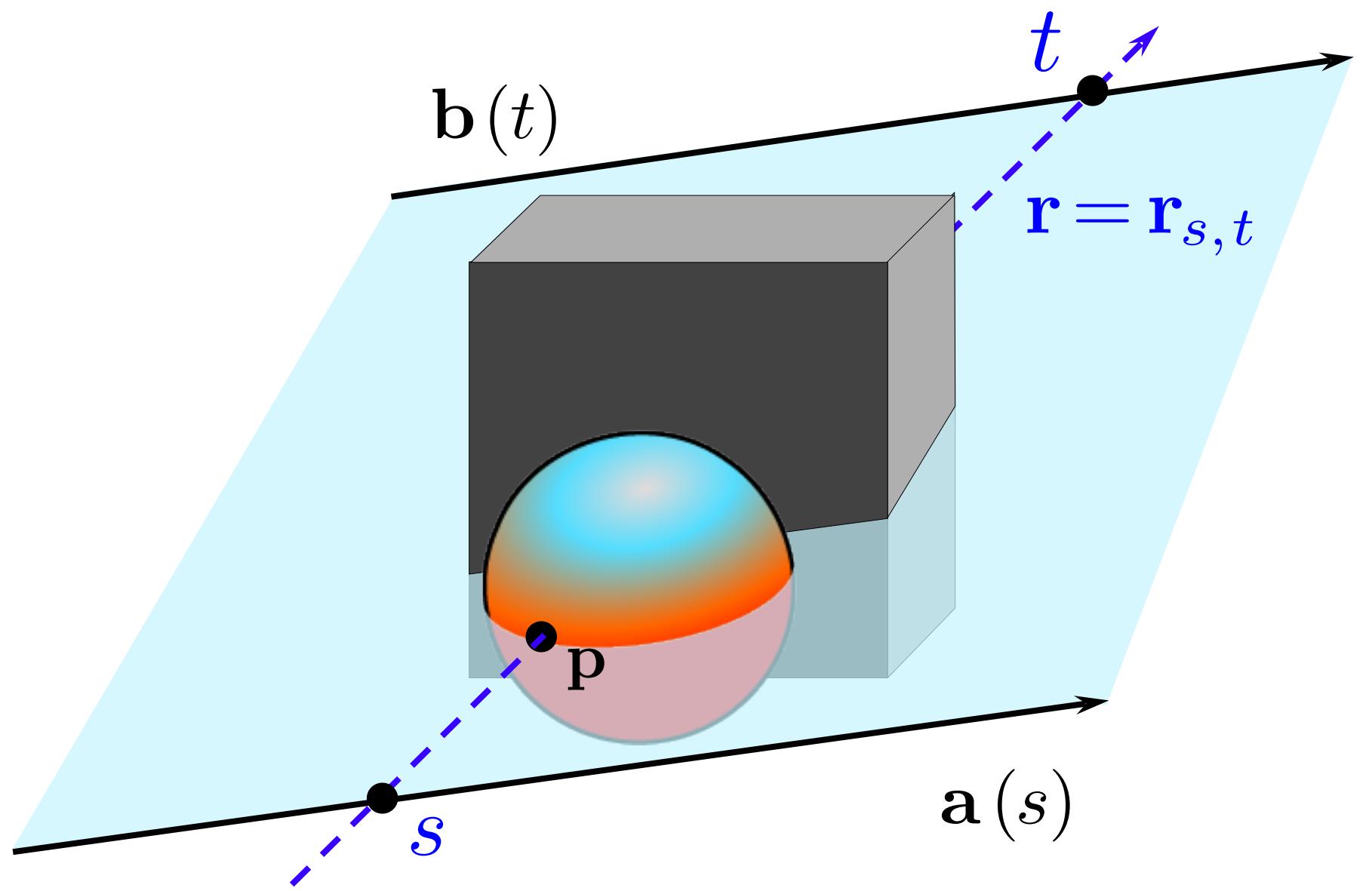
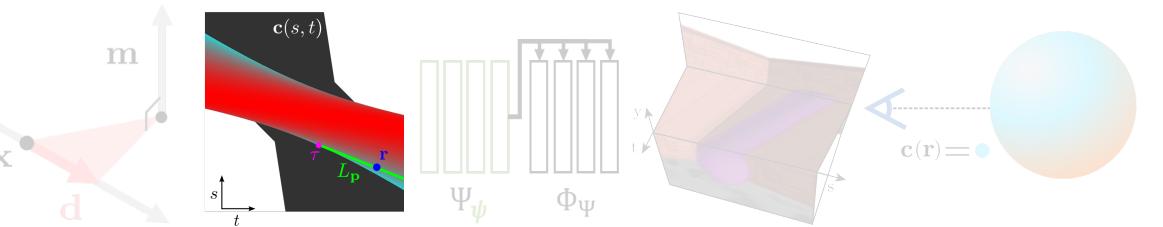
The geometry of LFNs



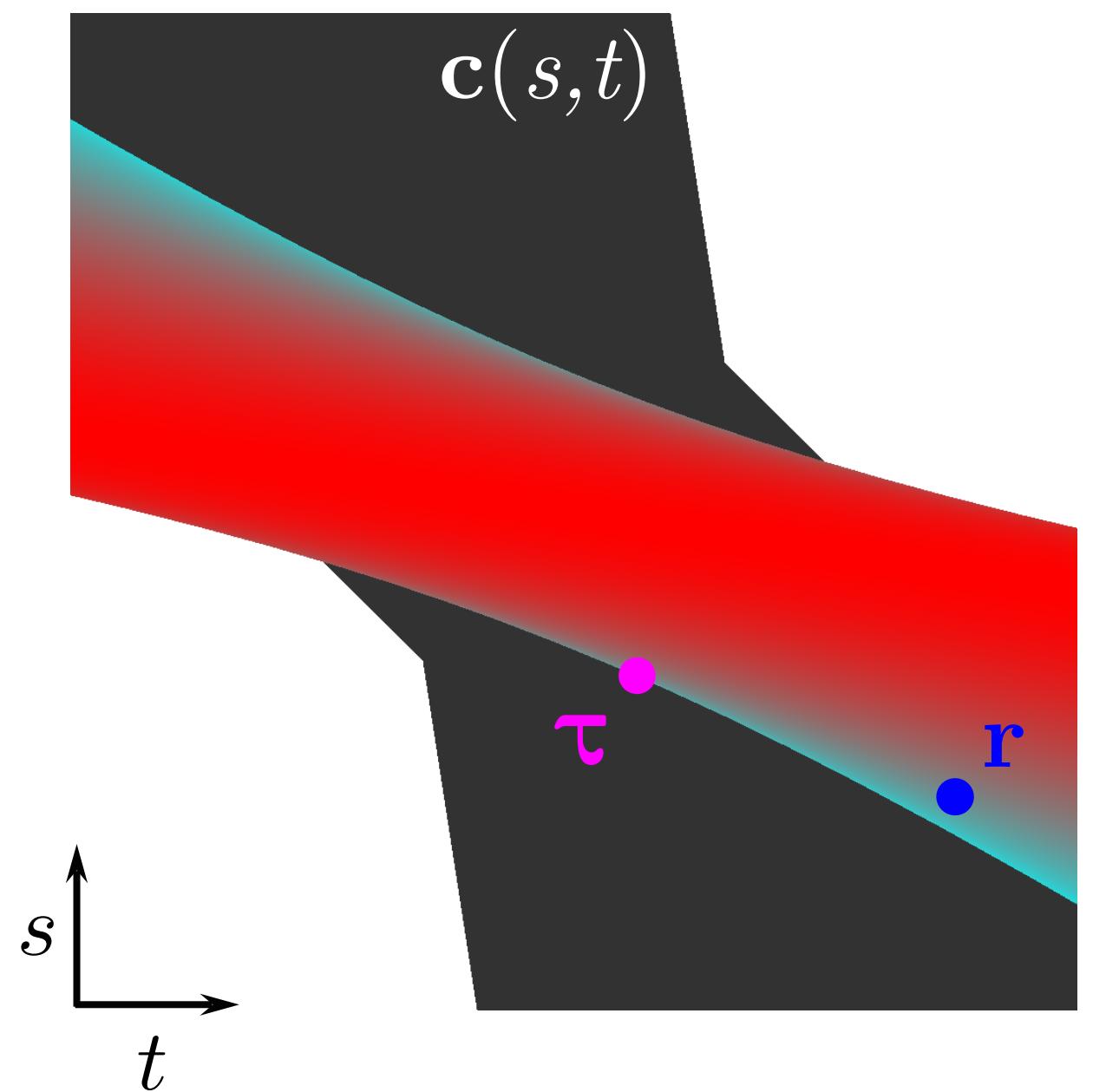
The geometry of LFNs



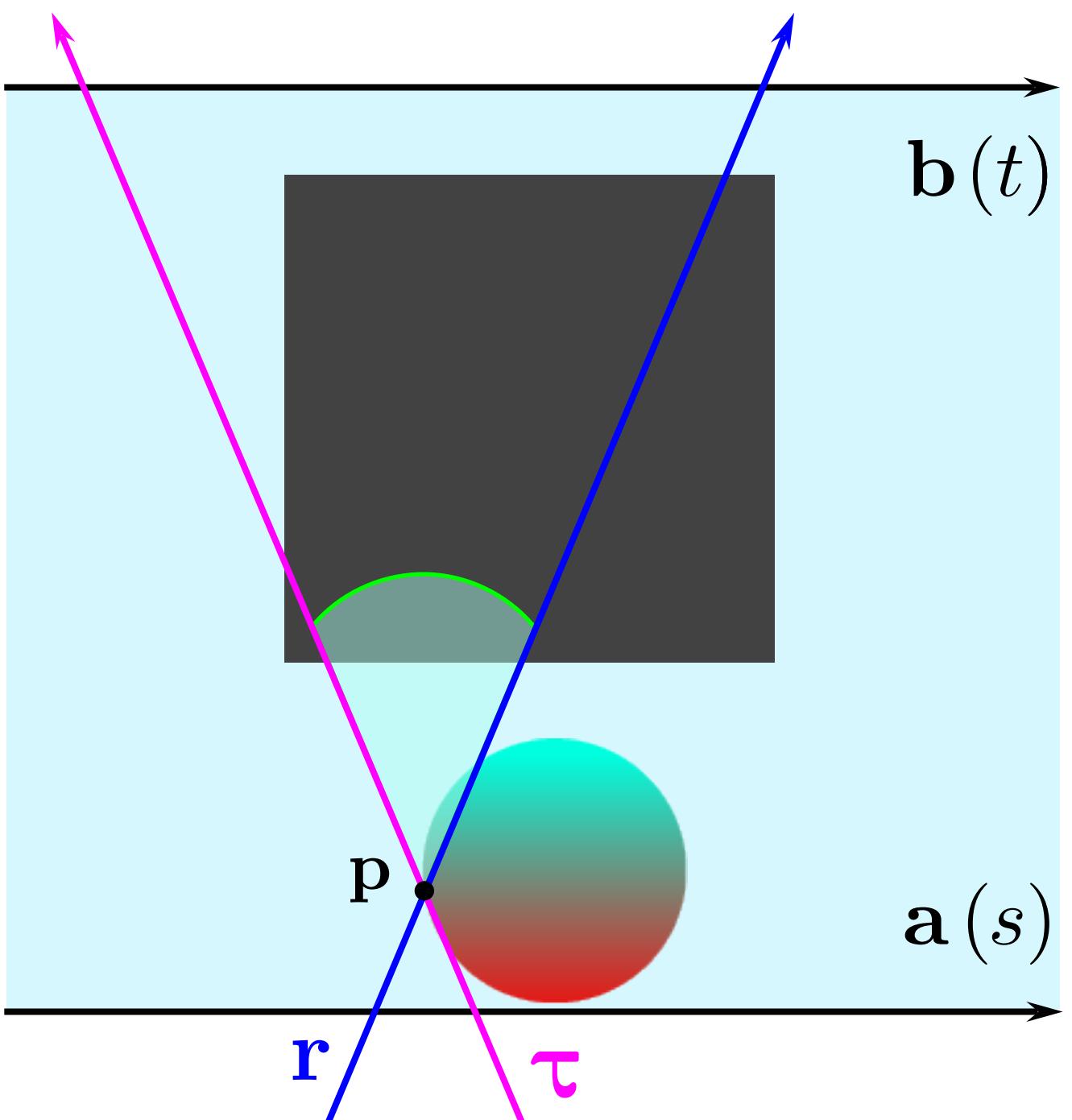
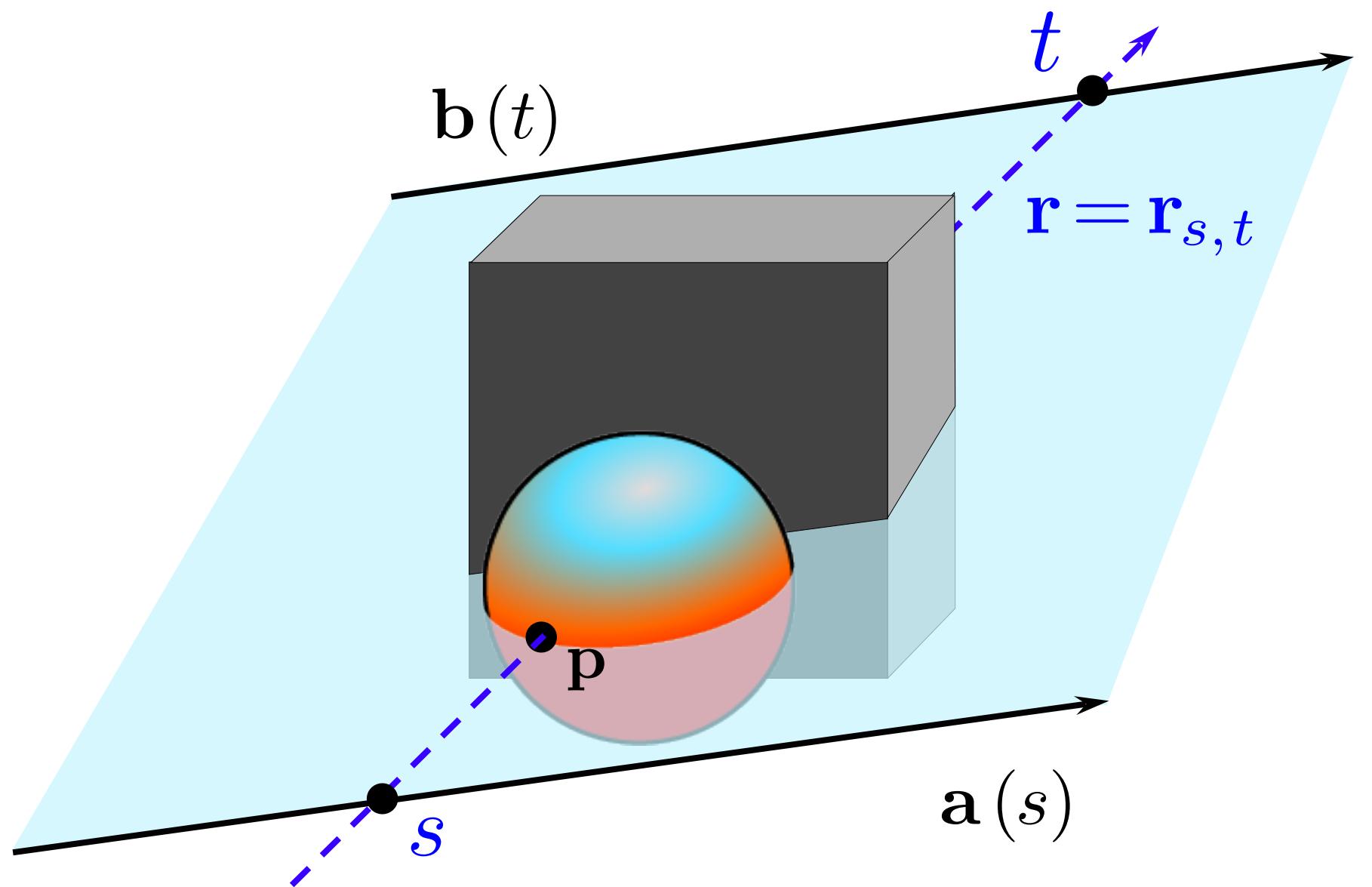
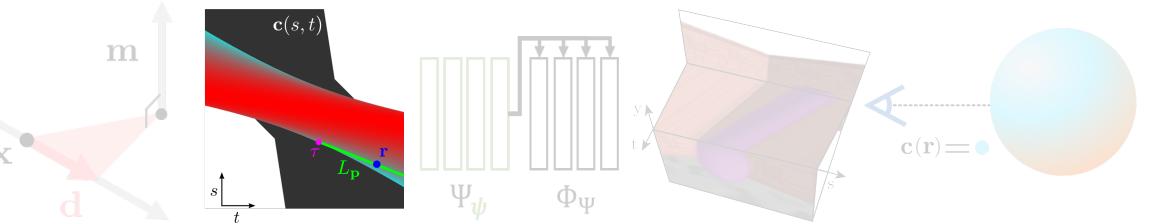
The geometry of LFNs



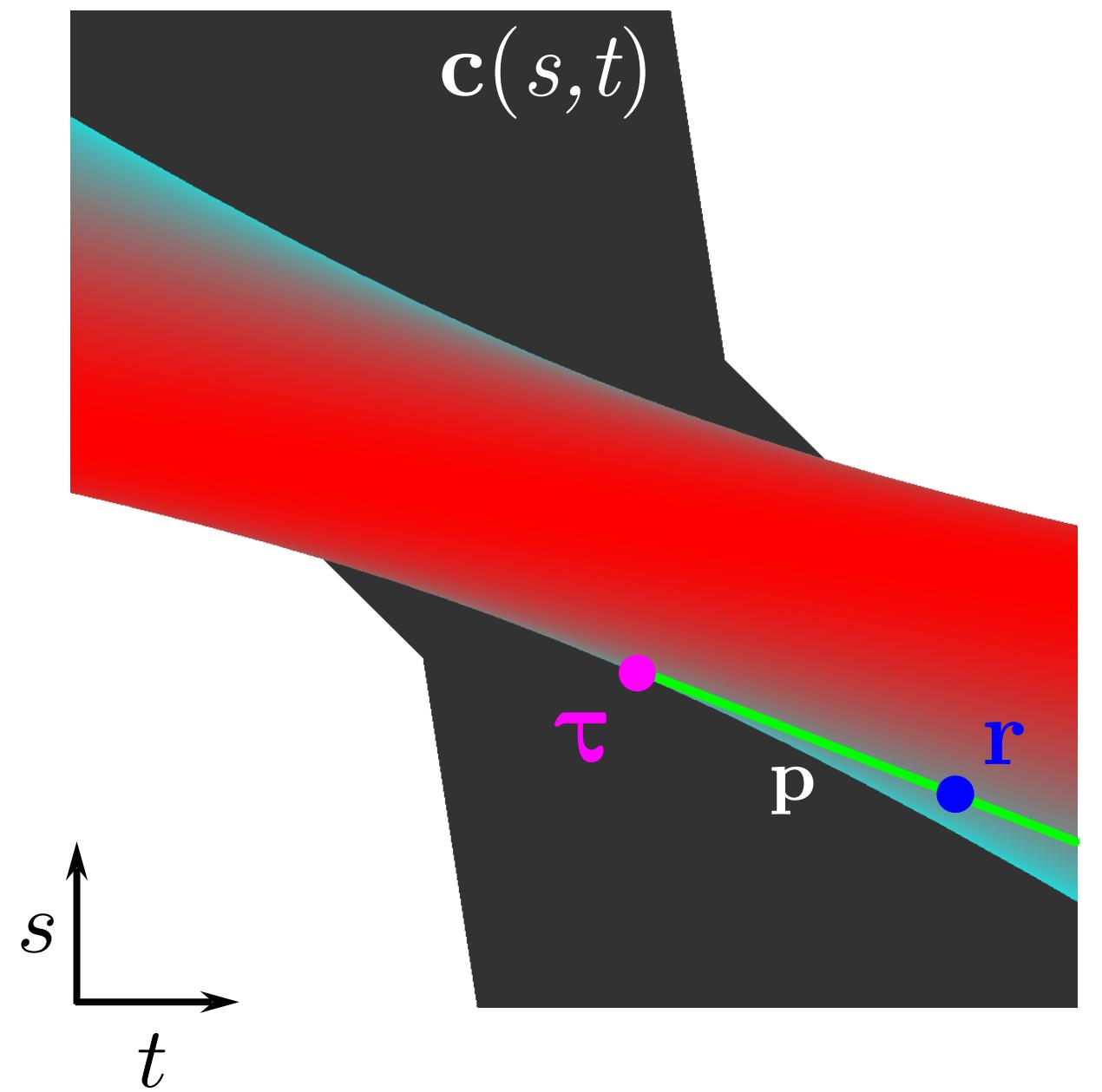
Epipolar Plane Image



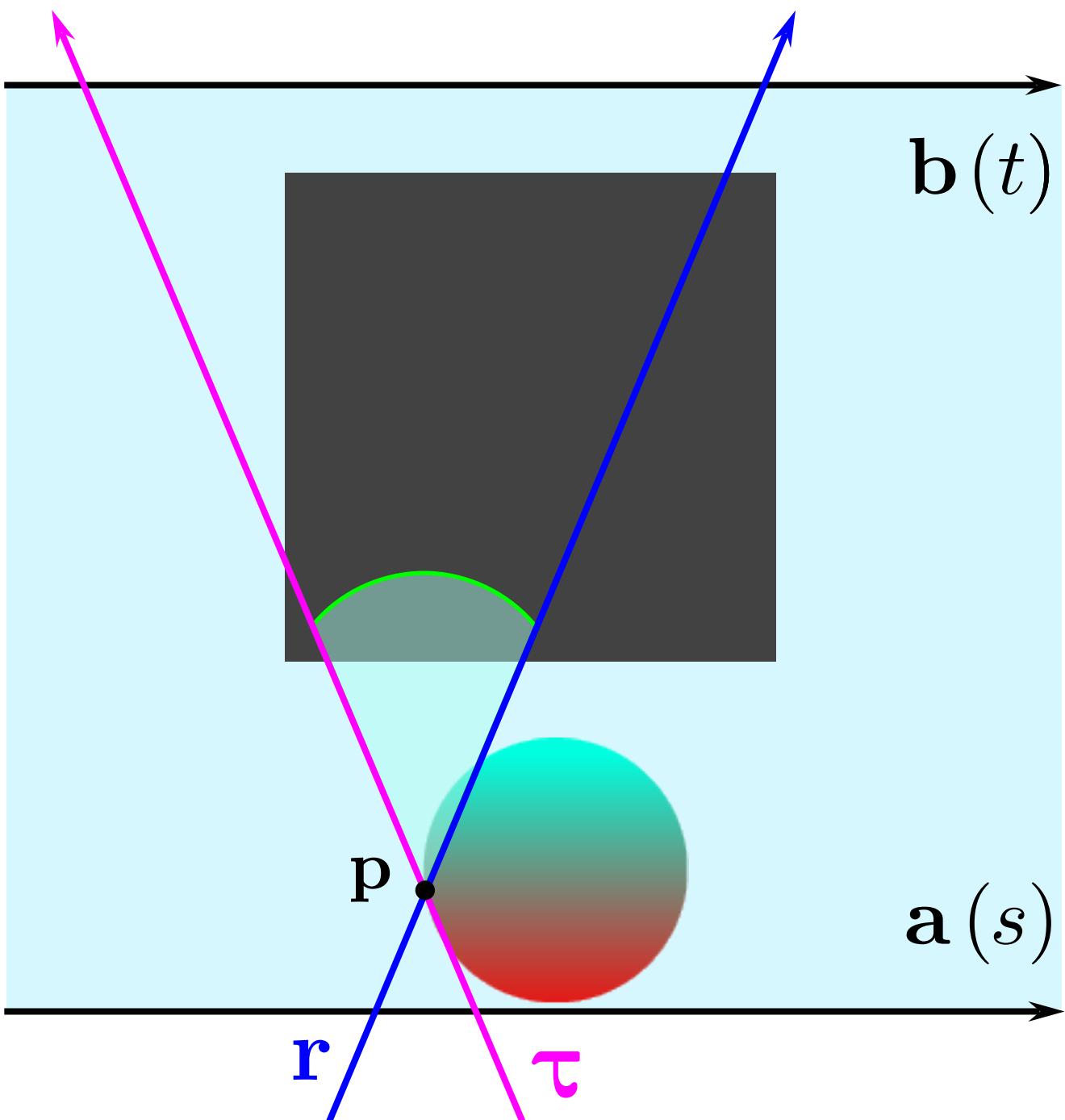
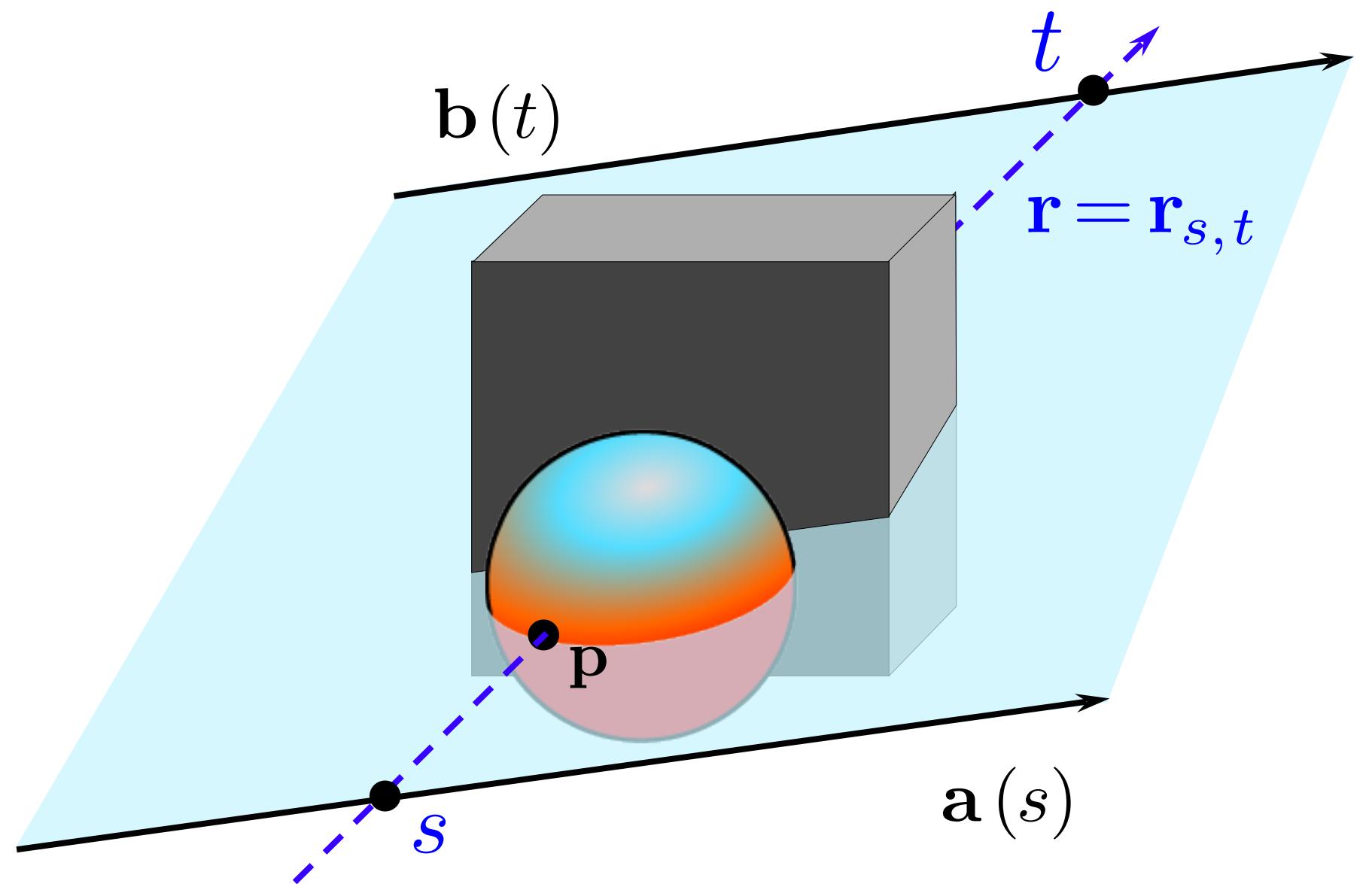
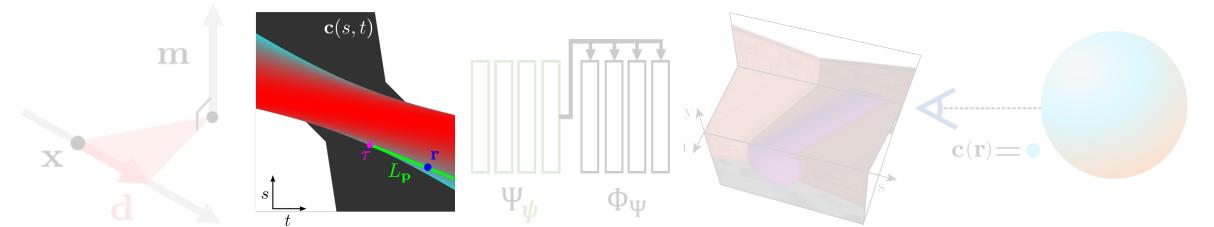
The geometry of LFNs



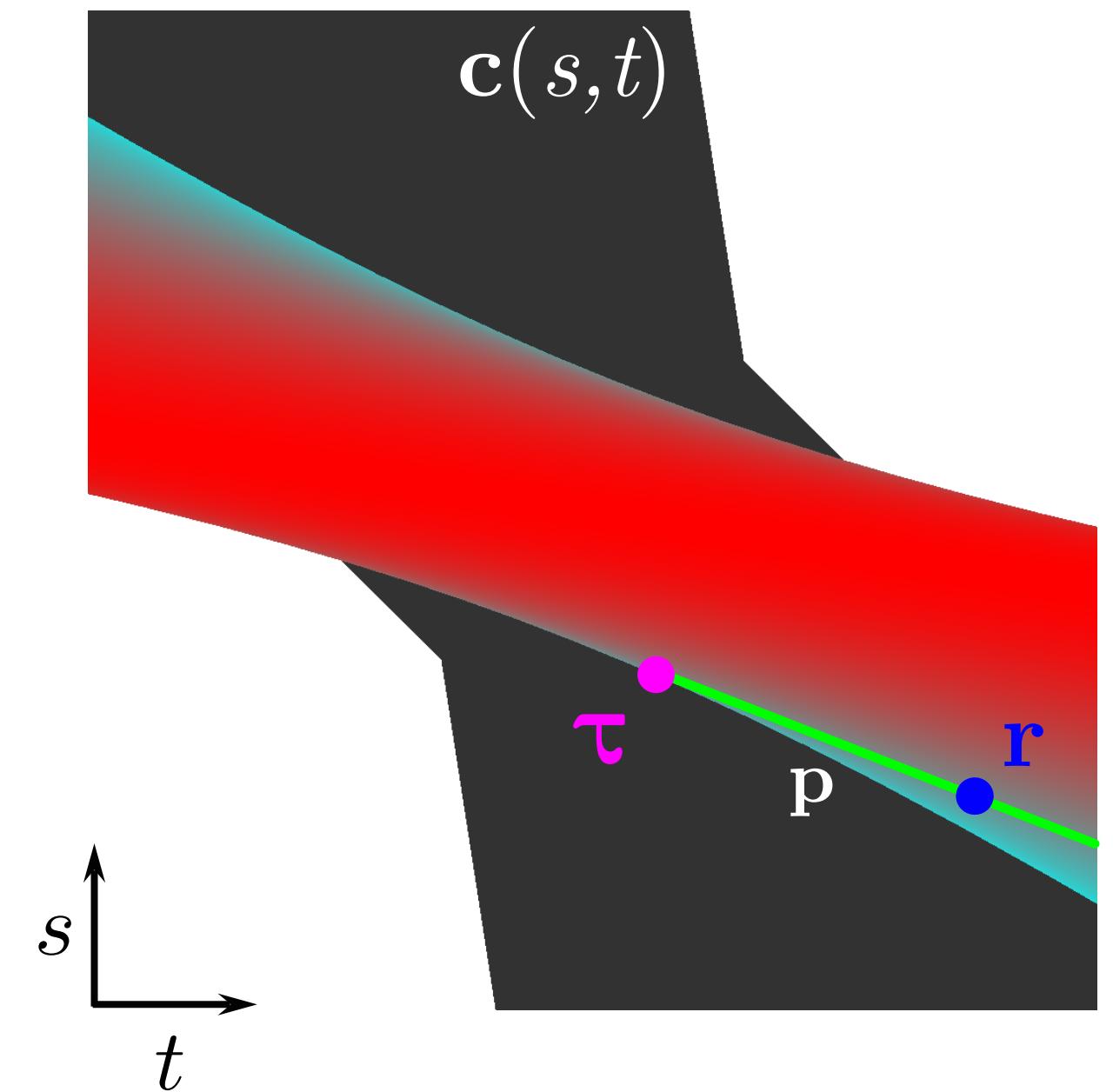
Epipolar Plane Image



The geometry of LFNs

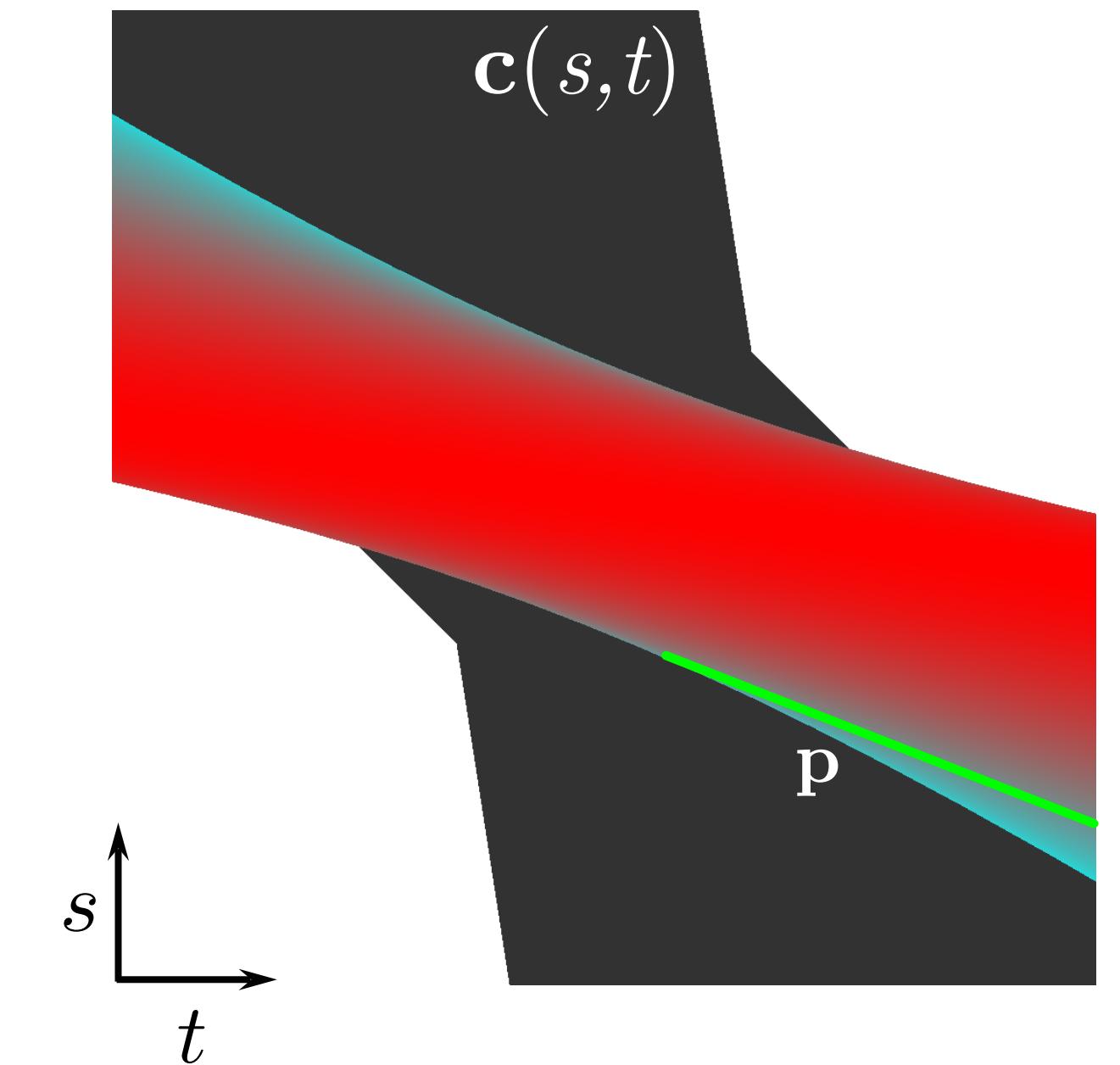
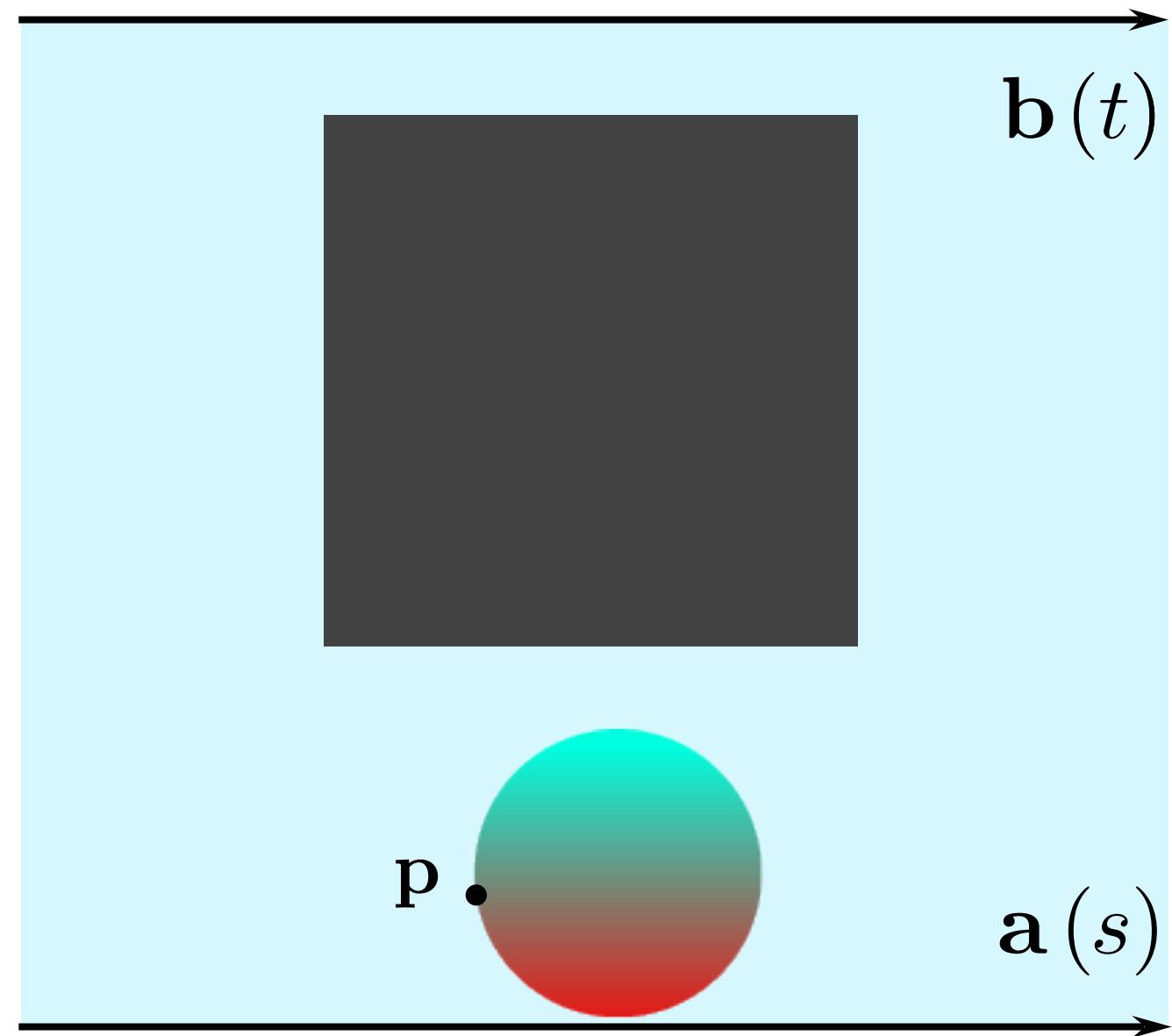
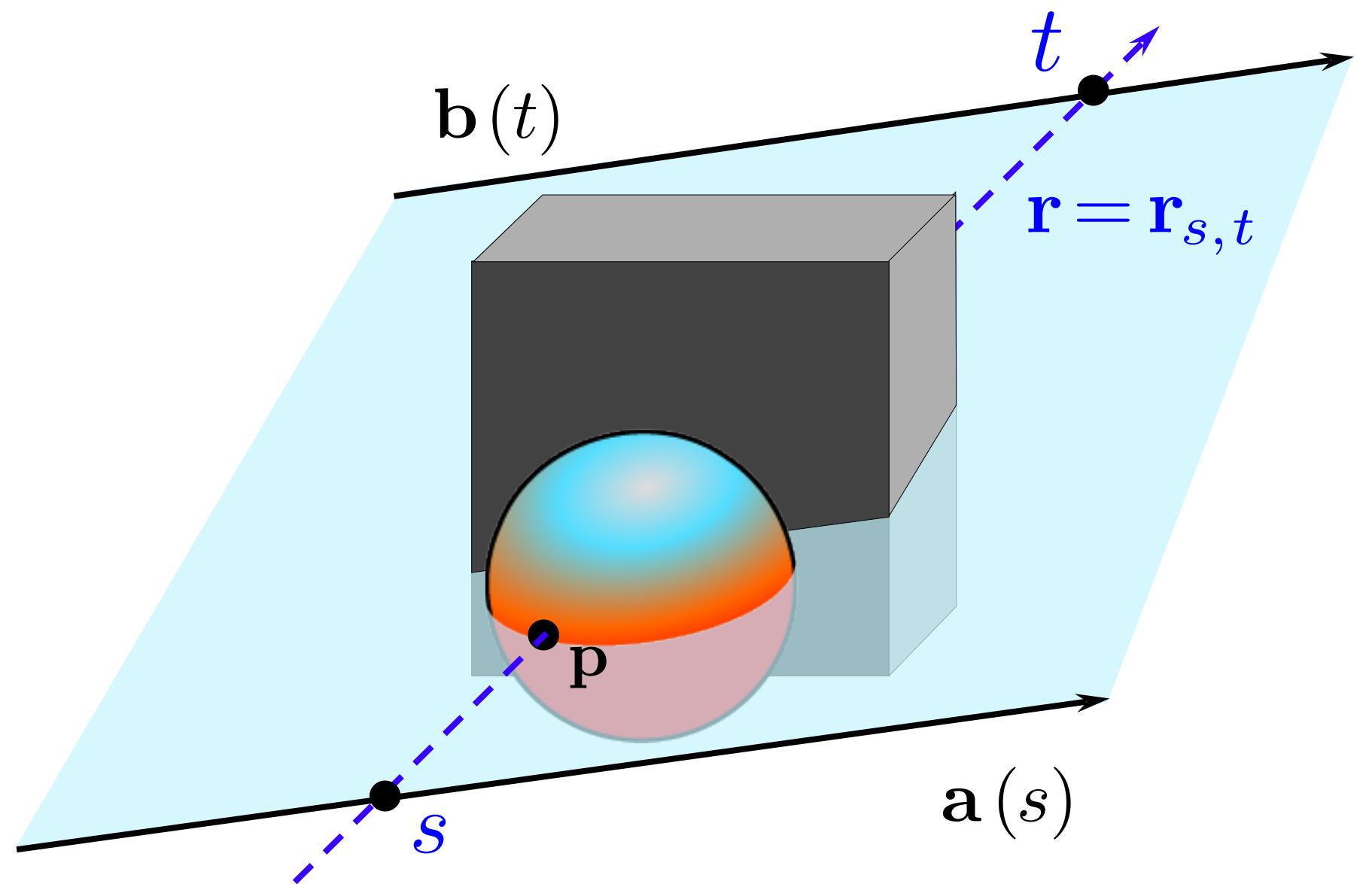
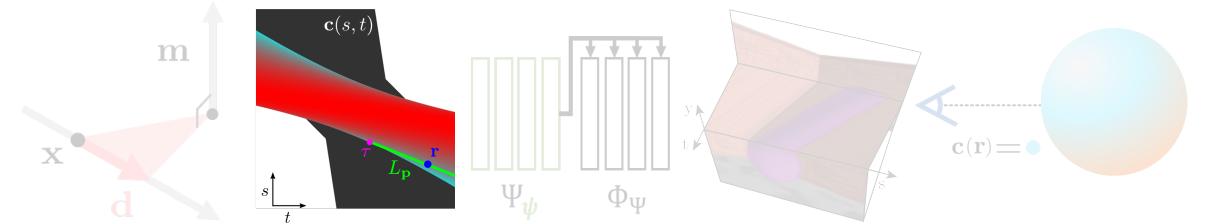


Epipolar Plane Image



Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

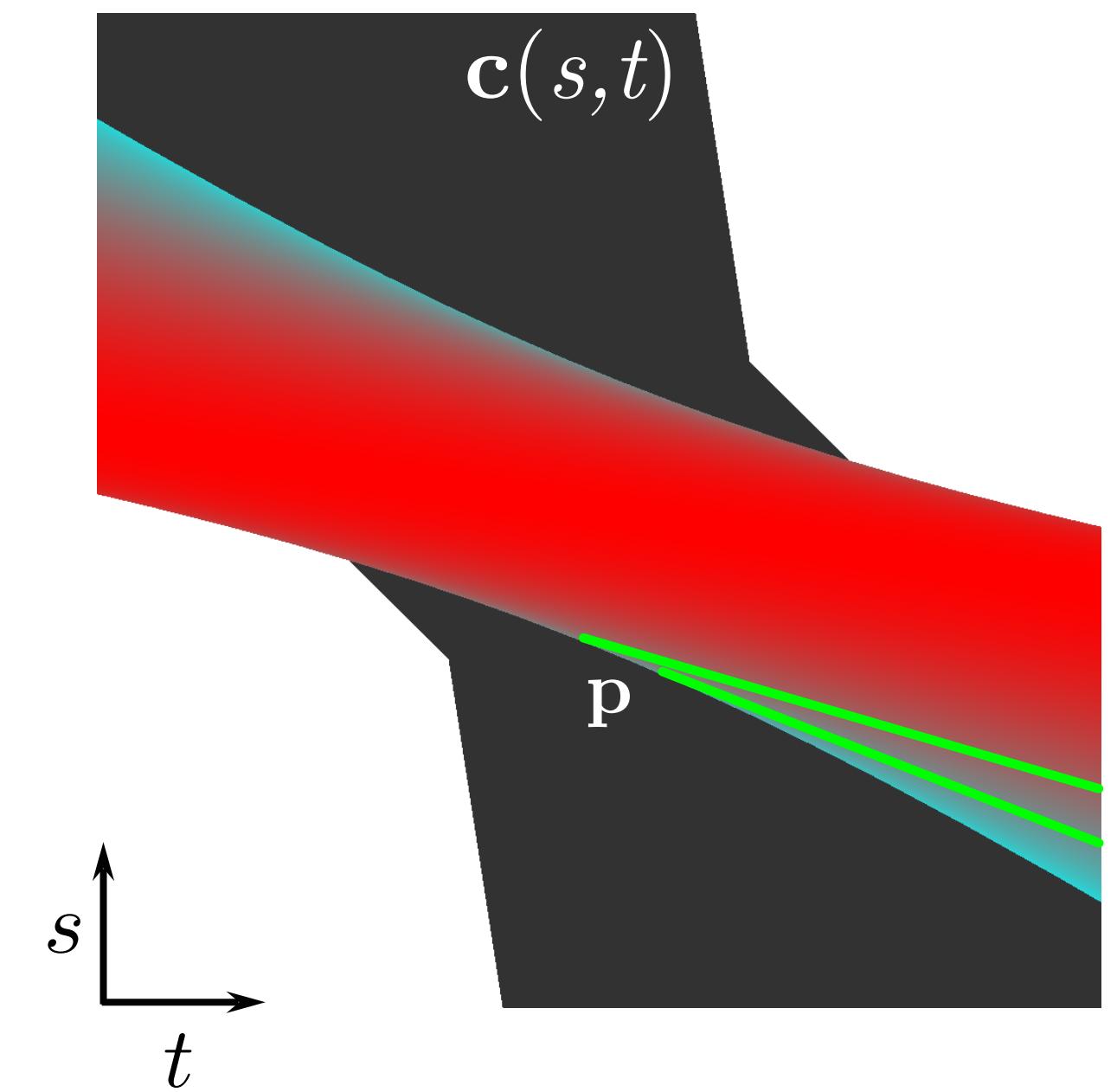
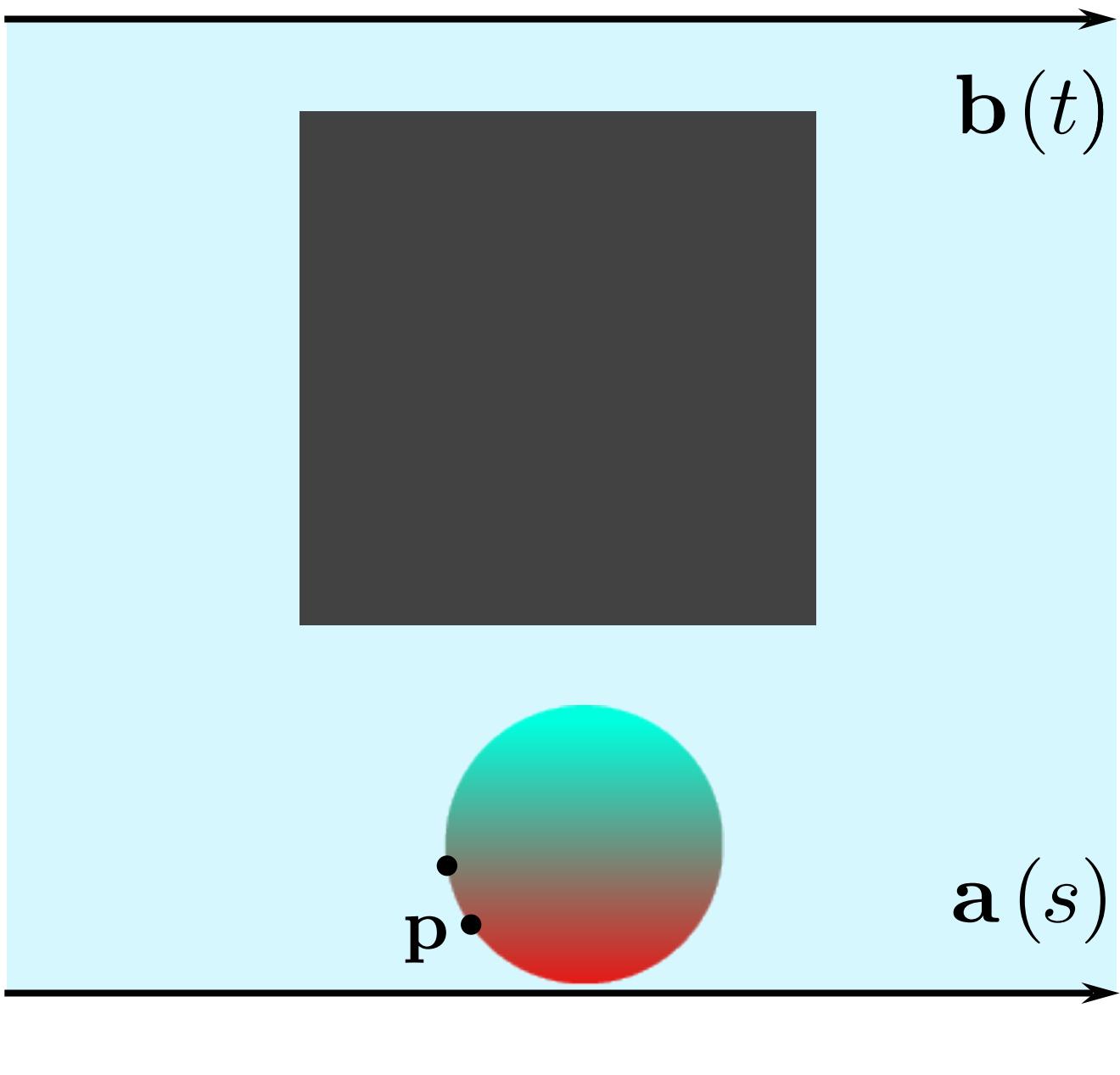
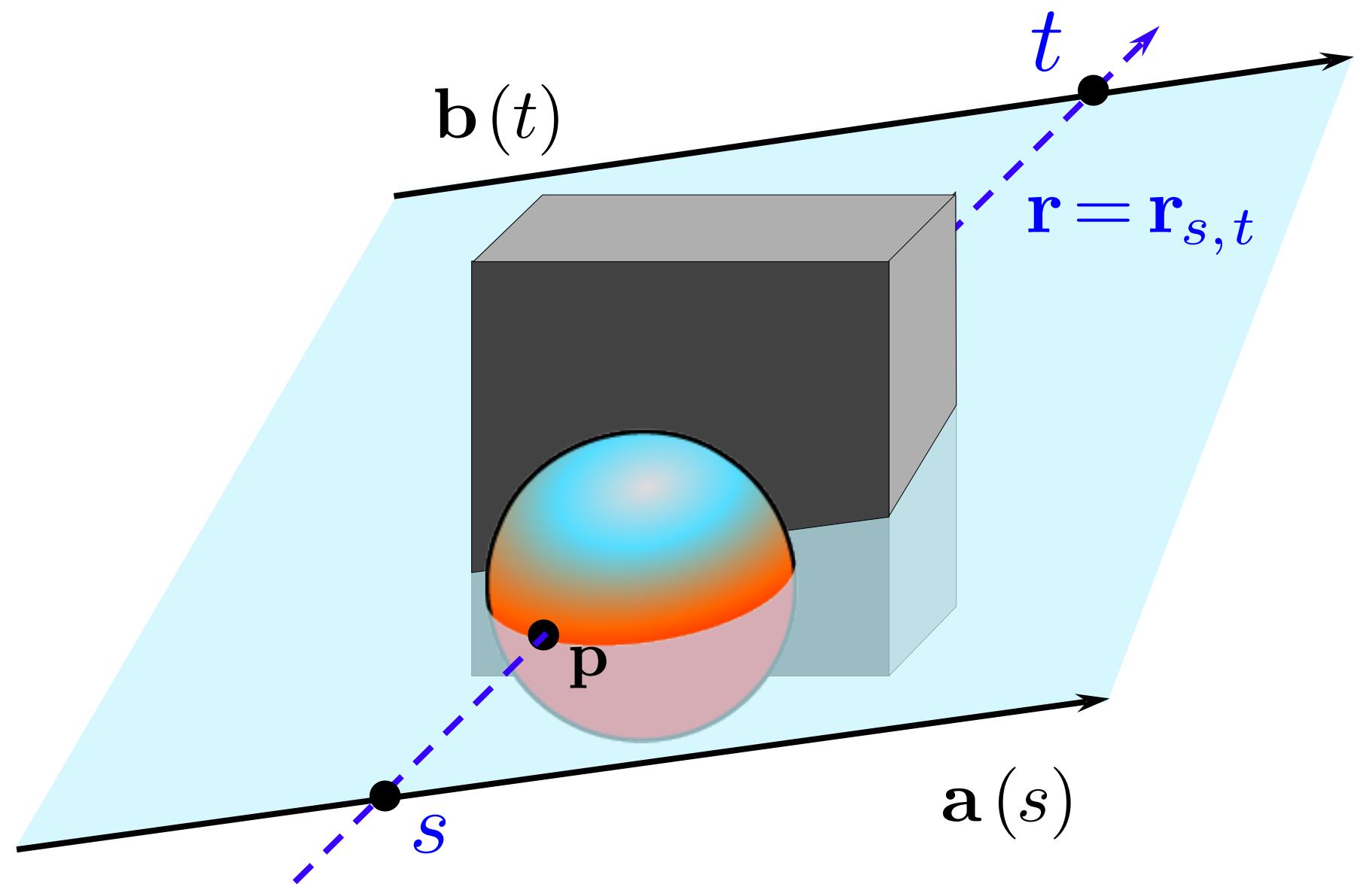
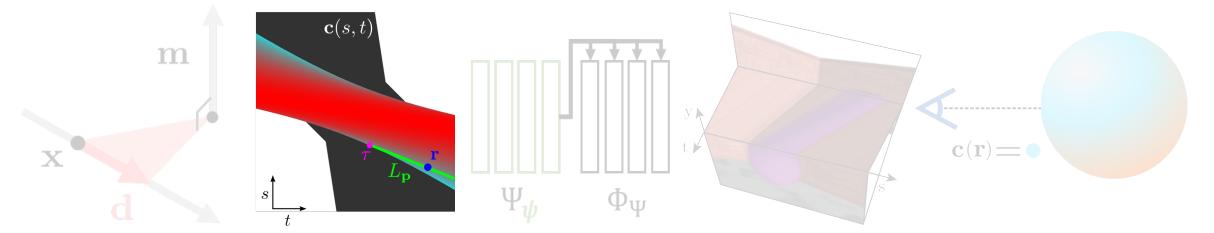
The geometry of LFNs



Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

Slope of line decreases as point moves closer.

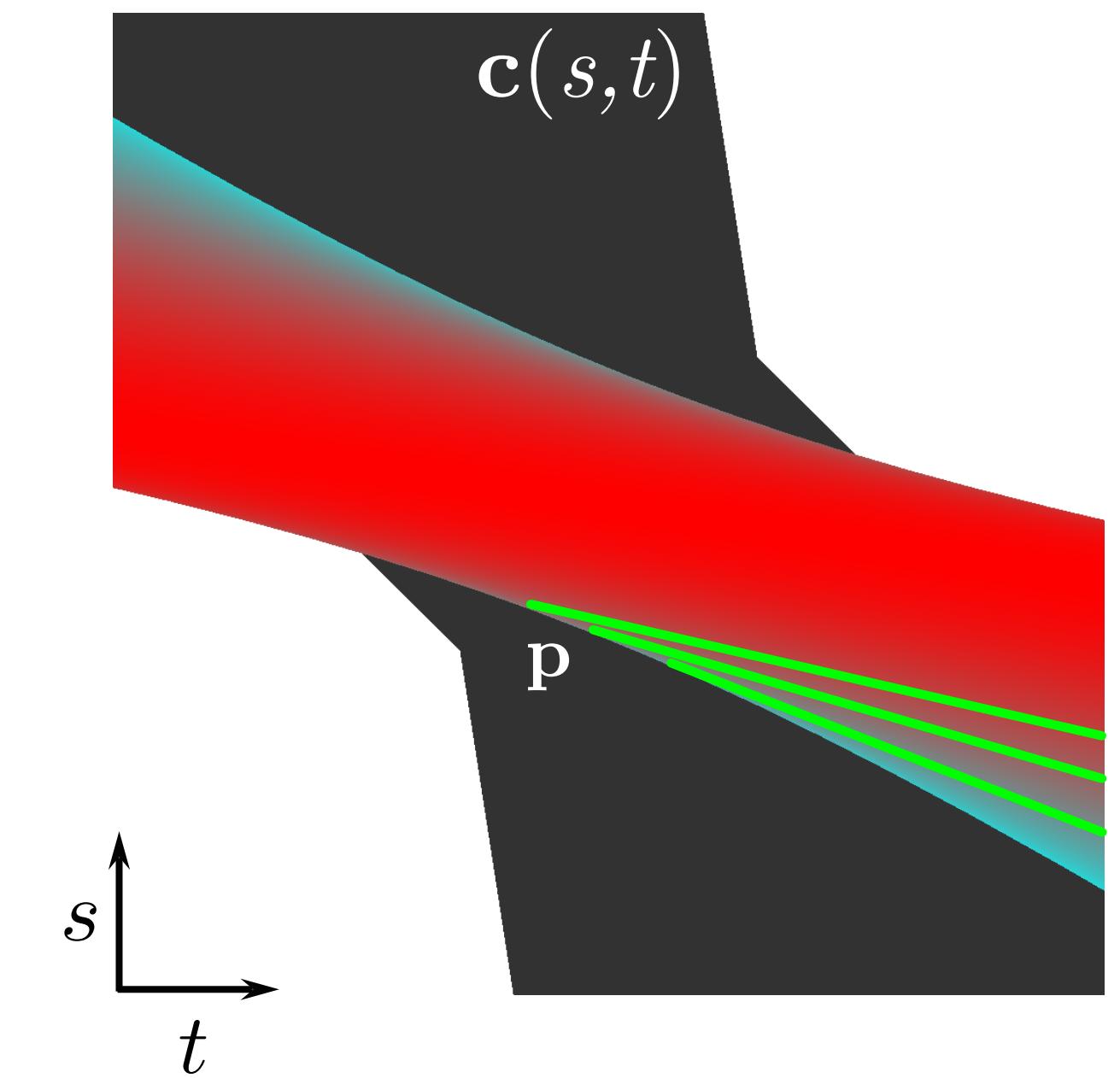
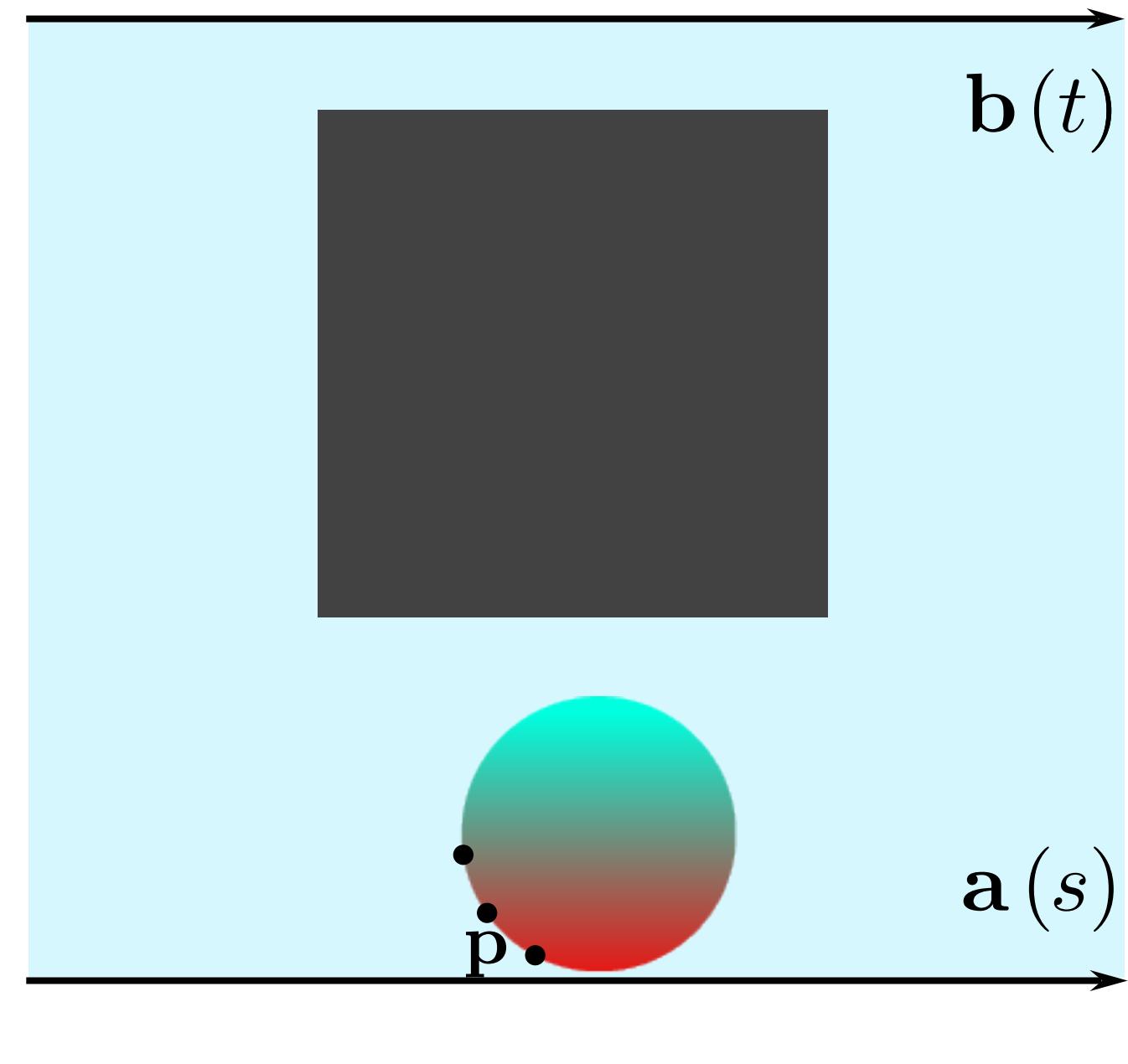
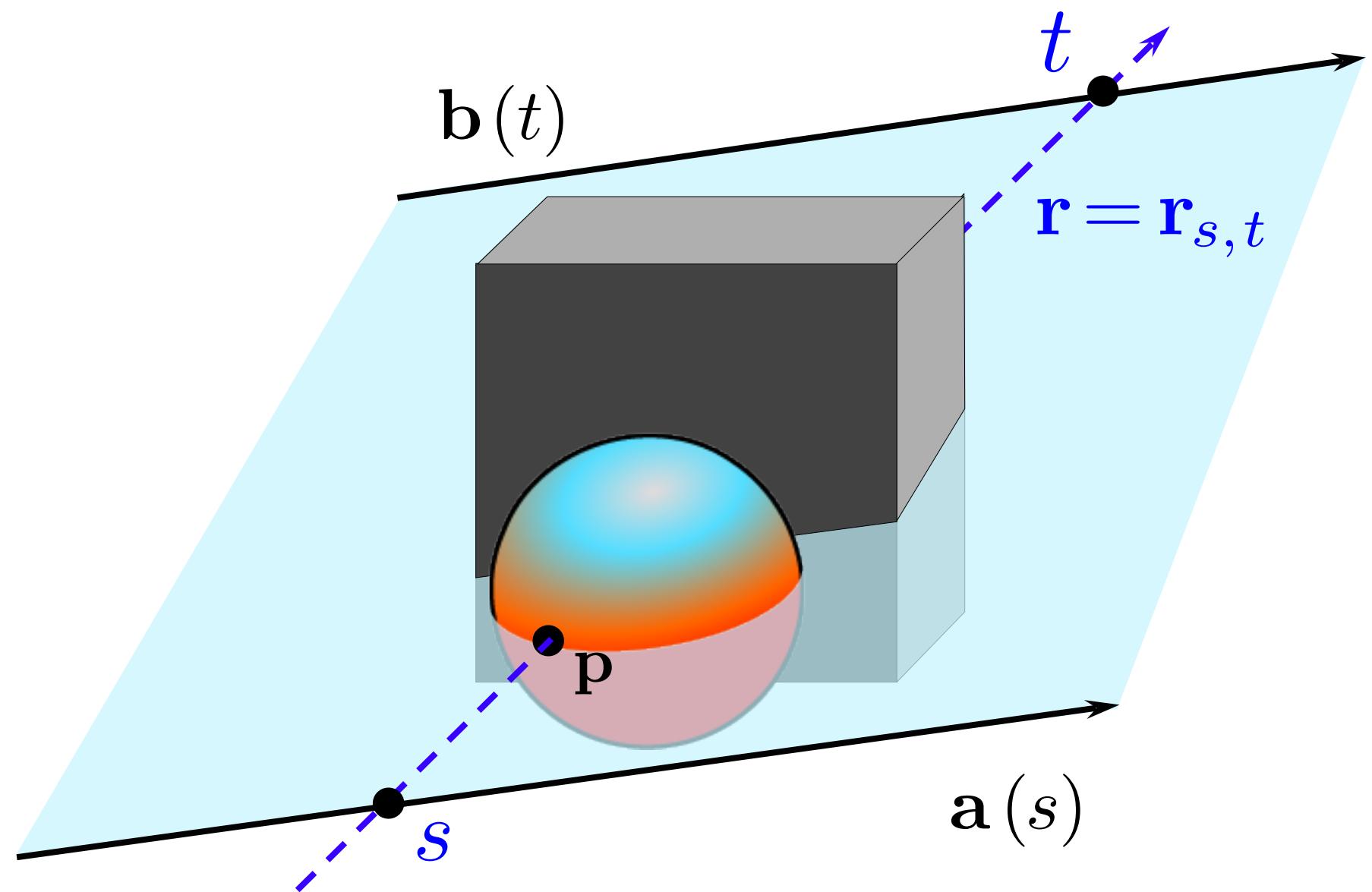
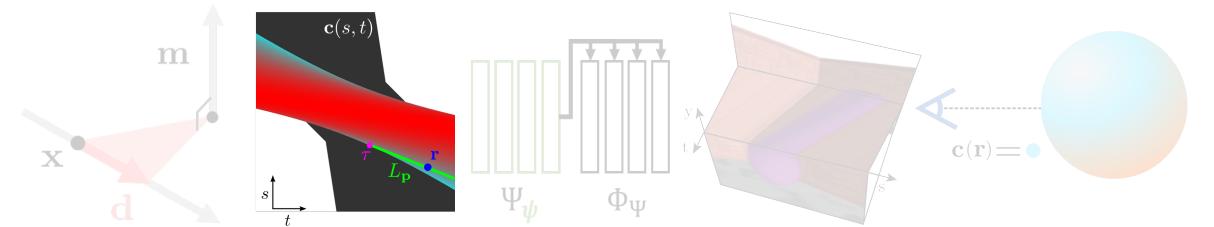
The geometry of LFNs



Points give lines of constant color in EPI $\mathbf{c}(s, t)$ – line is a levelset of the EPI.

Slope of line decreases as point moves closer.

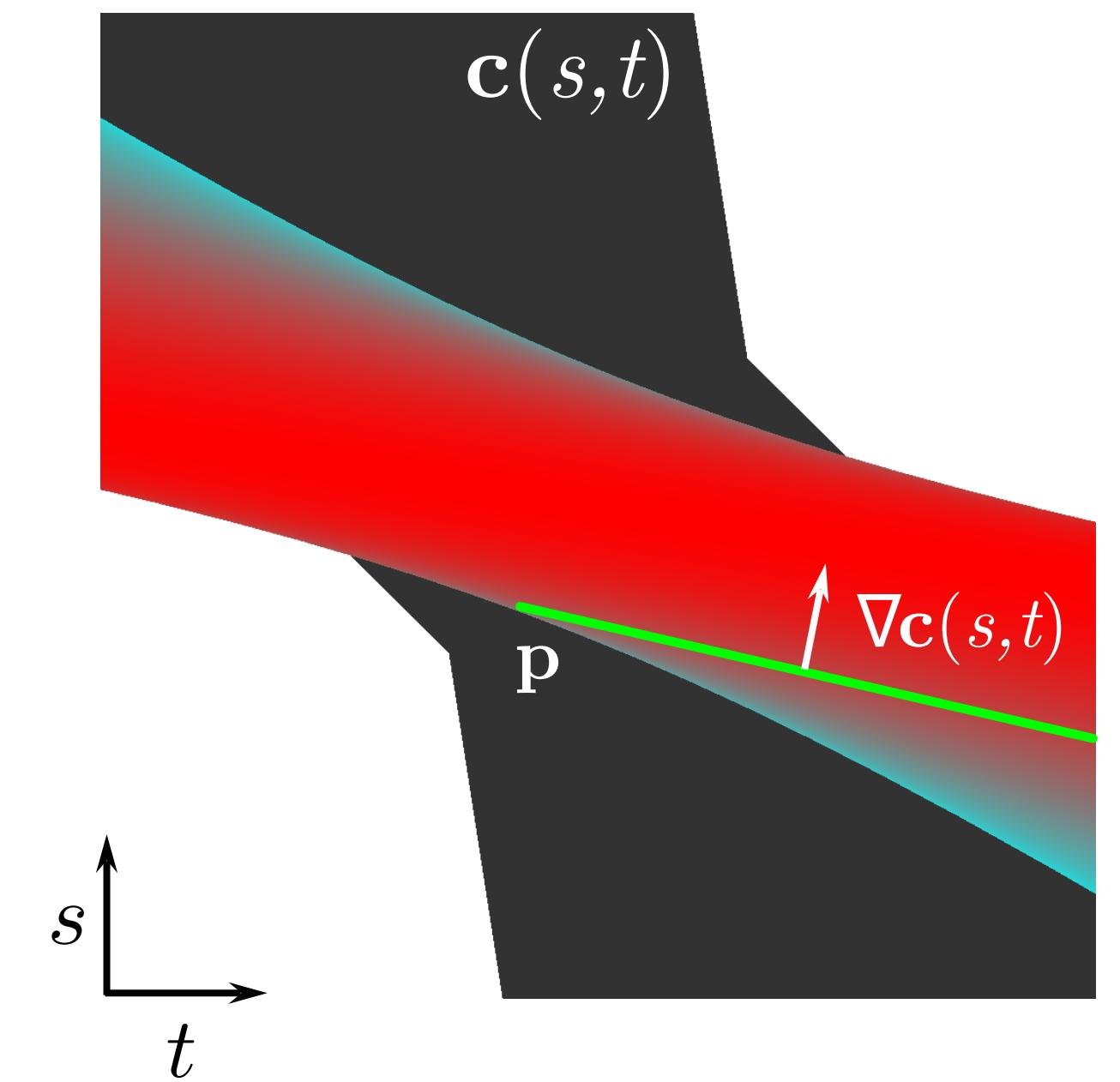
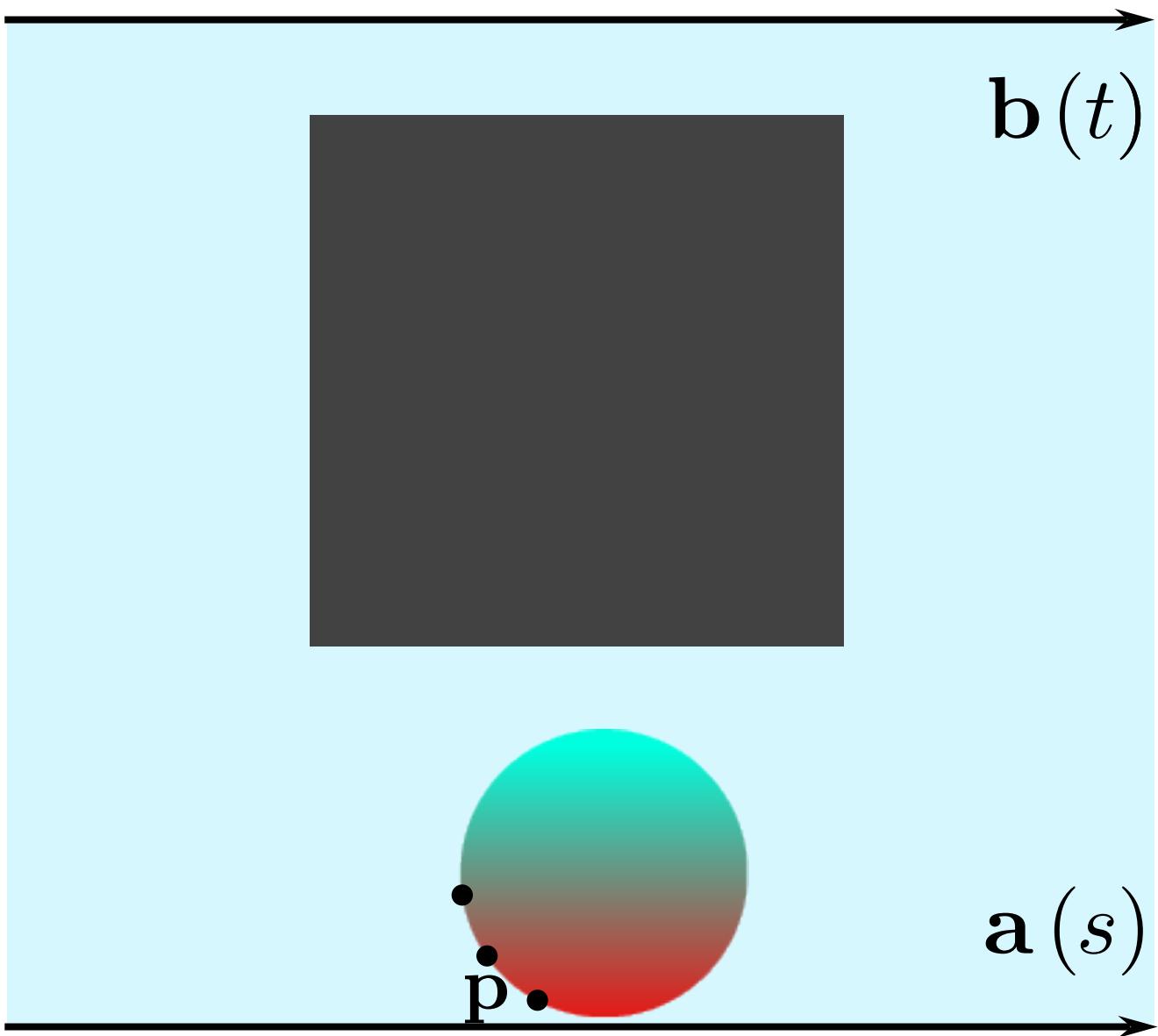
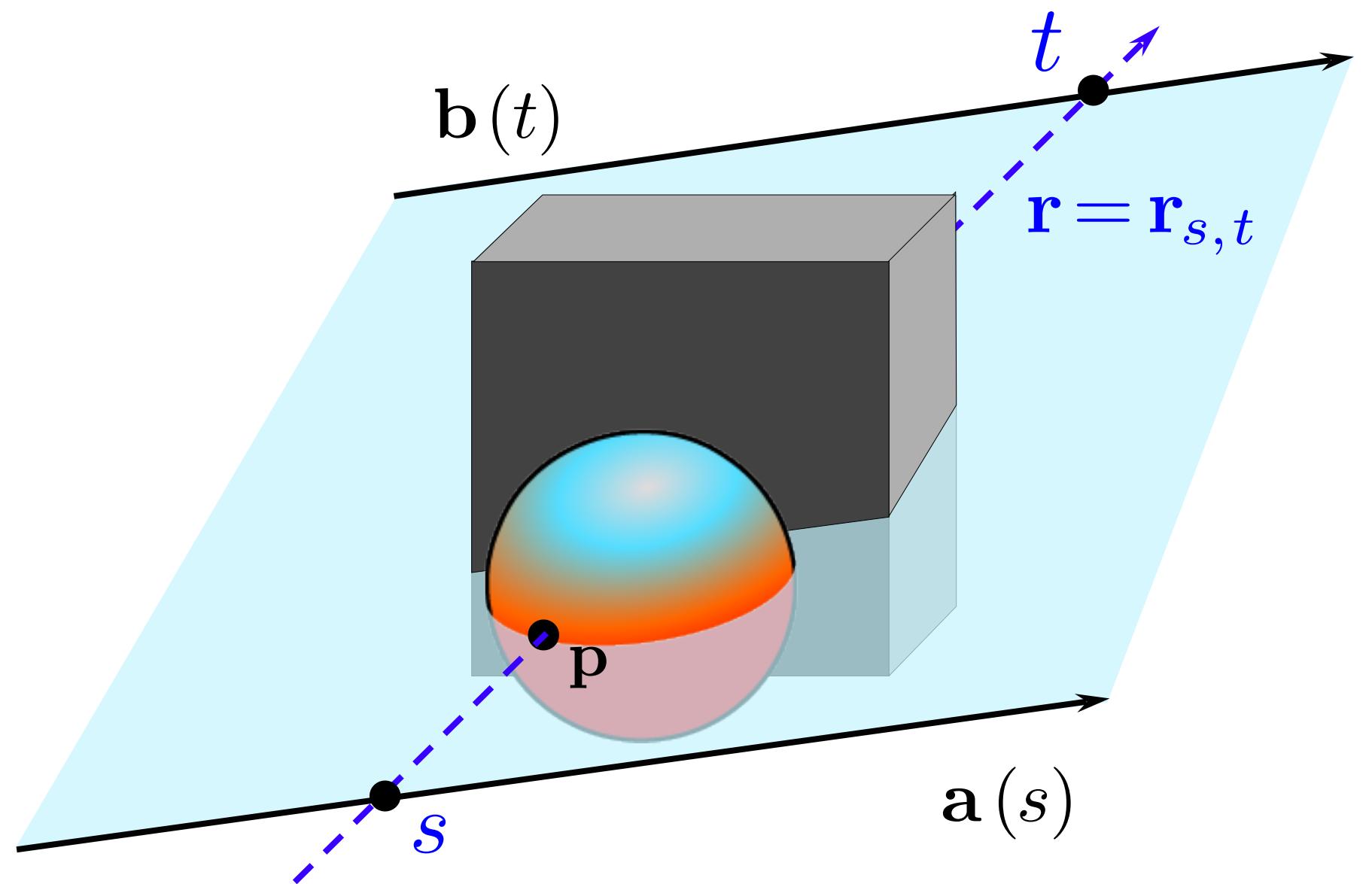
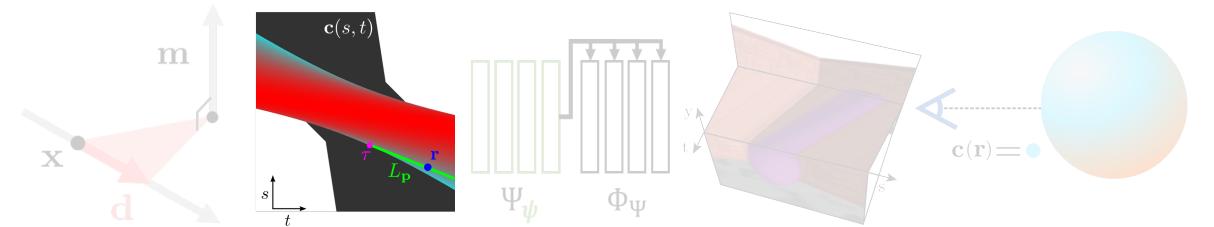
The geometry of LFNs



Points give lines of constant color in EPI $\mathbf{c}(s, t)$ – line is a levelset of the EPI.

Slope of line decreases as point moves closer.

The geometry of LFNs

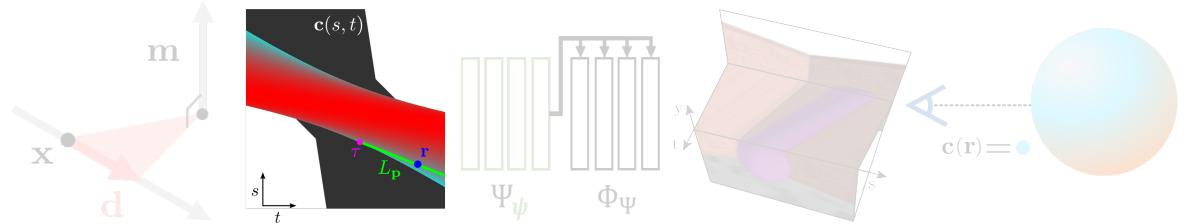


Points give lines of constant color in EPI $\mathbf{c}(s, t)$ – line is a levelset of the EPI.

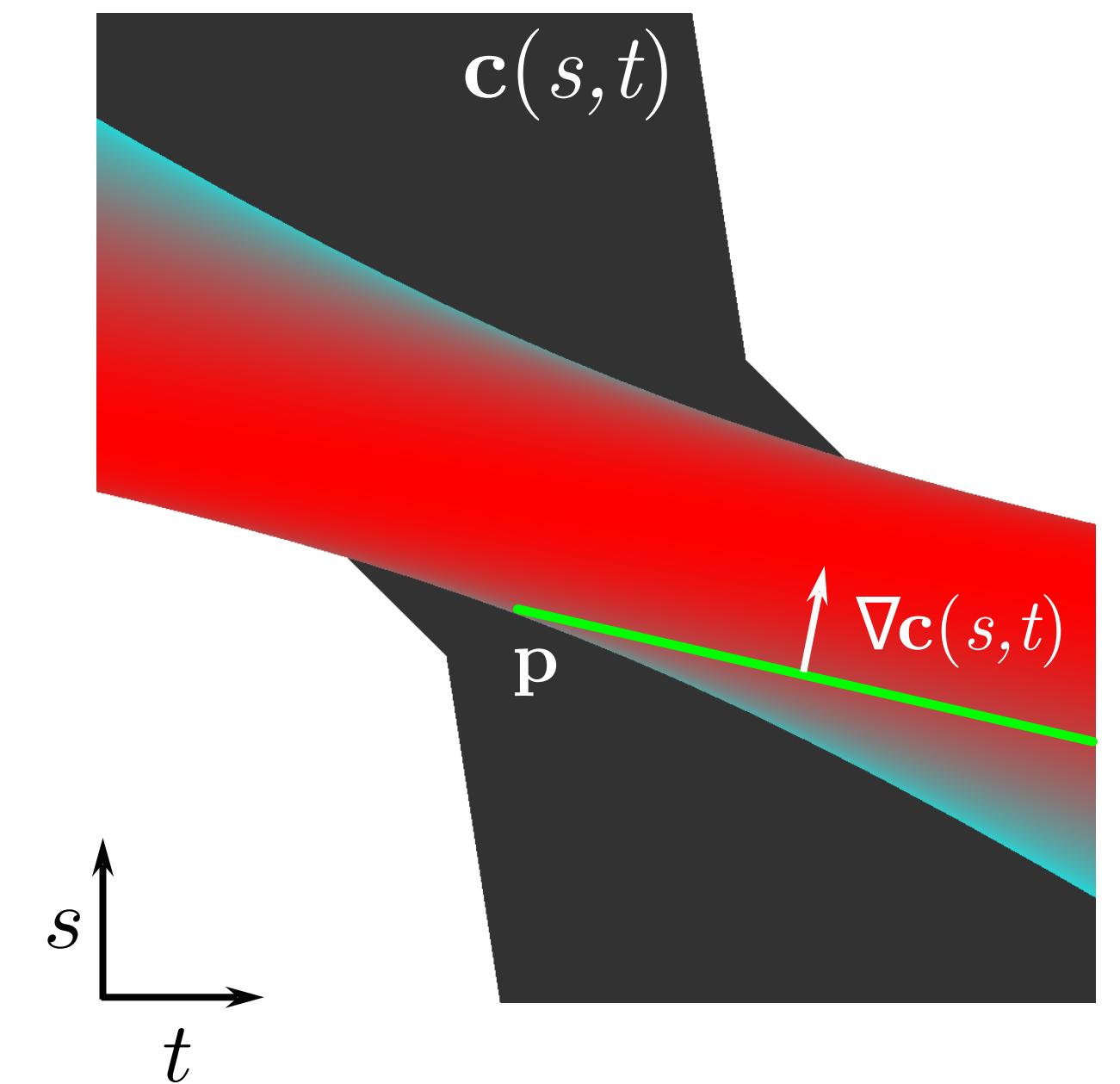
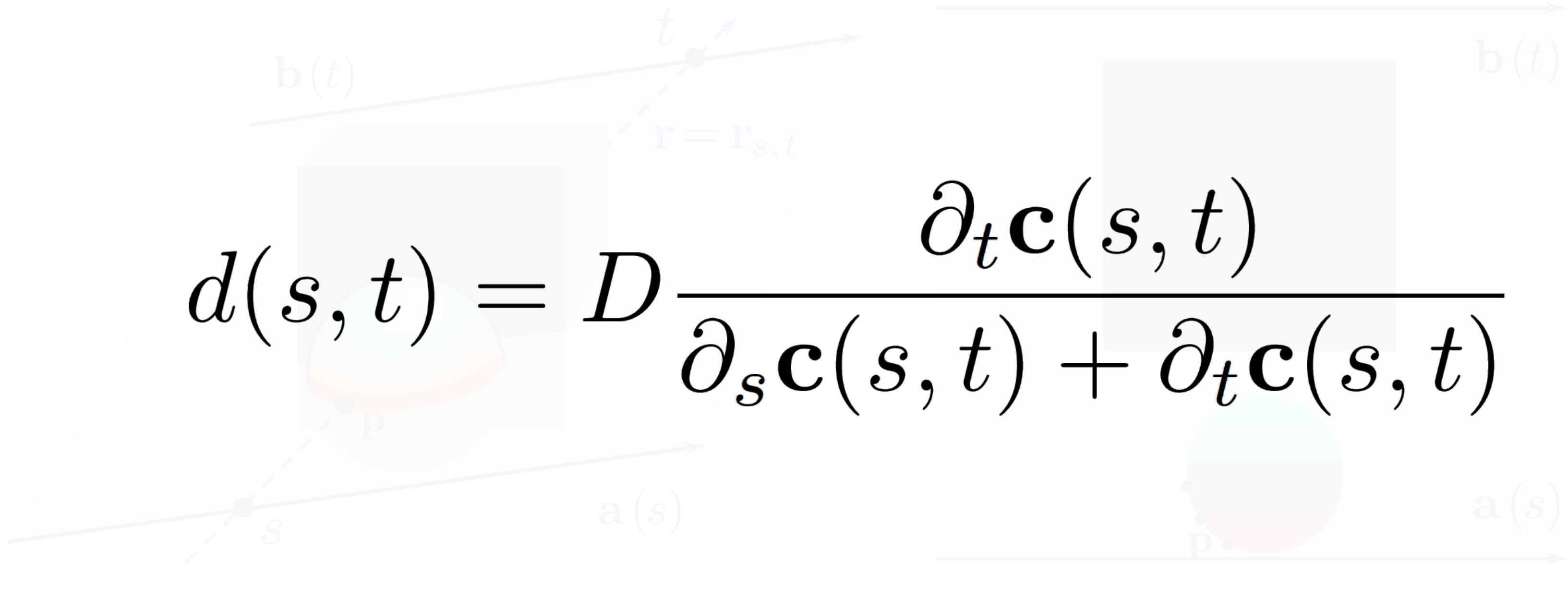
Slope of line decreases as point moves closer.

Gradient of $\mathbf{c}(s, t)$ is orthogonal to levelset -

The geometry of LFNs



Epipolar Plane Image

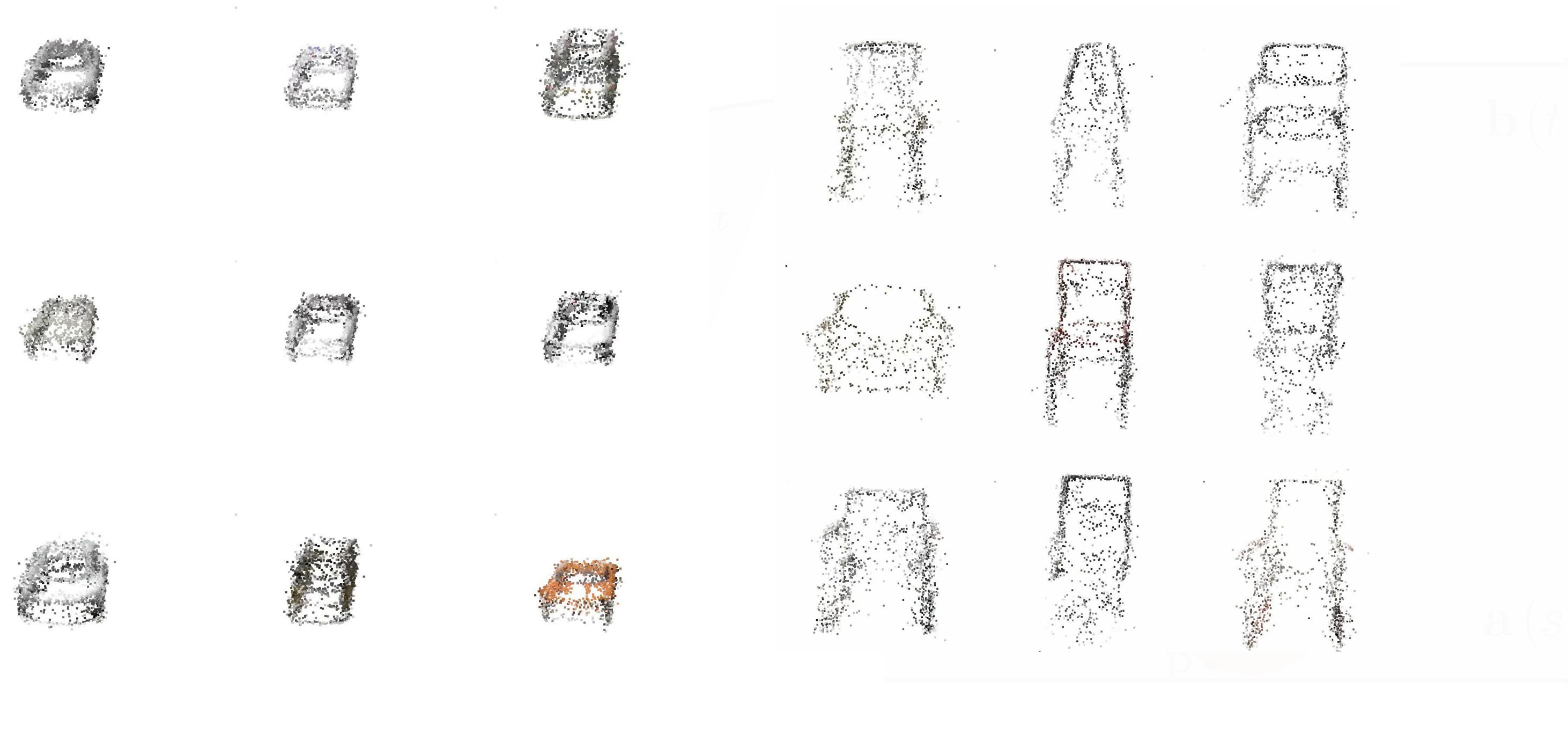
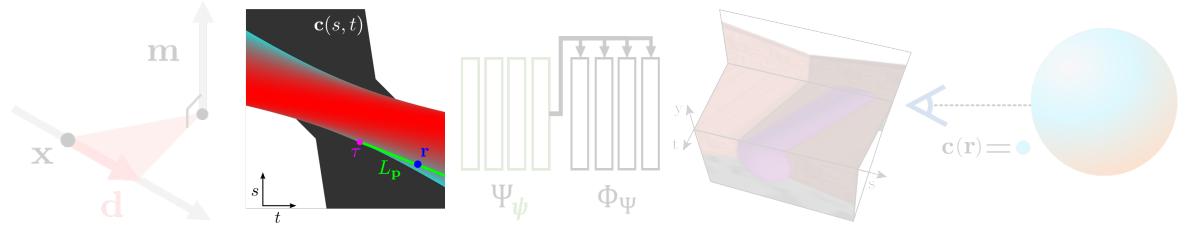


Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

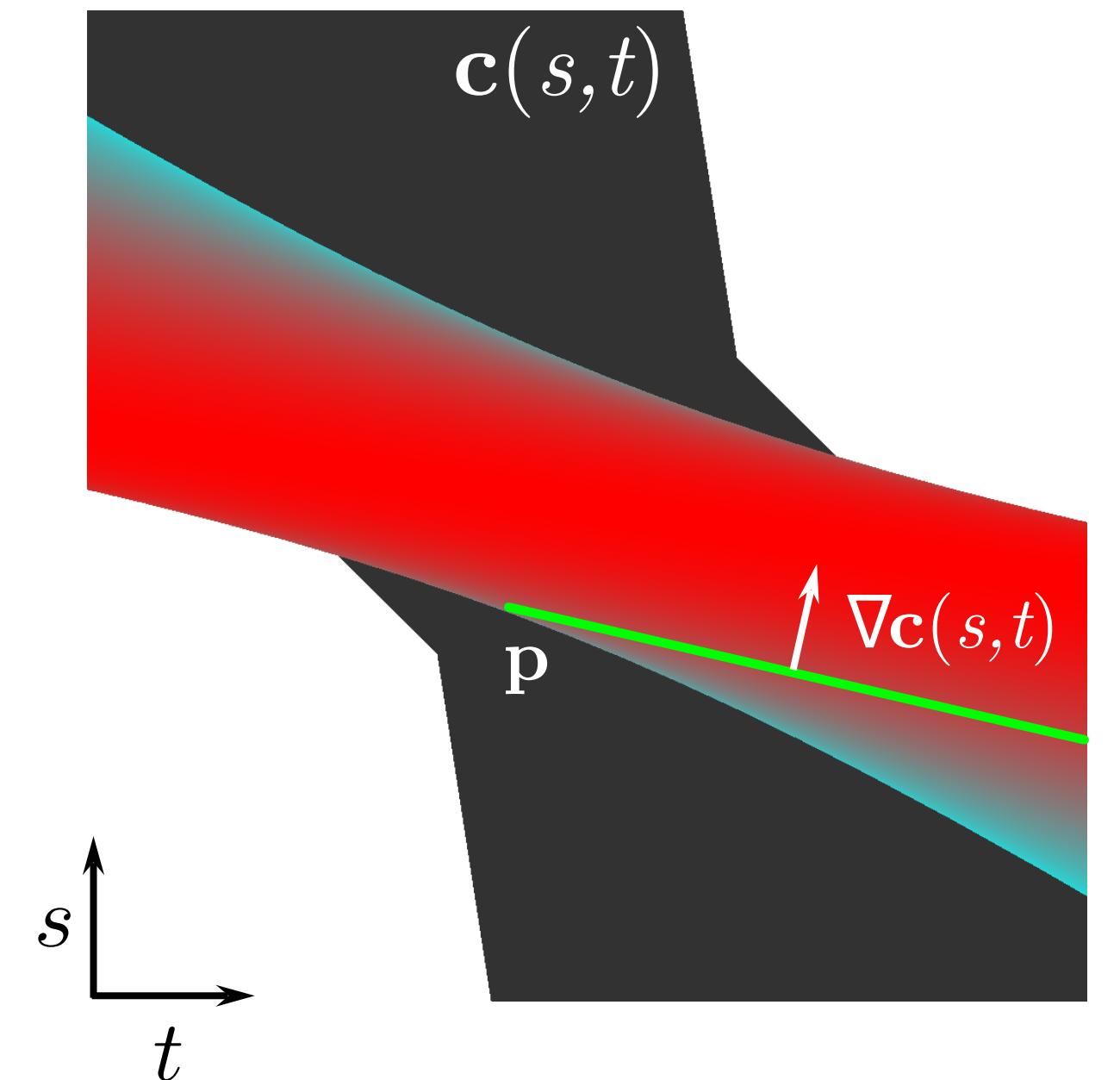
Slope of line decreases as point moves closer.

Gradient of $\mathbf{c}(s,t)$ is orthogonal to levelset - can extract depth from gradients of light field.

The geometry of LFNs



Epipolar Plane Image

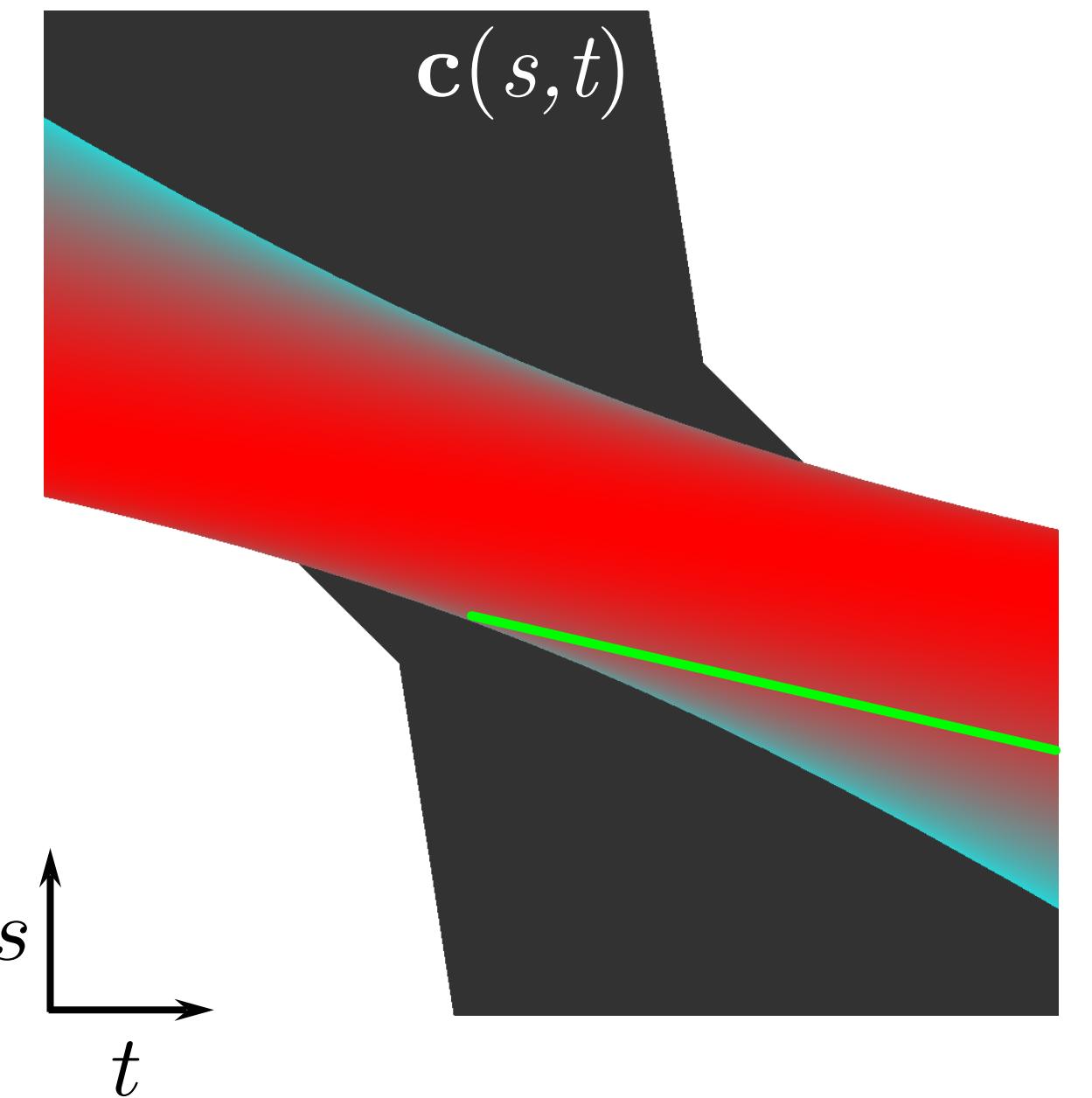
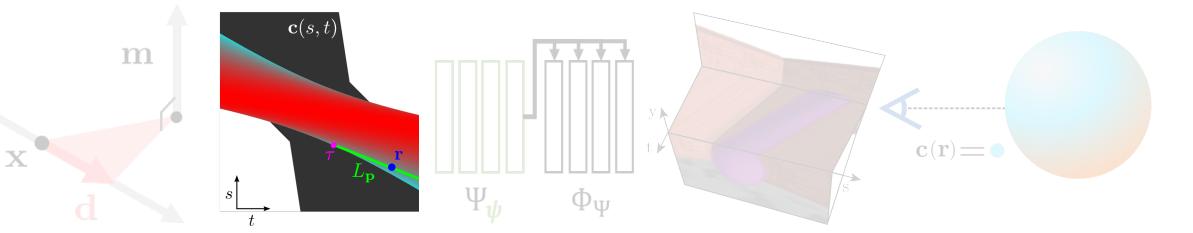


Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

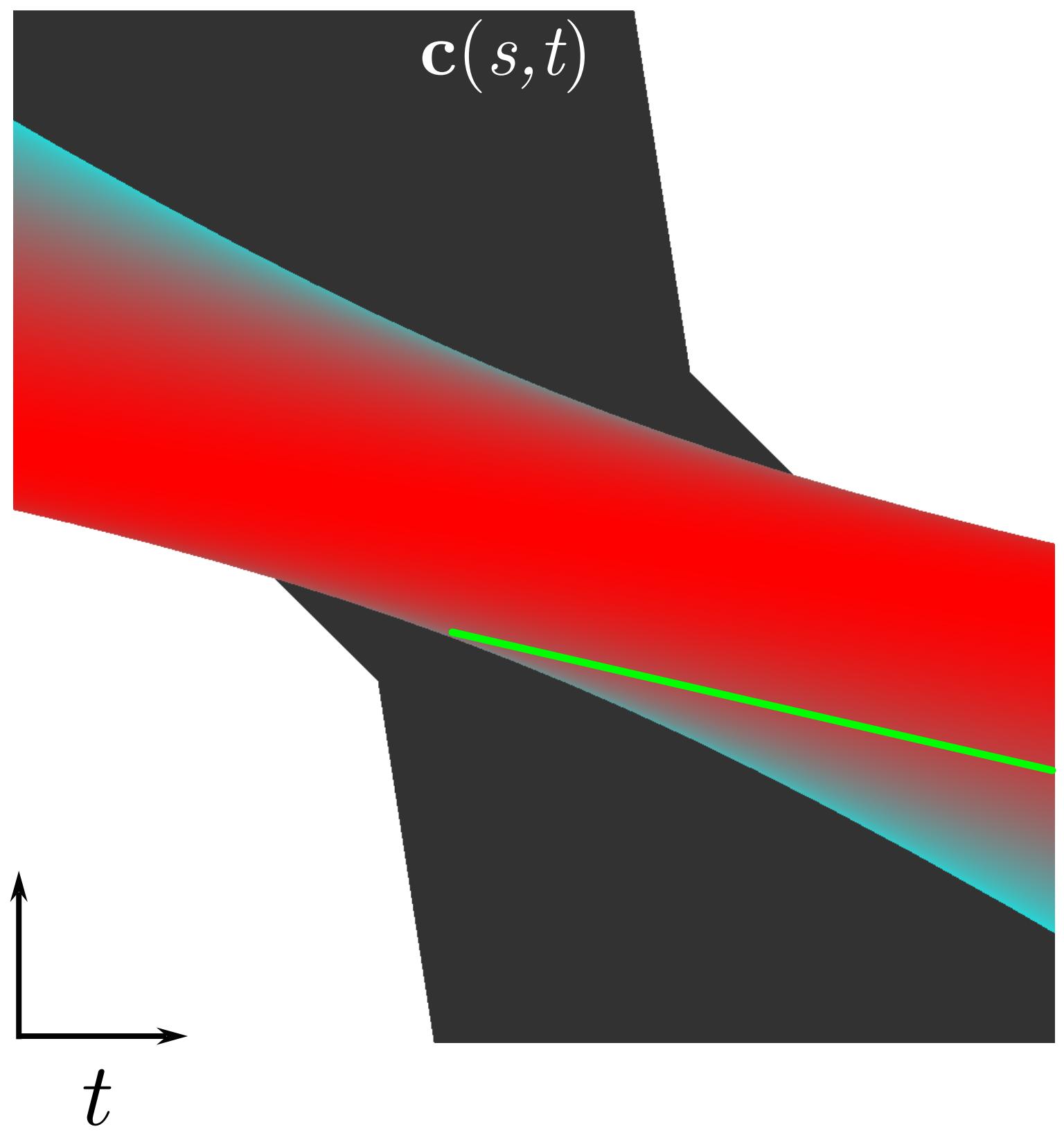
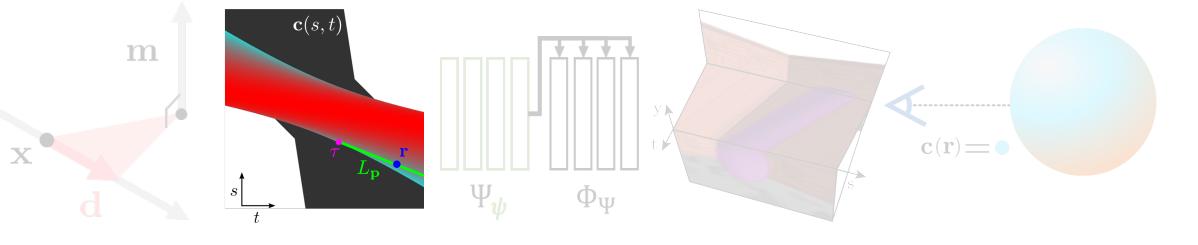
Slope of line decreases as point moves closer.

Gradient of $\mathbf{c}(s,t)$ is orthogonal to levelset - can extract depth from gradients of light field.

Multi-view consistency



Multi-view consistency



Light Field Networks
500 FPS
1 evaluation per ray

Volumetric Rendering (pixelNeRF)
0.033 FPS
196 evaluations per ray

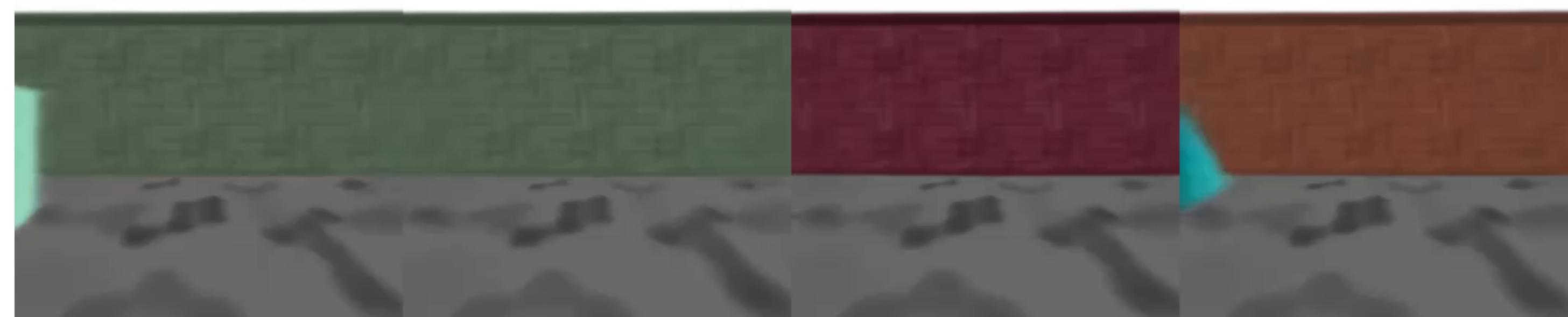


1x speed

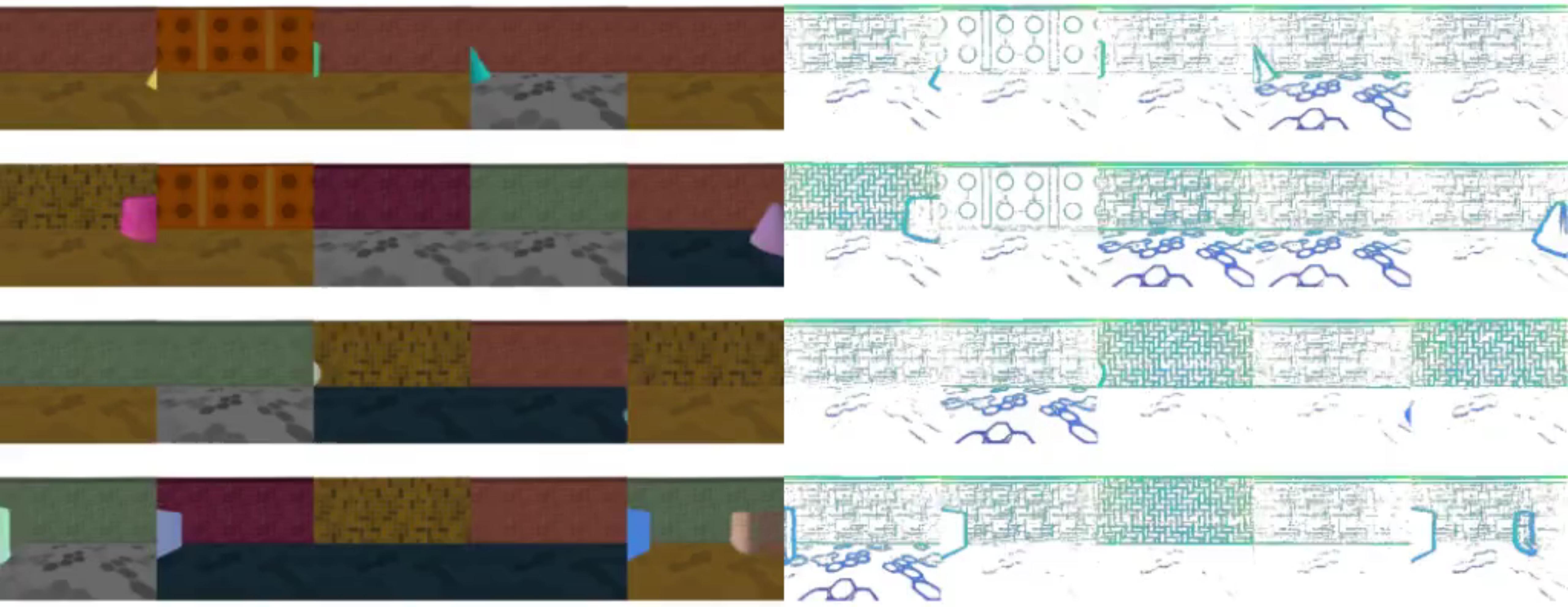


Real-time. No post-processing, no discrete data structures (octrees, voxelgrids, ...).
>100x reduction in memory: Can be trained on small GPUs!

Light Field Networks
500 FPS
1 evaluation per ray



Also Encode Depth in their 4D derivatives:
can be extracted via single evaluation of neural network and its gradient!



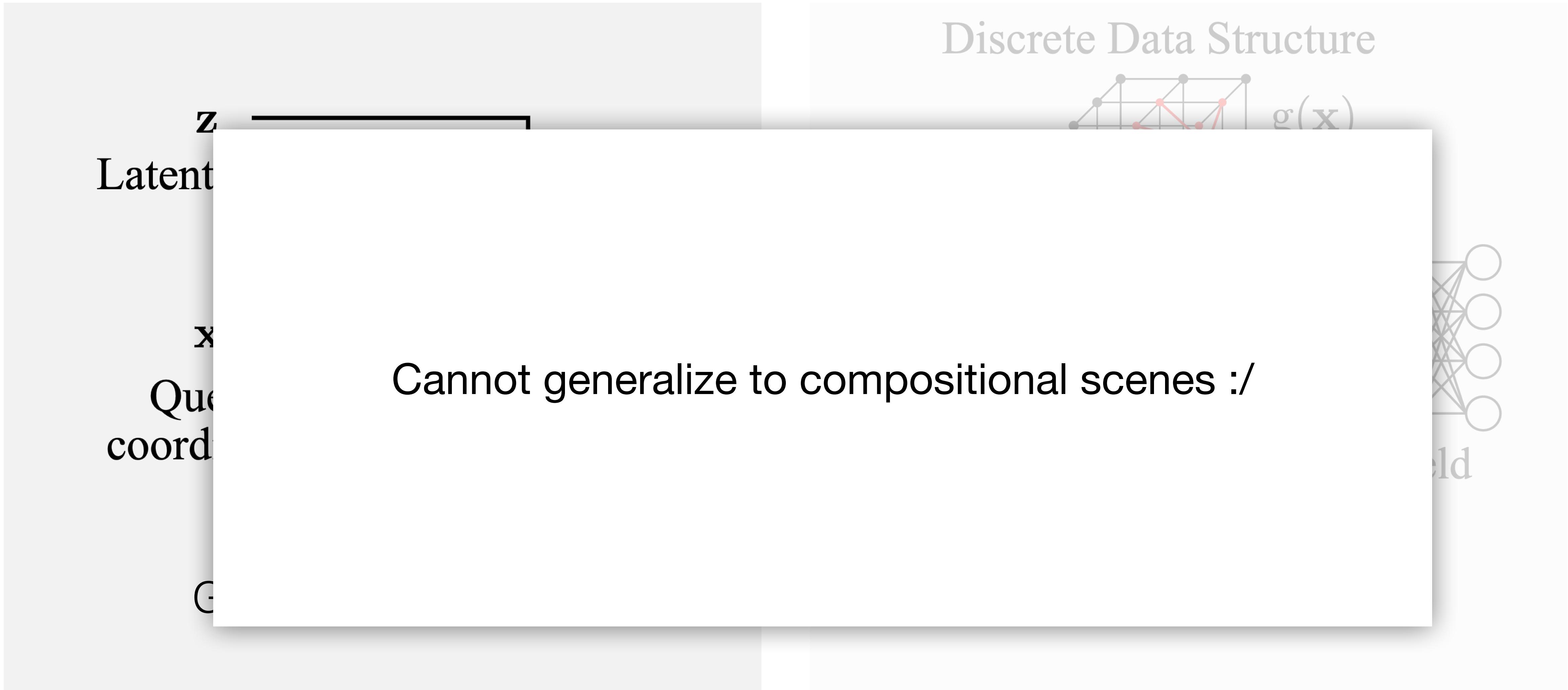
Outperforms global conditioning

	Chairs	Cars
SRNs [3]	22.89 / 0.89	22.25 / 0.89
LFN	22.26 / 0.90	22.42 / 0.89

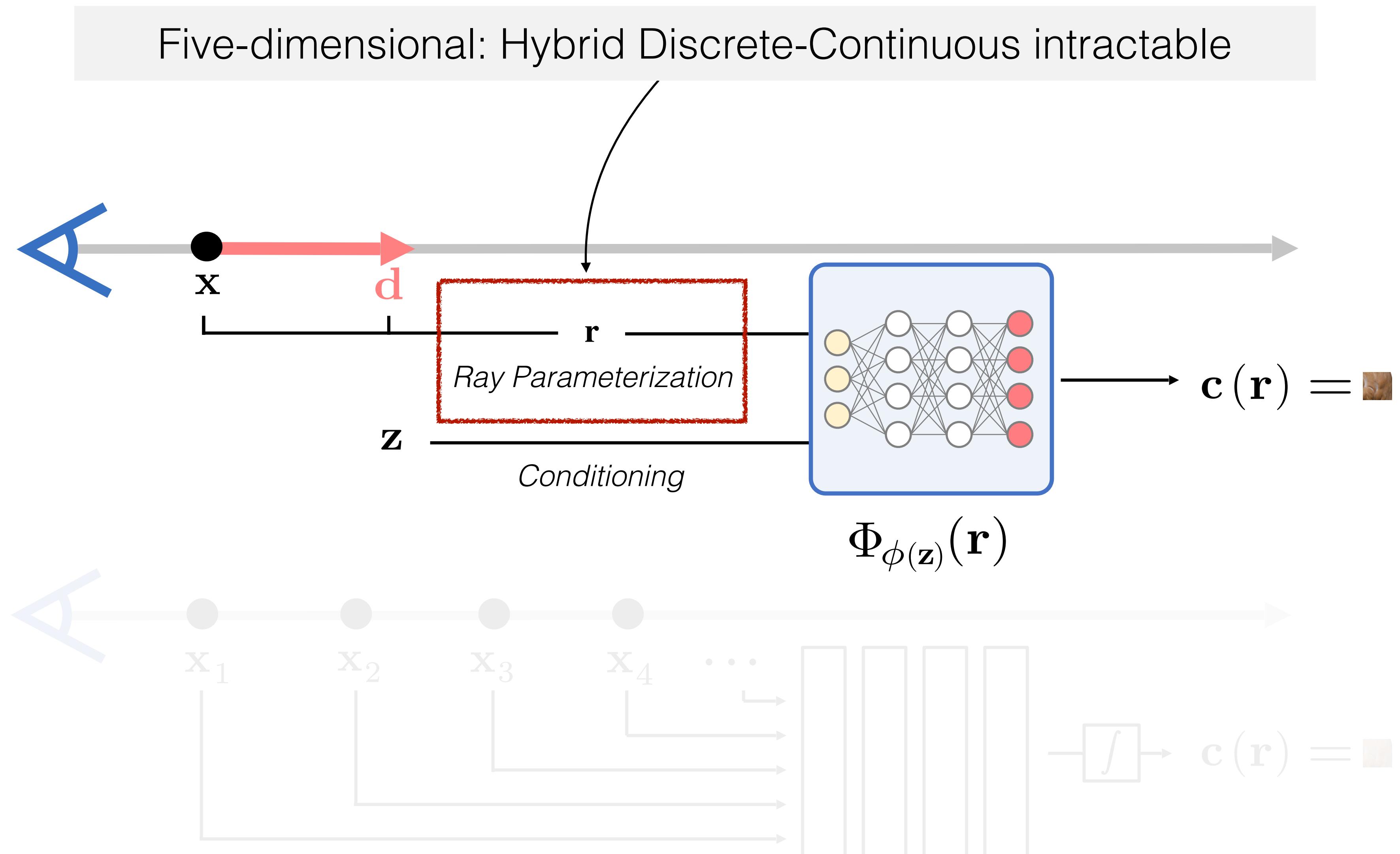
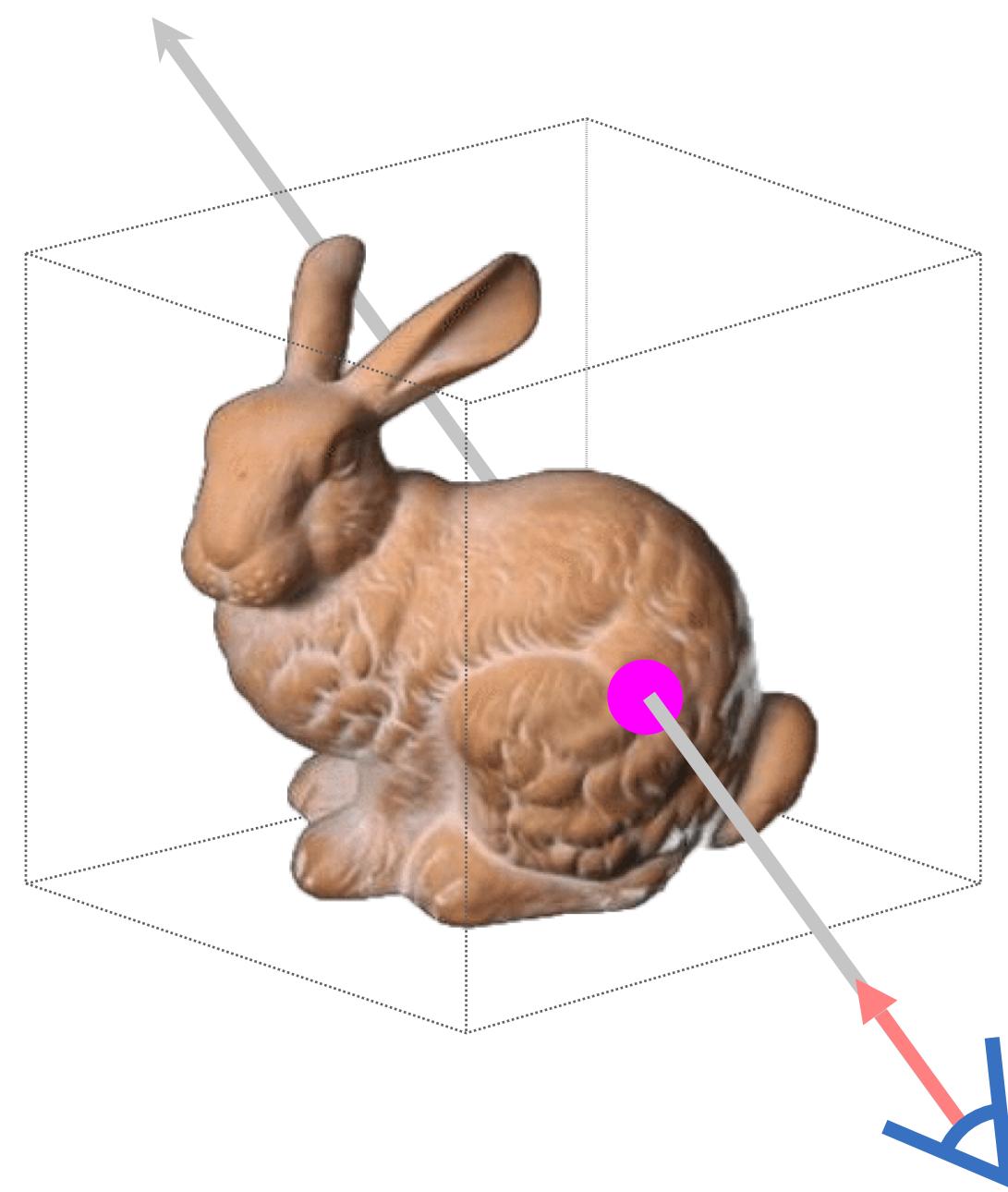
Does not outperform local conditioning

	Class-specific	Multi-Class
LFN	22.34 / 0.90	24.95 / 0.87
pixelNeRF	23.45 / 0.91	26.80 / 0.91

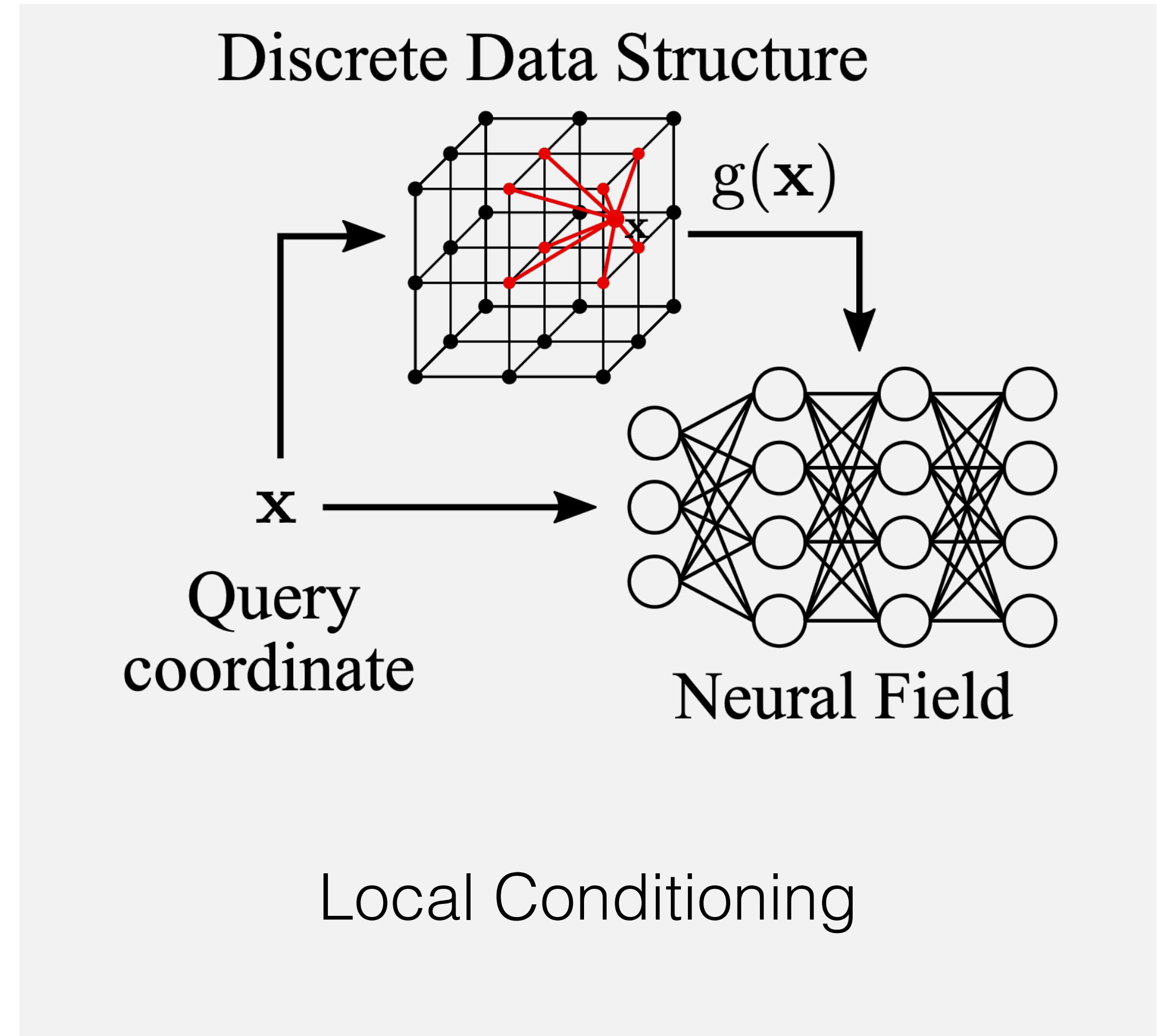
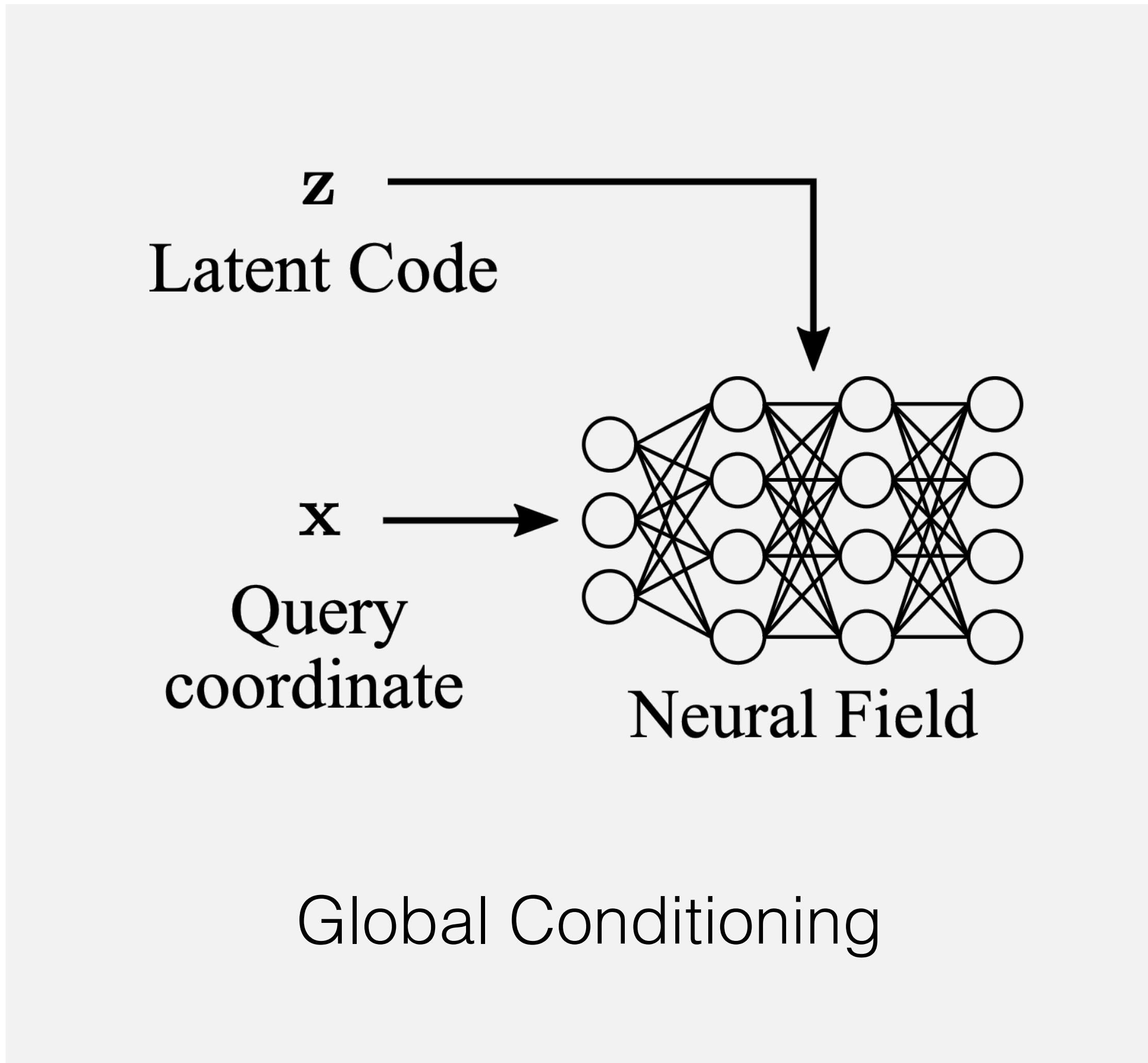
Global Conditioning: Single Latent Code for whole 3D Scene



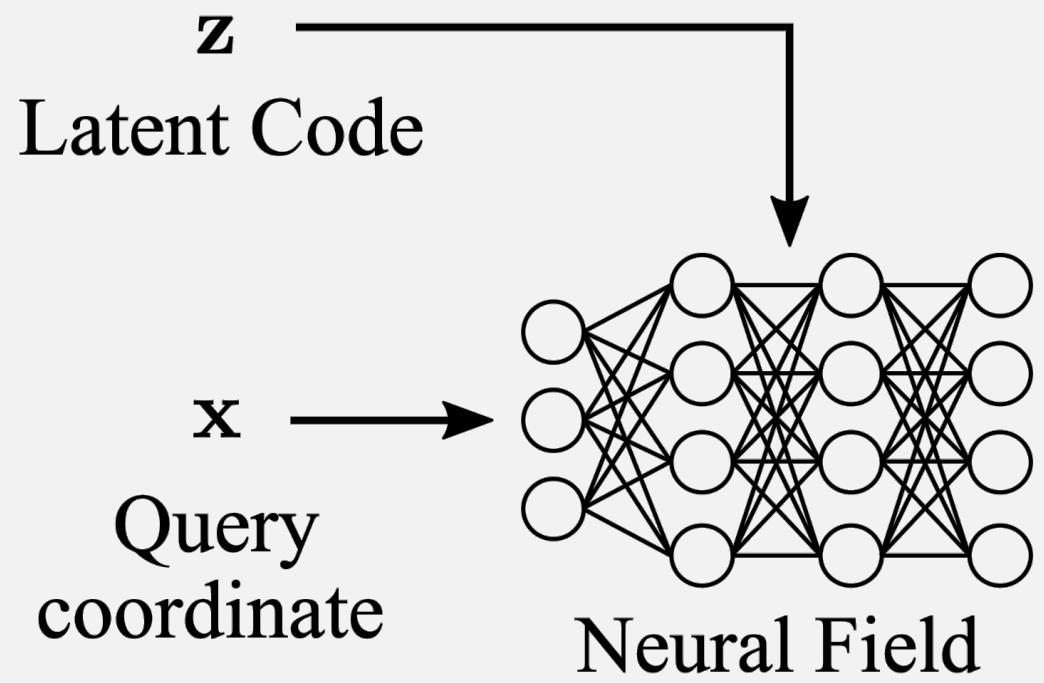
Now: Learn Prior over Light Fields!



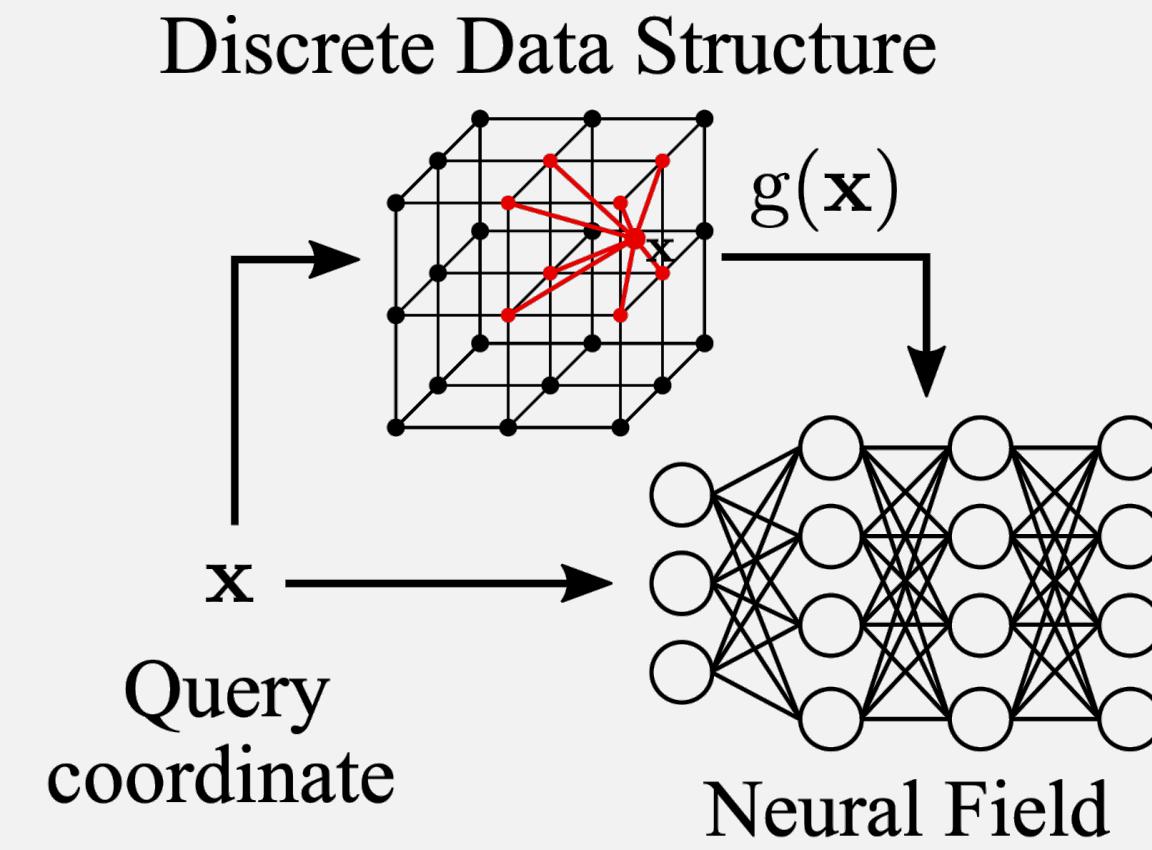
Recap: Global Conditioning and Local Conditioning



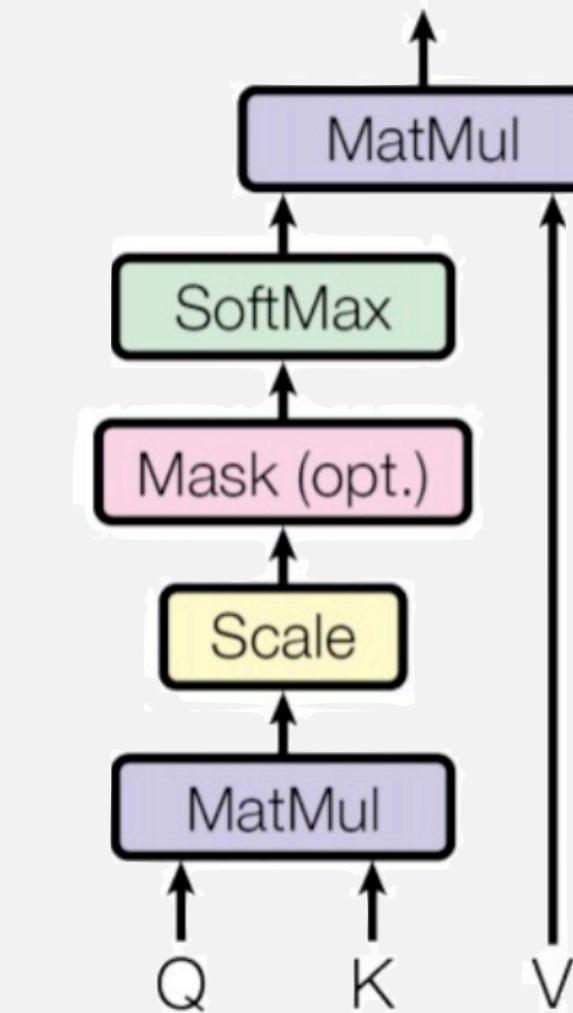
Other types of conditioning...?



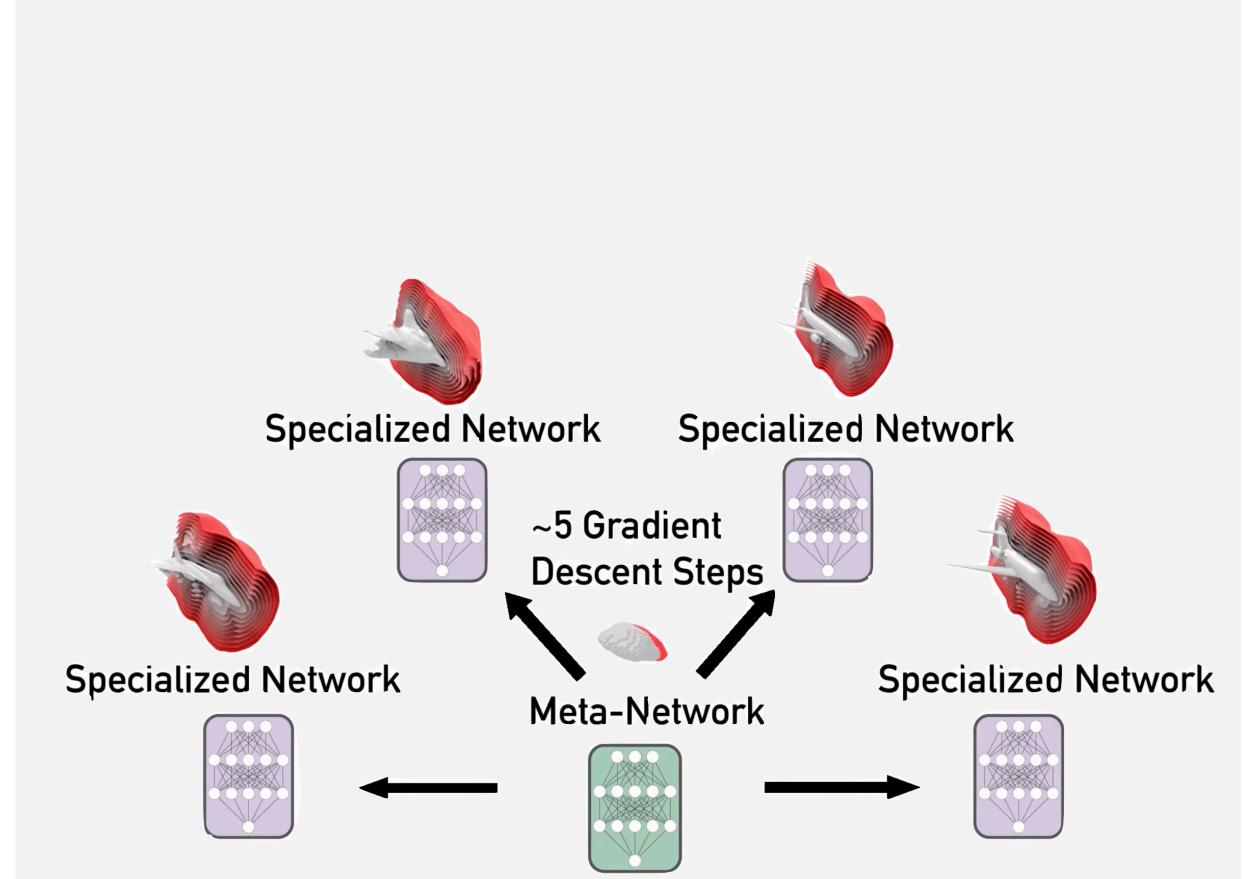
Global
Conditioning



Local
Conditioning

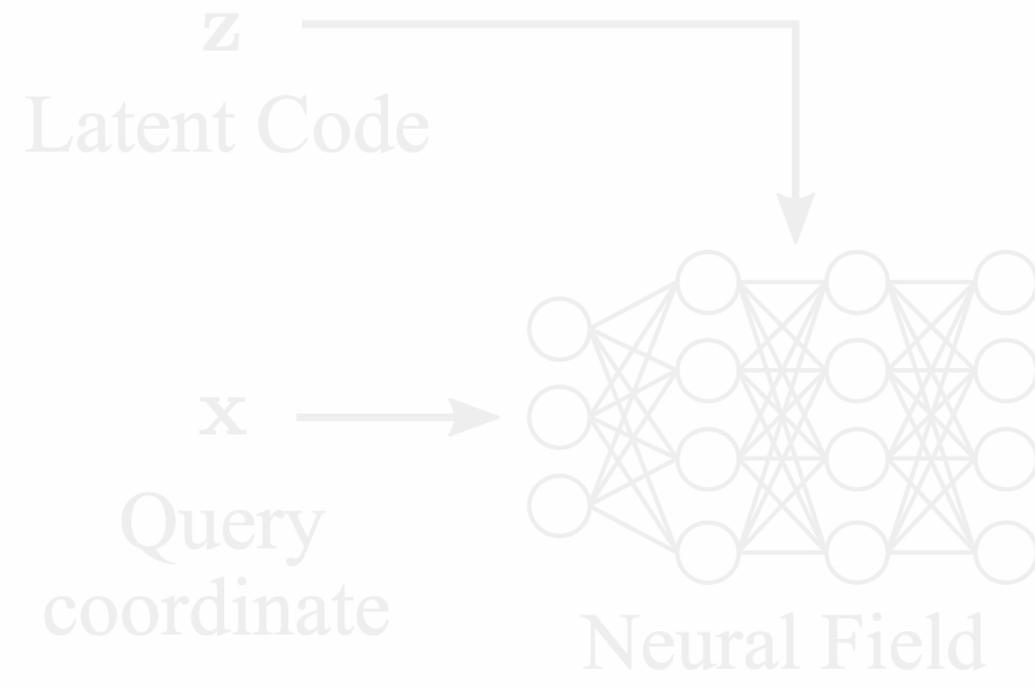


Attention-Based
Conditioning

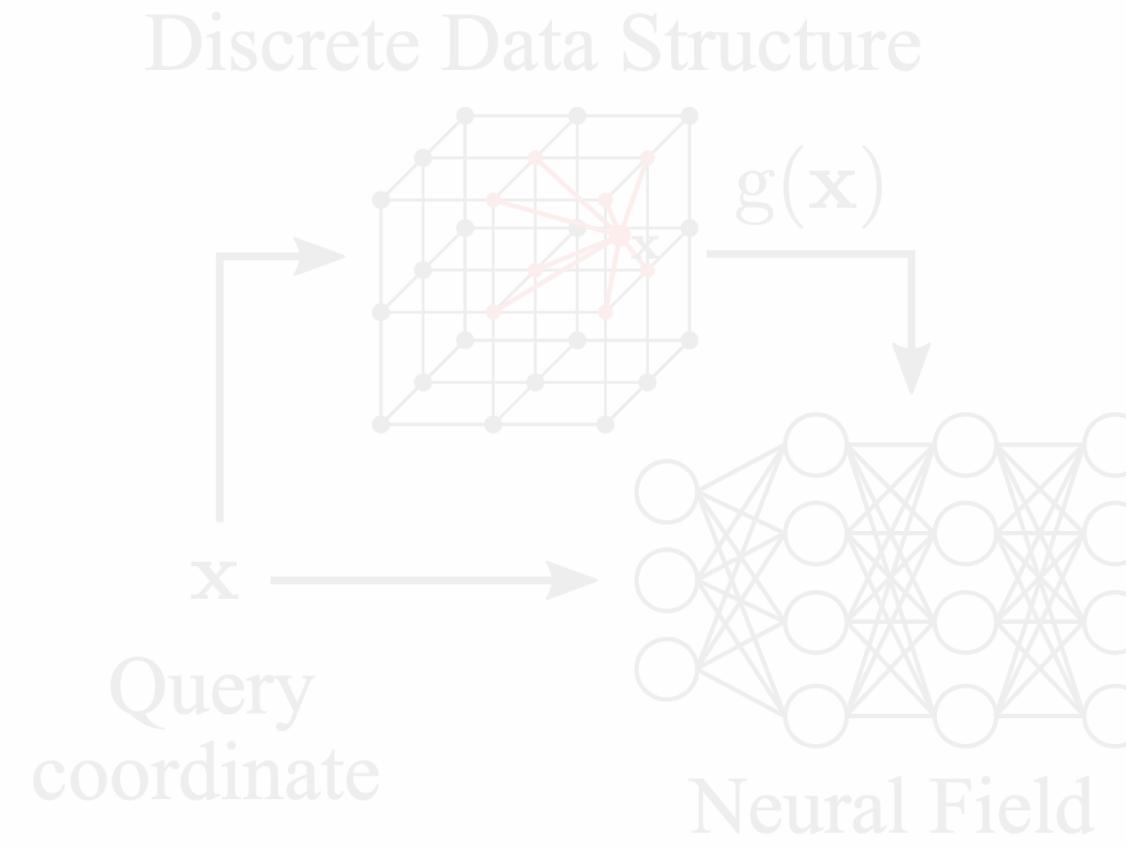


Gradient-Based
Meta-Learning

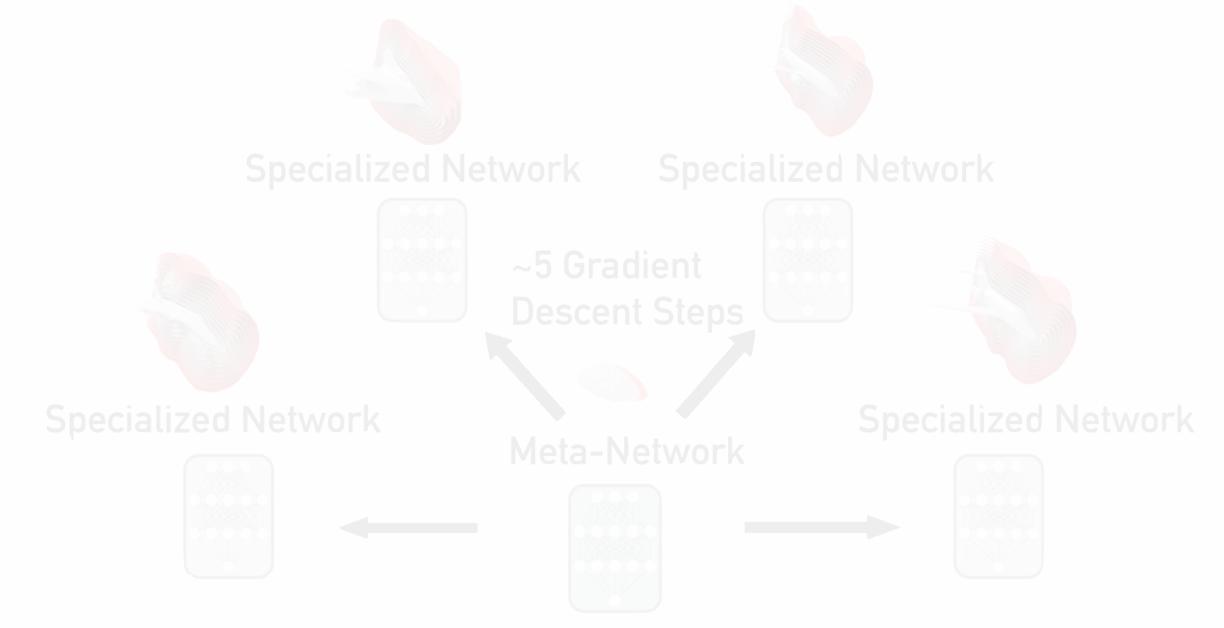
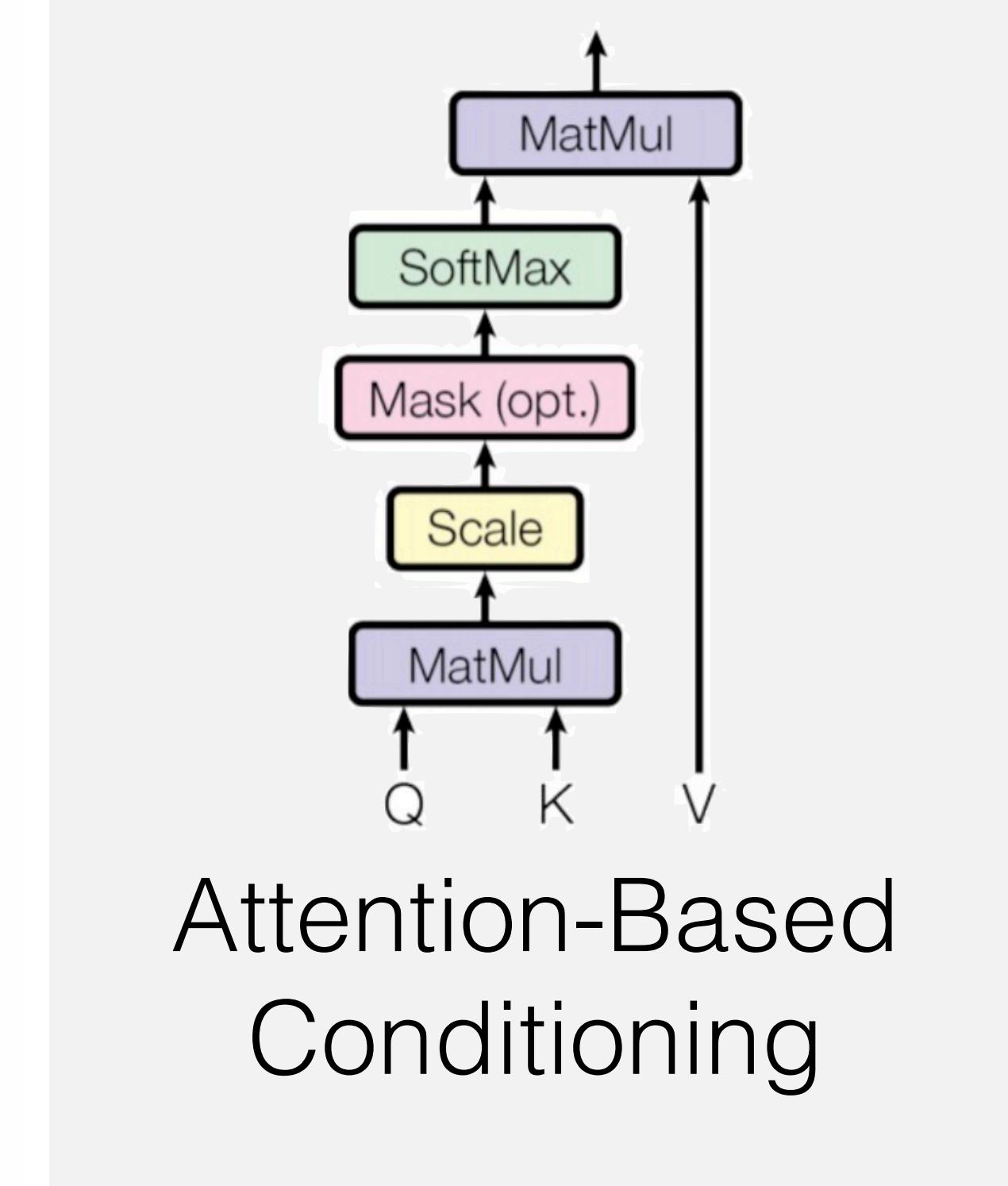
Other types of conditioning...?



Global
Conditioning



Local
Conditioning



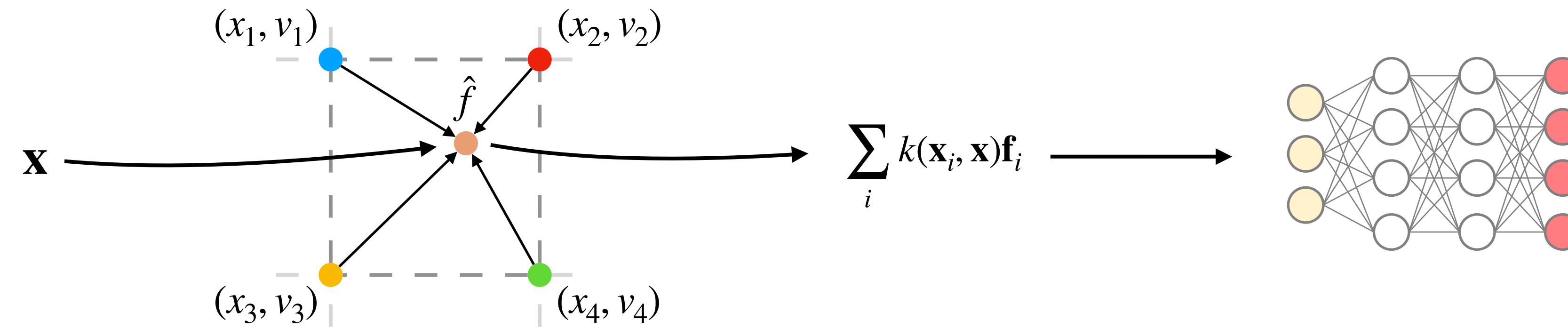
Gradient-Based
Meta-Learning

Recap: Scaled Dot-Product Attention

$$\text{softmax} \left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V} = \mathbf{O}$$

The diagram illustrates the Scaled Dot-Product Attention mechanism. It shows three input matrices: \mathbf{Q} (purple), \mathbf{K}^T (orange), and \mathbf{V} (blue). The \mathbf{Q} matrix is multiplied by the transpose of \mathbf{K} (\mathbf{K}^T) and then scaled by $1/\sqrt{d_k}$. The result is then multiplied by the \mathbf{V} matrix to produce the output \mathbf{O} (pink).

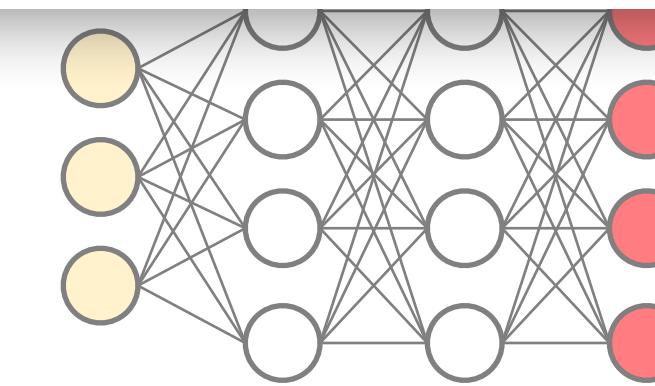
Attention as a form of local conditioning



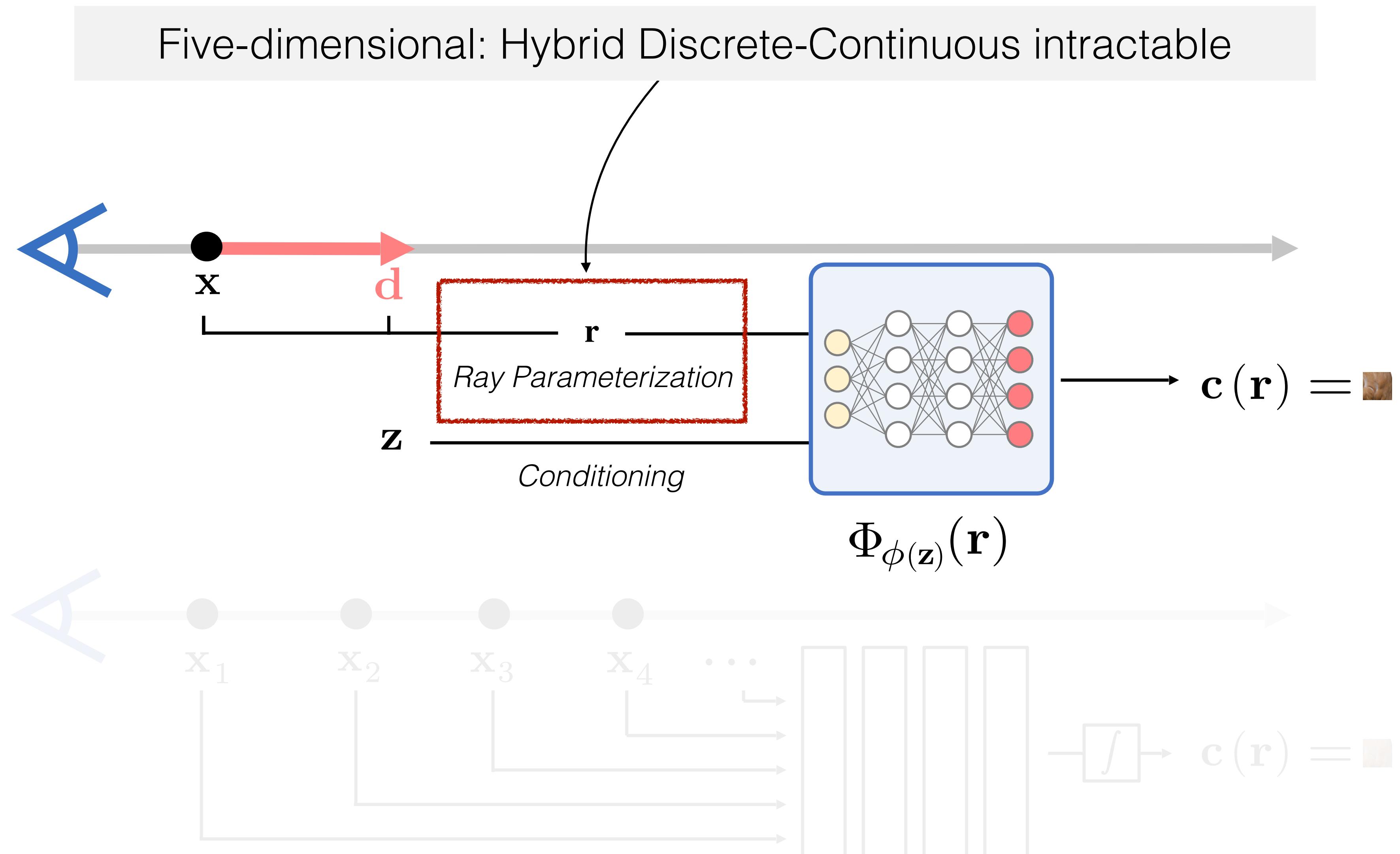
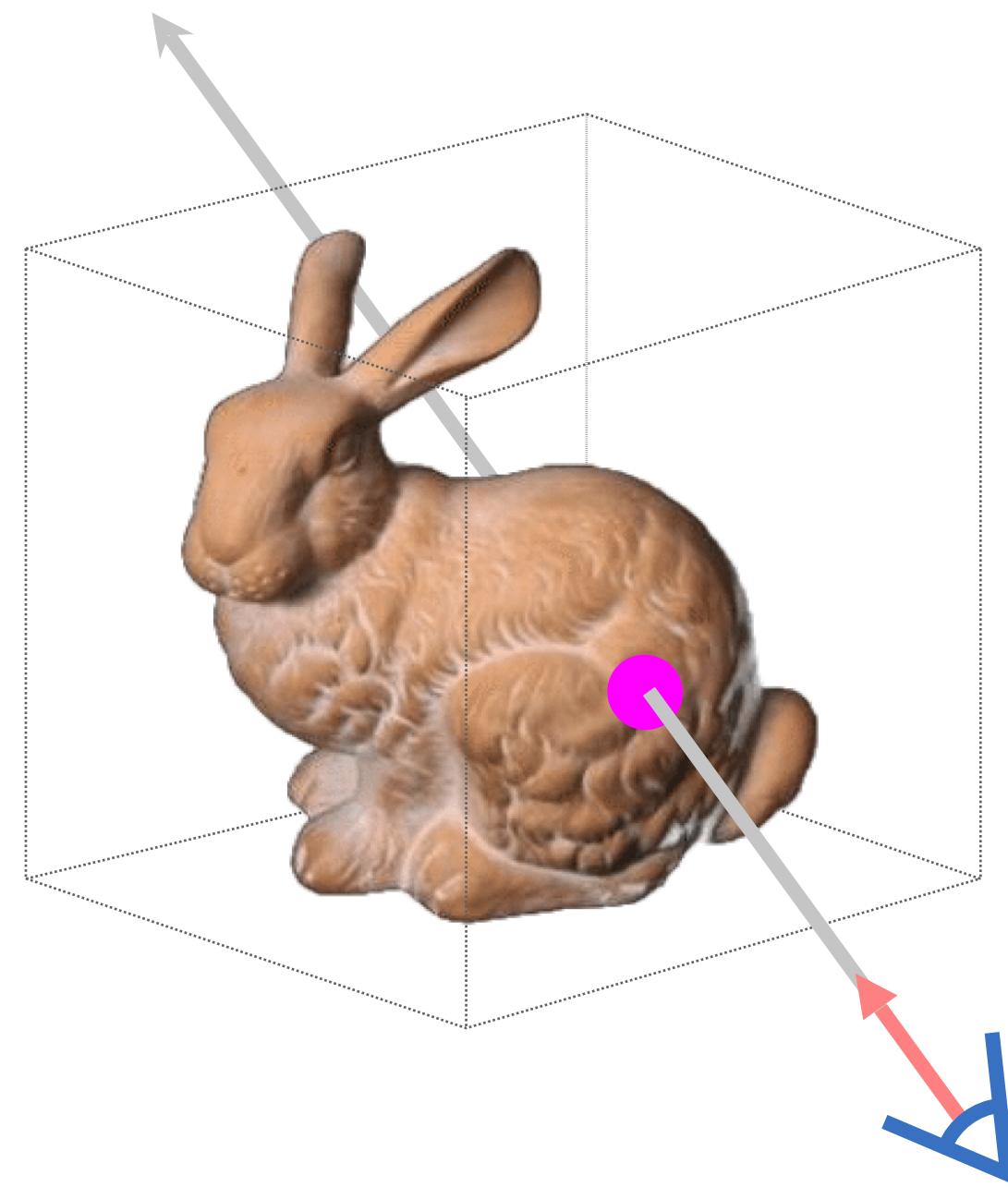
Attention as a form of local conditioning

It's just a different kernel!

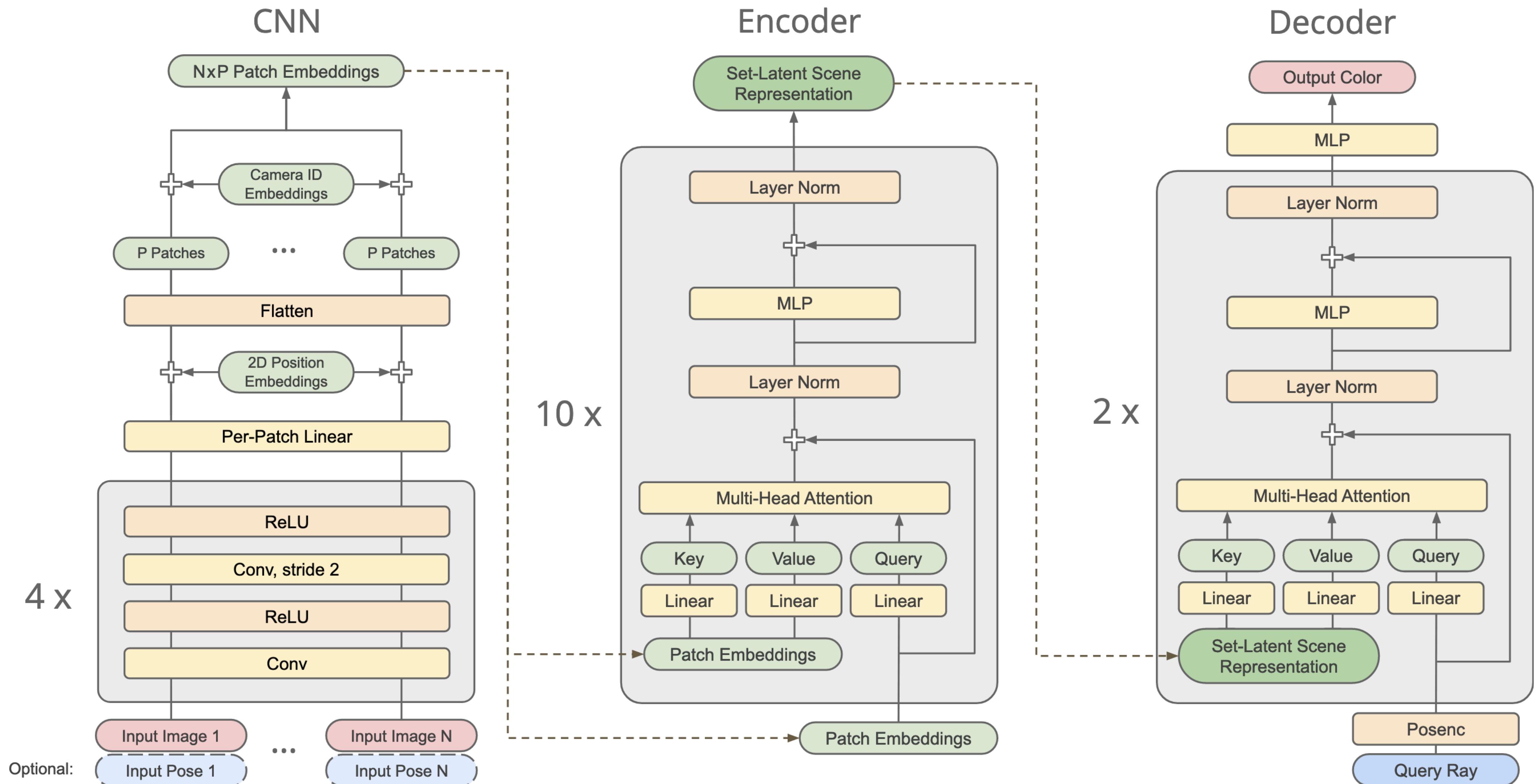
$$\mathbf{X} \quad \text{softmax}\left(\frac{\mathbf{K}^T \mathbf{V}}{\sqrt{d_k}}\right) \rightarrow \sum_i k(\mathbf{k}_i, \mathbf{x}) \mathbf{v}_i$$



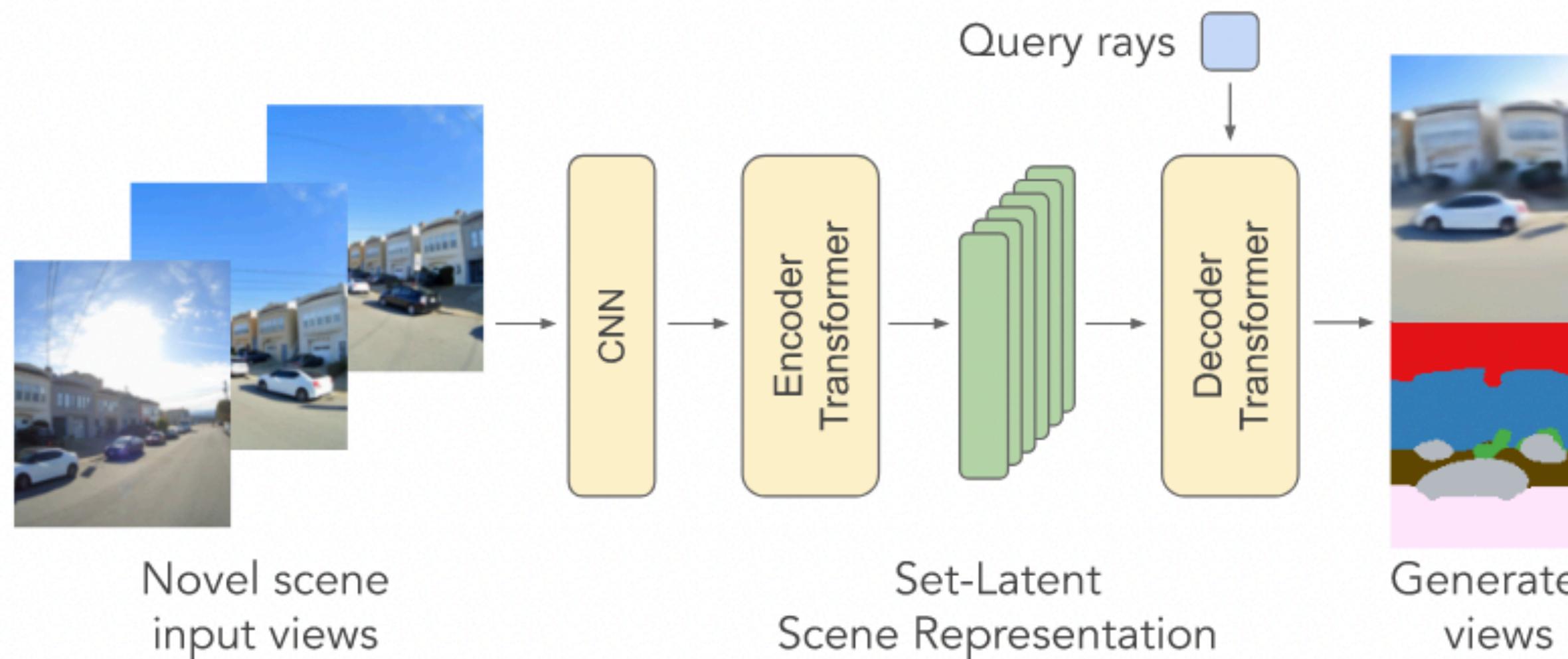
Now: Learn Prior over Light Fields!



Transformers for Light Fields: Scene Representation Transformer



Attention-conditioned Light Fields



**Scene Representation Transformer:
Geometry-Free Novel View Synthesis Through Set-Latent Scene Representations**

Mehdi S. M. Sajjadi

Noha Radwan

Jakob Uszkoreit*

Henning Meyer

Suhani Vora

Thomas Funkhouser

Etienne Pot

Mario Lučić

Andrea Tagliasacchi^{†‡}

Urs Bergmann

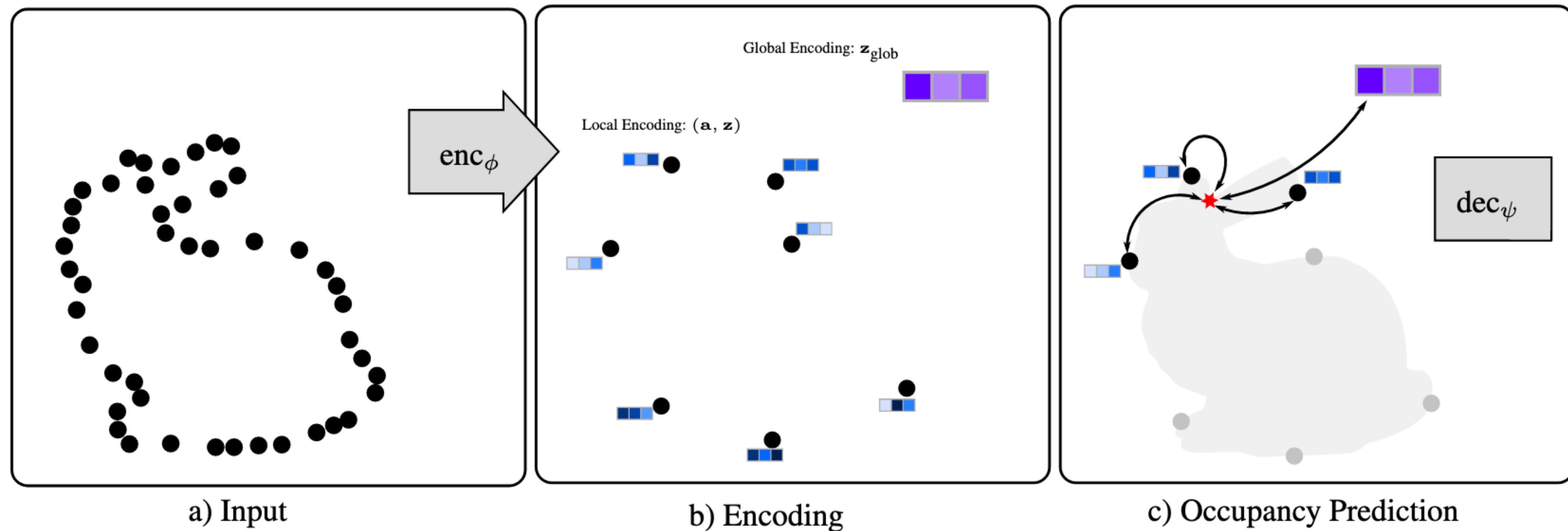
Daniel Duckworth

Andrea Tagliasacchi^{†‡}

Klaus Greff

Alexey Dosovitskiy*

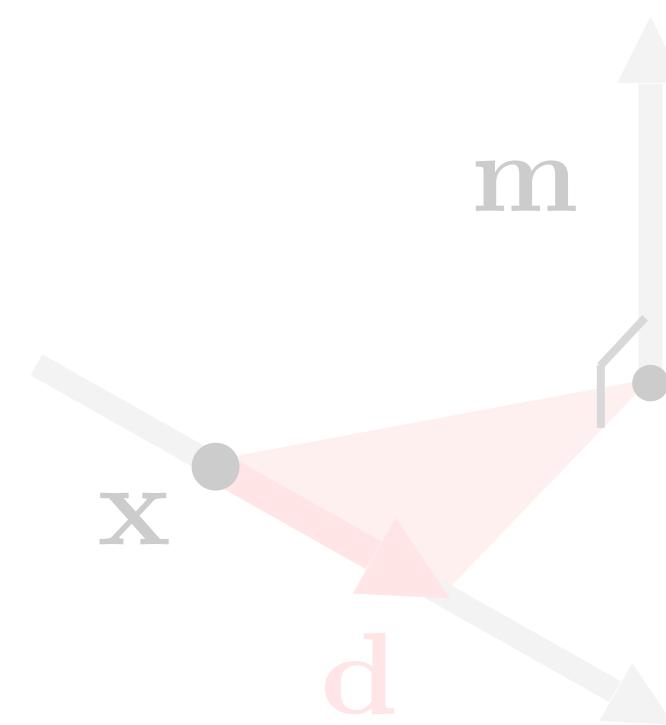
Transformer Conditioning for Point Clouds: AIR-Nets, Giebenhain et al. 2022



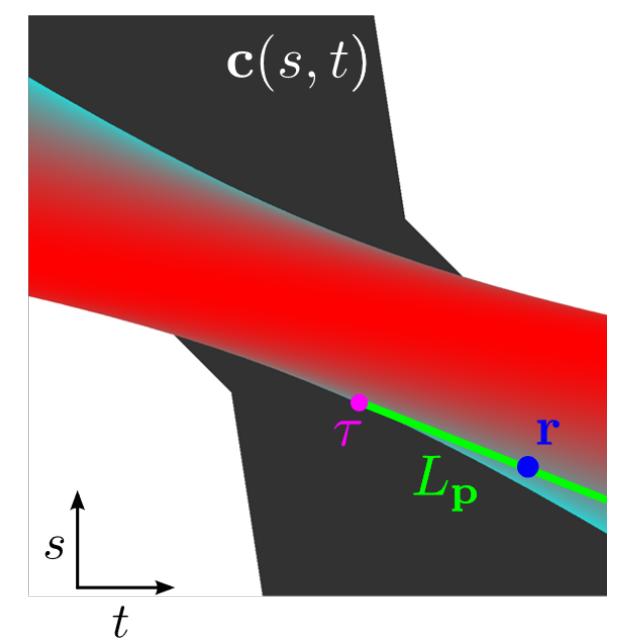
AIR-Nets, Giebenhain et al. 2022

Where's the geometry in a Light Field?

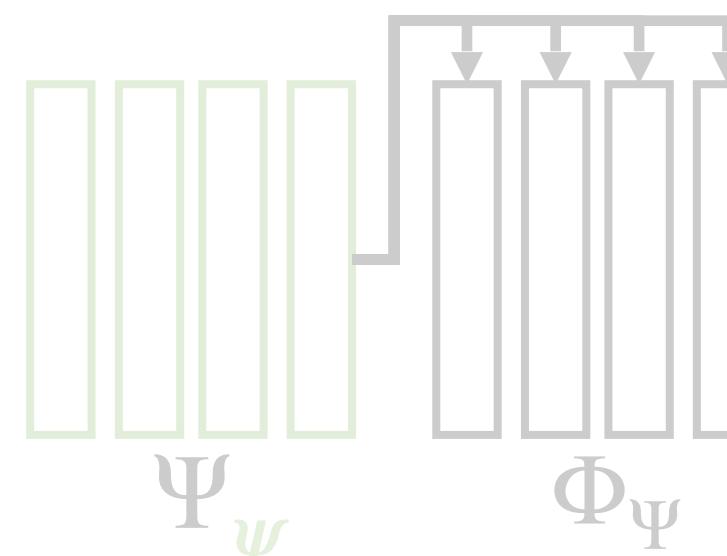
Parameterization



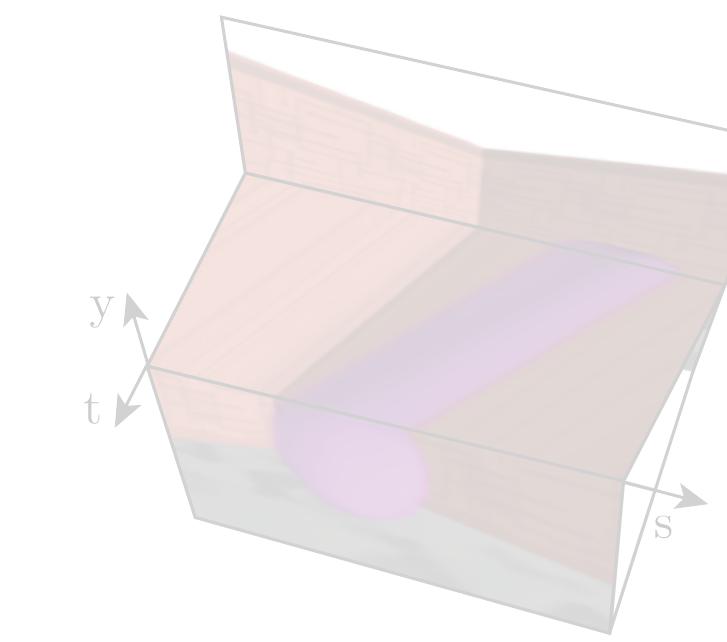
LFN Geometry



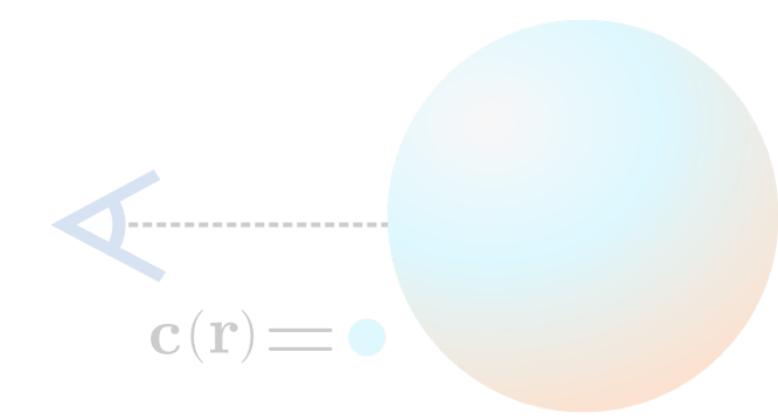
Meta-Learning



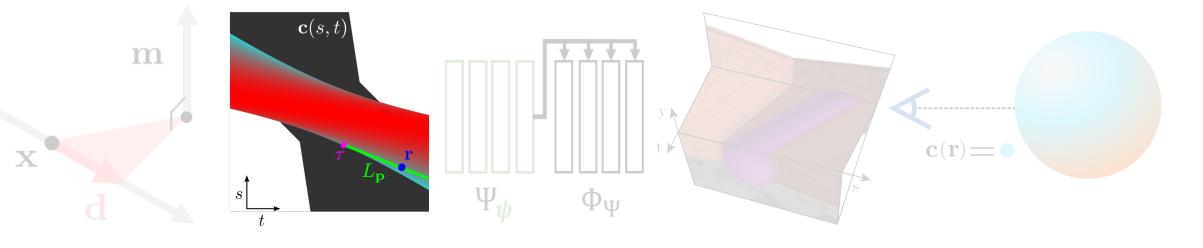
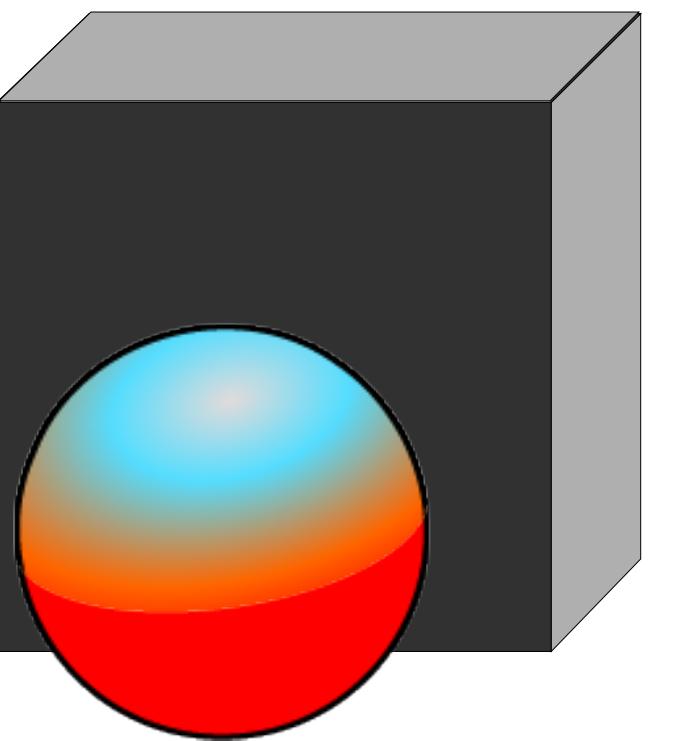
Results



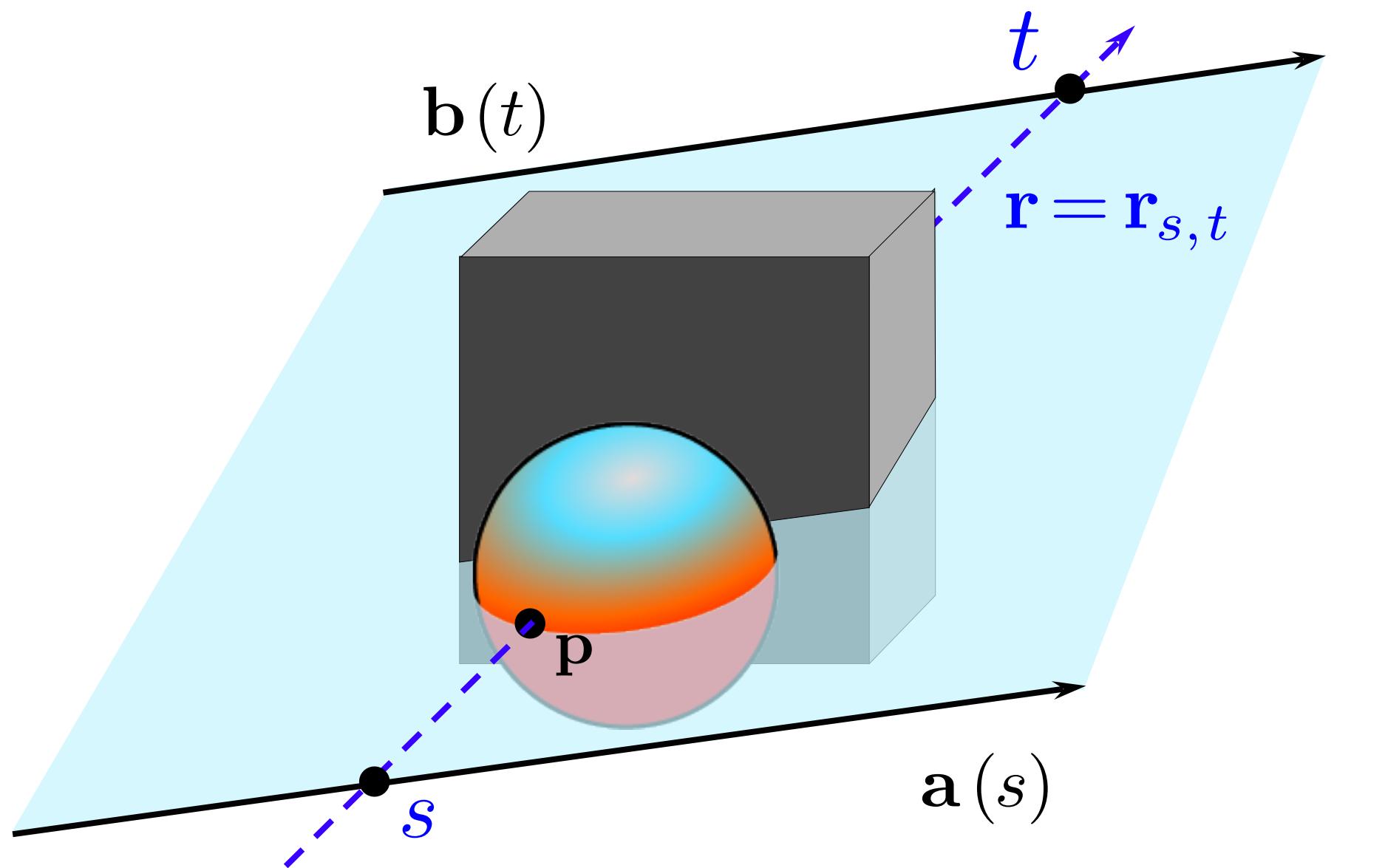
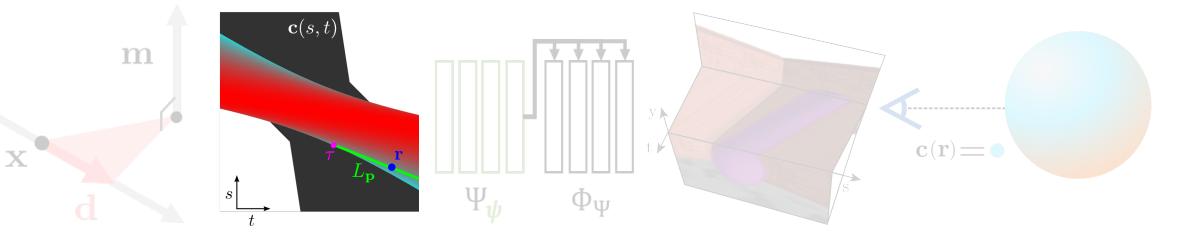
Limitations



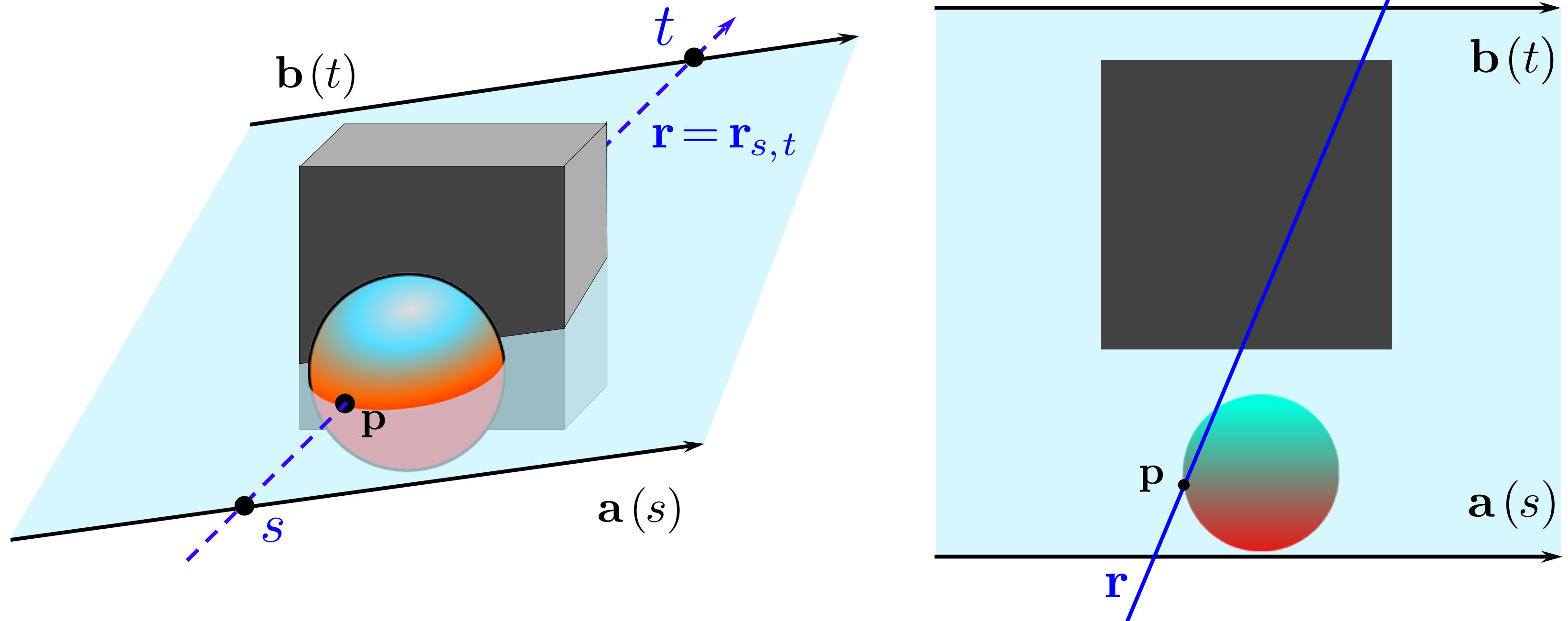
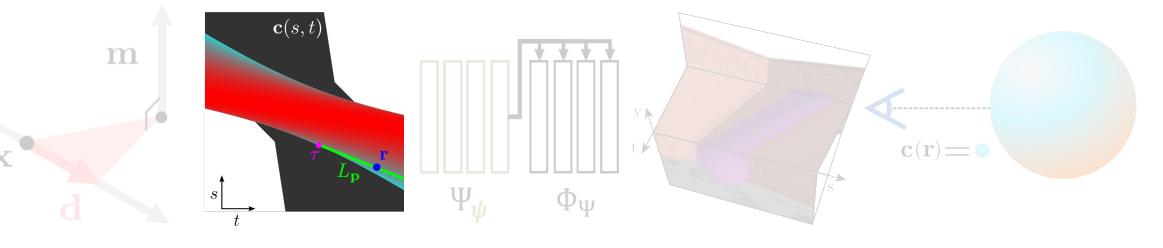
The geometry of LFNs



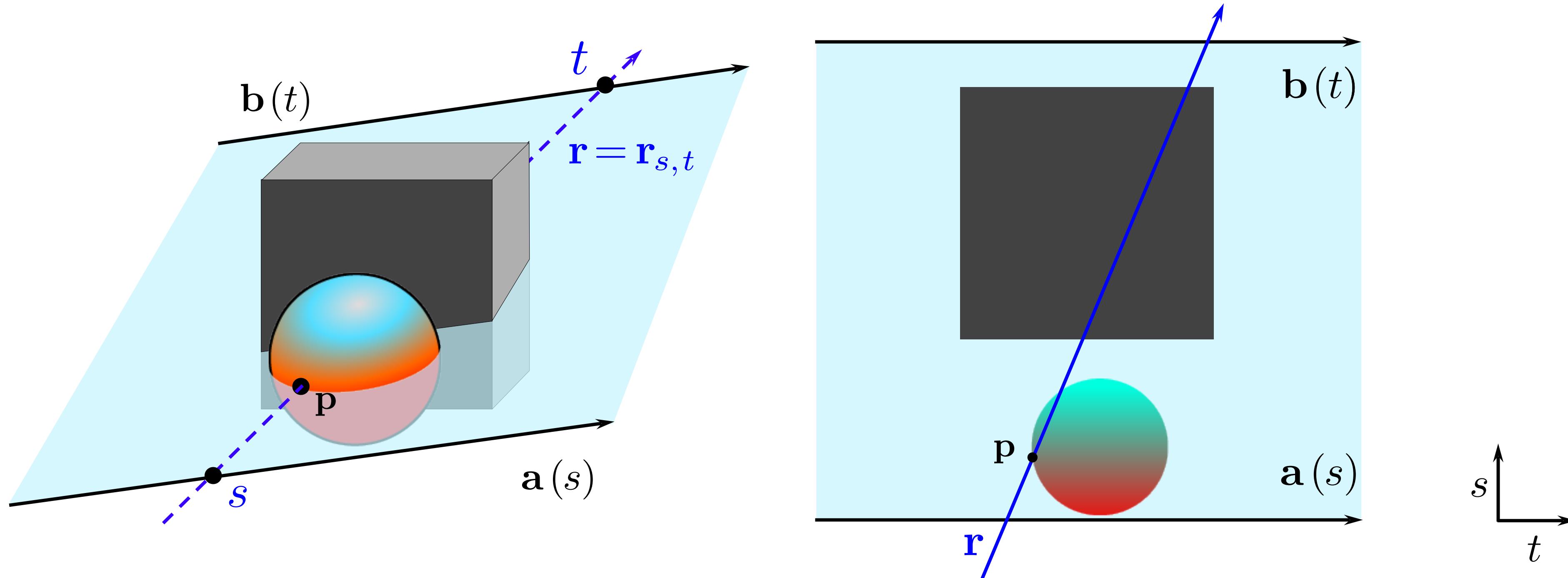
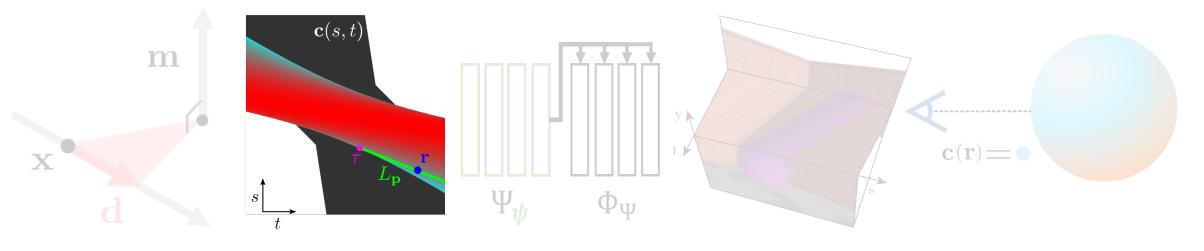
The geometry of LFNs



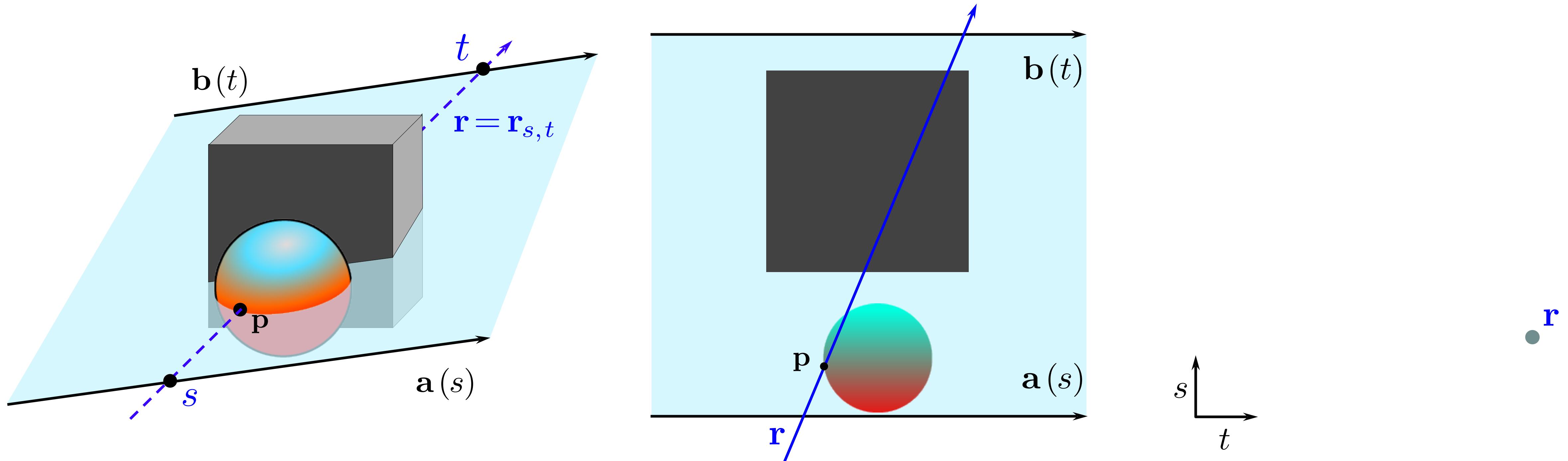
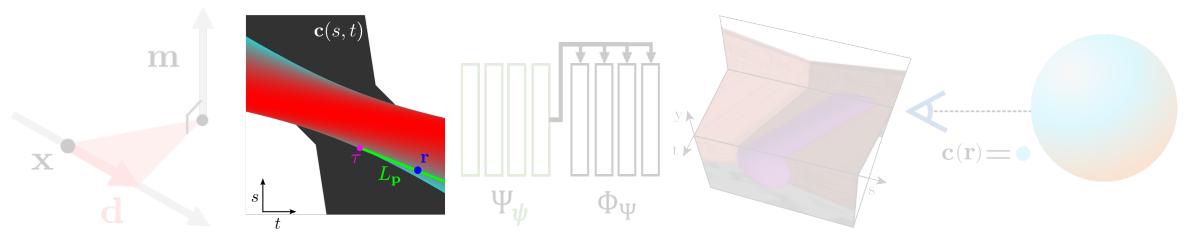
The geometry of LFNs



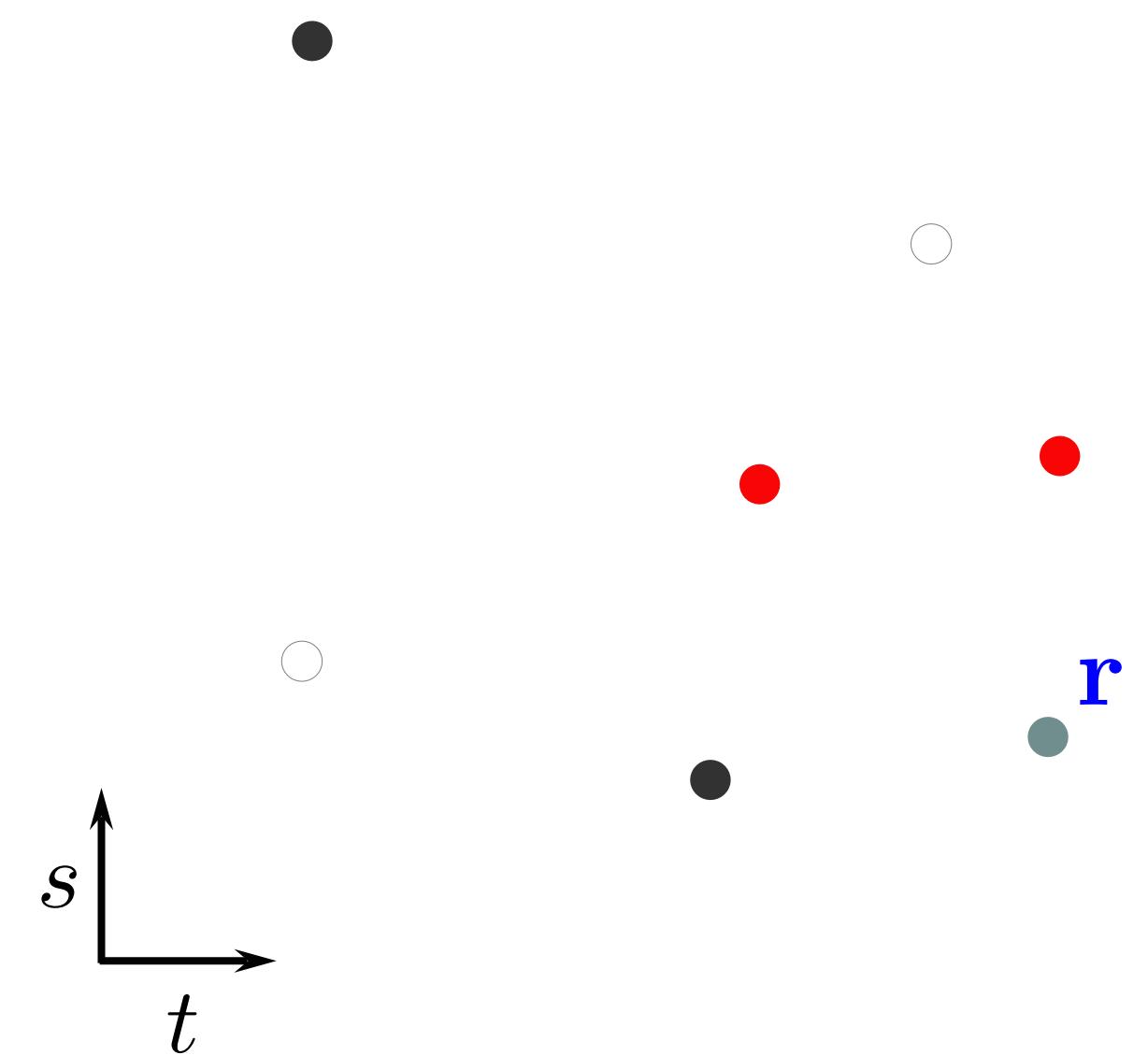
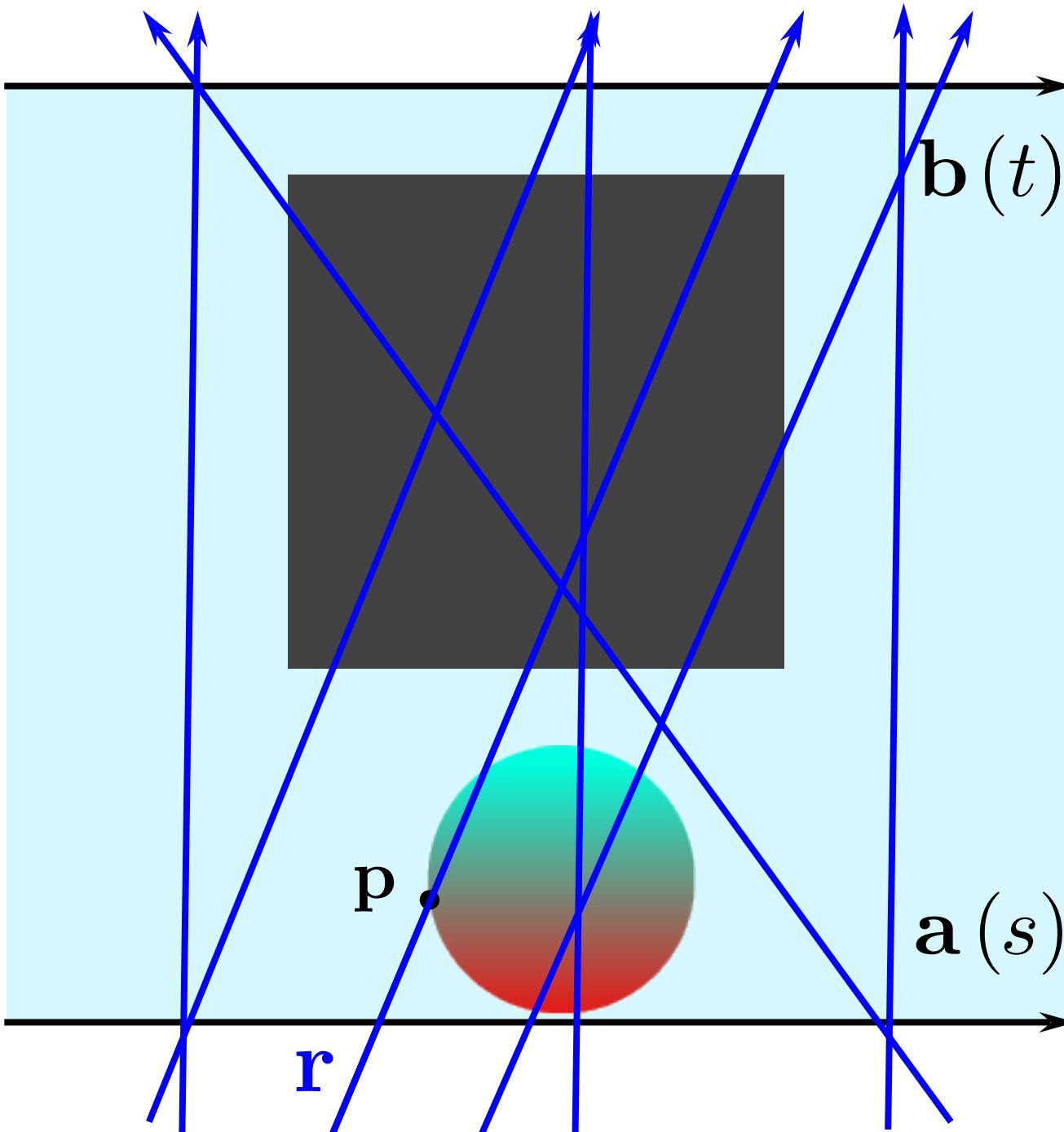
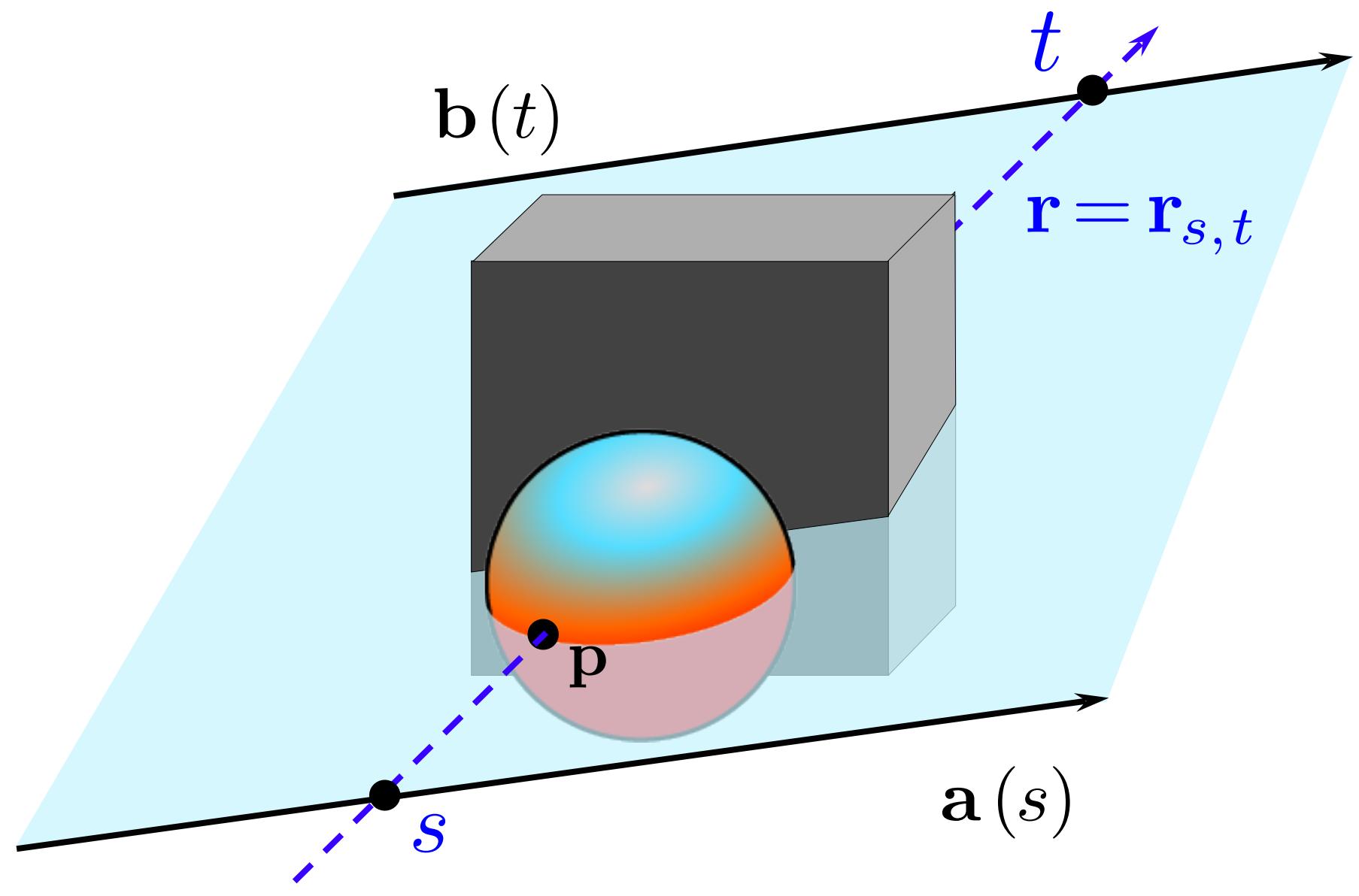
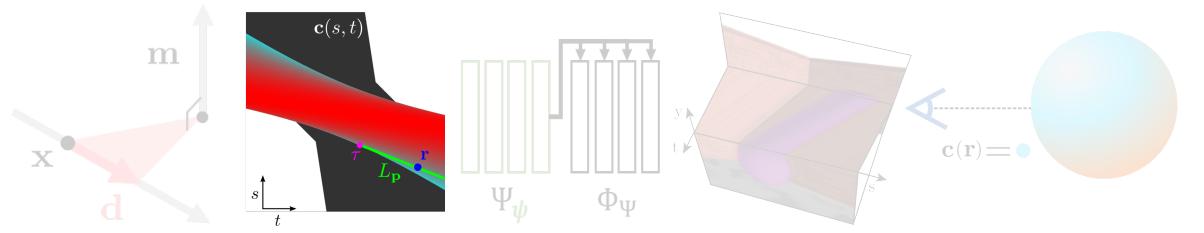
The geometry of LFNs



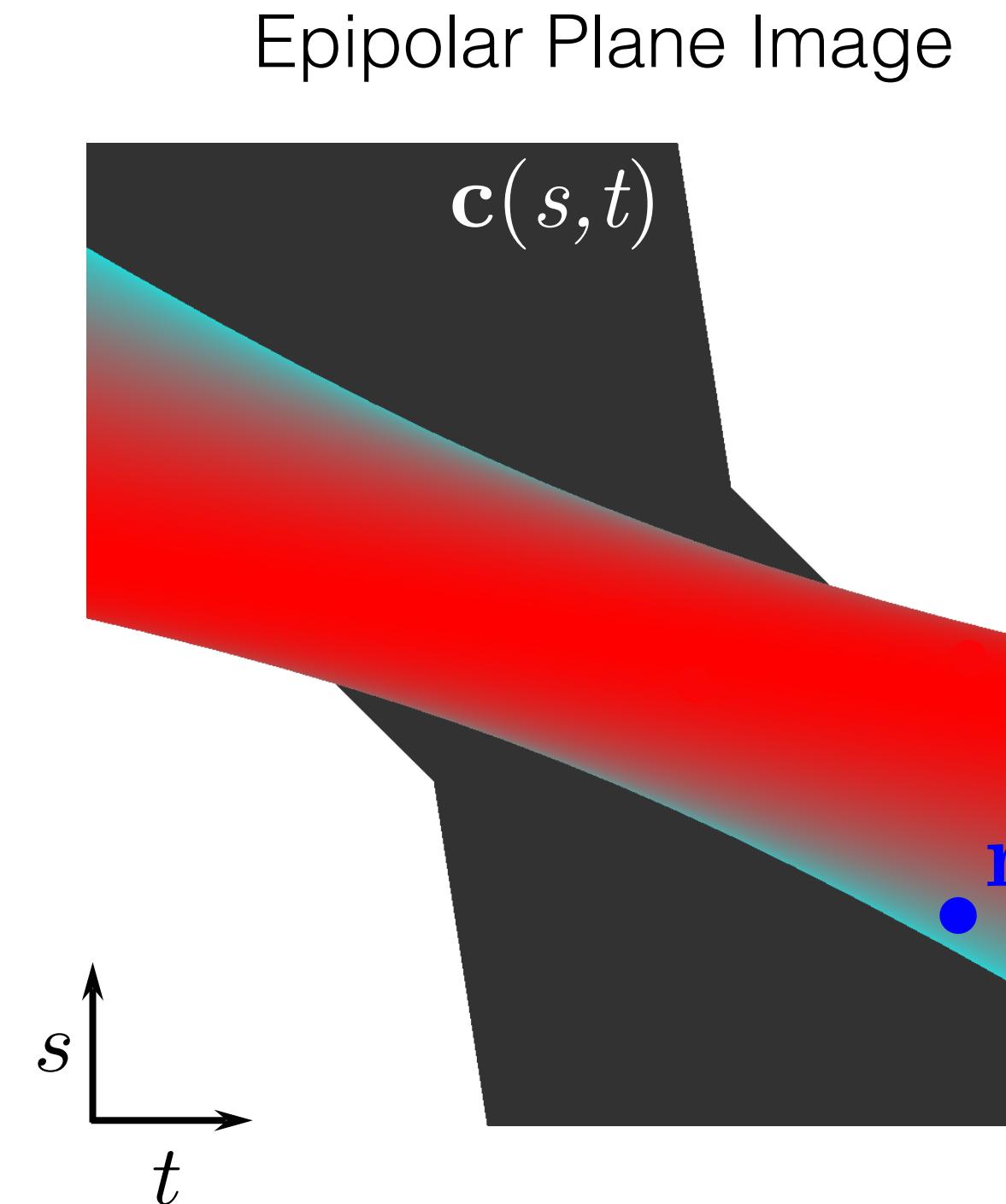
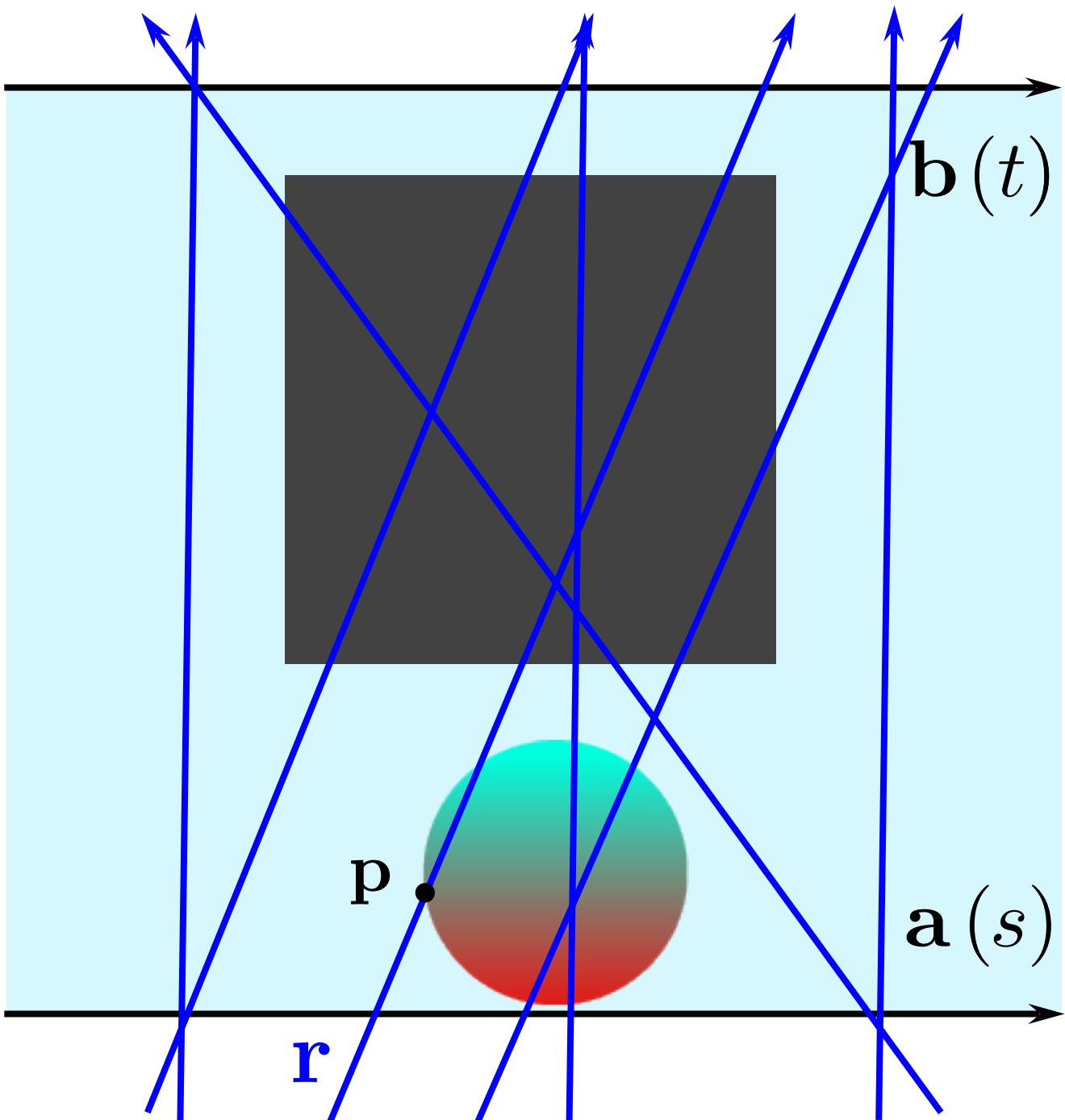
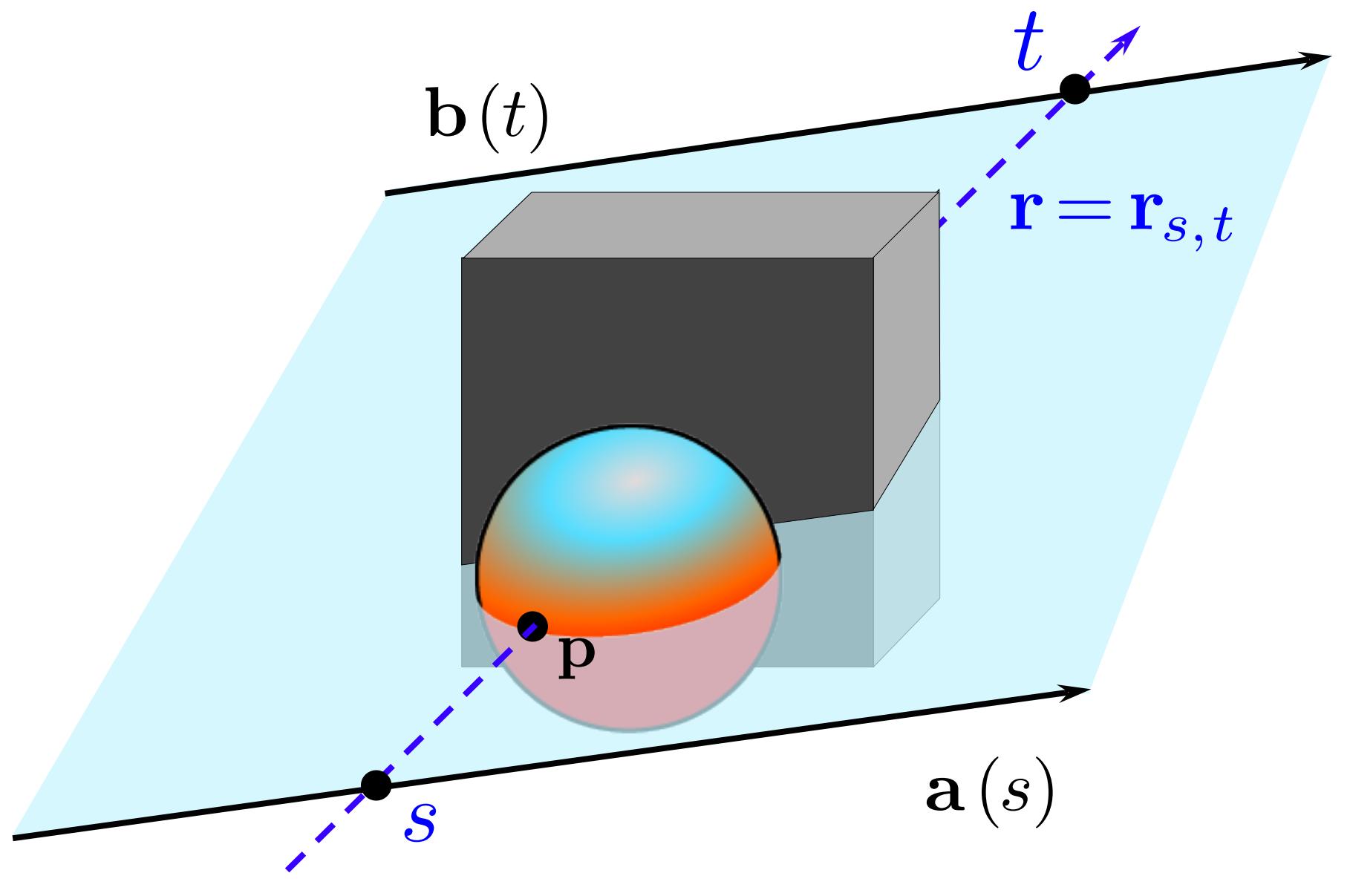
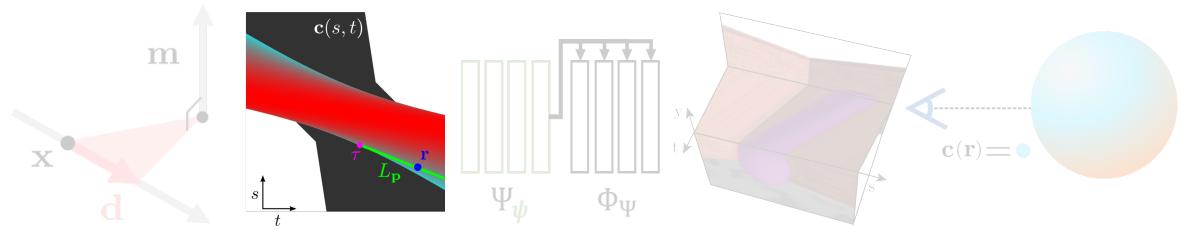
The geometry of LFNs



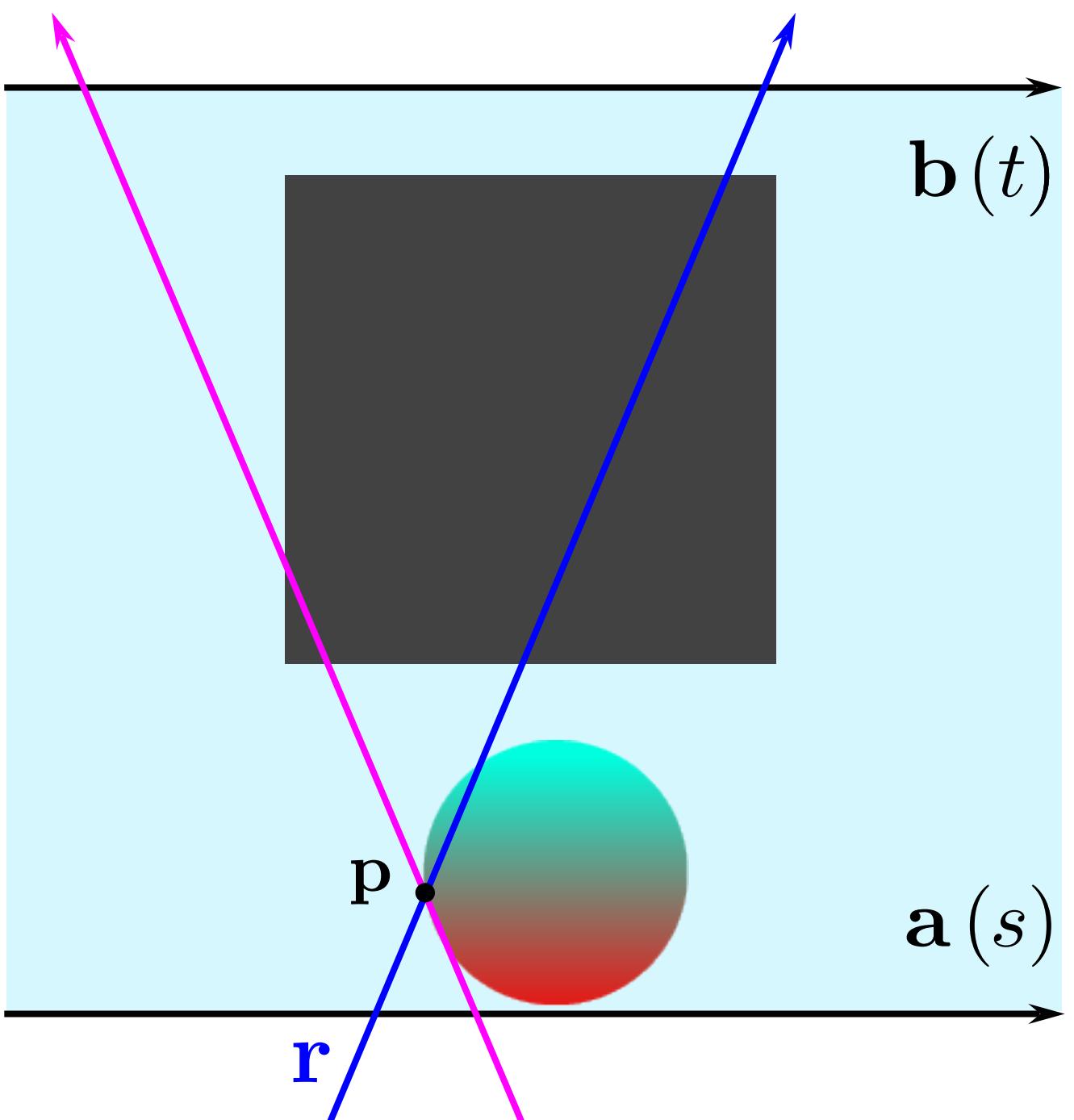
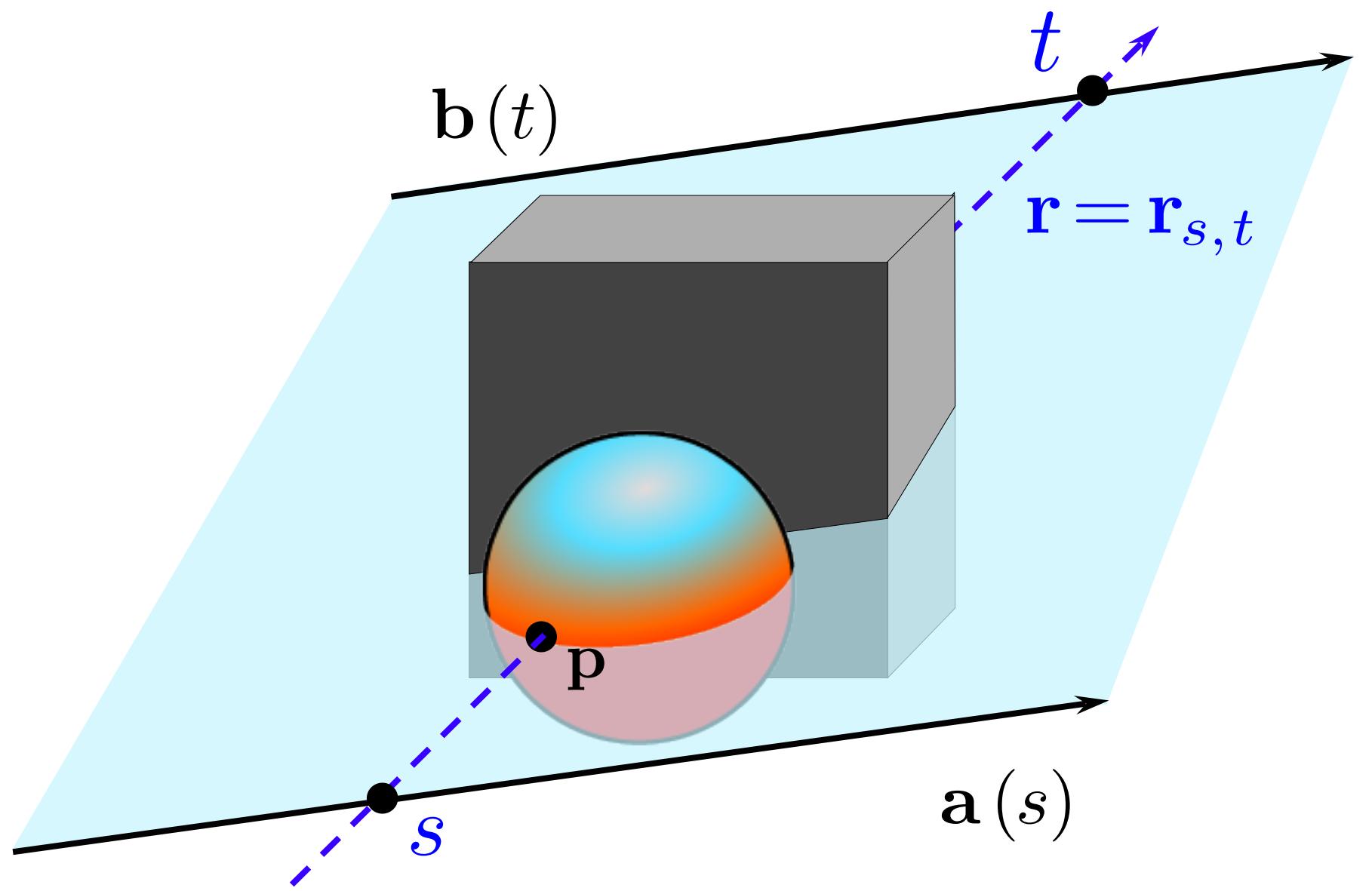
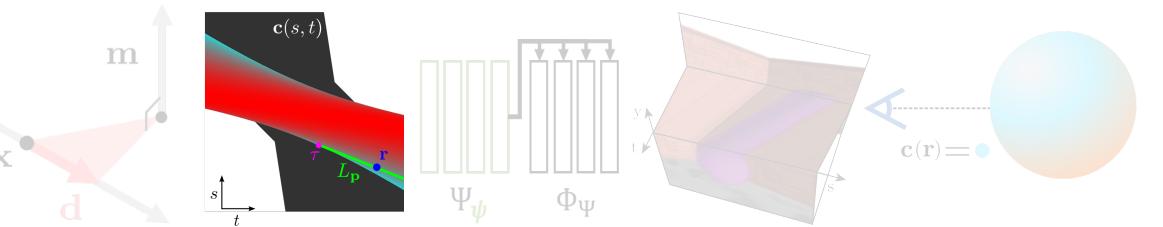
The geometry of LFNs



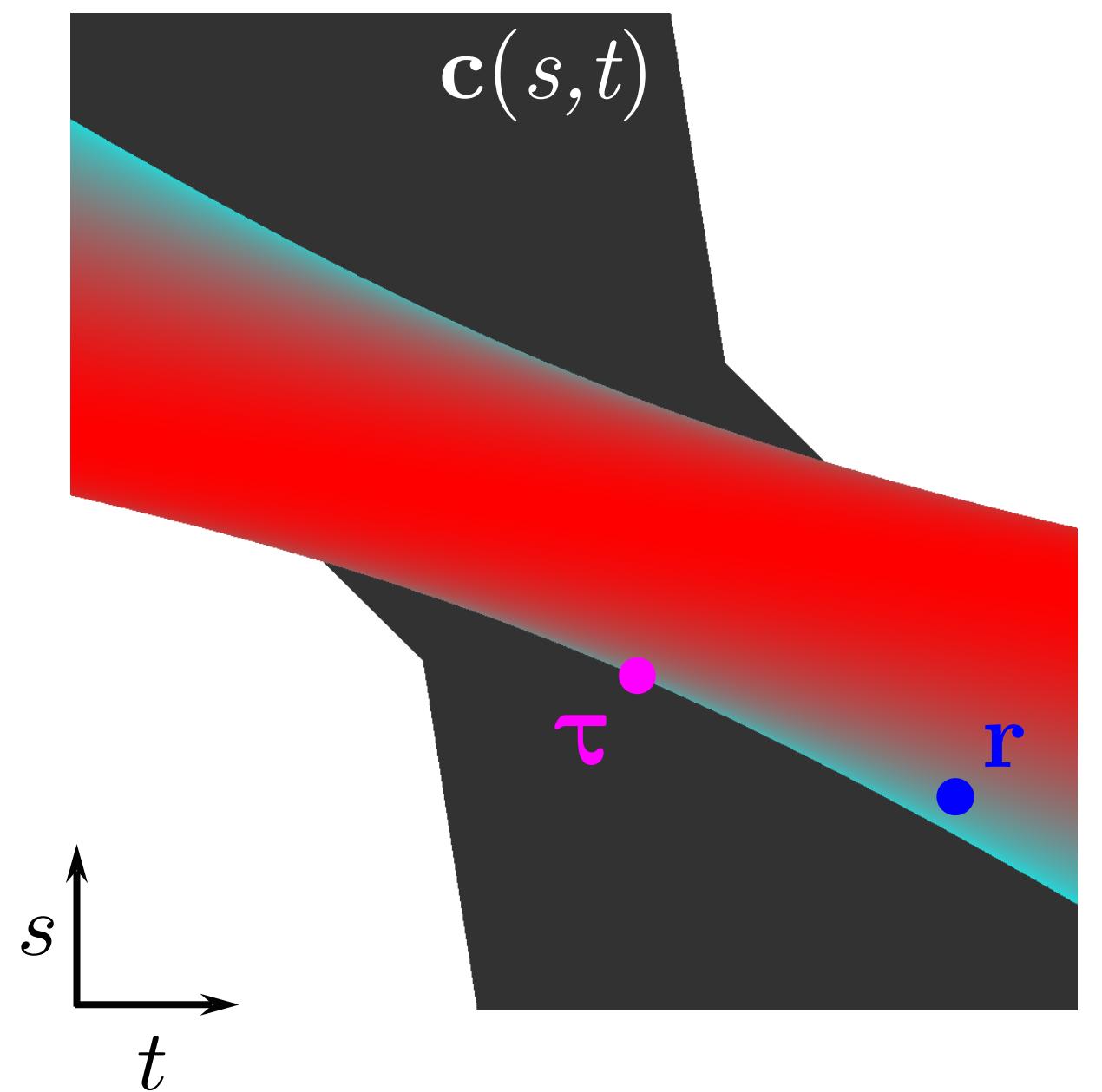
The geometry of LFNs



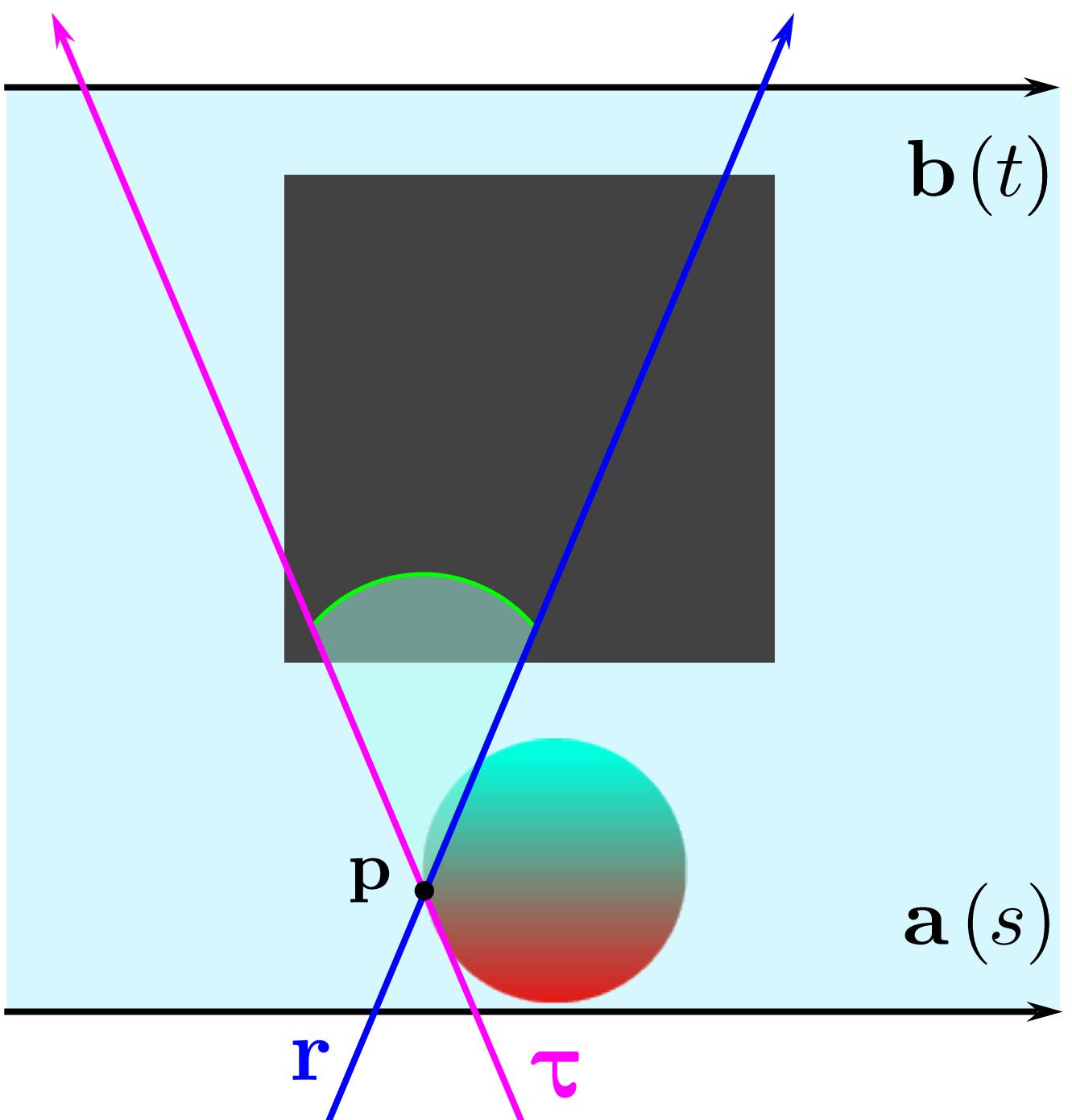
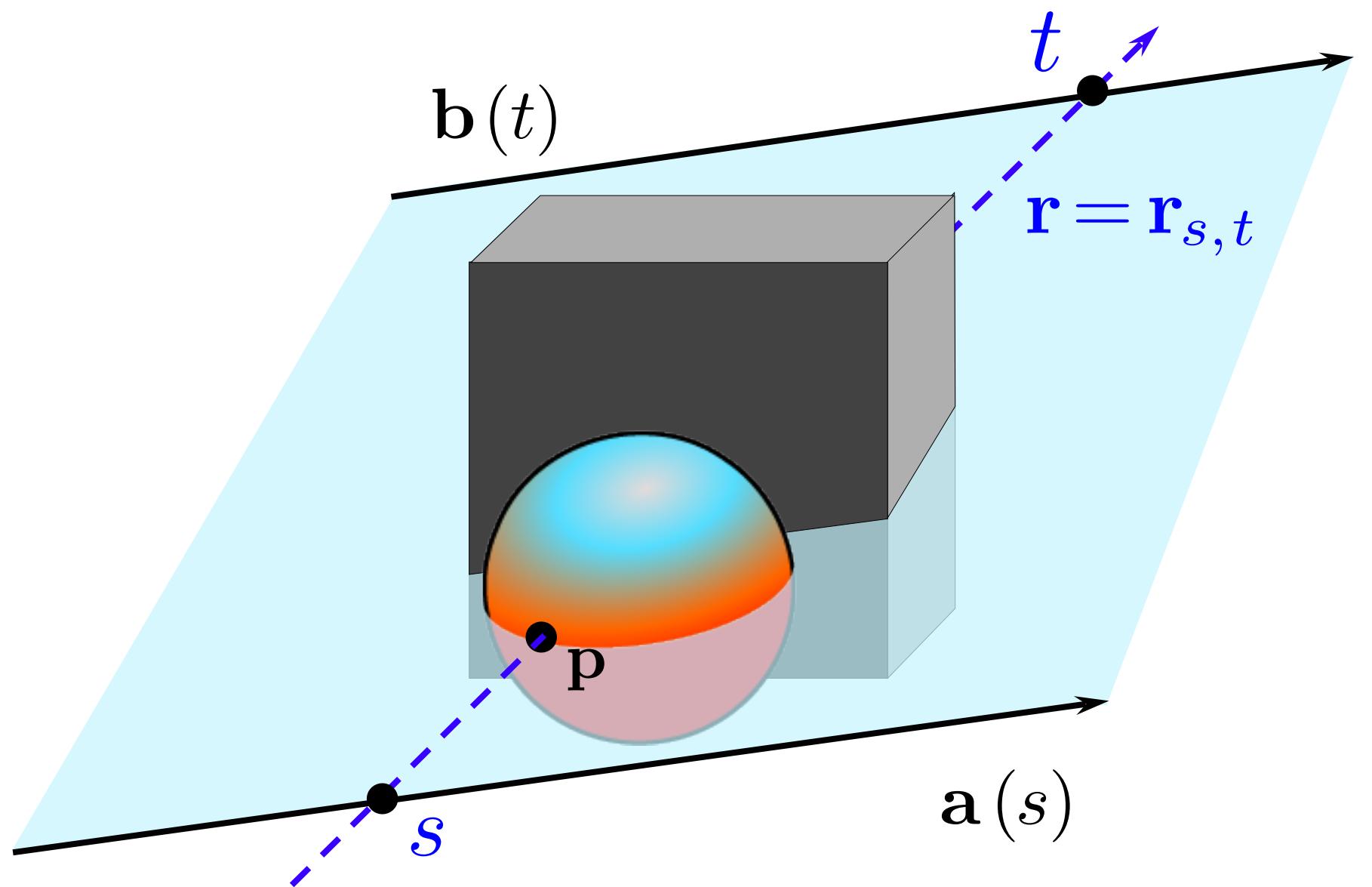
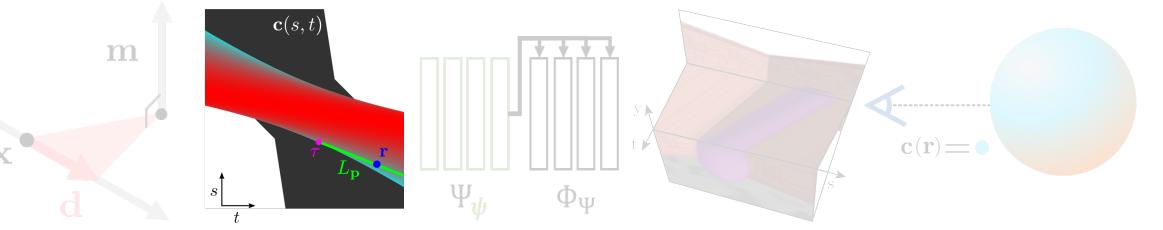
The geometry of LFNs



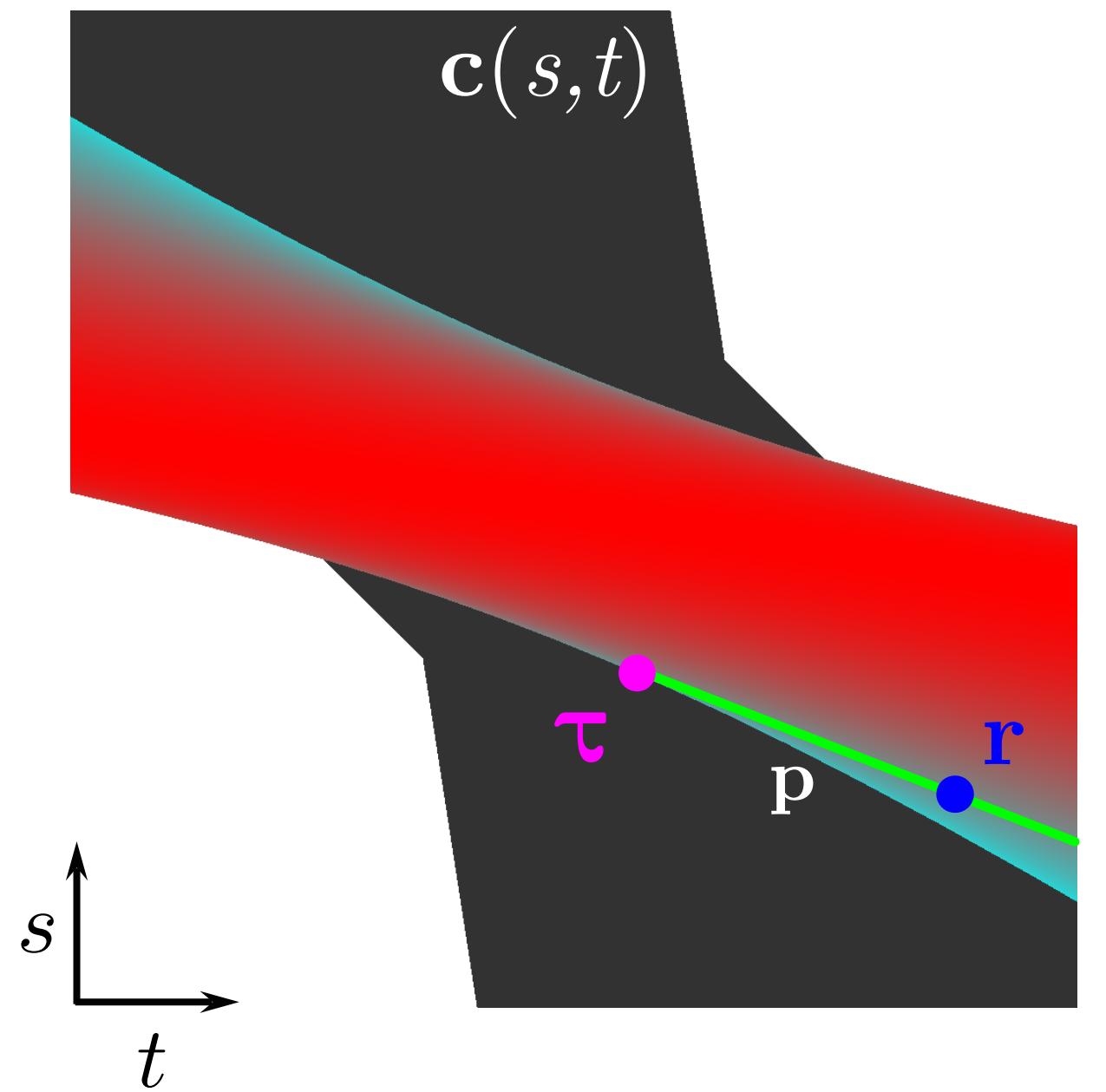
Epipolar Plane Image



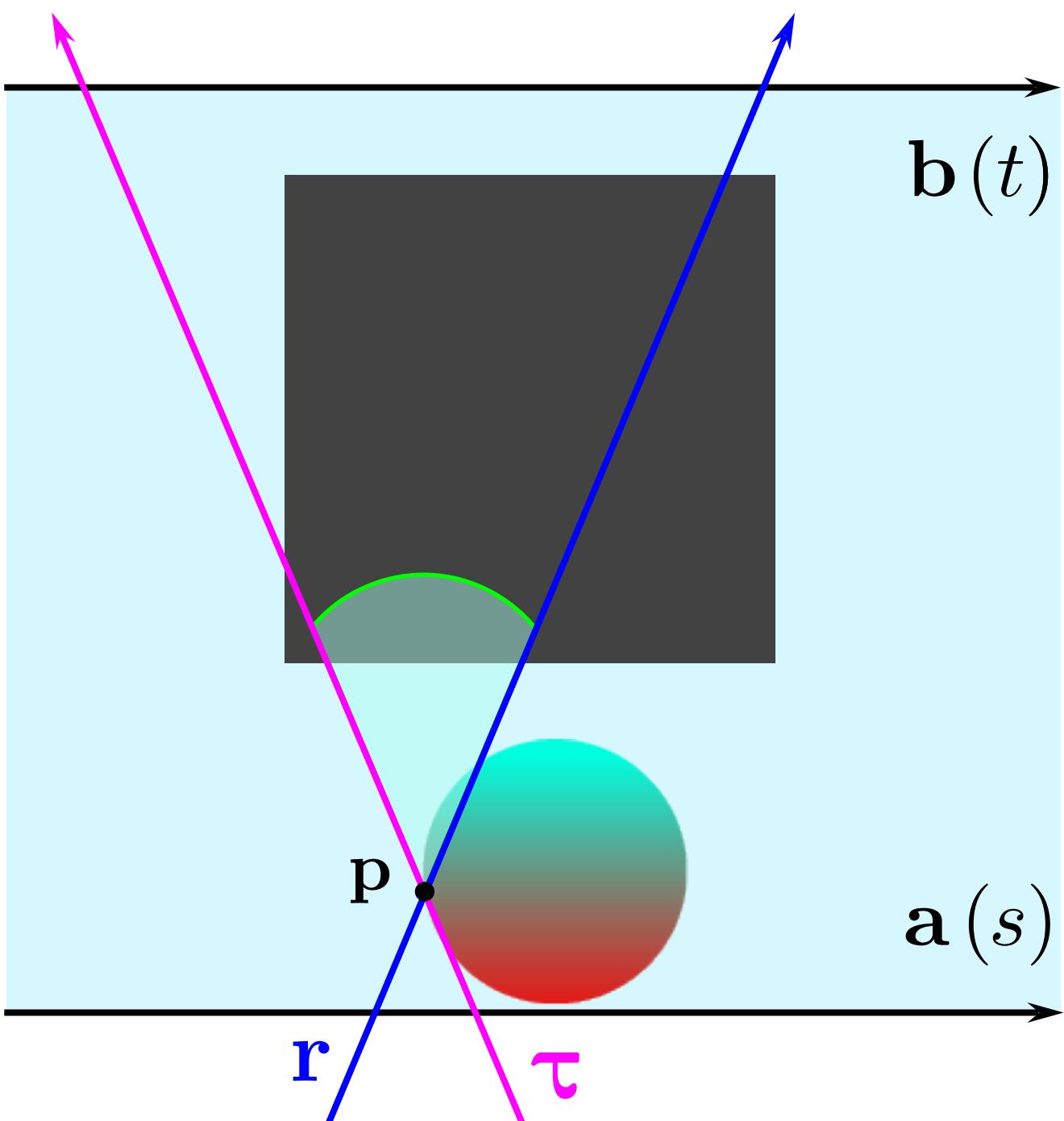
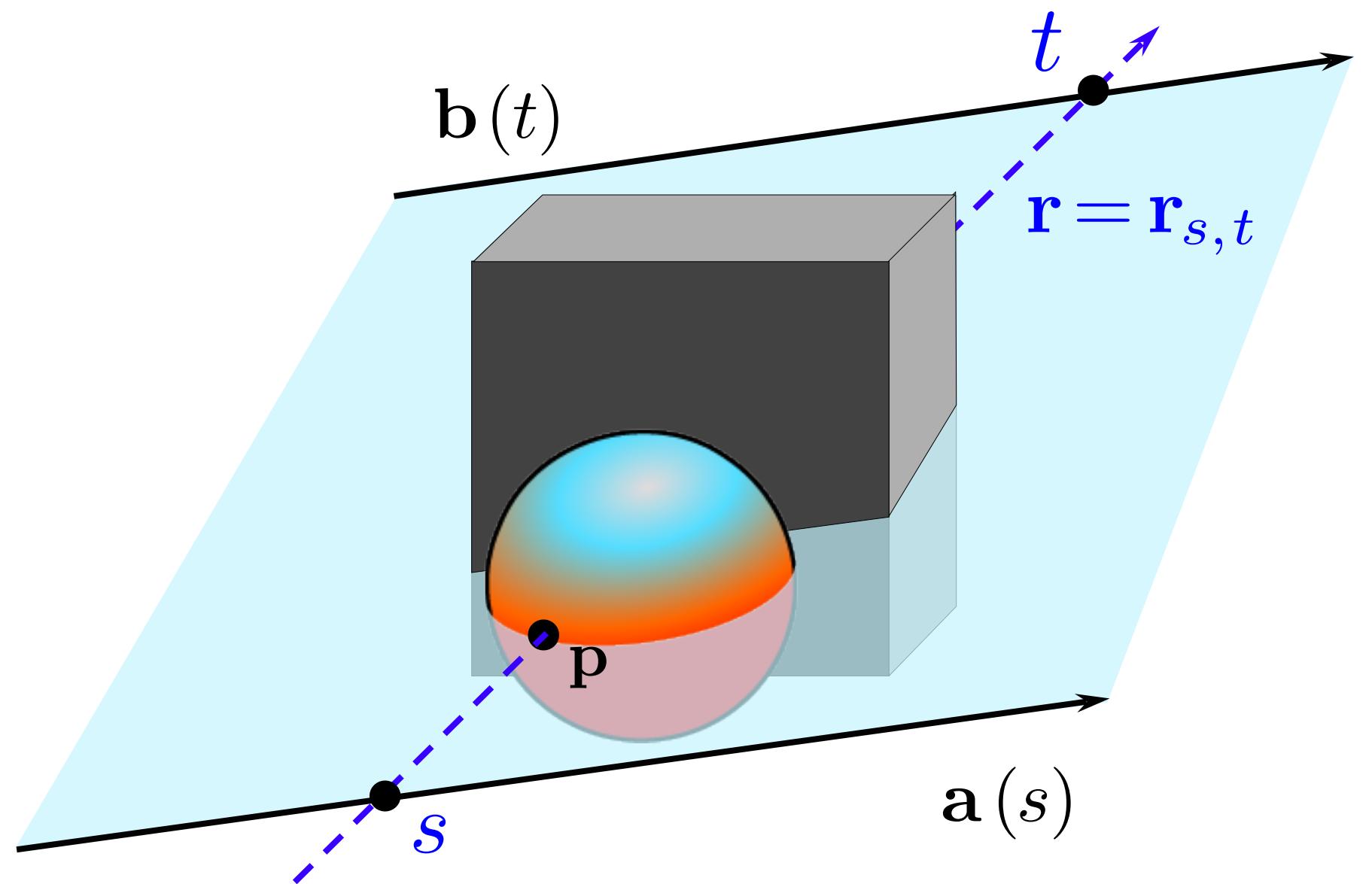
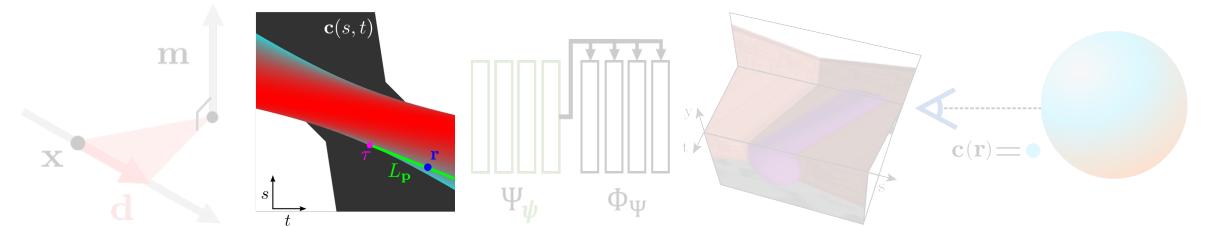
The geometry of LFNs



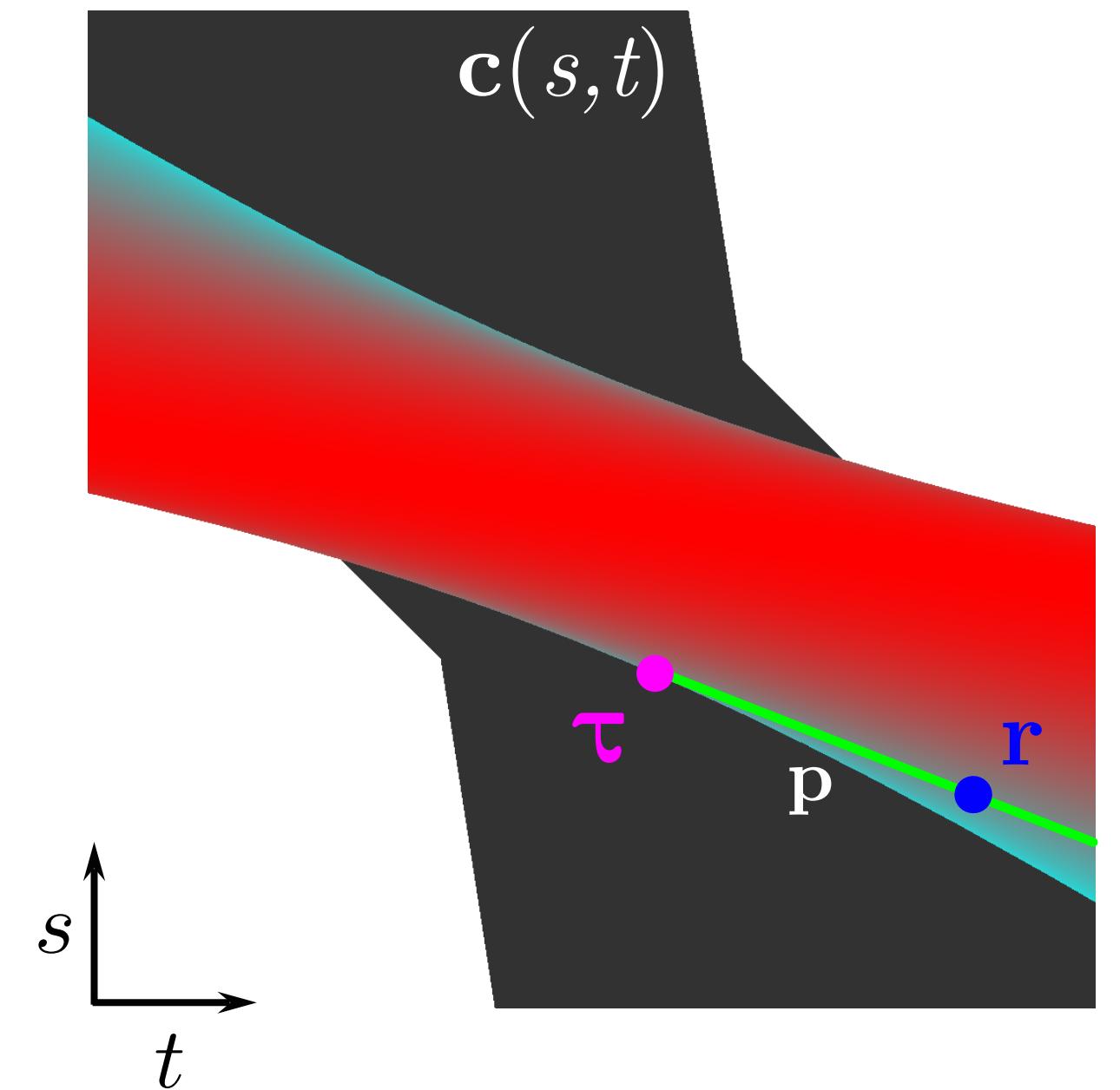
Epipolar Plane Image



The geometry of LFNs

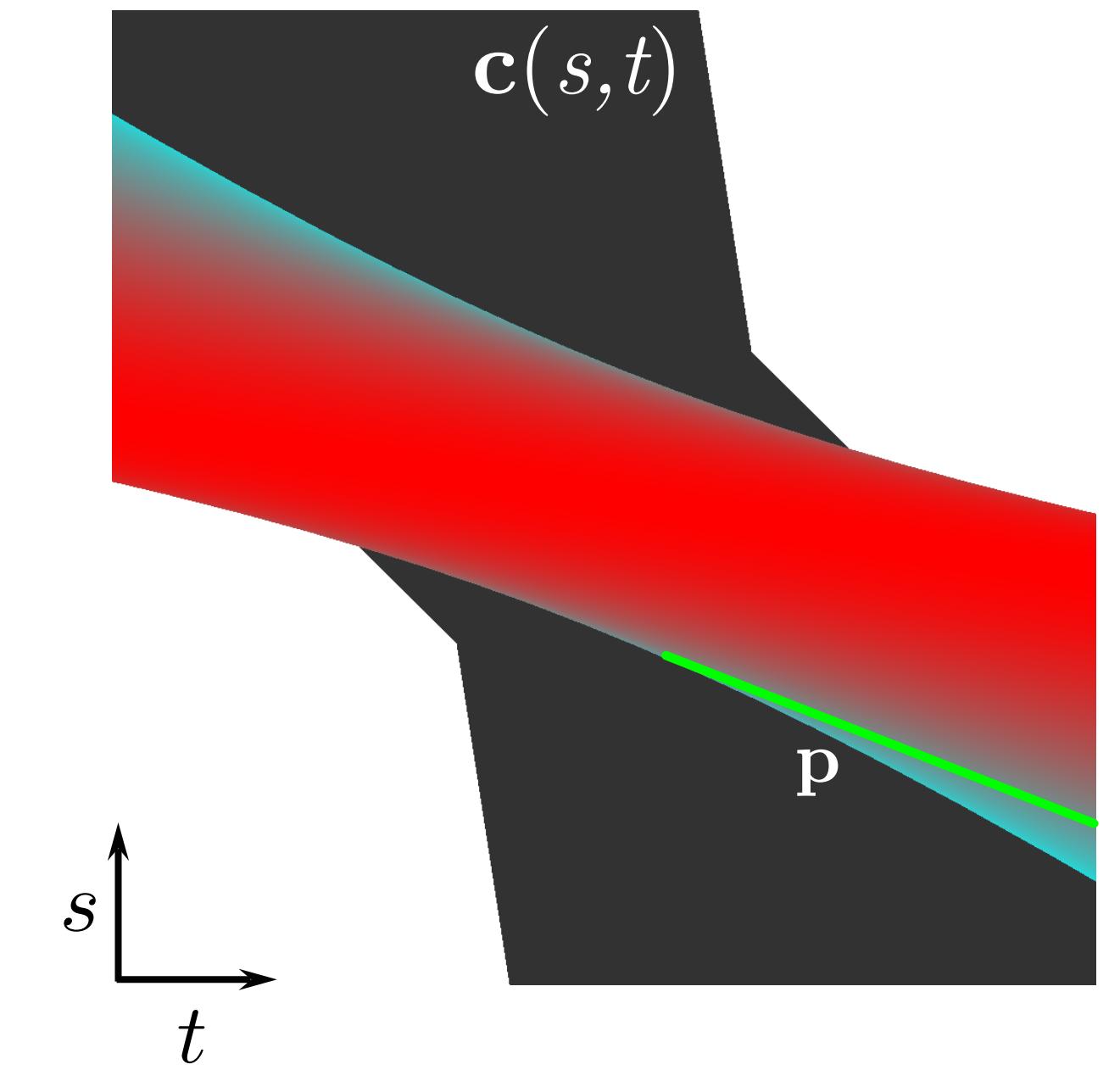
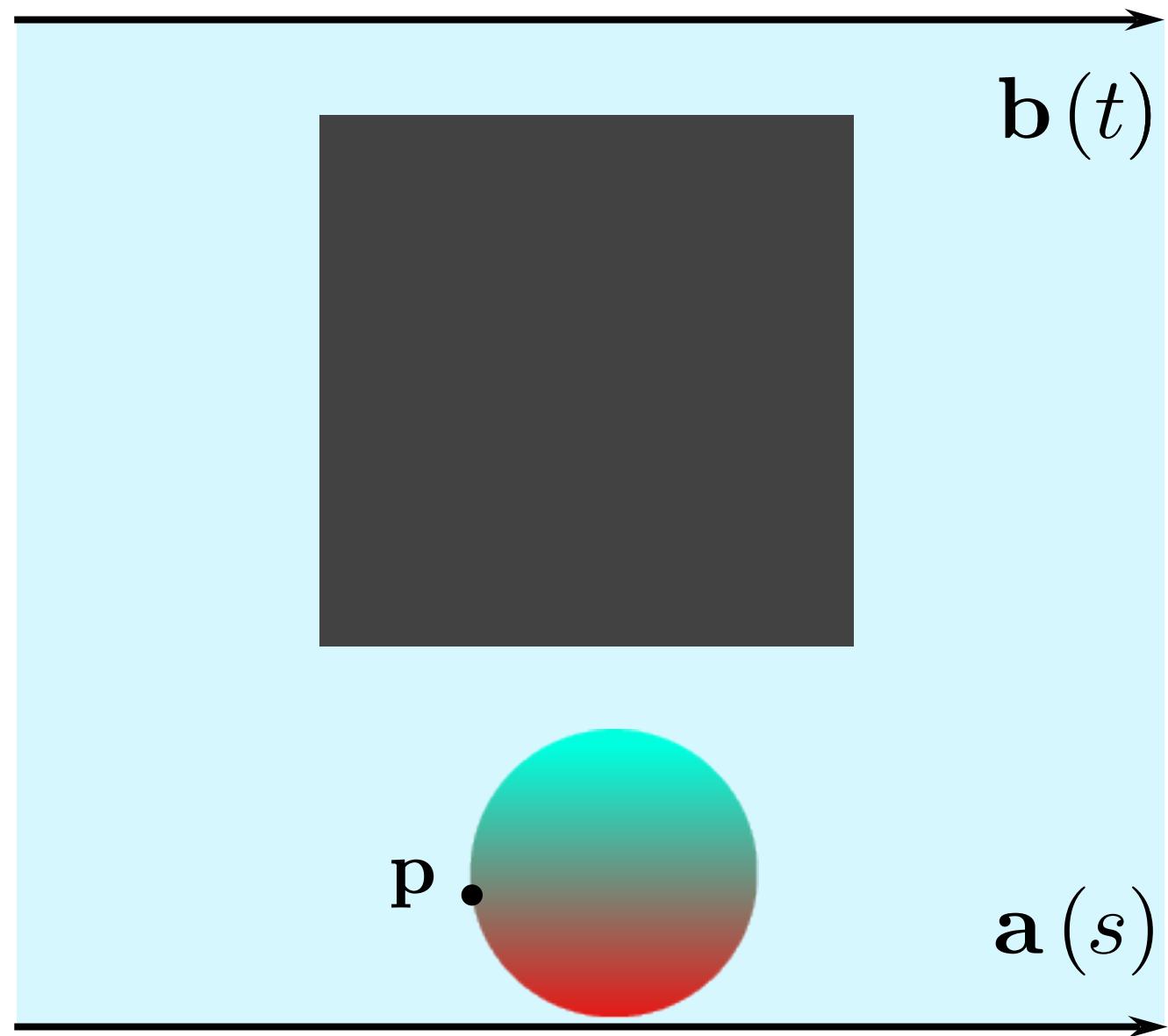
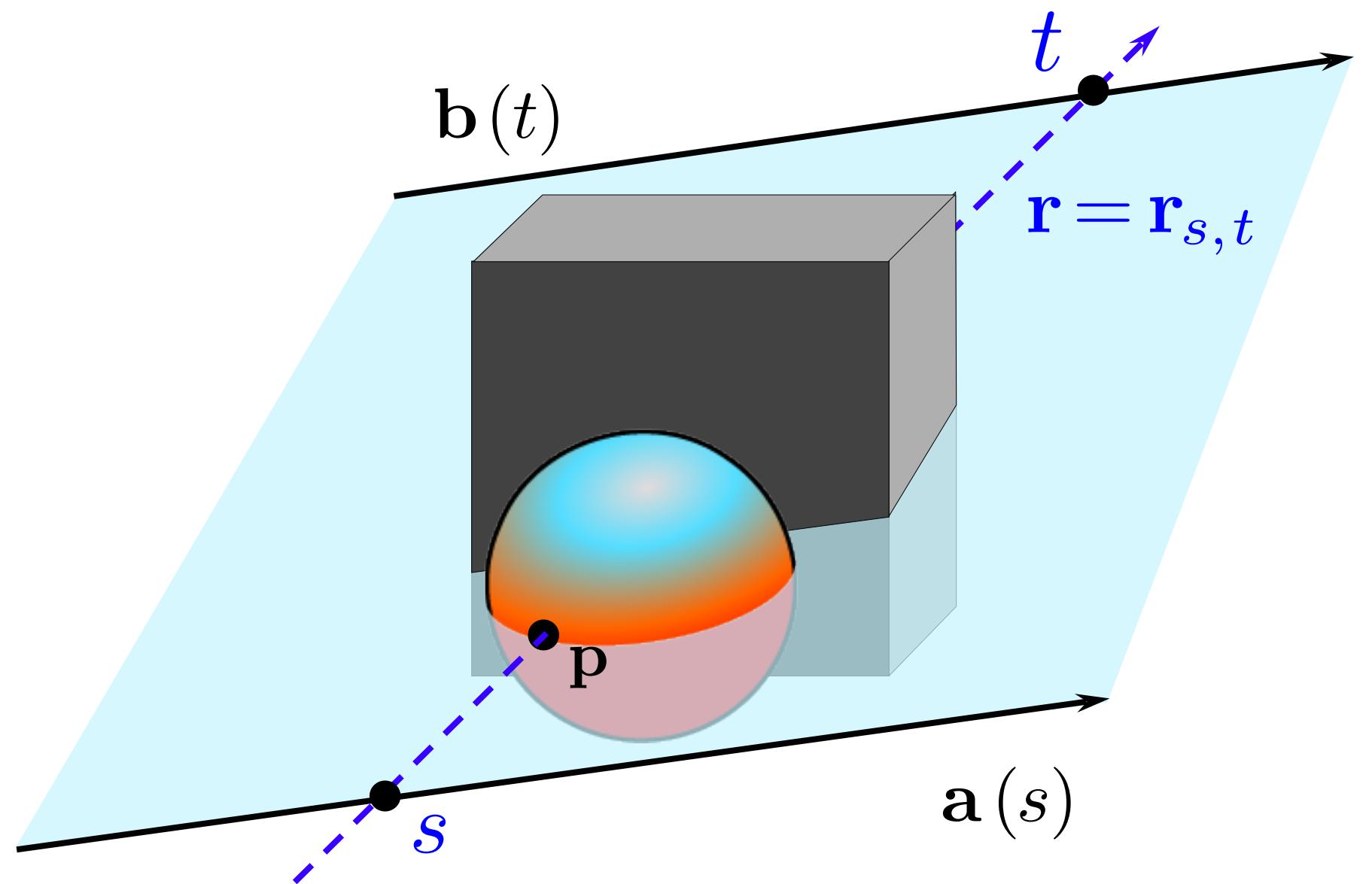
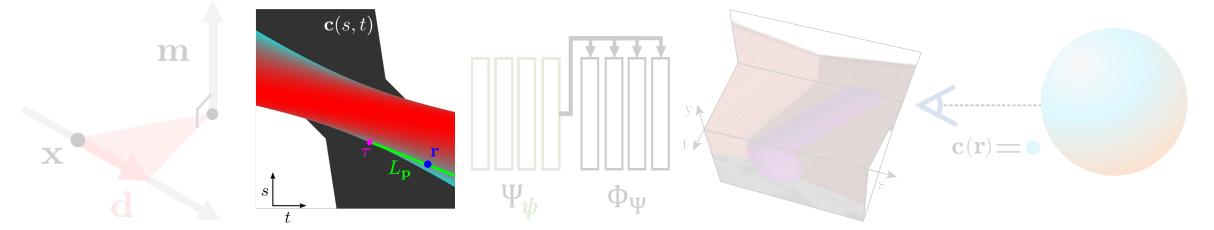


Epipolar Plane Image



Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

The geometry of LFNs

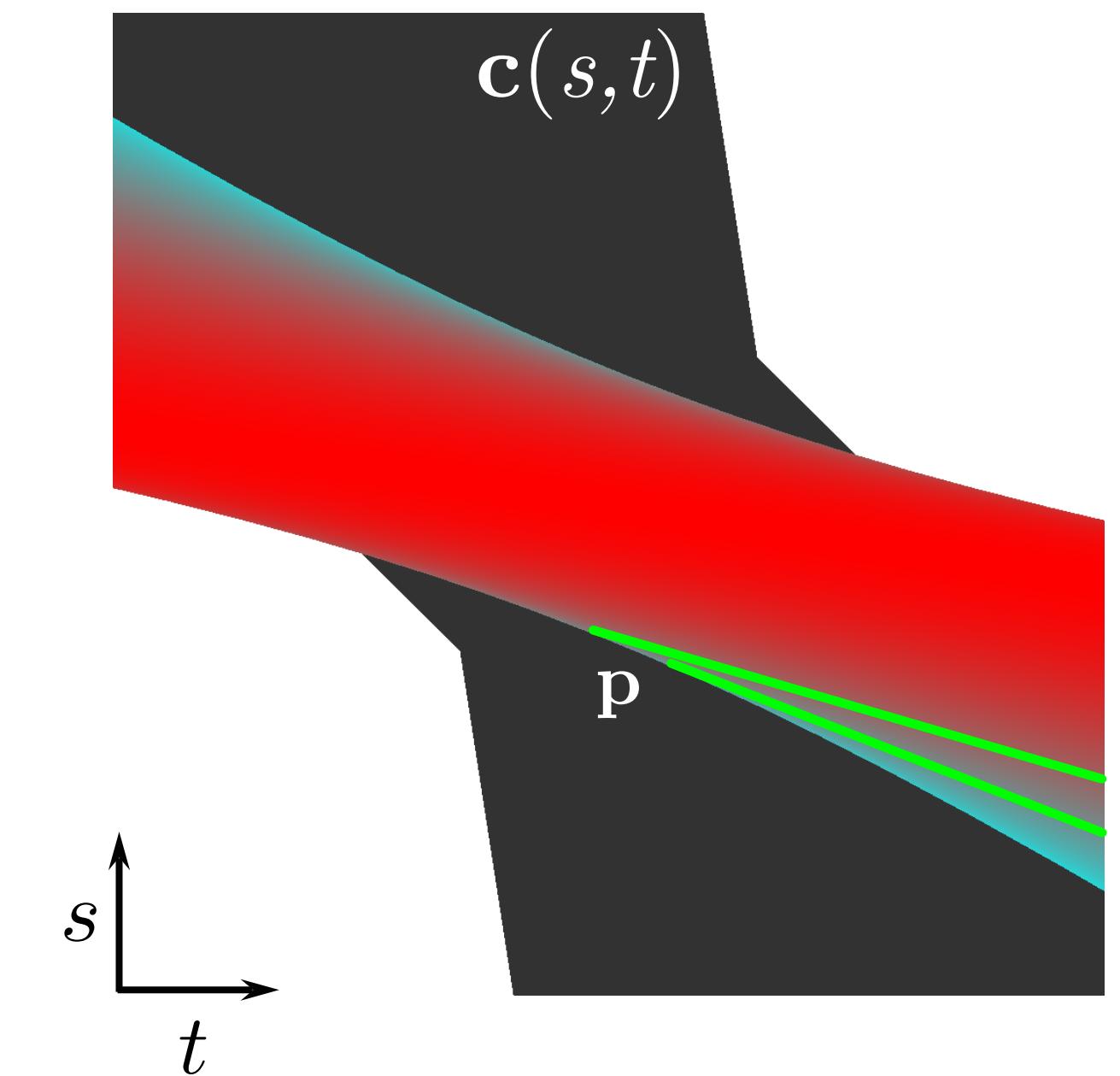
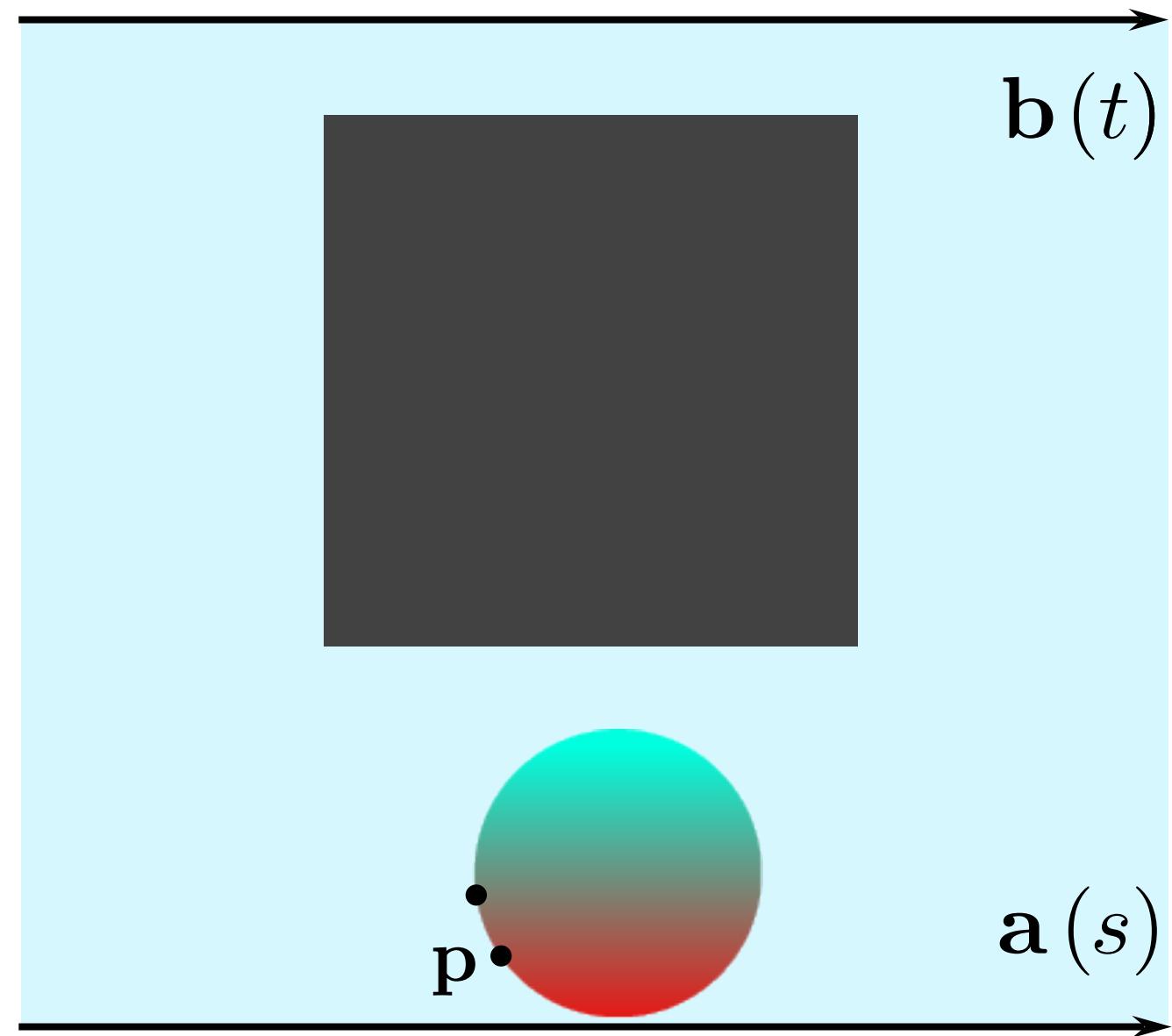
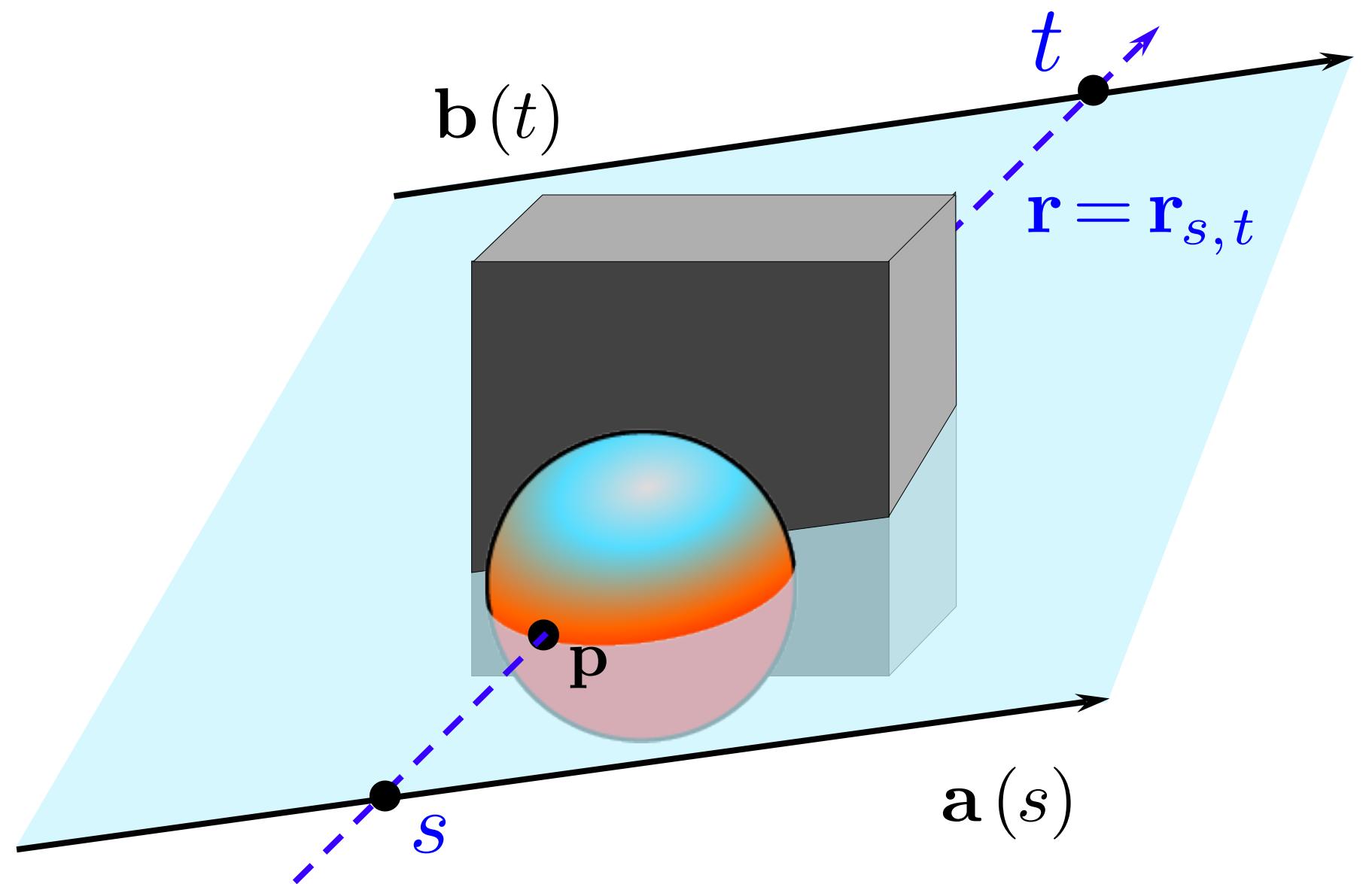
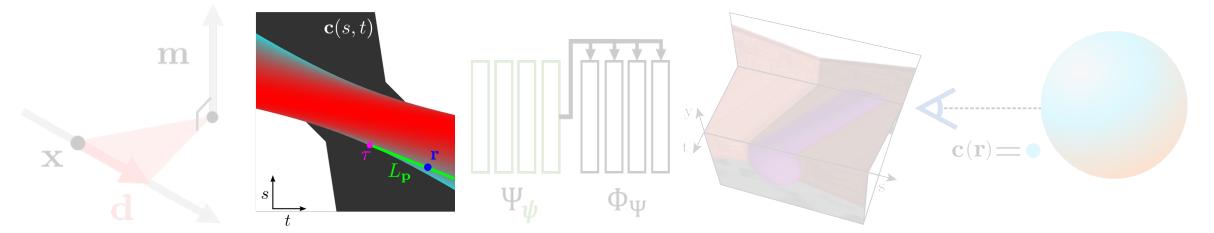


Epipolar Plane Image

Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

Slope of line decreases as point moves closer.

The geometry of LFNs

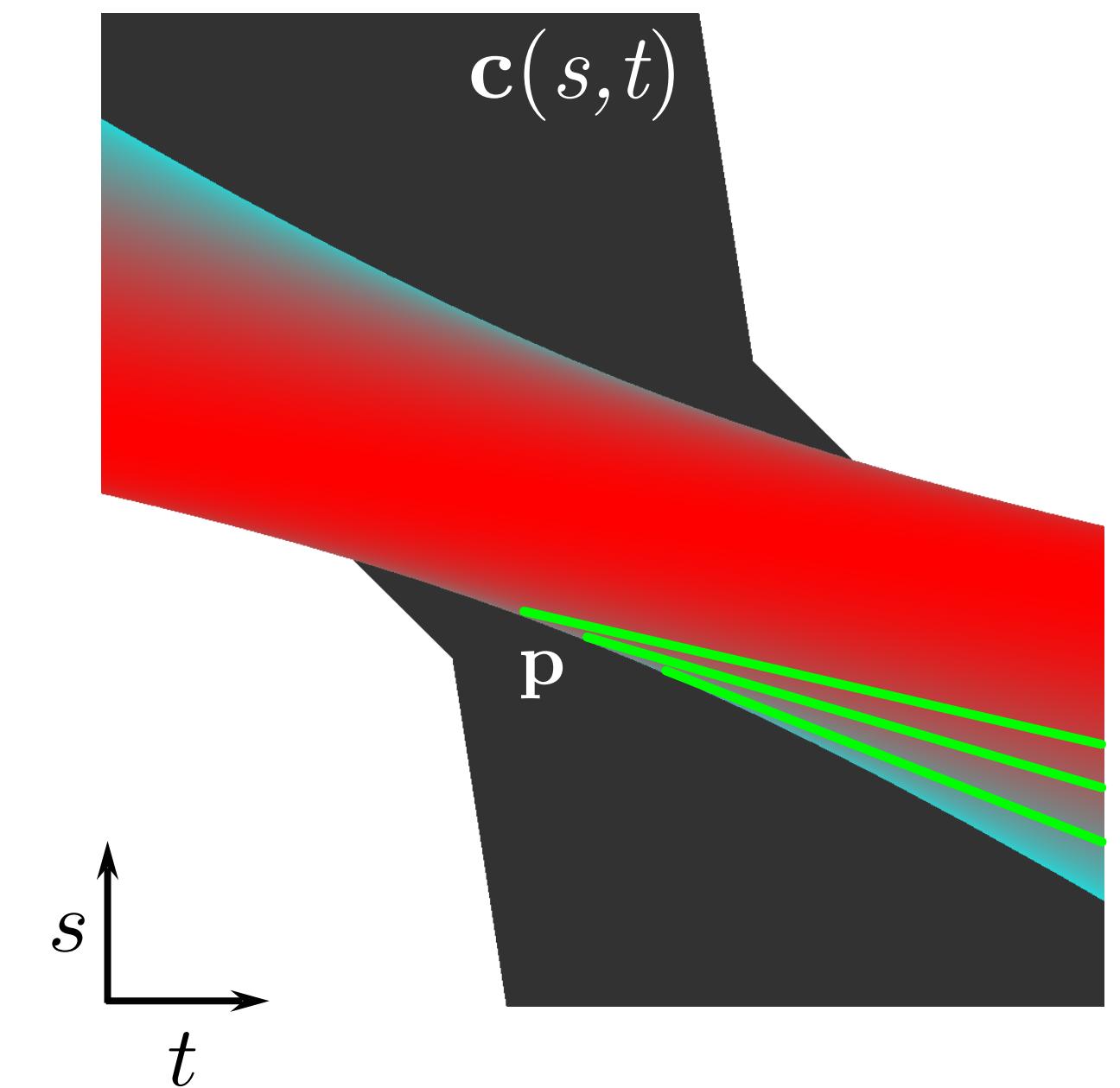
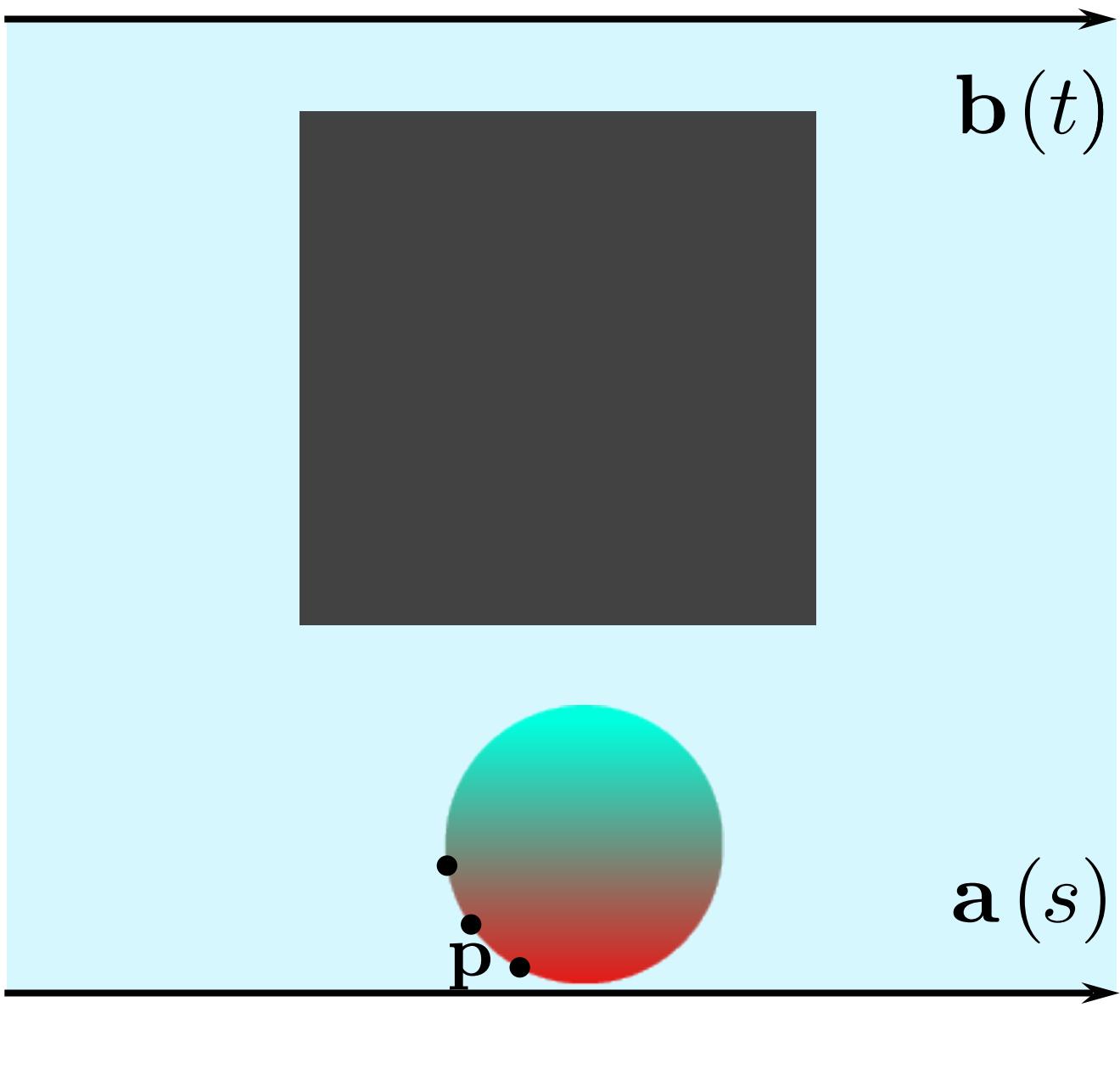
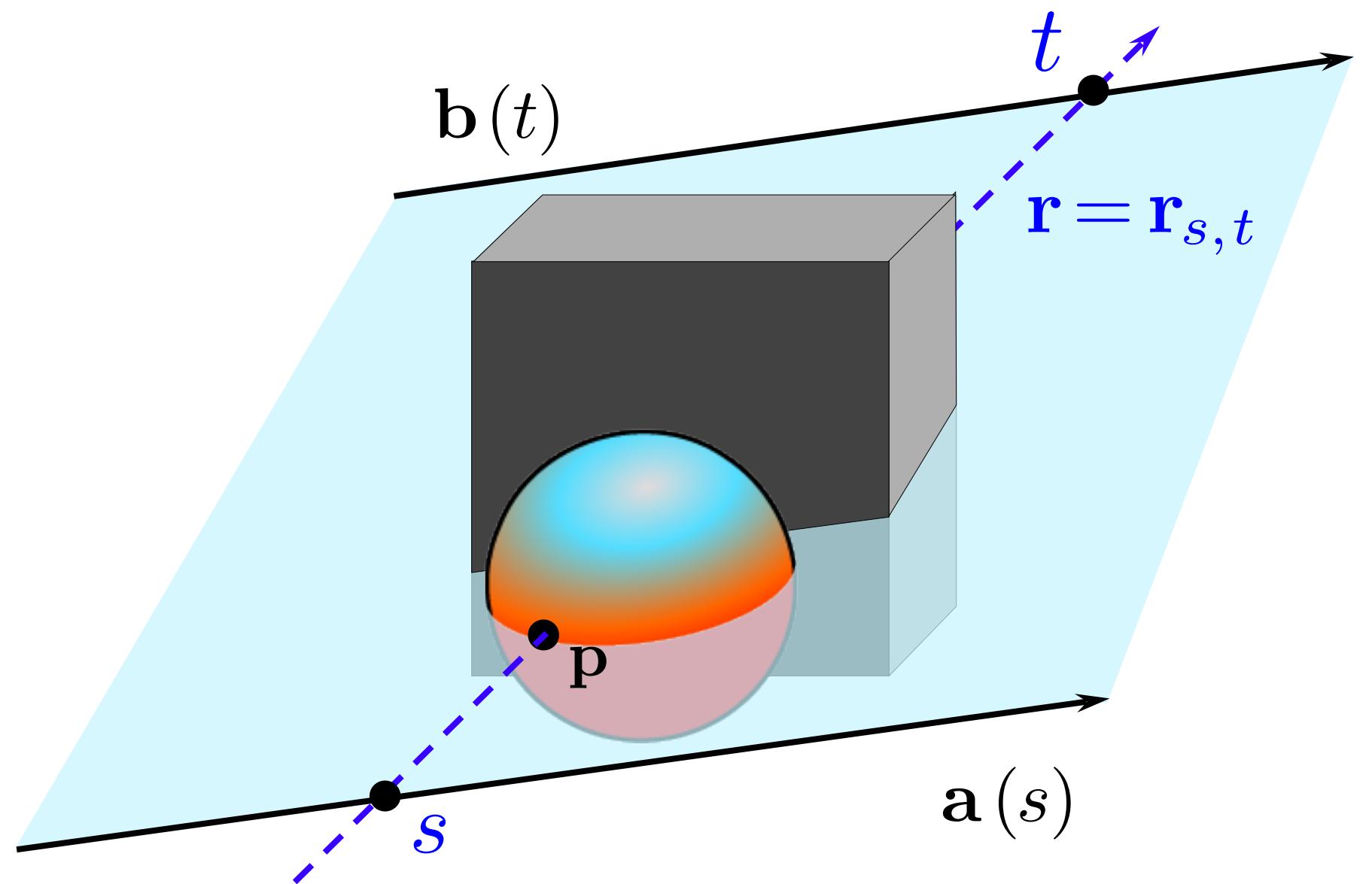
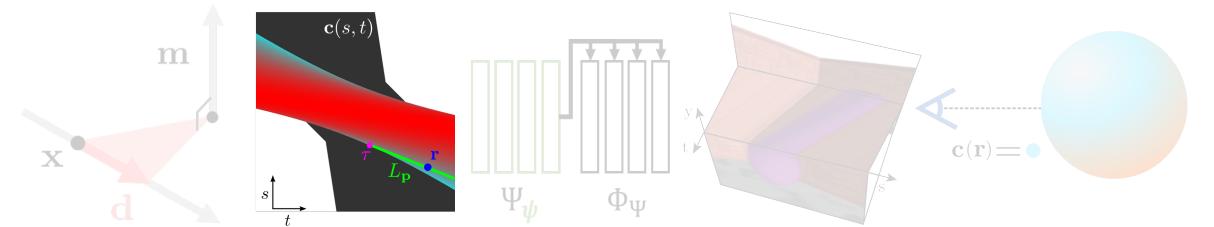


Epipolar Plane Image

Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

Slope of line decreases as point moves closer.

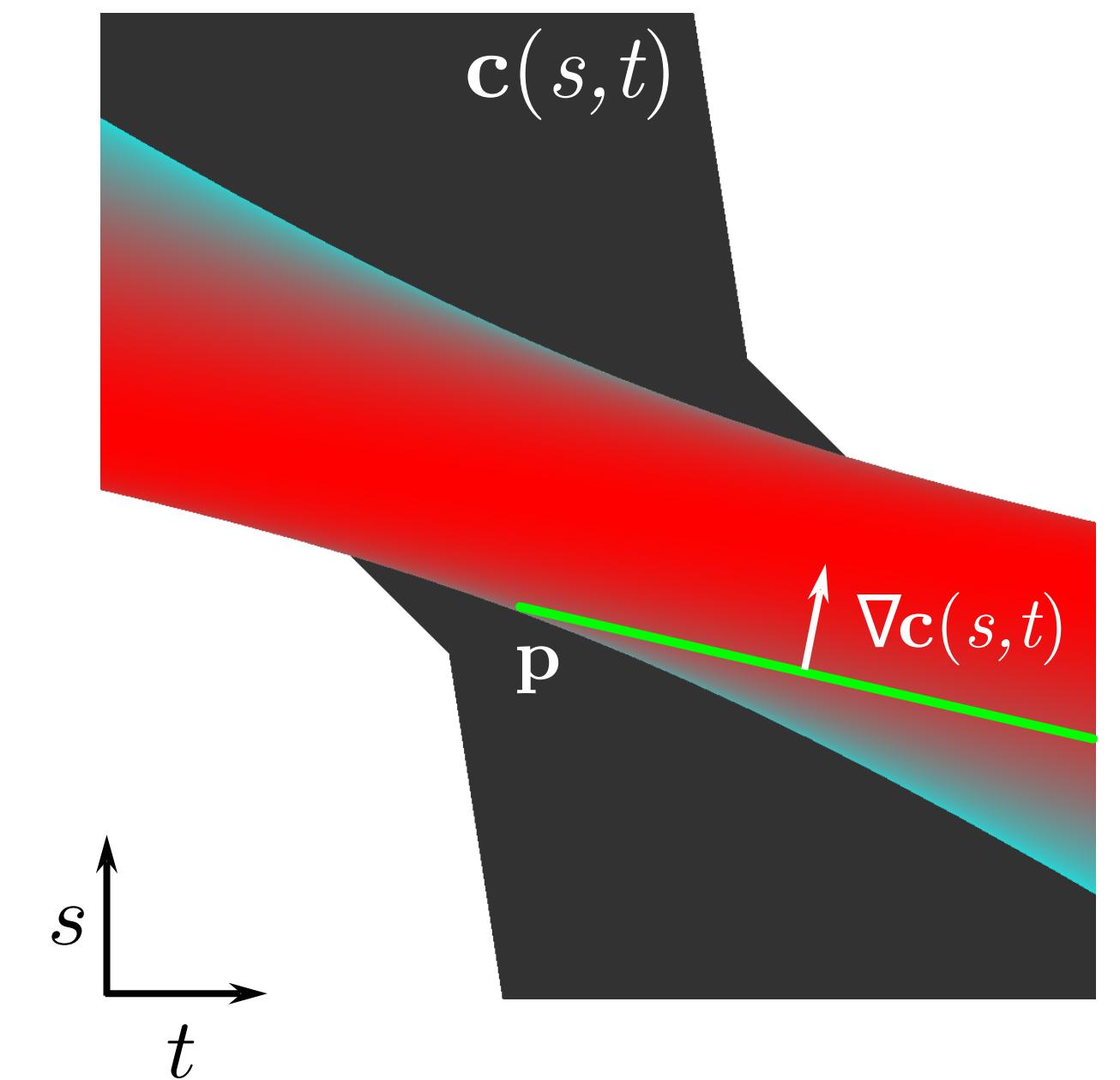
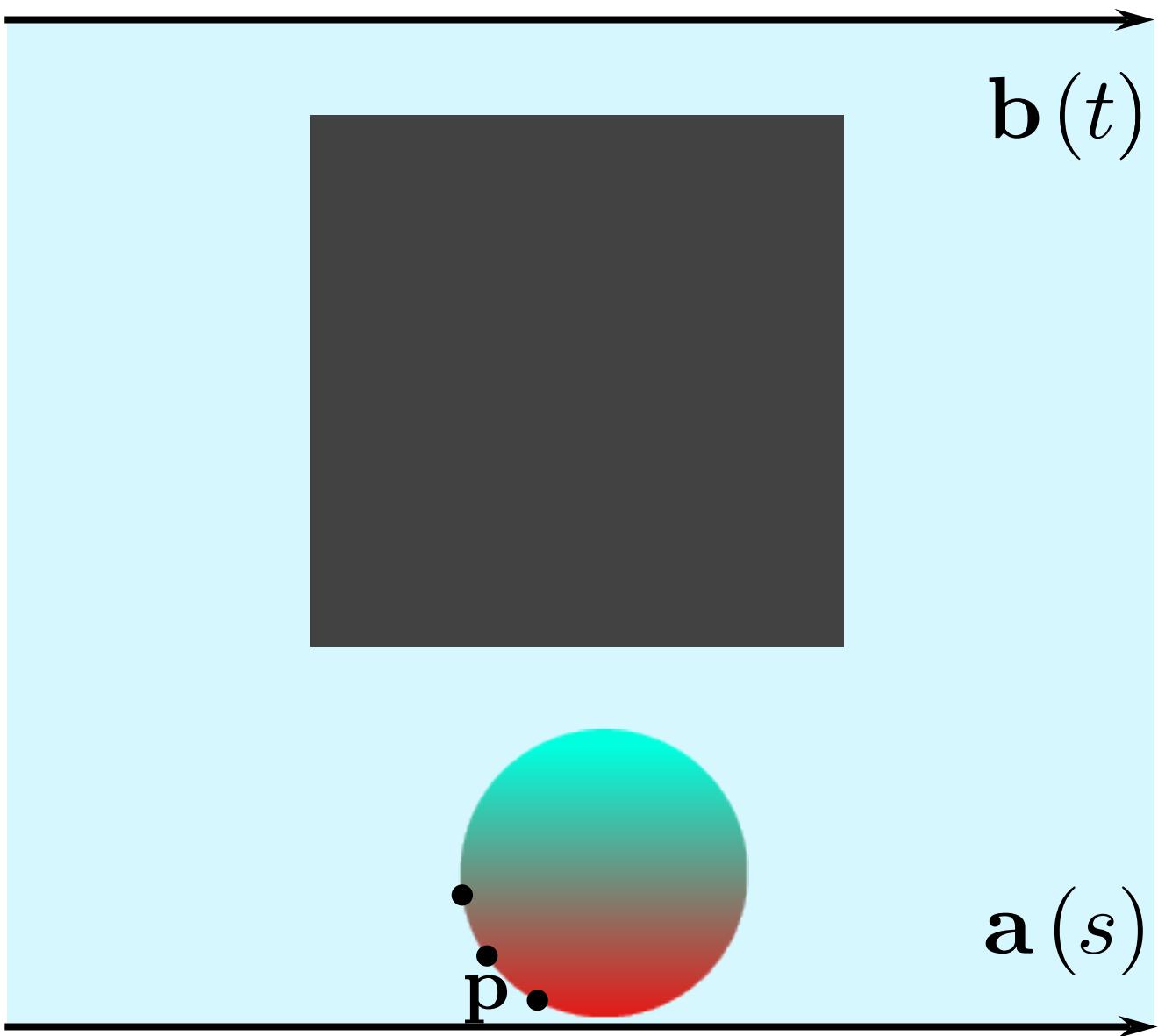
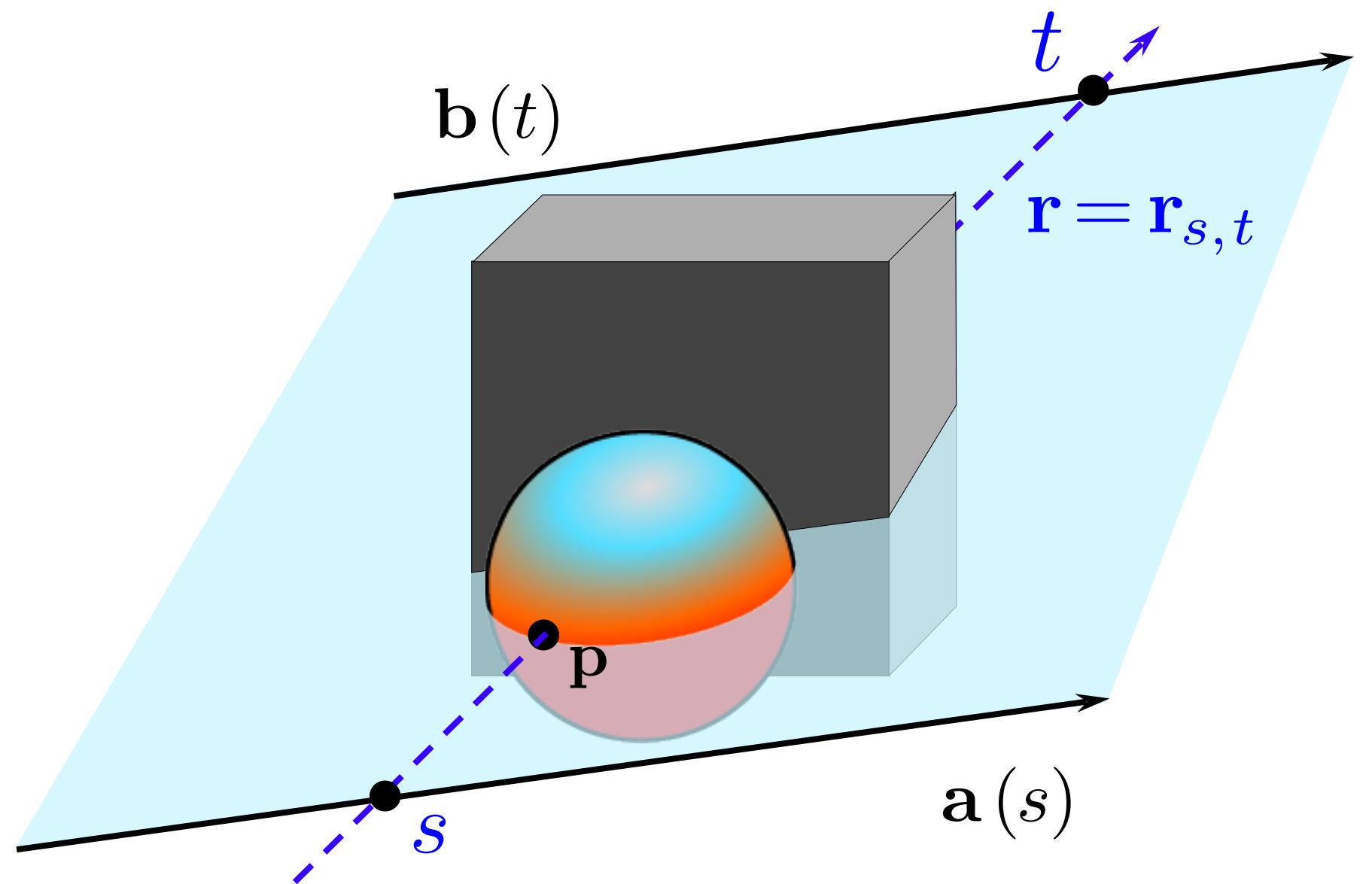
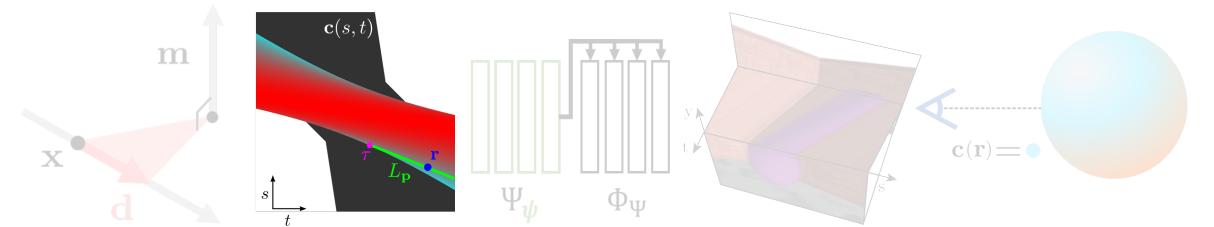
The geometry of LFNs



Points give lines of constant color in EPI $\mathbf{c}(s, t)$ – line is a levelset of the EPI.

Slope of line decreases as point moves closer.

The geometry of LFNs

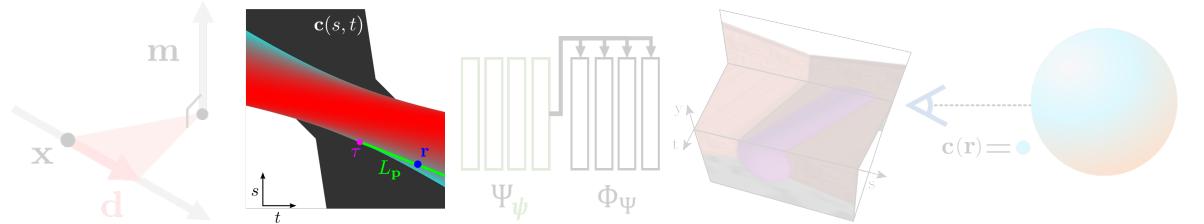


Points give lines of constant color in EPI $\mathbf{c}(s, t)$ – line is a levelset of the EPI.

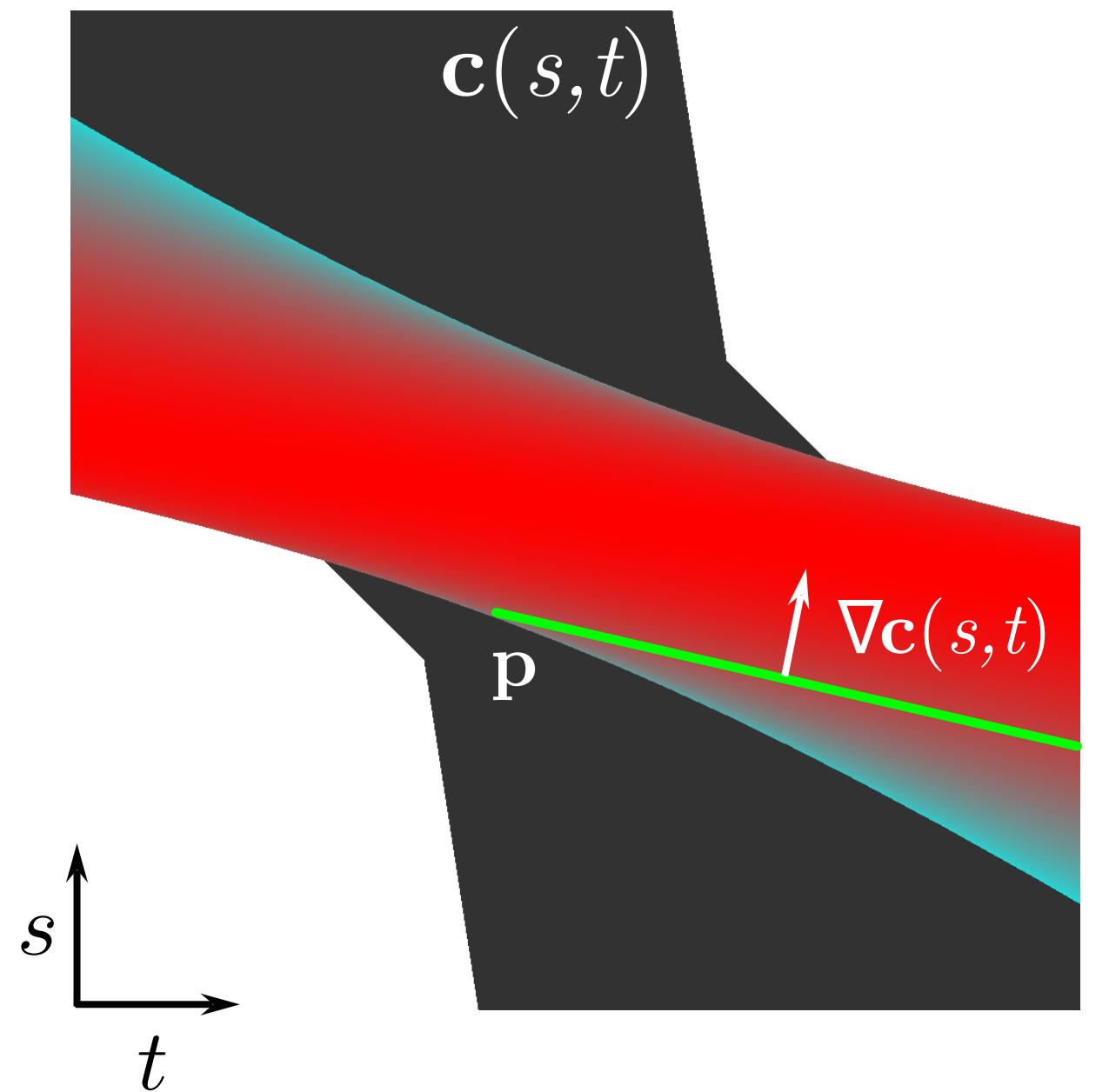
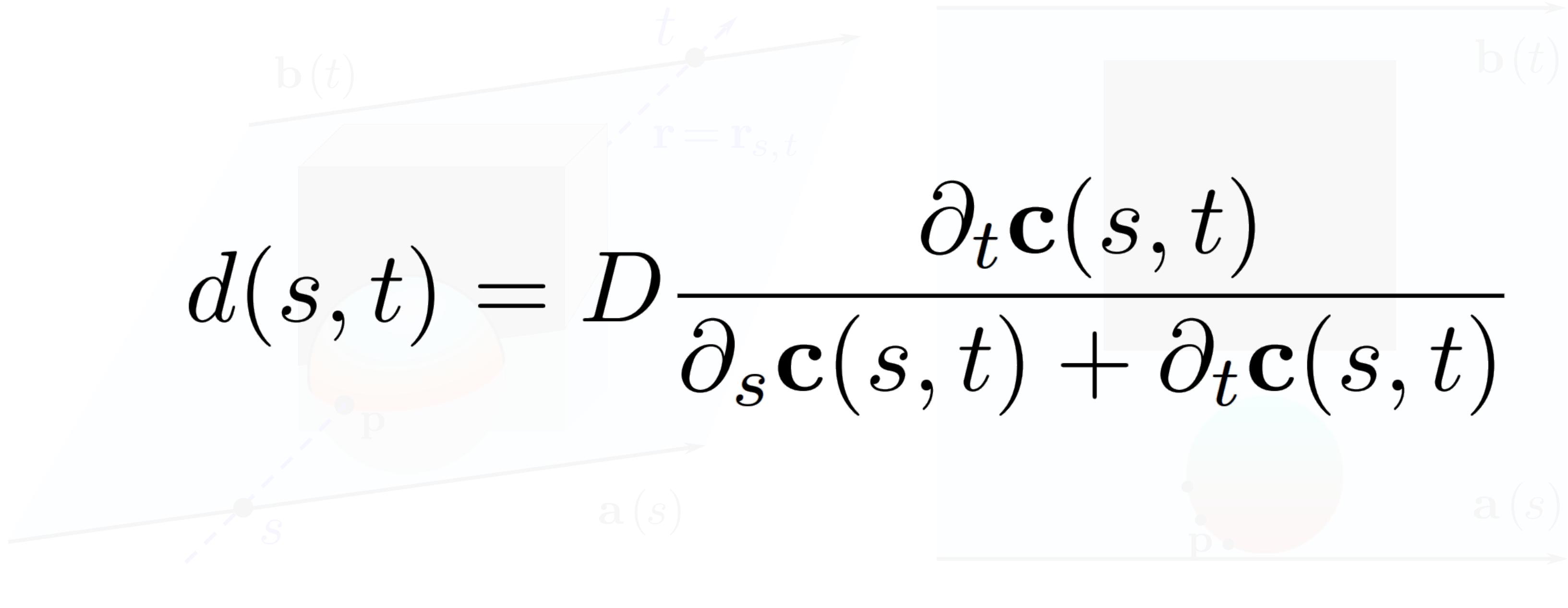
Slope of line decreases as point moves closer.

Gradient of $\mathbf{c}(s, t)$ is orthogonal to levelset -

The geometry of LFNs



Epipolar Plane Image

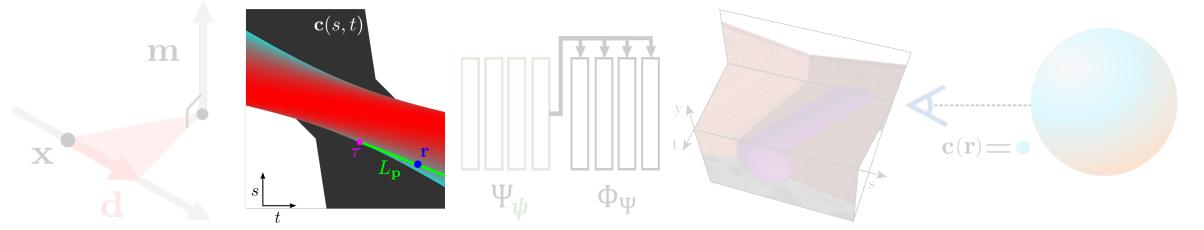


Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

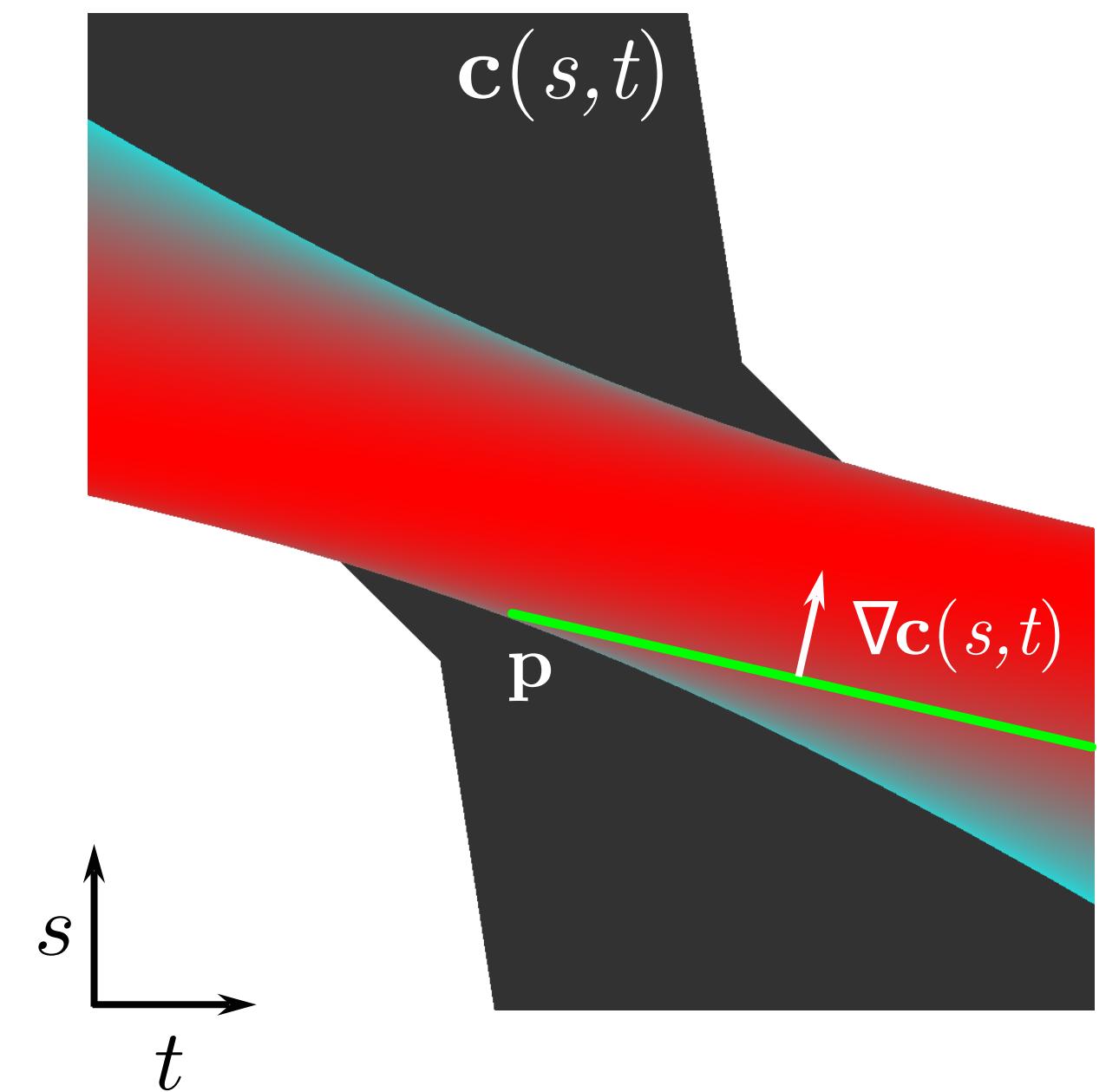
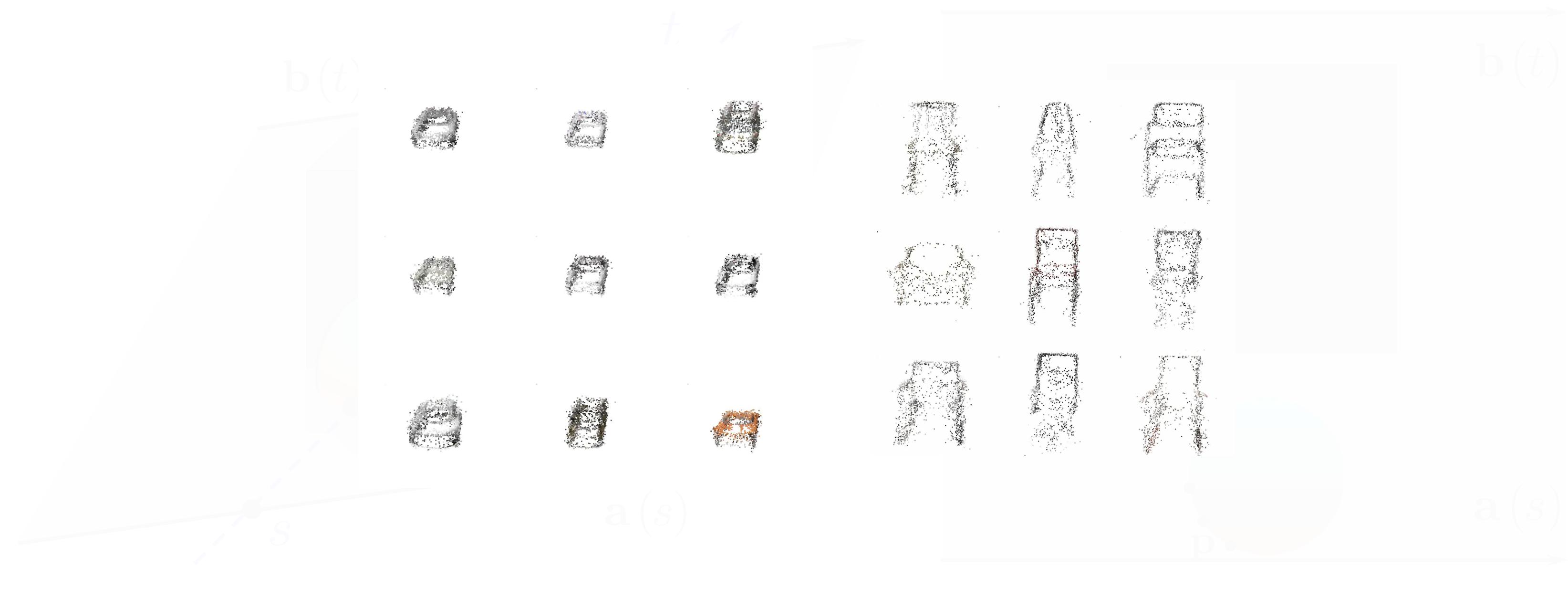
Slope of line decreases as point moves closer.

Gradient of $\mathbf{c}(s,t)$ is orthogonal to levelset - can extract depth from gradients of light field.

The geometry of LFNs



Epipolar Plane Image



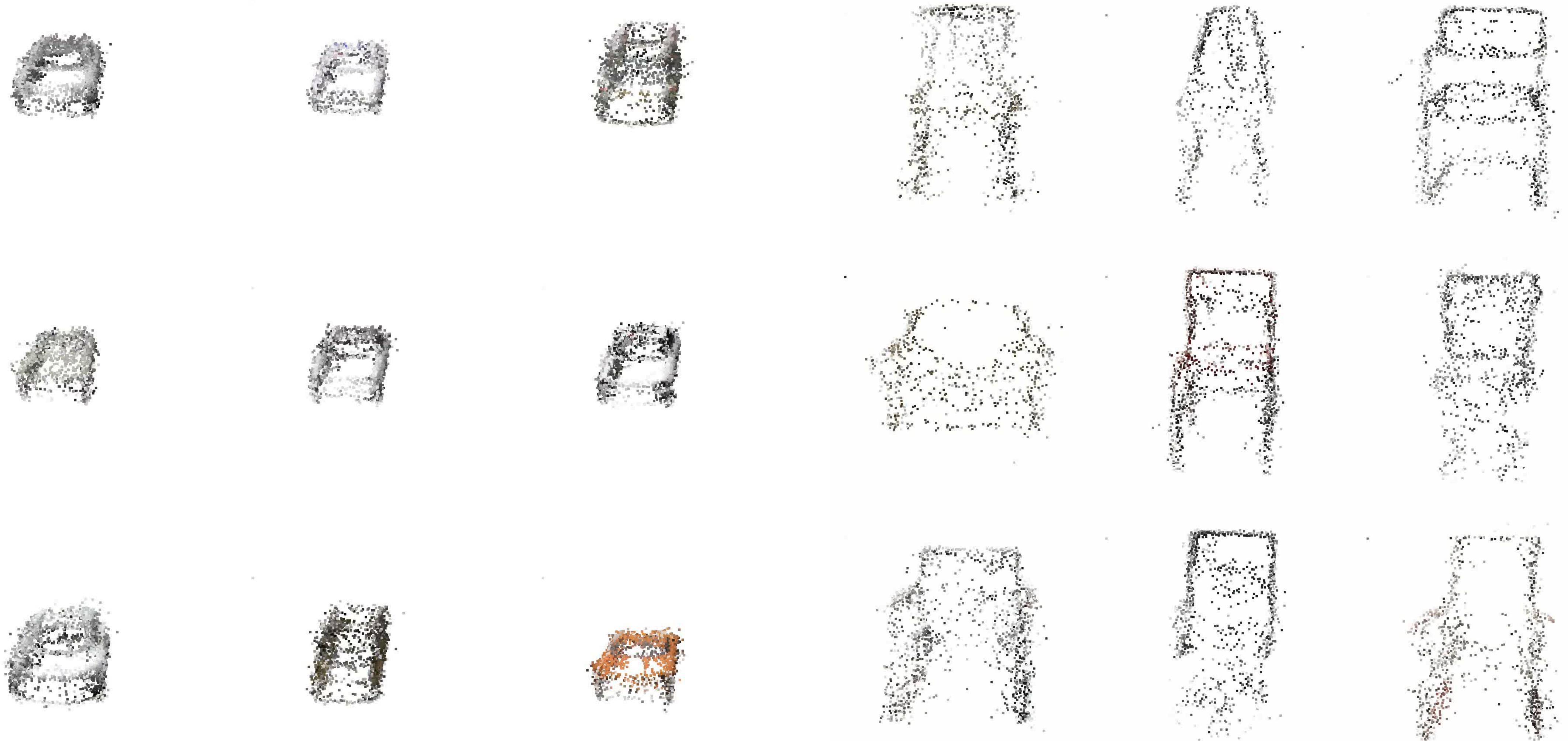
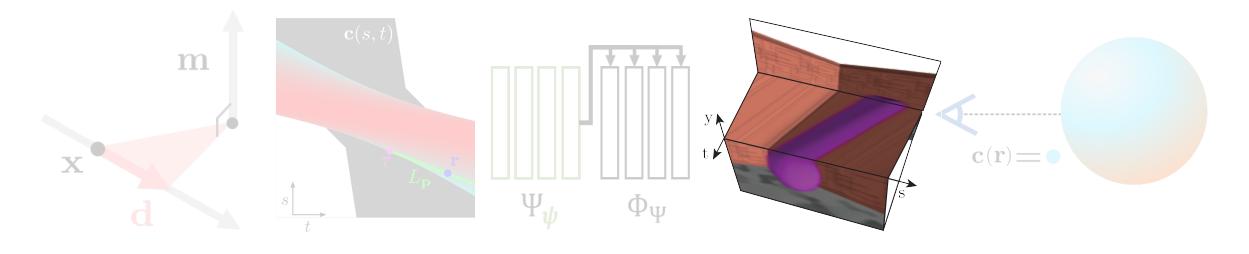
Points give lines of constant color in EPI $\mathbf{c}(s,t)$ – line is a levelset of the EPI.

Slope of line decreases as point moves closer.

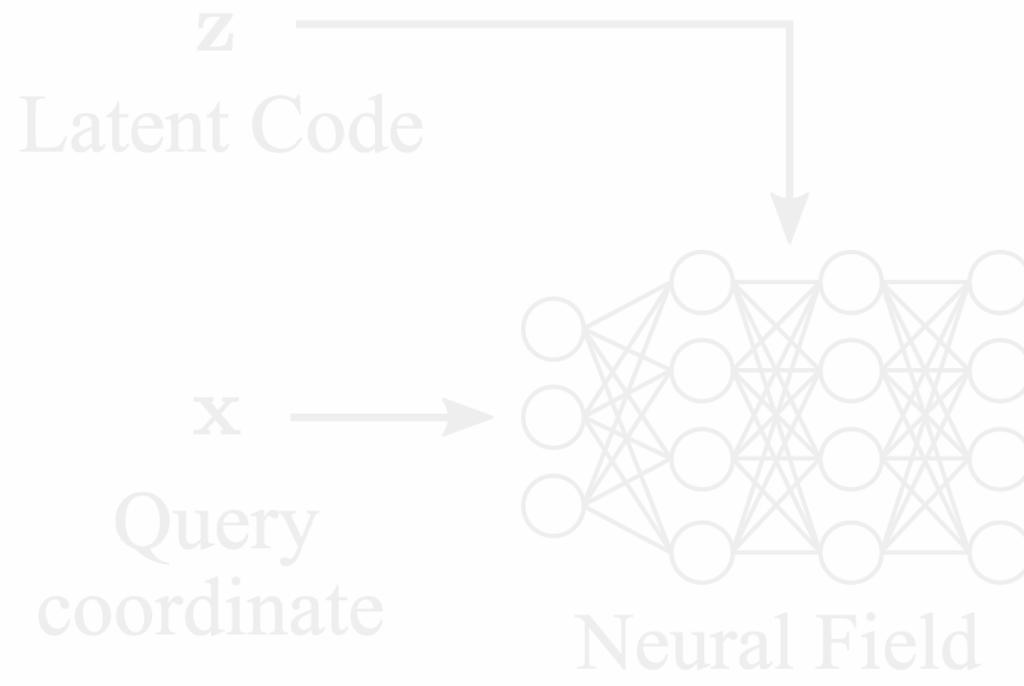
Gradient of $\mathbf{c}(s,t)$ is orthogonal to levelset - can extract depth from gradients of light field.

Pointclouds extracted from 4 views of cars & chairs

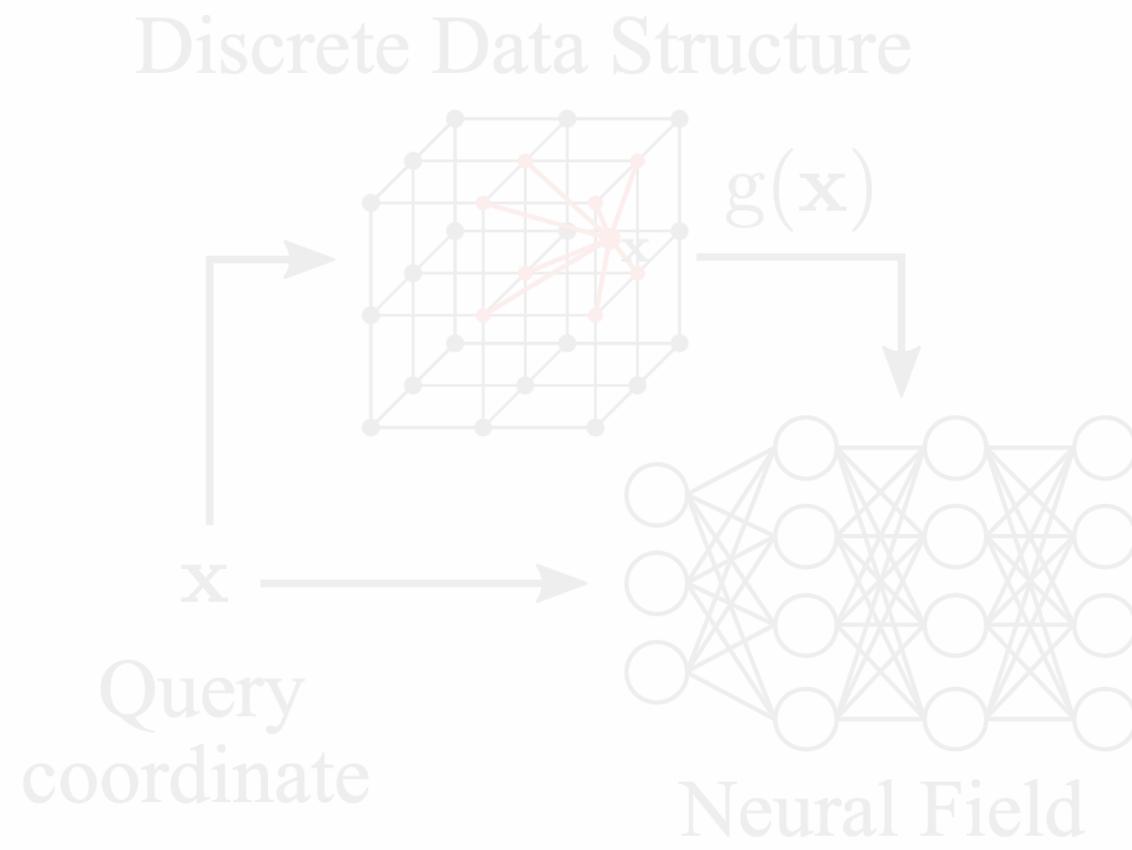
Single evaluation of network & its gradient per ray, constant complexity



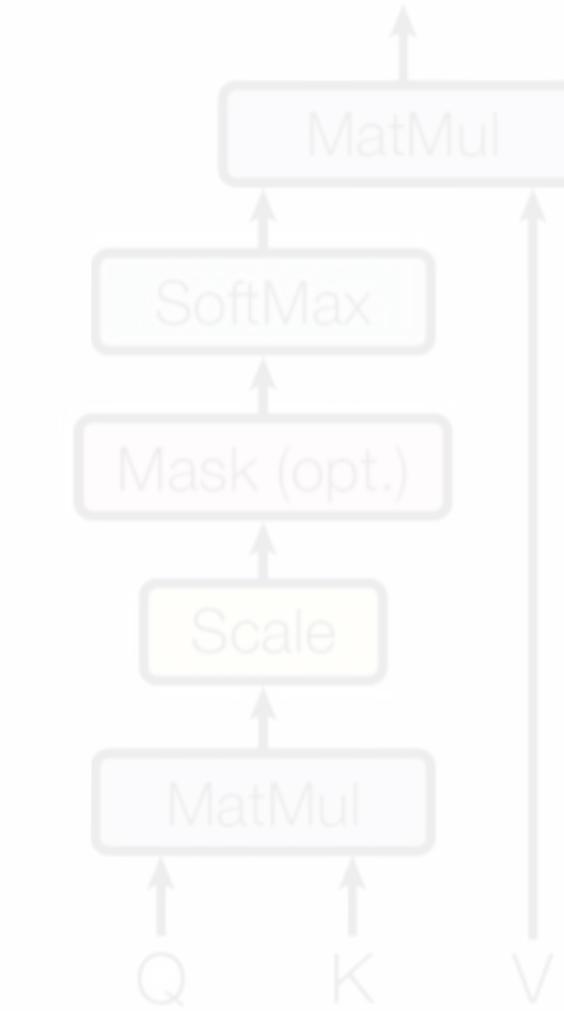
Other types of conditioning...?



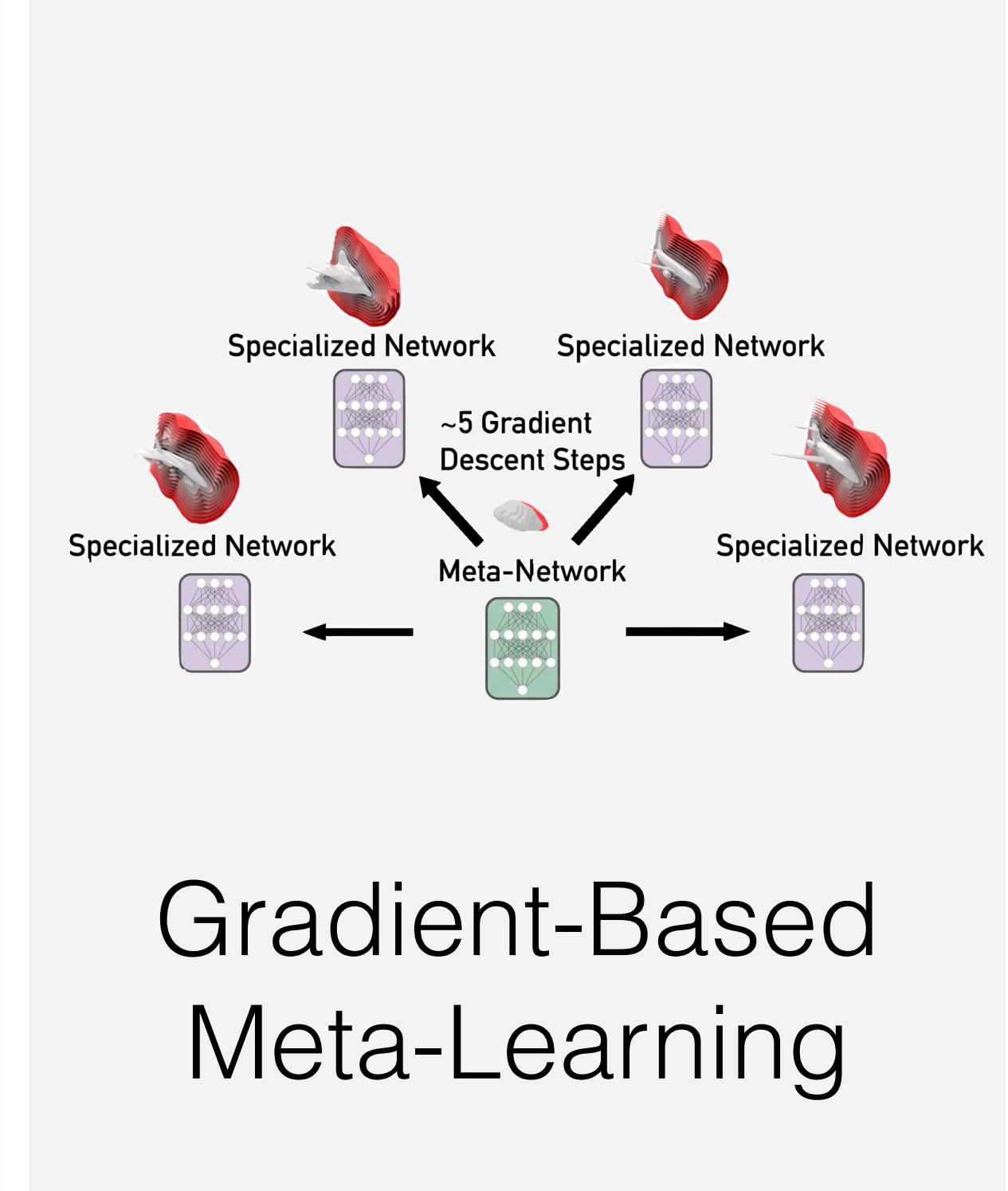
Global
Conditioning



Local
Conditioning

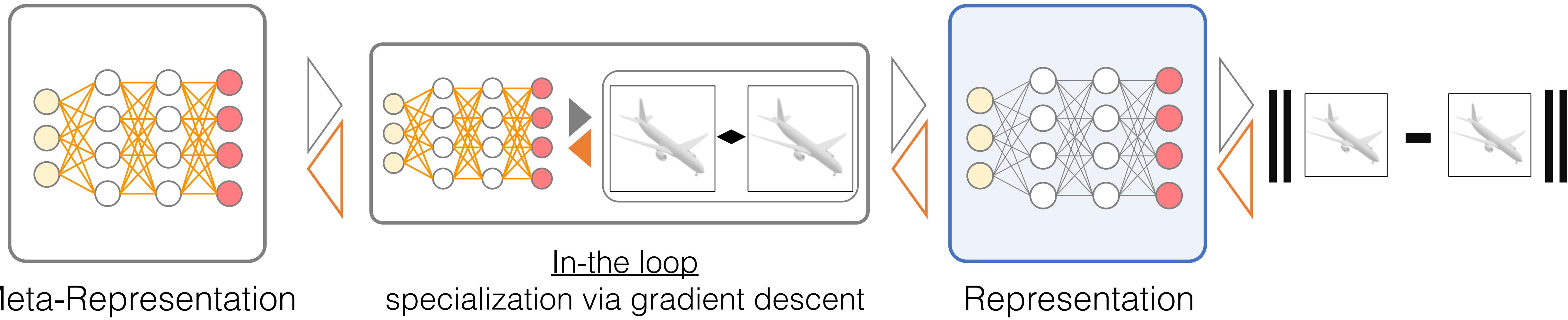


Attention-Based
Conditioning



Gradient-Based
Meta-Learning

Other forms of Generalization: Gradient-based meta-learning



Other forms of Generalization: Gradient-based meta-learning



Meta-Representations

Behaves somewhat similar to global conditioning.

gradient-based learning time.

Explains held-out observations when fit to context observation.

Next steps

- Oct 6th: Paper Session
- Oct 11th: Student Holiday
- Oct 18th: Paper Session
- Oct 20th: Lecture on Unconditional Generative Models of 3D Scenes!