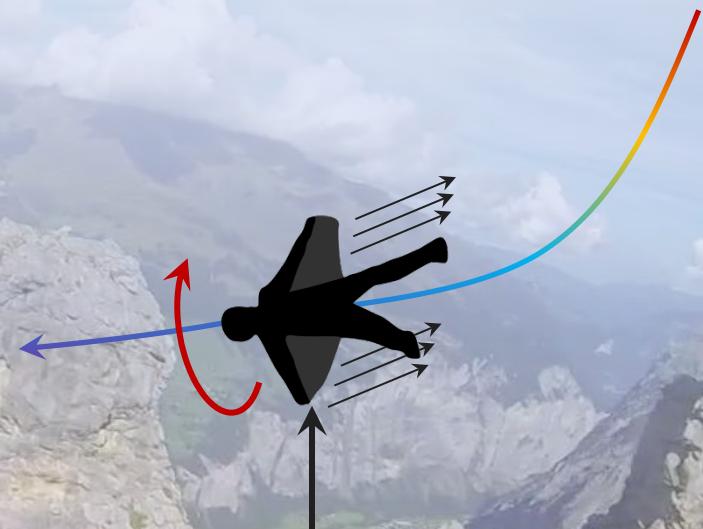
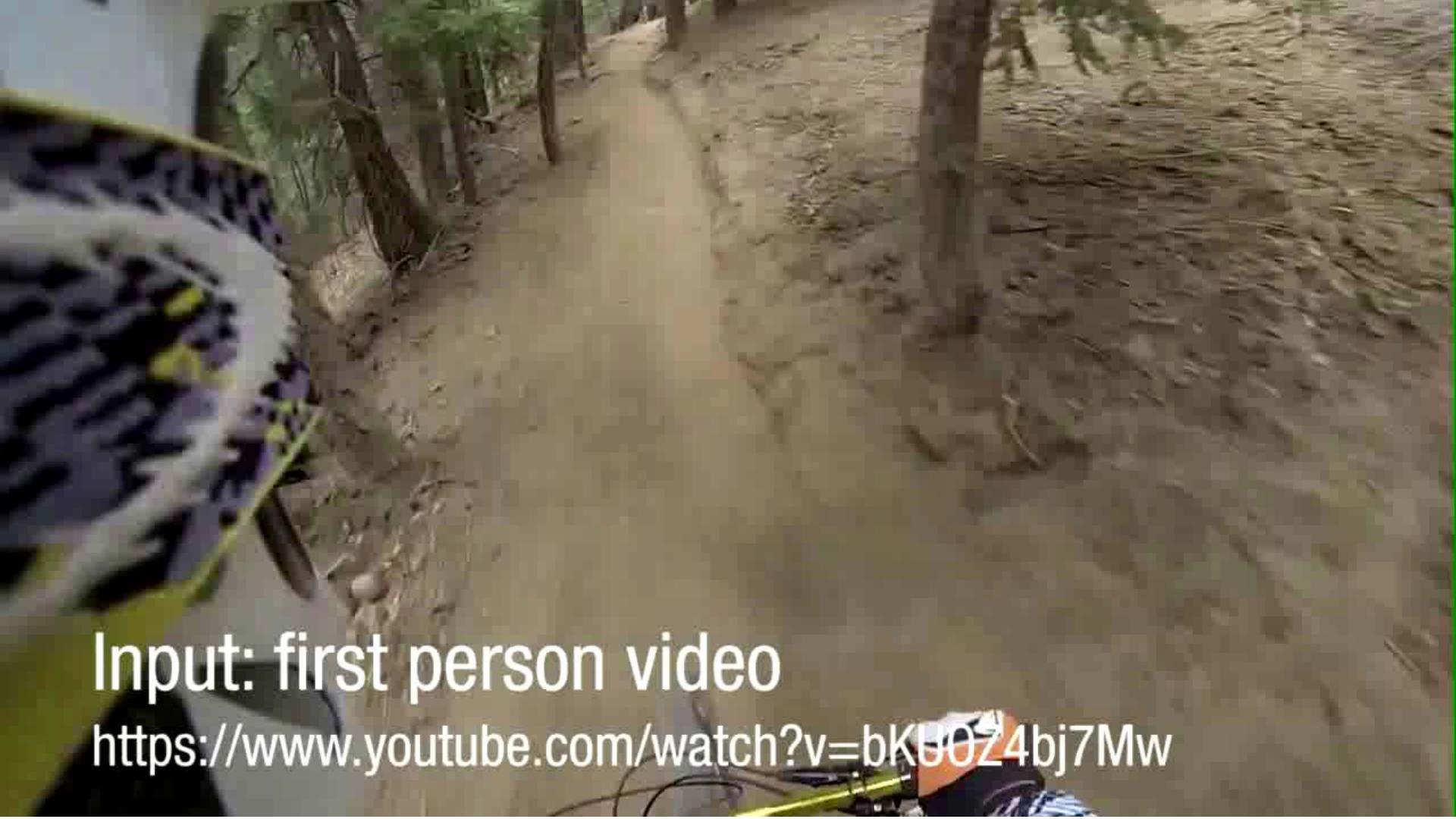


Visual Sensorimotor Behaviors II:

What can a first person video tell about how we control?





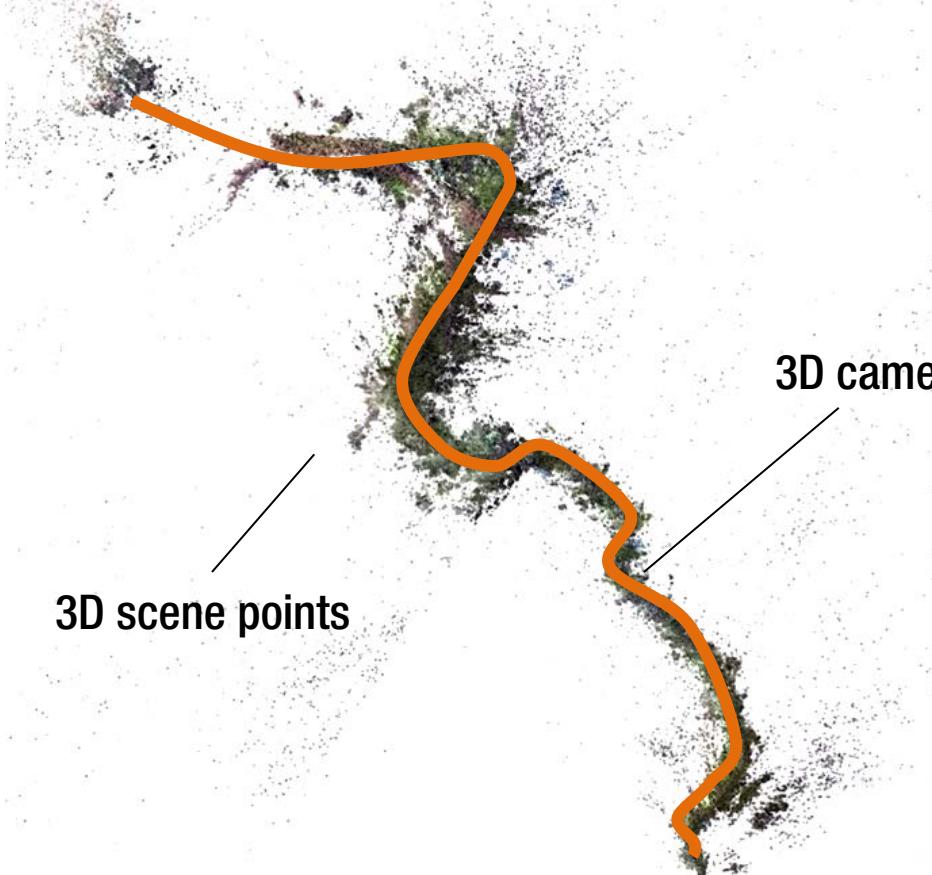
Input: first person video

<https://www.youtube.com/watch?v=bKs0z4bj7Mw>



3D reconstruction

3D reconstruction



3D scene points

3D camera trajectory



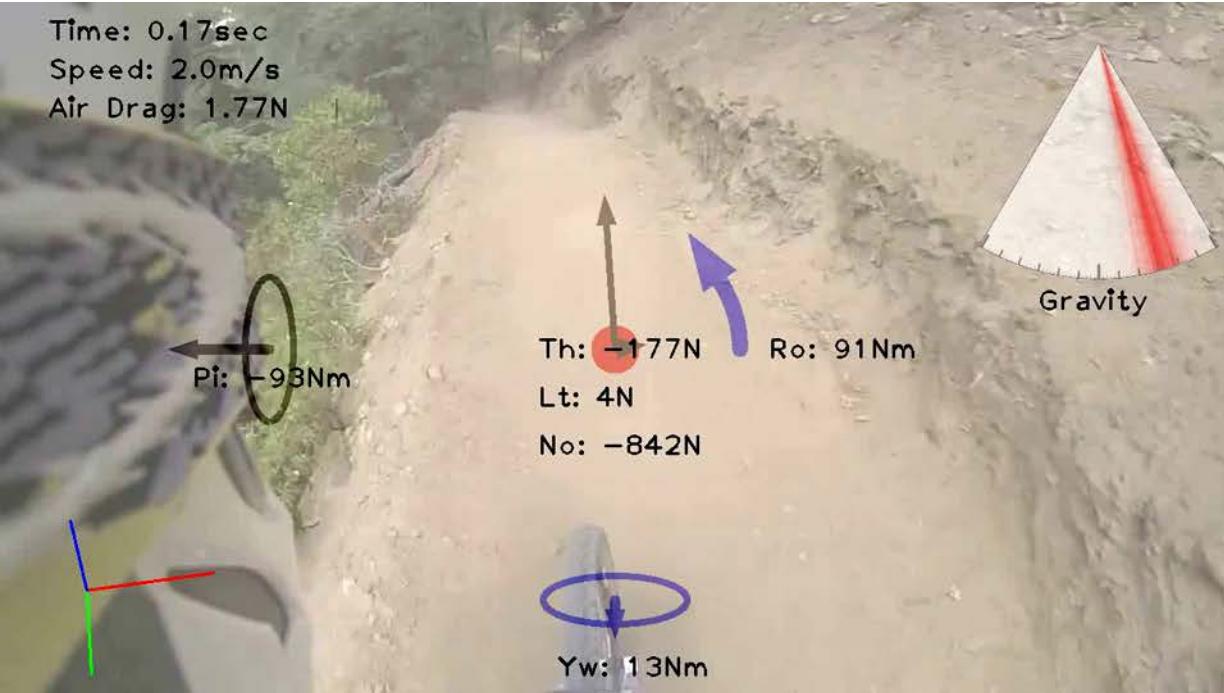


Is geometry or kinematics enough to
understand the biker's behaviors?

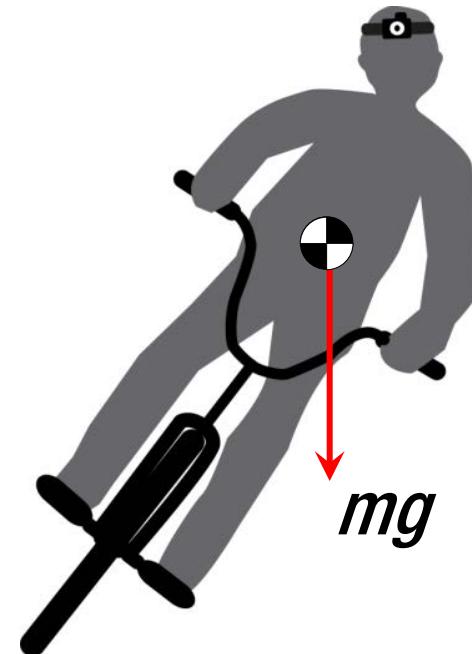
What causes motion?

3D reconstruction

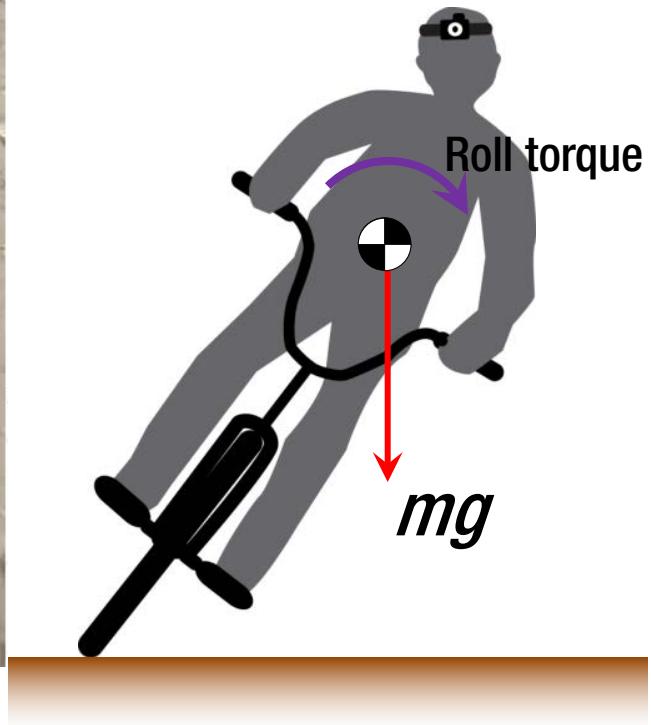
Time: 0.17sec
Speed: 2.0m/s
Air Drag: 1.77N



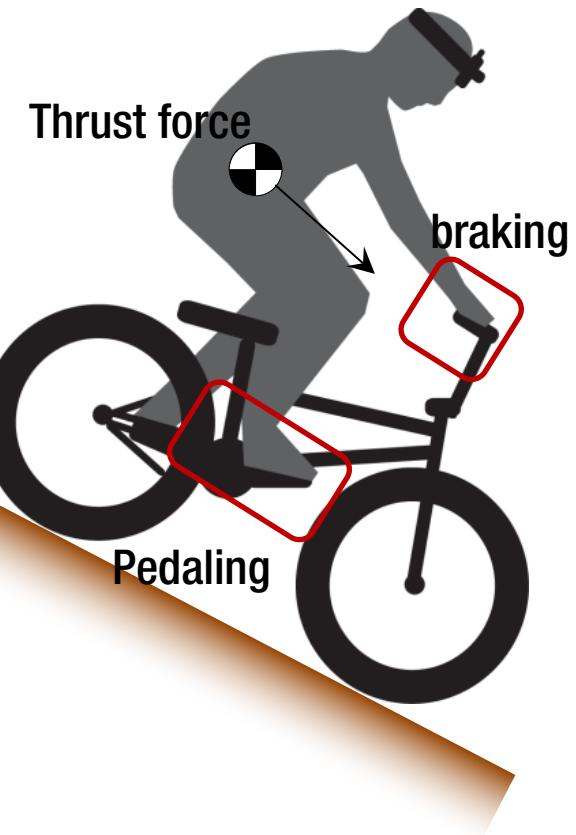
Time: 0.17sec
Speed: 2.0m/s
Air Drag: 1.77N



Time: 0.17sec
Speed: 2.0m/s
Air Drag: 1.77N



Time: 0.17sec
Speed: 2.0m/s
Air Drag: 1.77N

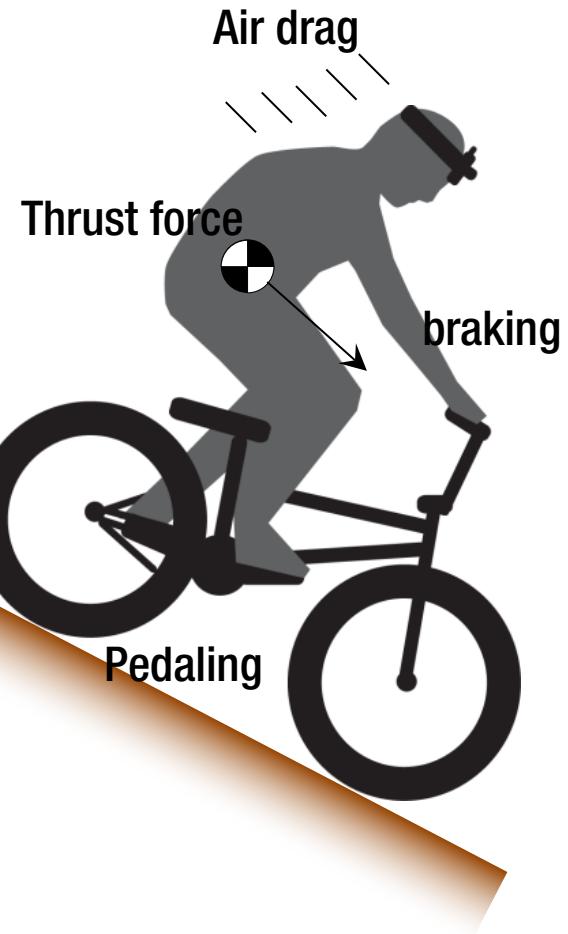


Time: 0.17sec
Speed: 2.0m/s
Air Drag: 1.77N

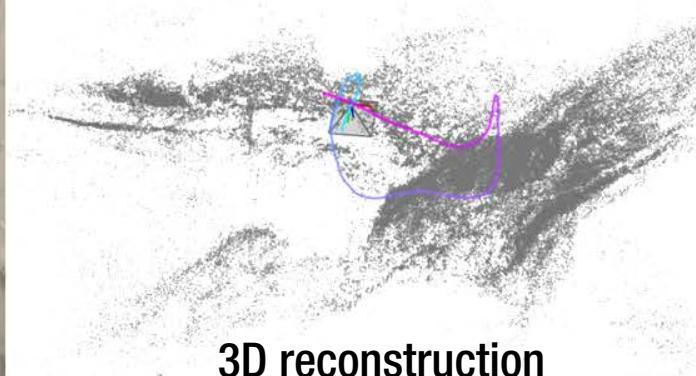
Pi: -93Nm

Th: -177N
Lt: 4N
No: -842N

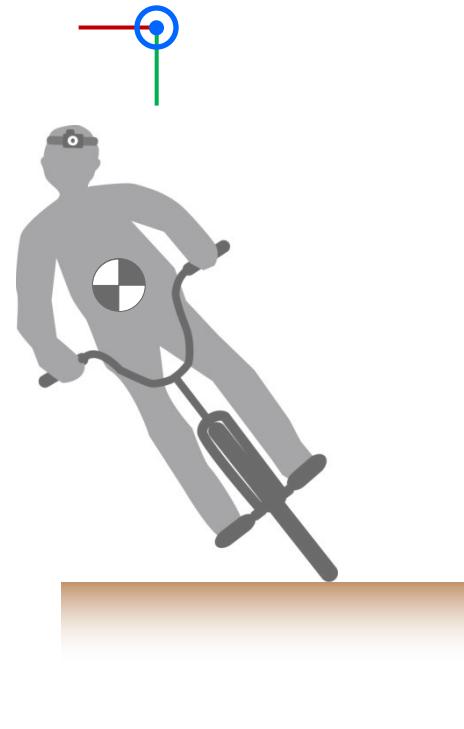
Yw: 13Nm

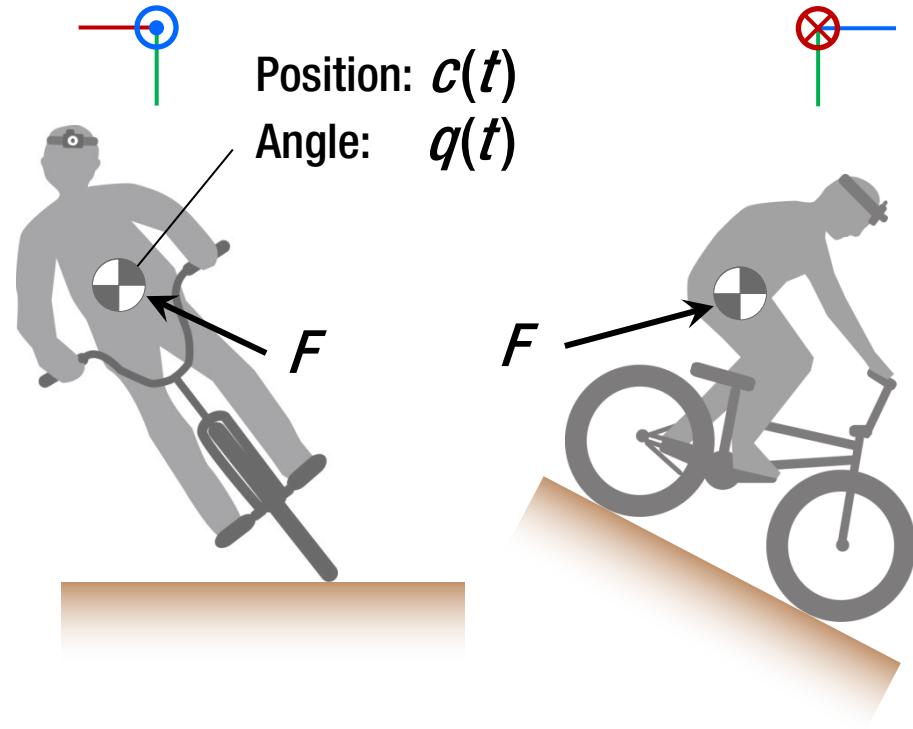


Time: 0.17sec
Speed: 2.0m/s
Air Drag: 1.77N



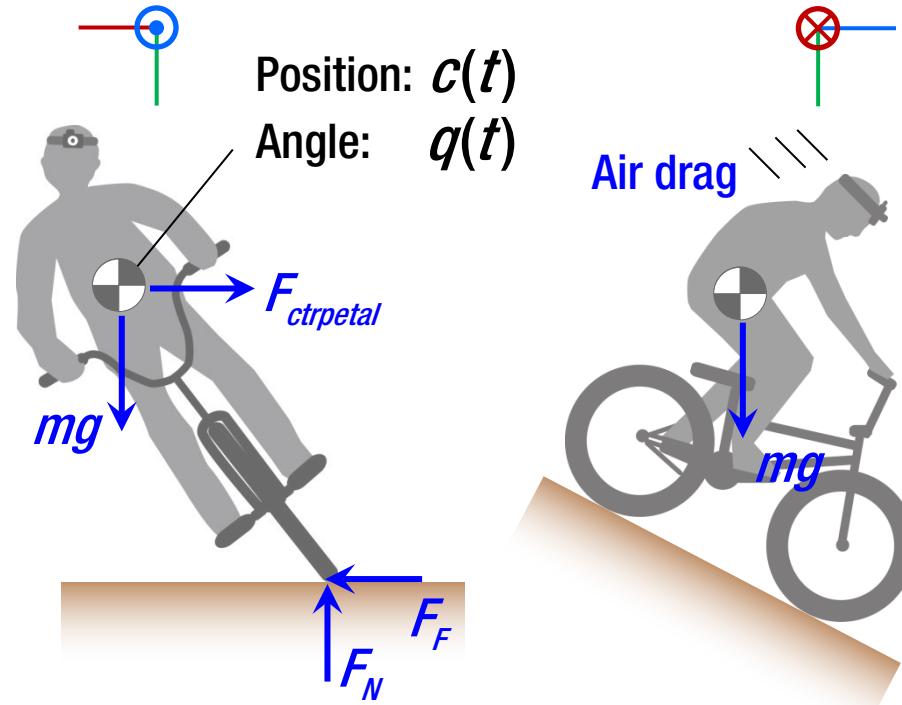
3D reconstruction





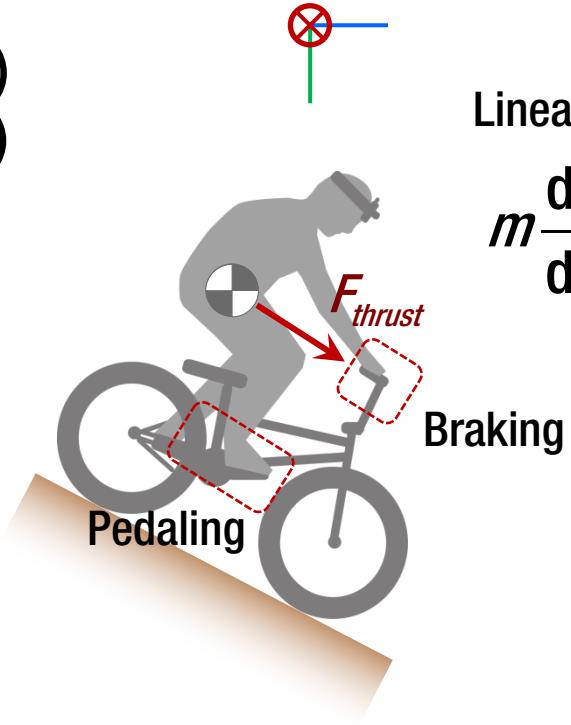
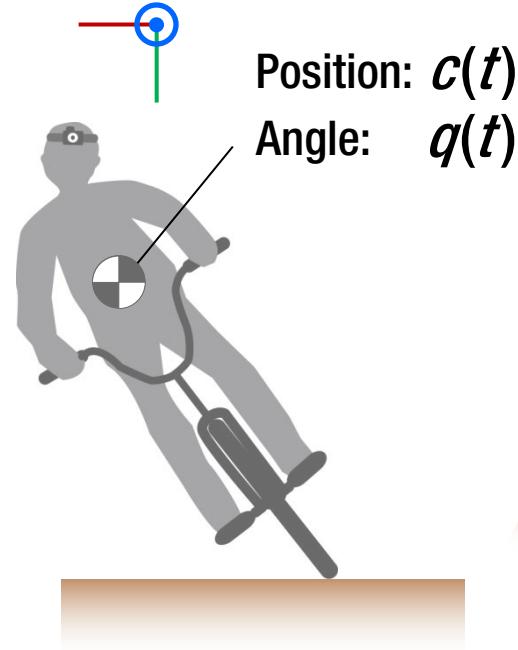
Linear force:

$$m \frac{d^2 c}{dt^2} = F$$



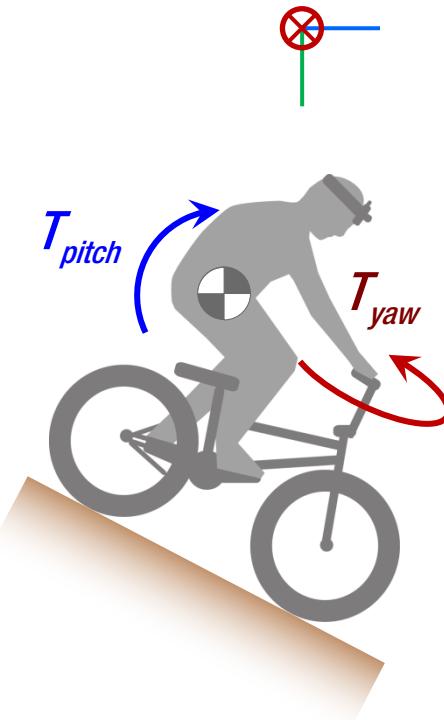
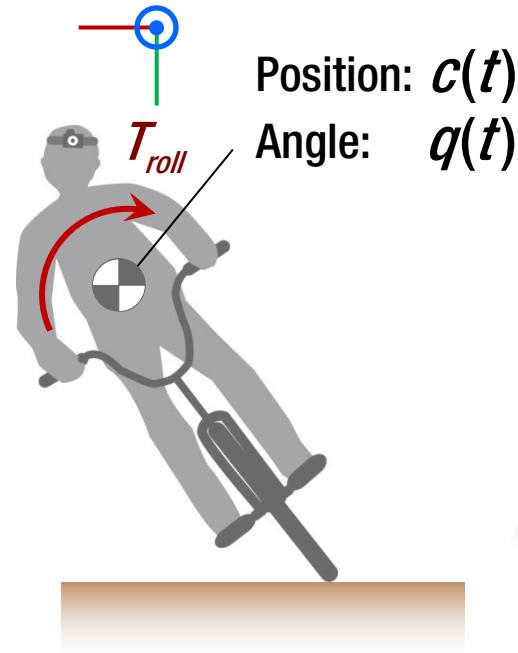
Linear force:

$$m \frac{d^2 c}{dt^2} = F = \sum \mathcal{F}_{\text{passive}} + \sum \mathcal{F}_{\text{active}}$$



Linear force:

$$m \frac{d^2 c}{dt^2} = F = \sum \mathcal{F}_{\text{passive}} + \sum \mathcal{F}_{\text{active}}$$



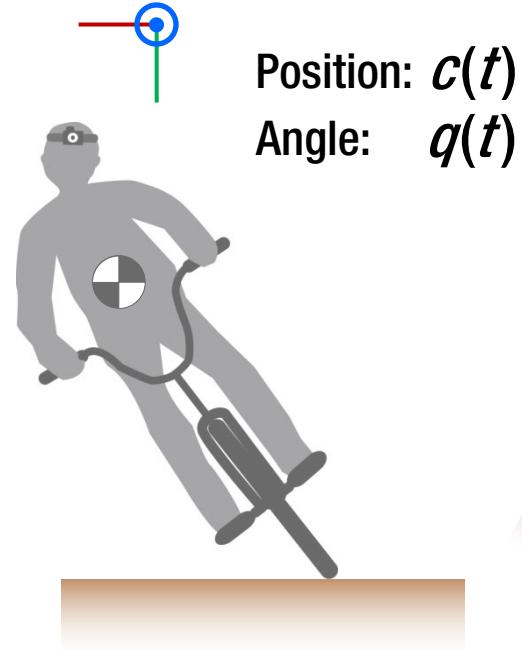
Linear force:

$$m \frac{d^2 c}{dt^2} = F = \sum \mathcal{F}_{\text{passive}} + \sum \mathcal{F}_{\text{active}}$$

Angular force (torque):

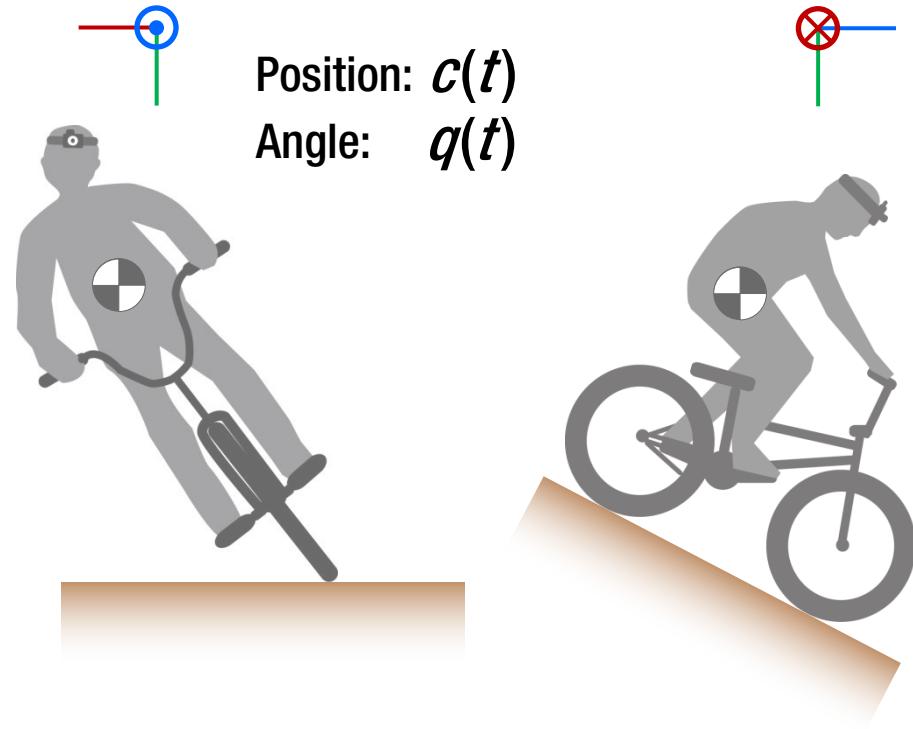
$$J \frac{d^2 q}{dt^2} + \frac{dq}{dt} \times J \frac{dq}{dt} = \sum \mathcal{T}_{\text{passive}} + \sum \mathcal{T}_{\text{active}}$$

J : Moment of inertia



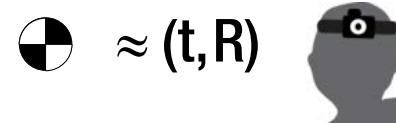
$$(c, q) = \text{ODE}(F_{\text{passive}}, F_{\text{active}}, T_{\text{passive}}, T_{\text{active}})$$



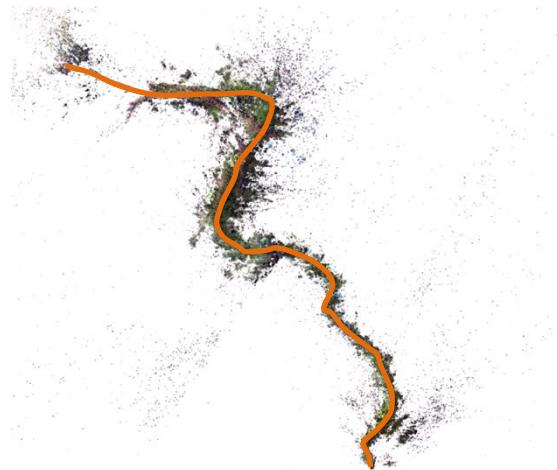


Position: $c(t)$
Angle: $q(t)$

$$(c, q) = \text{ODE}(F_{\text{passive}}, F_{\text{active}}, T_{\text{passive}}, T_{\text{active}})$$



$$\text{where } P = K[R \ t]$$

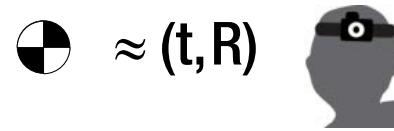


Inverse control:

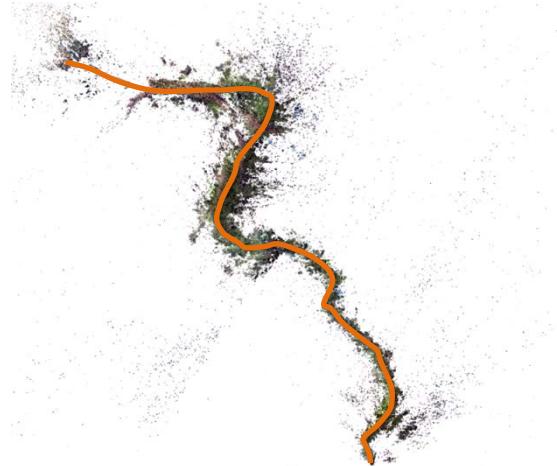
minimize E_{SfM}
 F, T, X Reprojection error

subject to $(t, R) = \text{ODE}(F, T)$

$$(c, q) = \text{ODE}(\mathcal{F}_{\text{passive}}, \mathcal{F}_{\text{active}}, \mathcal{T}_{\text{passive}}, \mathcal{T}_{\text{active}})$$



where $P = K[R \ t]$

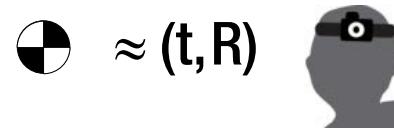


Inverse control:

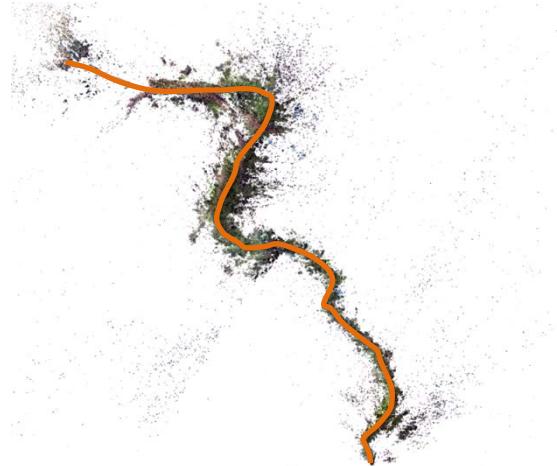
$$\underset{F, T, X}{\text{minimize}} \quad E_{\text{SfM}} + \lambda \frac{E_{\text{reg}}(F, T)}{\text{Temporal regularization}}$$

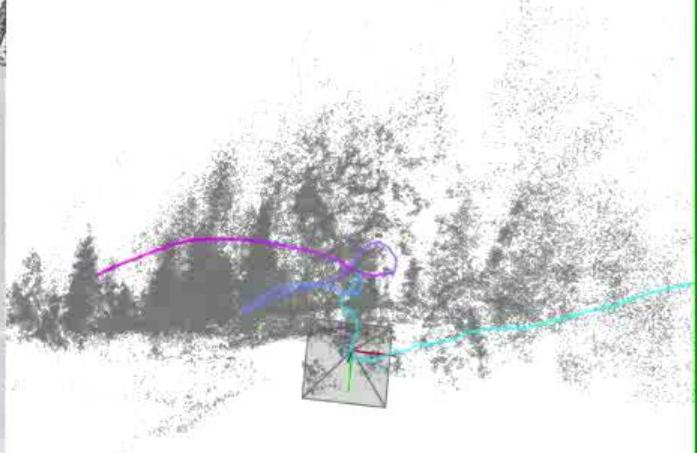
subject to $(t, R) = \text{ODE}(F, T)$

$$(c, q) = \text{ODE}(\textcolor{blue}{F}_{\text{passive}}, \textcolor{red}{F}_{\text{active}}, \textcolor{blue}{T}_{\text{passive}}, \textcolor{blue}{T}_{\text{active}})$$



where $P = K[R \ t]$





3D reconstruction

<https://www.youtube.com/watch?v=pCcuKCIUpLs>

Time: 5.17sec
Speed: 12.6m/s
Air Drag: 73.60N

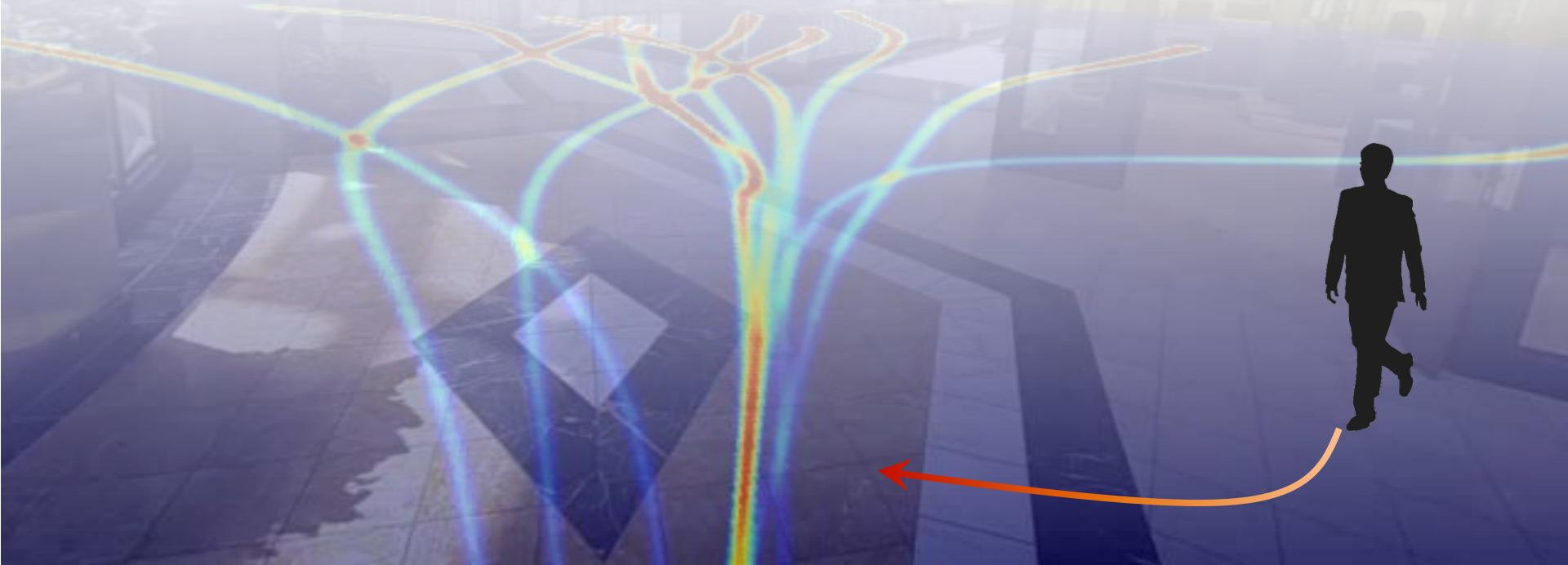


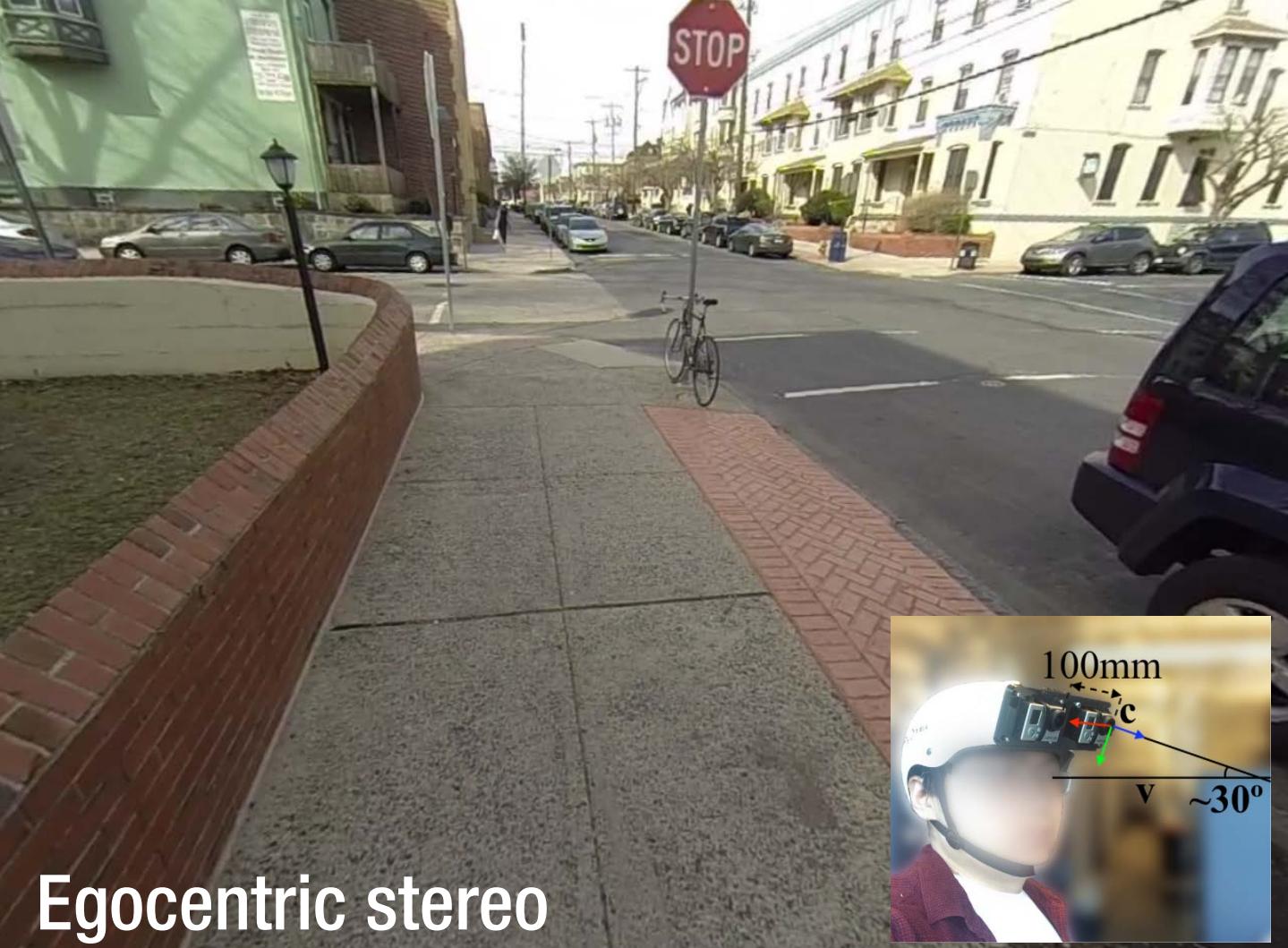
3D reconstruction

<https://www.youtube.com/watch?v=rnvvsjstveM>

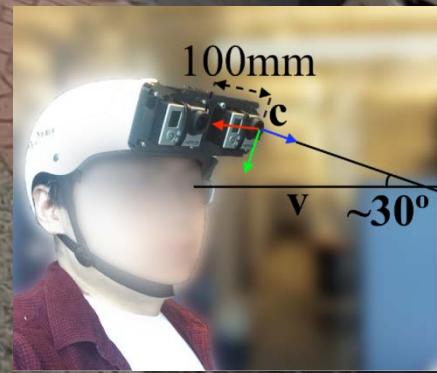
Visual Sensorimotor Behaviors II:

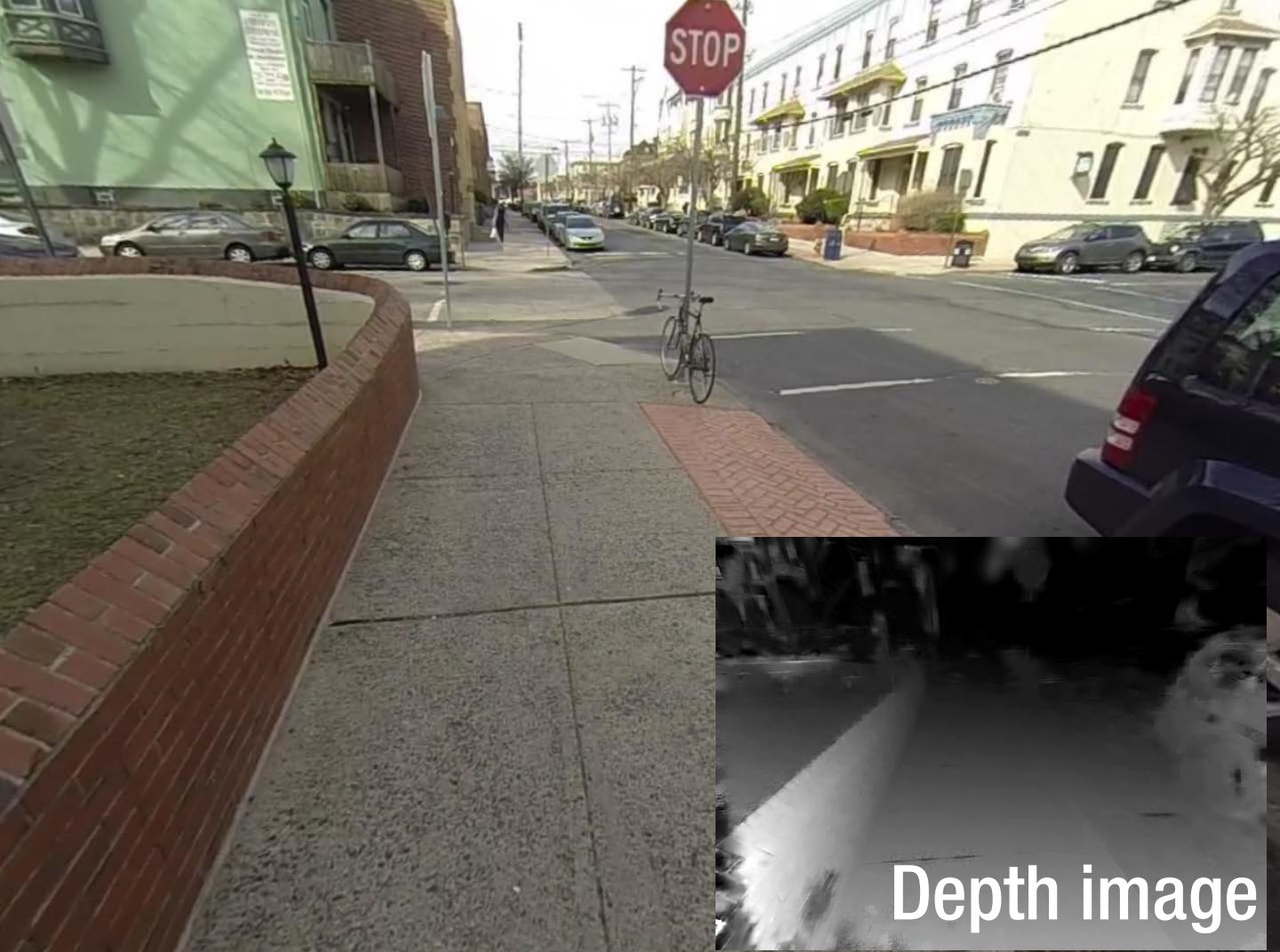
What can a first person video tell about my future?



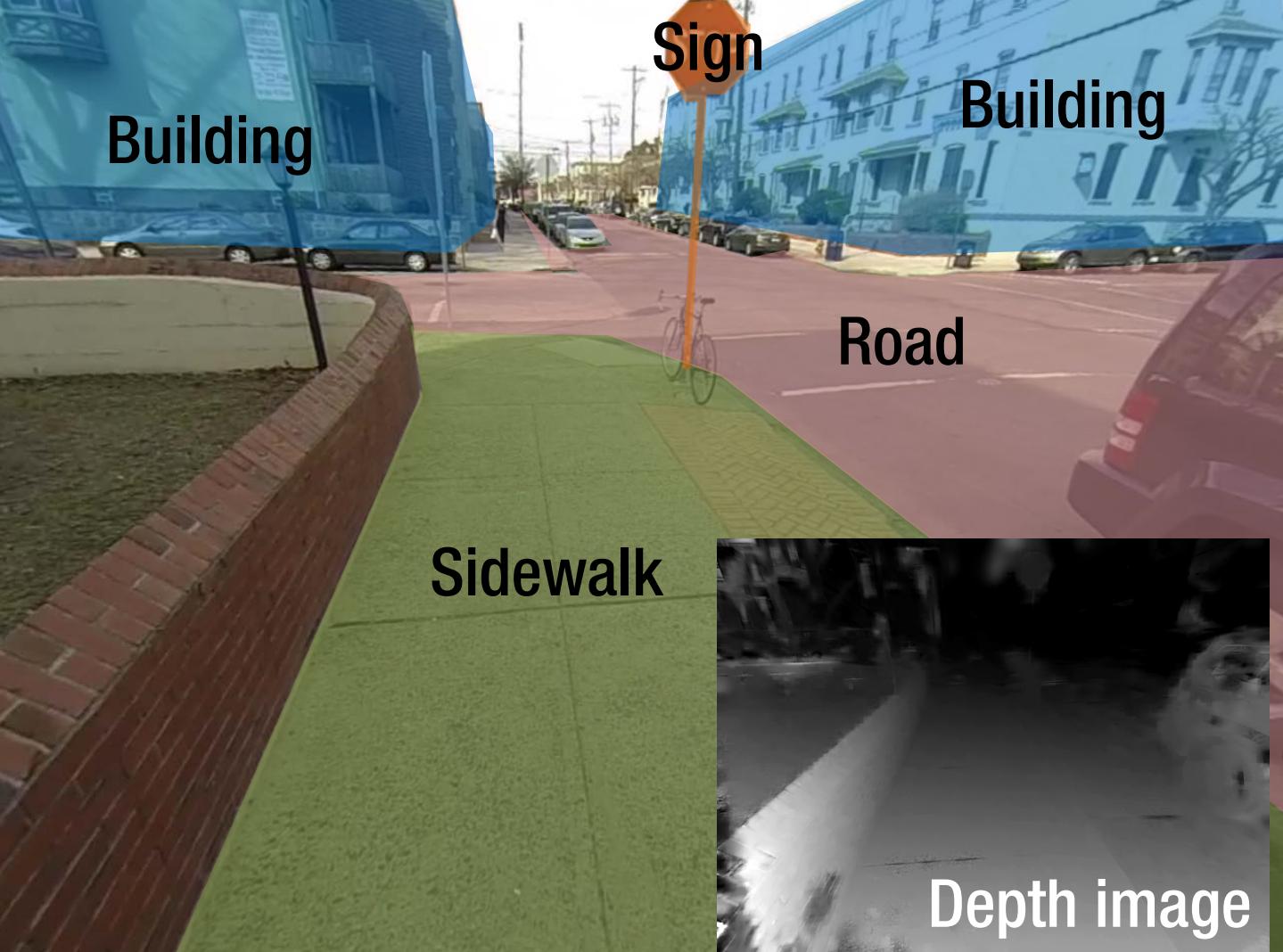


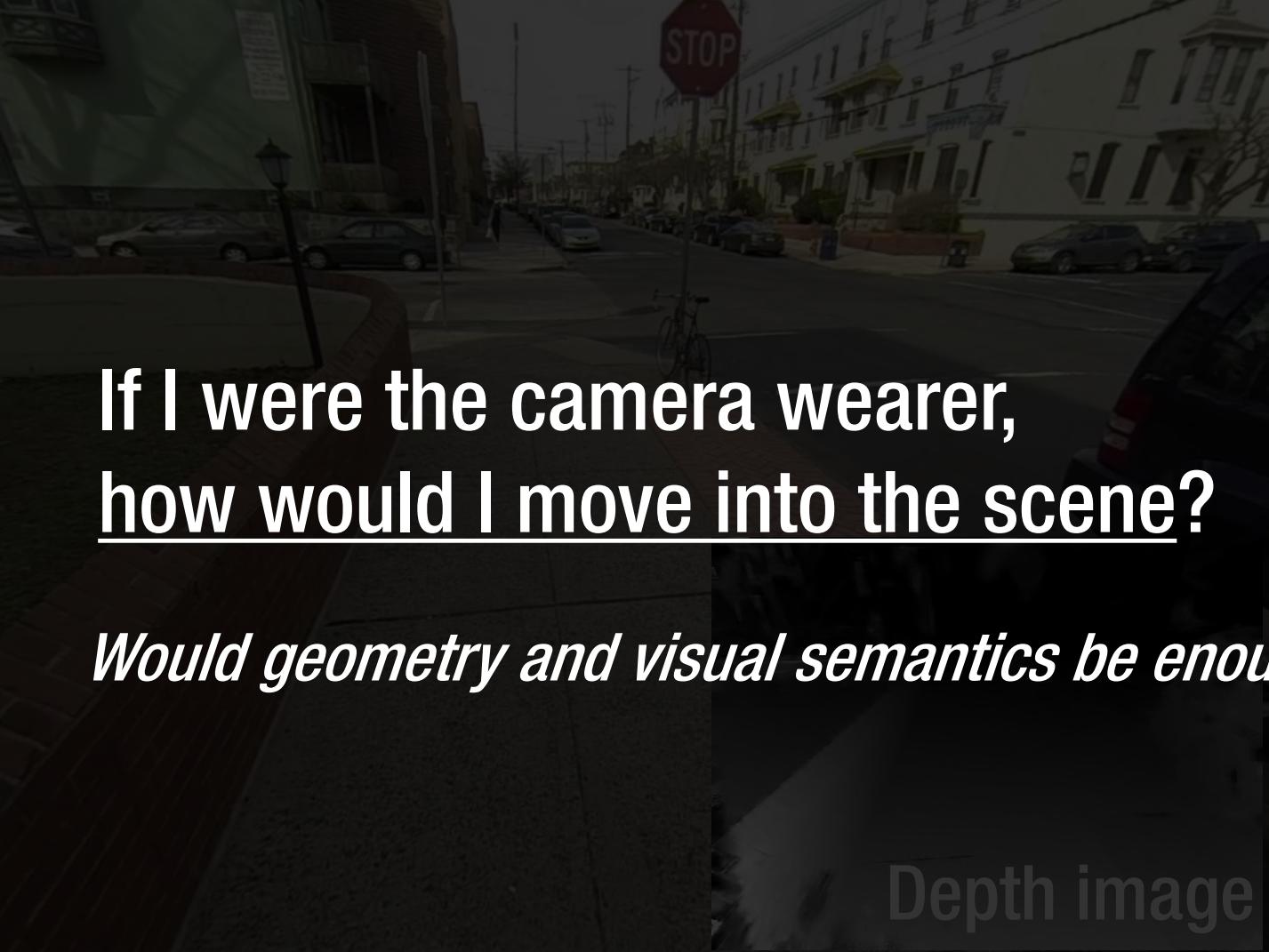
Egocentric stereo





Depth image

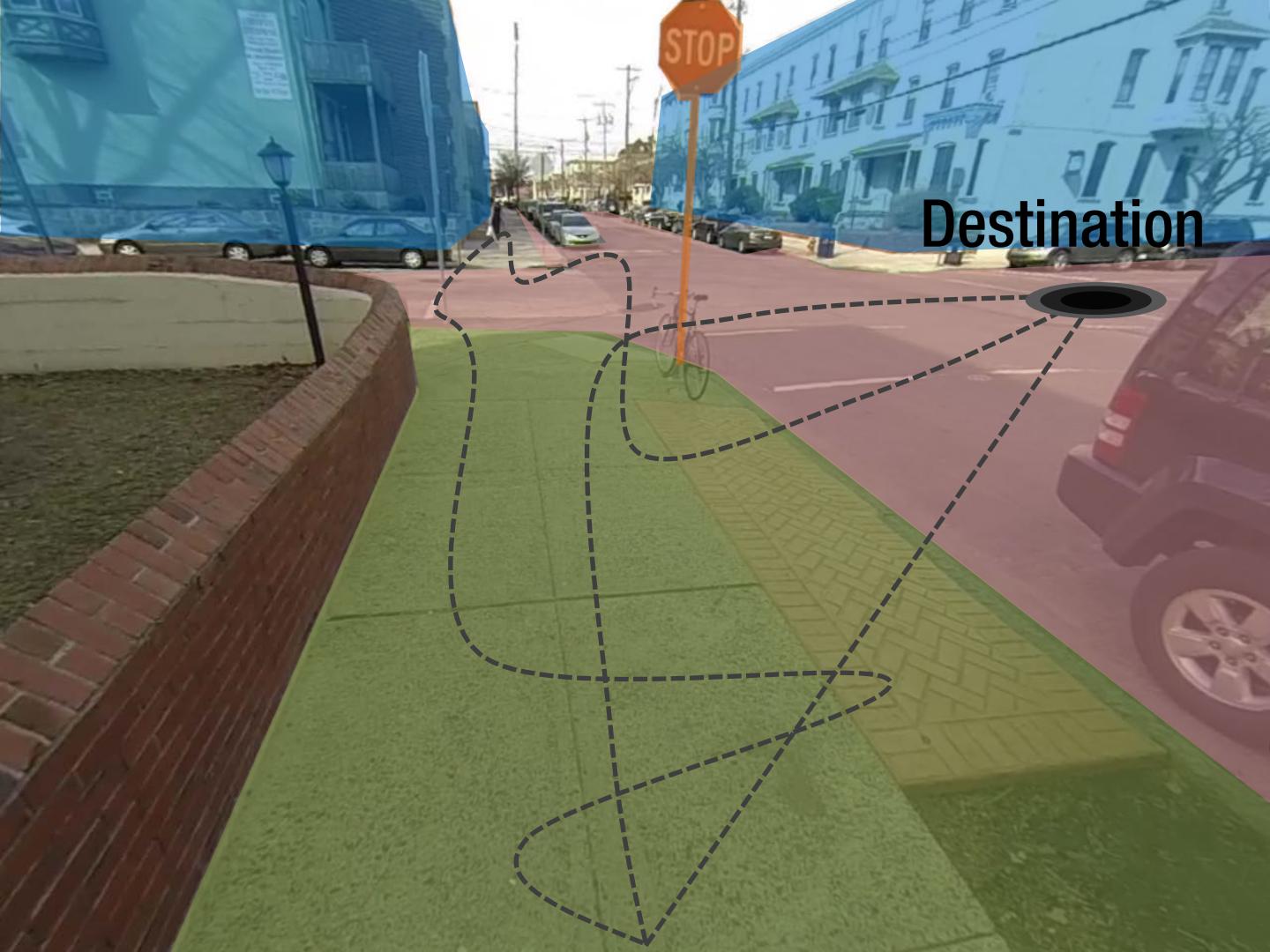




If I were the camera wearer,
how would I move into the scene?

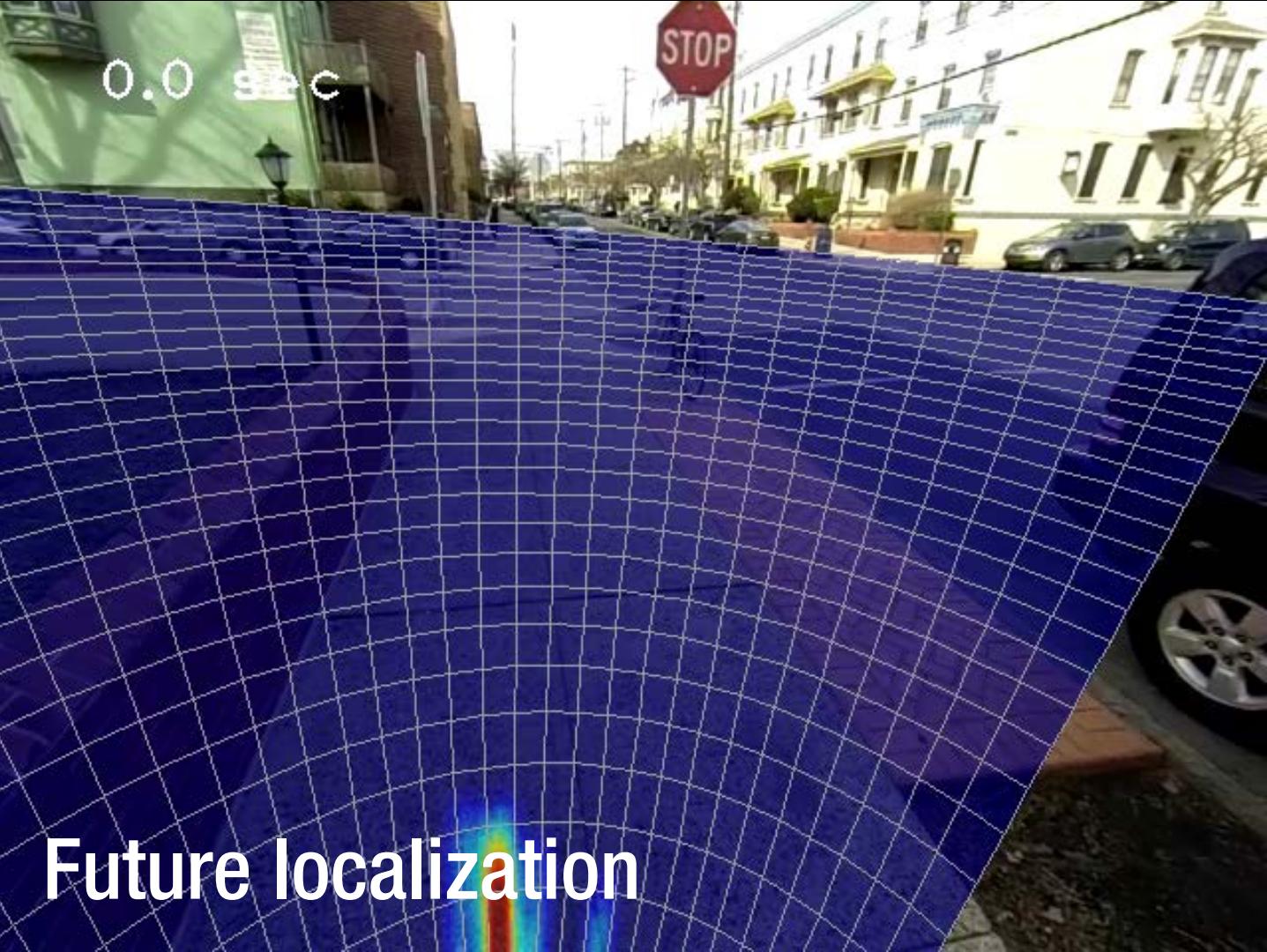
Would geometry and visual semantics be enough?

Depth image

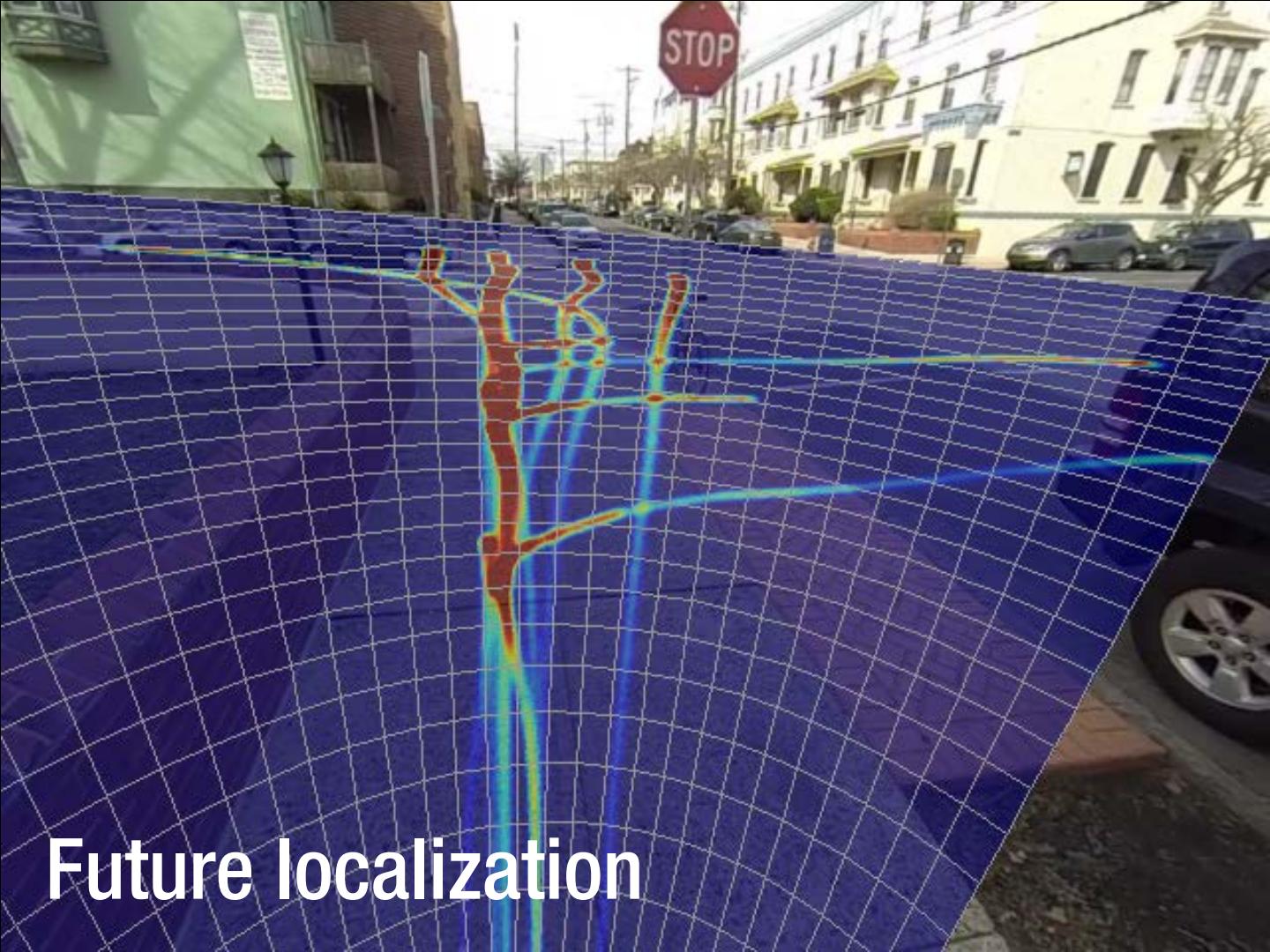


Destination

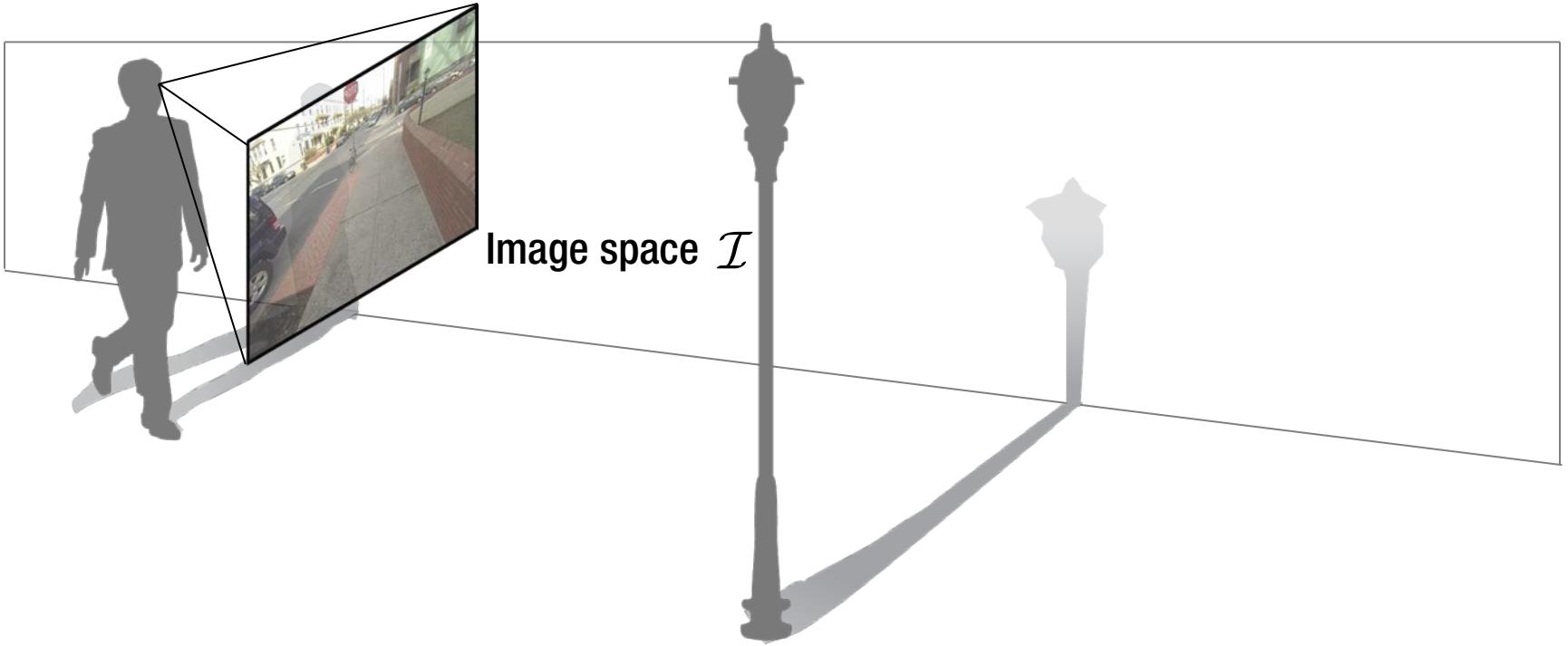
0.0 sec



Future localization



Future localization



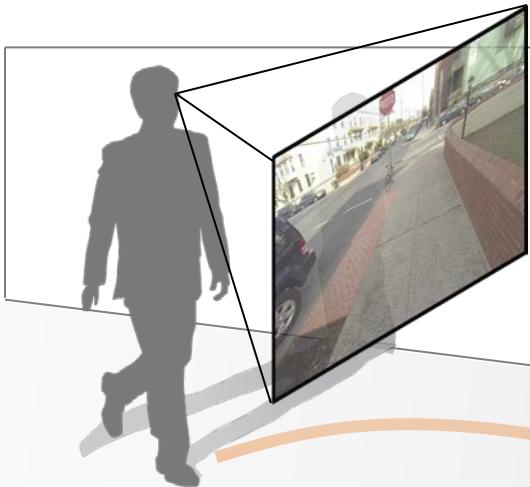
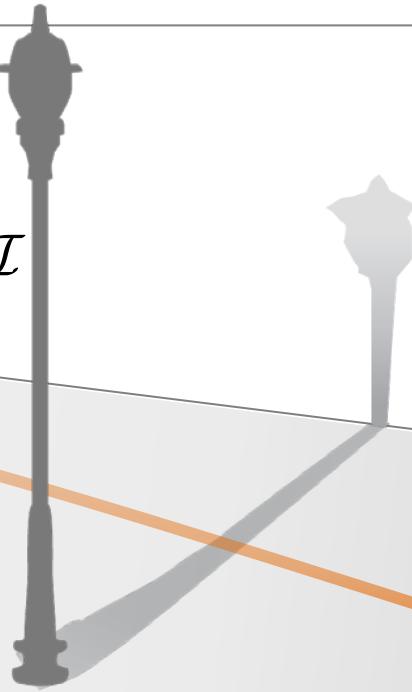


Image space \mathcal{I}



Prediction:
Configuration space (ground plane)



Ground plane

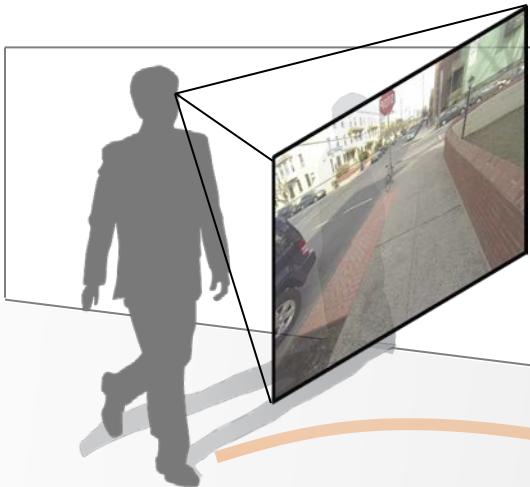


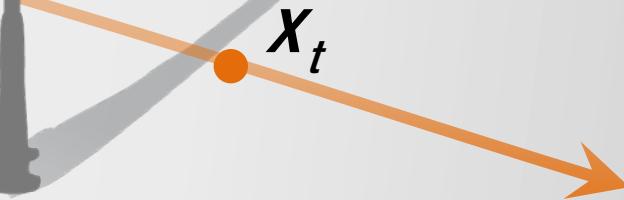
Image space \mathcal{I}



(x_1, \dots, x_F)
Predicted trajectory



Ground plane



Prediction:
Configuration space (ground plane)

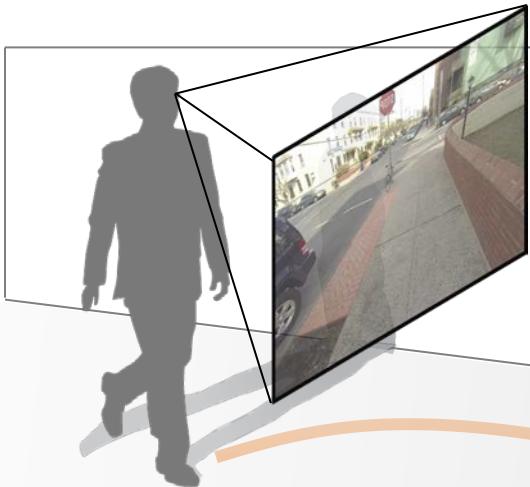
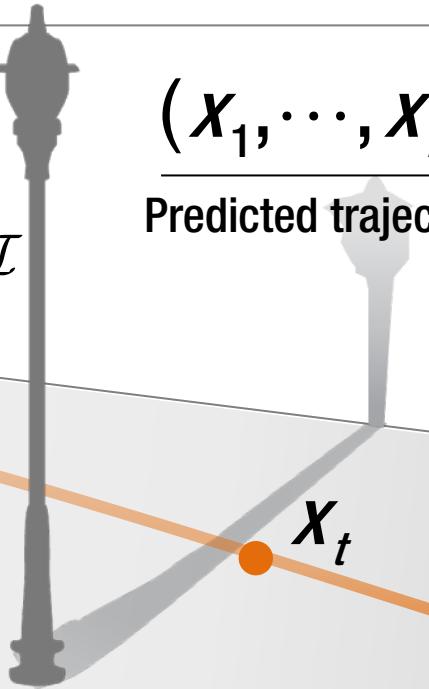


Image space \mathcal{I}



$$(x_1, \dots, x_F) = g(\mathcal{I})$$

Predicted trajectory



Ground plane

Prediction:
Configuration space (ground plane)

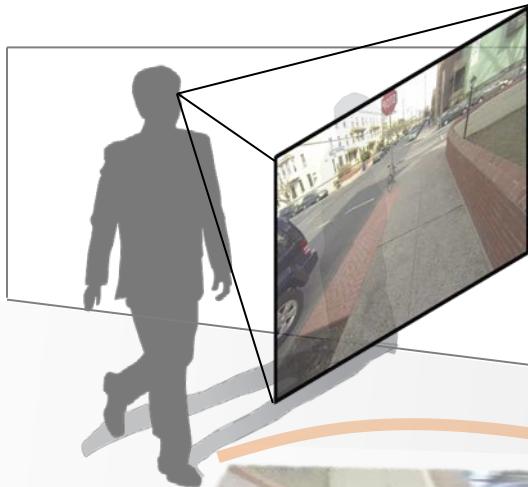


Image space \mathcal{I}

$$(x_1, \dots, x_F) = \underline{g(f(\mathcal{I}))}$$

Projection to cfg. space



x_t



Ground plane

Prediction:
Configuration space (ground plane)

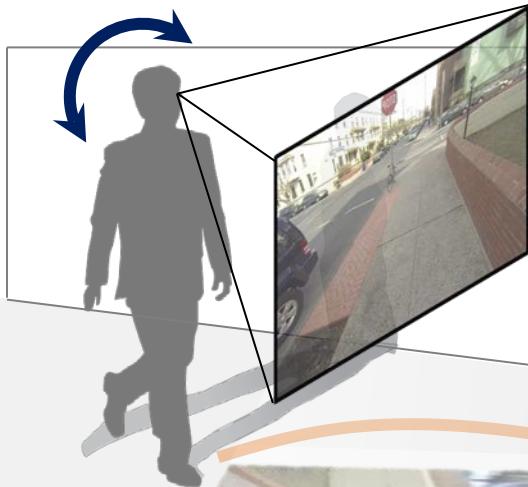
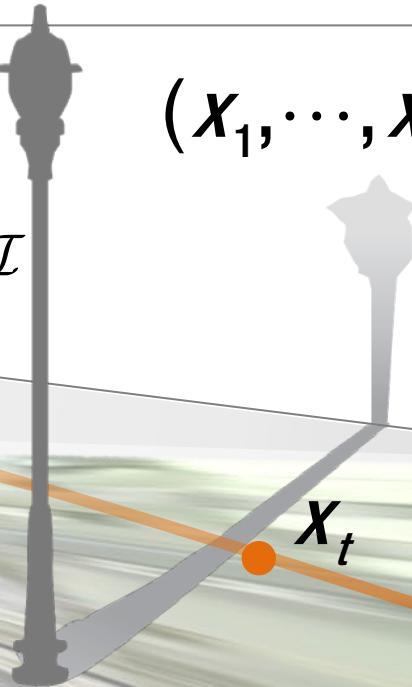


Image space \mathcal{I}

$$(x_1, \dots, x_F) = \underline{g(f(\mathcal{I}))}$$

Projection to cfg. space



Prediction:
Configuration space (ground plane)
Pitch angle invariant



Ground plane

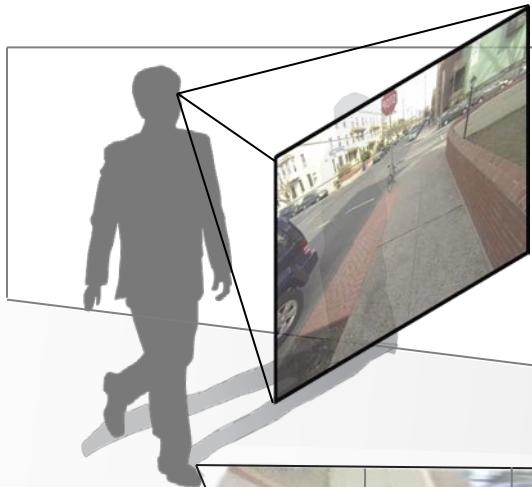


Image space \mathcal{I}

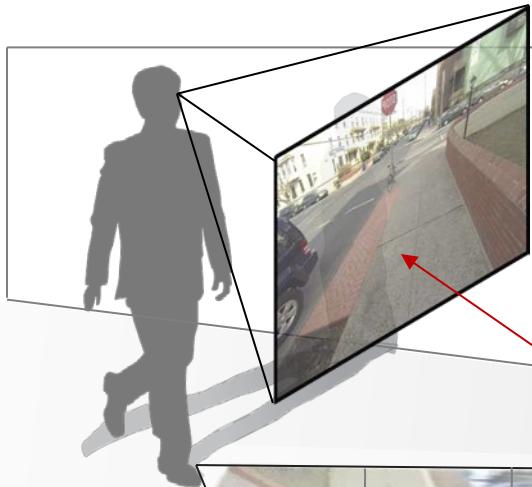
$$(x_1, \dots, x_F) = g(f(\mathcal{I}))$$

$$f(\mathcal{I})_{r,\theta}$$



Ground plane

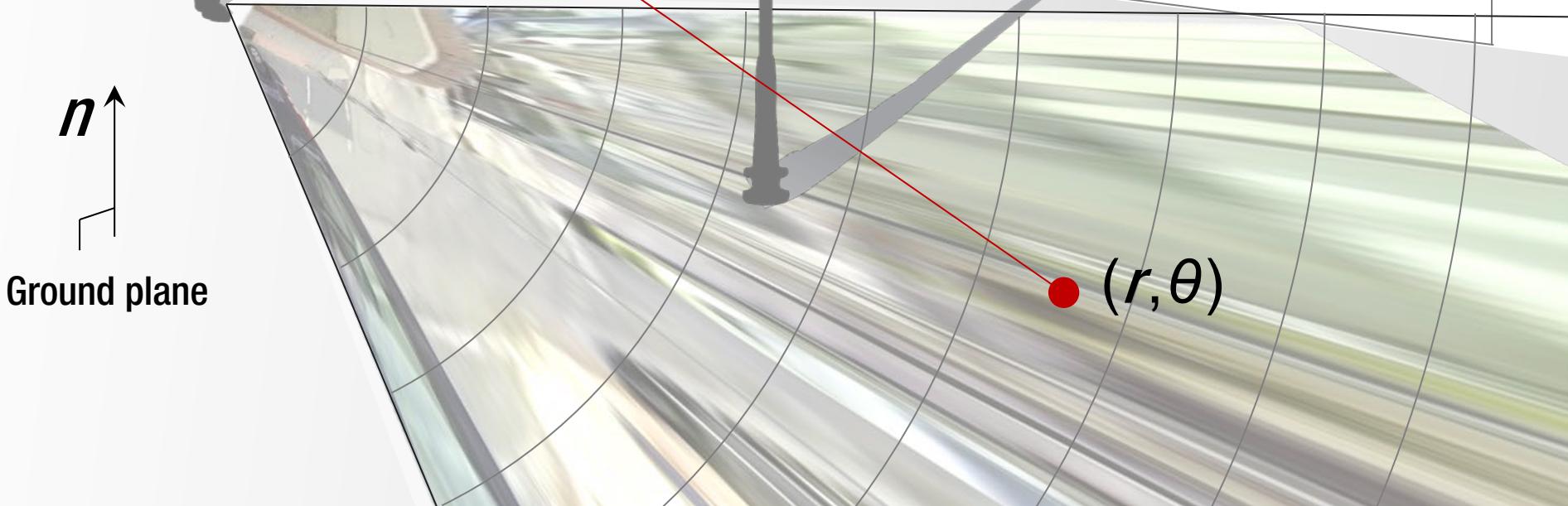
$$\bullet (r, \theta)$$



$$(x_1, \dots, x_F) = g(f(\mathcal{I}))$$

$$f(\mathcal{I})_{r,\theta} = \underline{\mathcal{I}_{\text{proj}(r,\theta)}}$$

Image projection



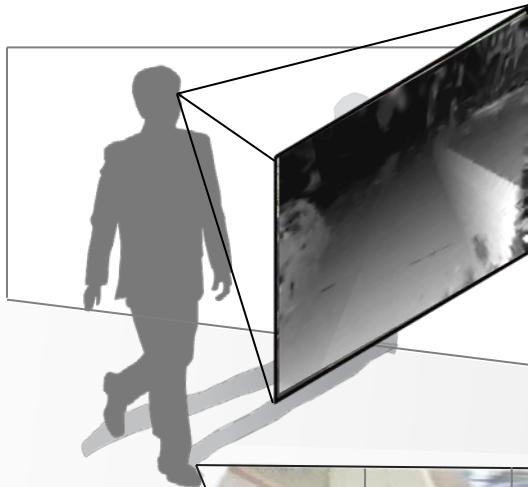
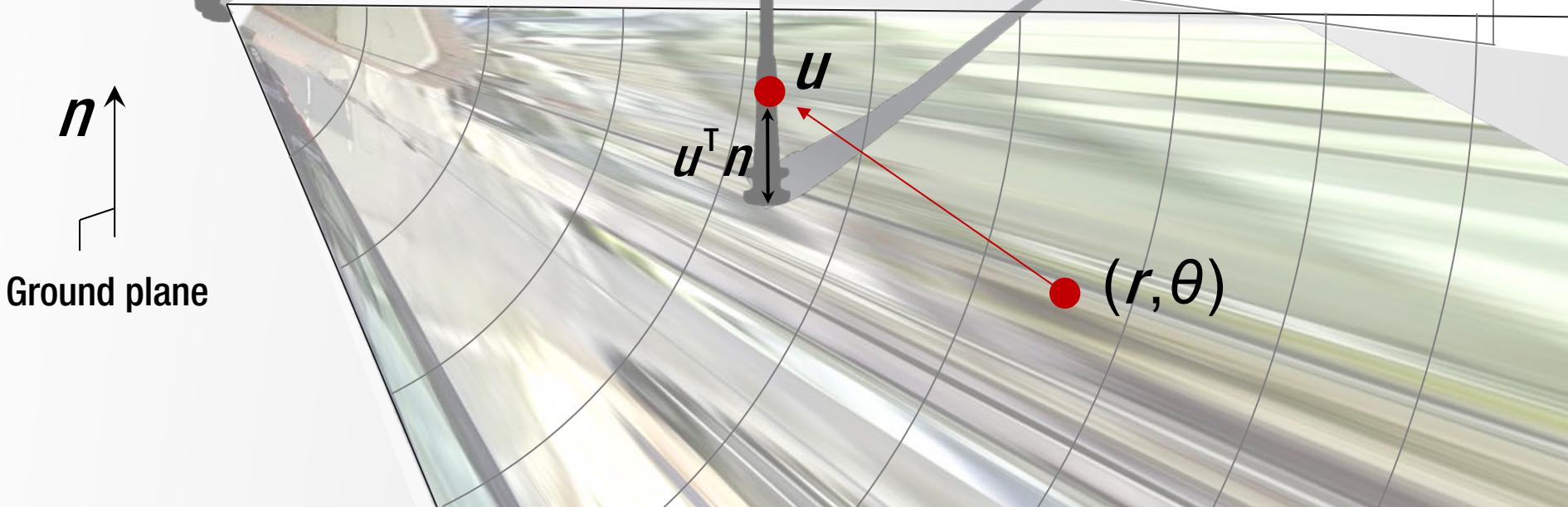


Image space \mathcal{I}

$$(x_1, \dots, x_F) = g(f(\mathcal{I}))$$

$$f(\mathcal{I})_{r,\theta} = (\mathcal{I}_{\text{proj}(r,\theta)}, \underline{u^T n})$$

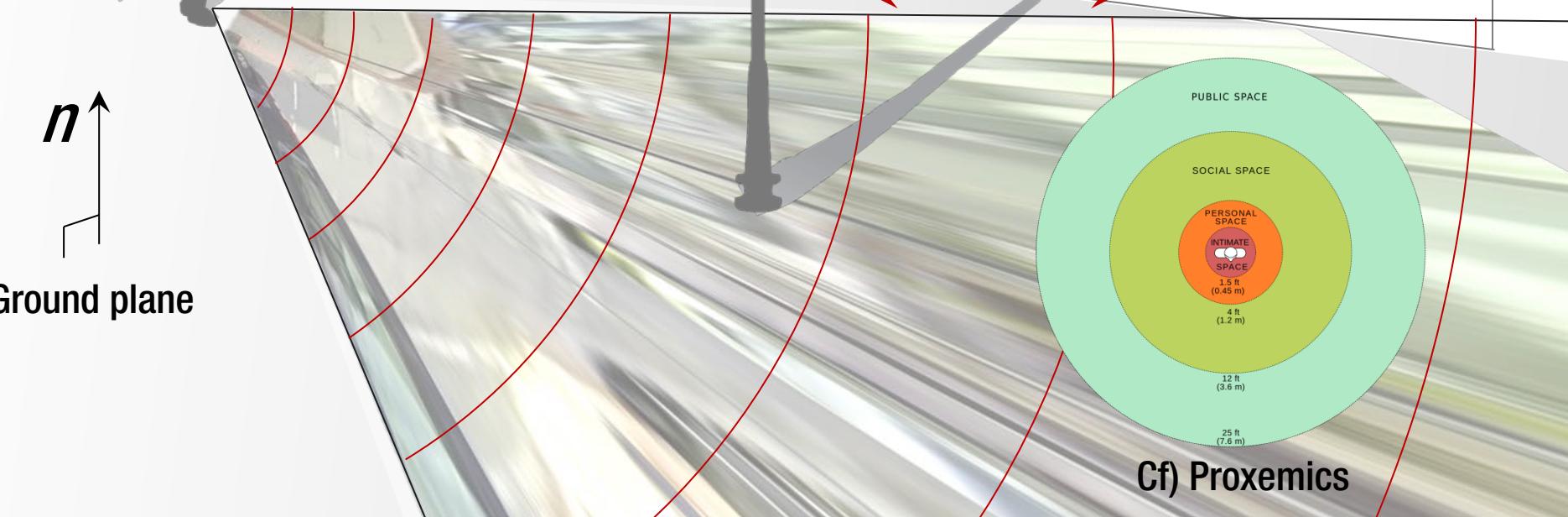
Height of occluding object

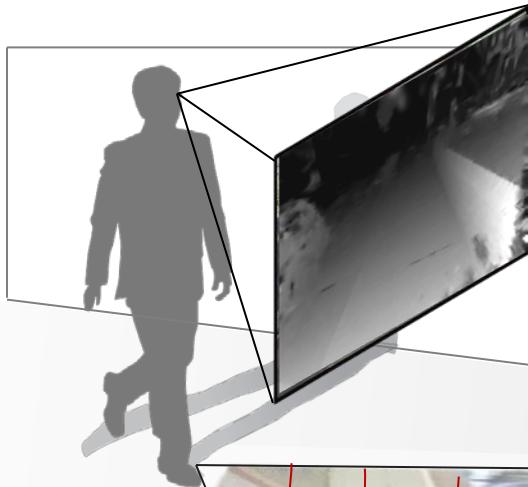


Ground plane

Retinal representation

$$\Delta r \propto \log \frac{1}{D} \text{ where } D \text{ is depth.}$$

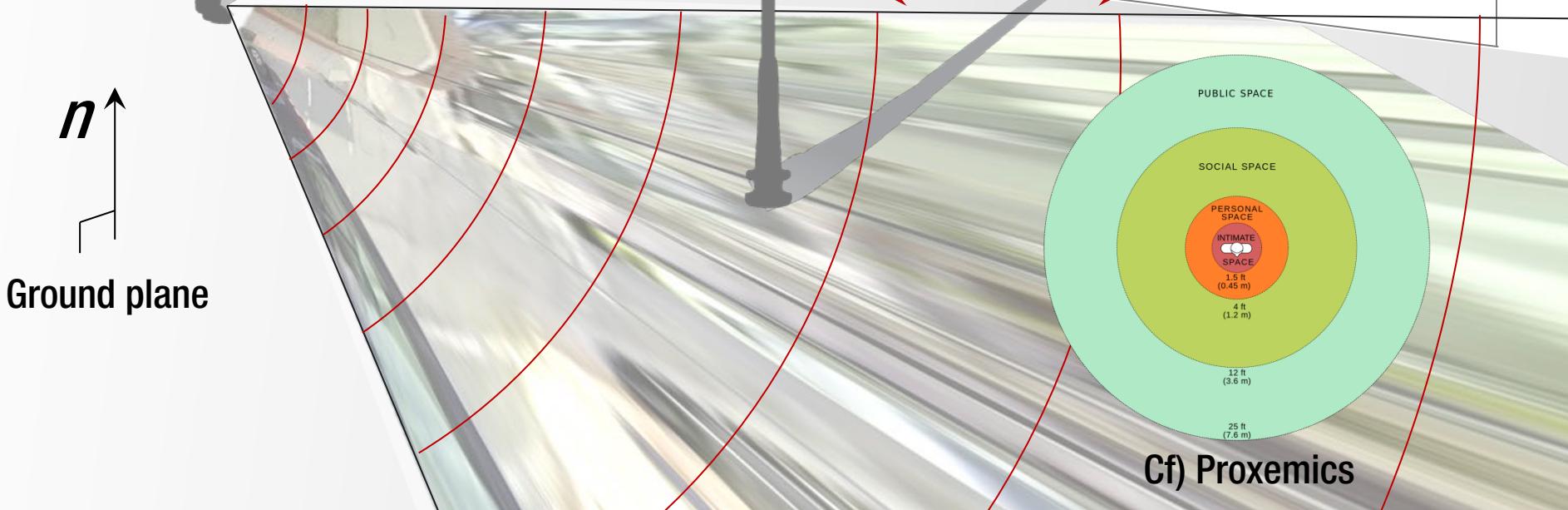


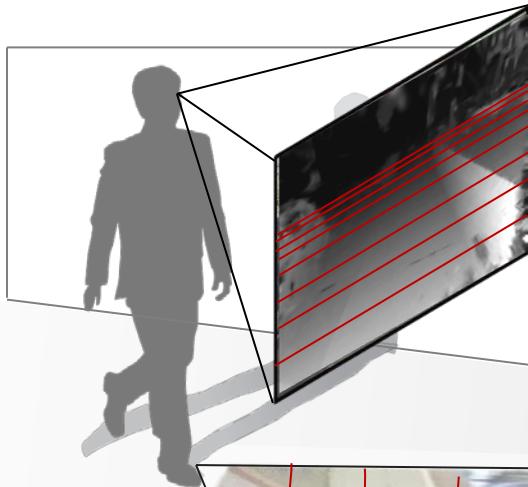


Retinal representation

Persistent to 2D and 3D distance

$$\Delta r \propto \log \frac{1}{D} \text{ where } D \text{ is depth.}$$



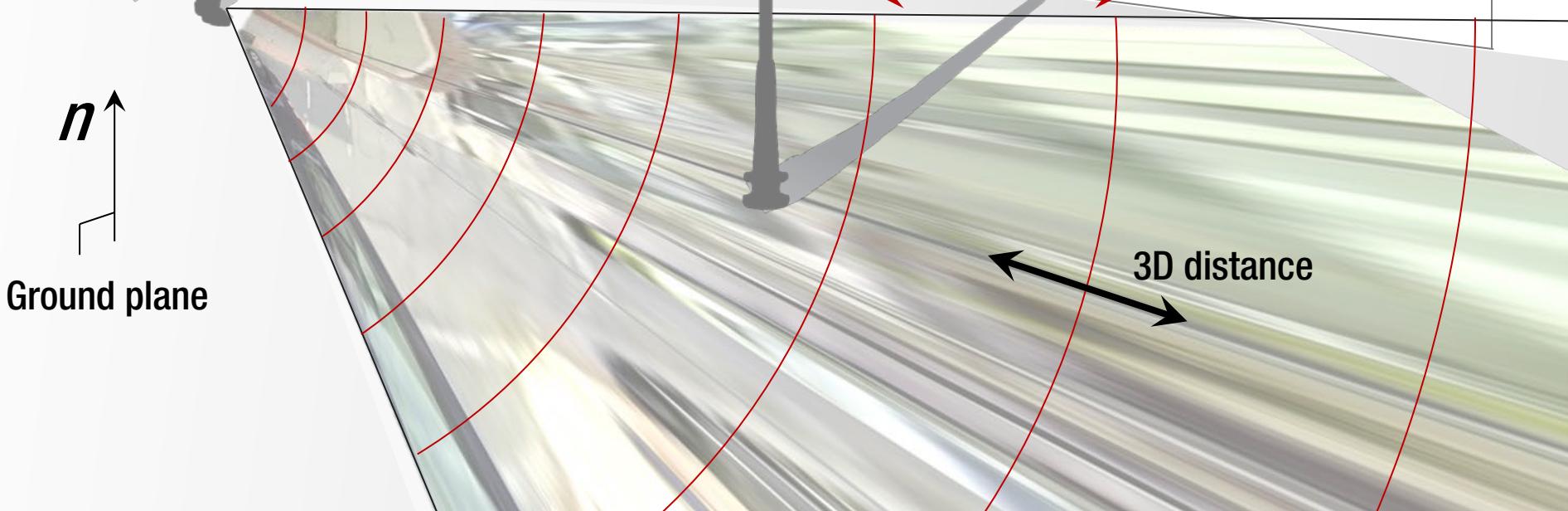


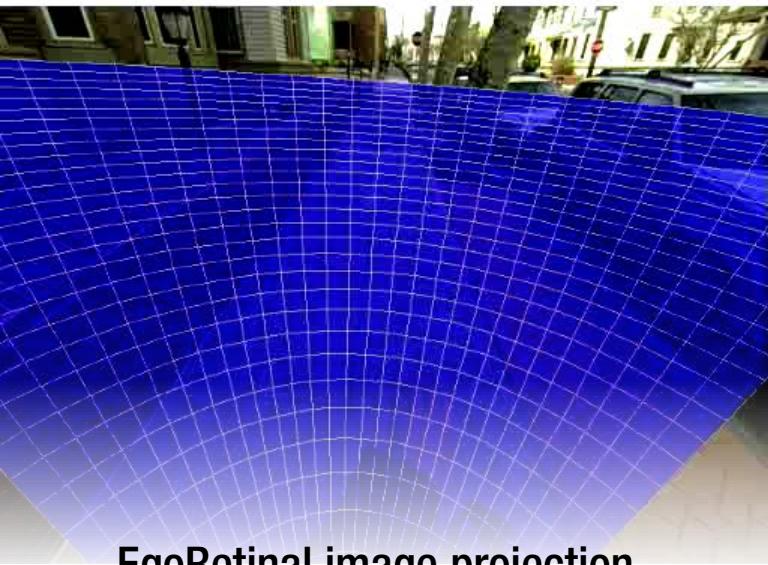
$\log p$

Retinal representation

Persistent to 2D and 3D distance

$$\Delta r \propto \log \frac{1}{D} \text{ where } D \text{ is depth.}$$

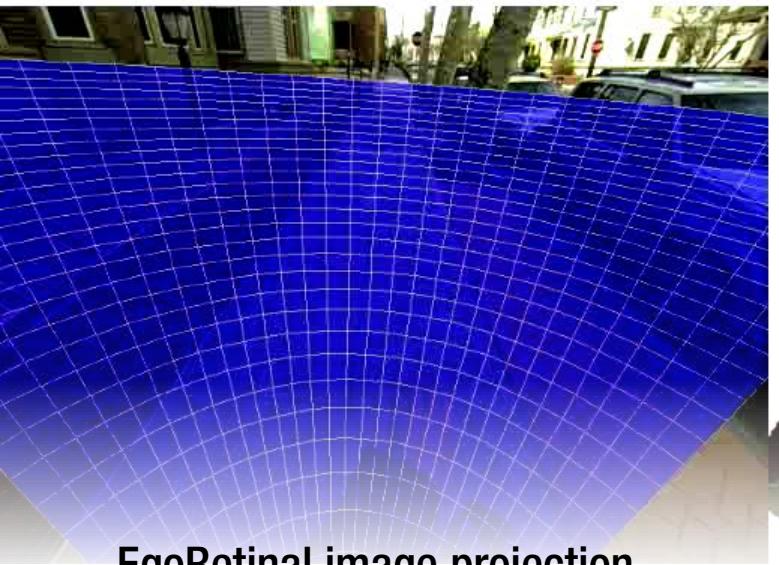




EgoRetinal image projection



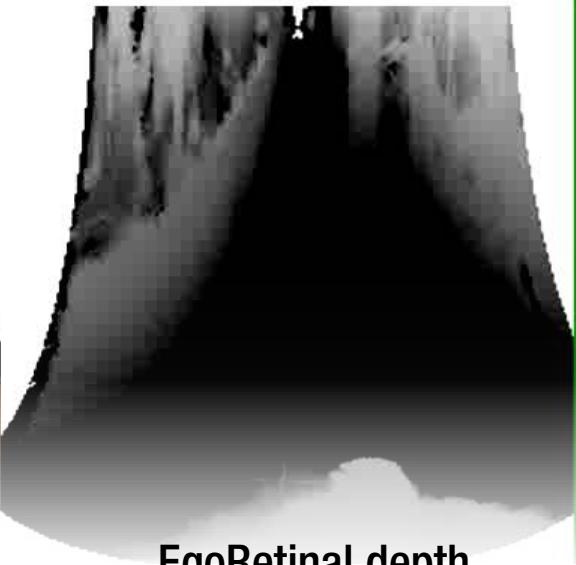
EgoRetinal RGB



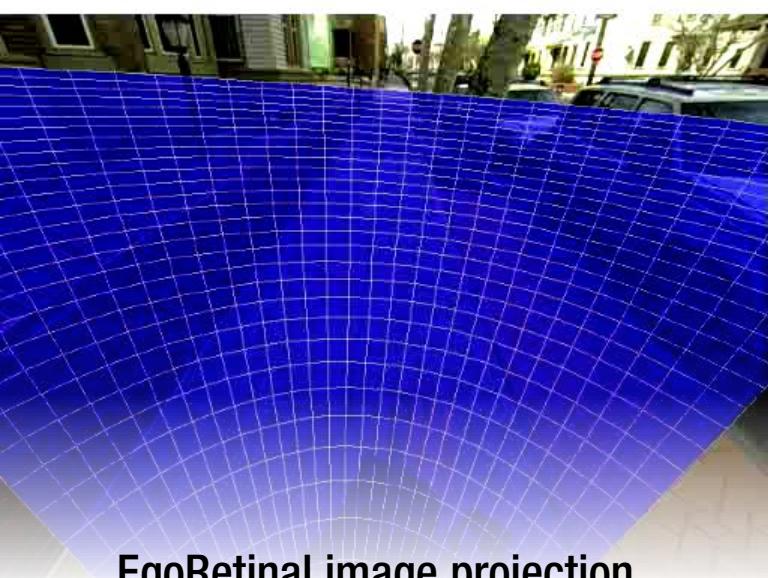
EgoRetinal image projection



EgoRetinal RGB



EgoRetinal depth



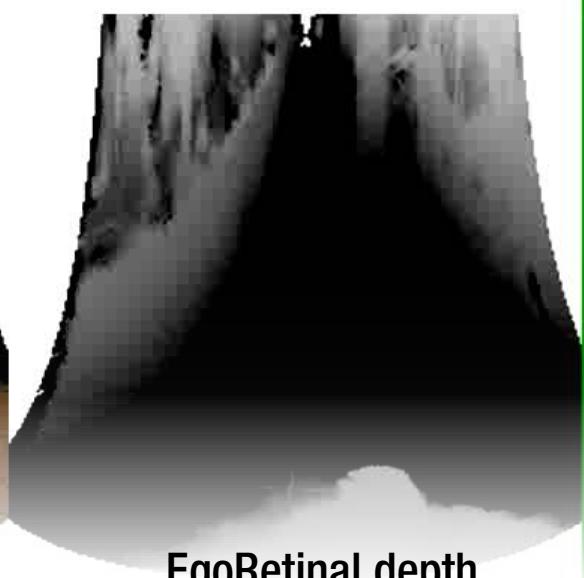
EgoRetinal image projection



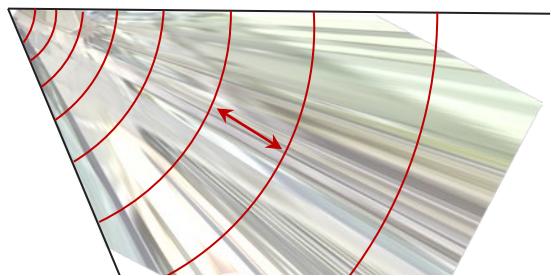
P1: Pitch angle invariant



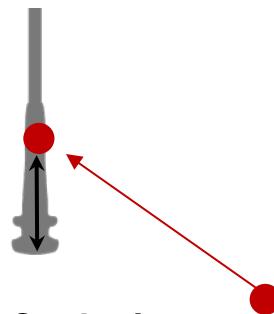
EgoRetinal RGB



EgoRetinal depth



P2: 2D and 3D persistent



P3: Occlusion reasoning

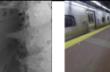
EgoMotion Dataset (outdoor)



EgoMotion

out
in

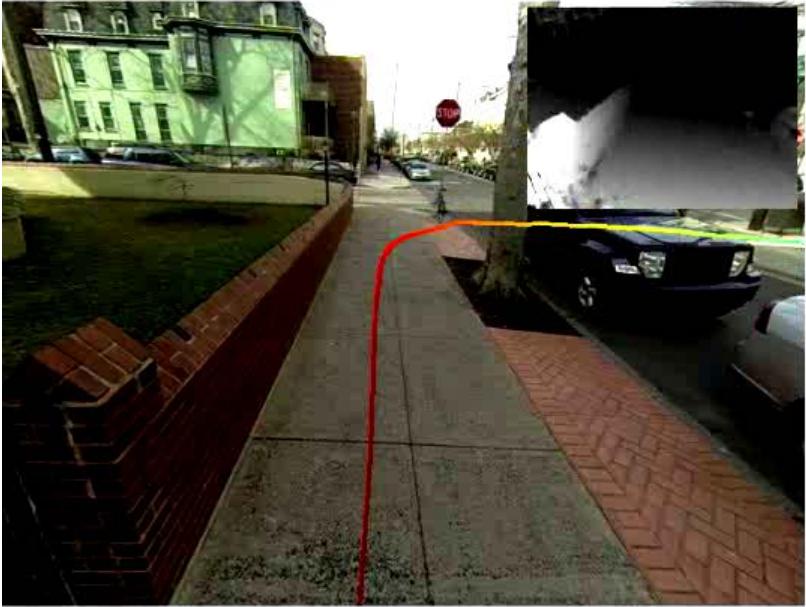
left
right

Image Disparity								
Scene	IKEA	Costco	Mall	Park	School1/2	Downtown1/2	Grocery1/2/3	Bus1/2
Frames	966	577	2683	3088	3754/3736	2856/3405	2858/2892/2834	2292/1850
Duration	08:03	04:49	22:22	25:44	31:17/31:08	23:48/28:23	23:49/24:06/23:37	19:06/15:25
Image Disparity								
Scene	Campus1/2/3	CVS1/2	Train Sta.1/2	River1/2	Dep. store	Library	Apartment	Caffe
Frames	2607/1884/1975	2359/3337	4034/2568	3378/2250	2250	1255	2050	1550
Duration	21:44/15:42/16:28	19:40/27:49	33:37/21:24	28:09/18:45	13:20	10:30	17:05	13:00

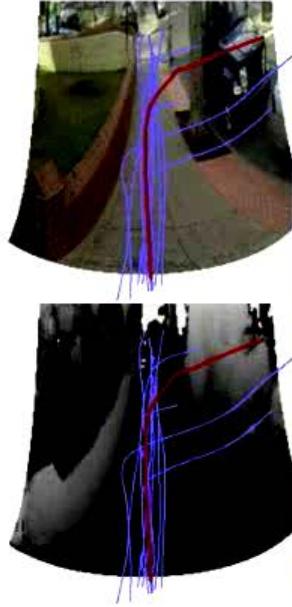
Dataset summary

- **1280x960 stereo (100mm baseline, ~15m depth resolution)**
- **26 scenes (13 indoor, 13 outdoor)**
- **65.5k frames (9.1 hours)**

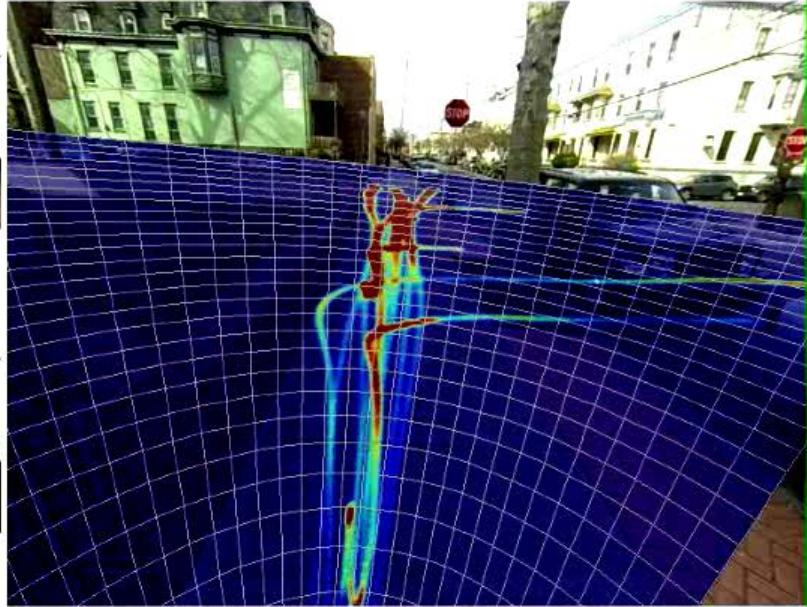
The trajectory is projected onto the ground plane and its time is color coded.



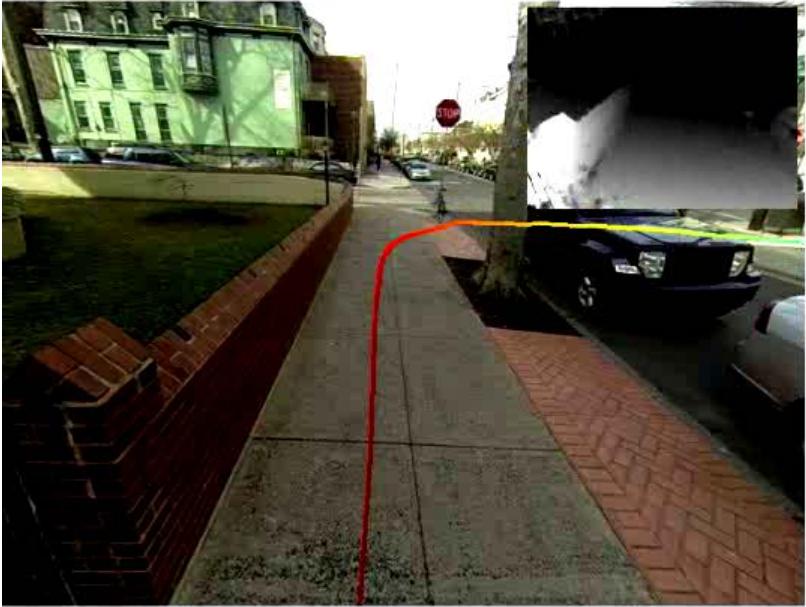
Ground truth



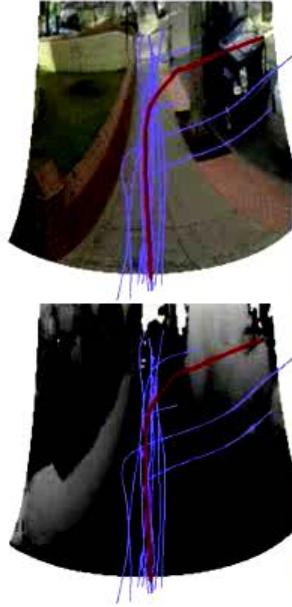
EgoRetinal map



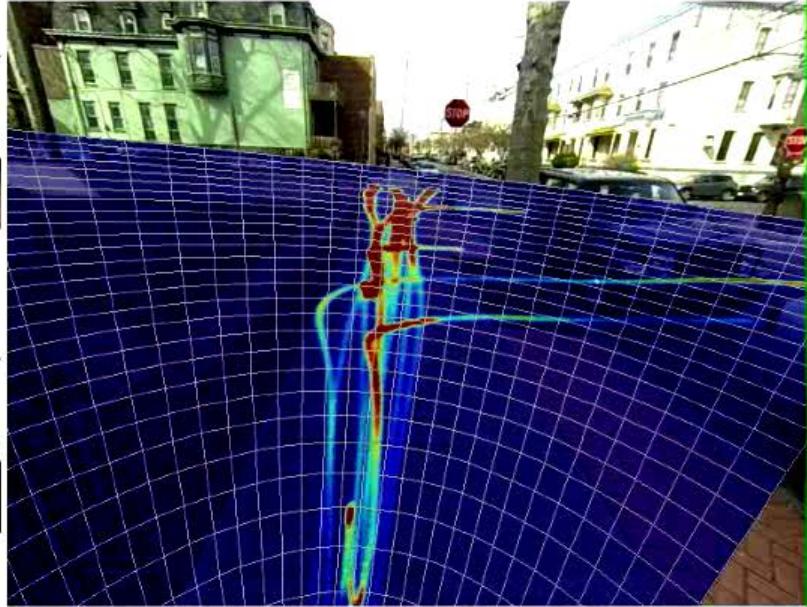
Predicted trajectories



Ground truth



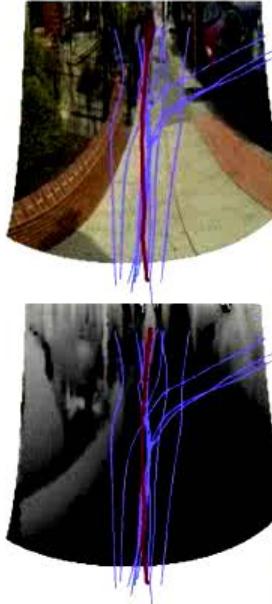
EgoRetinal map



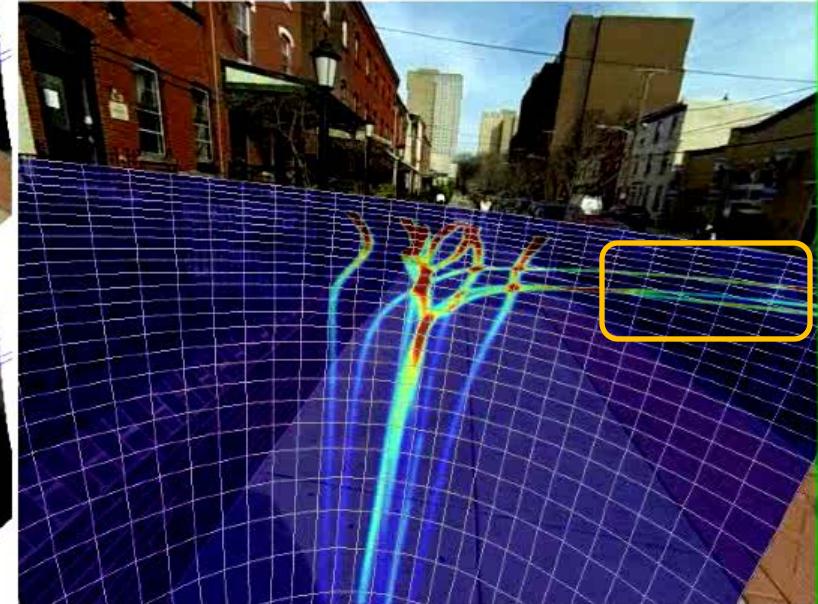
Predicted trajectories



Ground truth



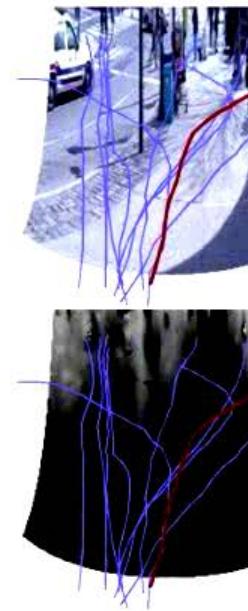
EgoRetinal map



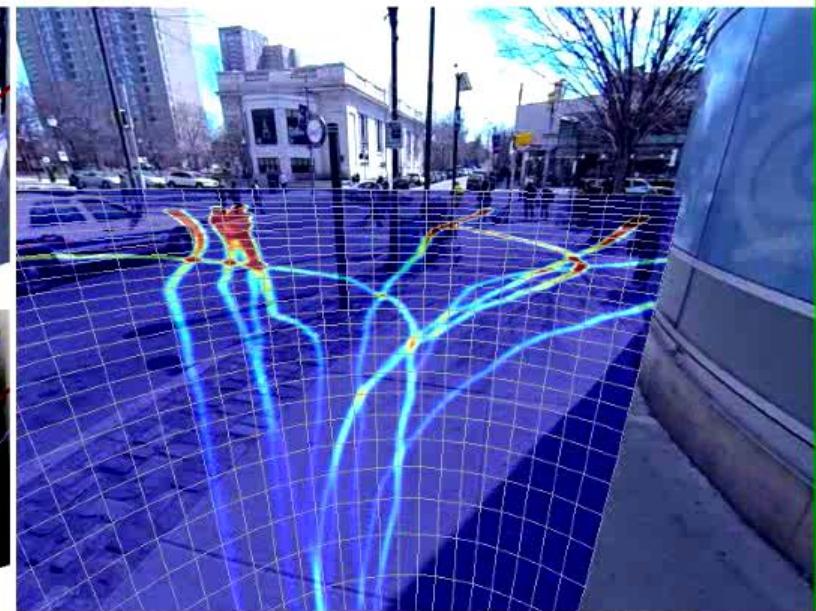
Predicted trajectories



Ground truth



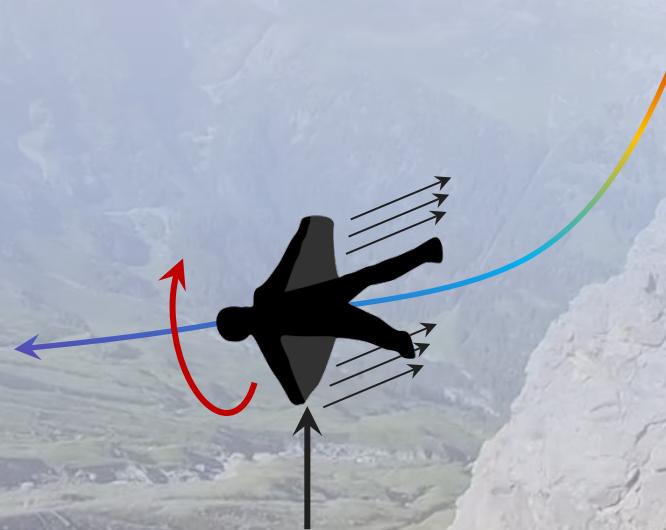
EgoRetinal map



Predicted trajectories

Visual Sensorimotor Behaviors II: Control and Planning

Force from Motion
(ORAL) Afternoon, Wed, June 29



Egocentric Future Localization
(ORAL) Morning, Thr, June 30

