

Event-based Robot Vision

Prof. Dr. Guillermo Gallego
Chair: Robotic Interactive Perception

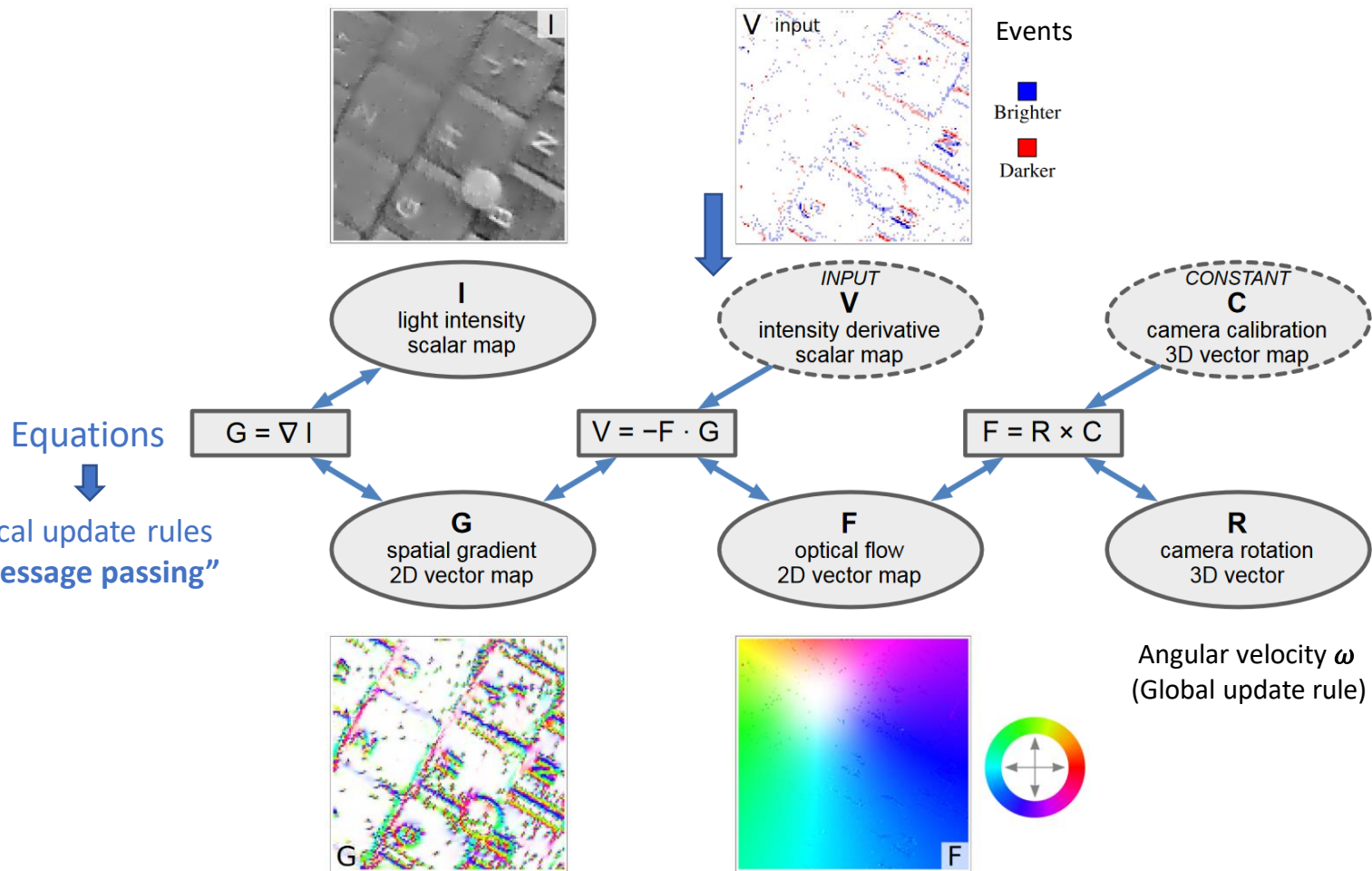
guillermo.gallego@tu-berlin.de

<http://www.guillermogallego.es>

Image (intensity)
Reconstruction
Literature Review

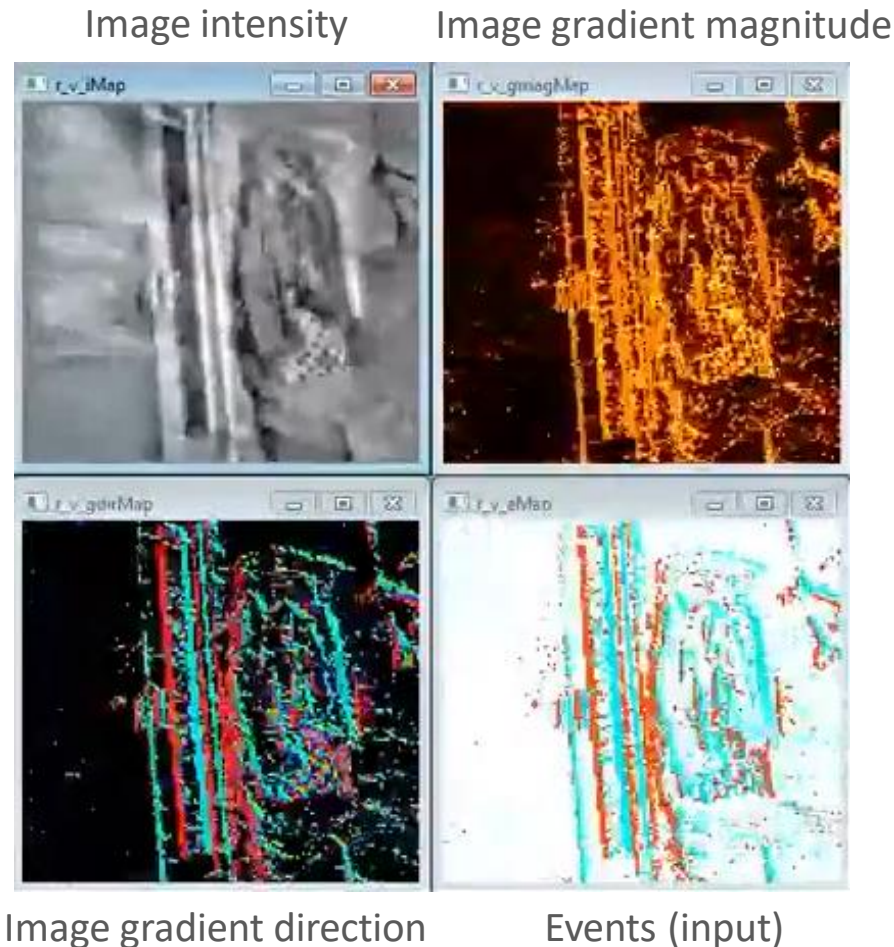
“Interacting Maps”

- Simultaneous estimation of multiple visual quantities with a *rotating* event camera



“Interacting Maps”

- Simultaneous estimation of multiple visual quantities with a *rotating* event camera



Inspired by the primary visual cortex.

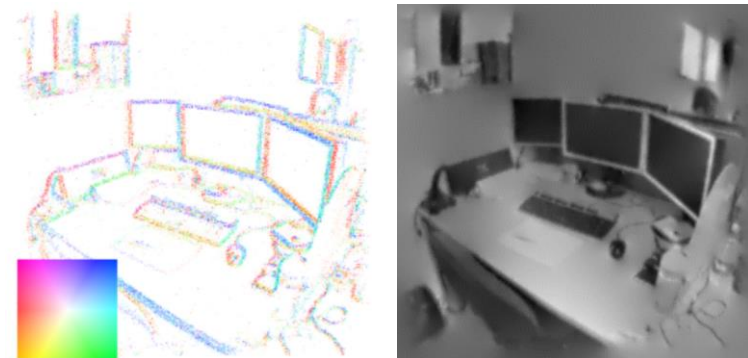
From the events, jointly estimate:

- Rotation (3 DOF ego-motion)
- Optical flow
- Image gradient
- Intensity reconstruction ←

Pure rotation: no translation or depth

Simultaneous mosaicing and tracking

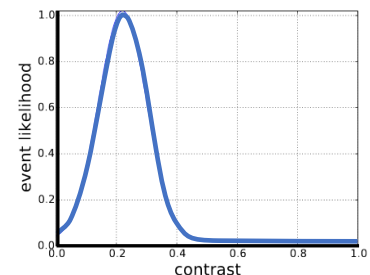
- Parallel tracking and mapping
- *Rotating* event camera (no translation or depth)



Intensity gradient

Image intensity

- **Mapping** = mosaicing (panoramic imaging)
 - Get intensity gradient g using pixel-wise EKF:
$$h(e|R) = \frac{g \cdot v}{c}$$
 should equal the pixel event rate $\frac{1}{\Delta\tau}$
(uses linearized event generation model)
 - Poisson reconstruction to obtain intensity
- **Tracking** (ego-motion estimation)
 - Random diffusion in motion space
 - Particle filter: particle weights updated using the map: $\text{contrast} = |\log I(t) - \log I(t - \Delta t)|$



Simultaneous mosaicing and tracking

Input Events



Gradient Map Estimation



Scene Reconstruction

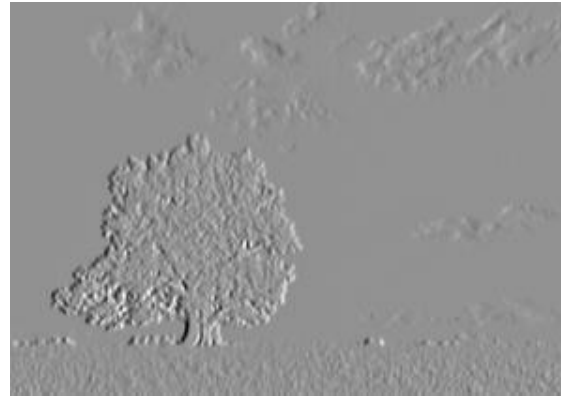


Image reconstruction by Poisson integration

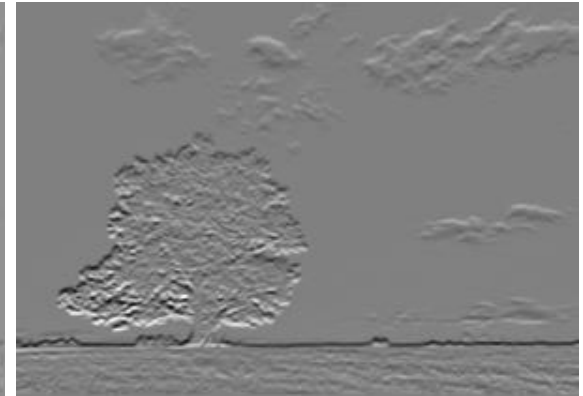
Integrate gradient map \mathbf{g} to get absolute brightness M



Original Image



Gradient in x direction
($g_x = \partial_x I$)



Gradient in y direction
($g_y = \partial_y I$)

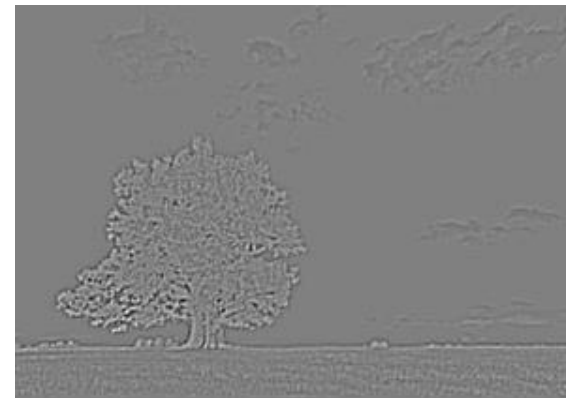


Reconstructed Image

2D integration



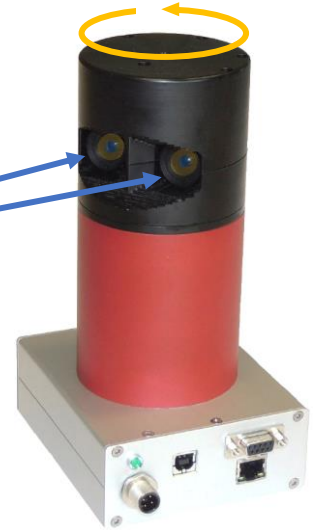
Solve Poisson eq:
($\Delta \tilde{I} = \text{div } \mathbf{g}$)
fast using the FFT



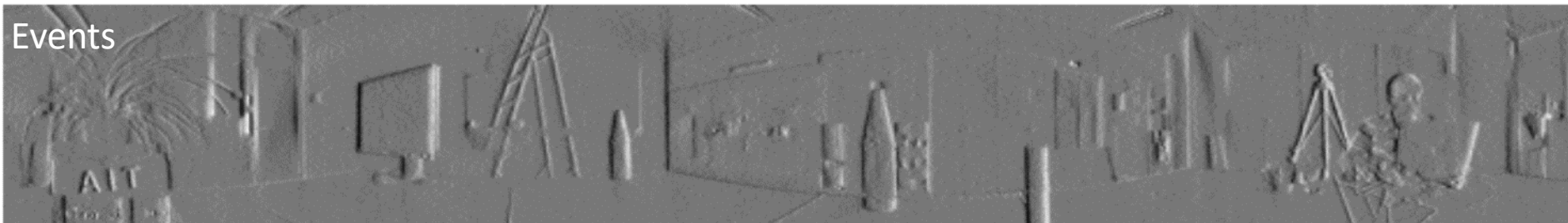
Divergence
($\text{div } \mathbf{g} = \partial_x g_x + \partial_y g_y$)

TUCO-3D Panoramic Imaging

- 360 deg Panoramic Imaging
- Exploit constrained motion
 - Two rotating 1D event cameras (stereo)
 - Event integration must satisfy periodic boundary conditions
- It also provides depth (3D)



TUCO-3D by AIT

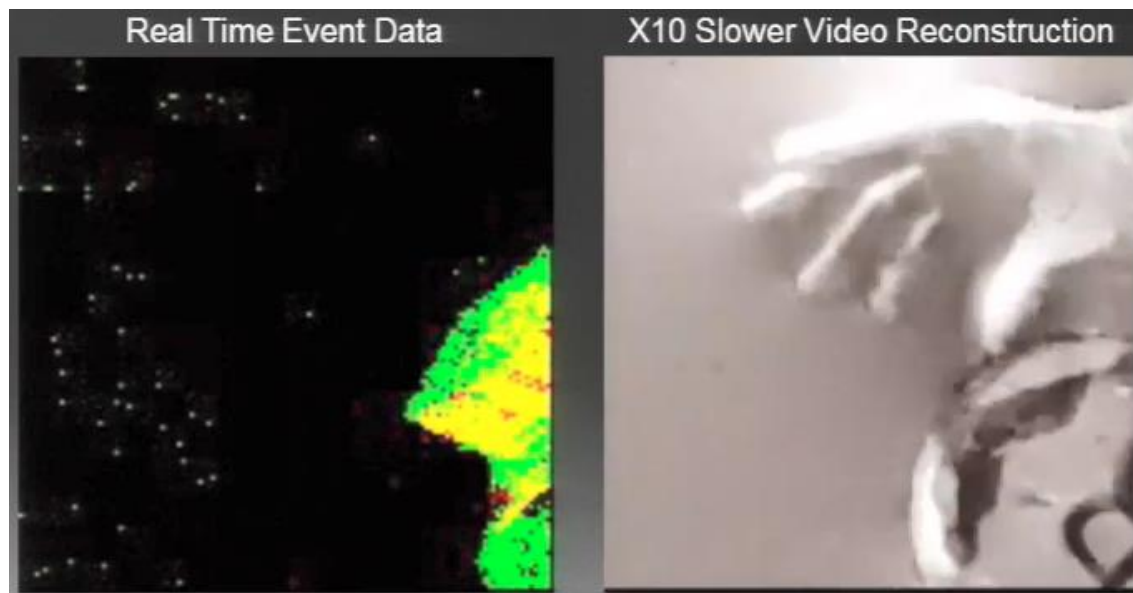
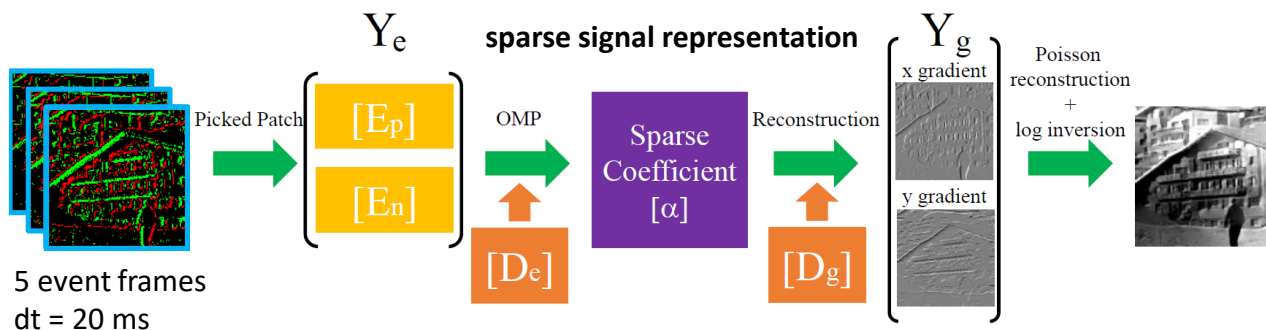


Integrate brightness changes, line by line



Image Reconstruction using a Sparse Dictionary

- Shows that there is **no need to estimate motion** for reconstruction
- **Patch-based dictionary of events** learnt from simulated data



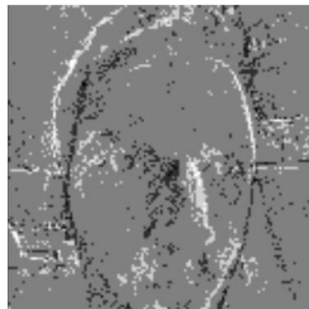
Green: positive
Red: negative
Yellow: both events

Video at **2 kHz**

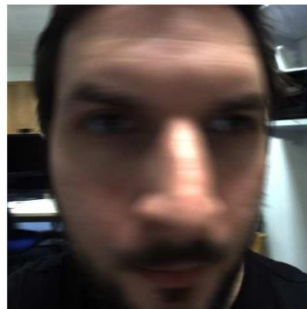
Patches of 9 x 9 pixels

SOFIE: Simultaneous Optical Flow & IE

- **Joint optimization** over image reconstruction and **optical flow** to explain a volume of events (voxel grid)
- More relaxed than SLAM methods; just need optical flow, i.e., works on **dynamic scenes** with **little assumptions about motion**.
- Solve a variational optimization problem:



(a) Raw event camera output



(b) Standard camera image



(c) Intensity estimate from events



(d) Optical flow from events

$$\min_{\mathbf{u}, L} \int_{\Omega} \int_T \left(\lambda_1 \|\mathbf{u}_{\mathbf{x}}\|_1 + \lambda_2 \|\mathbf{u}_t\|_1 + \lambda_3 \|L_{\mathbf{x}}\|_1 + \lambda_4 \|\langle L_{\mathbf{x}}, \delta_t \mathbf{u} \rangle + L_t\|_1 + \lambda_5 h_{\theta}(L - L(t_p)) \right) dt d\mathbf{x} + \int_{\Omega} \sum_{i=2}^{|P(\mathbf{x})|} \|L(t_i) - L(t_{i-1}) - \theta \rho_i\|_1 d\mathbf{x},$$

Smoothness terms

Optical flow term (brightness constancy)

No-event term

Event term

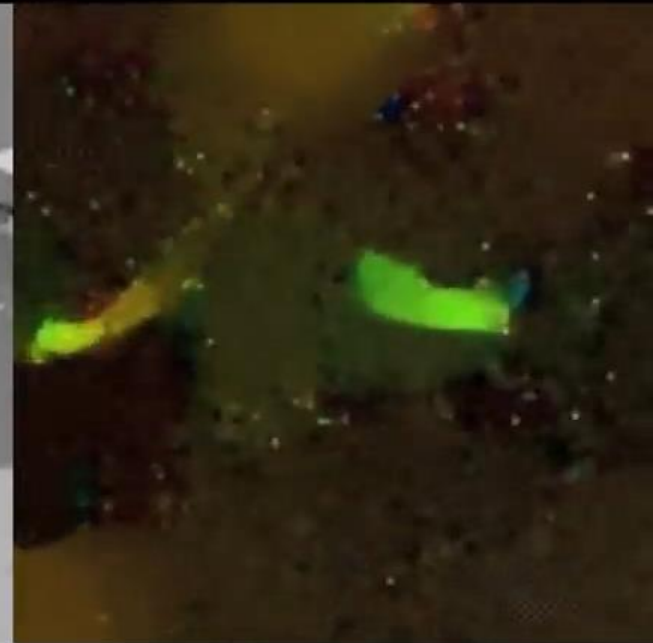
SOFIE: Simultaneous Optical Flow & IE



Camera Image



Reconstruction

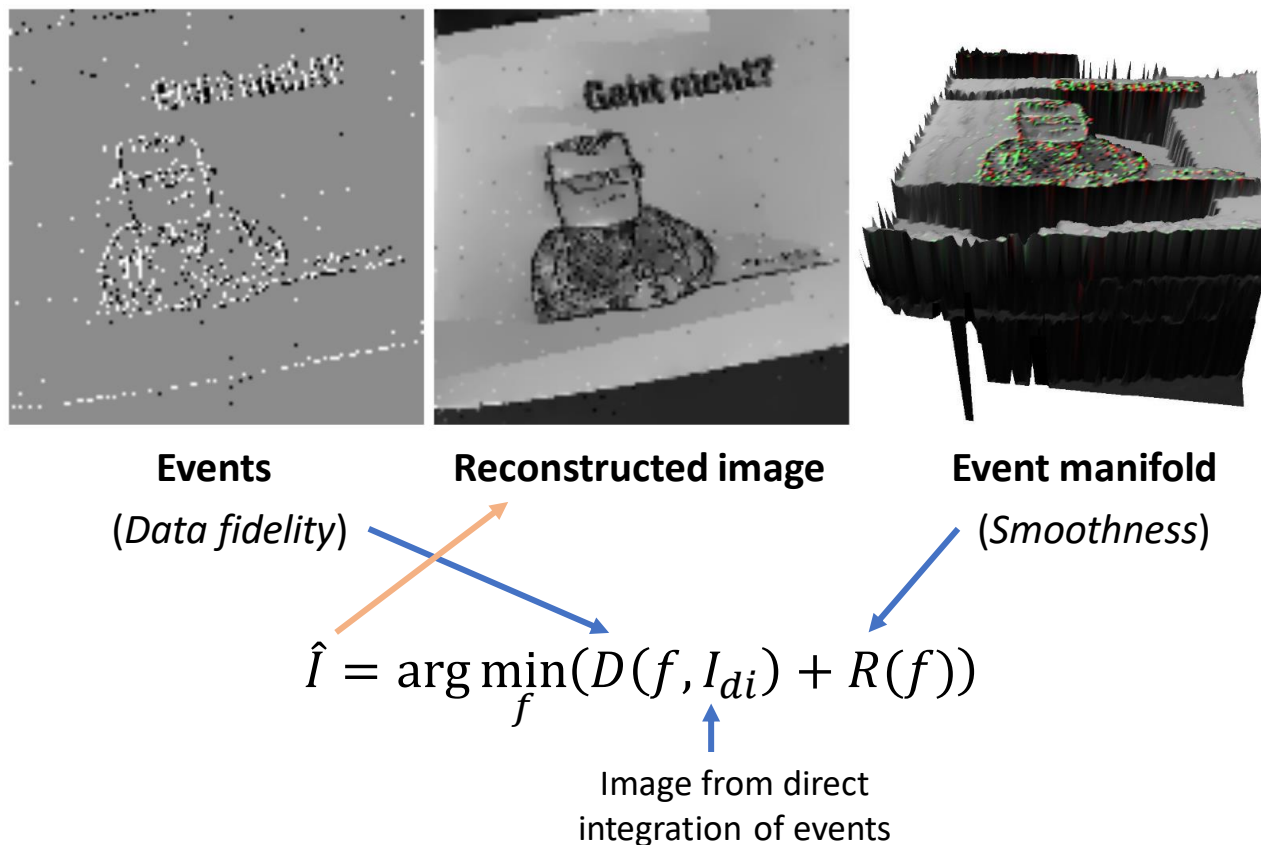


Optical Flow

1 frame = 4ms

Reconstruction using “Manifold Regularisation”

- Does not need to estimate motion
- Reconstruction is posed as **variational nonlinear image denoising**, using the time surface (event timestamps) to guide the denoising



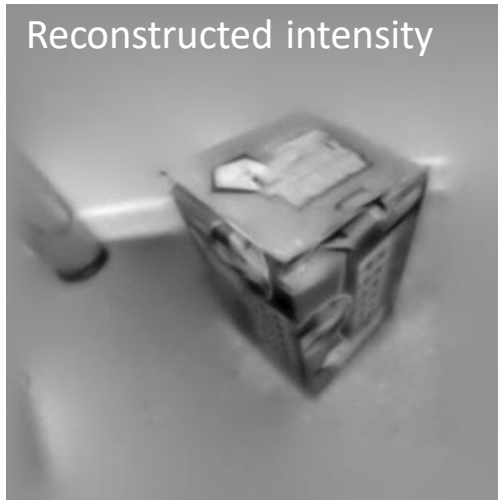
Reconstruction using “Manifold Regularisation”

- Does not need to estimate motion
- Reconstruction is posed as **variational nonlinear image denoising**, using the time surface (event timestamps) to guide the denoising

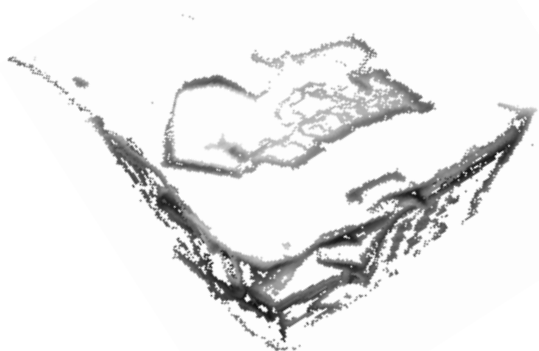


Event-based 6-DOF SLAM with 3 parallel filters

Parallel 6 DOF Tracking & Mapping in real time on a GPU



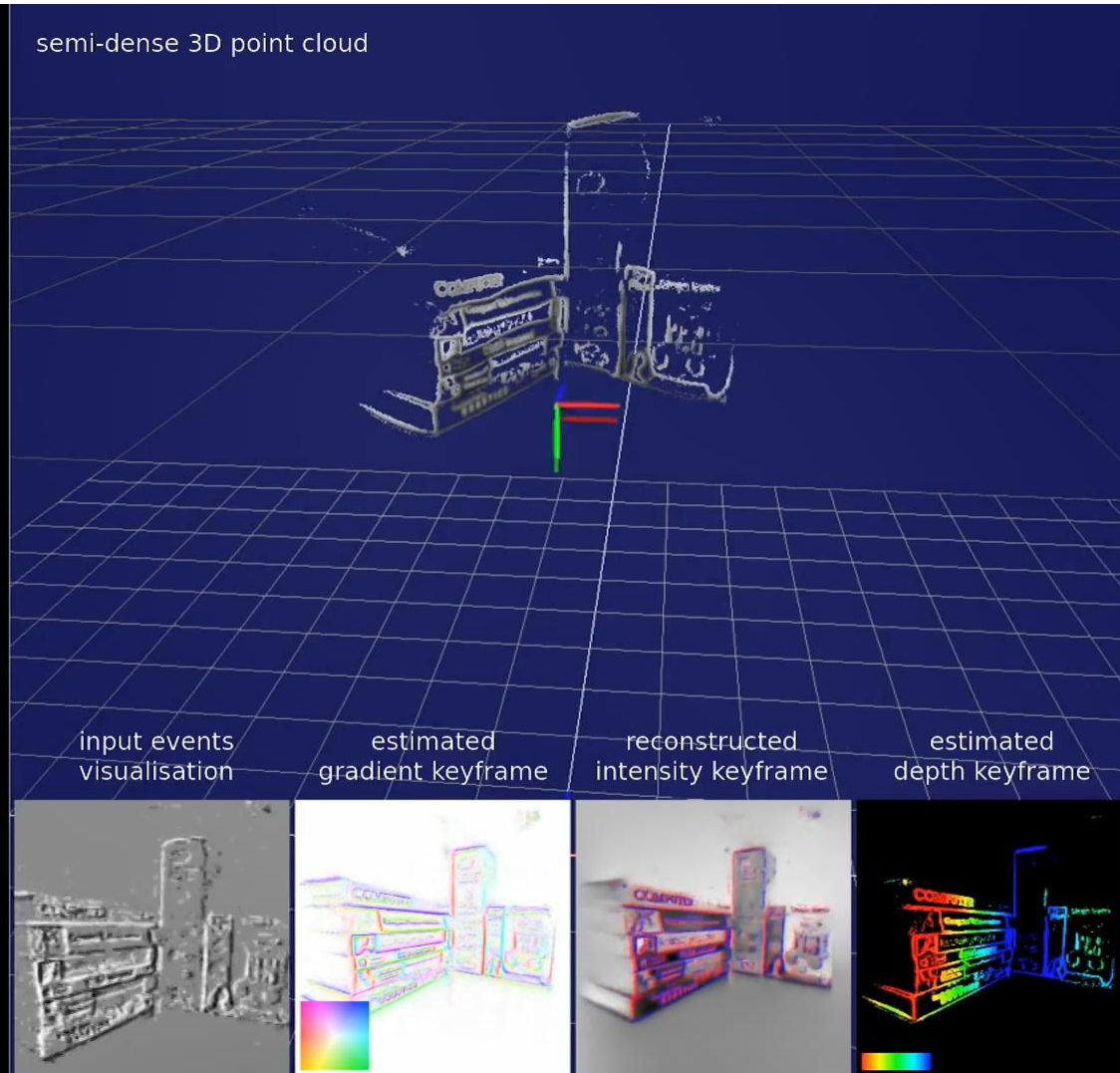
3D map of the scene



- **3 Kalman Filters** running in parallel.
Each filter needs the output from the others.
- **Tracking:** EKF in 6 DOF pose
 - Motion model: constant position
 - Use contrast $h_{\mathbf{x}}(\mathbf{x}^{(t|t-\tau)}) = \mathcal{I}_l(\mathbf{p}_w^{(t)}) - \mathcal{I}_l(\mathbf{p}_w^{(t-\tau_c)})$ to update pose (needs depth and intensity)
- **Intensity reconstruction:** pixel-wise EKF like Kim'14 & robust Poisson (Huber norm)
- **Mapping:** pixel-wise EKF on inverse depth, using contrast $h_{\rho} = \mathcal{I}_l(\mathbf{p}_w^{(t)}) - \mathcal{I}_l(\mathbf{p}_w^{(t-\tau_c)})$

Event-based 6-DOF SLAM with 3 parallel filters

Parallel 6 DOF Tracking & Mapping in real time on a GPU



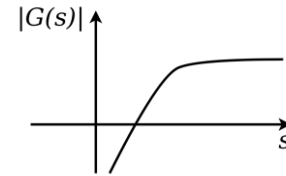
Event-based 6-DOF SLAM on a CPU

HDR image reconstruction from the output of SLAM



Reconstruction by Temporal Filtering

- Replace **pixel-wise** direct integration with a **high-pass temporal filter** to remove accumulated event noise
- No spatial filtering needed



Direct integration of events



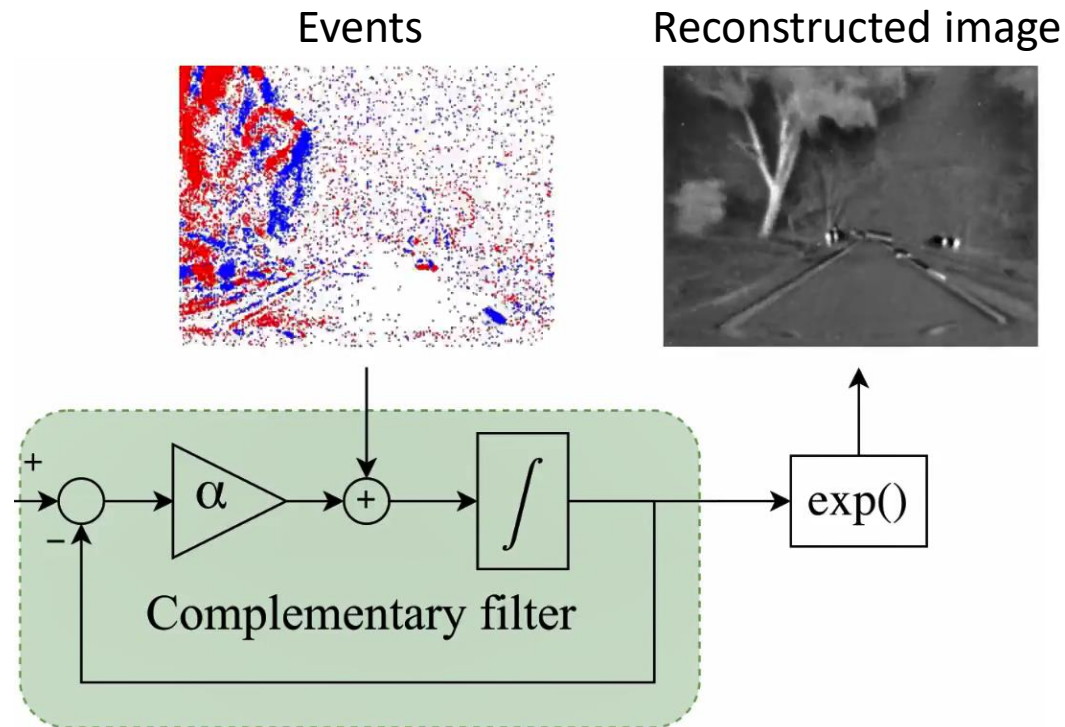
High-pass filtered (in time)



But slow-time information (low freq.)
(static background) is lost.
⇒ fuse with grayscale frames from DAVIS
(complementary filter)

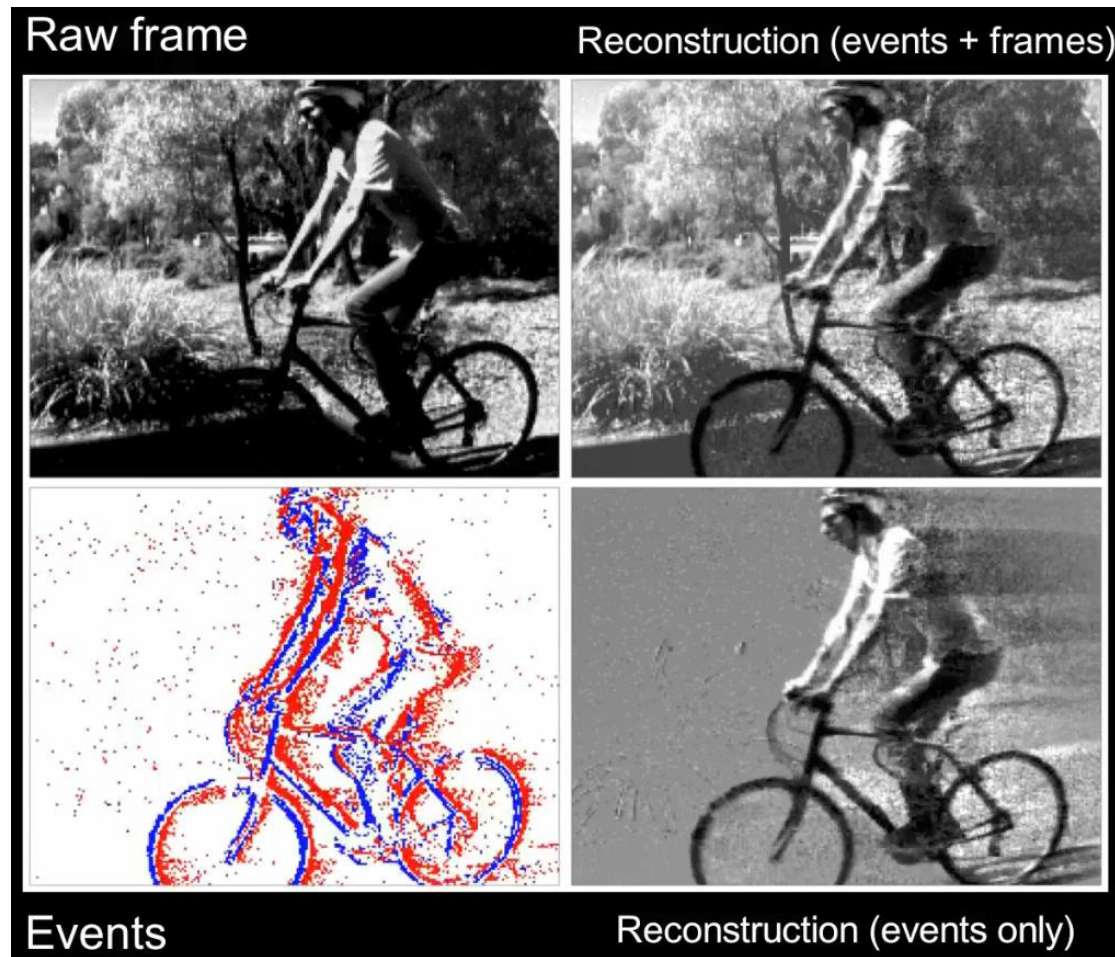
Reconstruction by Temporal Filtering

- Replace **pixel-wise** direct integration with a **high-pass filter** to remove accumulated event noise
- The internal **state of the filter** is an **image** that is updated **asynchronously, per-pixel** with each incoming event



Reconstruction by Temporal Filtering

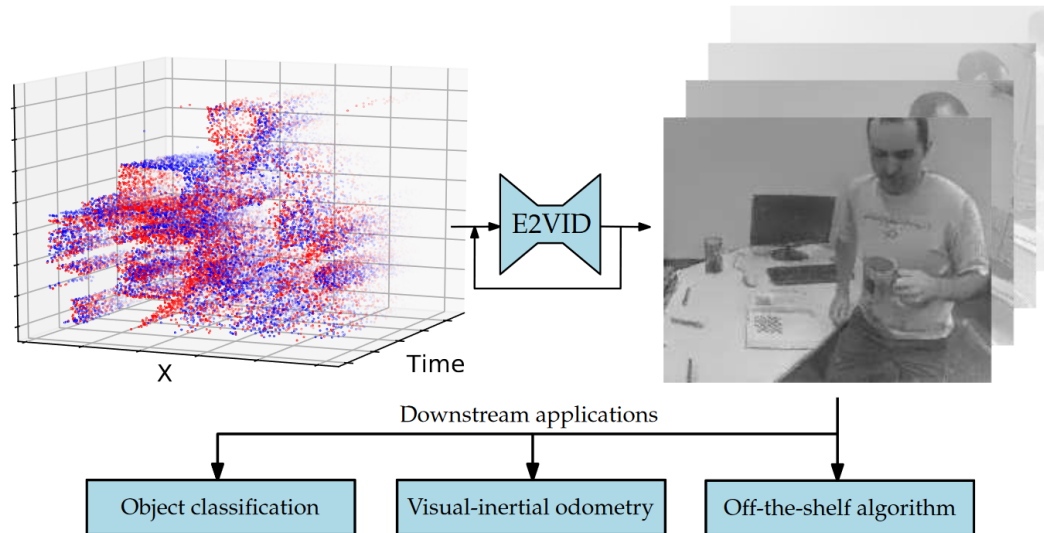
- Replace **pixel-wise** direct integration with a **high-pass filter** to remove accumulated event noise



Reconstruction using Deep-Learning

- **Events-to-Video (E2VID)**

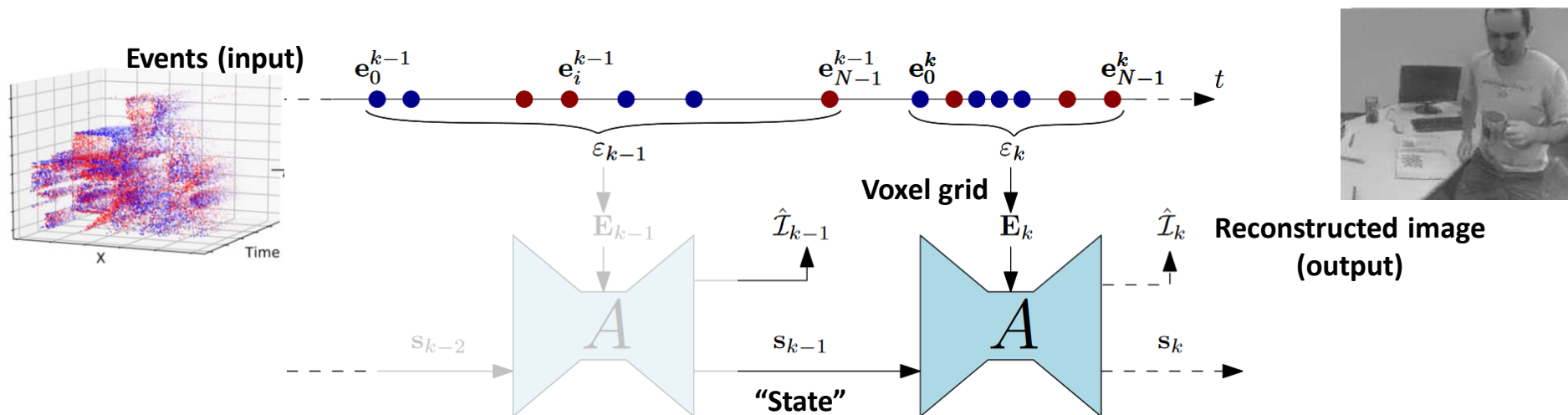
- Deep Learning method: **recurrent network** (with a U-Net)
- **Loss function: perceptual** (LPIPS) + temporal consistency
- Trained on simulation, transfers well to real-world data
- Shows a **big improvement** with respect to previous methods
- Shows reconstructed images can be used on **off-the-shelf computer vision methods** designed for image data



Reconstruction using Deep-Learning

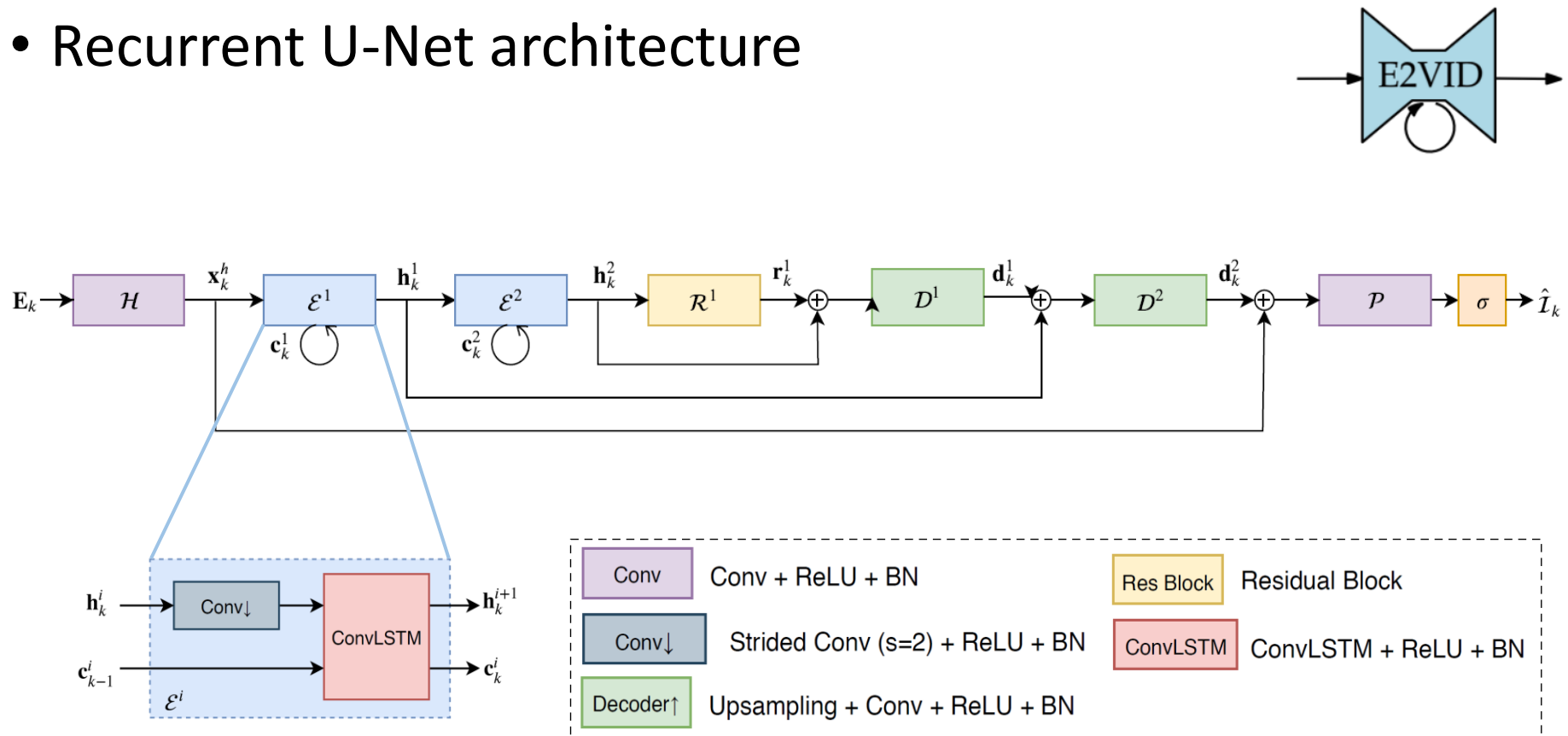
- **Events-to-Video (E2VID)**

- Deep Learning method: **recurrent network** (with a U-Net)
- **Loss function: perceptual** (LPIPS) + temporal consistency
- Trained on simulation, transfers well to real-world data
- Shows a **big improvement** with respect to previous methods
- Shows reconstructed images can be used on **off-the-shelf computer vision methods** designed for image data



Network architecture

- Recurrent U-Net architecture



E2VID - Results



Huawei P20 Pro (240 FPS)



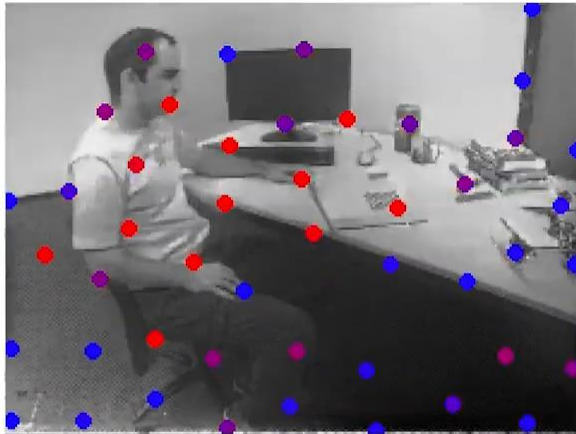
Our reconstruction (5400 FPS)

100 x slow motion

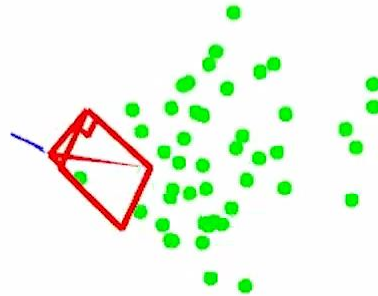
E2VID - Applications of Reconstructed images



Events



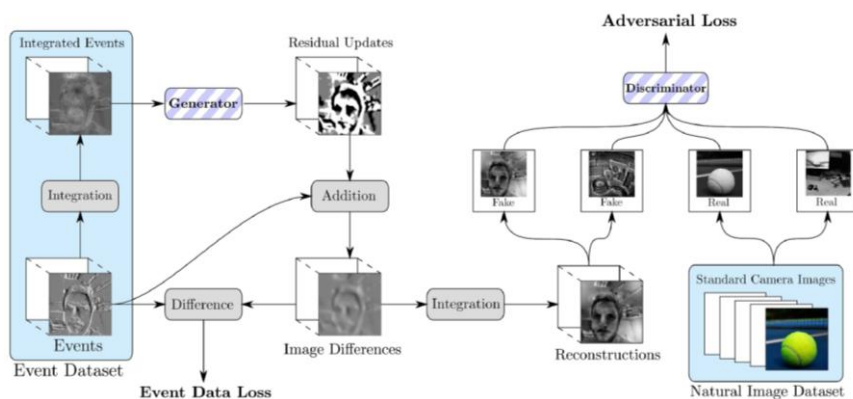
Our reconstruction
+ tracked features



VINS-Mono running on our reconstruction
from events

Unsupervised Deep Learning, using GANs

- Bardow (Ch. 6): Reconstruction using only Natural Image Priors



- Mostafavi CVPR 2019

