



Vision Algorithms for Mobile Robotics

Lecture 14 Event-based Vision

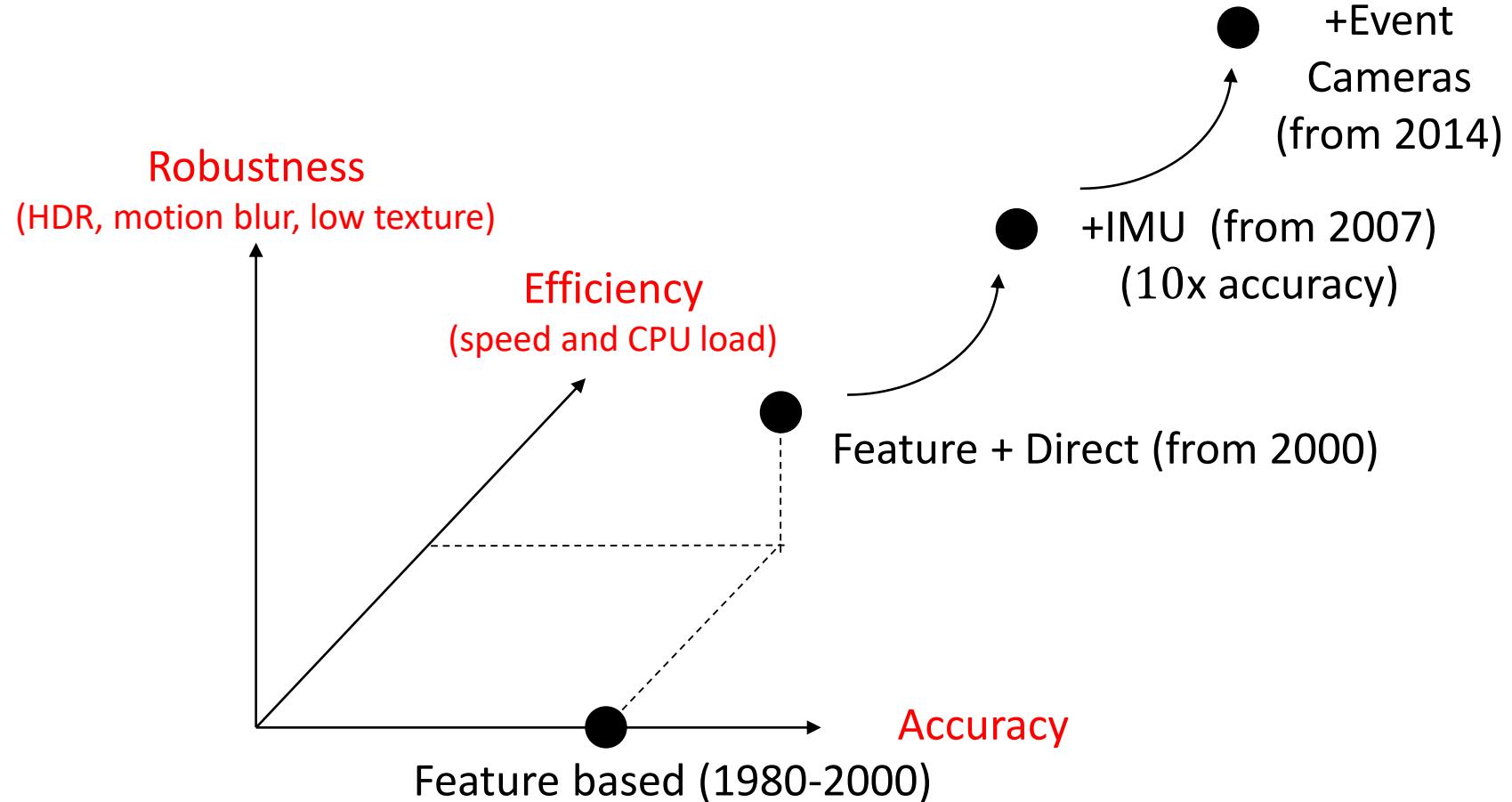
Davide Scaramuzza

<http://rpg.ifi.uzh.ch>

Lab Exercise – Today

Q&A on Exams followed by final VO integration

A Short Recap of the last 30 years of VIO



References

- **Tutorial paper:**
Gallego, Delbruck, Orchard, Bartolozzi, Taba, Censi, Leutenegger, Davison, Conradt, Daniilidis, Scaramuzza,
Event-based Vision: A Survey, IEEE Transactions of Pattern Analysis and Machine Intelligence, 2020. [PDF](#)
- List of event camera papers, codes, datasets, companies: https://github.com/uzh-rpg/event-based_vision_resources
- Event-camera simulator: <http://rpg.ifi.uzh.ch/esim.html>
- More on event camera research: http://rpg.ifi.uzh.ch/research_dvs.html

Open Challenges in Computer Vision

- The past 60 years of research have been devoted to frame-based cameras but they are not good enough

Latency & Motion blur



Dynamic Range



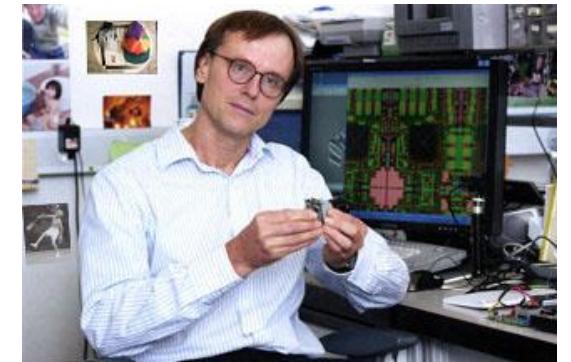
- **Event cameras** do not suffer from these problems

What is an Event Camera

First commercialized by Prof. T. Delbrück in 2008 at the Institute of Neuroinformatics of UZH & ETH under the name of Dynamic Vision Sensor (DVS)

Advantages

- **Low-latency** (~1 micro-seconds)
- **High-dynamic range (HDR)** (140 dB instead 60 dB)
- **High updated rate** (1 MHz)
- **Low power** (10mW instead 1W)



Prof. Tobi Delbrück, UZH & ETH Zurich

Challenges

- **Paradigm shift:** Requires totally **new vision algorithms** because:
 - **Asynchronous pixels**
 - **No intensity information** (only binary intensity changes)

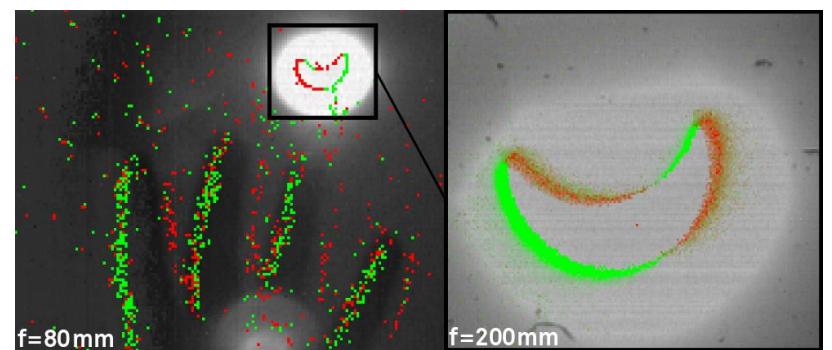
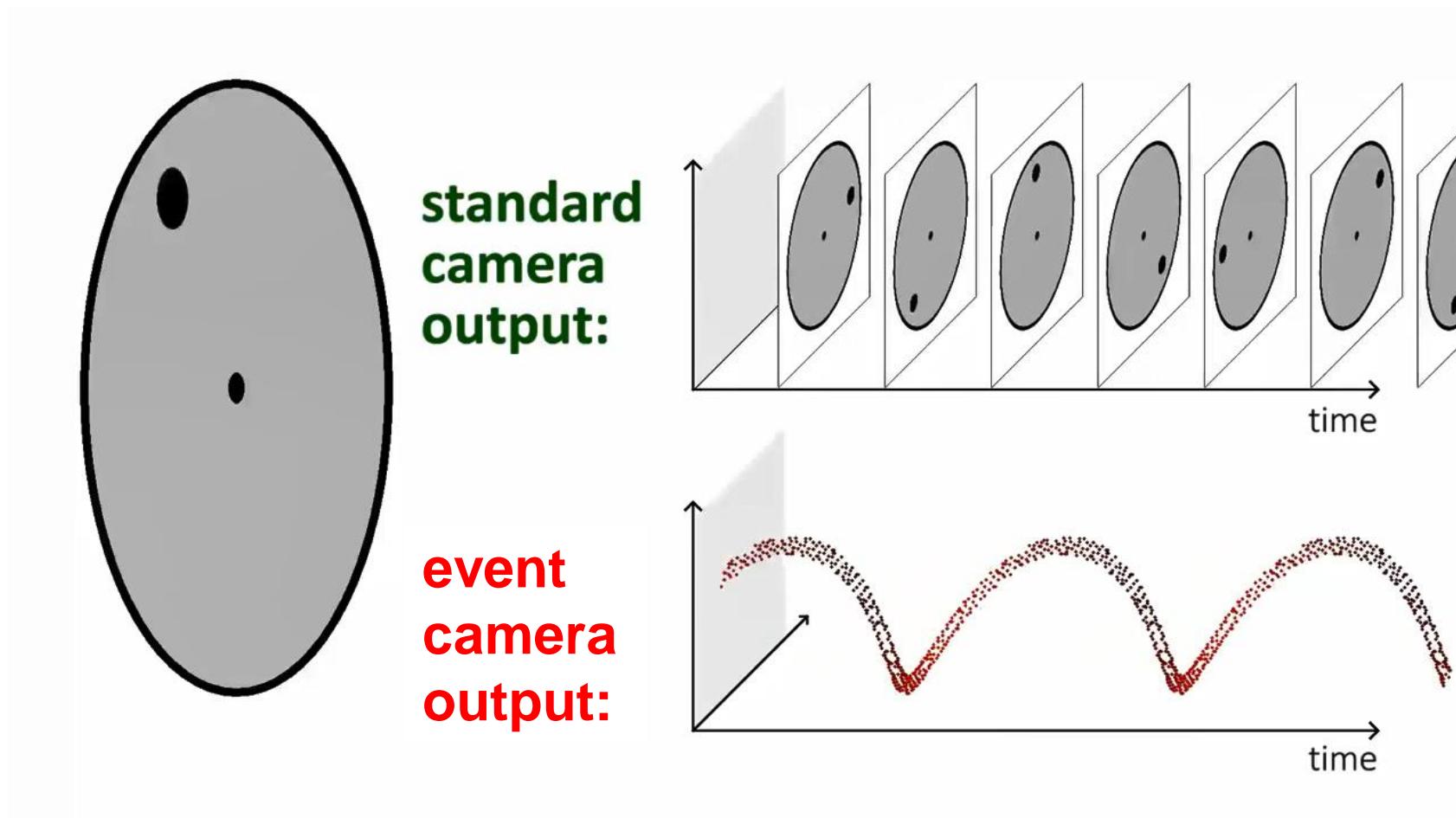


Image of solar eclipse captured by an event camera without black filter

Animation of an Event Camera Output



Video from here: <https://youtu.be/LauQ6LWTkxM?t=30>

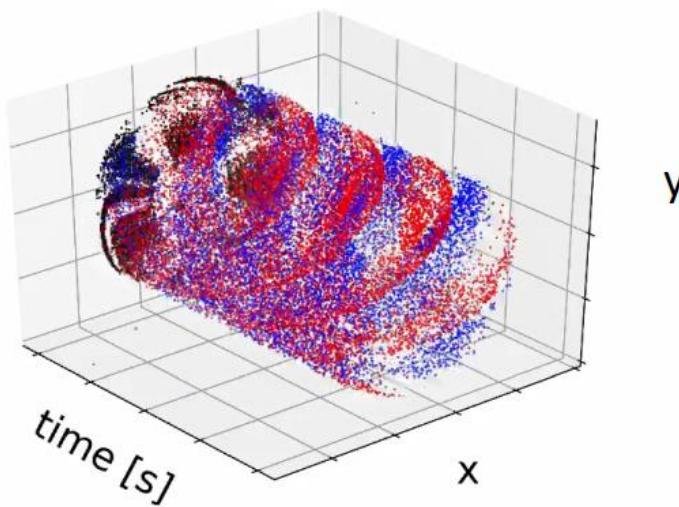
Conventional frames



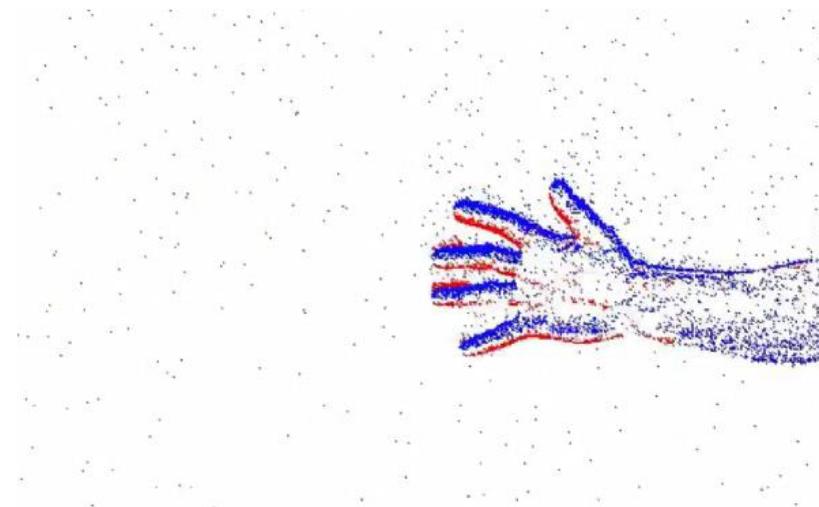
Conventional frames



Events



Events in the **space-time** domain (x, y, t)



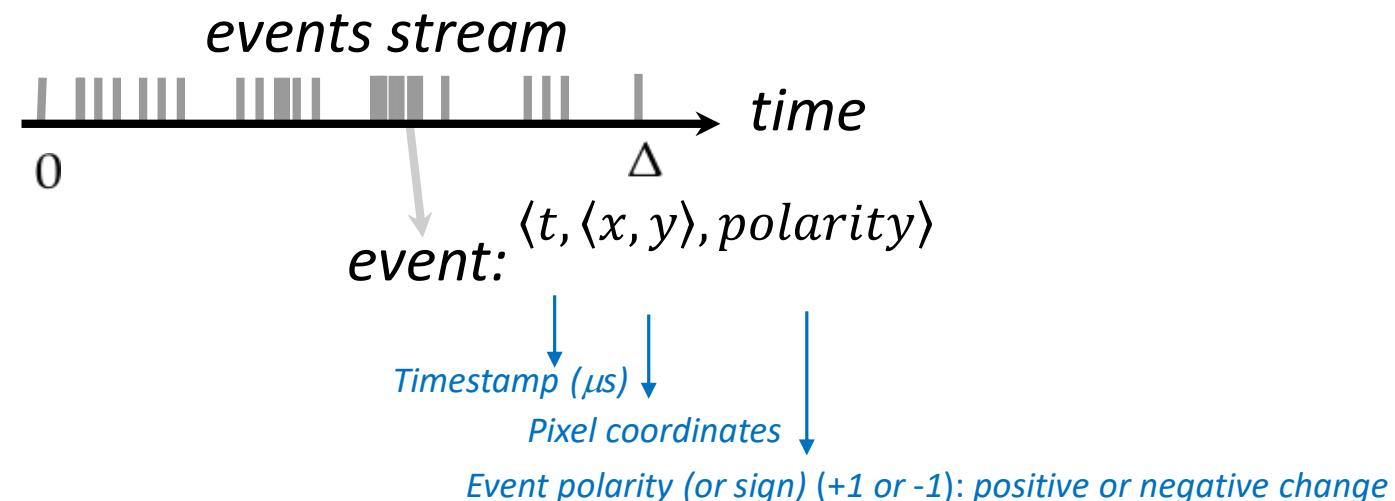
Events in the **image** domain (x, y)
Integration time but can be arbitrary: from 1 microsecond to infinity)

Standard Camera vs. Event Camera

- A traditional camera outputs frames at **fixed time intervals**:



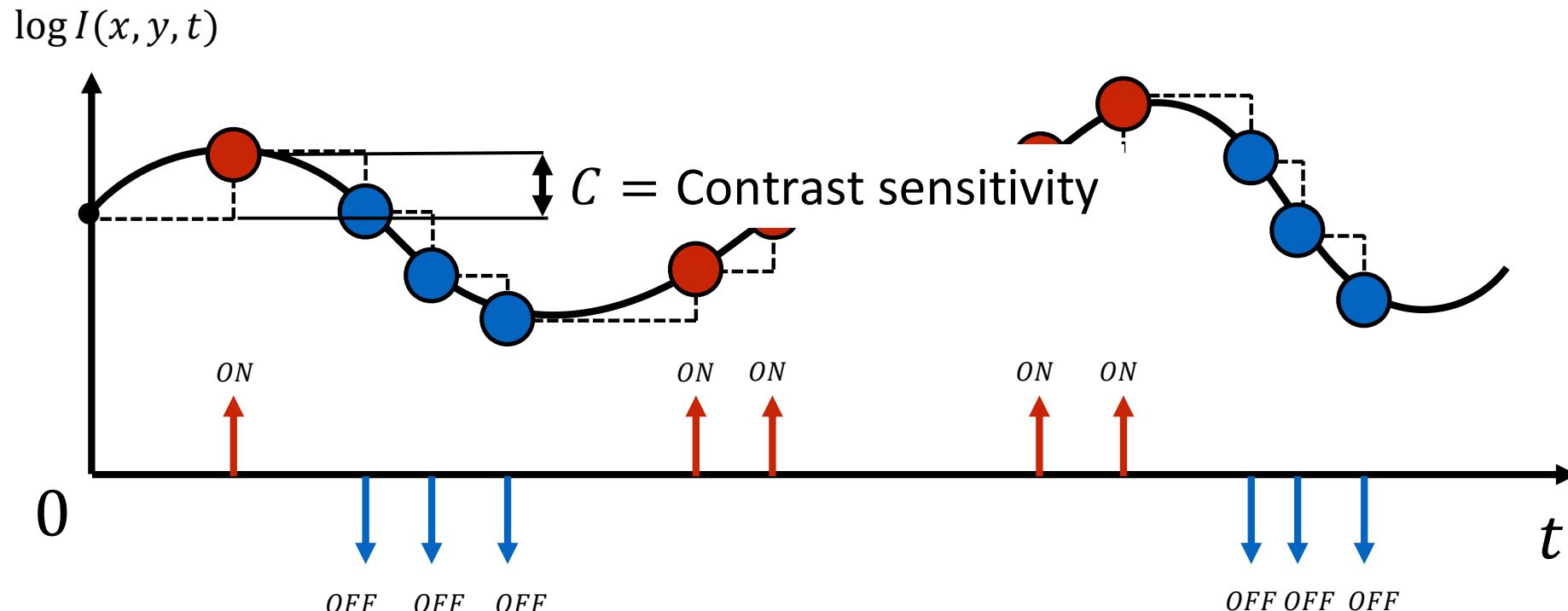
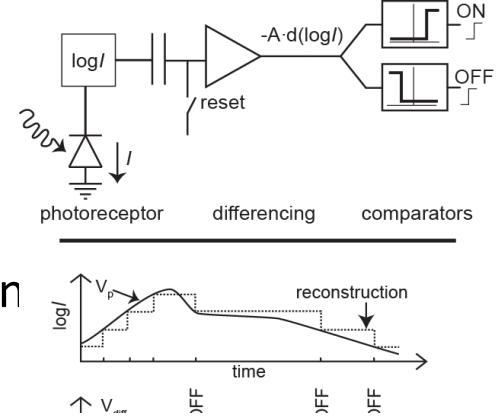
- By contrast, an event camera outputs **asynchronous events** at **microsecond resolution**. An event is generated each time a single pixel detects a change of intensity



Generative Event Model

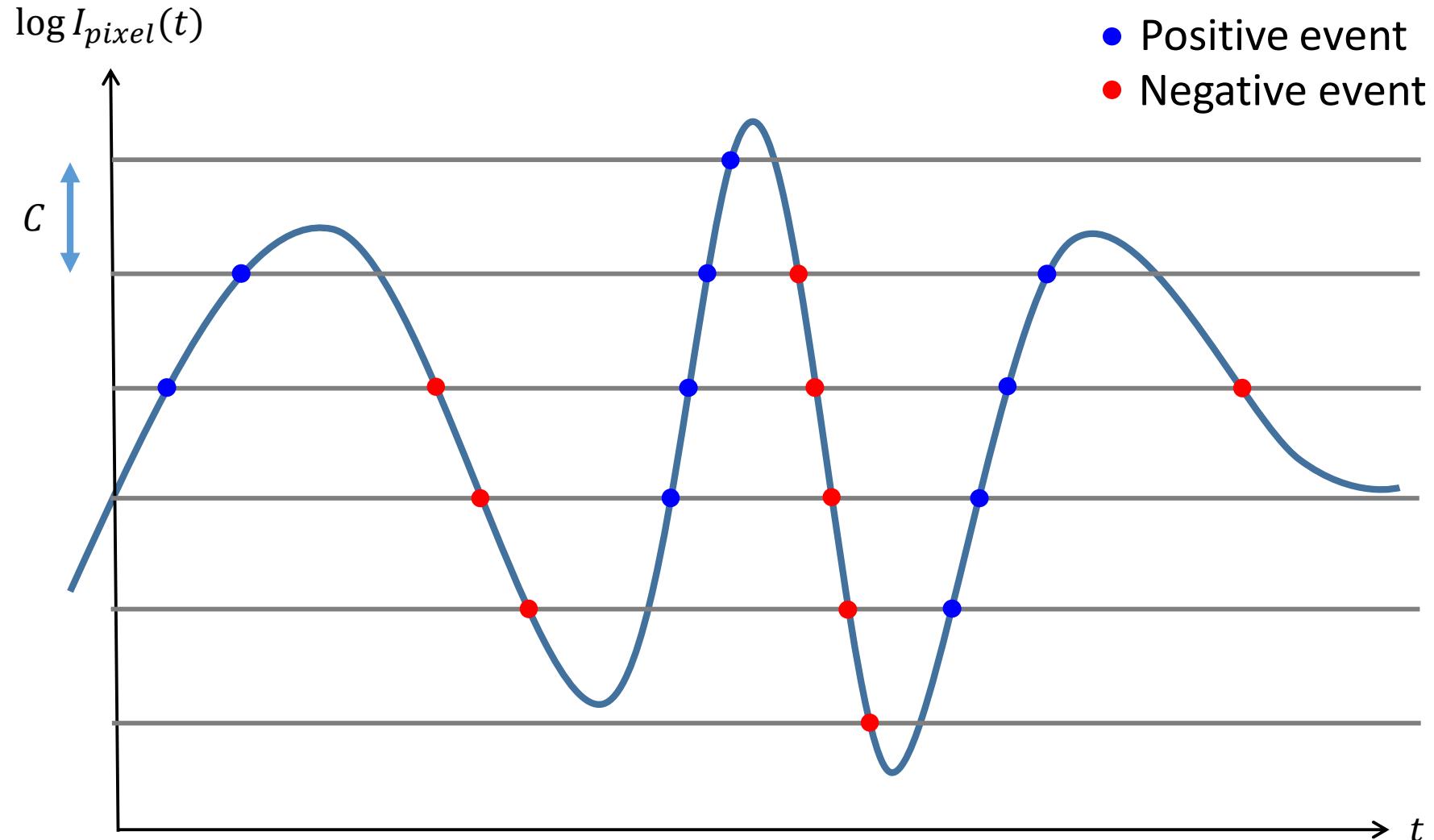
- Consider the intensity at a **single pixel** (x, y) . An event is generated when the following is satisfied:

$$\log I(x, y, t + \Delta t) - \log I(x, y, t) = \pm C$$

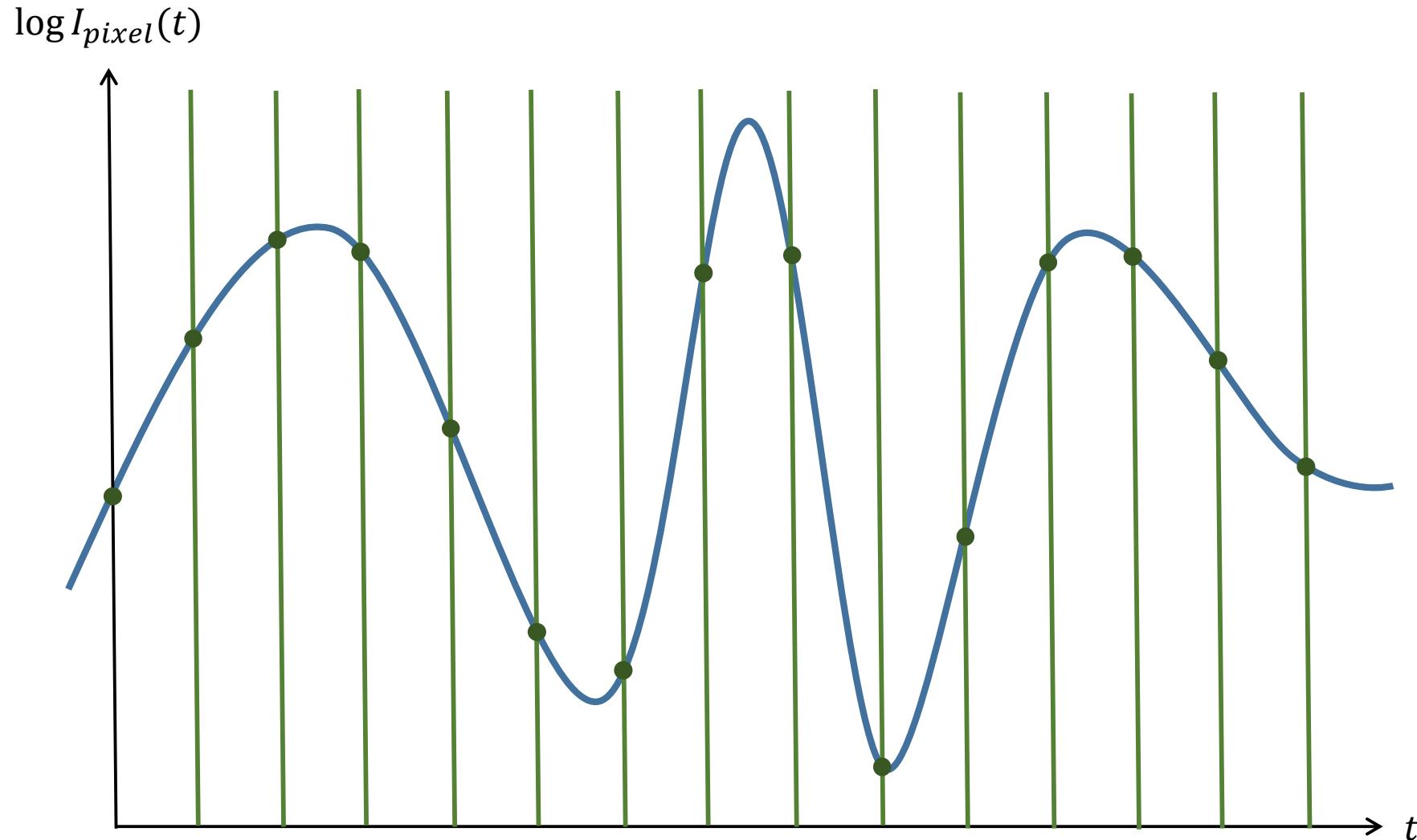


Events are triggered **asynchronously**

Event cameras sample intensity when this crosses a threshold (Level-crossing sampling)



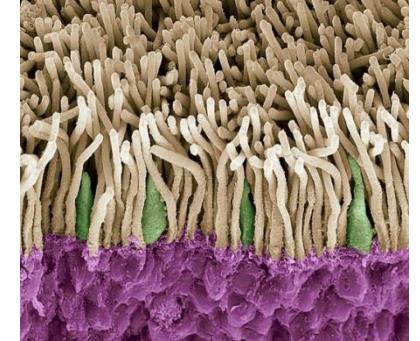
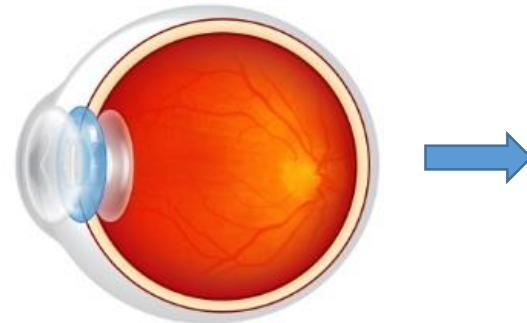
Standard cameras sample intensity at uniform time intervals (uniform time sampling)



Event cameras are inspired by the Human Eye

Human retina:

- 130 million **photoreceptors**
- But only 2 million **axons**!



Brain

Who sells event cameras and how much are they?

- Prophesee & SONY:
 - **ATIS sensor:** events, IMU, absolute intensity at the event pixel
 - Resolution: **1M pixels**
 - Cost: ~5,000 USD
- Inivation & Samsung
 - **DAVIS sensor:** frames, events, IMU.
 - Resolution: **VGA** (640x480 pixels)
 - Cost: ~5,000 USD
- CelePixel Technology & Omnivision:
 - **Celex One:** events, IMU, absolute intensity at the event pixel
 - Resolution: **1M pixels**
 - Cost: ~1,000 USD
- **Cost to sink to <5\$** when killer application found
(recall first ToF camera (>10,000 USD) today <5 USD)



SAMSUNG



Omnivision®

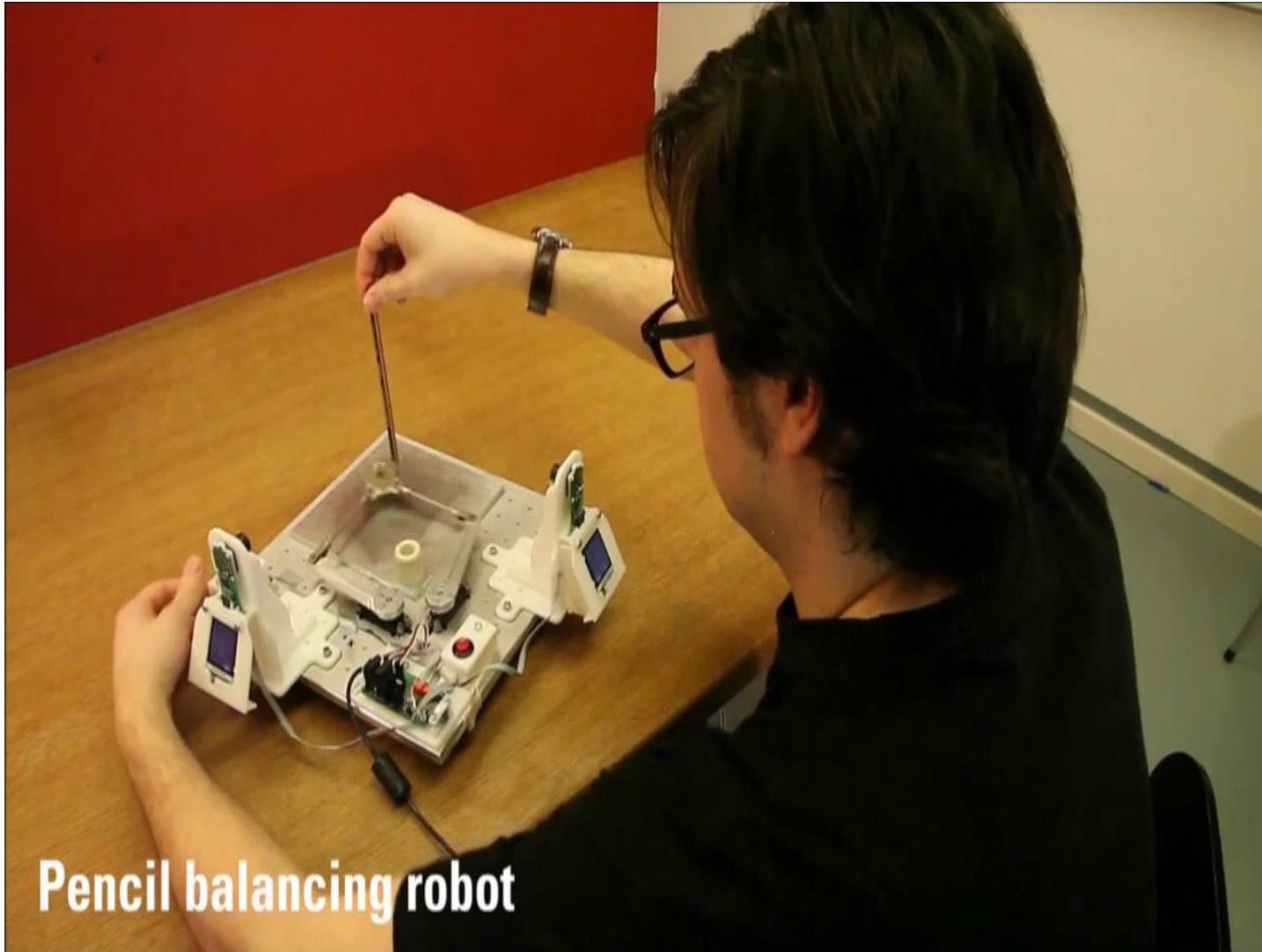


Event Camera Demo



<https://youtu.be/QxJ-RTbpNXw>

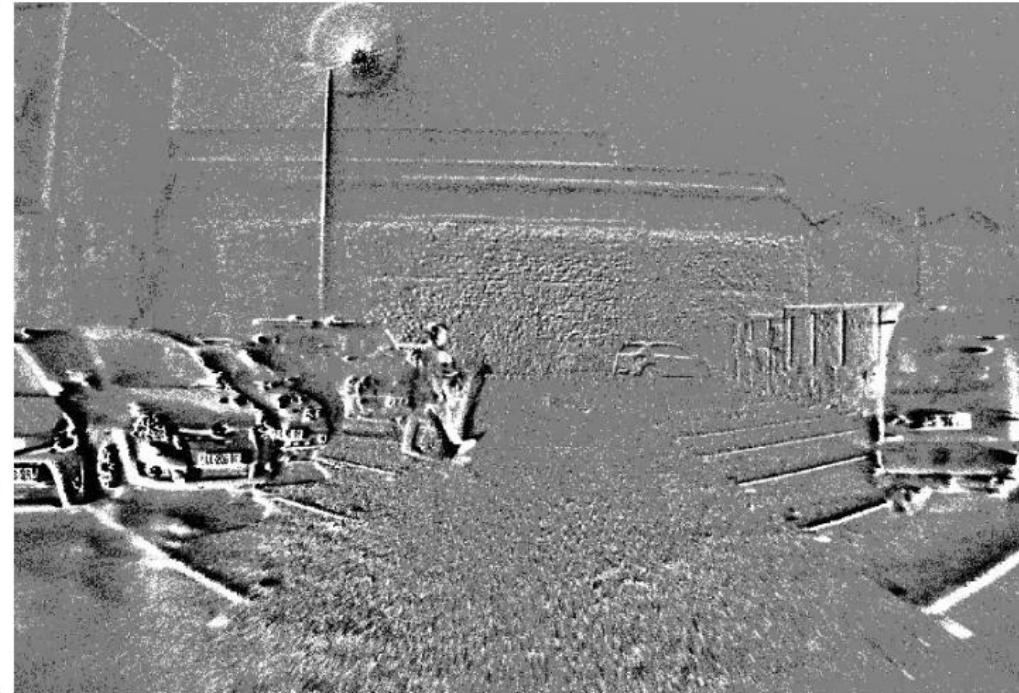
Event Camera Demo



Low-light Sensitivity (night drive)



GoPro Hero 6



Aggregated event image

(pixel intensity equal to the sum of positive (+1) and negative (-1) events in a given time interval)

High-speed Camera vs. Event Camera



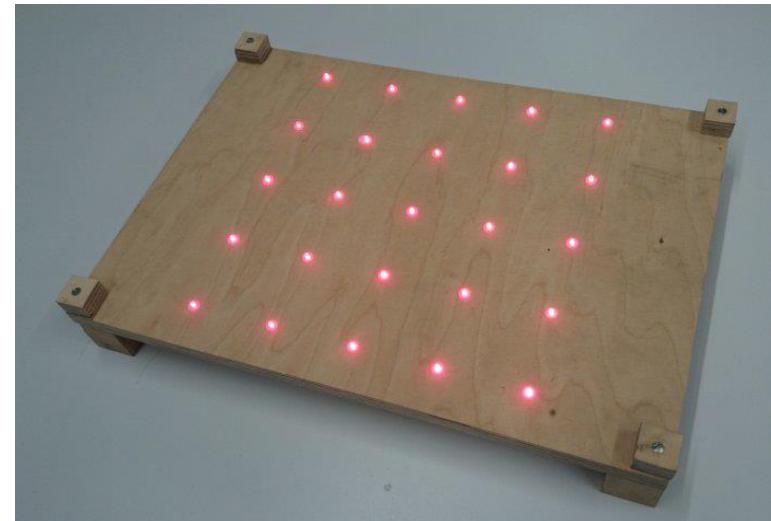
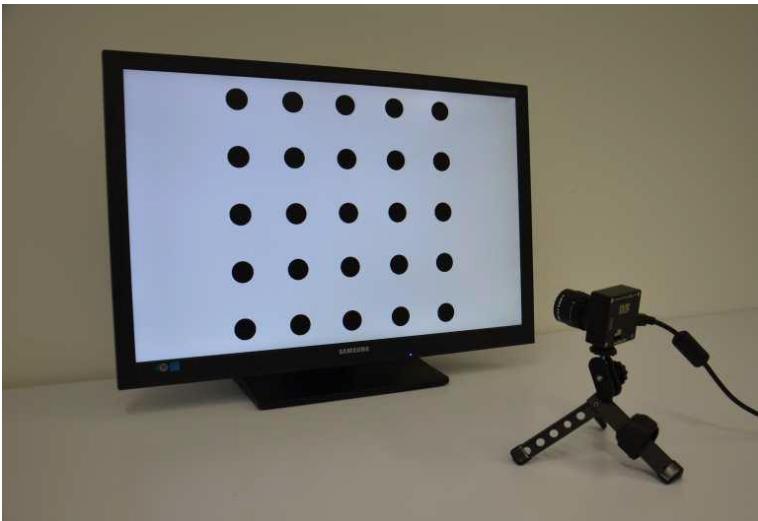
	High speed camera	Standard camera	Event Camera
Max fps or measurement rate	Up to 1MHz	100-1,000 fps	1MHz
Resolution at max fps	64x16 pixels	>1Mpxl	>1Mpxl
Bits per pixels (event)	12 bits	8-10 per pixel	~40 bits/event {t,(x,y),p)}
Weight	6.2 Kg	30 g	30 g
Active cooling	yes	No cooling	No cooling
Data rate	1.5 GB/s	32MB/s	~1MB/s on average (depends on dynamics & contrast threshold)
Mean power consumption	150 W + external light	1 W	1 mW
Dynamic range	not specified	60 dB	140 dB

Current commercial applications

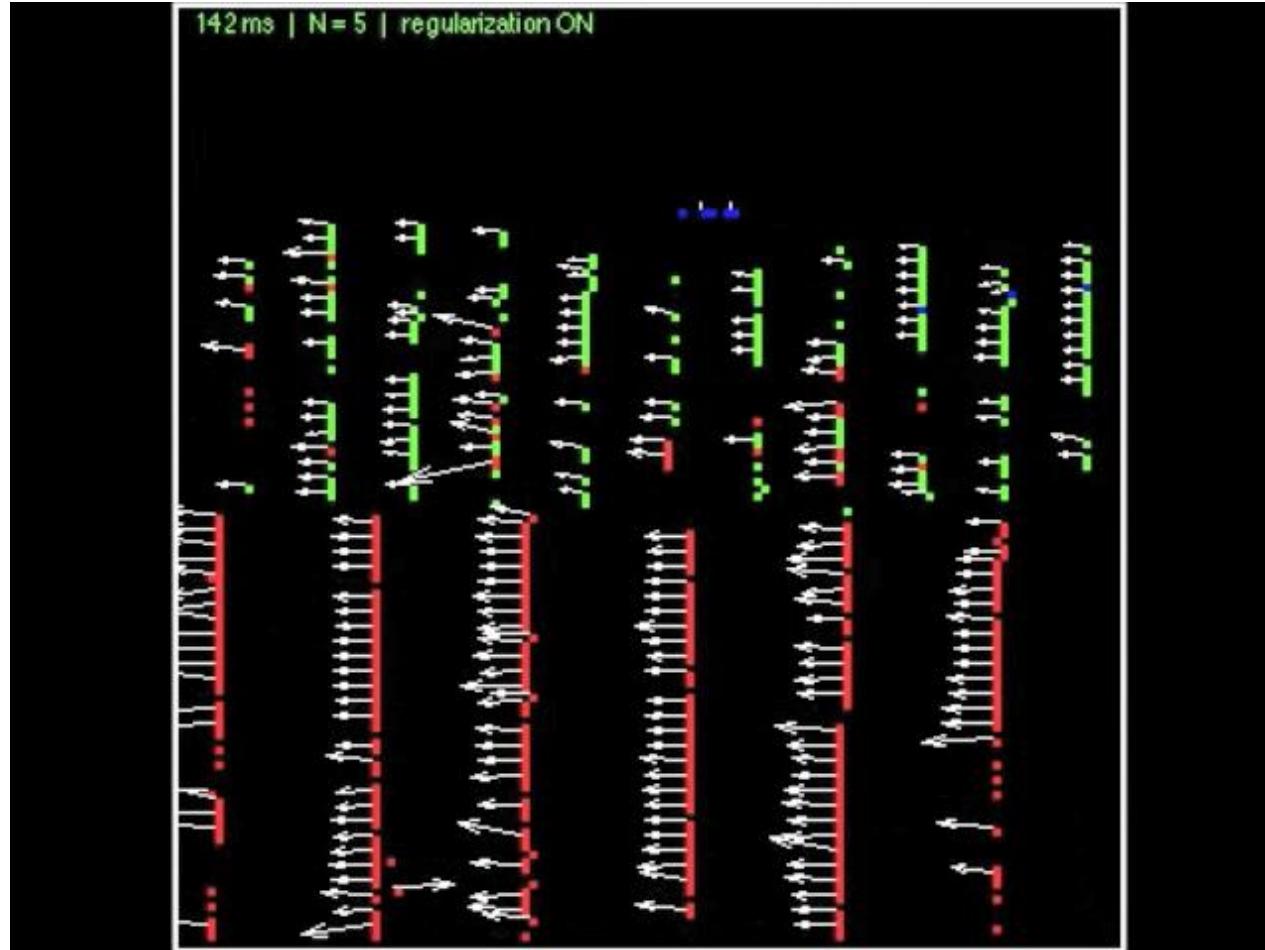
- **Monitoring and surveillance**
 - Action and gesture recognition in HDR scenes
- **Industrial automation**
 - Fast object counting
- **Computational photography**
 - Deblurring, super resolution, HDR, slow-motion video
- **Automotive:**
 - low-latency detection, object classification, low-power and low-memory storage

Calibration of an Event Camera

- Standard **pinhole camera model** still valid (same optics)
- Standard passive calibration patterns **cannot be used**
 - need to move the camera → inaccurate corner detection
- **Blinking patterns** (computer screen, LEDs)
- ROS DVS driver + intrinsic and extrinsic mono & stereo calibration: https://github.com/uzh-rpg/rpg_dvs_ros

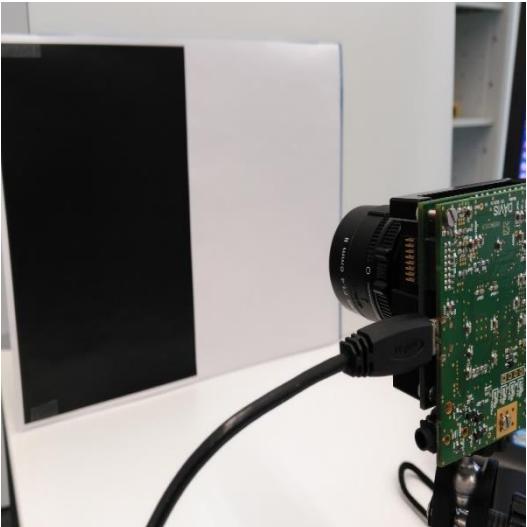


A Simple Optical Flow Algorithm

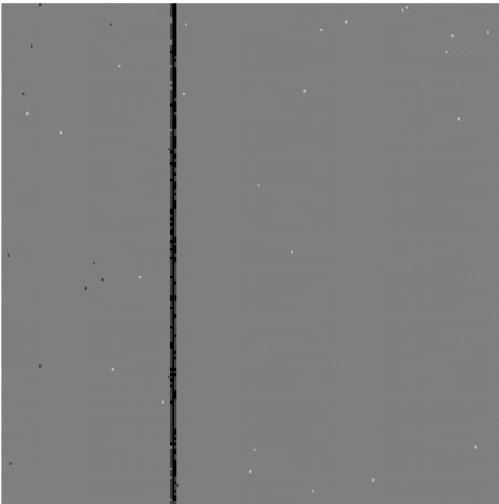


A Simple Optical Flow Algorithm

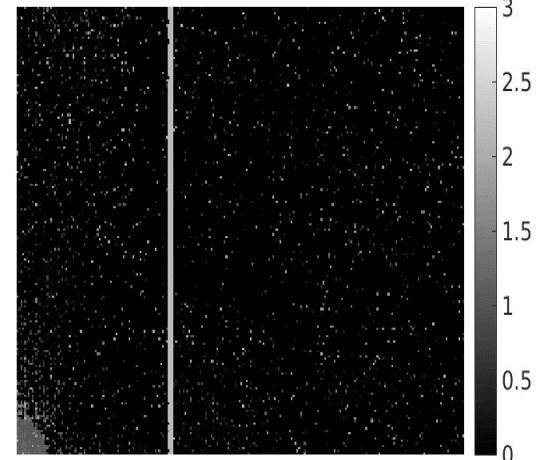
- Let's assume pure horizontal motion
- White pixels become black → brightness decrease → negative events (in black color)



Event image (1000 events). $t = 2.228$



Time of the last event

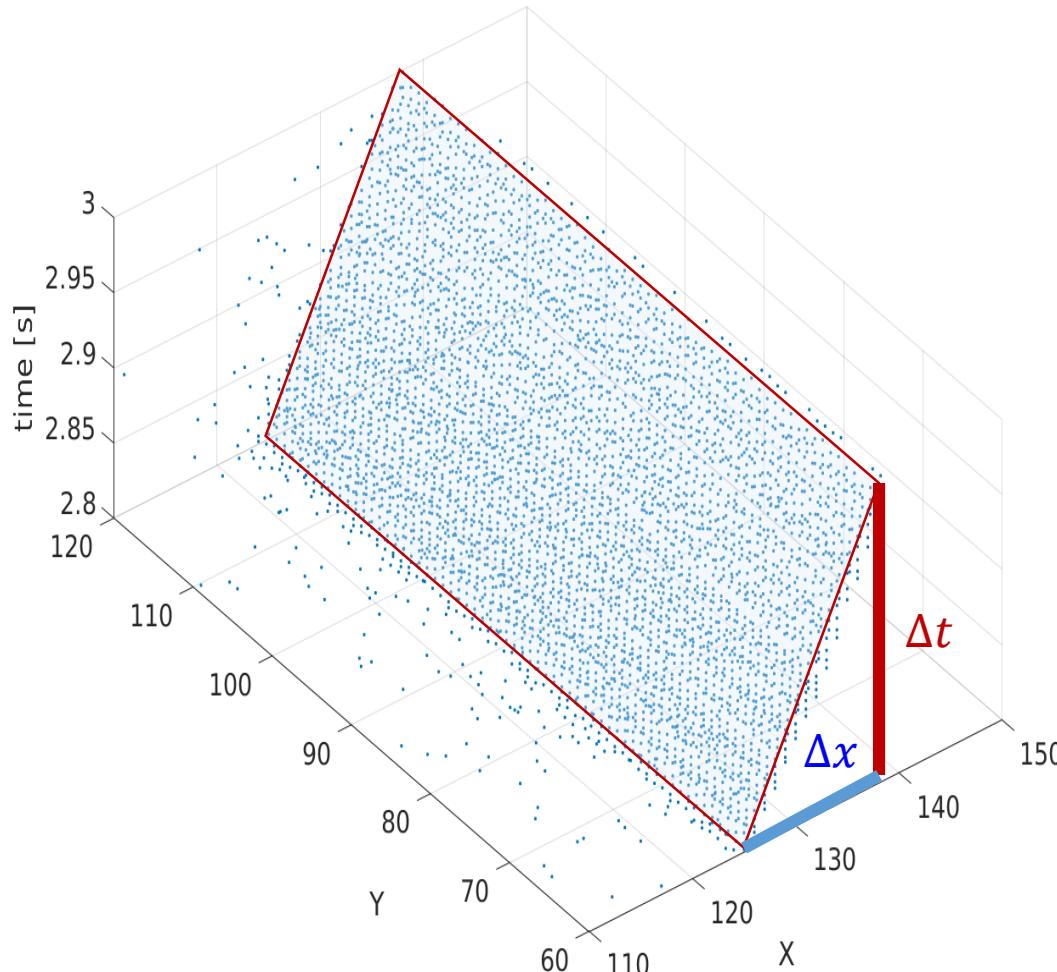


A Simple Optical Flow Algorithm

- The same edge, visualized in space-time
- Events are represented by dots

The edge is moving at
a speed of:

$$v = \frac{\Delta x}{\Delta t}$$



How do we unlock the outstanding potential of event cameras?

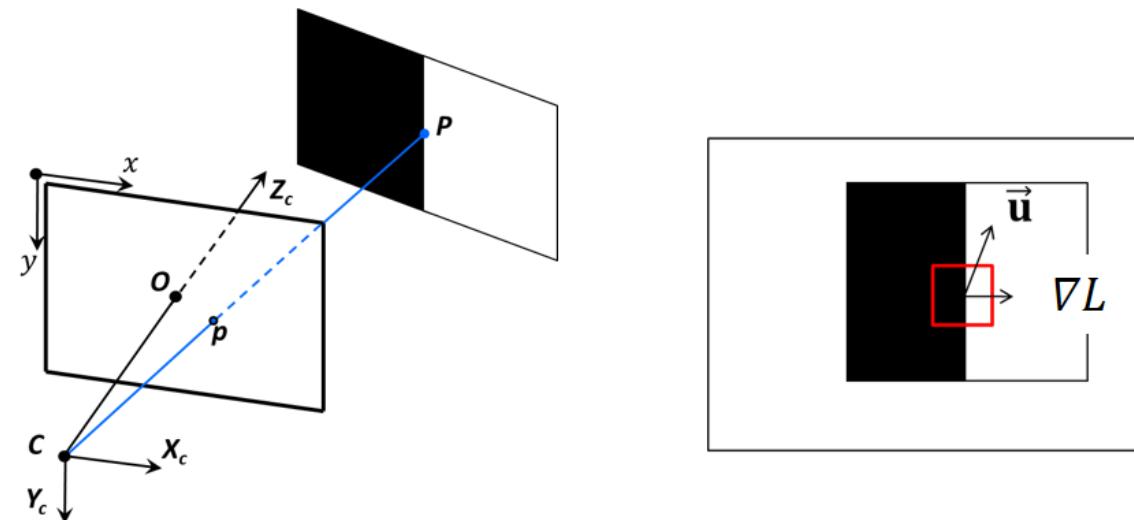
- Low latency
- High dynamic range
- No motion blur

1st order approximation of the Generative Event Model

- An event is generated when the following condition is satisfied:

$$\log I(x, y, t + \Delta t) - \log I(x, y, t) = \pm C$$

- For many applications, it is convenient to derive a 1st order approximation
- Let us define $L(x, y, t) = \log(I(x, y, t))$
- Consider a given pixel $p(x, y)$ with gradient $\nabla L(x, y)$ undergoing the motion $\mathbf{u} = (u, v)$ in pixels, induced by a moving 3D point P



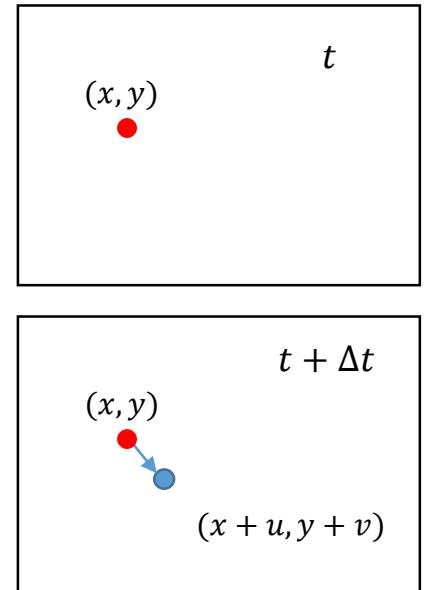
1st order approximation of the Generative Event Model

- Let's apply the **brightness constancy assumption**, which says that the intensity value of p before and after the motion must remain unchanged:

$$L(x, y, t) = L(x + u, y + v, t + \Delta t)$$

- By replacing the right-hand term with its 1st order approximation at $t + \Delta t$, we get:

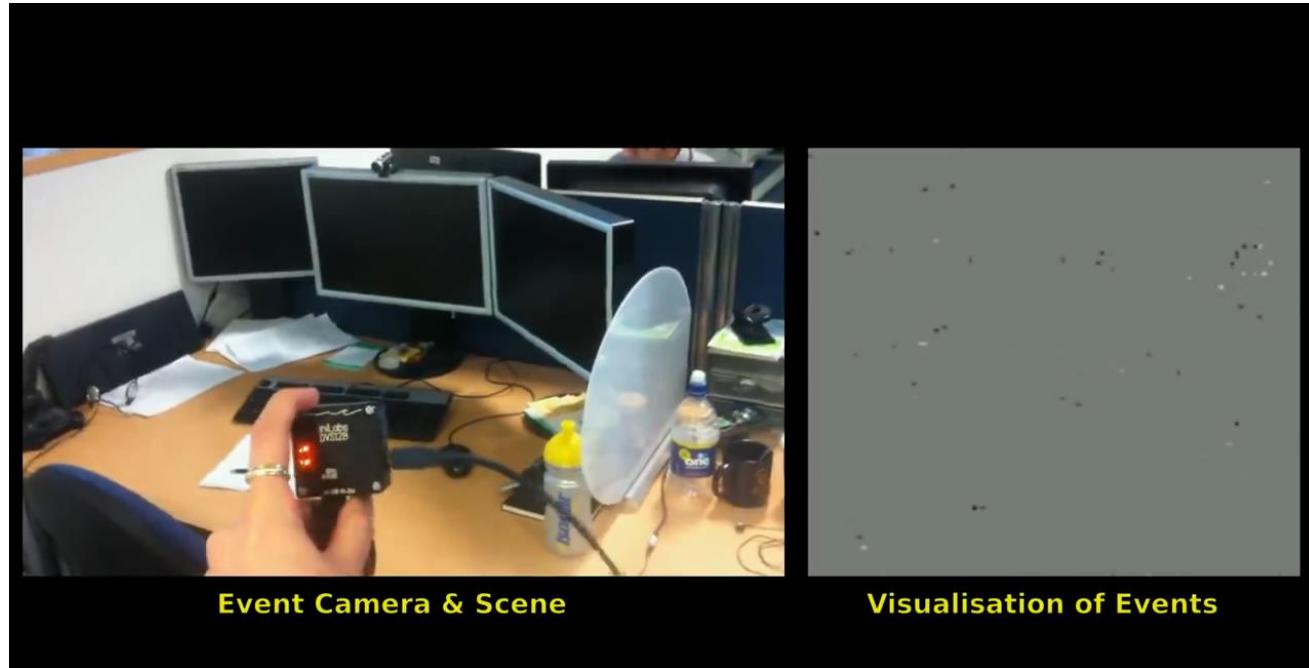
$$\begin{aligned} L(x, y, t) &= L(x, y, t + \Delta t) + \frac{\partial L}{\partial x} u + \frac{\partial L}{\partial y} v \\ \Rightarrow L(x, y, t + \Delta t) - L(x, y, t) &= -\frac{\partial L}{\partial x} u - \frac{\partial L}{\partial y} v \\ \Rightarrow \pm C &= -\nabla L \cdot \mathbf{u} \end{aligned}$$



- This formula shows that **maximum generation of events** (i.e., higher event rate) occurs when the **relative motion of the camera is perpendicular to the edge** and is **minimum when parallel to the edge**.

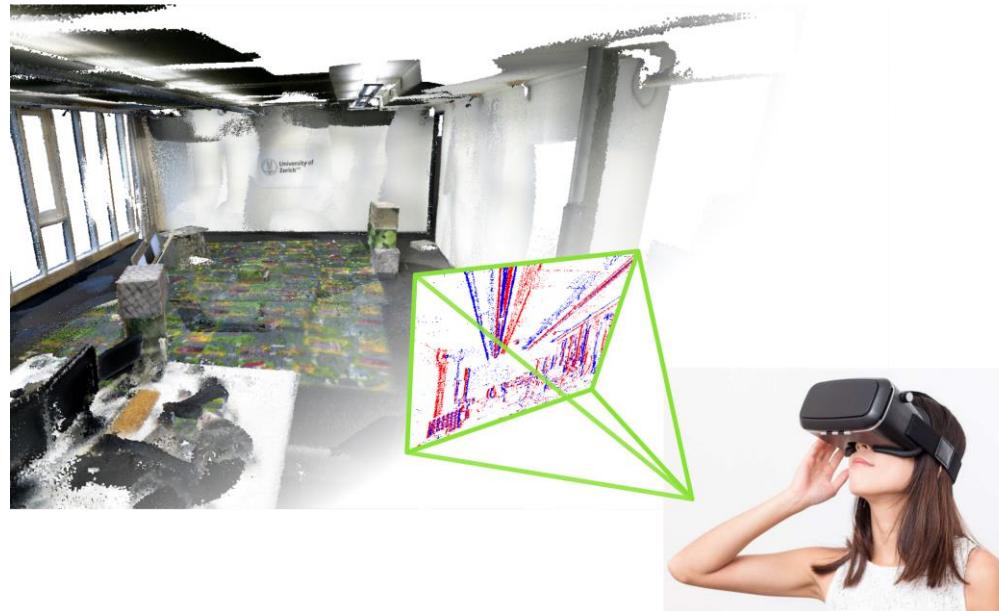
Application 1: Image Reconstruction from events

- Probabilistic **simultaneous gradient reconstruction and rotation estimation** from $\pm C = -\nabla L \cdot \mathbf{u}$
- Obtain **image intensity from gradient** via Poisson reconstruction
- The reconstructed image has **super-resolution and High Dynamic Range (HDR)**
- Can run in **real time on a GPU**



Application 2: 6DoF Tracking from Photometric Map

- Probabilistic **6DoF motion estimation** from $\pm C = -\nabla L \cdot \mathbf{u}$
- Assumes **photometric map** (x, y, z , grayscale Intensity) is **given**
- Useful for **VR/AR applications** (low-latency, HDR, no motion blur)
- Can run in **real time on a GPU**



Application 2: 6DoF Tracking from Photometric Map

Event camera

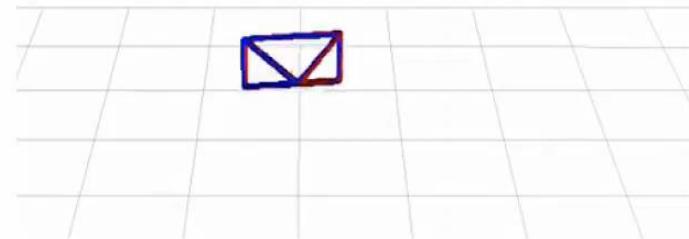


Standard camera

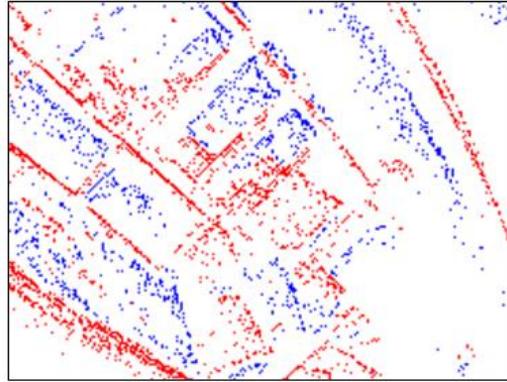


Motion estimation

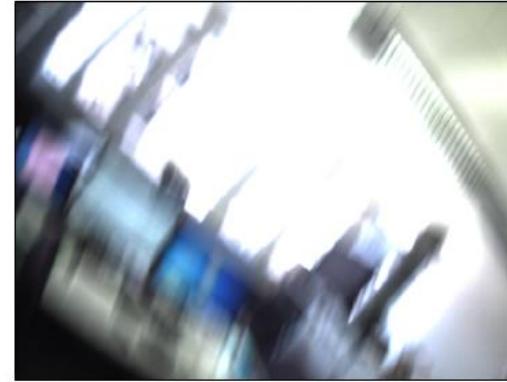
Event-based (EB)
Frame-based (FB)



Combining Standard Cameras with Event Cameras



Event Camera

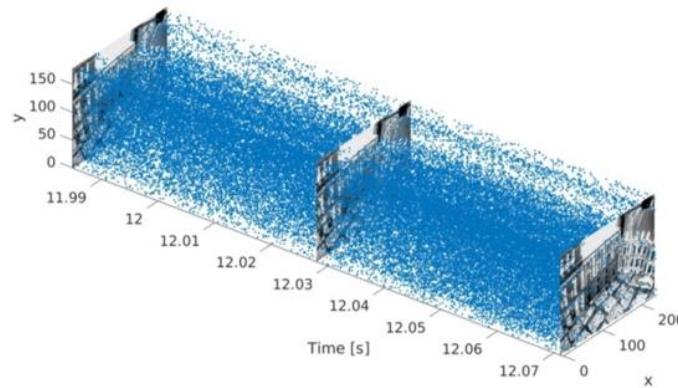


Standard Camera

Update rate	High (asynchronous): 1 MHz	Low (synchronous)
Dynamic Range	High (140 dB)	Low (60 dB)
Motion Blur	No	Yes
Static motion	No (event camera is a high pass filter)	Yes
Absolute intensity	No (but reconstructable up to a constant)	Yes
Maturity	< 10 years of research	> 60 years of research!

DAVIS sensor: Events + Images + IMU

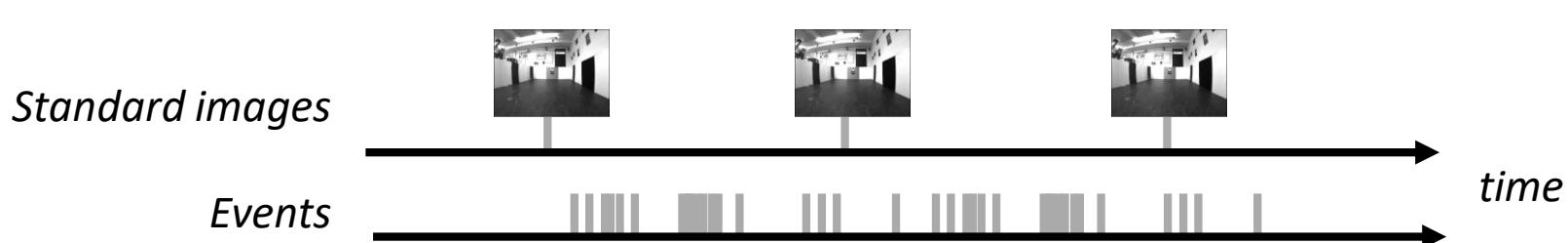
- Combines an **event and a standard camera** in the **same pixel array** (→ the same pixel can both trigger events and integrate light intensity).
- It also has an **IMU**



Spatio-temporal visualization
of the output of a DAVIS sensor



Temporal aggregation of events
overlaid on a DAVIS frame



Application 1: Deblurring a blurry video

- **Idea:** A **blurry image** can be regarded as the **integral of a sequence of latent images** during the exposure time, while the **events** indicate the **changes between the latent images**
- **Solution:** sharp image obtained by subtracting the double integral of event from input image

$$\log \text{ Input blur image} - \iint \text{ Input events} = \log \text{ Output sharp image}$$

The diagram illustrates the deblurring process. It shows three images: a blurry input image, a map of event data, and a sharp output image. The blurry input image is labeled "Input blur image". The event data map is labeled "Input events" and features a grid with red and blue dots representing event data. The sharp output image is labeled "Output sharp image". The mathematical notation indicates that the log of the input blur image minus the double integral of the input events results in the log of the output sharp image.

Application 1: Deblurring a blurry video

- **Idea:** A **blurry image** can be regarded as the **integral of a sequence of latent images** during the exposure time, while the **events** indicate the **changes between the latent images**
- **Solution:** sharp image obtained by subtracting the double integral of event from input image



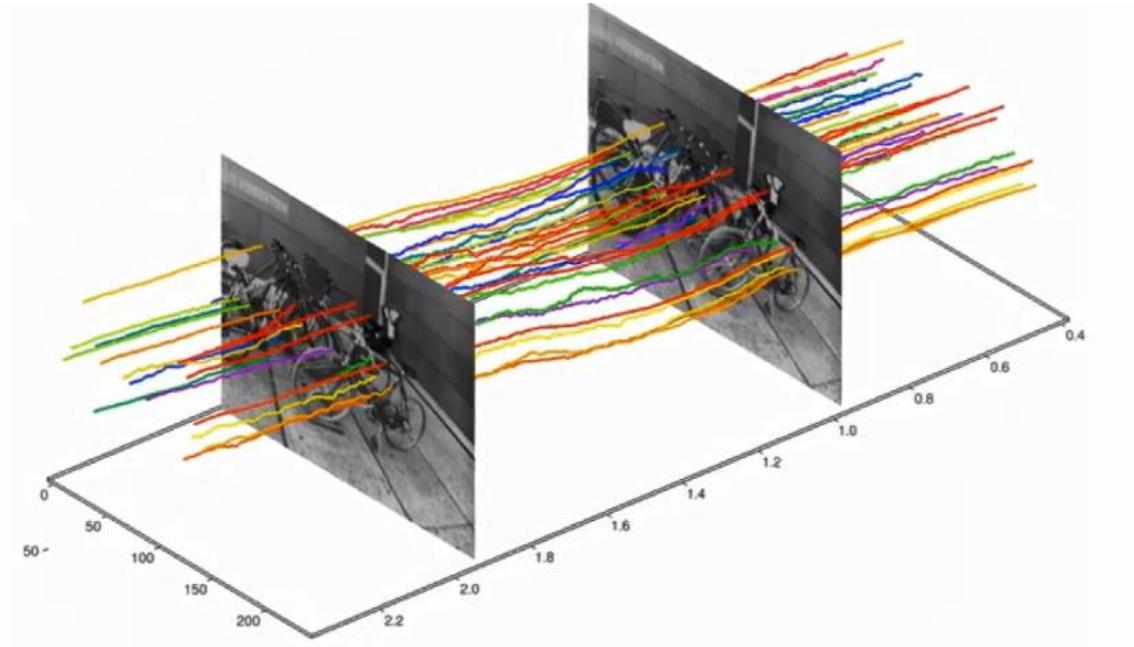
Input blur image



Output sharp video

Application 3: Event-based KLT Tracking

- **Goal:** Extract **features from standard frames** and track them using only **events** in the **blind time between two frames**
- Uses the event generation model via joint estimation of patch warping and optic flow



Source code: https://github.com/uzh-rpg/rpg_eklt

Recap

- All the approaches seen so far use the **generative event model**

$$\log I(x, y, t + \Delta t) - \log I(x, y, t) = \pm C$$

- or its 1st order approximation

$$\pm C = -\nabla L \cdot \mathbf{u}$$

which **requires knowledge of the contrast sensitivity C**

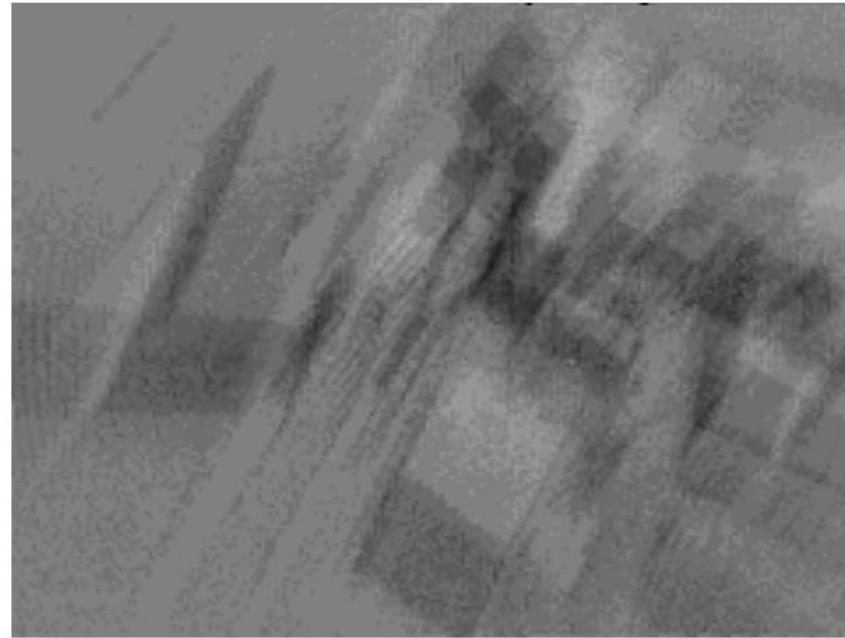
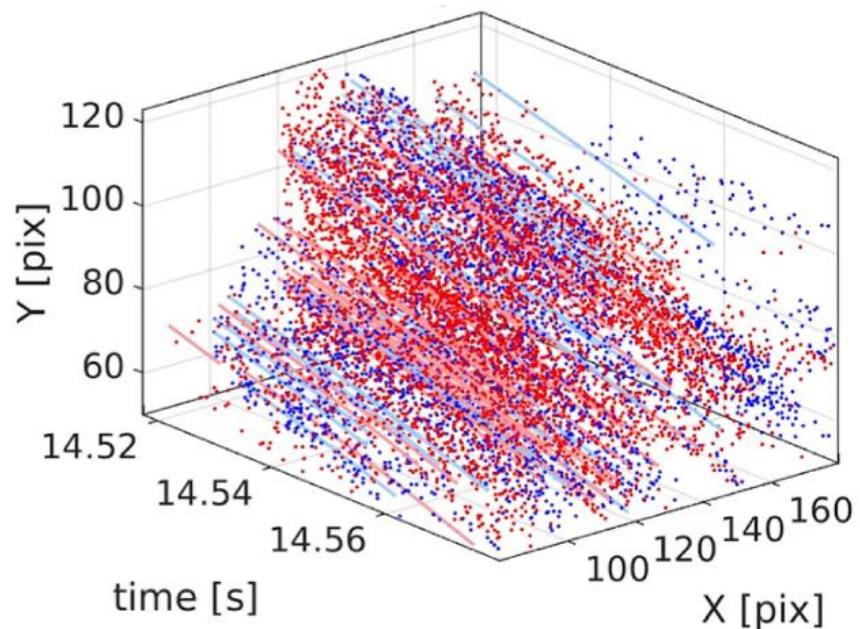
- Unfortunately, **C is scene dependent** and might **differ from pixel to pixel**
- **Alternative approach: Contrast maximization framework**

Contrast Maximization Framework

- Motion estimation
- 3D reconstruction
- SLAM
- Optical flow estimation
- Feature tracking
- Motion segmentation
- Unsupervised learning

Contrast Maximization Framework

Idea: Warp spatio-temporal volume of events to **maximize contrast** (e.g., sharpness) of the resulting image



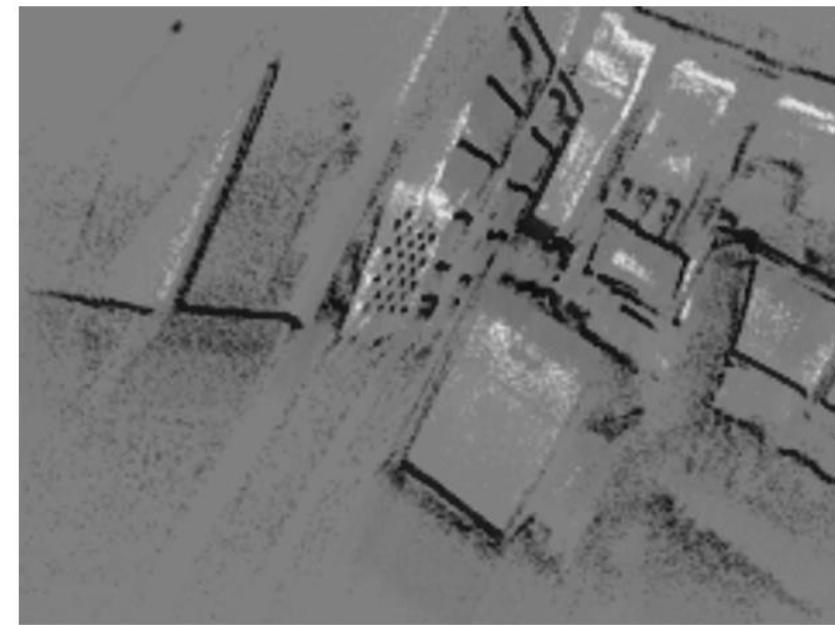
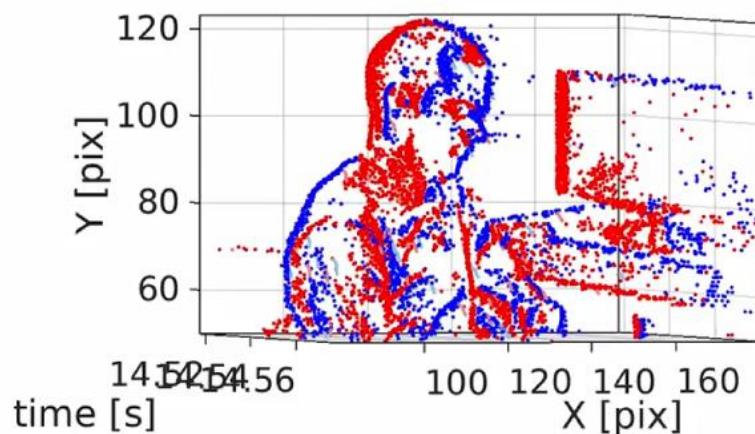
Aggregated image
without motion correction

Gallego, Rebecq, Scaramuzza, *A Unifying Contrast Maximization Framework for Event Cameras*, CVPR18, [PDF](#), [Video](#)

Gallego, Gehrig, Scaramuzza, *Focus Is All You Need: Loss Functions for Event-based Vision*, CVPR19, [PDF](#).

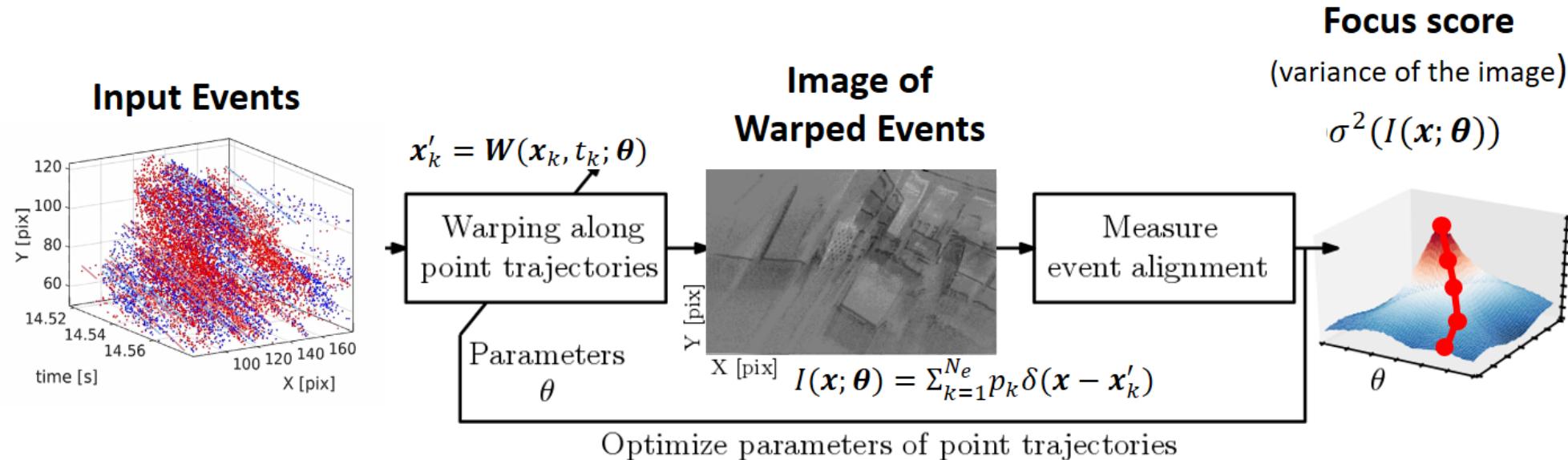
Contrast Maximization Framework

Idea: Warp spatio-temporal volume of events to **maximize contrast** (e.g., sharpness) of the resulting image



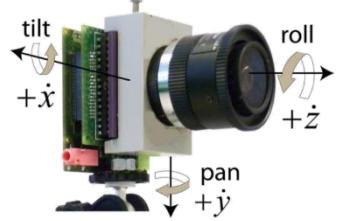
Aggregated image
with motion correction

Contrast Maximization Framework

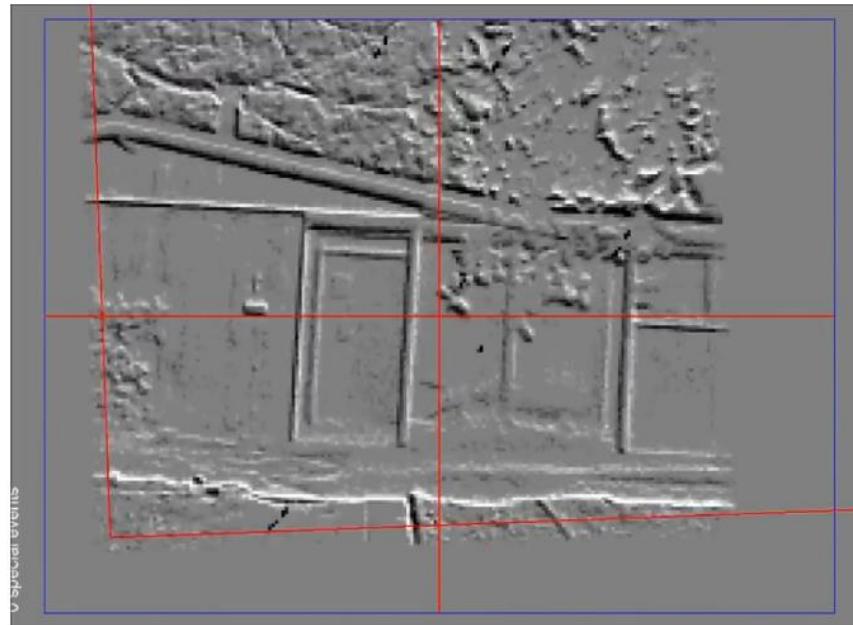


- $x'_k = W(x_k, t_k; \theta)$: This warps the (x, y) pixels coordinates of each event, not their time. Possible warps: roto-translation, affine, homography.
- $I(x; \theta) = \sum_{k=1}^{N_e} p_k \delta(x - x'_k)$: This builds a grayscale image, where the intensity of each pixel at the warped location (x', y') is equal to the summation of the polarity p (i.e., positive and negative events $(+1, -1)$)
- $\sigma^2(I(x; \theta))$: The assumption here is that if an image contains **high variance** then there is a wide **spread of responses, both edge-like and non-edge like**, representative of a normal, in-focus image. But if there is **very low variance**, then there is a tiny spread of responses, indicating there are very little edges in the image. As we know, the more an image is blurred, the less edges there are.

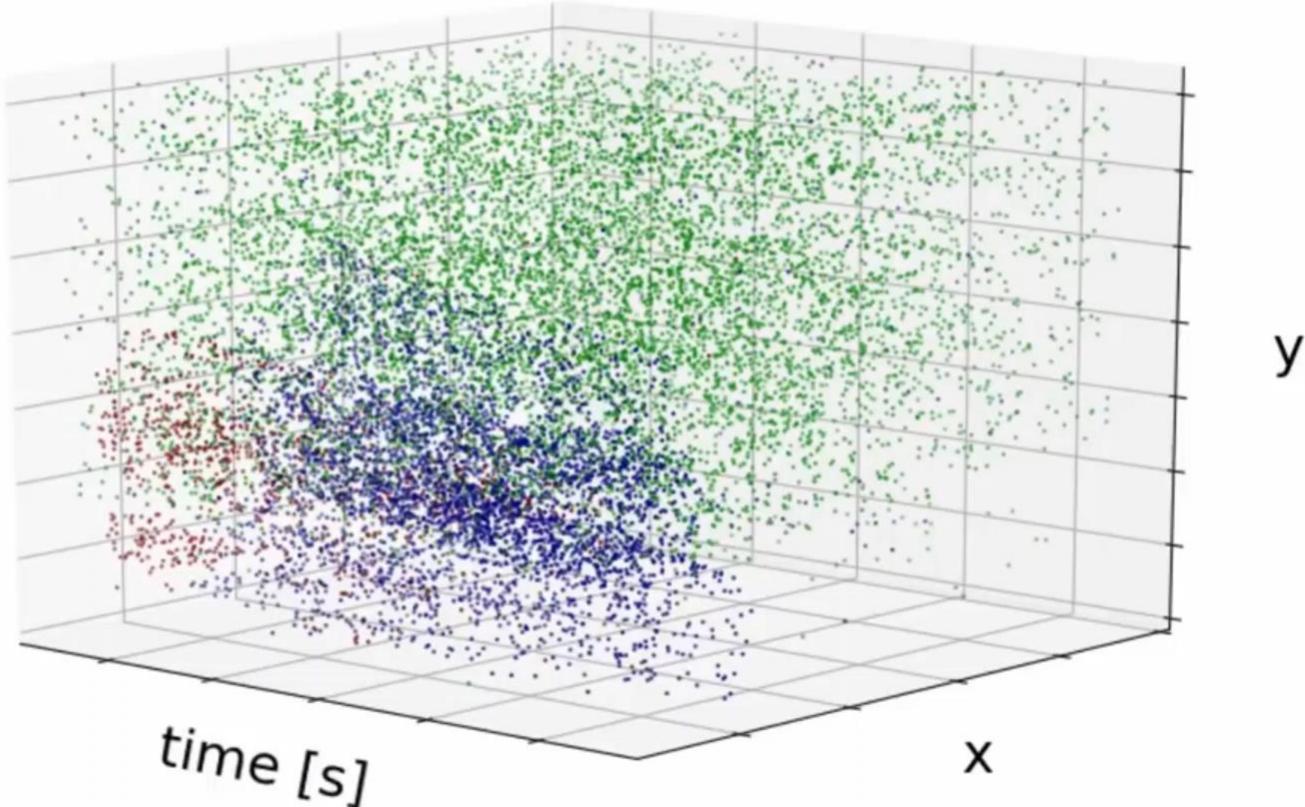
Application 1: Image Stabilization



- Goal: **Estimate rotational motion (3DoF)** of an event camera
- Can process millions of events per second in real time on a smartphone PC (e.g., OdroidXU4)
- Works up to over ~1,000 deg/s



Application 2: Motion Segmentation



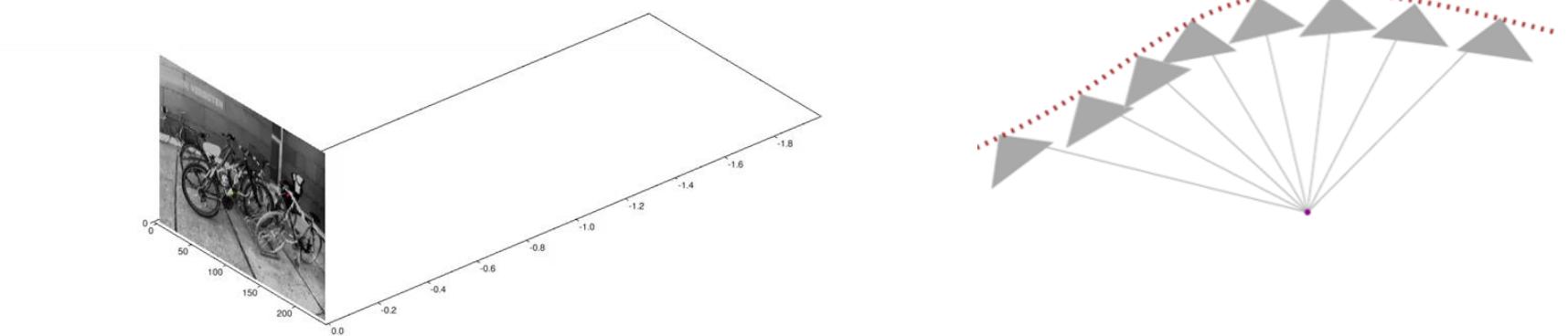
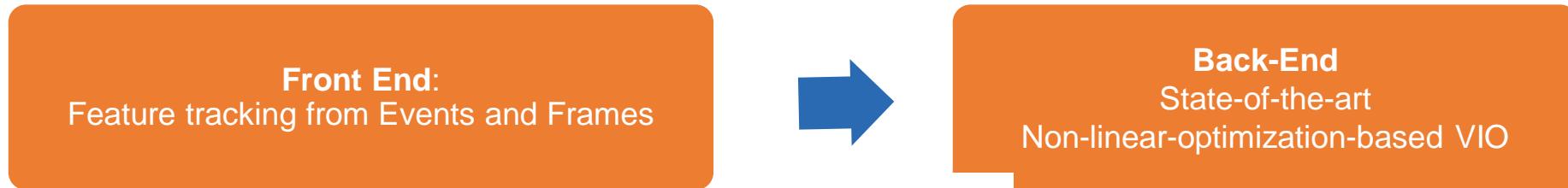
Application 3: Dynamic Obstacle Avoidance

- Works with relative speeds of up to **10 m/s**
- Perception **latency: 3.5 ms**



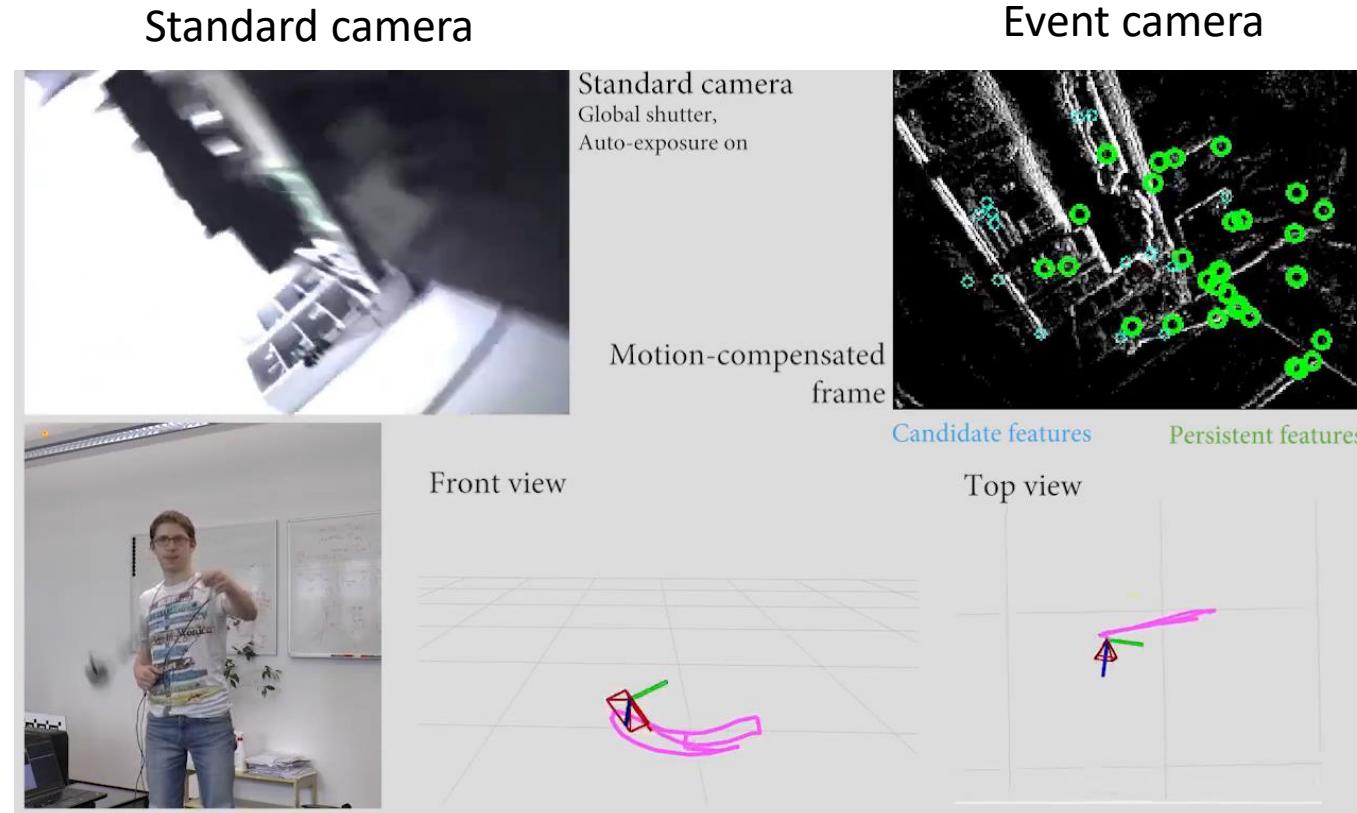
Application 4: “Ultimate SLAM”

Goal: combining **events**, **images**, and **IMU** for robust visual SLAM in HDR and high speed scenarios



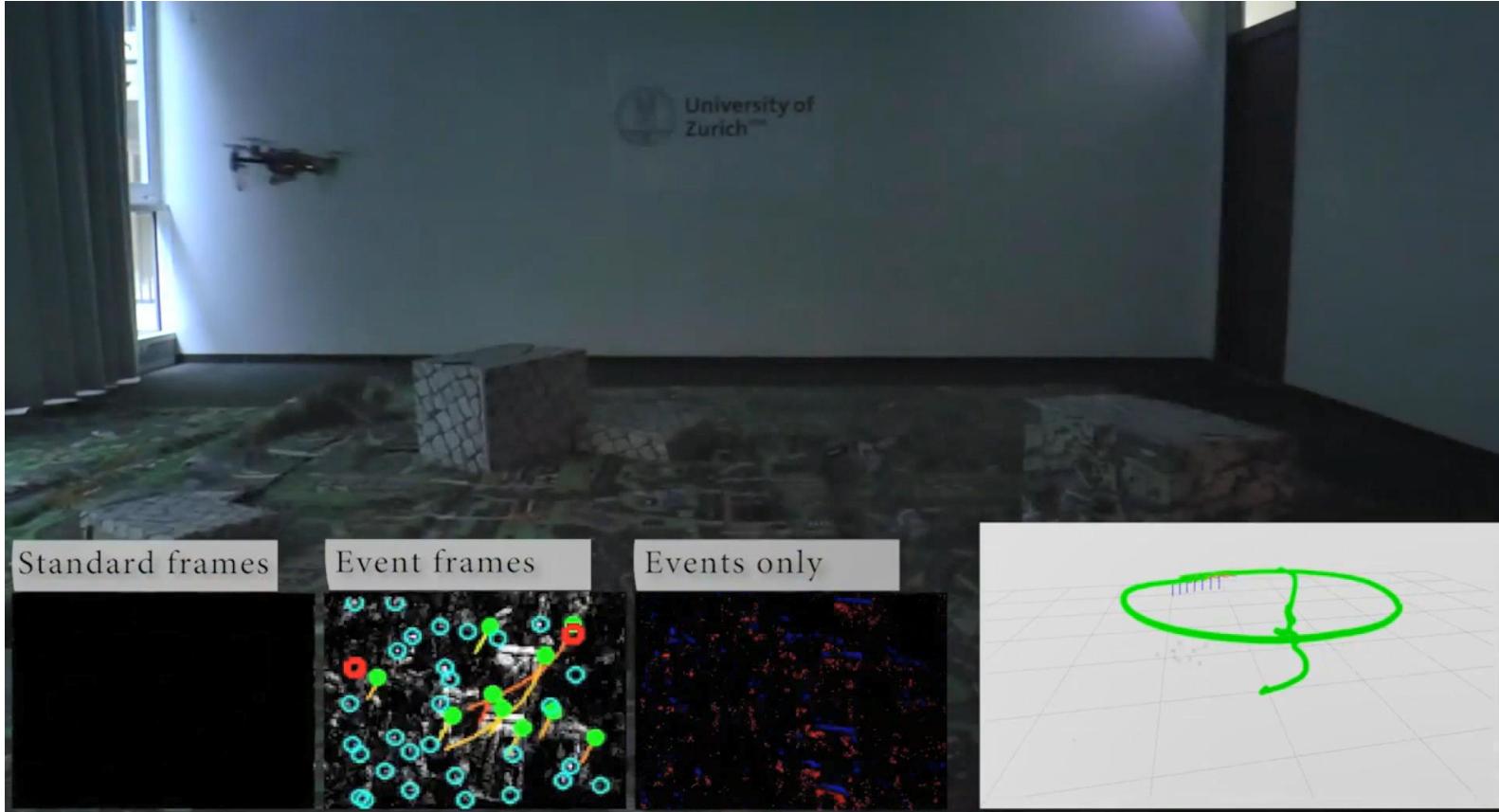
Application 4: “Ultimate SLAM”

- 85% accuracy gain over standard VIO in **HDR and high speed scenarios**



Application 5: Autonomous Navigation in Low Light

- UltimateSLAM running on board (CPU: Odroid XU4)

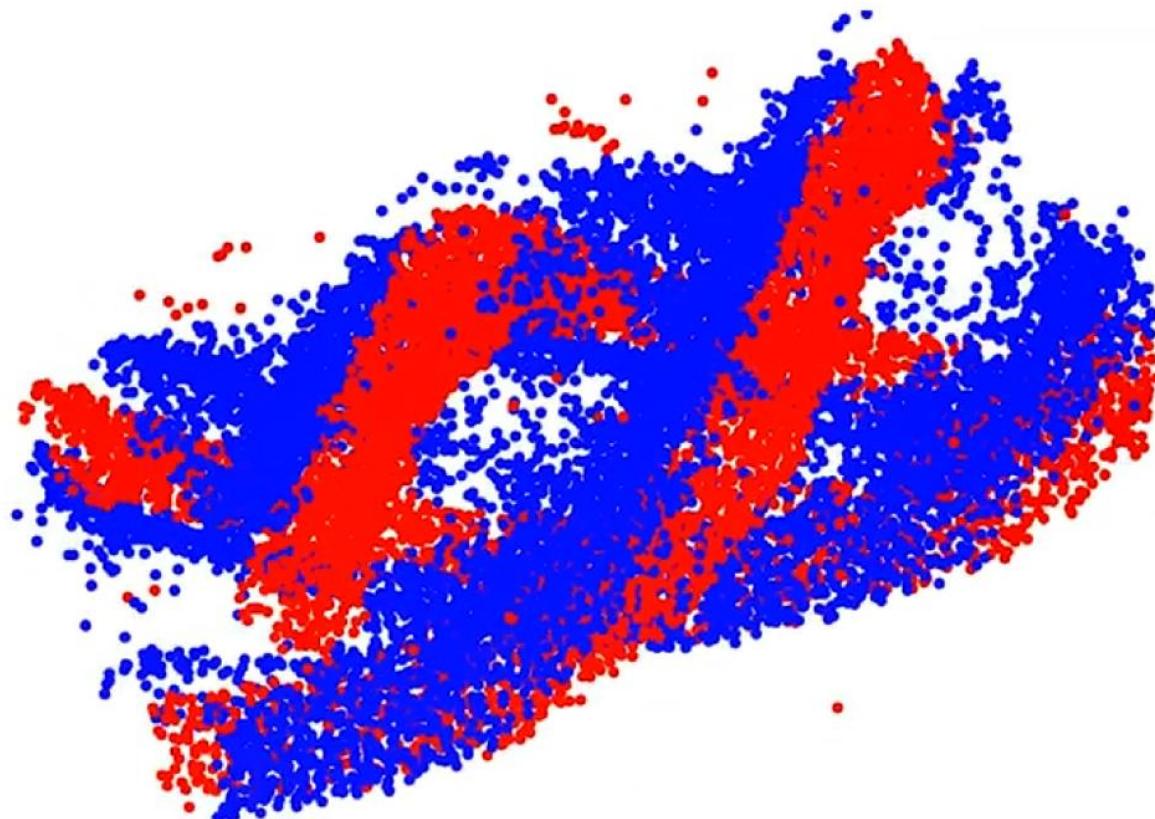


Learning with Event Cameras

- Approaches using synchronous, Artificial Neural Networks (ANNs) designed for standard images
- Asynchronous, Sparse ANNs
- Approaches using asynchronous, Spiking neural networks (SNNs)

Input representation

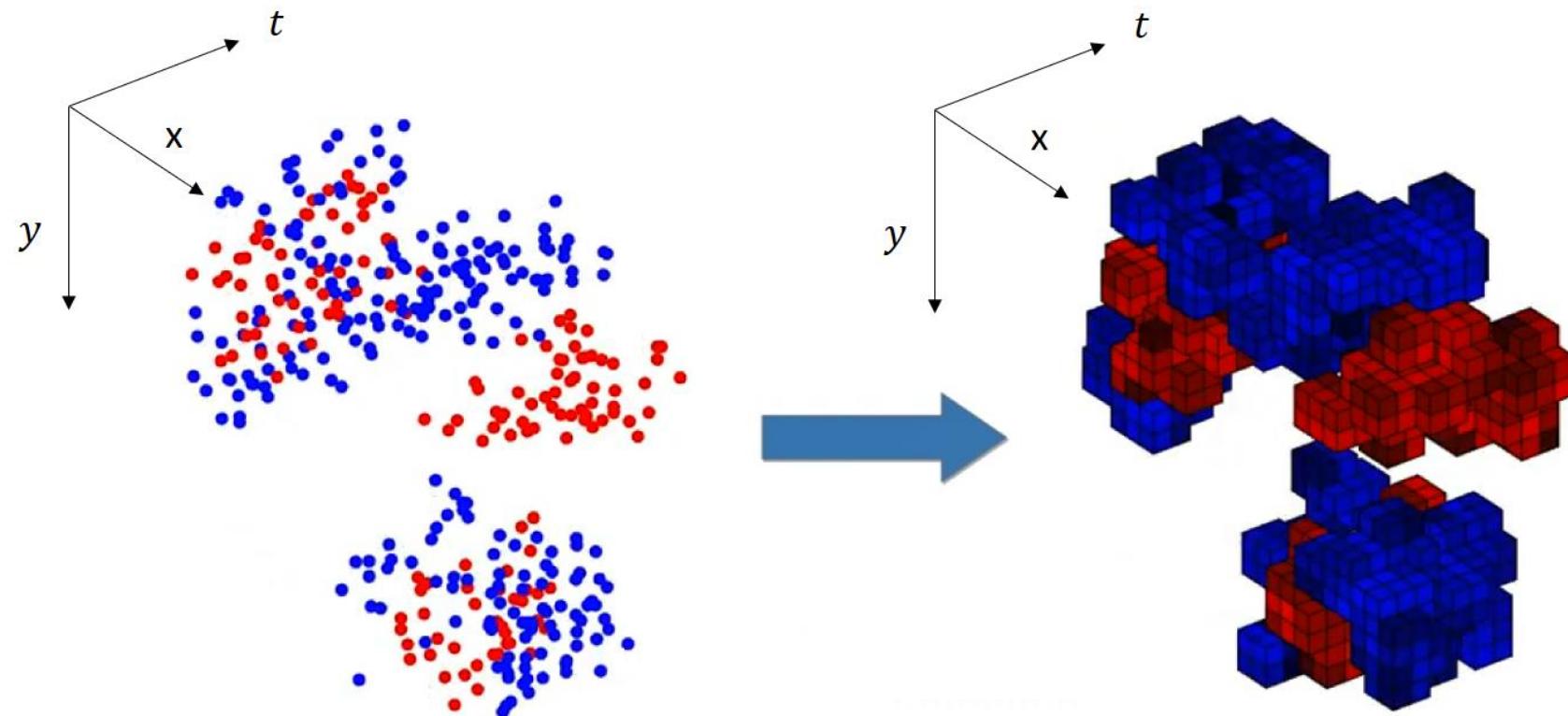
How do we pass sparse events into a convolutional neural network designed for standard images?



Video from [here](#)

Input representation

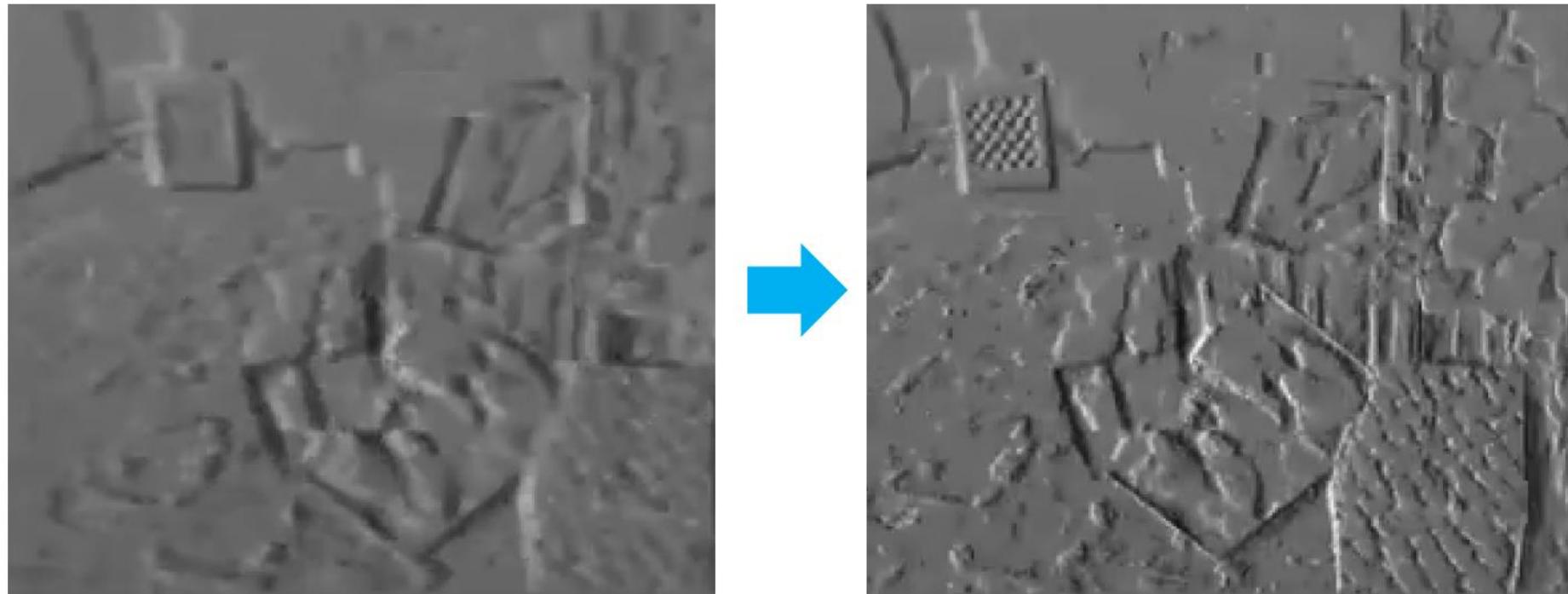
Represent events in space-time into a 3D voxel grid (x, y, t): each voxel contains sum of positive and negative events falling within the voxel



Video from [here](#)

Contrast as Loss Function for Unsupervised Learning

- Idea: maximize sharpness of the aggregated event image



Video from [here](#)

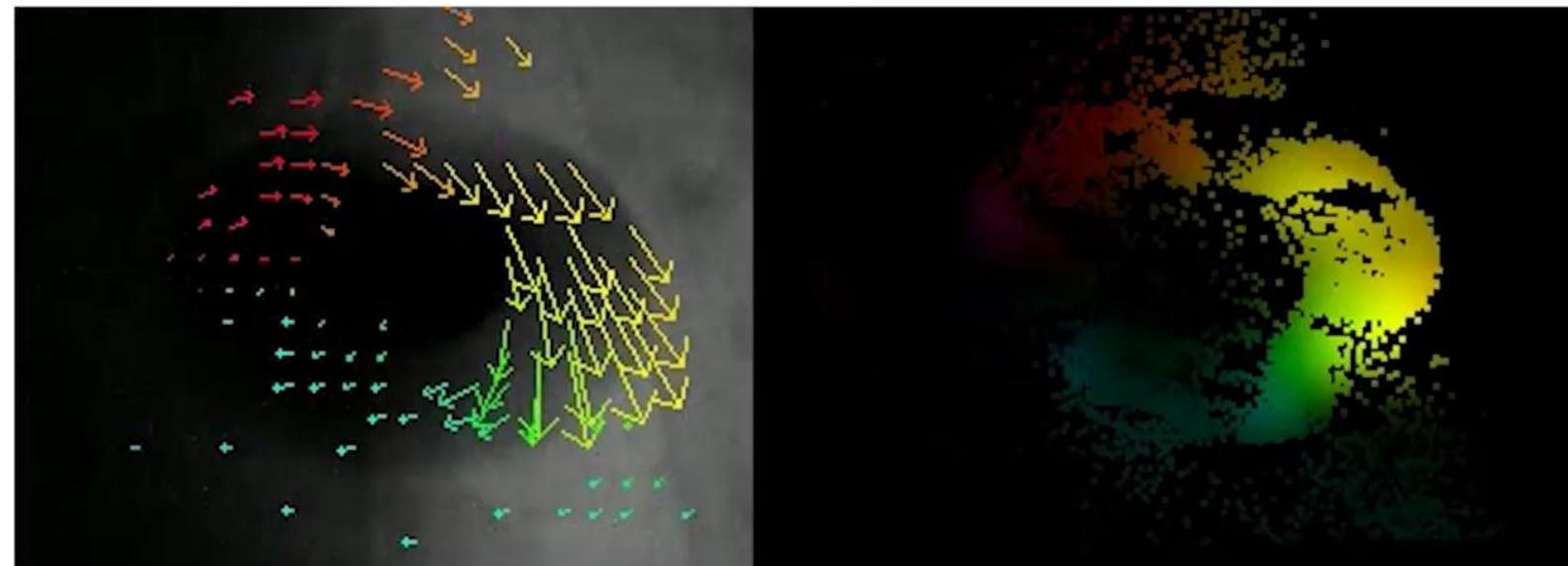
Zhu, Yuan, Chaney, Daniilidis, *Unsupervised Event-based Learning of Optical Flow, Depth and Egomotion*, CVPR 19. [PDF](#).

Gallego, Gehrig, Scaramuzza, *Focus Is All You Need: Loss Functions for Event-based Vision*, CVPR19, [PDF](#).

Application1: Unsupervised Learning of Optical Flow

- Idea: maximize sharpness of the aggregated event image

Fidget Spinner w/ Challenging Lighting



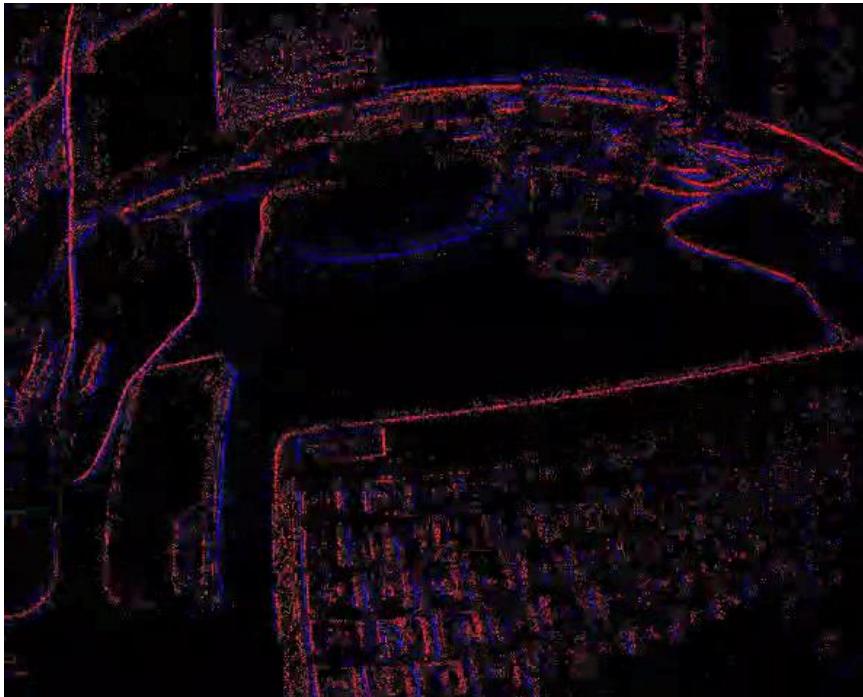
Grayscale Image w/ Sparse Flow Quiver

Dense Flow Output

1x realtime

Application 2: Image Reconstruction from Events

Events



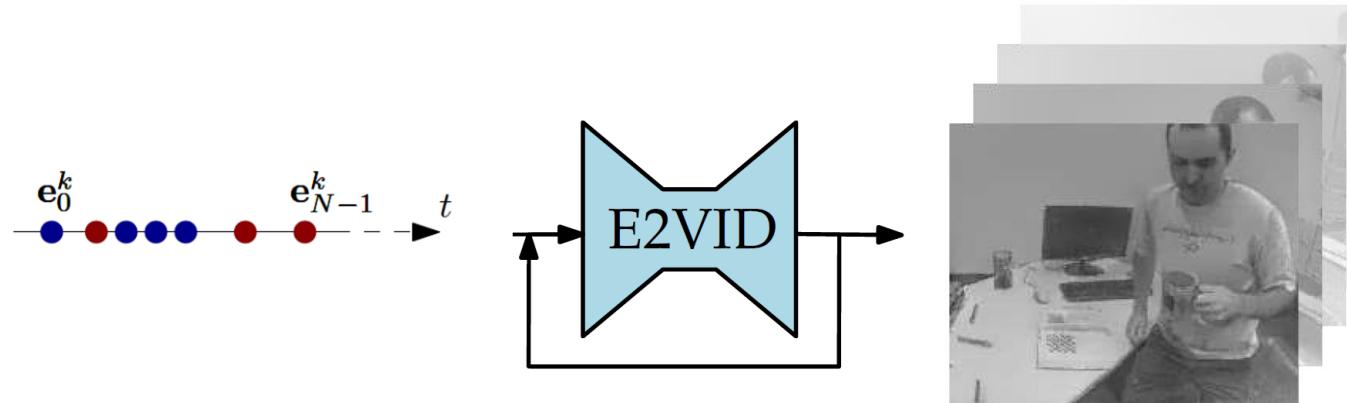
Reconstructed image from events



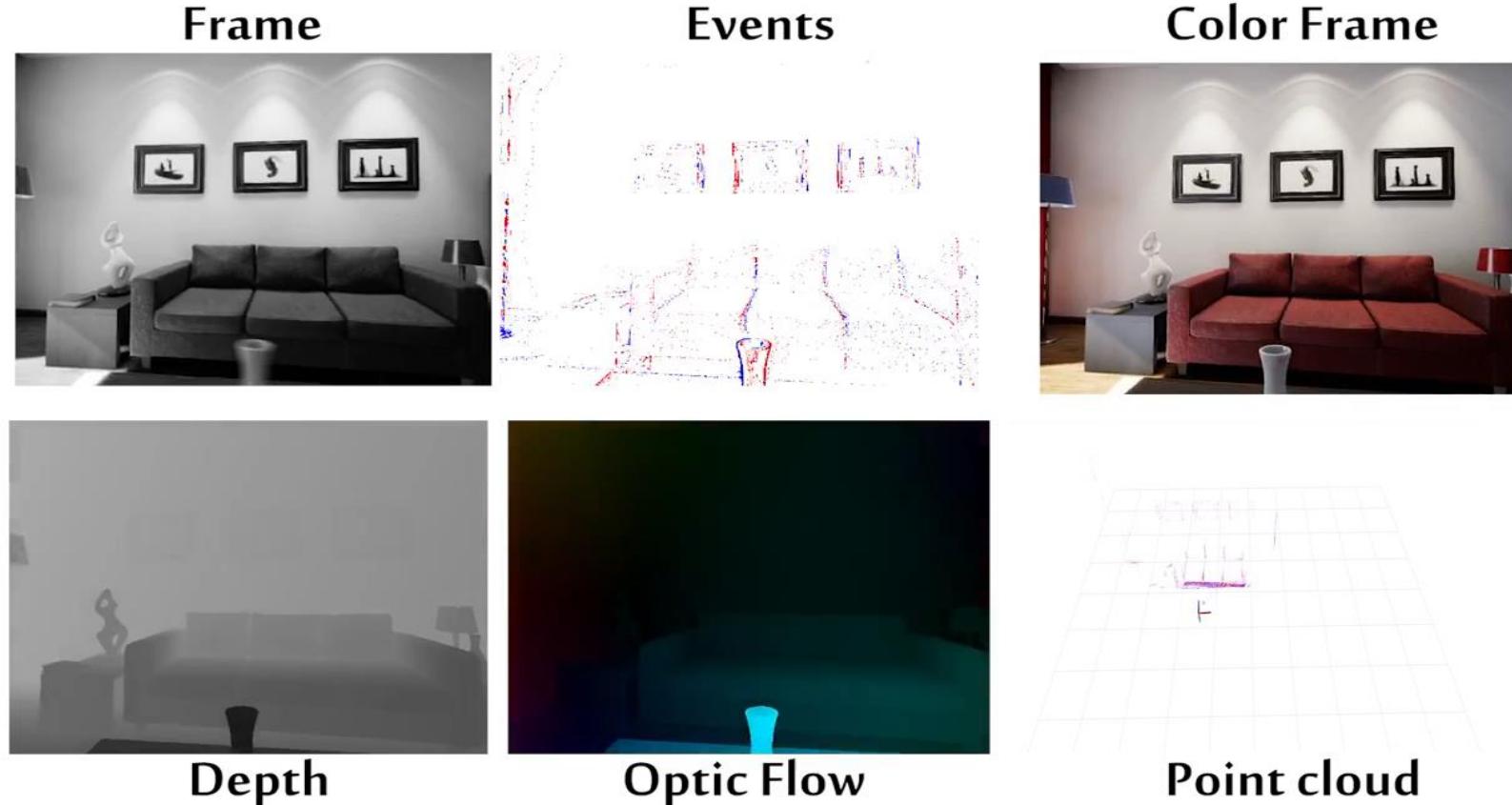
Code & datasets: https://github.com/uzh-rpg/rpg_e2vid

Overview

- **Recurrent neural network** (main module: Unet)
- Input: sequences of *event tensors* (3D spatio-temporal volumes of events^[3])
- **Trained in simulation only**, without seeing a single real image
- To improve robustness **we randomize the contrast sensitivity** during simulation.
- Event camera simulator (ESIM): <http://rpg.ifi.uzh.ch/esim.html>



ESIM: Event Camera Simulator



Open Source: <http://rpg.ifi.uzh.ch/esim.html>

Bullet shot by a gun (1,300 km/h)

Recall: trained in simulation only!



Huawei P20 Pro (240 FPS)



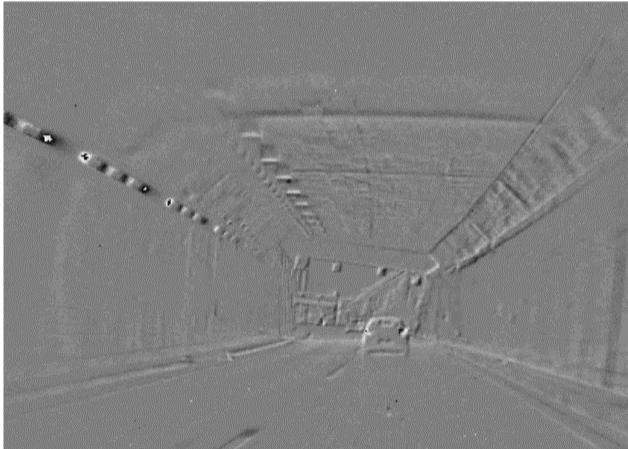
Our reconstruction (5400 FPS)

Code & datasets: https://github.com/uzh-rpg/rpg_e2vid 100 x slow motion

HDR Video: Driving out of a tunnel

Recall: trained in simulation only!

Driving out of a tunnel



Events



Our reconstruction



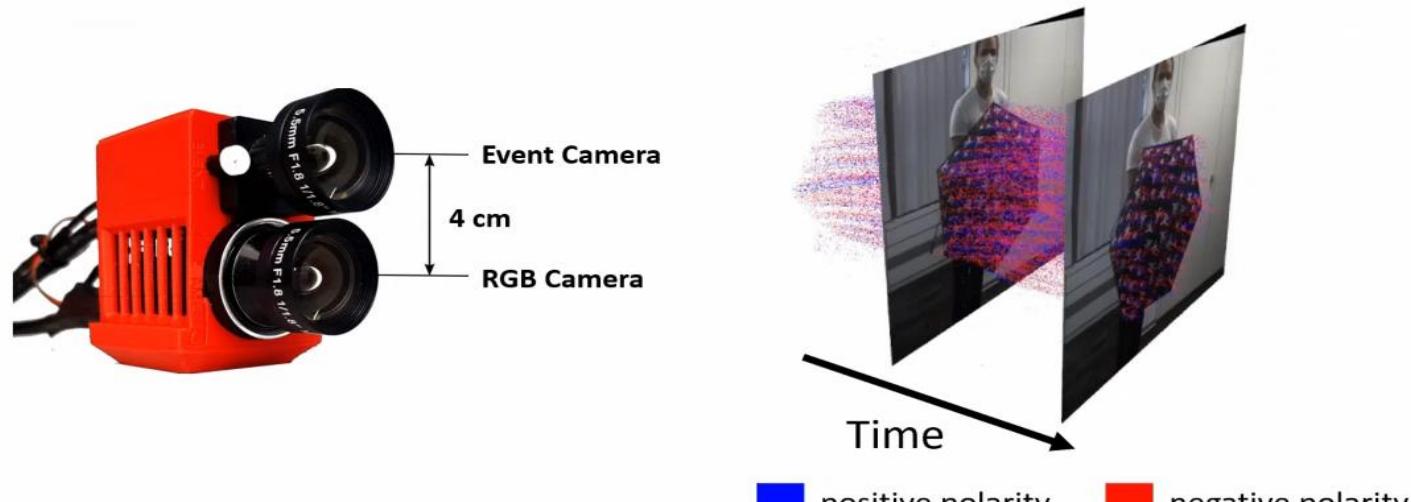
Phone camera

Code & datasets: https://github.com/uzh-rpg/rpg_e2vid

Application 3: Slow Motion Video

- We can combine an event camera with an HD RG camera
- We use events to **upsample low-framerate video by over 50 times with only 1/40th of the memory footprint!**

It does this by leveraging event cameras which provide a compressed stream of visual information in the blind-time between frames.

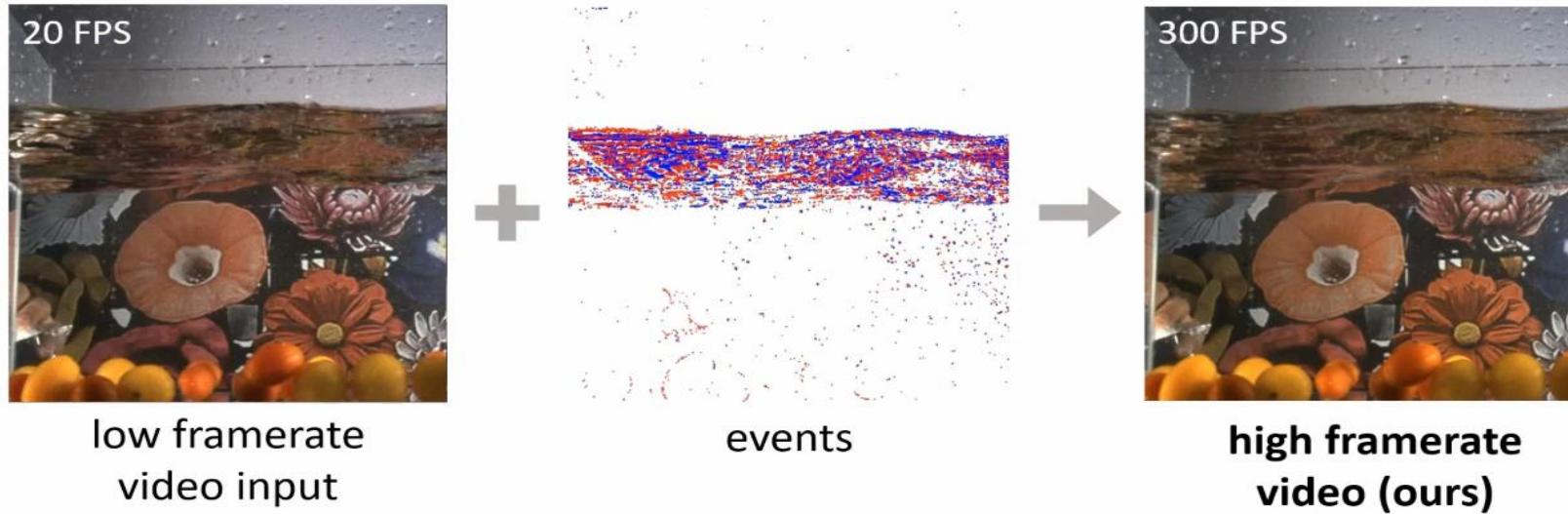


Code & Datasets: <http://rpg.ifi.uzh.ch/timelens>

Application 3: Slow Motion Video

- We can combine an event camera with an HD RG camera
- We use events to **upsample low-framerate video by over 50 times with only 1/40th of the memory footprint!**

stream of visual information in the blind-time between frames.



Code & Datasets: <http://rpg.ifi.uzh.ch/timelens>

Application 3: Slow Motion Video

- We can combine an event camera with an HD RG camera
- We use events to **upsample low-framerate video by over 50 times with only 1/40th of the memory footprint!**



low framerate video input



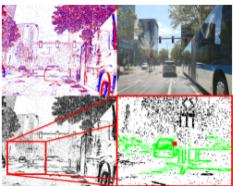
Time Lens (this work)

Code & Datasets: <http://rpg.ifi.uzh.ch/timelens>

We have student projects on event cameras

http://rpg.ifi.uzh.ch/student_projects.php

Asynchronous Processing for Event-based Deep Learning - [Available](#)



Description: Event cameras such as the Dynamic Vision Sensor (DVS) are recent sensors with large potential for high-speed and high dynamic range robotic applications. Since their output is sparse traditional algorithms, which are designed for dense inputs such as frames, are not well suited. The goal of this project is explore ways to adapt existing deep learning algorithms to handle sparse asynchronous data from events. Applicants should have experience in C++ and python deep learning frameworks (tensorflow or pytorch), and have a strong background in computer vision.

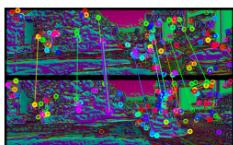
Goal: The goal of this project is explore ways to adapt existing deep learning algorithms to handle sparse asynchronous data from events.

Contact Details: Daniel Gehrig (dgehrig@ifi.uzh.ch)

Thesis Type: Semester Project / Master Thesis

[See project on SIROP](#)

Data-driven Keypoint Extractor for Event Data - [Available](#)



Description: Neuromorphic cameras exhibit several amazing properties such as robustness to HDR scenes, high-temporal resolution, and low power consumption. Thanks to these characteristics, event cameras are applied for camera pose estimation for fast motions in challenging scenes. A common technique for camera pose estimation is the extraction and tracking of keypoints on the camera plane. In the case of event cameras, most existing keypoint extraction methods are handcrafted manually. As a new promising direction, this project tackles the keypoint extraction in a data-driven fashion based on recent advances in frame-based keypoint extractors.

Goal: The project aims to develop a data-driven keypoint extractor, which computes interest points in event data. Based on the current advances of learned keypoint extractors for traditional frames, the approach will leverage neural network architectures to extract and describe keypoints in an event stream. The student should have prior programming experience in a deep learning framework and completed at least one course in computer vision.

Contact Details: Contact Details: Nico Messikommer [nmessi@ifi.uzh.ch], Mathias Gehrig [mgehrig@ifi.uzh.ch]

Thesis Type: Semester Project / Master Thesis

[See project on SIROP](#)

Designing a New Event Camera with Events and Images - [Available](#)



Description: Event cameras such as the Dynamic Vision Sensor (DVS) are recent sensors with a lot of potential for high-speed and high dynamic range robotic applications. They have been successfully applied in many applications, such as high speed video and high speed visual odometry. Due to their high speed and

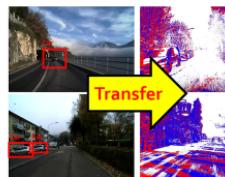
Goal: The goal of this project is to design a new event camera that combines events and standard images.

Contact Details: Daniel Gehrig (dgehrig@ifi.uzh.ch), Mathias Gehrig (mgehrig@ifi.uzh.ch)

Thesis Type: Semester Project / Master Thesis

[See project on SIROP](#)

Domain Transfer between Events and Frames - [Available](#)



Description: During the last years, a vast collection of frame-based datasets was collected for countless tasks. In comparison, event-based datasets represent only a tiny fraction of the available datasets. Thus, it is highly promising to use labelled frame datasets to train event-based networks as current data-driven approaches heavily rely on labelled data.

Goal: In this project, the student extends current advances from the UDA literature for traditional frames to event data in order to transfer multiple tasks from frames to events. The approach should be validated on several tasks (segmentation, object detection, etc.) in challenging environments (night, high-dynamic scenes) to highlight the benefits of event cameras. As several deep learning methods are used as tools for the task transfer, a strong background in deep learning is required. If you are interested, we are happy to provide more details.

Contact Details: Nico Messikommer [nmessi@ifi.uzh.ch], Daniel Gehrig (dgehrig@ifi.uzh.ch)

Thesis Type: Semester Project / Master Thesis

[See project on SIROP](#)

Understanding Check

Are you able to answer the following questions?

- What is an event camera and how does it work?
- What are its pros and cons vs. standard cameras?
- Can we apply standard camera calibration techniques?
- How can we compute optical flow with a DVS?
- Could you intuitively explain why we can reconstruct the intensity?
- What is the generative model of an event camera (formula). Can you derive its 1st order approximation?
- What is a DAVIS sensor?
- What is the focus maximization framework and how does it work? What is its advantage compared with the generative model?