



# **EE5110/EE6110**

## **Special Topics in Automation and Control**

### **Segment B**

**From Frames to Events:  
Theory and Applications of Event-based Vision**

Jiang Rui, Ph.D., BEng  
[rui.j@nus.edu.sg](mailto:rui.j@nus.edu.sg)

Adjunct Lecturer  
Department of Electrical and Computer Engineering

Semester 1, AY2021/2022

# Course Structure

- **Lecture 1 – Event-based Vision System**
  - Conventional machine vision systems
  - Event-based sensors and cameras
  - Event data representation and compression
  - Event data simulation and dataset

# Course Structure

- **Lecture 2 – Algorithms and Applications**
  - Detection and tracking
  - Optical flow estimation
  - Pose estimation and SLAM
  - Image reconstruction
  - Motion segmentation

# Learning Objectives

- In this lecture, you are expected to:
  1. Describe and explain the basic concepts and methods in event-based computer vision algorithms and applications;
  2. Understand the pros and cons by comparing the event-based methods with the frame-based counterparts.

# Event-based Detection and Tracking

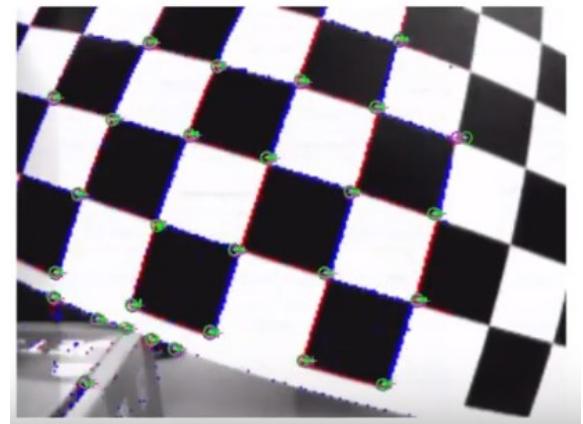
- Connection between detection and tracking

- **Feature detection:**

- Extract those robust, representative events ("features").

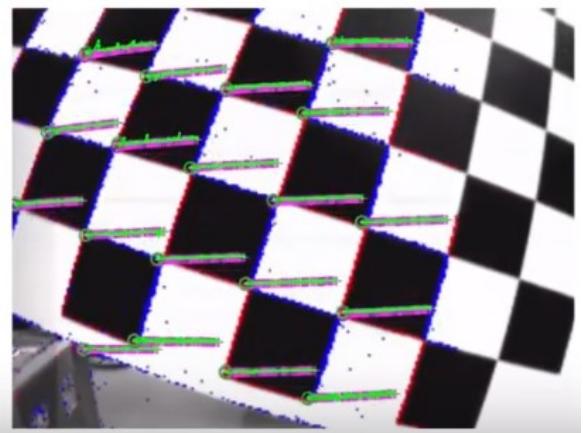
- **Object detection:**

- Identify and find the interested objects.



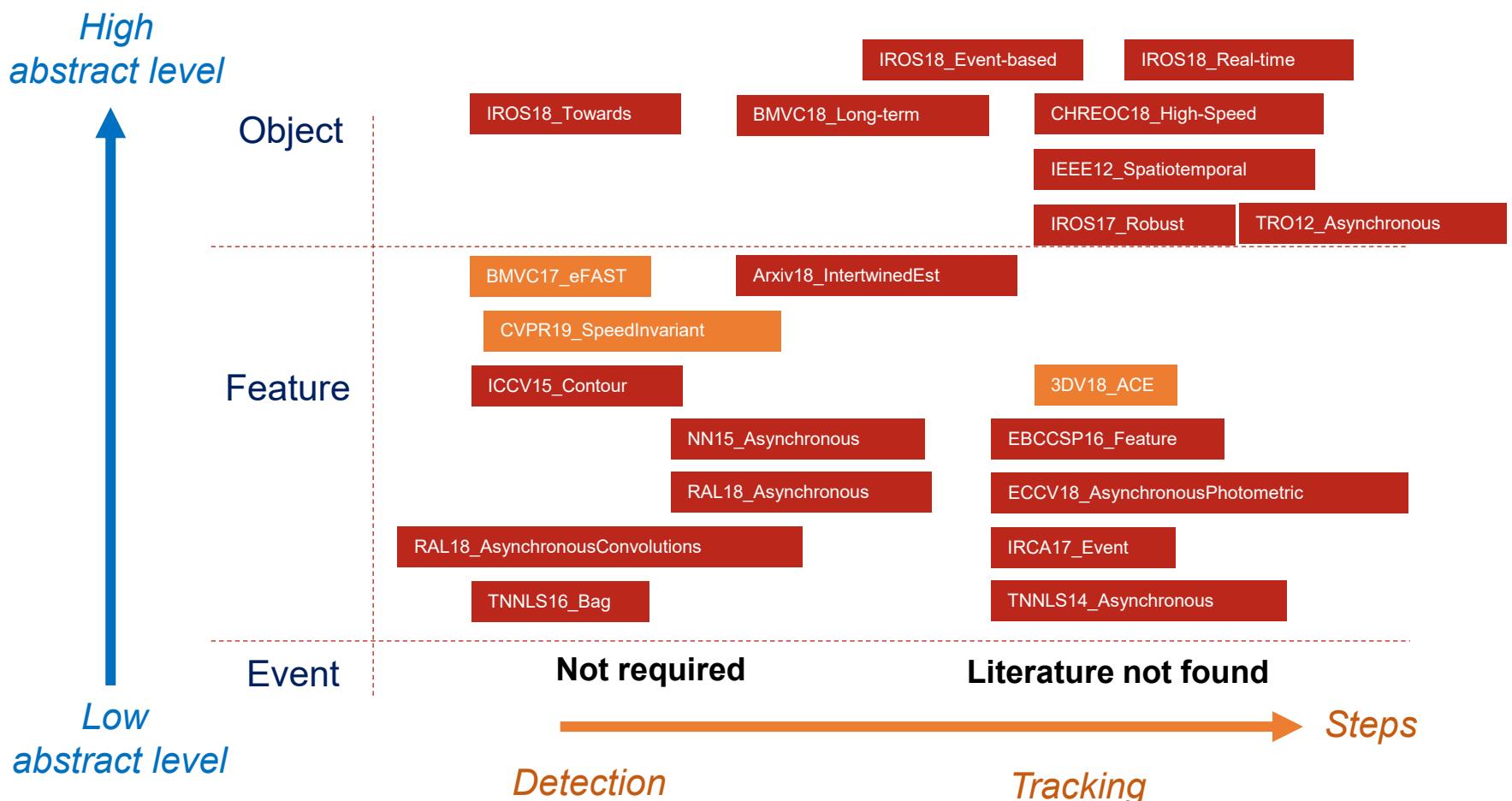
- **Feature/object tracking:**

- Establish correspondences between features or objects.



# Event-based Detection and Tracking

- Connection between detection and tracking



\*Above references are for demonstration only, and do not reflect the recent progress of this field.

# Event-based Detection and Tracking

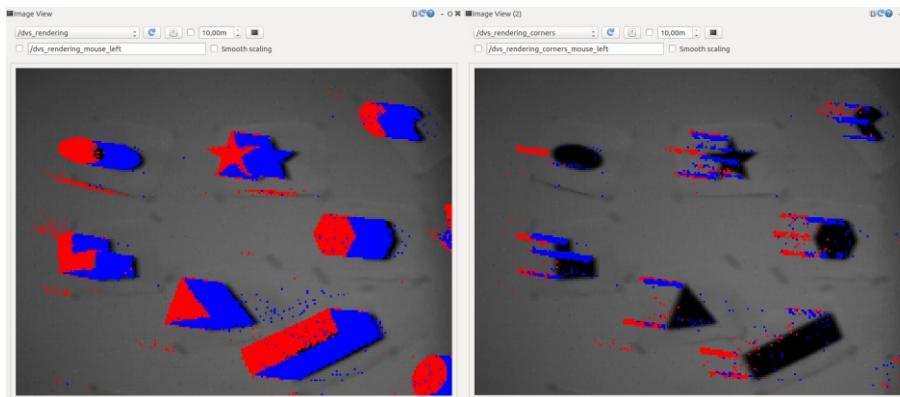
- References' details

- 3DV18\_ACE\_An Efficient Asynchronous Corner Tracker for Event Cameras
- Arxiv18\_Event-Based Features Selection and Tracking from Intertwined Estimation of Velocity and Generative Contours
- BMVC17\_Fast Event-based Corner Detection
- BMVC18\_Long-term object tracking with a moving event camera
- CHREOC18\_High-Speed Object Tracking with Dynamic Vision Sensor
- CVPR19\_Speed Invariant Time Surface for Learning to Detect Corner Points with Event-Based Cameras
- EBCCSP16\_Feature detection and tracking with the dynamic and active-pixel vision sensor (DAVIS)
- ECCV18\_Asynchronous\_Photometric Feature Tracking using Events and Frames
- ICCV15\_Contour Detection and Characterization for Asynchronous Event Sensors
- IEEE12\_Spatiotemporal Multiple Persons Tracking Using Dynamic Vision Sensor
- IRCA17\_Event-based Feature Tracking with Probabilistic Data Association
- IROS17\_Robust Visual Tracking with a Freely-moving Event Camera
- IROS18\_Event-based Moving Object Detection and Tracking
- IROS18\_Real-time clustering and multi-target tracking using event-based sensors
- IROS18\_Towards Event-Driven Object Detection with Off-the-Shelf Deep Learning
- NN15\_Asynchronous event-based corner detection and matching
- RAL18\_Asynchronous Corner Detection and Tracking for Event Cameras in Real Time
- RAL18\_Asynchronous Spatial Image Convolutions for Event Cameras
- TNNLS14\_Asynchronous Event-Based Multi-kernel Algorithm for High Speed Visual Features Tracking
- TNNLS16\_Bag of Events\_An Efficient Probability-Based Feature Extraction Method for AER Image Sensors
- TR012\_Asynchronous Event-based Visual Shape Tracking for Stable Haptic Feedback in Microrobotics

\*Above references are for demonstration only, and do not reflect the recent progress of this field.

# Paper: Fast Event-based Corner Detection

- “Event camera generates much less data than conventional camera.”  
**Is this always true?**
- Not always! **Sometimes, more data from event camera...**
- That’s why we need feature detection.
- Paper highlights:
  - A fast method for corner detection in an event stream
  - Efficient computing using only comparison operations

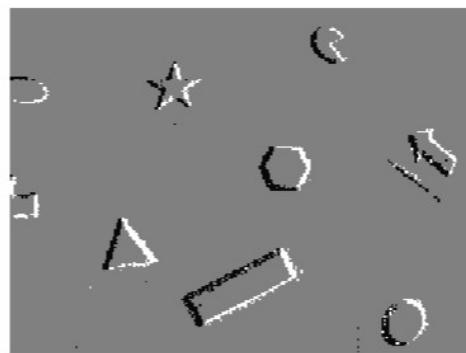


Reduce the event stream to a **corner event stream!**

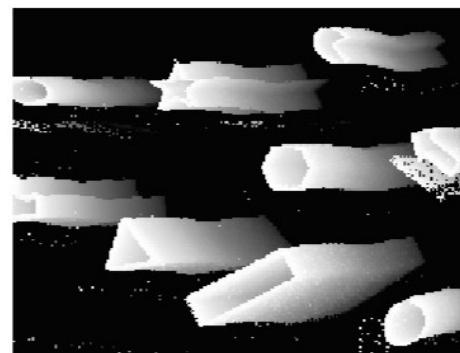
PowerPoint Presentation (uzh.ch)  
[uzh-rpg/rpg\\_corner\\_events: Fast Event-based Corner Detection \(github.com\)](https://uzh-rpg/rpg_corner_events)

# Paper: Fast Event-based Corner Detection

- Approach
  - Data representation: Surface of Active Events (a map with the timestamp of the latest event at each pixel)
  - Analyze timestamp distribution around current event
  - Detect corners by **searching for contiguous pixels with higher timestamps than the rest**



(a) Events over the last 33 ms. Positive events (white) and negative events (black)



(b) Surface of Active Events (SAE). Brightness represents time, from past (dark) to present (bright)



(c) Intensity Image (absolute brightness  $I(x,y)$  in (1))

Figure 2: Signal used for corner detection: the Surface of Active Events (SAE).

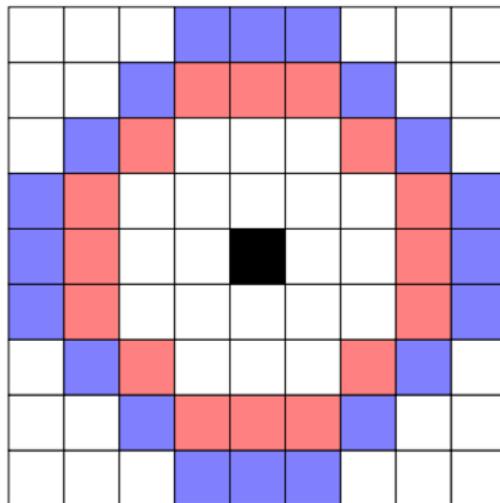
# Paper: Fast Event-based Corner Detection

- Idea: The corner will create below local pattern. **How to detect the pattern?**
  - Along each circle, the algorithm searches for segments of contiguous pixels (arcs) that are higher (i.e., they have a more recent timestamp) than all other pixels on the circle.

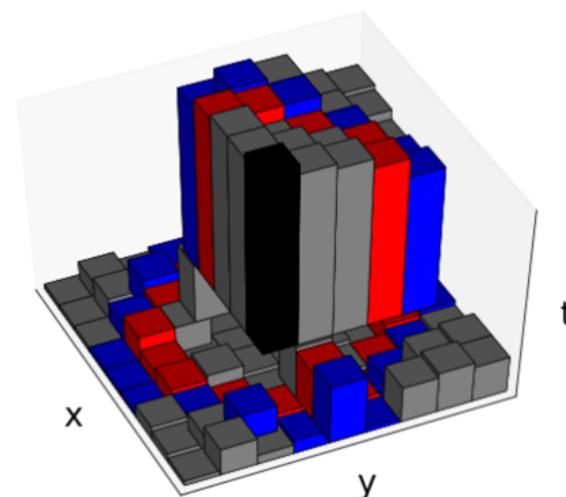
Outer circle: search for a segment of length 4-8

Inner circle: search for a segment of length 3-6

} The center is a corner event!



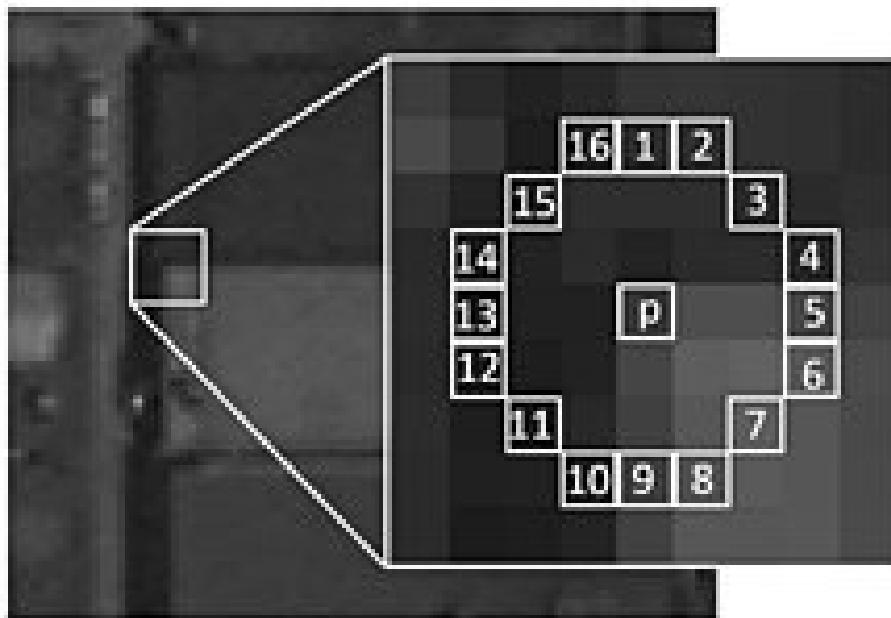
(a) Circles of radius 3 and 4 pixels



(b) Visualization of Surface of Active Events (SAE)

# Paper: Fast Event-based Corner Detection

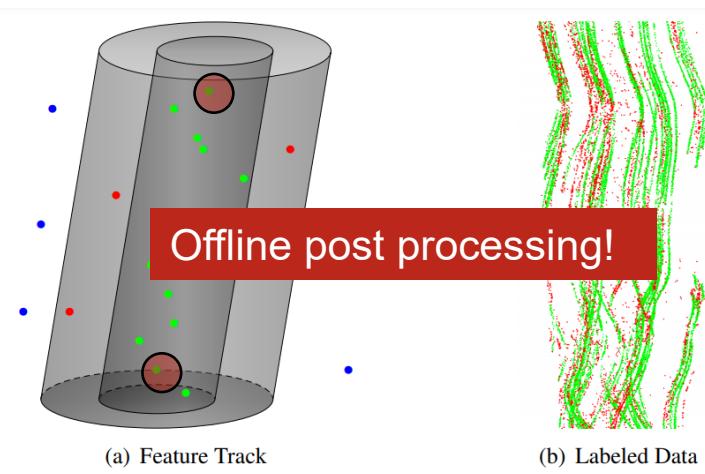
- Idea similar to “FAST” feature in conventional computer vision.



So it is also called “eFAST”!

# Paper: Fast Event-based Corner Detection

- Find feature tracks
  - A feature track is composed of an inner and an outer oblique cylinder in space-time.
  - Exhaustively search for feature tracks by creating hypotheses **using two corner events** and checking whether there are **enough corner events in the inner cylinder** and **few corner events in the outer cylinder**.
  - Use 3 and 5 pixels for the inner and outer radius, respectively.
  - Minimum of 30 inner events, and maximum ratio of outer-to-inner events of 25%.



# Paper: Fast Event-based Corner Detection

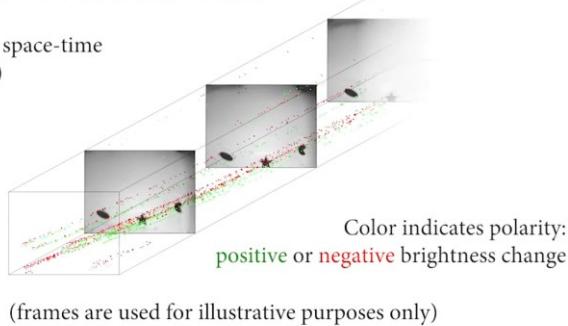
- How fast/accurate is eFAST corner detection?

event  
reduction rate

success rate of  
corner tracking

## Corner Event Stream

Real data in space-time  
(10x slower)



Texture	Dataset	Harris [17]		Ours	
		Red.	FT	Red.	FT
low	shapes_rotation	92.7	74.1	88.9	74.7
	shapes_translation	91.7	78.3	87.8	77.5
	shapes_6dof	90.6	77.1	87.0	76.8
medium	dynamic_rotation	95.1	53.3	96.4	46.4
	dynamic_translation	95.3	62.1	96.7	52.1
	dynamic_6dof	95.4	55.9	96.4	49.4
high	poster_rotation	92.6	35.3	95.7	30.5
	poster_translation	92.3	39.5	95.8	35.9
	poster_6dof	92.4	35.3	95.6	32.1
high	boxes_rotation	92.1	32.9	96.7	25.2
	boxes_translation	92.4	37.0	96.7	30.5
	boxes_6dof	92.7	34.4	96.8	26.7

Table 1: Performance of Harris and our detector expressed as reduction rate (Red.) of the event stream and matched Feature Tracks (FT). Values are given in percentages.

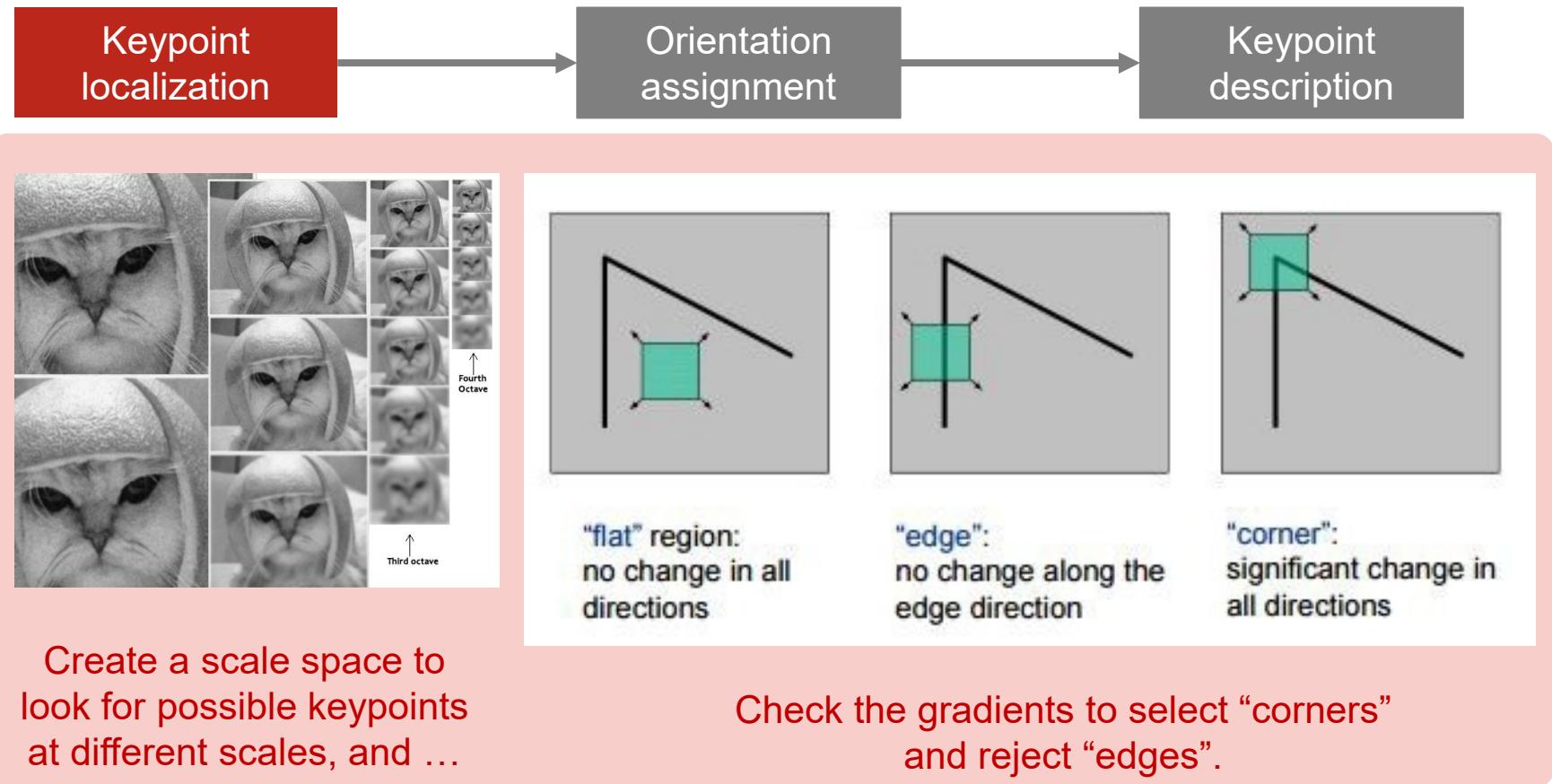
Method	Time per event [μs]	Max. event rate [e/s]
Harris [17]	11.6	86,230
Ours	0.78	1,275,000

Table 2: Runtime comparison per event and corresponding maximum event rate.

# Recap: Feature Detection and Description

- In conventional computer vision, a feature usually comes with a **descriptor**.

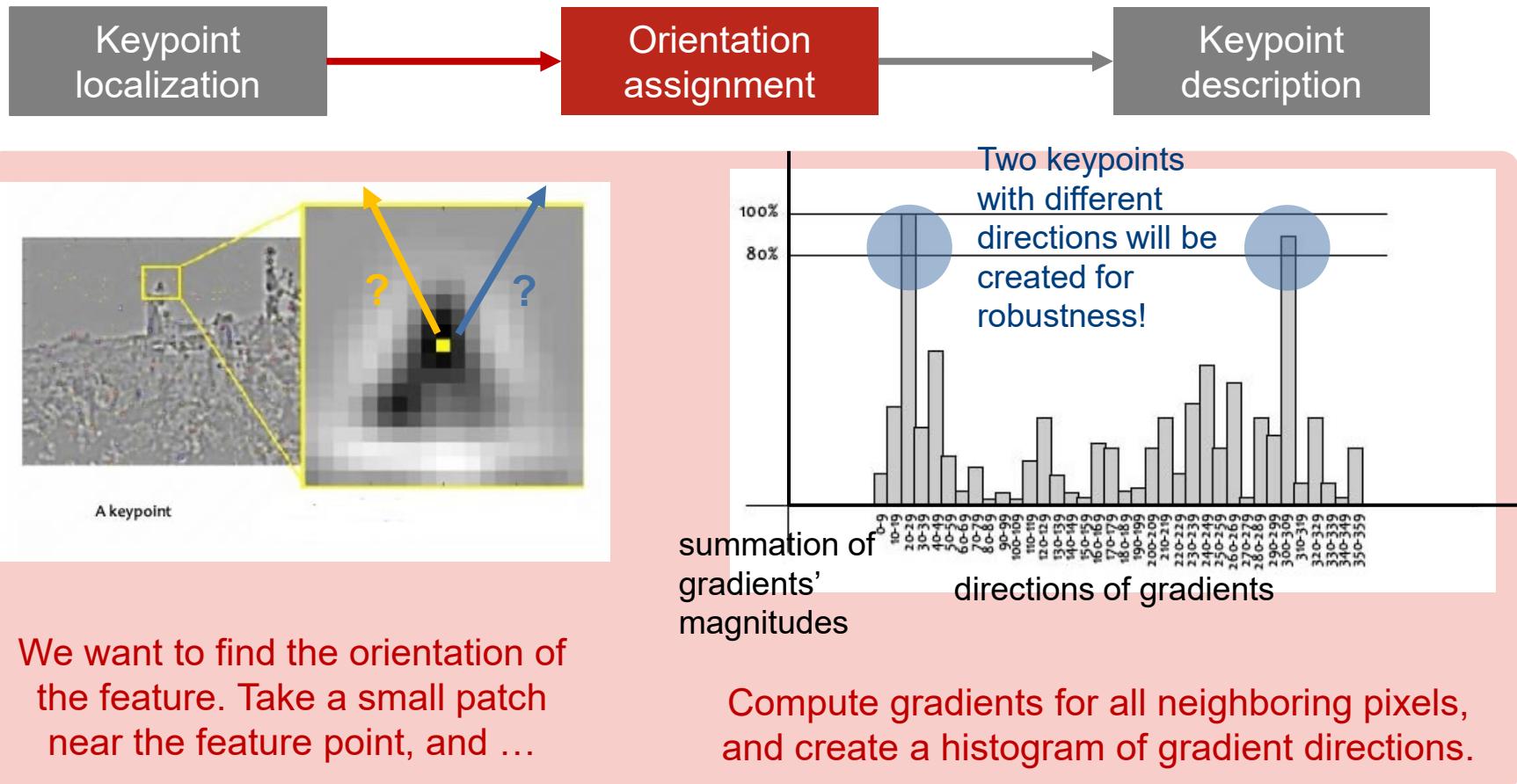
## Step 1: localize the feature



From feature descriptors to deep learning: 20 years of computer vision Eyes of Things  
[Introduction to SIFT\( Scale Invariant Feature Transform\)](#) | by Deepanshu Tyagi | Data Breach | Medium

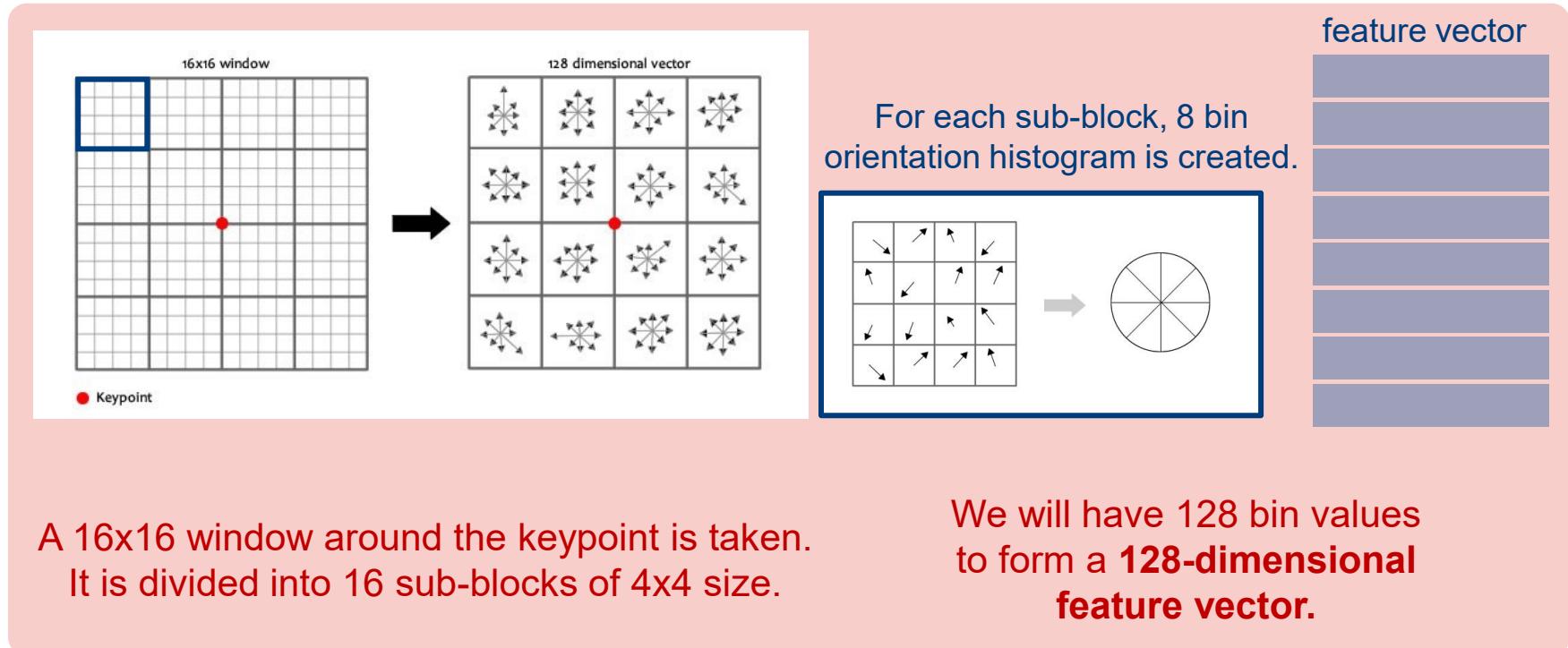
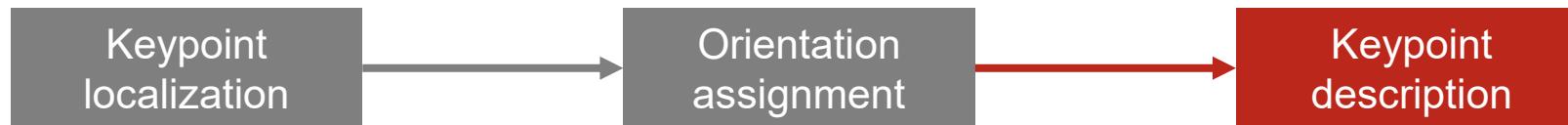
# Recap: Feature Detection and Description

- **Step 2: find the orientation of the feature**



# Recap: Feature Detection and Description

- Step 3: compute the descriptor of the feature



From feature descriptors to deep learning: 20 years of computer vision Eyes of Things  
[Introduction to SIFT\( Scale Invariant Feature Transform\) | by Deepanshu Tyagi | Data Breach | Medium](#)

# Recap: Feature Detection and Description

- We have just discussed **SIFT ( Scale Invariant Feature Transform)**
  - **Scale invariance:**

In the 1<sup>st</sup> step, we extract feature candidates in different scales.
  - **Rotation invariance:**

In the 2<sup>nd</sup> step, direction of the feature is obtained, so that the descriptor contains relative rotation angle, which is compared to the direction of the feature, not the direction of the image.
  - **Translation invariance:**

Obvious.
  - **Illumination invariance:**

Normalization used to reduce the influence from lighting conditions and illuminance.

# Paper: DART: Distribution Aware Retinal Transform for Event-based Cameras

- Highlights
  - Introduced a generic **visual descriptor** “DART”, that encodes the structural context using log-polar grids for event cameras
  - DART is **robust to scale, rotation, and view-point variations.**

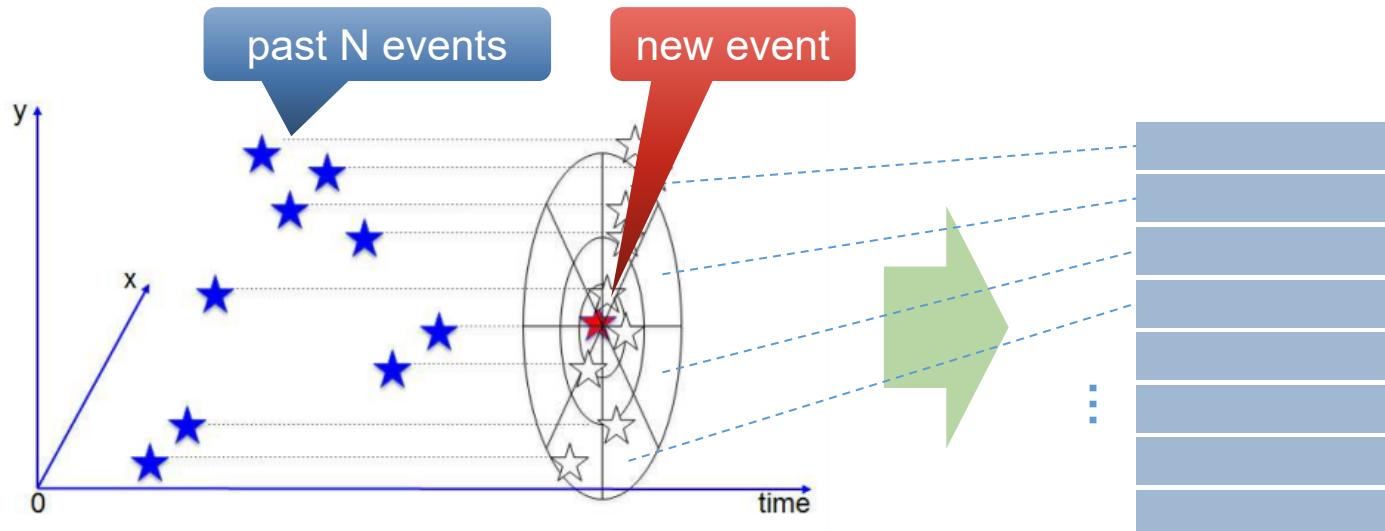
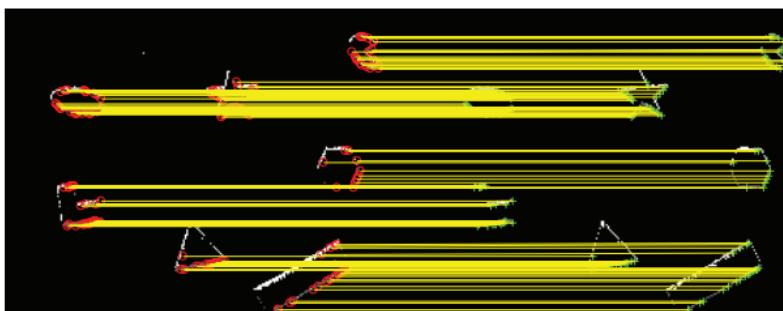


Fig. 1. The most recent event is indicated in red and the previous events are indicated in blue. White stars indicate the position of previous events mapped onto the log-polar bins at the current time.

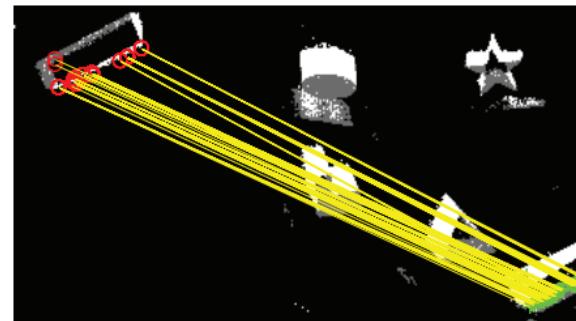
feature vector

# Paper: DART: Distribution Aware Retinal Transform for Event-based Cameras

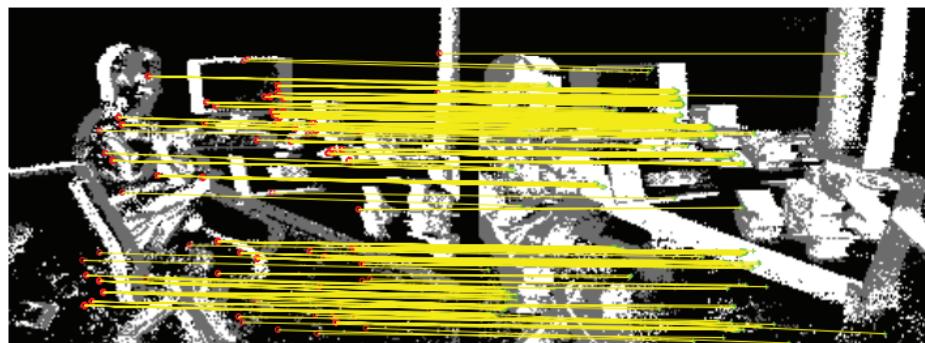
- Feature matching results



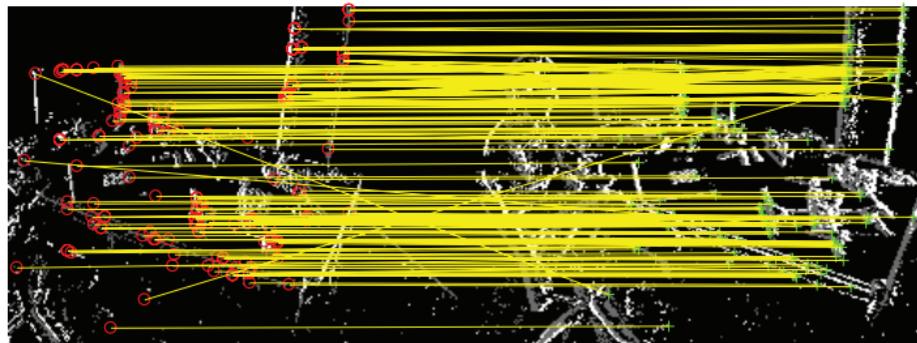
(a) Shapes data of temporally close sequences.



(b) Shapes data of temporally far sequences with occlusion.



(c) Indoor scene with different camera motion speeds.

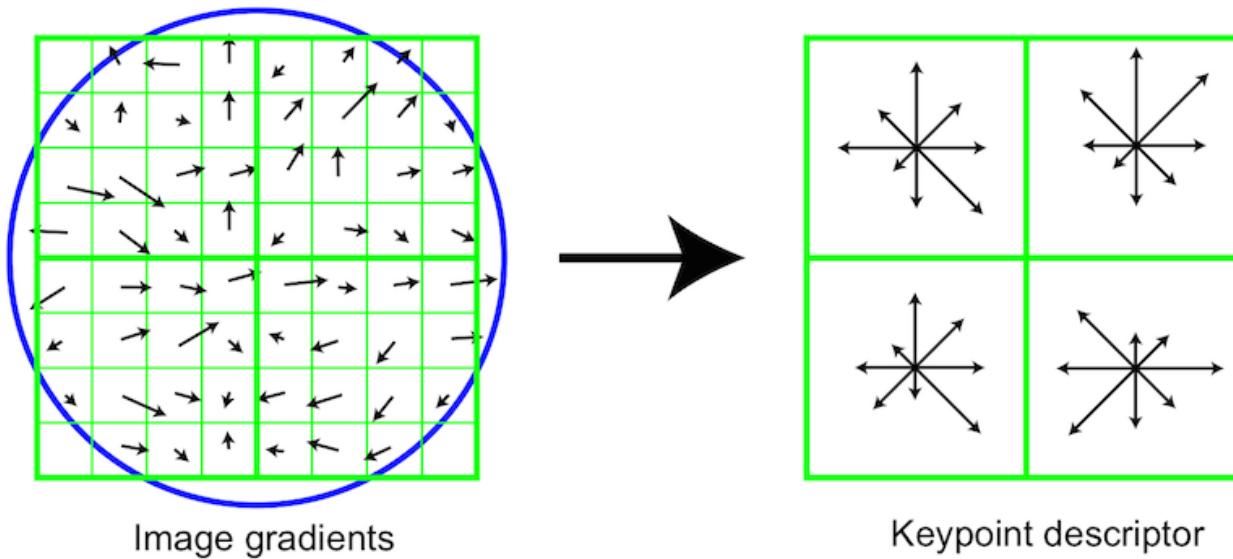


(d) Indoor scene with scale and rotation changes.

Red circles indicate the features in the left time-slice and the green crosses indicate the corresponding match in the right time-slice connected by yellow lines.

# History and Future for Feature Detection and Description

- **Manually-crafted** features/descriptors v.s. **Learned** features/descriptors
  - **The rise of the local feature descriptors: ~1995 to ~2000**  
e.g., SIFT (1999, [David Lowe](#) @ University of British Columbia)



**From comparing raw data (intensities) to comparing features!**

# History and Future for Feature Detection and Description

- **The modern dataset: ~2000 to ~2005**

Geometric problems (homography estimation, ground-plane estimation, robotic vision, SfM) greatly benefited from robust image features such as SIFT.

Images were going online. Datasets were being created.  
e.g., Caltech-101 (2003, Fei-Fei Li, et al. @ Caltech)



From feature descriptors to deep learning: 20 years of computer vision Eyes of Things  
[Introduction to SIFT\( Scale Invariant Feature Transform\) | by Deepanshu Tyagi | Data Breach | Medium](#)

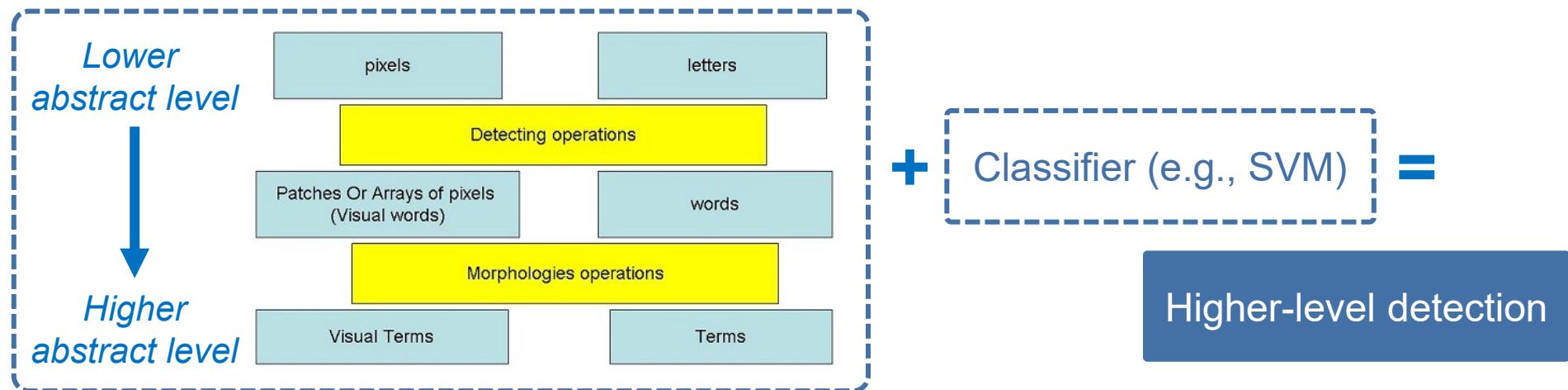
# History and Future for Feature Detection and Description

- **Bins, Grids, and Visual Words (aka Machine Learning meets descriptors): ~2000 to ~2005**

Raw SIFT: not powerful enough to provide object-level metric for matching

Statistical way: To combine more features and descriptors and compare the **feature set** instead of a single feature!

- e.g., Visual words: analogy text - image



From feature descriptors to deep learning: 20 years of computer vision Eyes of Things  
Visual Word - Wikipedia

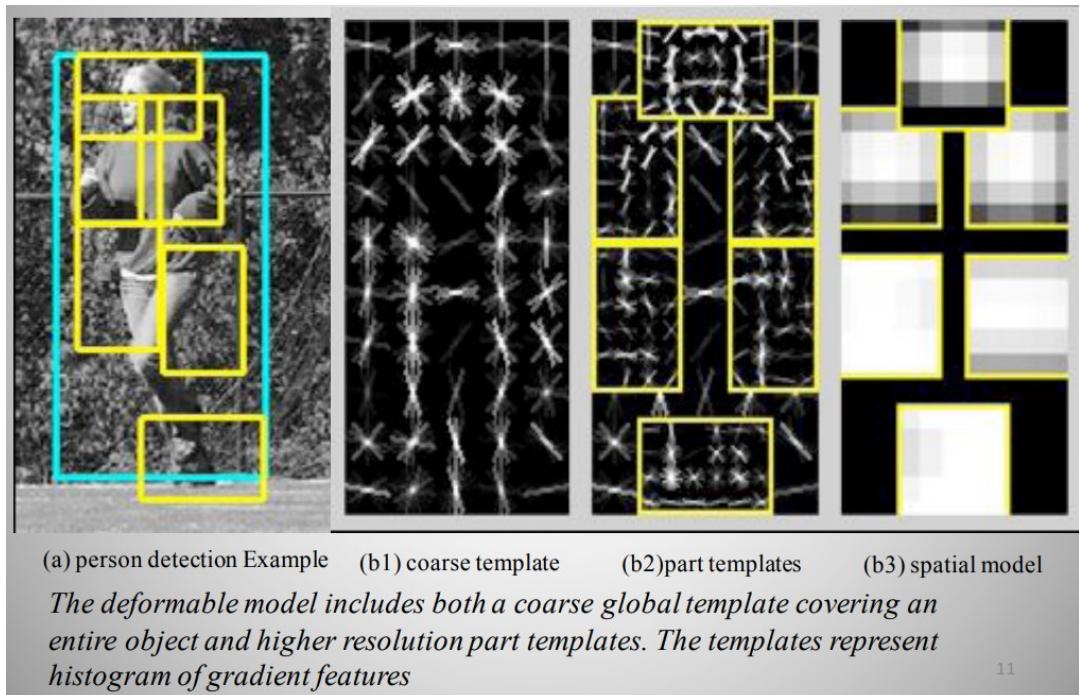
# History and Future for Feature Detection and Description

- **Object Templates (aka the reign of HOG and DPM): ~2005 to ~2010**

HOG: Histogram of oriented gradients  
describe local object appearance and shape by the **distribution of intensity gradients** or edge directions.

Based on HOG, later we have DPM: **Deformable Part Model**.

Idea: To separate object into parts and make them movable!



# History and Future for Feature Detection and Description

- **Big data, Convolutional Neural Networks and the promise of Deep Learning: ~2010 to ~2015 (until present)**

Bigger and bigger data!

+

Learning algorithms (since 1980s)

+

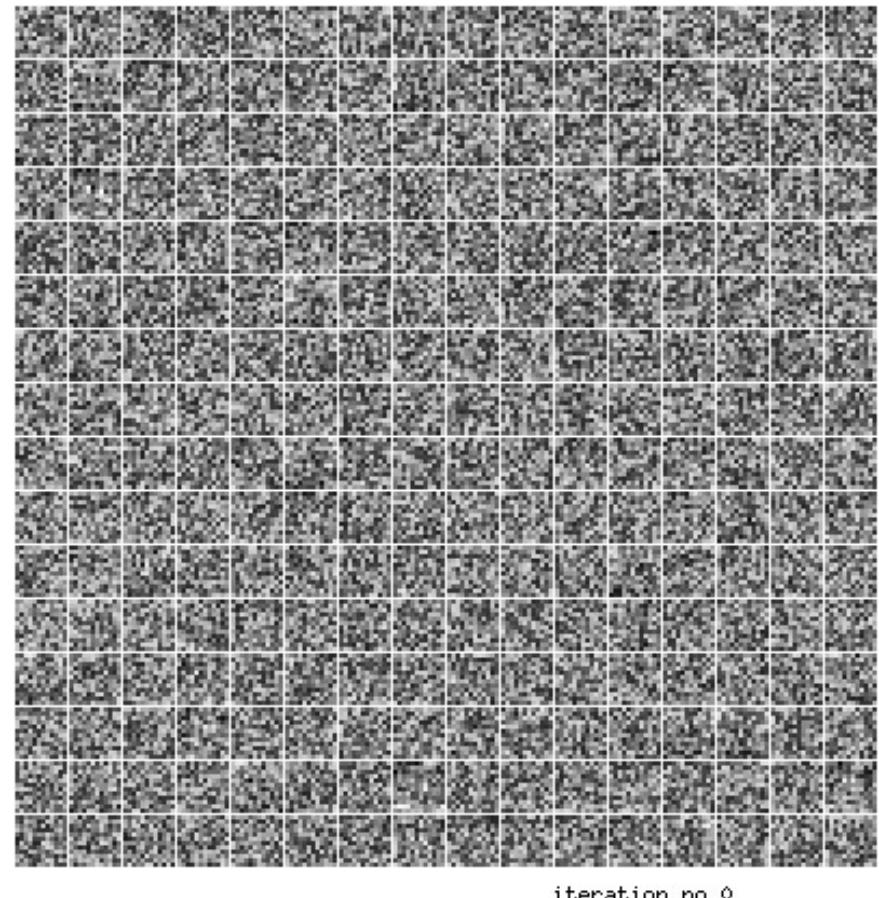
Hardware (GPU)

= ?

**Simplified feature engineering,  
better performance...**

**What about for event-based vision?**

Difficulties exist due to  
fundamentally different data type  
compared to conventional cameras!

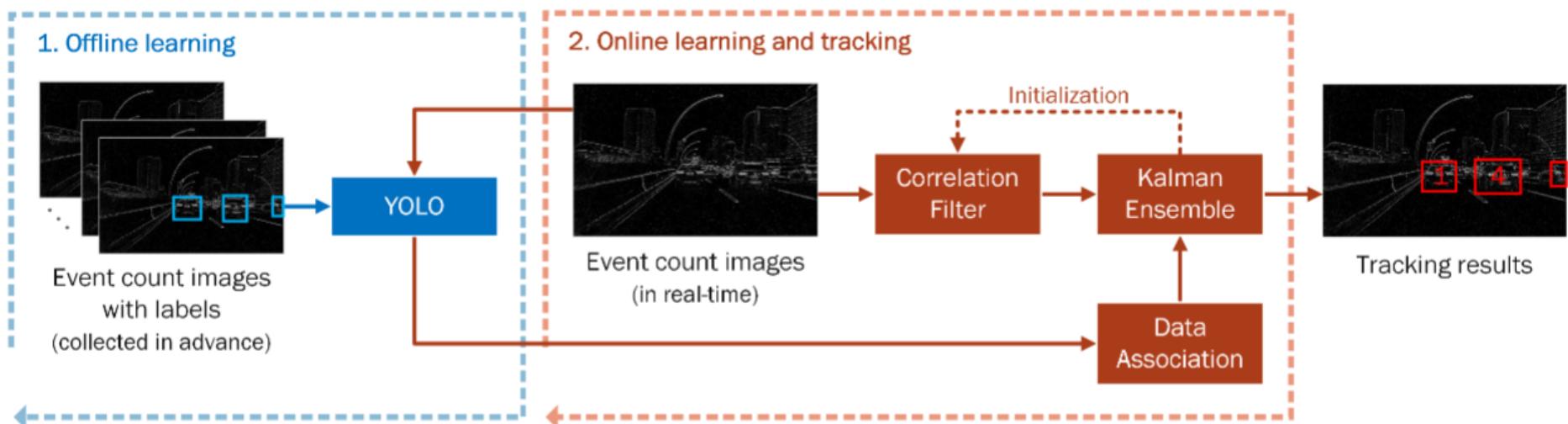


# Paper: Object Tracking on Event Cameras with Offline-Online Learning

- We have just discussed feature detection.  
**What about object-level detection and tracking?**
- Possible ways for object detection/tracking based on event cameras
  - **Detect and track objects on event frames – straightforward and easy**
  - Detect and track object in point clouds
- Highlights
  - Contains YOLO-Based **offline-trained** deep detector + Correlation Filter (CF)-based **online-trained** tracker, which complements each other;
  - A **real-time fusion scheme** designed based on the Kalman filter;
  - It is verified by experiments that: **event count images** contains enough information for object detection and tracking.

# Paper: Object Tracking on Event Cameras with Offline-Online Learning

- Object tracking framework

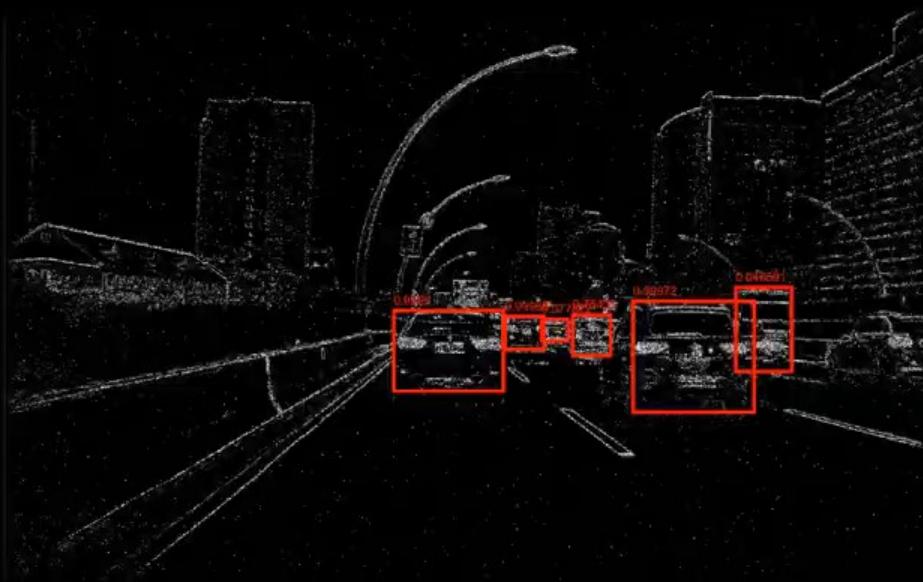


- The 1<sup>st</sup> stage needs to be completed off-line before tracking such that a YOLO detector is trained;
- The correlation filter, which tracks single object, is initialized/reinitialized using results from the YOLO detector.
- The Kalman filter fuses measurements from both the detector and the tracker to achieve multi-object tracking.

(PDF) Object Tracking on Event Cameras with Offline-Online Learning (researchgate.net)

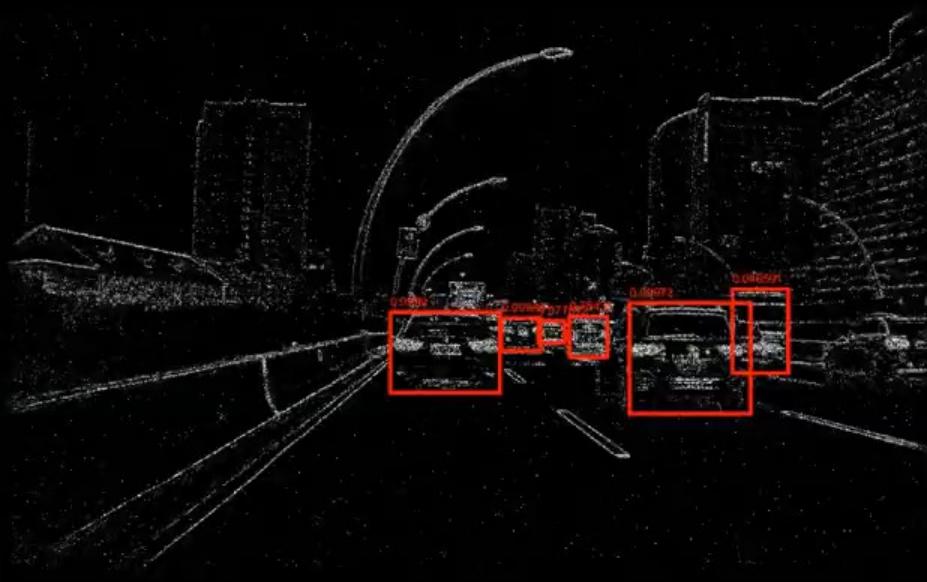
# Paper: Object Tracking on Event Cameras with Offline-Online Learning

## Object Tracking on Event Cameras with Offline-Online Learning Urban Scenario



Partial Configuration

without CF, tracking at 459 fps



Full Configuration

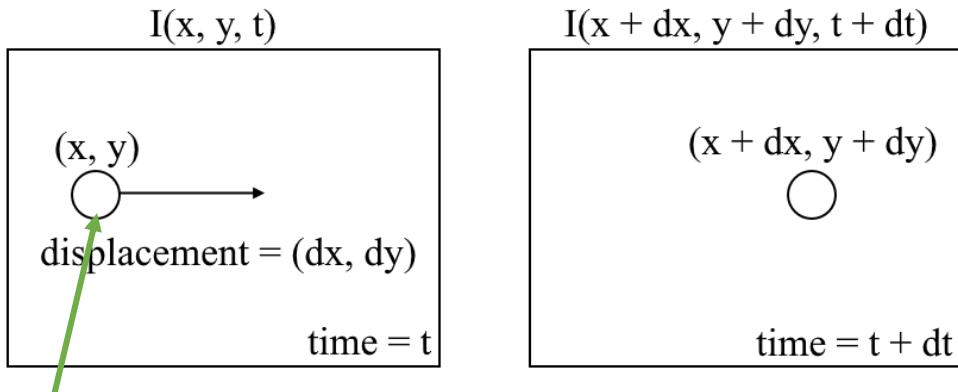
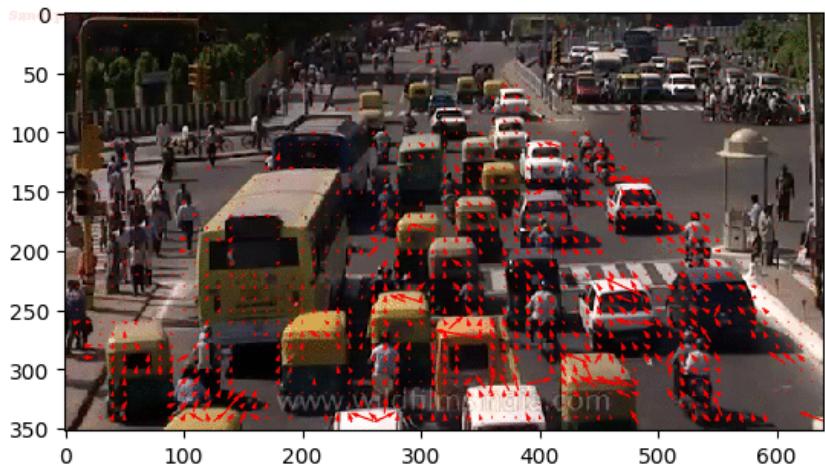
with offline-online tracker, at ~50 fps

## Summary – Detection and tracking

- Relation between detection and tracking
- Different abstraction levels of detection and tracking:
  - Event level → feature level → object level
- Why we need event-based feature detection?
  - To reduce the data rate, aiming at fast event data processing
- Feature detection and description
  - Descriptors required to “identify” features
  - From hand crafted to learned features
- Object-level detection and tracking
- Examples and papers

# Optical Flow Estimation

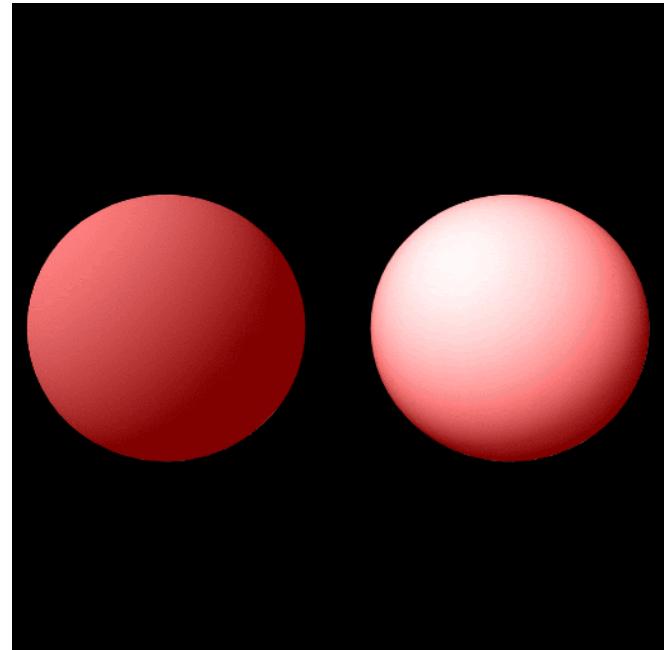
- **Optical flow field vs. motion field**
  - Motion field reflects **object's real motion** in 3-D space
  - Optical flow field is the **distribution of apparent velocities** of movement of **brightness pattern** in an image (2-D space)



“brightness  
pattern”

# Optical Flow Estimation

- Does the following scenario exist?
  - There is no motion field or optical flow;
  - There is both motion field and optical flow;
  - There is motion field but no optical flow;
  - There is optical flow but no motion field.



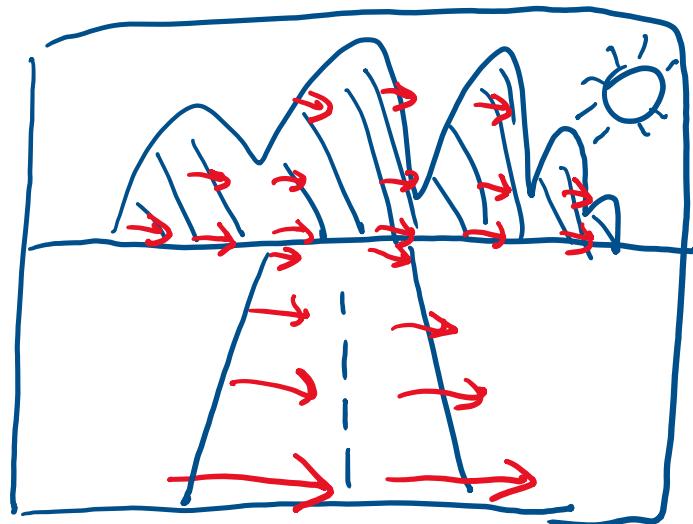
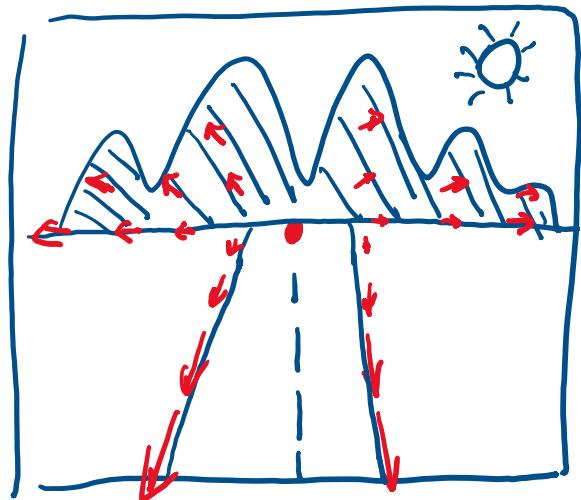
**Information loss may be due to:**

1. 3D to 2D projection
2. Non-informative texture on objects

Can you tell whether it is because of the motion of the sphere, or the light source?

# Information in Optical Flow

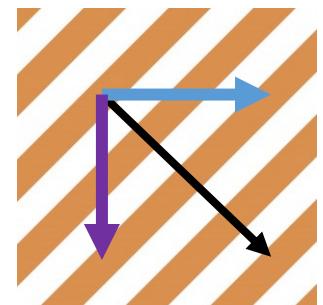
- Motion information: translation and rotation



**We always want to recover motion (2D/3D)!**

# Information in Optical Flow

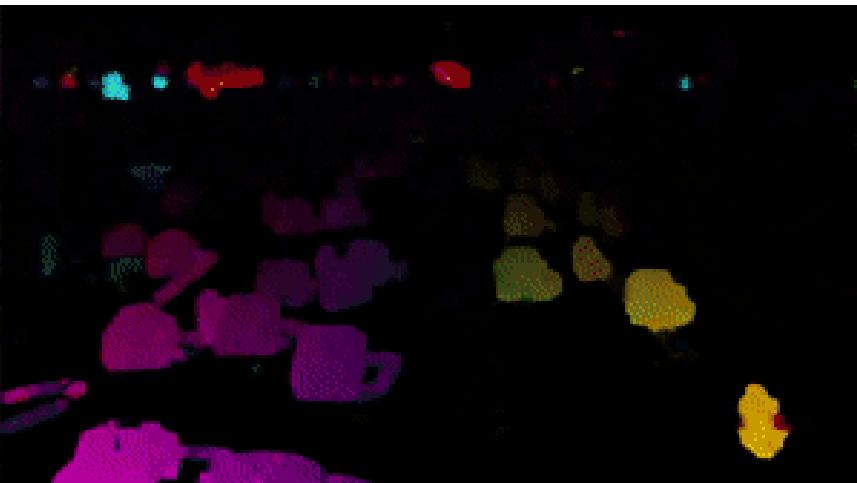
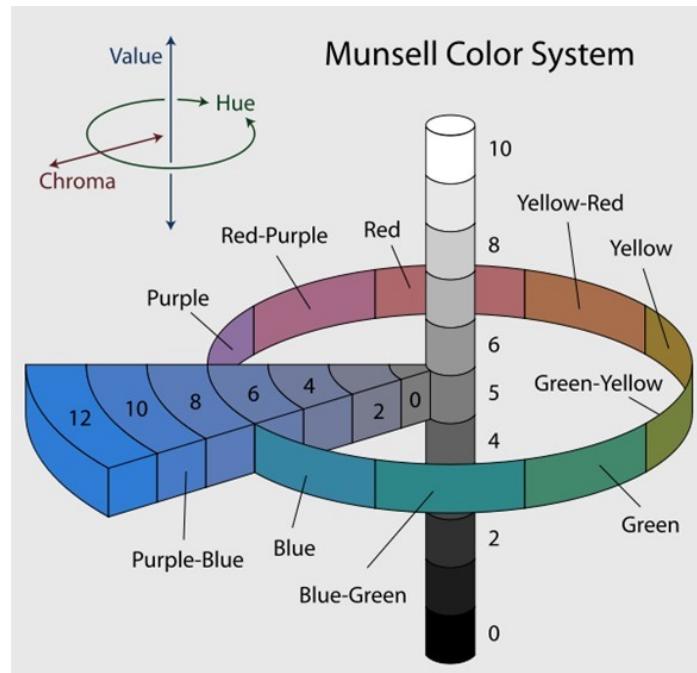
- Why motion recovery not easy?
  - Flow depending on scene/background, but **scene unknown**  
→ usually recover motion together with scene (reconstruction)
  - Ego-motion coupled with translation and rotation, **different motion may lead to same flow**  
→ internal measurements (IMU, other sensors)  
→ build theoretically-robust model for motion and scene
  - 3-D to 2-D projection, **irreversible information loss**  
→ camera needs to move, observe from multi-angles to reconstruct 3-D



**Estimating flow is the first step.**

# Visualization of Optical Flow

- HSV instead of RGB representation
  - Hue
  - Saturation/Chroma
  - Value



# Frame-based Flow Estimation

- The problem formulation:

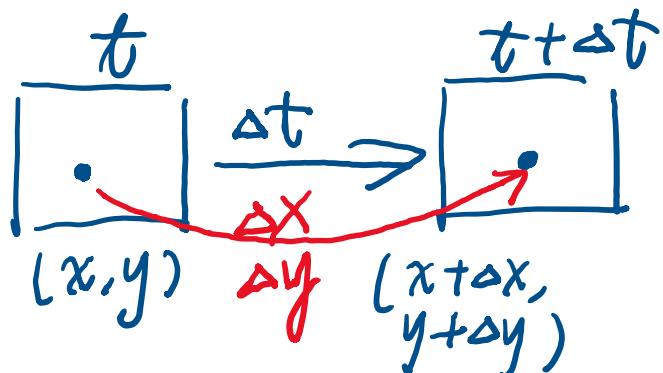
Optimization of a global energy function

$$E_{\text{Global}} = E_{\text{Data}} + \lambda E_{\text{Prior}}.$$

- **Data term:** how consistent the optical flow is with the input images;
- **Prior term:** it favors certain flow fields over others (e.g. smoothly varying flow fields).

# Flow Estimation: Data Term

- Brightness consistency → **optical flow constraint equation**



$$\begin{aligned} I(x, y, t) &= I(x + \Delta x, y + \Delta y, t + \Delta t) \\ &= I(x, y, t) + \frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t \\ &\quad + \text{H.o.t.} \end{aligned}$$

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

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

~~~~~

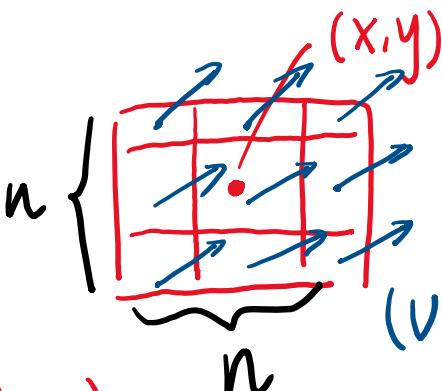
$$\begin{aligned} \frac{\partial I}{\partial x} V_x + \frac{\partial I}{\partial y} V_y + \frac{\partial I}{\partial t} V_t &= 0 \\ I_x V_x + I_y V_y - I_t &= 0 \end{aligned}$$

Above equation leads to one error per pixel...

**How to get the global constraints for the whole image?**

# Flow Estimation: Data Term

- Aggregating error over the local patch (**Lucas-Kanade** method)



Constant velocity assumption in neighborhood

$$m = n^2 \text{ eqns}$$
$$\begin{pmatrix} \frac{\partial I(x_1, y_1)}{\partial x} & \frac{\partial I(x_1, y_1)}{\partial y} \\ \vdots & \vdots \\ \frac{\partial I(x_m, y_m)}{\partial x} & \frac{\partial I(x_m, y_m)}{\partial y} \end{pmatrix} \begin{pmatrix} v_x \\ v_y \end{pmatrix} = \begin{pmatrix} -I_{t1} \\ \vdots \\ -I_{tm} \end{pmatrix}$$

Use the **Least Squares Method** to solve above **overdetermined system**.  
Only **local** and **sparse flow** can be obtained.

# Flow Estimation: Data Term

- Aggregating error over the whole image
  - L2 norm (used in **Horn and Schunck** algorithm)

$$E_{\text{Data}} = \sum_{x,y} \left[ u \frac{\partial I}{\partial x} + v \frac{\partial I}{\partial y} + \frac{\partial I}{\partial t} \right]^2.$$

- L1 norm

$$\|\mathbf{E}\|_1 = \sum_{x,y} |E_{x,y}| \approx \sum_{x,y} \sqrt{\|E_{x,y}\|^2 + \epsilon^2}.$$

# Flow Estimation: Prior Term

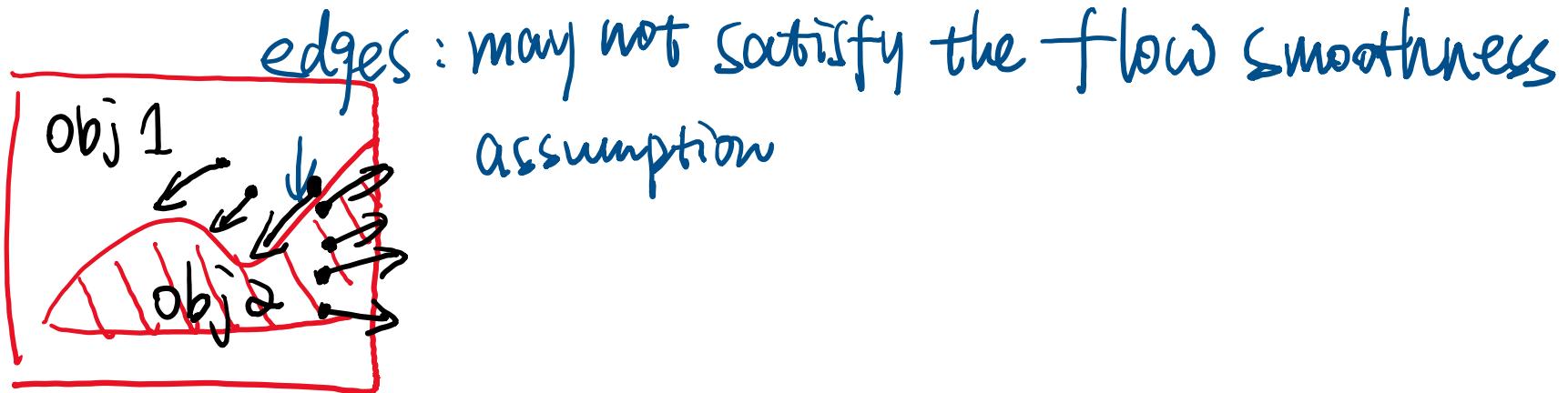
- Spatial and temporal smoothness
  - Favor small first-order derivatives (gradients) of the flow field (used in **Horn and Schunck** algorithm)

$$E_{\text{Prior}} = \sum_{x,y} \left[ \left( \frac{\partial u}{\partial x} \right)^2 + \left( \frac{\partial u}{\partial y} \right)^2 + \left( \frac{\partial v}{\partial x} \right)^2 + \left( \frac{\partial v}{\partial y} \right)^2 \right].$$

- Also, possible to add temporal smoothness terms  $\frac{\partial u}{\partial t}$  and  $\frac{\partial v}{\partial t}$ .

# Flow Estimation: Prior Term

- Different weights for different locations. **Why?**
  - Flow may not be smooth at edges.
  - We may want to reduce the weight of the prior at edges (with high  $|\nabla I|$ ).

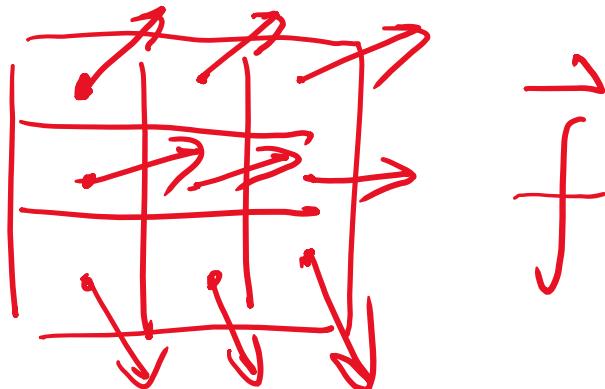


To do that, design the prior term as

$$E_{\text{Prior}} = \sum_{x,y} w(\nabla I) \left[ \left( \frac{\partial u}{\partial x} \right)^2 + \left( \frac{\partial u}{\partial y} \right)^2 + \left( \frac{\partial v}{\partial x} \right)^2 + \left( \frac{\partial v}{\partial y} \right)^2 \right].$$

# Flow Estimation: Optimization

- Recap: **How to formulate the flow estimation problem?**



Recall the flow est. prob. :  
 $\arg\min \{ E_{\text{Global}} \}$

$$\underset{\vec{f}}{\circlearrowleft}$$

vector field

The calculus of variations will help to solve the above problem.

- Hard to find the analytical solution due to various  $E_{\text{Global}}$
- Numerical and iterative solution

# Event-based Flow Estimation

- **Lucas-Kanade variants**

No frames/intensities... How to obtain spatial and temporal gradients?

1. For each neighborhood, use events to compute temporal gradients and approximate spatial gradients
2. Find the local flow based on L-K method

- + Straightforward and easy to implement
- Noisy results, needs further filtering
- Approximation on spatial gradients not accurate
- Only sparse flow available

$$\begin{cases} \frac{\partial e(x,y,t)}{\partial x} \sim \sum_{t-\Delta t}^t e(x,y,t) - \sum_{t-\Delta t}^t e(x-1,y,t) \\ \frac{\partial e(x,y,t)}{\partial y} \sim \sum_{t-\Delta t}^t e(x,y,t) - \sum_{t-\Delta t}^t e(x,y-1,t) \end{cases}$$

$$\text{grad}^t(I) \begin{pmatrix} v_x \\ v_y \end{pmatrix} + \frac{\partial I}{\partial t} = 0$$

$$\frac{\partial e(x,y,t)}{\partial t} \sim \frac{\sum_{t-\Delta t}^{t_1} e(x,y,t) - \sum_{t-\Delta t}^t e(x,y,t)}{t-t_1}$$

Benosman, Ryad, et al. "Asynchronous frameless event-based optical flow." Neural Networks 27 (2012): 32-37.

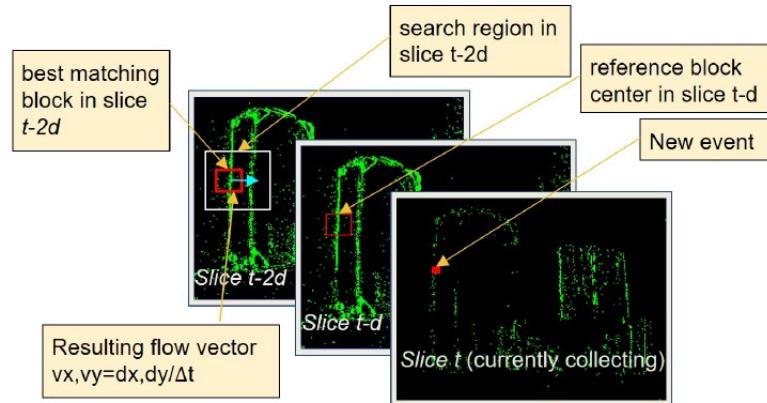
# Event-based Flow Estimation

- **Block matching**

- Idea: match **block**, instead of a single pixel

- + More robust as block may contain more information

- Not good for complex scenes due to repetitive textures and contours
  - Defining block features is much more difficult than defining point features!

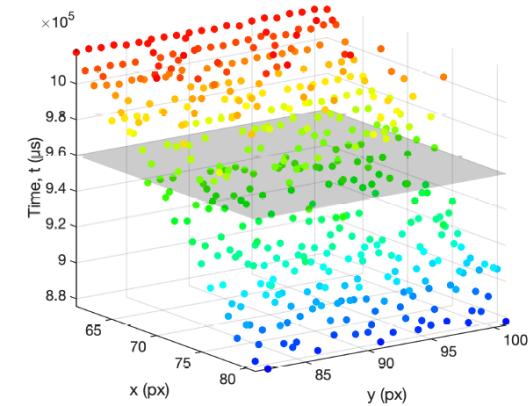


<http://www.bmva.org/bmvc/2018/contents/papers/0280.pdf>

Liu, Min, and Tobi Delbrück. "Adaptive time-slice block-matching optical flow algorithm for dynamic vision sensors." (2018).

# Event-based Flow Estimation

- **Local plane fitting**
  1. Create **event time surface**
  2. Use a local plane to fit the time surface
  3. Gradient (normal vector) of the plane indicates the flow
    - Noisy results that requires further filtering
- **Energy minimization/optimization** (deep learning)
  - + Global method
  - + Able to provide dense flow
  - High computational complexity
  - Data-driven for deep learning



$$\min_{\mathbf{u}, L} \int_{\Omega} \int_T \left( \lambda_1 \|\mathbf{u}_x\|_1 + \lambda_2 \|\mathbf{u}_t\|_1 + \lambda_3 \|L_x\|_1 + \lambda_4 \|\langle L_x, \delta_t \mathbf{u} \rangle + L_t\|_1 + \lambda_5 h_\theta(L - L(t_p)) \right) dt dx$$

$$+ \int_{\Omega} \sum_{i=2}^{|P(\mathbf{x})|} \|L(t_i) - L(t_{i-1}) - \theta \rho_i\|_1 dx,$$

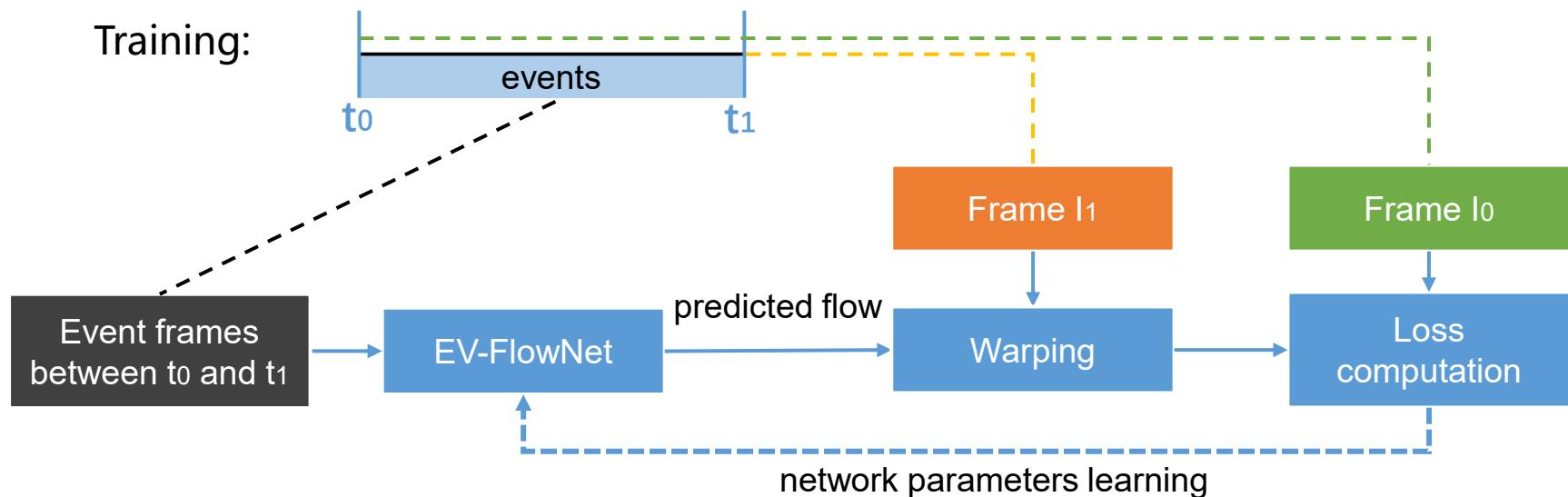
Ieng, Sio Hoi, João Carneiro, and Ryad Benjamin Benosman. "Event-Based 3D Motion Flow Estimation Using 4D Spatio Temporal Subspaces Properties." Frontiers in neuroscience 10 (2017): 596.  
Benosman, Ryad, et al. "Event-based visual flow." IEEE transactions on neural networks and learning systems 25.2 (2013): 407-417.

Low, Weng Fei, et al. "SOFEA: A Non-iterative and Robust Optical Flow Estimation Algorithm for Dynamic Vision Sensors." CVPR Workshops. 2020.

Pan, Liyuan, Miaomiao Liu, and Richard Hartley. "Single image optical flow estimation with an event camera." 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 2020.  
Bardow, Patrick, Andrew J. Davison, and Stefan Leutenegger. "Simultaneous optical flow and intensity estimation from an event camera." CVPR. 2016.

# Paper: EV-FlowNet

- **Self-Supervised** Optical Flow Estimation for Event-based Cameras
  - Self-supervised: no ground-truth flow required
  - Groundtruth frames required **only in training stage**



Inference: The input only contains events.

# Paper: EV-FlowNet

- Event representation
  - **Image form representation** instead of the point cloud form
  - Input to the network: **4 channel image**, same resolution as the camera



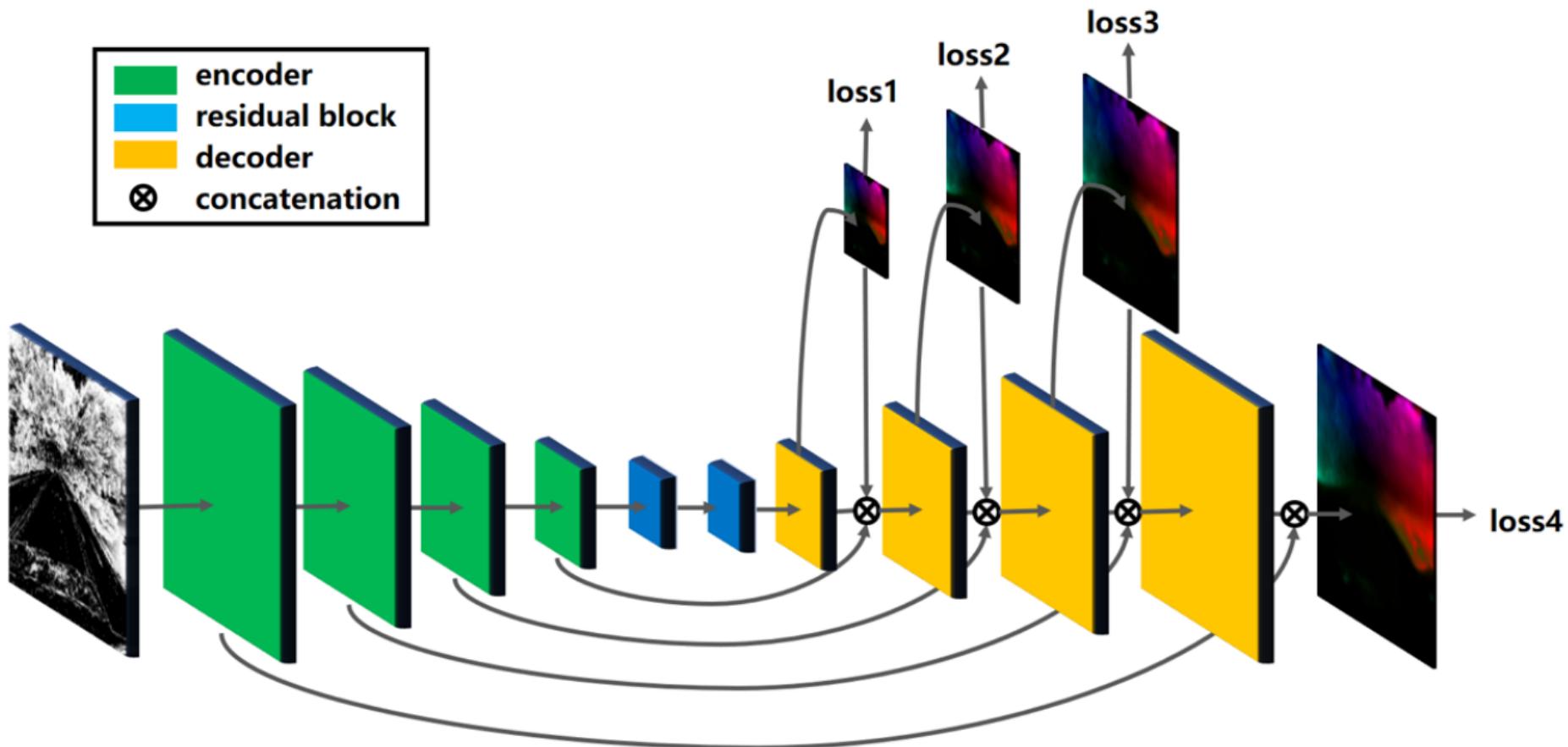
# Paper: EV-FlowNet

- Self-supervised loss
  - **Photometric loss:** minimize the difference in intensity between the warped second image and the first image
$$\ell_{\text{photometric}}(u, v; I_t, I_{t+1}) = \sum_{x,y} \rho(I_t(x, y) - I_{t+1}(x + u(x, y), y + v(x, y)))$$
where  $\rho$  is the Charbonnier loss function for outlier rejection.
  - **Smoothness loss:** regularize the output flow by minimizing the difference in flow between neighboring pixels horizontally, vertically and diagonally.
$$\ell_{\text{smoothness}}(u, v) = \sum_{x,y} \sum_{i,j \in N(x,y)} \rho(u(x, y) - u(i, j)) + \rho(v(x, y) - v(i, j))$$
where  $N$  is the set of neighbors around  $(x, y)$ .

[\[1802.06898\] EV-FlowNet: Self-Supervised Optical Flow Estimation for Event-based Cameras \(arxiv.org\)](https://arxiv.org/abs/1802.06898)

# Paper: EV-FlowNet

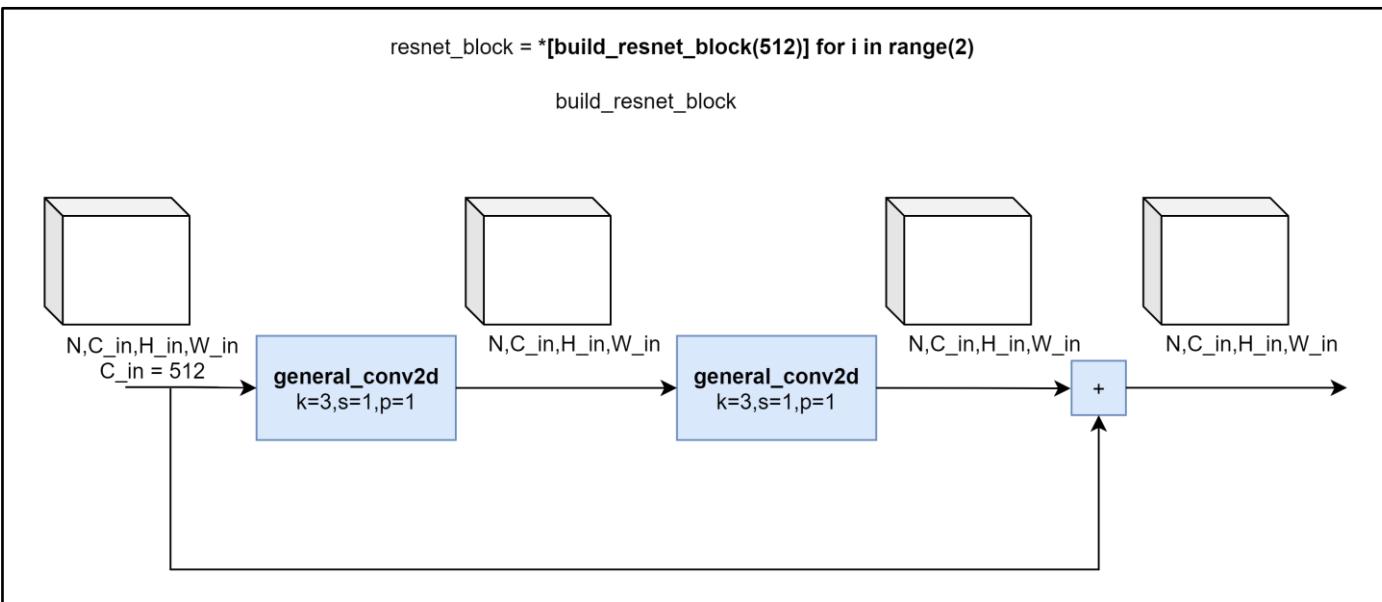
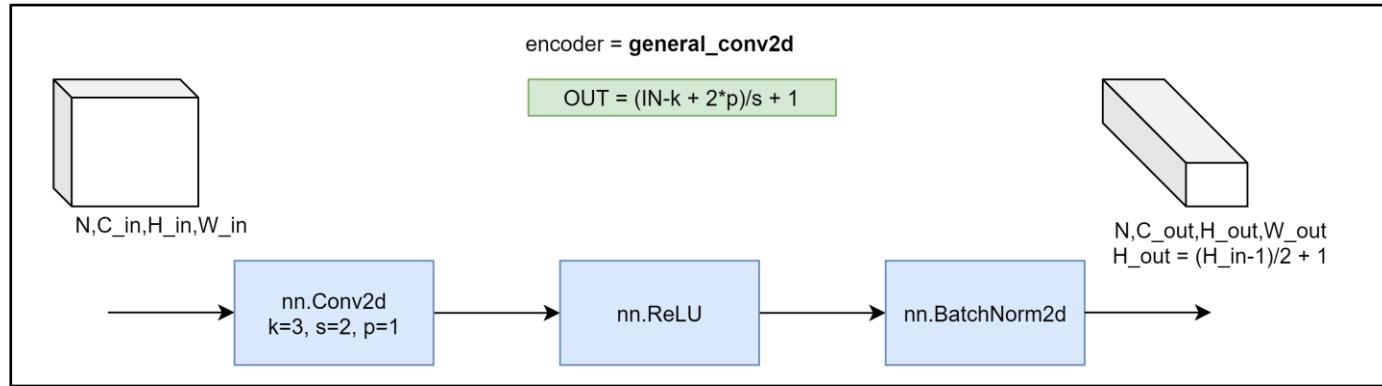
- EV-FlowNet architecture



[1802.06898] EV-FlowNet: Self-Supervised Optical Flow Estimation for Event-based Cameras (arxiv.org)

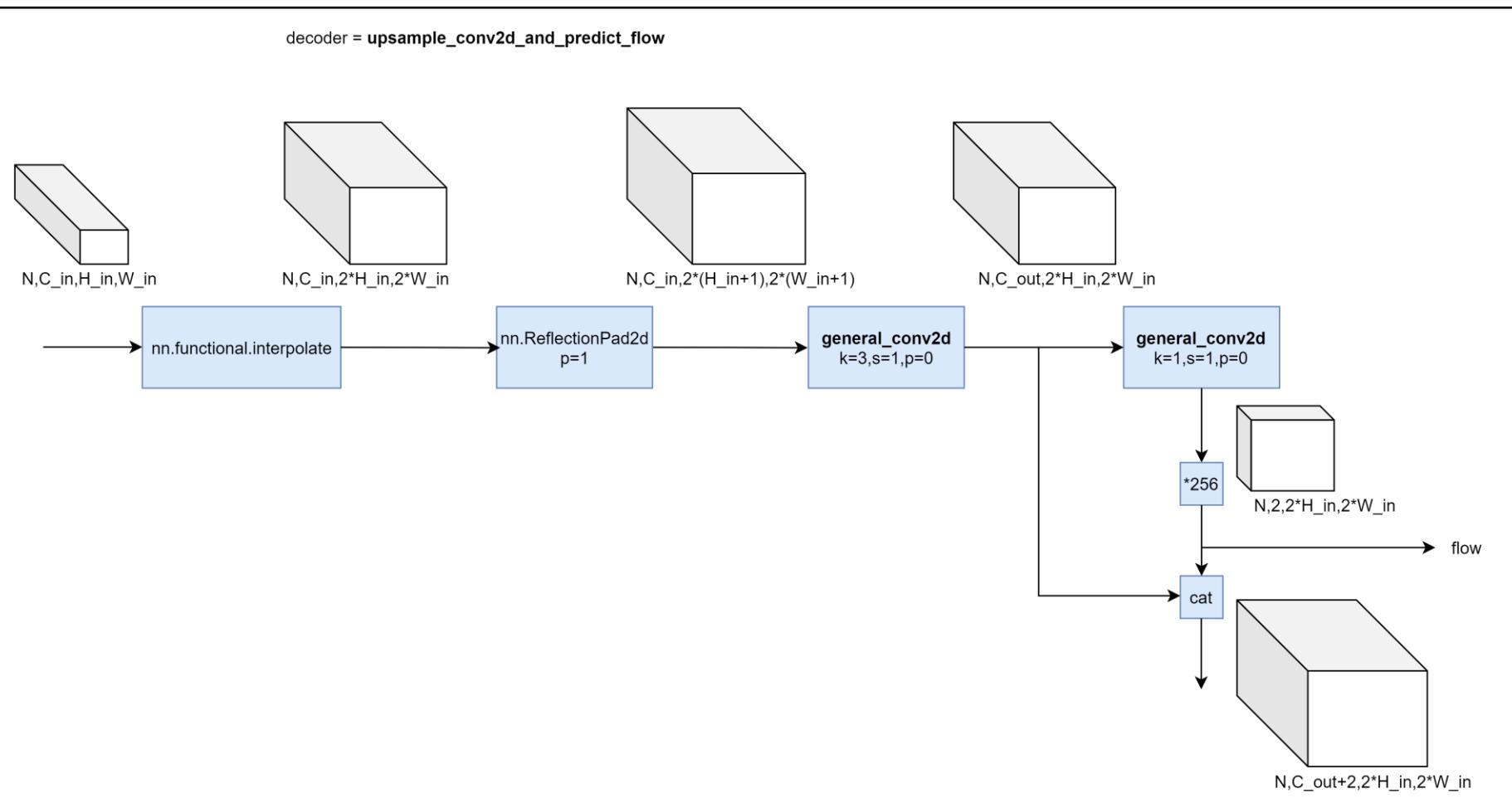
# Paper: EV-FlowNet

- EV-FlowNet architecture in details: sub-modules



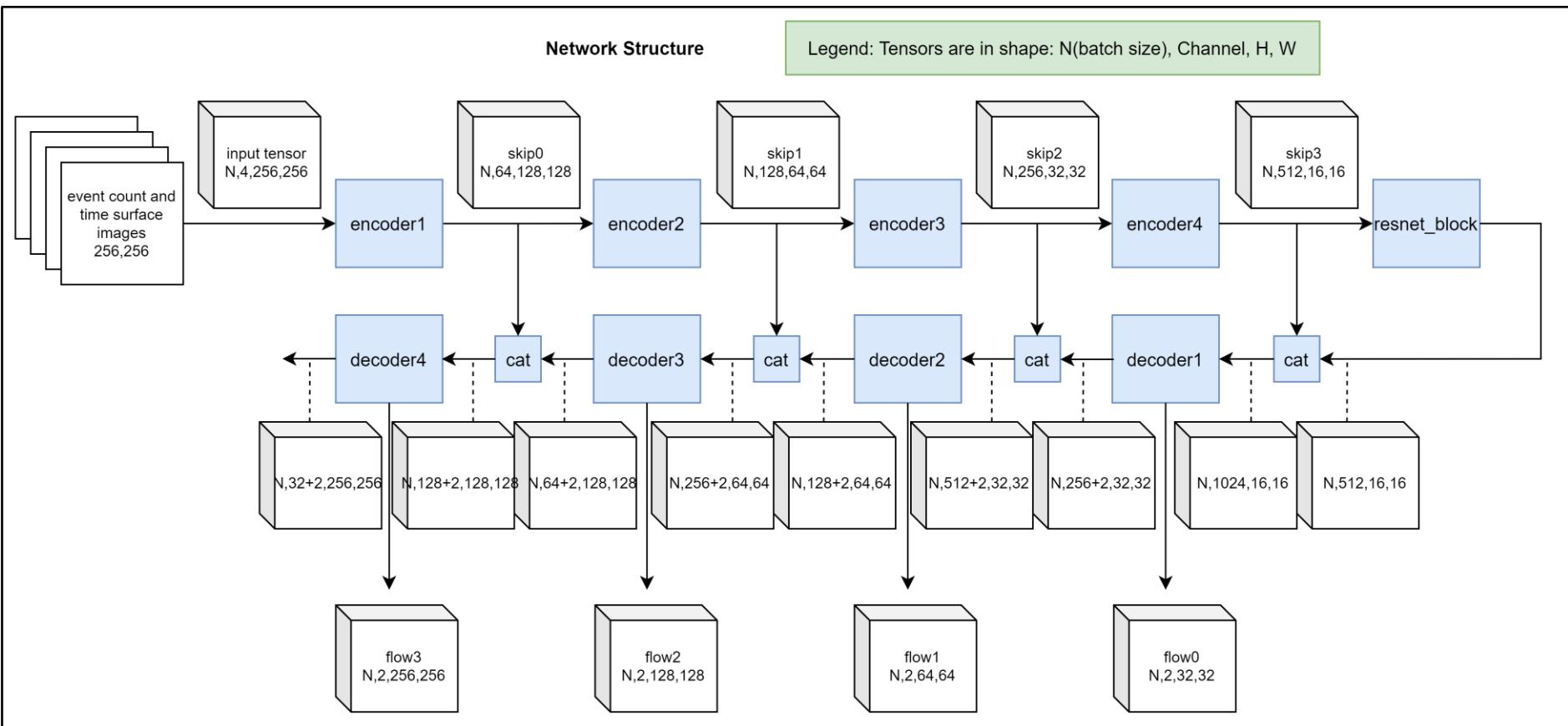
# Paper: EV-FlowNet

- EV-FlowNet architecture in details: sub-modules



# Paper: EV-FlowNet

- EV-FlowNet architecture in details: framework



# Paper: EV-FlowNet

## EV-FlowNet: Self-Supervised Optical Flow Estimation for Event-based Cameras

Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, Kostas Daniilidis

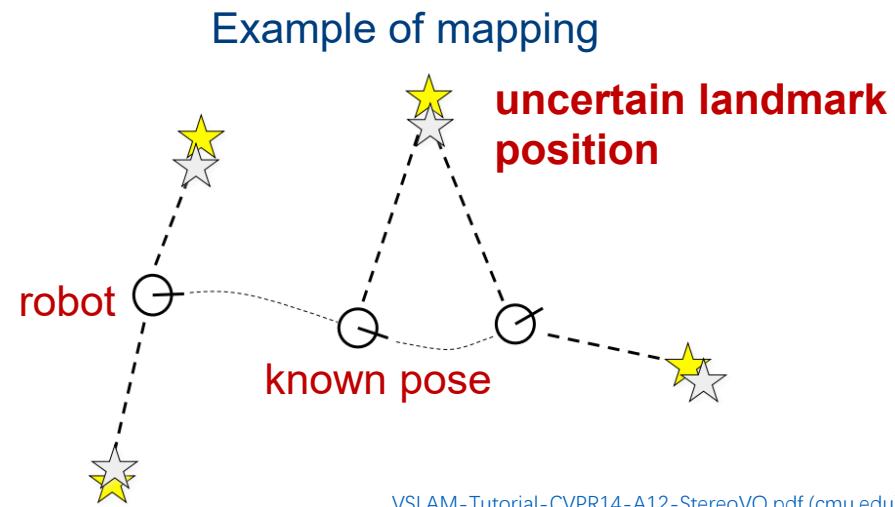
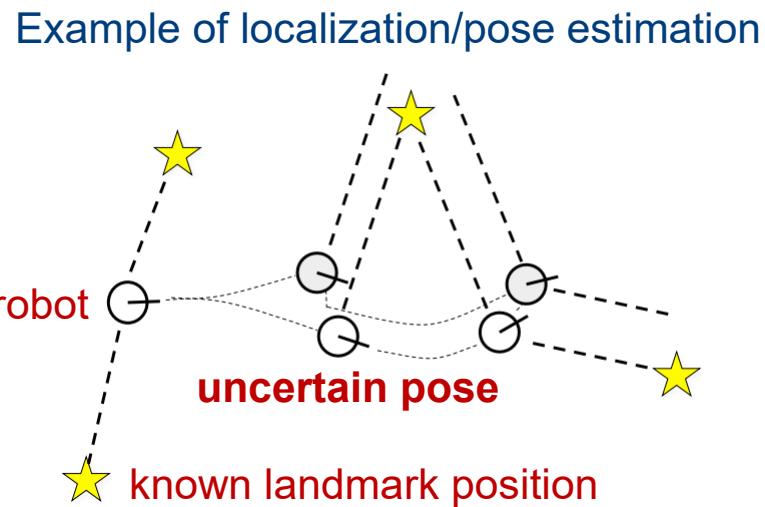


## Summary – Optical flow estimation

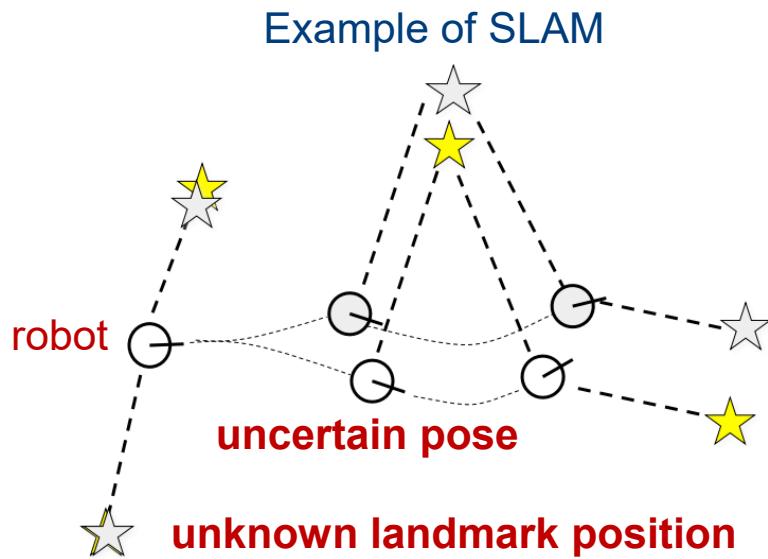
- Optical flow field vs. motion field
- Information in optical flow: Why it is necessary to obtain optical flow?
- Visualization of optical flow
- Frame-based optical flow estimation
- Event-based flow estimation
  - Lucas-Kanade variants
  - Block matching
  - Local plane fitting
  - Energy minimization/optimization (with/without deep learning)

# Pose Estimation and SLAM

- Problem
  - Pose estimation:  
Estimating **robot's poses only**.
  - SLAM (Simultaneous Localization And Mapping):  
Computing the **robot's poses** and the **map of the environment** at the same time.



# Pose Estimation and SLAM



- SLAM is a chicken-or-egg problem:
  - a map is needed for localization and
  - a pose estimate is needed for mapping

# Pose Estimation and SLAM

- Definition of the SLAM Problem

## Given

- The robot's controls

$$u_{1:T} = \{u_1, u_2, u_3 \dots, u_T\}$$

- Observations

$$z_{1:T} = \{z_1, z_2, z_3 \dots, z_T\}$$

## Wanted

- Map of the environment

$m$

- Path of the robot

$$x_{0:T} = \{x_0, x_1, x_2 \dots, x_T\}$$

# Pose Estimation and SLAM

- Definition of the SLAM Problem

## Given

- The robot's controls

$$u_{1:T} = \{u_1, u_2, u_3 \dots, u_T\}$$

- Observations

$$z_{1:T} = \{z_1, z_2, z_3 \dots, z_T\}$$

## Wanted

- Map of the environment

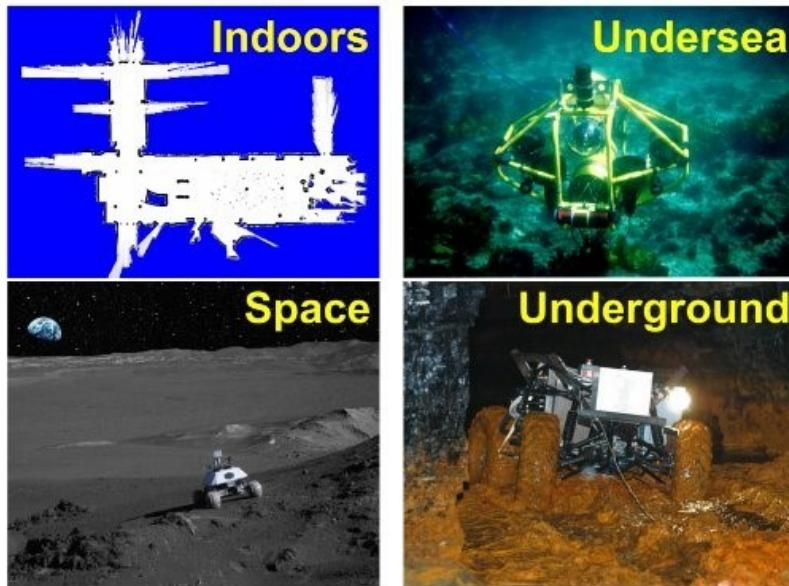
$m$

- Path of the robot

$$x_{0:T} = \{x_0, x_1, x_2 \dots, x_T\}$$

# Pose Estimation and SLAM

- Why we need SLAM since we have the GPS and the Google Map?



- What we want for SLAM solutions?
  - Fast control feedback;
  - Environment flexibility, e.g., working in high dynamic range scenes;
  - Low power consumption and hardware complexity.

# Pose Estimation and SLAM

- What sensors have been used?

- Passive sensors

- Monocular



UEye Camera

- Stereo



Point Grey Stereo Cam

- Omnidirectional



Bubl Omnicam

- Active sensors

- Lidar



Velodyne Lidar

- Time-of-flight



Mesa TOF Cam

- RGB-Depth



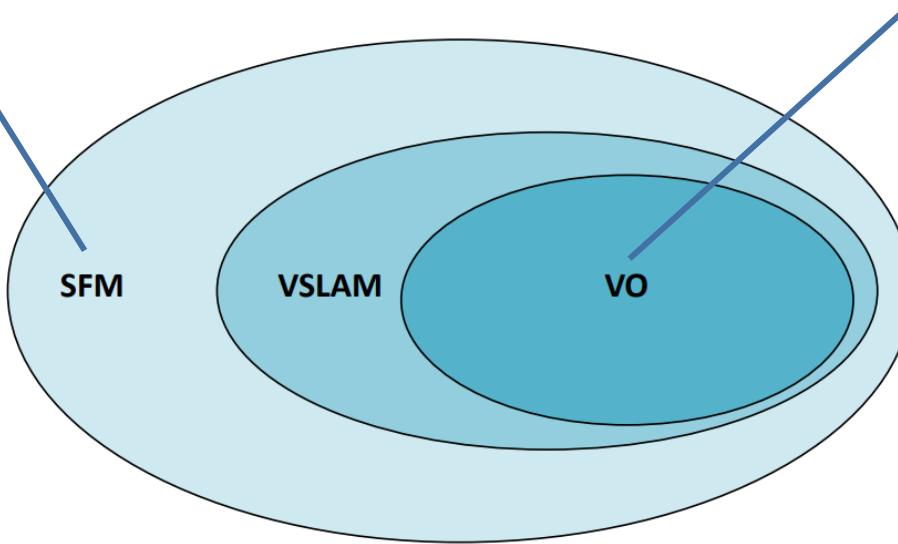
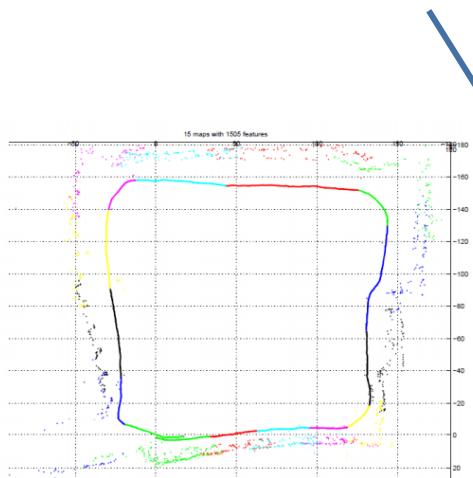
Kinect RGB-D

- Hybrid sensors

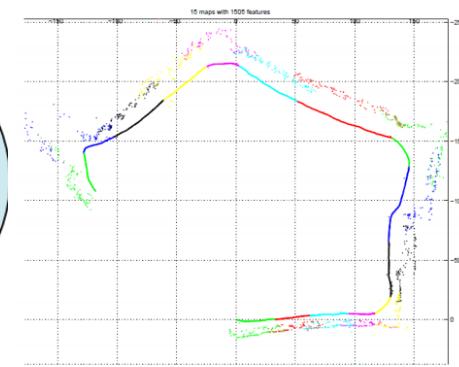
- Uses multiple sensors

# Visual Odometry (VO), VSLAM, and SFM

3-D reconstruction and 6-DOF pose estimation from **unordered image sets**



Estimating the 3D motion of the camera **sequentially and in real time**



- VO: Focusing on **local consistency** of the trajectory
- SLAM: Focusing on **global consistency** of the trajectory

VO can be used as a building block in a SLAM system.

[Robots that know what they do \(upenn.edu\)](http://Robots that know what they do (upenn.edu))  
[Robots that know what they do \(uzh.ch\)](http://Robots that know what they do (uzh.ch))

# Visual Odometry (VO)

- VO Working Principle

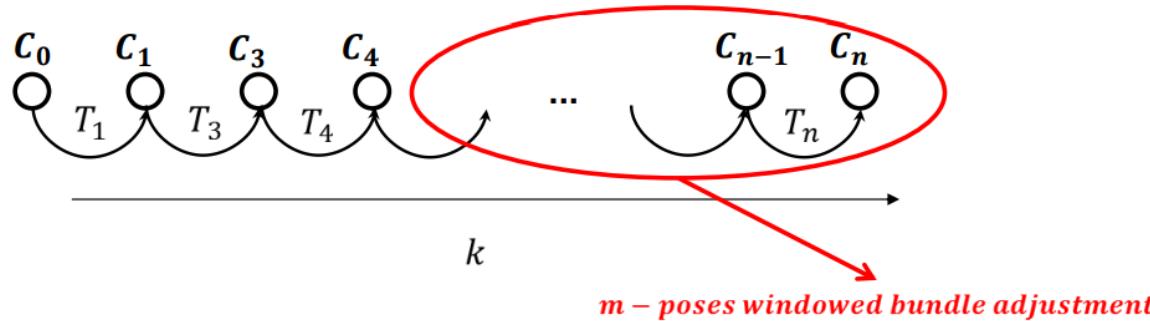
1. Compute the relative motion  $T_k$  from images  $I_{k-1}$  to image  $I_k$

$$T_k = \begin{bmatrix} R_{k,k-1} & t_{k,k-1} \\ 0 & 1 \end{bmatrix}$$

2. Concatenate them to recover the full trajectory

camera pose  $C_n = C_{n-1} T_n$  transformation

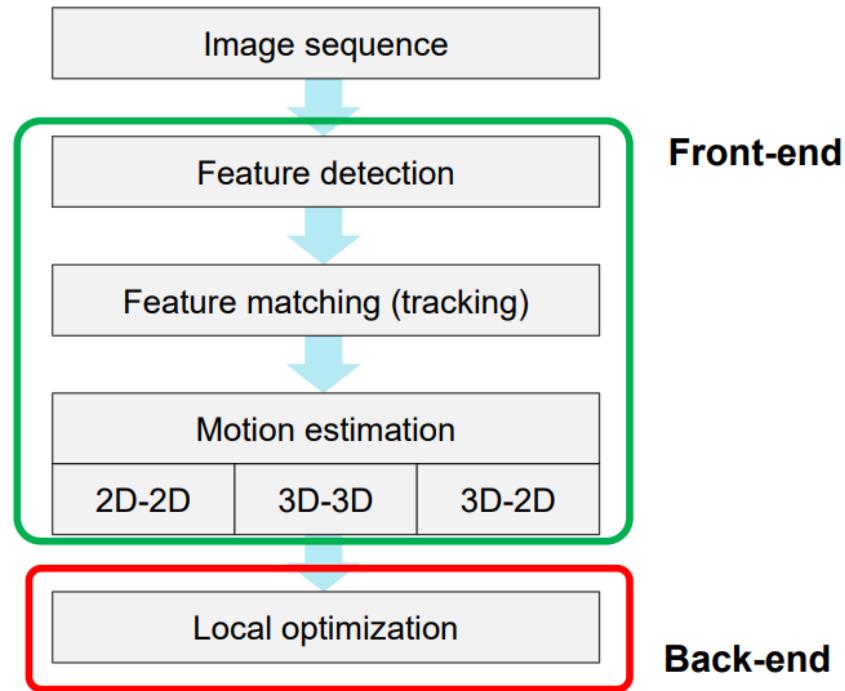
3. An optimization over the last  $m$  poses can be done to refine locally the trajectory (Pose-Graph or Bundle Adjustment)



(5) (PDF) Visual Odometry [Tutorial] (researchgate.net)

# Front-End and Back-End

- VO systems contains the front-end and the back-end.
- The **front-end** is responsible for
  - Feature extraction, matching, and outlier removal
  - Loop closure detection
  - Motion estimation between two frames
- The **back-end** is responsible for the pose and structure optimization (e.g., with libraries iSAM, g2o, Google Ceres)

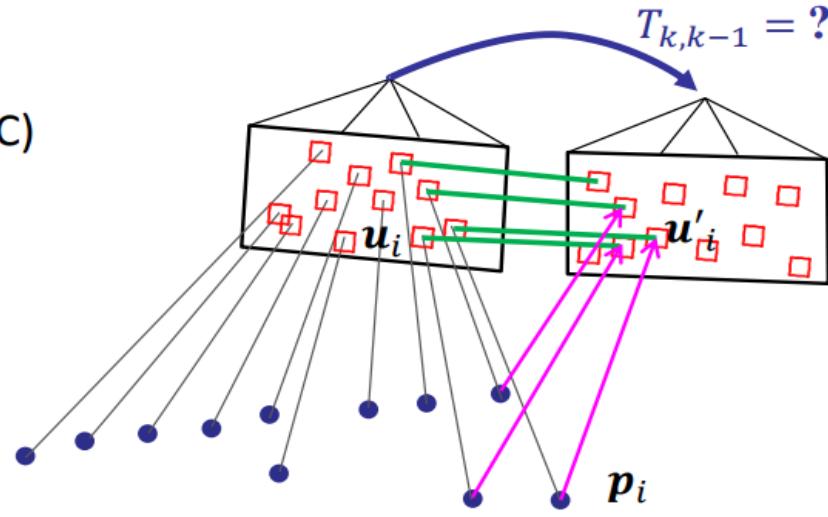


# Estimating the Relative Motion $T_k$

## Feature-based methods

1. Extract & match features (+RANSAC)
2. Minimize **Reprojection error** minimization

$$T_{k,k-1} = \arg \min_T \sum_i \| \mathbf{u}'_i - \pi(\mathbf{p}_i) \|_{\Sigma}^2$$



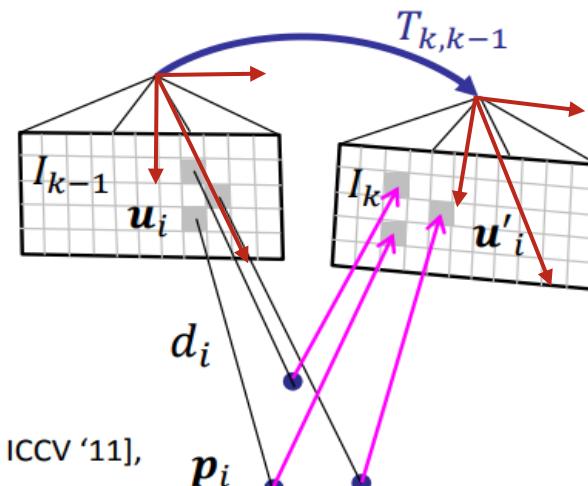
## Direct methods

1. Minimize **photometric error**

$$T_{k,k-1} = \arg \min_T \sum_i \| I_k(\mathbf{u}'_i) - I_{k-1}(\mathbf{u}_i) \|_{\sigma}^2$$

where  $\mathbf{u}'_i = \pi(T \cdot (\pi^{-1}(\mathbf{u}_i) \cdot d))$

[Jin,Favaro,Soatto'03] [Silveira, Malis, Rives, TRO'08], [Newcombe et al., ICCV '11],  
[Engel et al., ECCV'14], [Forster et al., ICRA'14]



[Robots that know what they do \(upenn.edu\)](http://Robots that know what they do (upenn.edu))

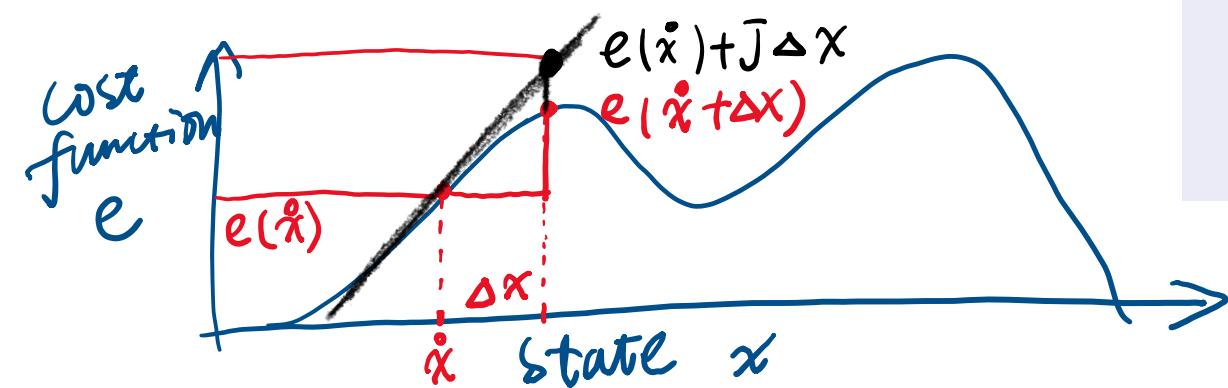
# Solving Relative Motion $\mathbf{T}_k$ by Optimization

- Graph Representation

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \sum_{<i,j> \in \mathcal{C}} \underbrace{\mathbf{e}_{ij}(\mathbf{z}_{ij}, \mathbf{x}_{ij})^\top \boldsymbol{\Omega}_{ij} \mathbf{e}_{ij}(\mathbf{z}_{ij}, \mathbf{x}_{ij})}_{\mathbf{F}_{ij}}$$

- To linearize:

$$\mathbf{e}_{ij}(\dot{\mathbf{x}}_{ij} + \Delta \mathbf{x}_{ij}) \simeq \mathbf{e}_{ij}(\dot{\mathbf{x}}_{ij}) + \mathbf{J}_{ij} \Delta \mathbf{x}_{ij}$$



error function

$$\mathbf{e}_{ij}(\mathbf{z}_{ij}, \mathbf{x}_i, \mathbf{x}_j) = \mathbf{z}_{ij} \ominus \mathbf{h}_{ij}(\mathbf{x}_i, \mathbf{x}_j)$$

$\mathbf{x}_i$ : state of vertex  $i$ ;

$\mathbf{h}_{ij}(\mathbf{x}_i, \mathbf{x}_j)$ : measurement equation;

$\mathbf{z}_{ij}$ : actual measurement.

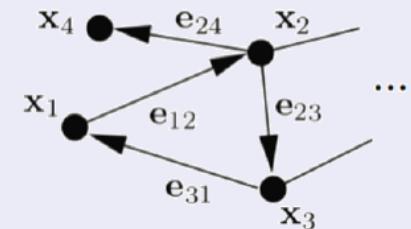
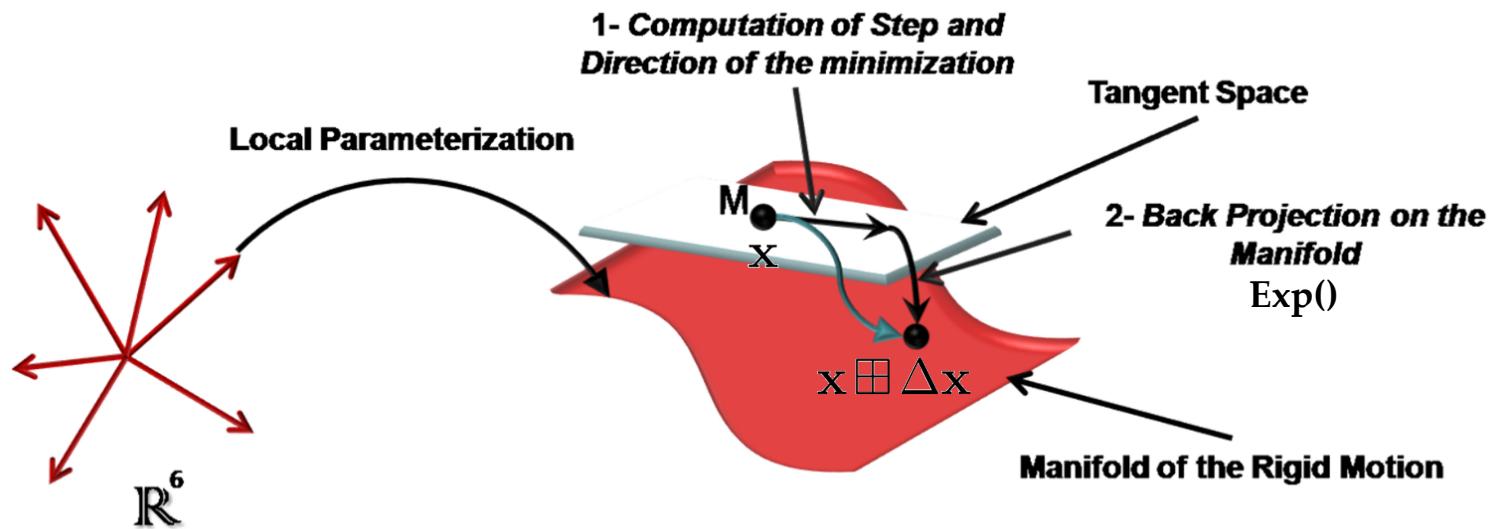
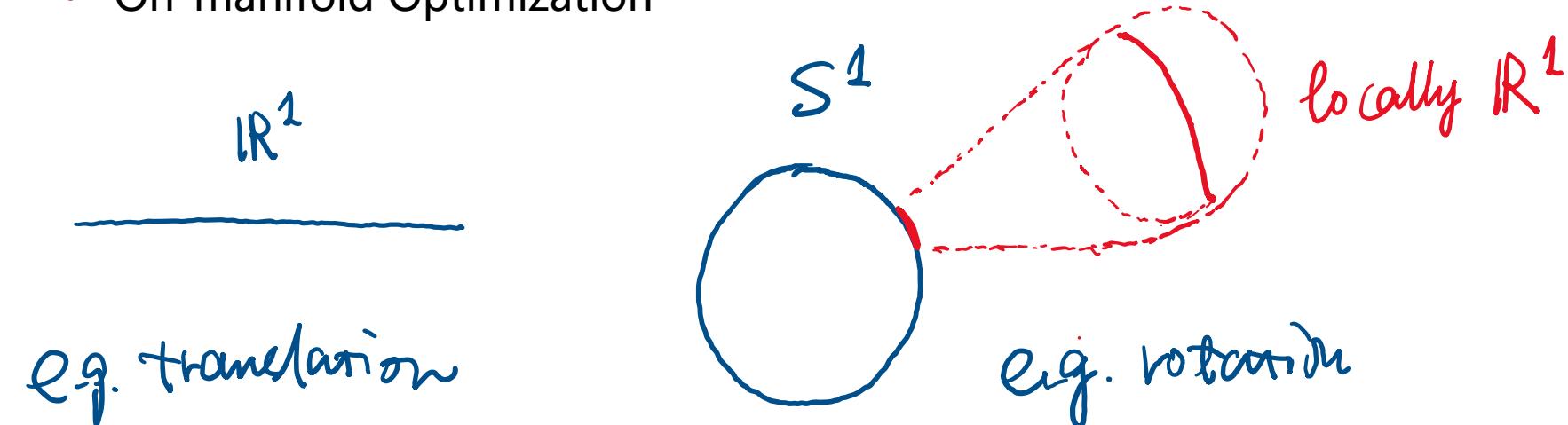


Figure: Pose graph representation [1].

# Solving Relative Motion $T_k$ by Optimization

- On-manifold Optimization

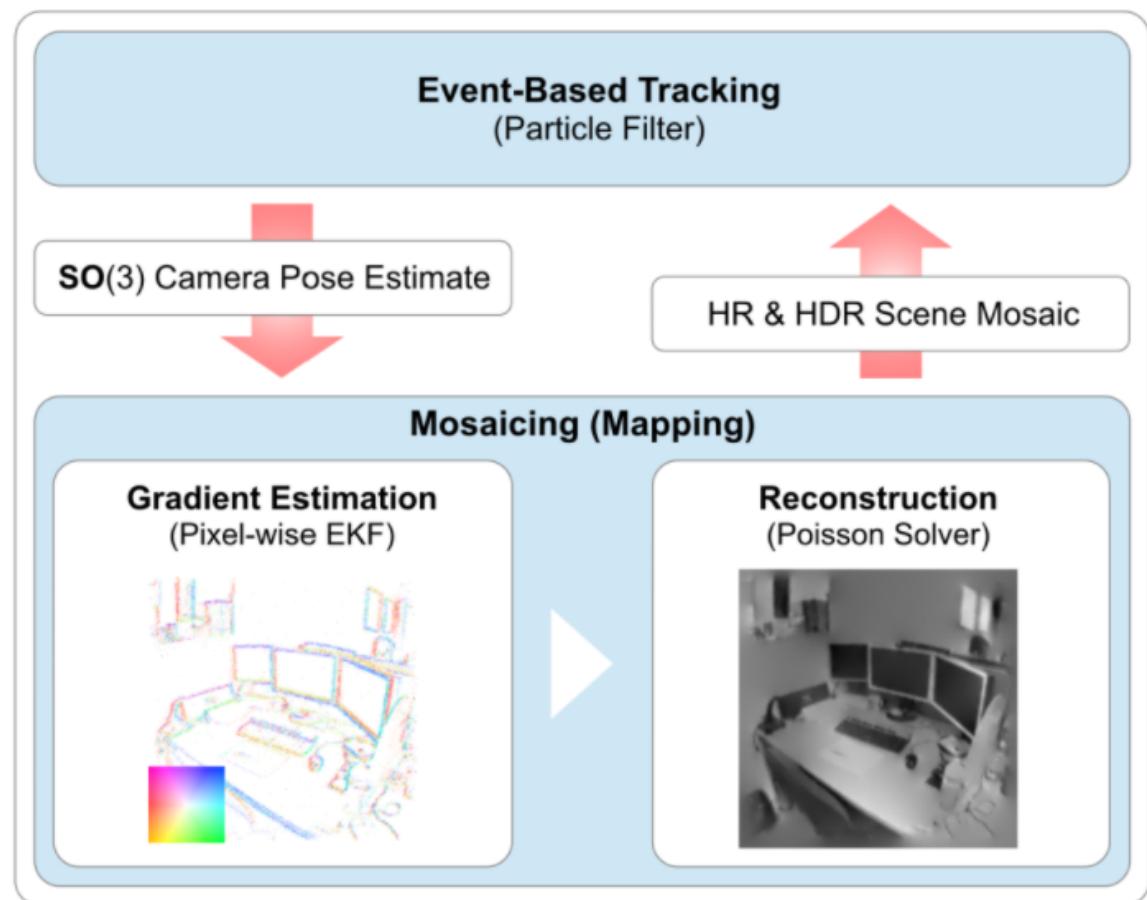
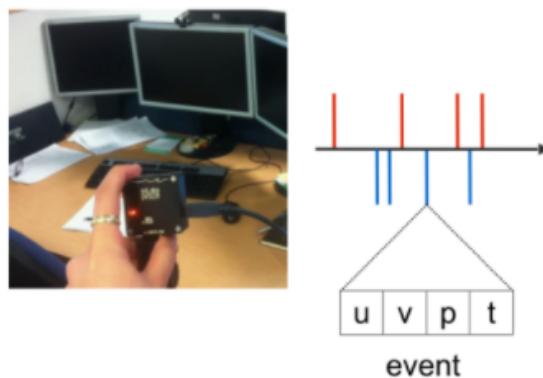


# Pose Estimation and SLAM with an Event Camera

- Challenges
  - Event streams cannot be fed directly to existing methods designed for standard cameras!
- Ideas
  - Classification by data type:  
**event frames v.s. event point clouds**
  - Classification by way of processing:  
**batch-based v.s. event-based**
- Practical considerations (besides accuracy)
  - Latency
  - Computation complexity / power consumption

# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

- Proposed algorithm
  - Simplified problem: only considers 3-DoF camera ego-motion
  - No frame involved



[Simultaneous Mosaicing and Tracking with an Event Camera \(bmva.org\)](http://bmva.org/research/papers/2018/0001.pdf)

# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

- **Event-based camera tracking based on a particle filter**

Use  $N$  particles to represent camera pose:  $\{\mathbf{p}_1^{(t)}, \mathbf{p}_2^{(t)}, \dots, \mathbf{p}_N^{(t)}\}$

3D rotation group  $\boxed{\mathbf{R}_i^{(t)} \in \mathbf{SO}(3)}$   $\boxed{w_i^{(t)}}$

Two-step process:

- 1) **Motion prediction:** constant pose motion model

$$\mathbf{R}_i^{(t)} = \mathbf{R}_i^{(t-\tau)} \exp \left( \sum_{k=1}^3 \mathbf{n}_k G_k \right)$$

Lie group generators  
Gaussian distributed noise vector

# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

- **Event-based camera tracking based on a particle filter**

## 2) Measurement update using Bayes rule

$$w_i^{(t)} = P(z|R_i^{(t)})w_i^{(t-\tau)}$$

How to define measurement  $z$ ?

Difference on log intensity at the  
same triggered pixel

$$z = \log(M(\mathbf{p}_m^{(t)})) - \log(M(\mathbf{p}_m^{(t-\tau_c)})) ,$$

$$\text{where } \mathbf{p}_m^{(t)} = \pi\left(R_i^{(t)} K^{-1} \dot{\mathbf{p}}_c\right) .$$

$K$  is the camera intrinsics;

$\pi(\mathbf{p}) = \frac{1}{p_2}(p_0, p_1)^\top$  is the projection function.

$$\dot{\mathbf{p}}_c = (u, v, 1)^\top$$

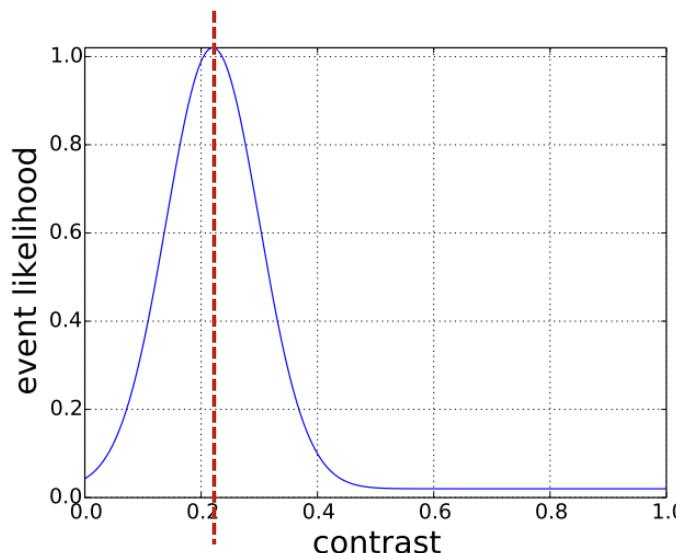
# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

- Event-based camera tracking based on a particle filter

## 2) Measurement update using Bayes rule

$$w_i^{(t)} = P(z|\mathcal{R}_i^{(t)})w_i^{(t-\tau)}$$

How to define conditional probability  $P$ ?



Theoretically, should be triggering threshold

# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

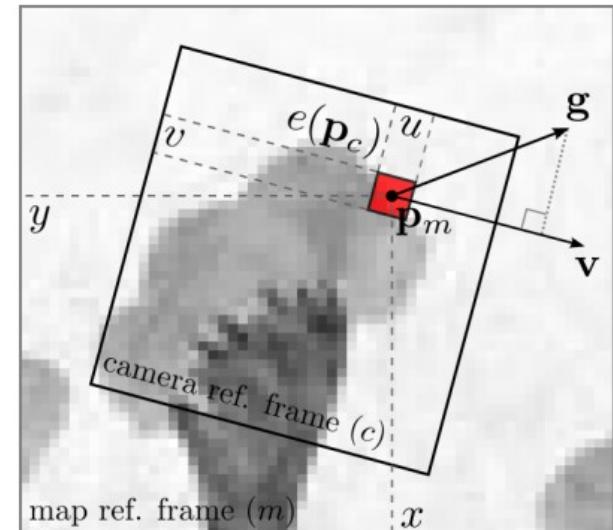
- **Mapping - Gradient estimation based on pixel-wise EKF**

Two-step process:

## 1) Process equation

Define the **gradient vector** as state vector

Apply the **constant gradient process model**



## 2) Measurement equation

$$\left( \mathbf{g}^{(t)} \cdot \mathbf{v}^{(t)} \right) \tau_c = \pm C$$

Define the measurement  $z$  as the **instantaneous event rate** at this pixel,  
we have the measurement equation:

$$h^{(t)} = \frac{\mathbf{g}^{(t)} \cdot \mathbf{v}^{(t)}}{C}$$

$$z^{(t)} = \frac{1}{\tau_c}$$

# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

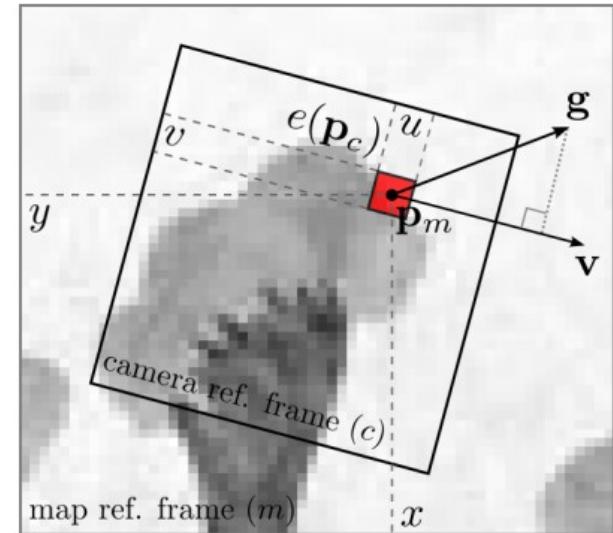
- **Mapping - Gradient estimation based on pixel-wise EKF**

Two-step process:

## 2) Measurement equation

How to find the velocity of the camera at a pixel  $\mathbf{p}_m$ ?

$$\mathbf{v}^{(t)} = \begin{pmatrix} v_x \\ v_y \end{pmatrix} = \frac{\mathbf{p}_m^{(t)} - \mathbf{p}_m^{(t-\tau_c)}}{\tau_c}$$



We can now iterate to estimate the pixel gradient based on EKF process.

# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

- **Mapping - Reconstruction from gradients**

Since we already obtained the gradients, how to reconstruct the intensity?

By optimization! Design the cost function as:

$$J(M) = \int \int (M_x - g_x)^2 + (M_y - g_y)^2 dxdy .$$

where  $M_x, M_y$  are **true gradients** in  $x$  and  $y$  directions (we don't have);  
 $g_x, g_y$  are **estimated gradients** in  $x$  and  $y$  directions (we have);

$M$  denotes the intensity that we aim to solve.

Minimizing  $J(M)$  leads to the well known Poisson equation:

$$\nabla^2 M = \frac{\partial}{\partial x} g_x + \frac{\partial}{\partial y} g_y .$$

where  $\nabla^2 M = \frac{\partial^2 M}{\partial x^2} + \frac{\partial^2 M}{\partial y^2}$  is the Laplacian.

**Fortunately, many methods are available to solve the above Poisson equation!**

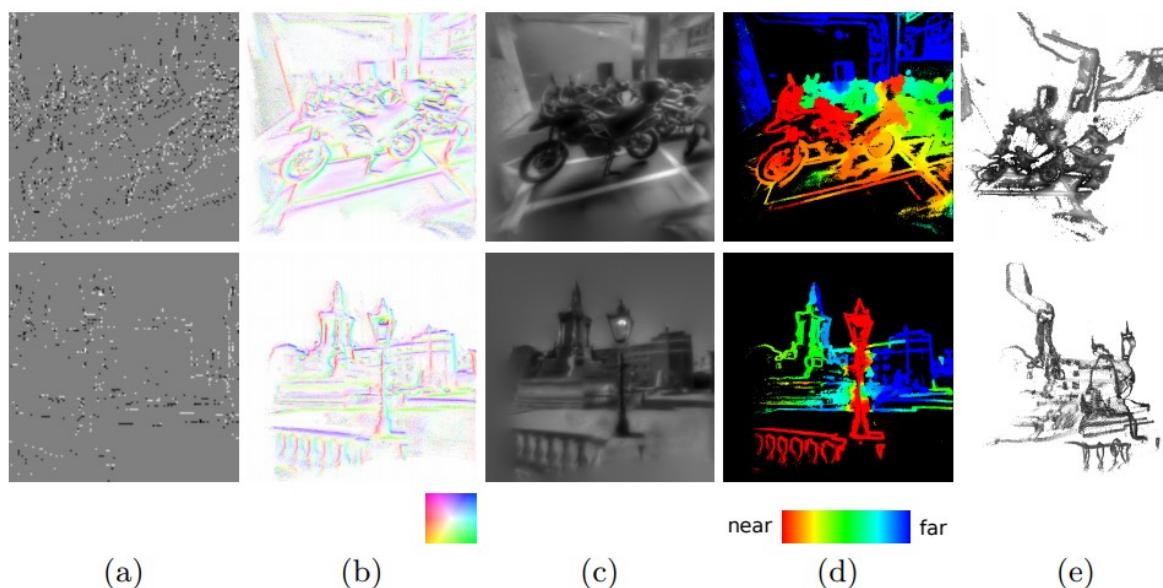
# Simultaneous Mosaicing and Tracking with an Event Camera

Hanme Kim [1], Ankur Handa [2], Ryad Benosman [3],  
Sio-Hoi Ieng [3], Andrew J. Davison [1]

- [1] Imperial College London
- [2] University of Cambridge
- [3] UPMC Univ Paris 06

# Paper: Simultaneous Mosaicing and Tracking with an Event Camera

- The authors have expanded their work into 6-DoF pose estimation
  - Based on 3 decoupled probabilistic filters, estimating **6-DoF camera motion**, scene log **intensity gradient**, and scene **inverse depth**, respectively.



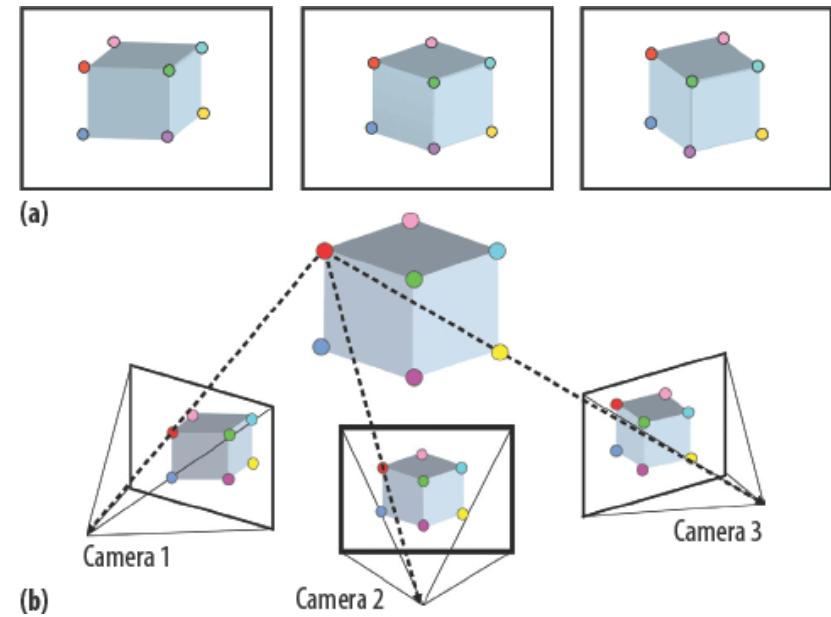
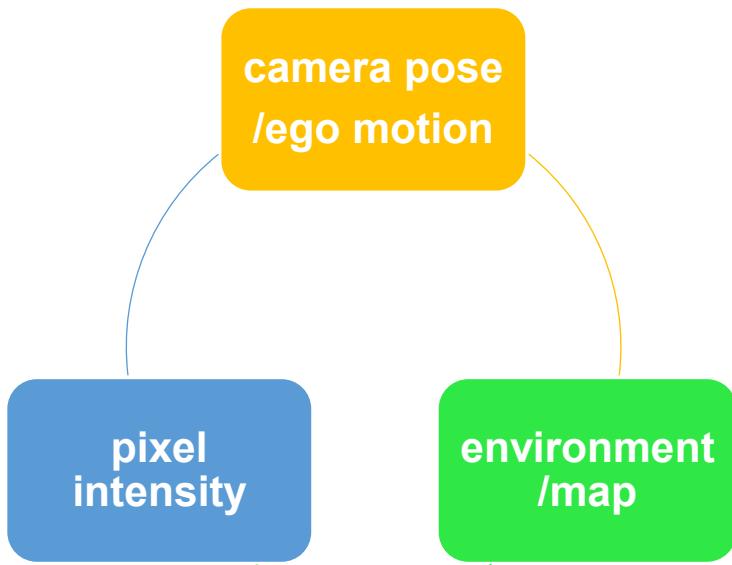
(a) visualisation of the input event stream; (b) estimated gradient keyframes; (c) reconstructed intensity keyframes with super resolution and high dynamic range properties; (d) estimated depth maps; (e) semi-dense 3D point clouds.

## Summary – Pose estimation and SLAM

- Problem formulation
- Challenges for SLAM using conventional cameras
- Two main streams for event-based SLAM
  - Taking events as frames (frame-based)
  - Taking events as point clouds (pure event-based)
- Visual odometry (VO) - the core in a SLAM system
- Examples and related work

# Image Reconstruction

- Close related to pose estimation, SLAM, and optical flow estimation.



- Here we define the **image reconstruction problem** as:
  - Estimate the **intensity** of each pixel on images
  - **Without explicitly** estimating camera's pose or environment

Usually focusing on image quality instead of the camera motion / environment modelling.

# Image Reconstruction

- Why brings event cameras into image reconstruction?
  - High speed imaging



- High dynamic range imaging



- High resolution imaging  
by reconstructing sub-pixel intensity.

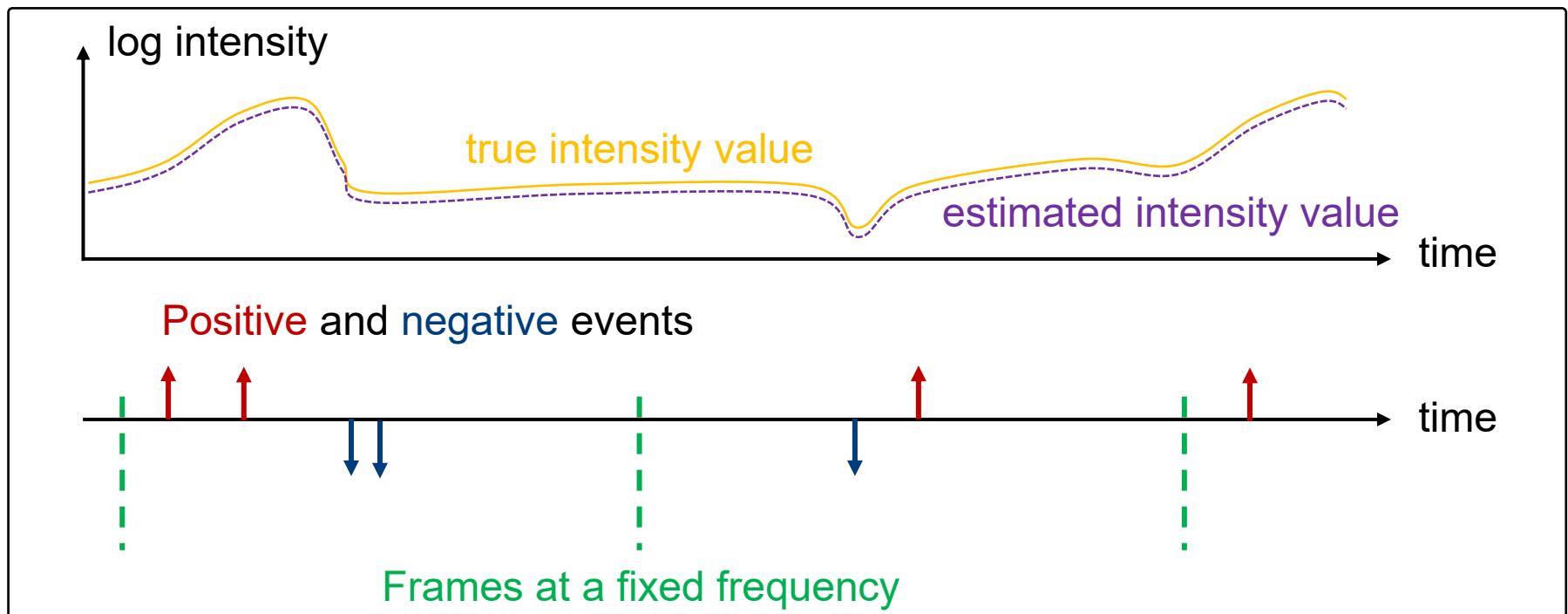


[paper066.pdf \(bmva.org\)](http://paper066.pdf (bmva.org))

# Image Reconstruction

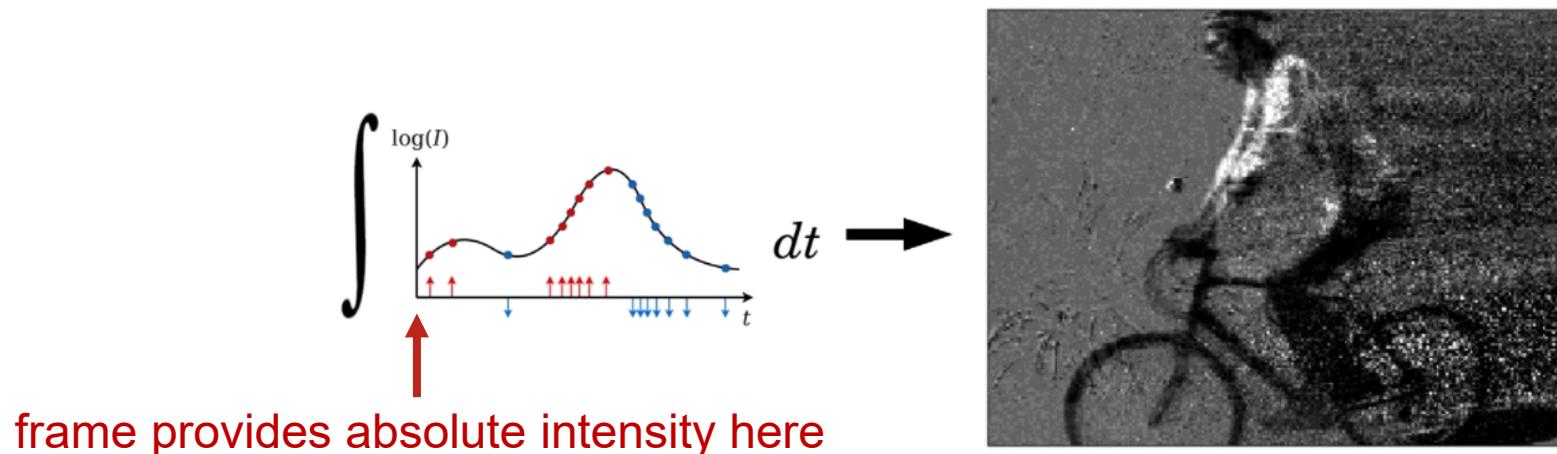
- Ideas
  - Pure event-based reconstruction
  - **Frame + event fusion**

Single-pixel frame + event fusion as a **state estimation problem**



# Paper: Continuous-time Intensity Estimation Using Event Cameras

- A naïve and simple idea would be pure integration



- Note that:
  - An event camera essentially outputs **high-frequency information**.
  - A conventional camera outputs **low-frequency information**.

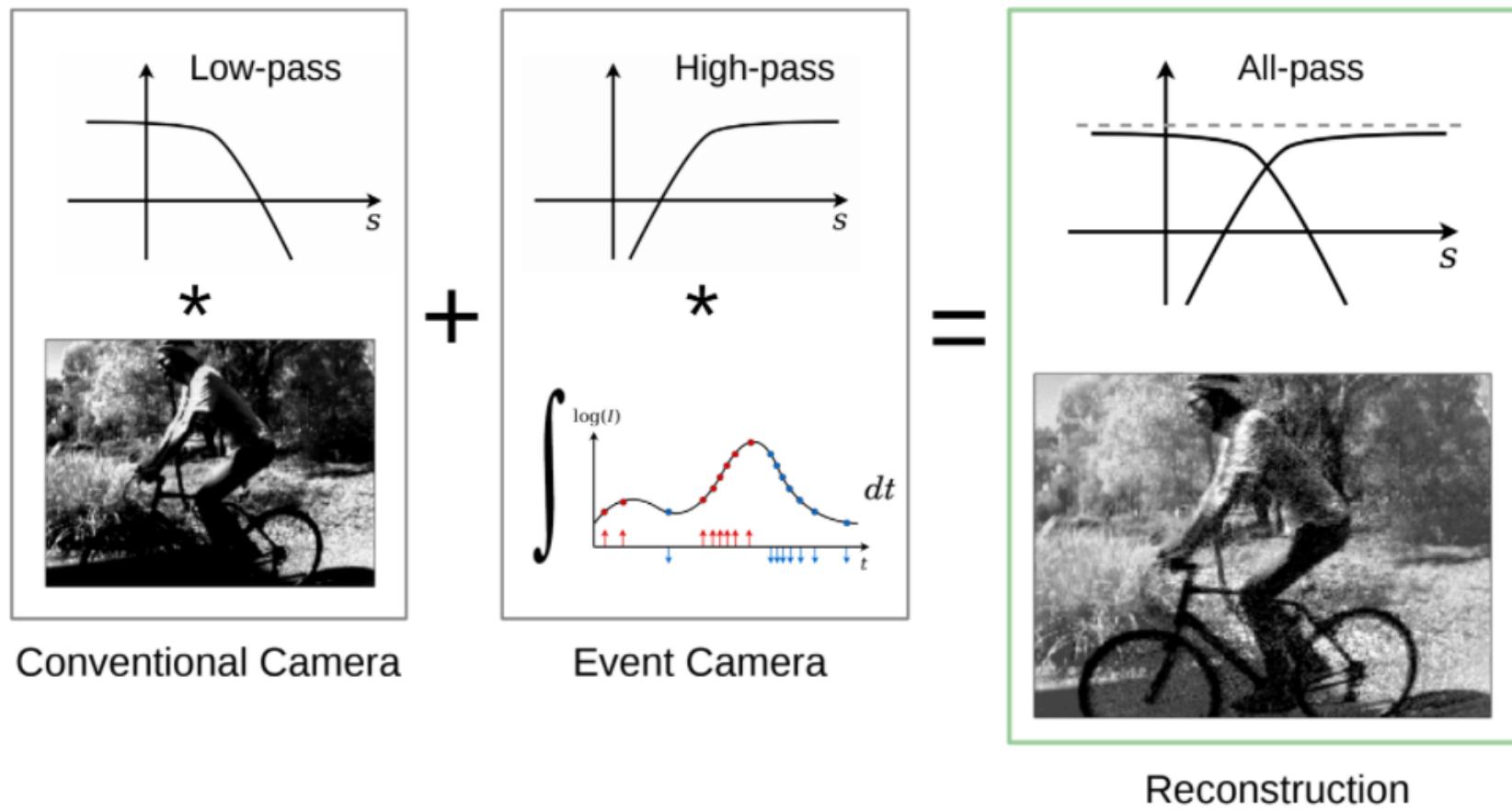
Can we fuse low-frequency information from frames  
with high-frequency information from events? And how?

<https://arxiv.org/pdf/1811.00386.pdf>

Continuous-time Intensity Estimation Using Event Cameras - Cedric Scheerlinck

# Paper: Continuous-time Intensity Estimation Using Event Cameras

- The complementary filter



[lecture\\_15.pdf \(mit.edu\)](#)

<https://arxiv.org/pdf/1811.00386.pdf>

Continuous-time Intensity Estimation Using Event Cameras - Cedric Scheerlinck



# Continuous-time Intensity Estimation Using Event Cameras

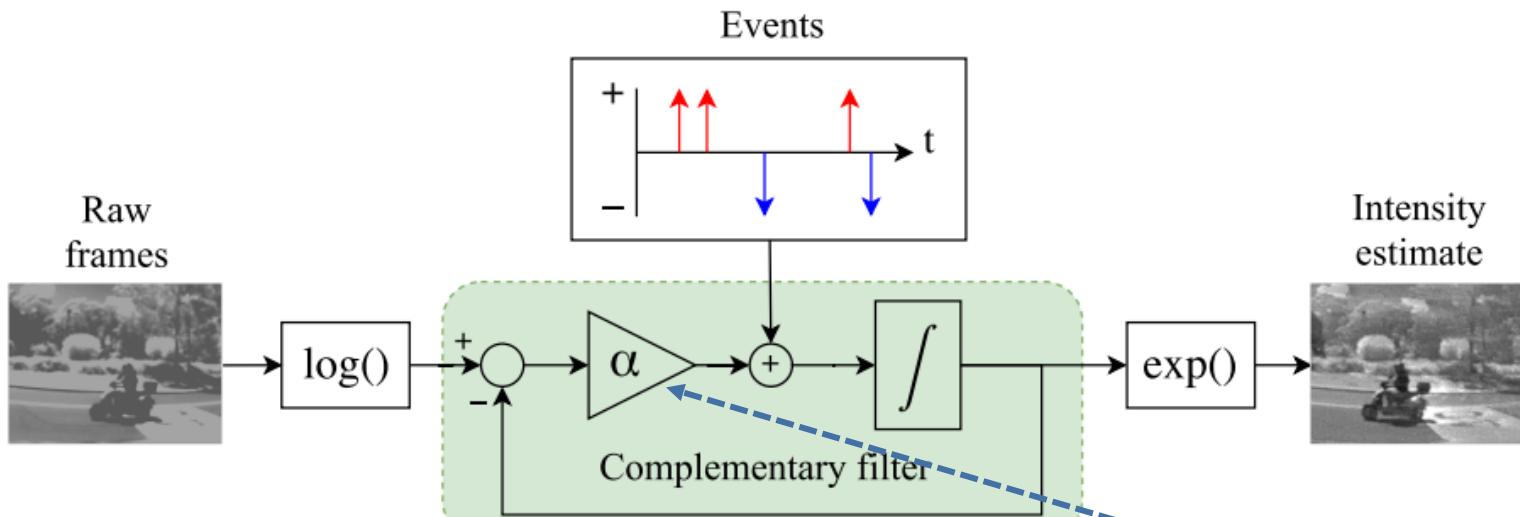
Cedric Scheerlinck, Nick Barnes, Robert Mahony



Australian  
National  
University



# Paper: An Asynchronous Kalman Filter for Hybrid Event Cameras

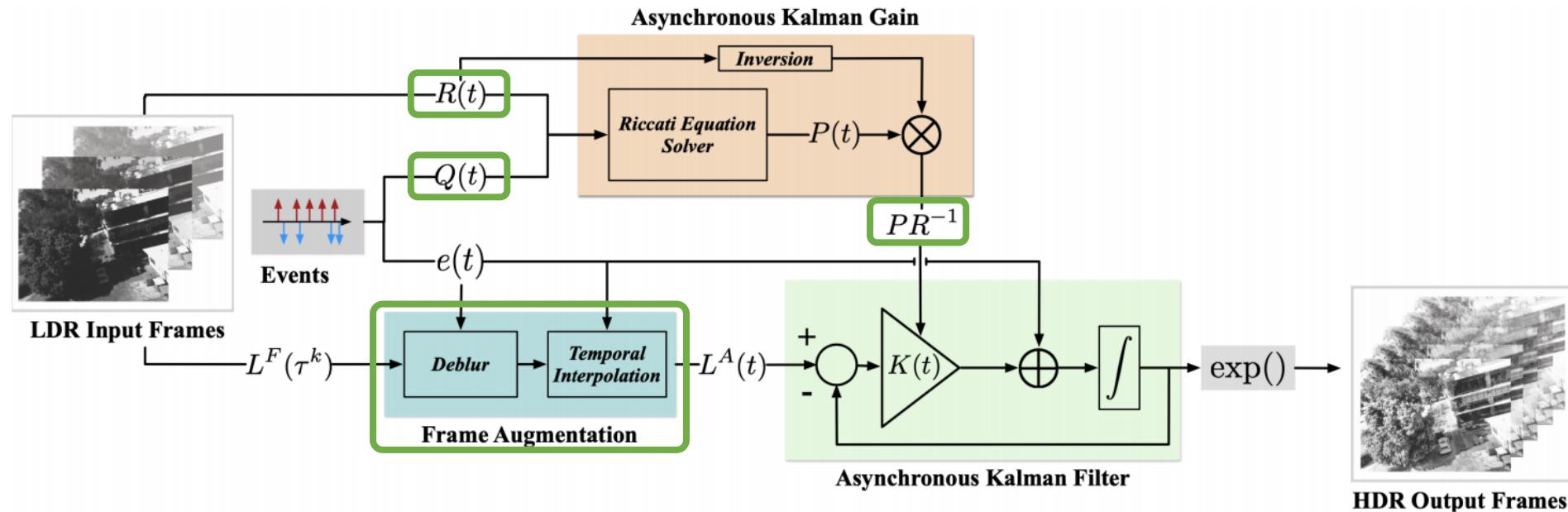


- A Complementary Filter has been proposed. There is a **filter parameter** that requires manual selection. Practical issues:
  - Hard to decide the parameter for different scenarios.
  - Not adaptive to scene changes.

Is it possible to set the filter parameter automatically?

# Paper: An Asynchronous Kalman Filter for Hybrid Event Cameras

- Framework



- Contributions

- Varying and adaptive filter gain
- EVS + CIS uncertainty models
- Deblur and temporal interpolation for frame augmentation

# Paper: An Asynchronous Kalman Filter for Hybrid Event Cameras

- Event Camera (Noise) Model

$$e_{\mathbf{p}}(t) = \sum_{i=1}^{\infty} (c\sigma_{\mathbf{p}}^i + \eta_{\mathbf{p}}^i) \delta(t - t_{\mathbf{p}}^i)$$

$\eta_{\mathbf{p}}^i \sim \mathcal{N}(0, Q_{\mathbf{p}}(t))$

Noise term

$$Q_{\mathbf{p}}(t) := \sum_{i=1}^{\infty} (Q_{\mathbf{p}}^{\text{proc.}}(t) + Q_{\mathbf{p}}^{\text{iso.}}(t) + Q_{\mathbf{p}}^{\text{ref.}}(t)) \delta(t - t_{\mathbf{p}}^i)$$

Process noise

Isolated pixel noise

Refractory period noise

# Paper: An Asynchronous Kalman Filter for Hybrid Event Cameras

- Event Camera (Noise) Model
  - Process noise

$$Q_{\mathbf{p}}^{\text{proc.}}(t_{\mathbf{p}}^i) = \sigma_{\text{proc.}}^2(t_{\mathbf{p}}^i - t_{\mathbf{p}}^{i-1})$$

- Isolated pixel noise

$$Q_{\mathbf{p}}^{\text{iso.}}(t_{\mathbf{p}}^i) = \sigma_{\text{iso.}}^2 \min\{t_{\mathbf{p}}^i - t_{N(\mathbf{p})}^*\}$$

timestamp of neighborhood events

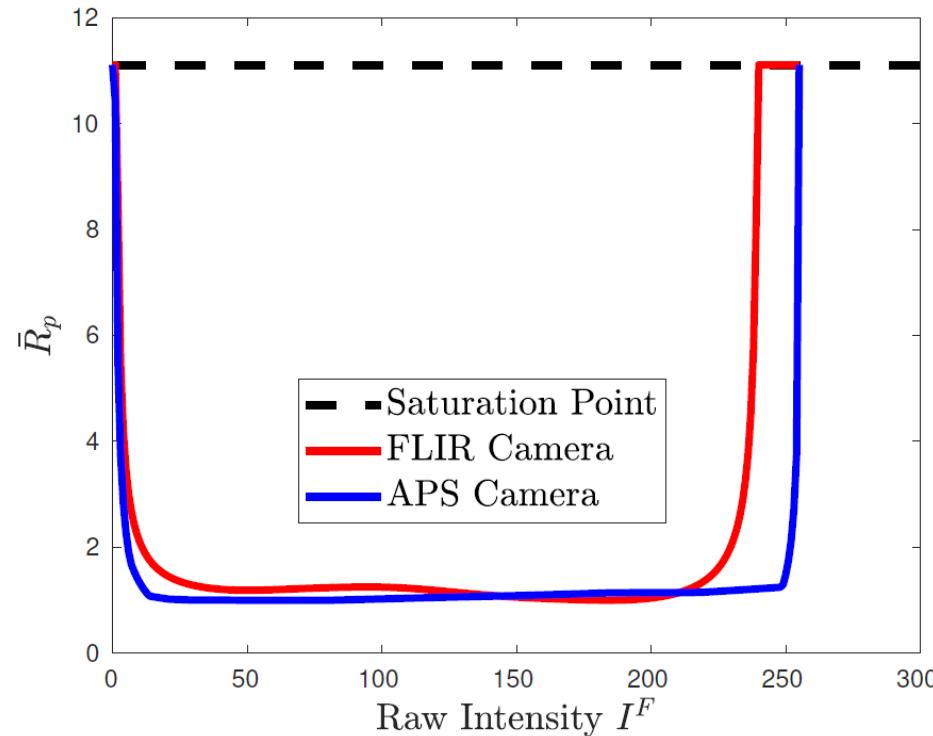
- Refractory period noise

$$Q_{\mathbf{p}}^{\text{ref.}}(t_{\mathbf{p}}^i) = \begin{cases} 0 & \text{if } t_{\mathbf{p}}^i - t_{\mathbf{p}}^{i-1} > \rho \\ \sigma_{\text{ref.}}^2 & \text{if } t_{\mathbf{p}}^i - t_{\mathbf{p}}^{i-1} \leq \bar{\rho} \end{cases}$$

refractory time

# Paper: An Asynchronous Kalman Filter for Hybrid Event Cameras

- Conventional Camera (Noise) Model
  - Gaussian assumption, covariance related to output intensity:  
Over-exposed/under-exposed areas → Larger covariance  
Normally-exposed areas → Smaller covariance





# An Asynchronous Kalman Filter for Hybrid Event Cameras

Ziwei Wang, Yonhon Ng, Cedric Scheerlinck,  
Robert Mahony

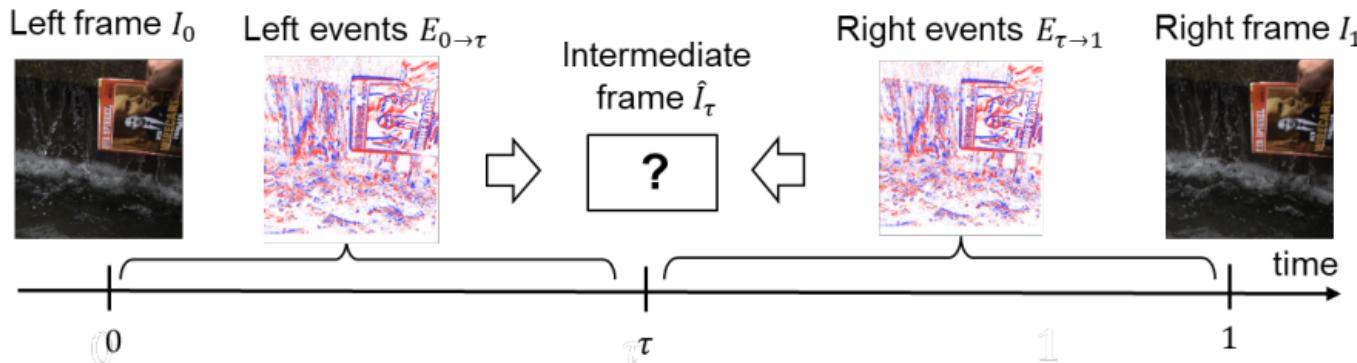


Australian  
National  
University

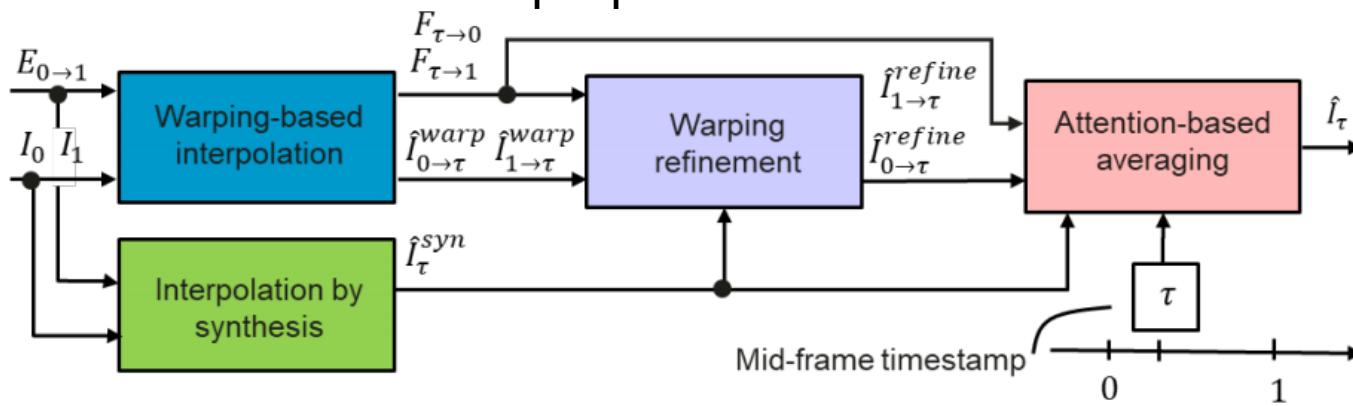
Systems Theory &  
Robotics Group

# Paper: TimeLens: Event-based Video Frame Interpolation

- Deep learning-based approach.
- Proposed event-based interpolation approach utilizes **events and frames from two directions**.

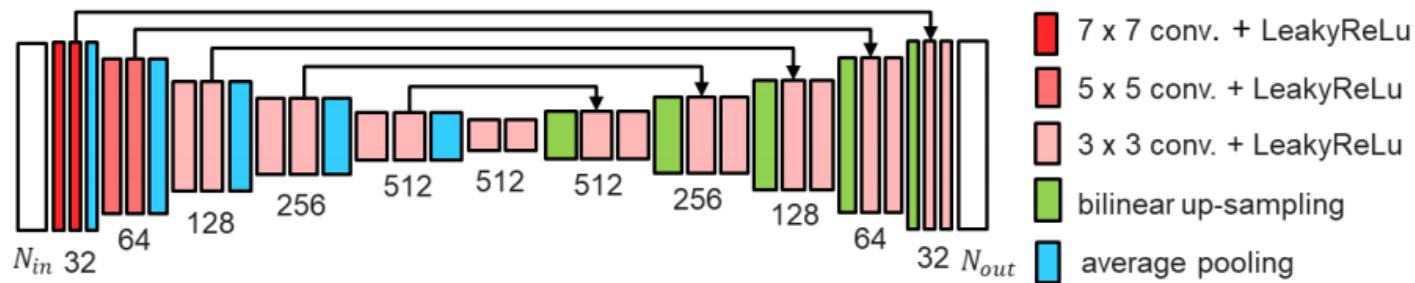


- Overview framework of the proposed method:

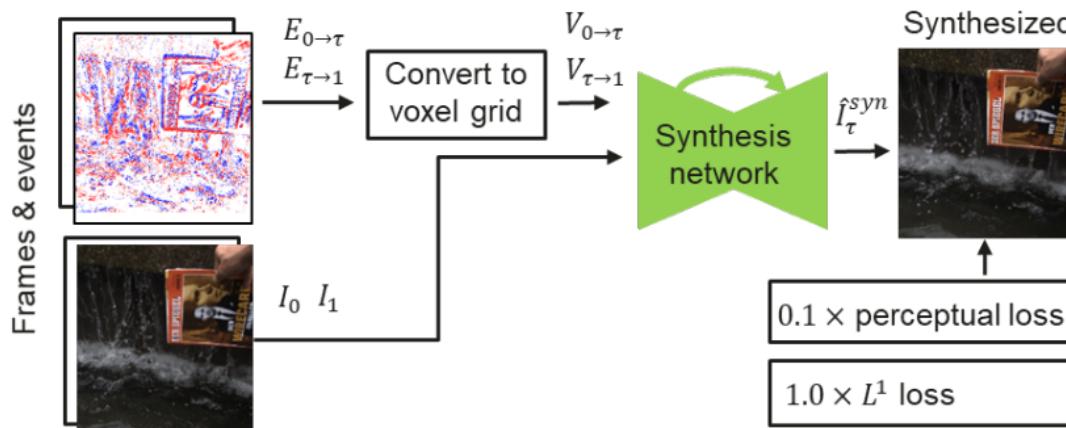


# Paper: TimeLens: Event-based Video Frame Interpolation

- Backbone network architecture: **hourglass network**

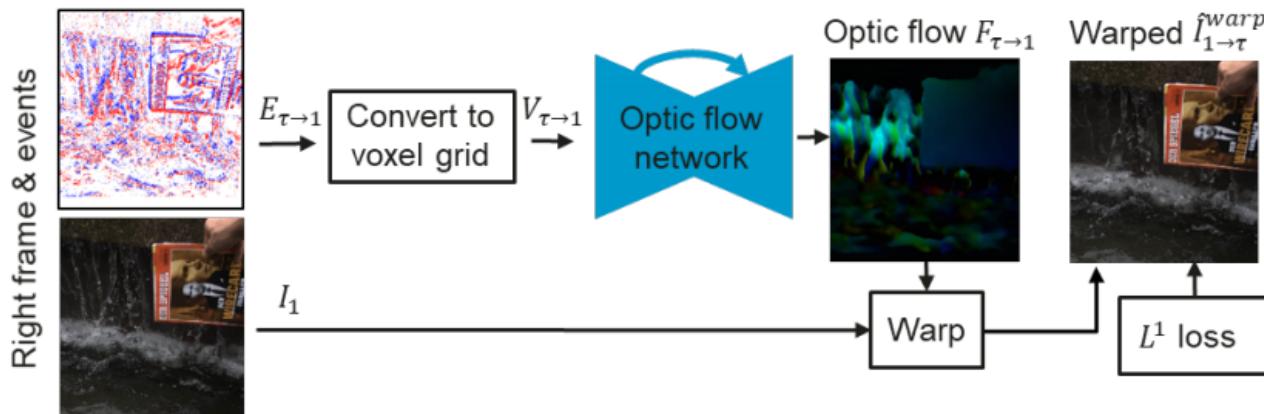


- Interpolation by synthesis module

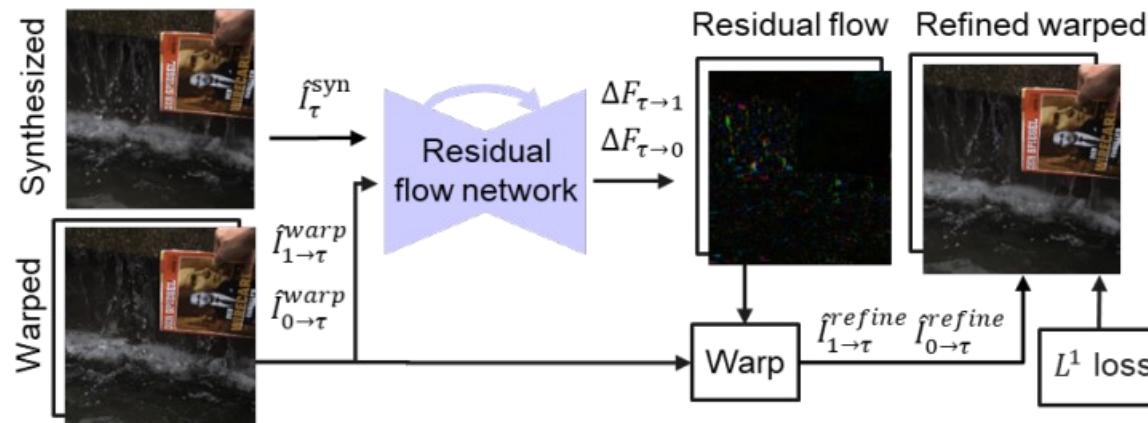


# Paper: TimeLens: Event-based Video Frame Interpolation

- Warping-based interpolation module

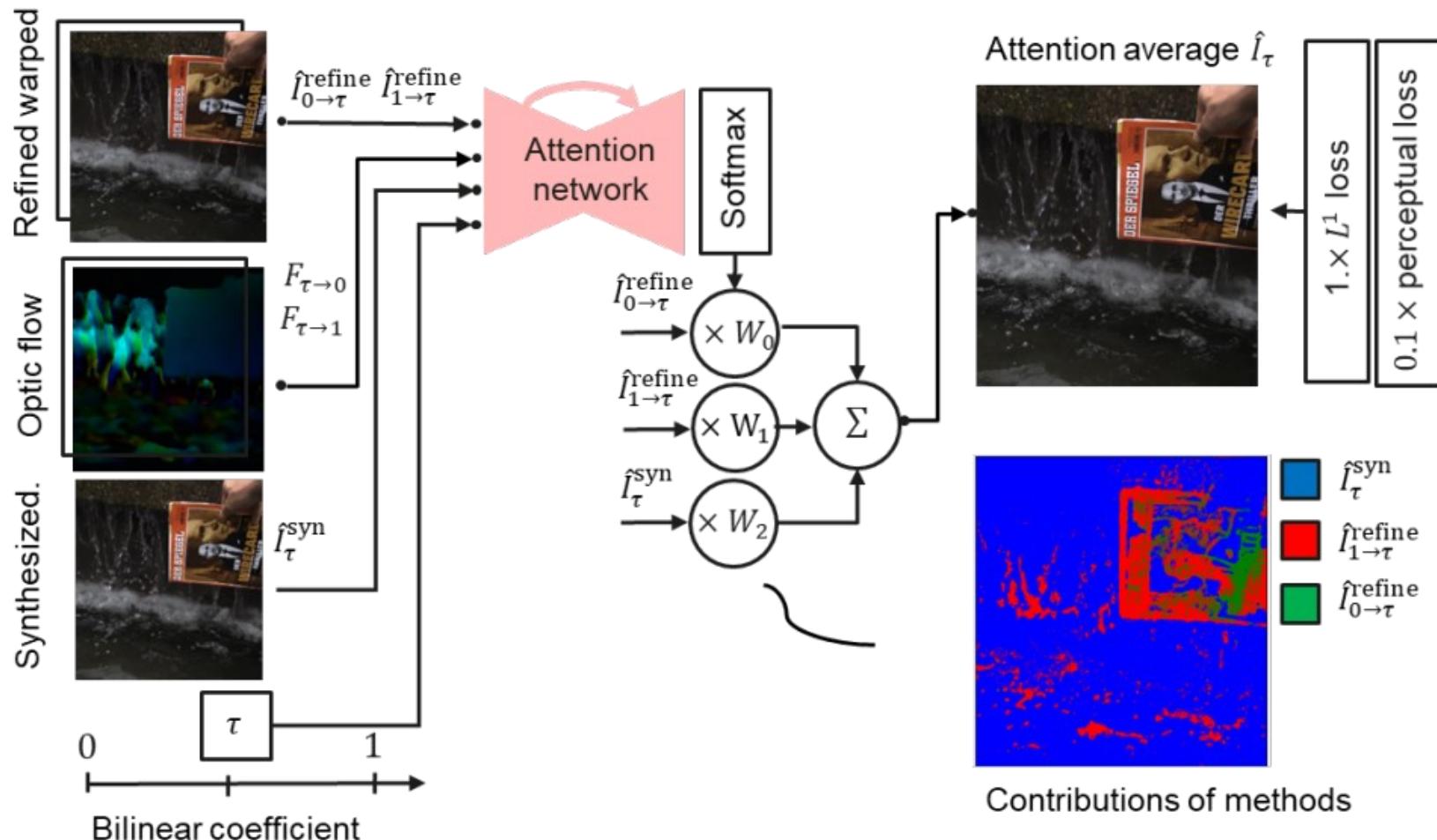


- Warping refinement module



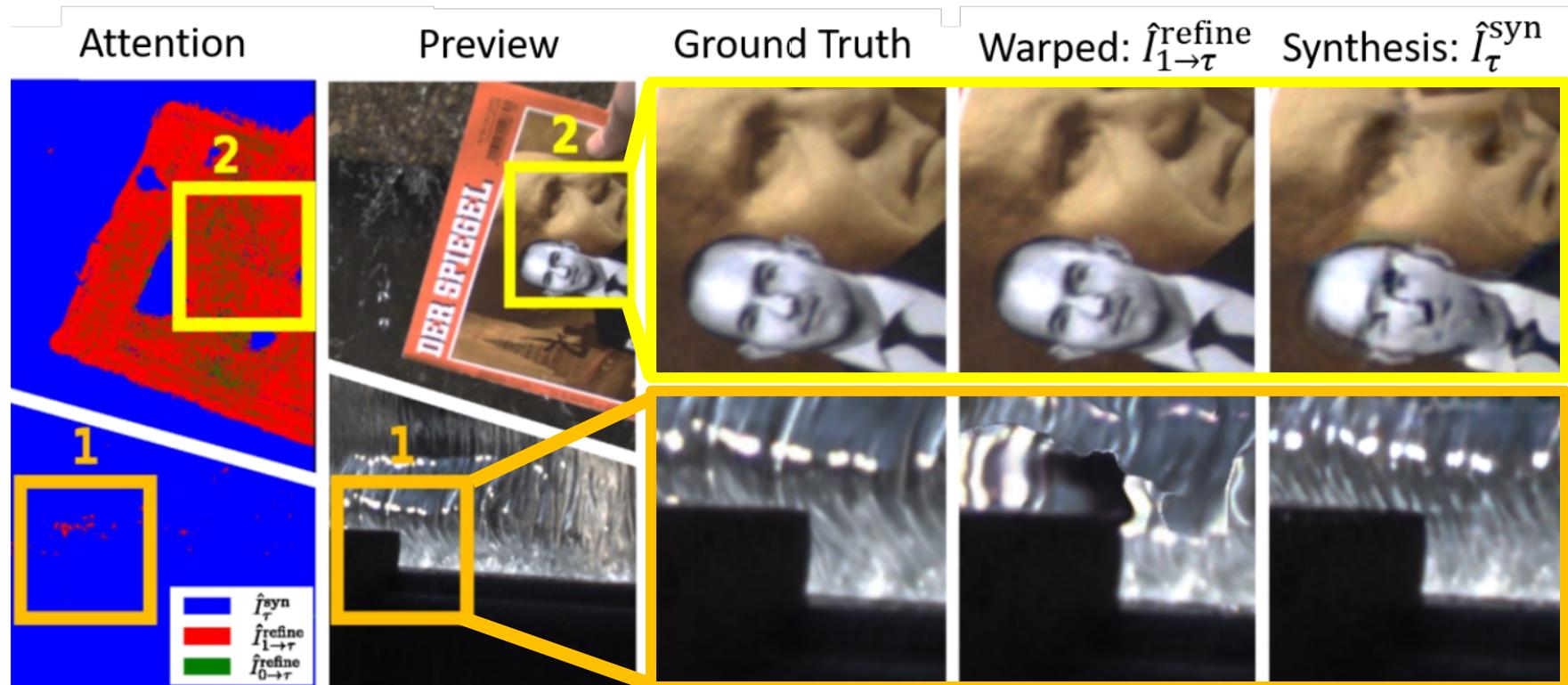
# Paper: TimeLens: Event-based Video Frame Interpolation

- Attention-based averaging module



# Paper: TimeLens: Event-based Video Frame Interpolation

- Complementarity of warping- and synthesis-based interpolation.
  - The attention network learns where to trust synthesized interpolation results, and where to trust warped interpolation results.





## TimeLens: Event-based Video Frame Interpolation

Stepan Tulyakov\*, Daniel Gehrig\*, Stamatios Geourgoulis,  
Julius Erbach, Mathias Gehrig, Yuanyou Li, Davide Scaramuzza

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



University of  
Zurich<sup>UZH</sup>



\* these authors contributed equally

# Paper: Back to Event Basics: Self-Supervised Learning of Image Reconstruction for Event Cameras via Photometric Constancy

- One issue for learning-based methods: **lack of data**. Possible reasons:
  - The **cost** of event sensor is high;
  - The **data format** of event sensor is not unified.
  - ...?
- How to reduce the dependence of data in learning-based methods?
  - Self-supervised learning/unsupervised learning.  
e.g., EV-FlowNet for optical flow estimation:  
No ground-truth optical flow is used in training; Ground-truth CIS images are required.

Can we train a network using events only?



## Back to Event Basics: Self-Supervised Learning of Image Reconstruction for Event Cameras via Photometric Constancy

Federico Paredes-Valles and Guido C. H. E. de Croon  
(Poster Session Three, ID: 8305)

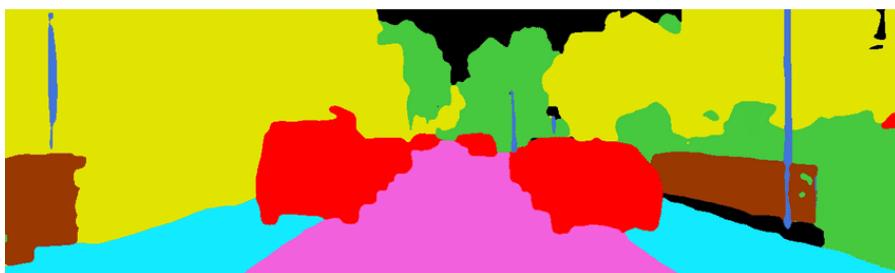


## Summary – Image reconstruction

- Image reconstruction: closely related to pose estimation, SLAM, and optical flow estimation.
  - Modelling-based approaches
  - Learning-based approaches

# Motion Segmentation

- **Image segmentation** is
  - the process of partitioning a digital image into multiple segments.
  - the process of assigning a label to every pixel in an image such that pixels with the same label share certain characteristics.



|      |            |          |         |
|------|------------|----------|---------|
| Road | Sidewalk   | Building | Fence   |
| Pole | Vegetation | Vehicle  | Unlabel |



(4) (PDF) Towards a Meaningful 3D Map Using a 3D Lidar and a Camera (researchgate.net)  
Deploying robotics applications — Using image segmentation to drive the Turtlebot3 Burger - Qualcomm Developer Network

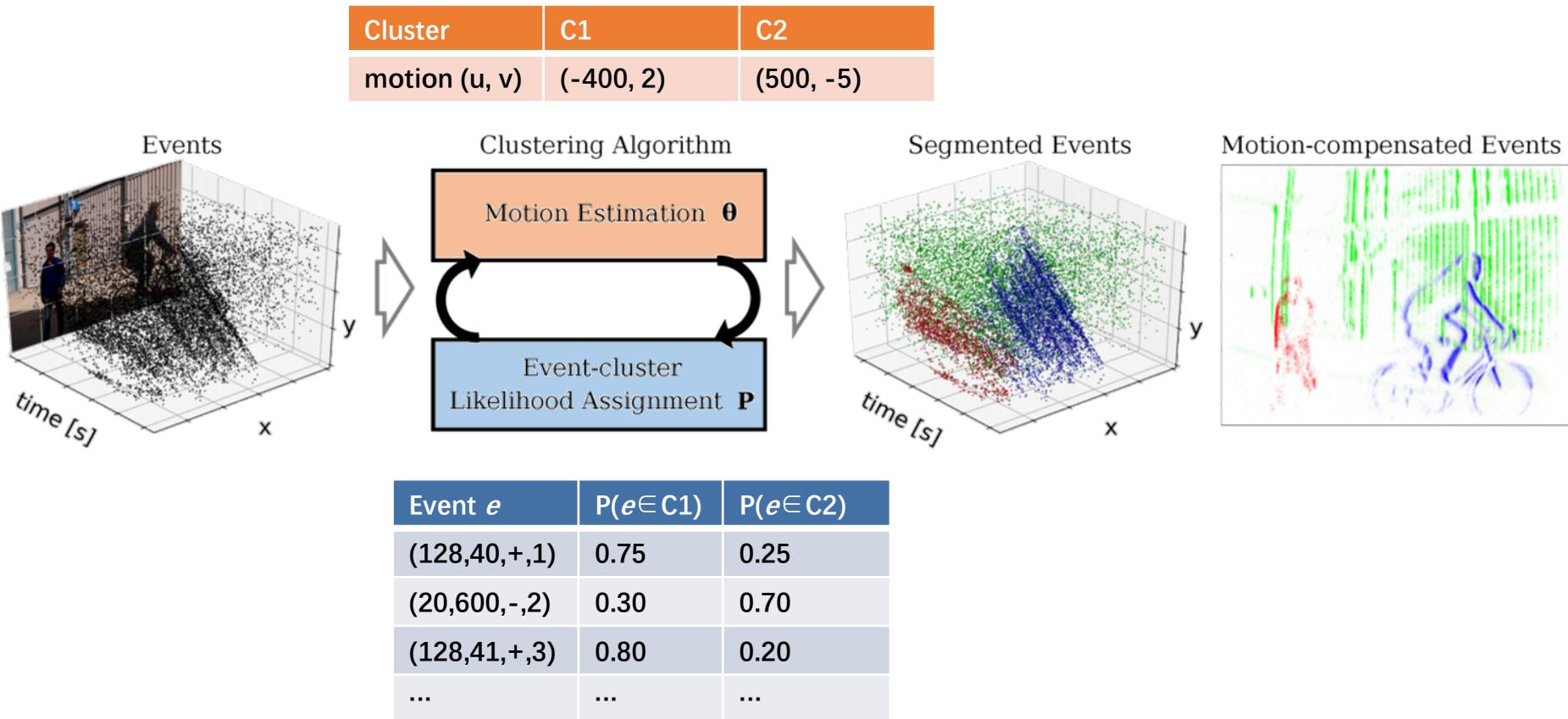
# Motion Segmentation

- Image segmentation with an event camera
  - Idea 1: **fuse events and frames**: combine motion patterns and texture/geometry patterns
  - Idea 2: use **events only**: segmenting objects based on motion patterns  
→ “**motion segmentation**”

# Paper: Event-Based Motion Segmentation by Motion Compensation

- Framework of motion segmentation

Idea: maximizing event alignment, i.e., maximizing the sharpness of motion compensated images (one per cluster) of warped events.



# Paper: Event-Based Motion Segmentation by Motion Compensation

- Mathematical formulation
  - **Soft event-cluster associations:**

$$p_{kj} = P(e_k \in \ell_j)$$

is the probability of the  $k$ -th event belonging to the  $j$ -th cluster.

Let  $\mathbf{P} \equiv (p_{kj})$  be an  $N_e \times N_\ell$  matrix with all event-cluster associations. Define the weighted IWE (Image of Warped Events) of the  $j$ -th cluster as

$$I_j(\mathbf{x}) \doteq \sum_{k=1}^{N_e} p_{kj} \delta(\mathbf{x} - \mathbf{x}'_{kj}),$$

with  $\mathbf{x}'_{kj} = \mathbf{W}(\mathbf{x}_k, t_k; \boldsymbol{\theta}_j)$  the warped event location, and  $\delta$  the Dirac delta. Above equation states that events are warped:

$$e_k \doteq (\mathbf{x}_k, t_k, s_k) \mapsto e'_k \doteq (\mathbf{x}'_k, t_{\text{ref}}, s_k)$$

and probabilities are accumulated at the warped locations.

# Paper: Event-Based Motion Segmentation by Motion Compensation

We want to maximize the sharpness, but how to define the sharpness?

- Define the sharpness as the variance of the image:

$$\text{Var}(I_j) \doteq \frac{1}{|\Omega|} \int_{\Omega} (I_j(\mathbf{x}) - \mu_{I_j})^2 d\mathbf{x}$$

where  $\mu_{I_j}$  is the mean of IWE over the image plane  $\Omega$ .

**Problem:** to find the associations  $\mathbf{P}$  and cluster parameters  $\boldsymbol{\theta}$  that maximize the sum of contrasts of all clusters:

$$(\boldsymbol{\theta}^*, \mathbf{P}^*) = \arg \max_{(\boldsymbol{\theta}, \mathbf{P})} \sum_{j=1}^{N_\ell} \text{Var}(I_j)$$

No closed-form solution → Only iterative, alternating optimization!

# Paper: Event-Based Motion Segmentation by Motion Compensation

- Alternating Optimization
  - Each iteration contains 2 steps

The proposed alternating method **converges locally** (i.e., there is no guarantee of convergence to a global solution).

If associations  $\mathbf{P}$  are fixed, update the motion parameters using

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mu \nabla_{\boldsymbol{\theta}} \left( \sum_{j=1}^{N_\ell} \text{Var}(I_j) \right) \quad (6)$$

“gradient ascent”

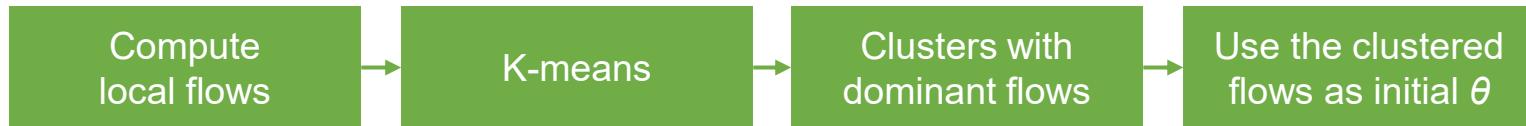
Fix motion parameters  $\boldsymbol{\theta}$ , refine the associations based on probability partitioning law:

$$p_{kj} = \frac{c_j(\mathbf{x}'_k(\boldsymbol{\theta}_j))}{\sum_{i=1}^{N_\ell} c_i(\mathbf{x}'_k(\boldsymbol{\theta}_i))} \quad (7)$$

where  $c_j(\mathbf{x}) \doteq I_j(\mathbf{x})$ .

# Paper: Event-Based Motion Segmentation by Motion Compensation

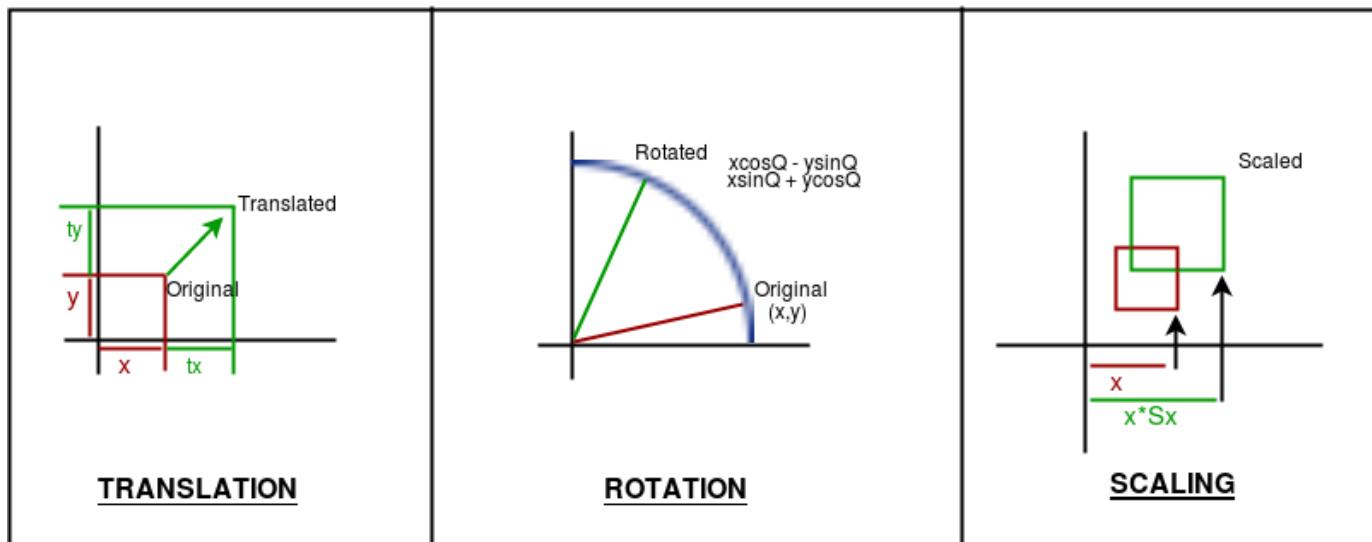
- Initialization
  - Requires initialization of  $\theta, P$  to start the iteration.
  - Several initialization schemes are possible, depending on the motion models.
    1. K-means initialization for **optical flow-type motion**;



2. Greedy initialization for **all kinds of motion**.  
No details/codes released from the authors.

# Paper: Event-Based Motion Segmentation by Motion Compensation

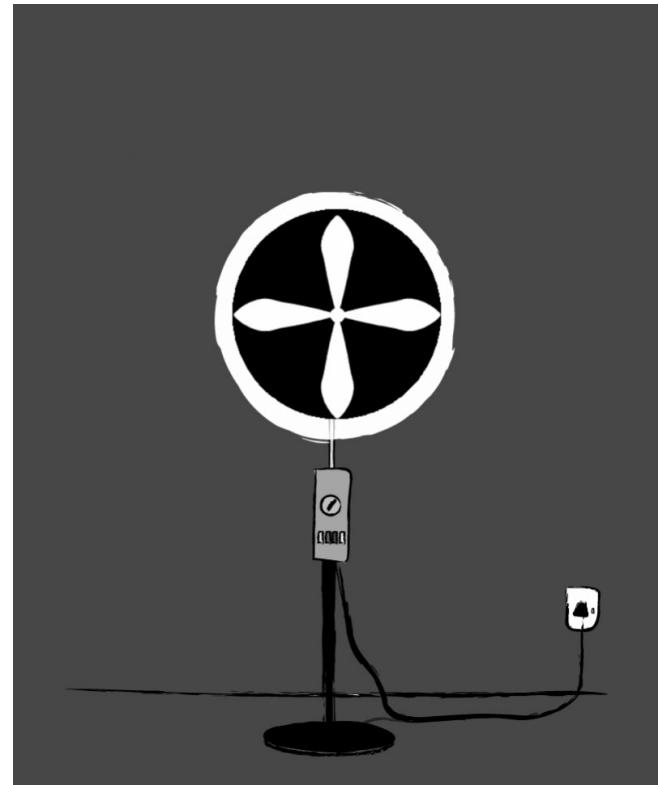
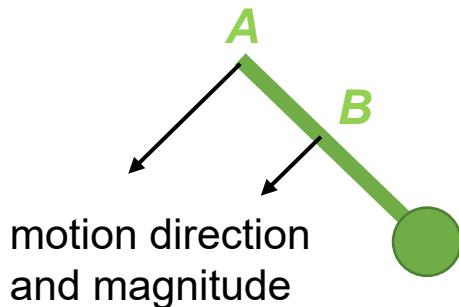
- Discussion on **motion model**  $\mathbf{x}'_{kj} = \mathbf{W}(\mathbf{x}_k, t_k; \theta_j)$   
 $e_k \doteq (\mathbf{x}_k, t_k, s_k) \mapsto e'_k \doteq (\mathbf{x}'_k, t_{\text{ref}}, s_k)$ 
  - Optical flow type motion (translation only)
  - Other types of motion (both translation and rotation)



An appropriate motion model helps to **reduce the cluster number!**

# Paper: Event-Based Motion Segmentation by Motion Compensation

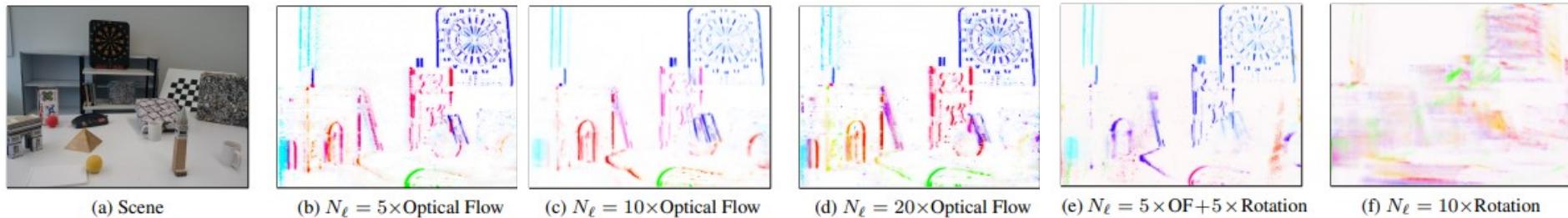
- Discussion on **motion model**  
Consider the rotating fan scenario.  
**What motion model should be used?**
- Do points A and B share the same motion parameter, if using an optical flow type model?
- What if using a rotation model instead?



All points on the edge of blade have **different motion parameters**, if using an flow type model.  
However, they do share the **same motion parameters**, if using a rotation model!

# Paper: Event-Based Motion Segmentation by Motion Compensation

- Discussion on **limit of object number**
  - Algo requires to pre-define number of clusters  $N_l$ ;
  - It is noted that  $N_l$  is not a particularly important parameter;
  - If  $N_l$  too large, unnecessary clusters end up not having any events allocated to them and thus “die”;
  - $N_l$  can simply be chosen to be large.

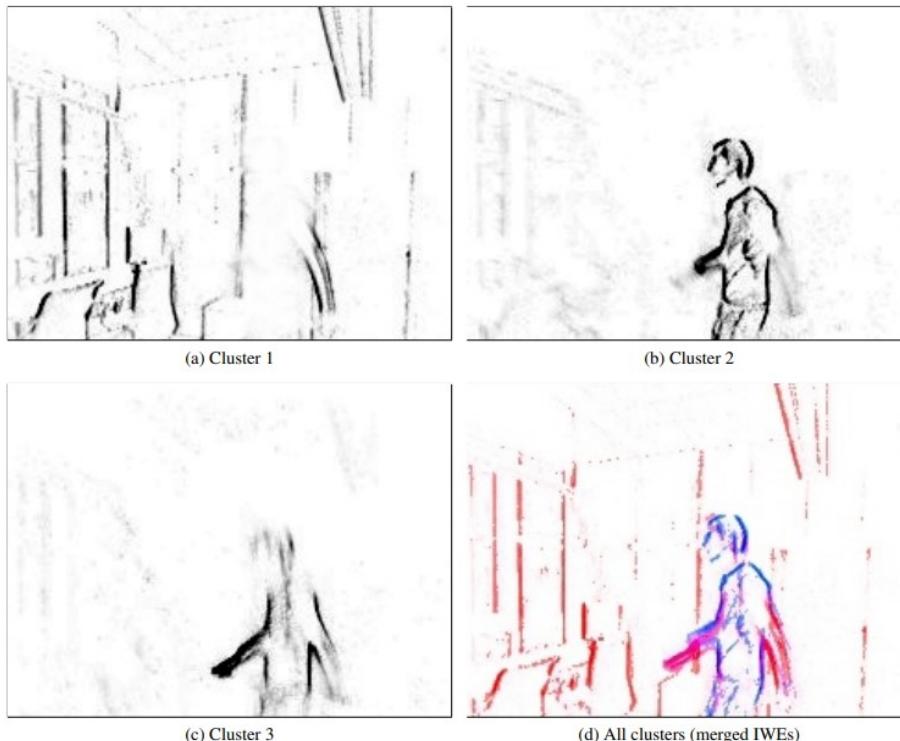


Experiment on *slider depth* sequence. Motion-compensated images in (b) to (f) show events colored by cluster.

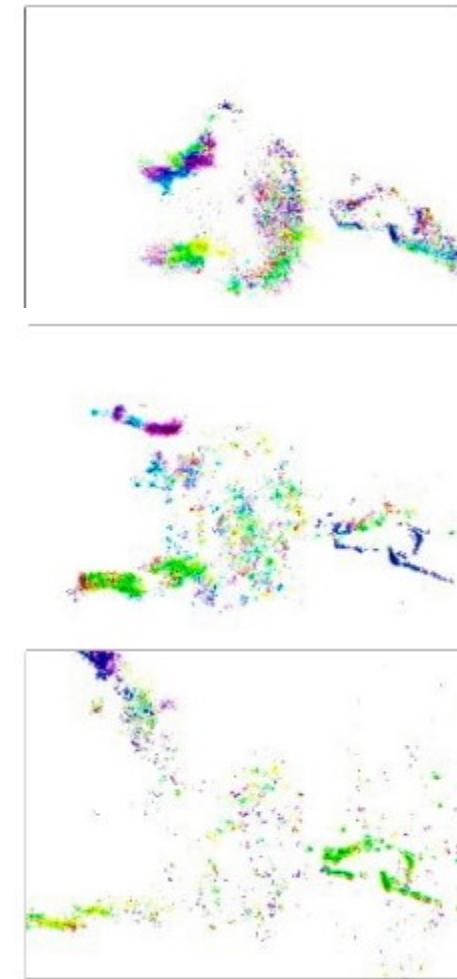
**The real motion for all clusters should be in “Optical Flow” type.**

# Paper: Event-Based Motion Segmentation by Motion Compensation

- Discussion on **non-rigid objects' motion**



A person walks across a room, arms swinging.



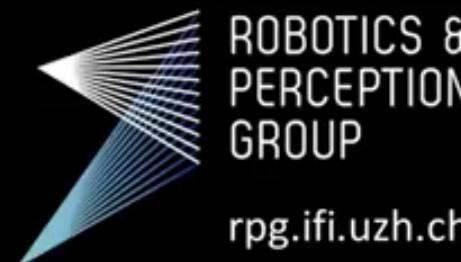
Snapshots of segmentation of balloon popping.  
Num. of clusters = 4,  
Events colored by cluster membership.

# Event-Based Motion Segmentation by Motion Compensation

Timo Stoffregen, Guillermo Gallego, Tom Drummond, Lindsay Kleeman,  
Davide Scaramuzza



Department of Neuroinformatics



## Summary – Motion segmentation

- Image segmentation and motion segmentation
- Problem formulation
- Ideas
  - Events + frame fusion-based
  - Pure event-based
- Selected papers
  - Motion segmentation towards **contrast maximization**

# Supplementary Materials for Self-Learning

# Paper: Simultaneous Estimation for Flow and Intensity

- Method: Sliding window **variational optimisation**
  - Idea on how to design the energy function:

$\vec{u}$  &  $L$ : functions  $\vec{u}(\vec{x}, t), \vec{L}(\vec{x}, t)$

*joint estimation on flow & intensity*

$$\min_{\mathbf{u}, L} \int_{\Omega} \int_T \left( \underbrace{\lambda_1 \|\mathbf{u}_x\|_1 + \lambda_2 \|\mathbf{u}_t\|_1 + \lambda_3 \|L_x\|_1 + \lambda_4 \|\langle L_x, \delta_t \mathbf{u} \rangle + L_t\|_1}_{\text{smoothness of the flow & intensities}} + \lambda_5 h_\theta(L - L(t_p)) \right) dt dx$$

*brightness consistency*

$$+ \int_{\Omega} \sum_{i=2}^{|P(x)|} \left\{ \|L(t_i) - L(t_{i-1}) - \theta \rho_i\|_1 \right\} dx,$$

*all events fired at  $\vec{x}$*

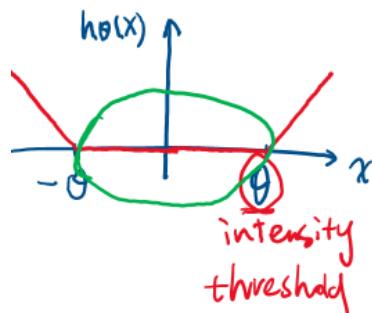
*polarity*

# Paper: Simultaneous Estimation for Flow and Intensity

- The first 3 terms: regularise the smoothness of the flow (both spatially and temporally) and the smoothness of the intensities.
- The 4<sup>th</sup> term: first order Taylor approximation of the brightness consistency equation.
- The last 2 terms: the **no-event data term** and the **event data term**.

No-event data term:

$$\lambda_5 h_\theta(L - L(t_p))$$



$$h_\theta(x) = \begin{cases} |x| - \theta, & \text{if } |x| > \theta \\ 0, & \text{otherwise,} \end{cases}$$

Event data term:

$$\int_{\Omega} \sum_{i=2}^{|P(\mathbf{x})|} \|L(t_i) - L(t_{i-1}) - \theta \rho_i\|_1$$

We want event data to satisfy

$$|L(\mathbf{x}, t) - L(\mathbf{x}, t_p(\mathbf{x}, t))| \geq \theta$$

as much as possible.

## **Conference on Computer Vision and Pattern Recognition (CVPR 2016)**

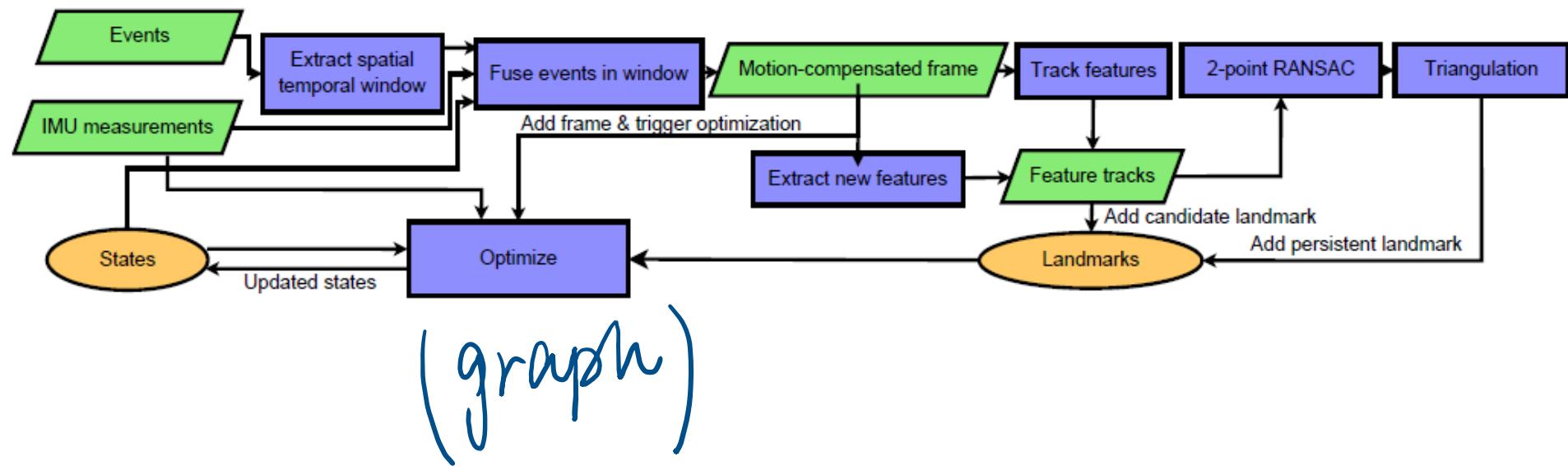
Simultaneous Optical Flow and Intensity Estimation  
from an Event Camera

Patrick Bardow, Andrew Davison and Stefan Leutenegger

Dyson Robotics Laboratory,  
Imperial College London

# Paper: Real-time Visual-Inertial Odometry for Event Cameras using Keyframe-based Nonlinear Optimization

- Overview



# Paper: Real-time Visual-Inertial Odometry for Event Cameras using Keyframe-based Nonlinear Optimization

- Creating Motion-Corrected Event Frames

$$\mathbf{x}'_j = \pi\left(T_{t_k^f, t_j}(Z(\mathbf{x}_j)\pi^{-1}(\mathbf{x}_j))\right)$$

To obtain  $T_{t_k^f, t_j}$  and  $Z(\mathbf{x}_j)$

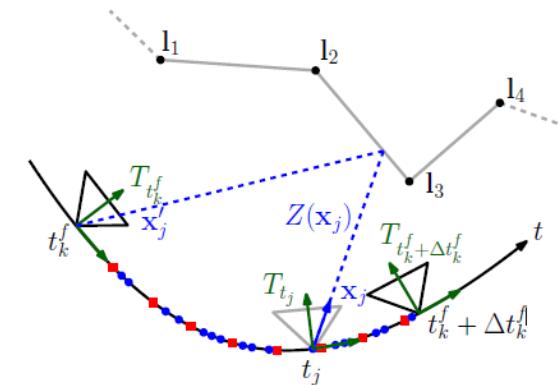
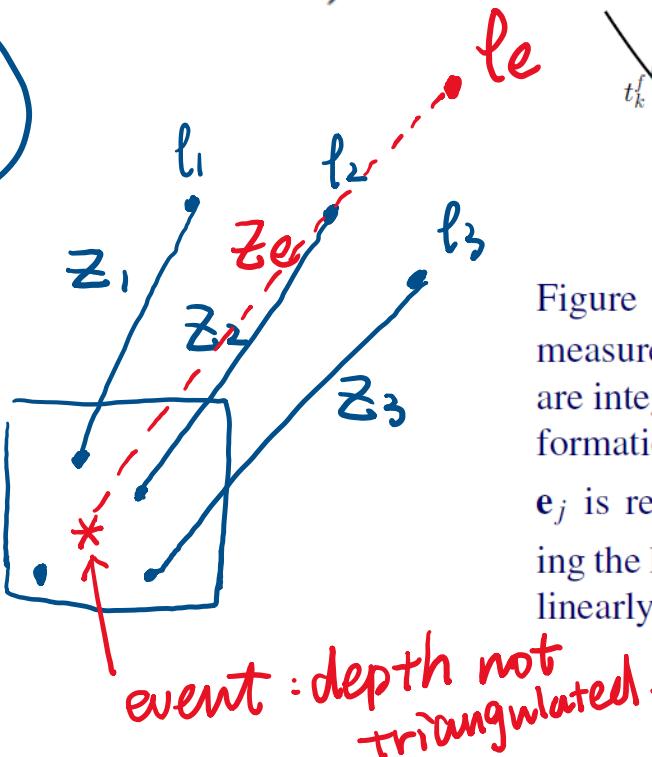
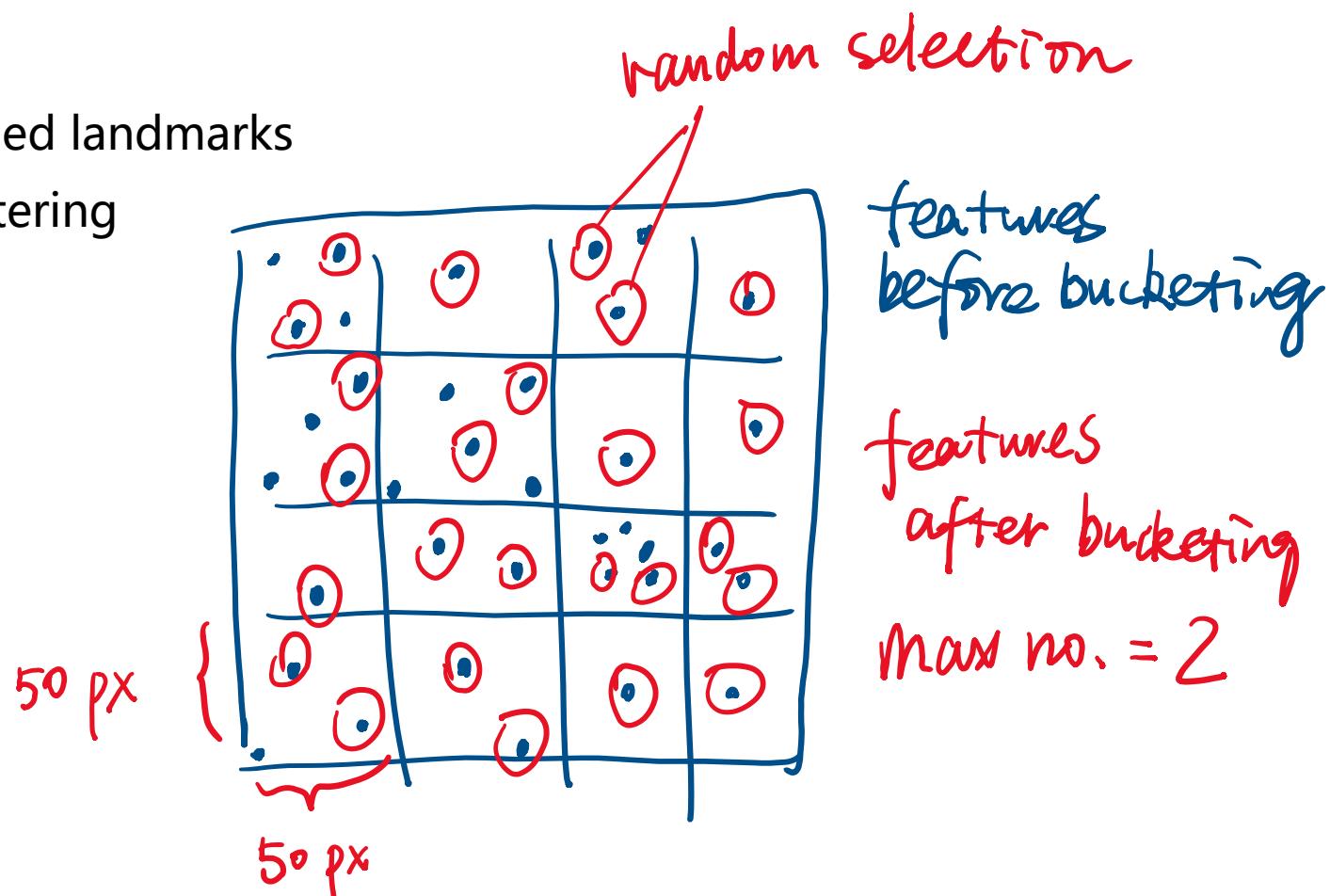


Figure 3: Motion correction: the inertial measurements in  $[t_k^f, t_k^f + \Delta t_k^f]$  (red squares) are integrated to compute the relative transformation  $T_{t_k^f, t_k^f + \Delta t_k^f}$ . Each event (blue dot)  $\mathbf{e}_j$  is reprojected to camera frame  $C_{t_k^f}$  using the linearly interpolated pose  $T_{t_j}$  and the linearly interpolated depth  $Z(\mathbf{x}_j)$ .

# Paper: Real-time Visual-Inertial Odometry for Event Cameras using Keyframe-based Nonlinear Optimization

- Feature Tracking and Landmark Triangulation

- FAST corner detector
- Bucketing
- Dispose failed landmarks
- RANSAC filtering



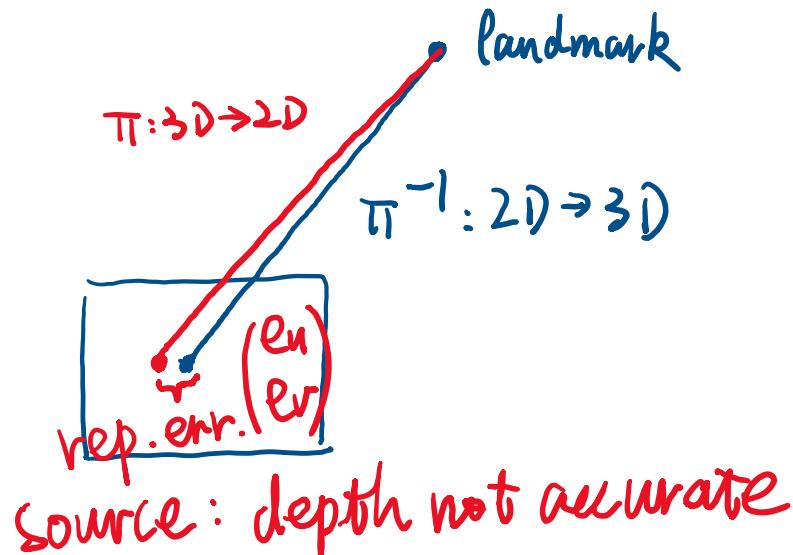
# Paper: Real-time Visual-Inertial Odometry for Event Cameras using Keyframe-based Nonlinear Optimization

- Optimization on Reprojection Error and Inertial Error:

$$J := \sum_{k=1}^K \sum_{j \in \mathcal{J}(k)} \mathbf{e}^{j,k T} \mathbf{W}_r^{j,k} \mathbf{e}^{j,k} + \sum_{k=1}^{K-1} \mathbf{e}_s^{k T} \mathbf{W}_s^k \mathbf{e}_s^k$$

j: landmark index

k: frame index

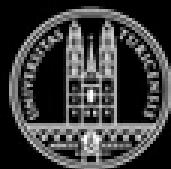


eg. in 2D

$$\begin{aligned} x_1 & \xrightarrow{\text{pred.}} \tilde{x}_2 = \hat{x}_1 + \hat{v}_1 \Delta t + \frac{1}{2} \hat{a}_1 \Delta t^2 \\ & \xrightarrow{\text{est.}} \hat{x}_2 \\ e &:= \frac{\tilde{x}_2 - \hat{x}_2}{\text{meas.}} \end{aligned}$$

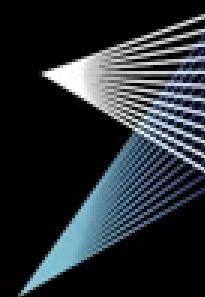
# Real-time Visual-Inertial Odometry for Event Cameras using Keyframe-based Nonlinear Optimization

Henri Rebecq, Timo Horstschaefer, Davide Scaramuzza



University of  
Zurich<sup>UZH</sup>

Department of Informatics

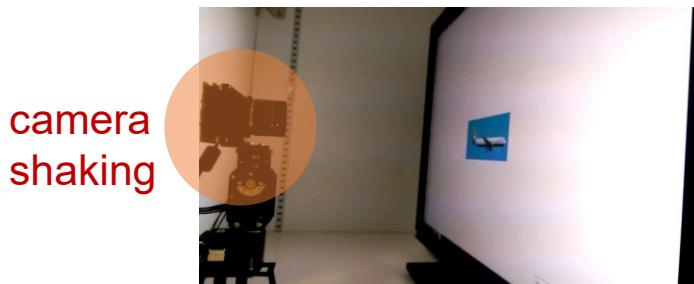


ROBOTICS &  
PERCEPTION  
GROUP

rpg.ifi.uzh.ch

# Recognition

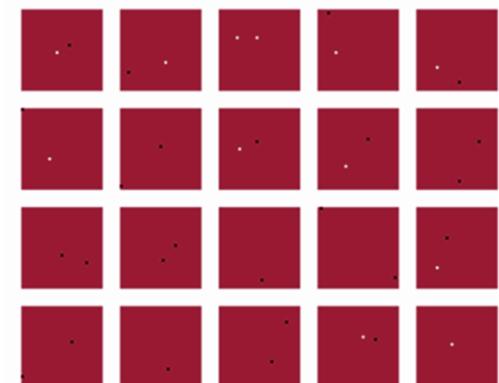
- Tasks (from easy to hard)
    - Detecting the presence of a simple shape (such as a circle)
    - Classification of more complex shapes, such as card pips, block letters and faces
    - More difficult objects:  
similar to conventional Caltech-101 and Caltech-256 datasets
    - Gestures and Motions
- Recording Neuromorphic-Caltech101



## Card symbols recognition



## Event-based MNIST datasets

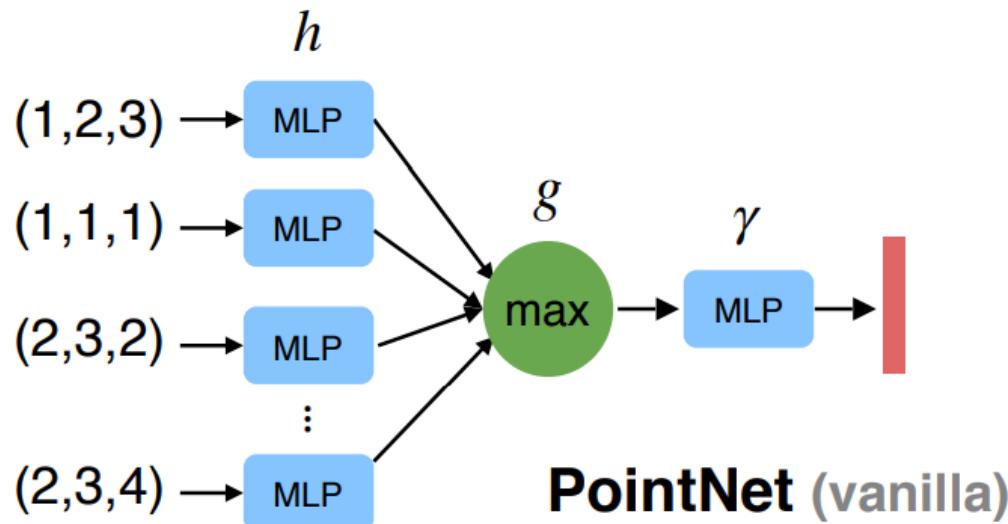


Datasets | Neuromorphic Engineering ([greg-cohen.com](http://greg-cohen.com))  
[www2.imse-cnm.csic.es/cavier/POKERDVS.html](http://www2.imse-cnm.csic.es/cavier/POKERDVS.html)  
Garrick Orchard - N-Caltech101

# Recognition

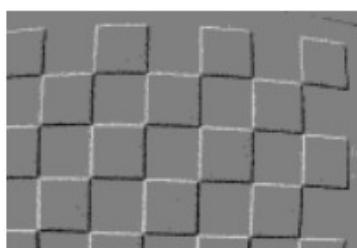
- Ideas
  - **Event representation:** event frames, point clouds, ...
  - **Recognition algorithm:** mostly using neural networks, including Convolutional Neural Network (CNN), Multi-layer Perceptron (MLP), and Spiking Neural Network (SNN)

The Pointnet (which is based on MLP), can be used to process event data.

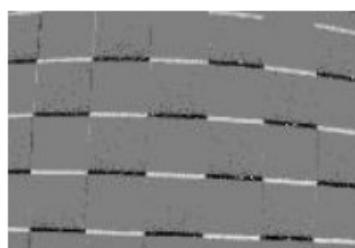


# Recognition

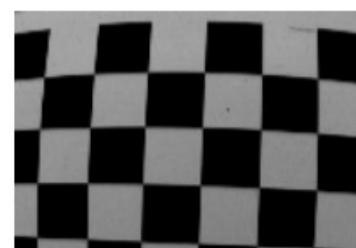
- Advantages
  - Recognizing fast-moving objects / in HDR scene
  - Sparse data: may lead to less computational complexity / power consumption
  - Good at capturing motion within a time interval, instead of a timestamp
- Challenges
  - Still yet to find the most appropriate application scenarios
  - Inadequate dataset, compared to frame-base recognition
  - Requires relative motion to generate data, and data pattern depends on the motion



(a)



(b)

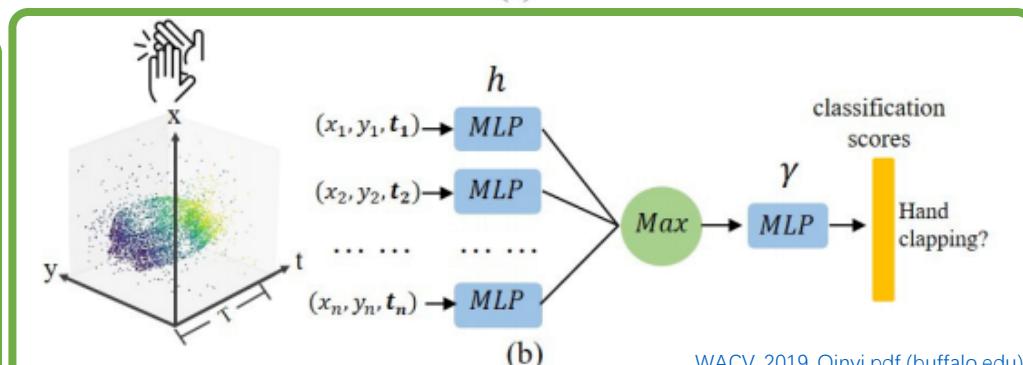
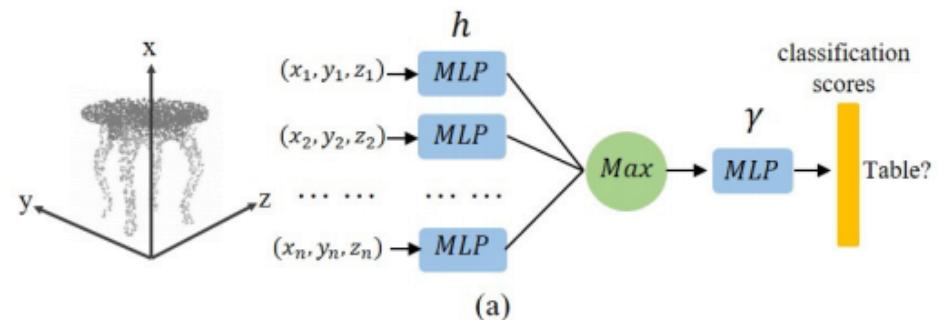
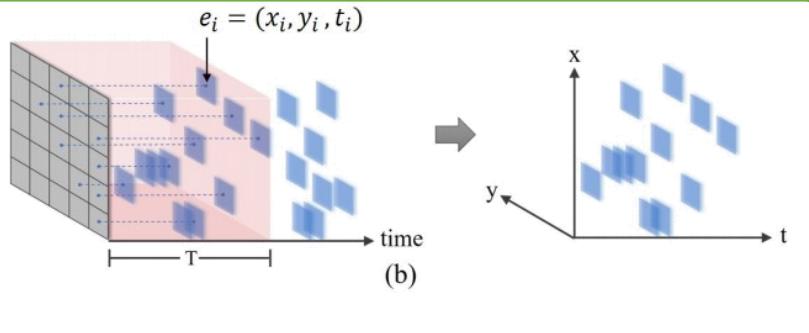
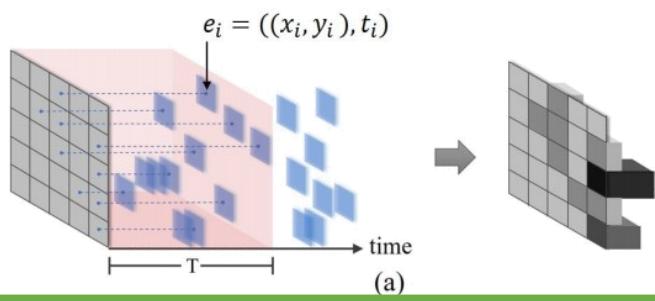


(c)

Different relative motions cause different event patterns, which brings difficulties in recognition!

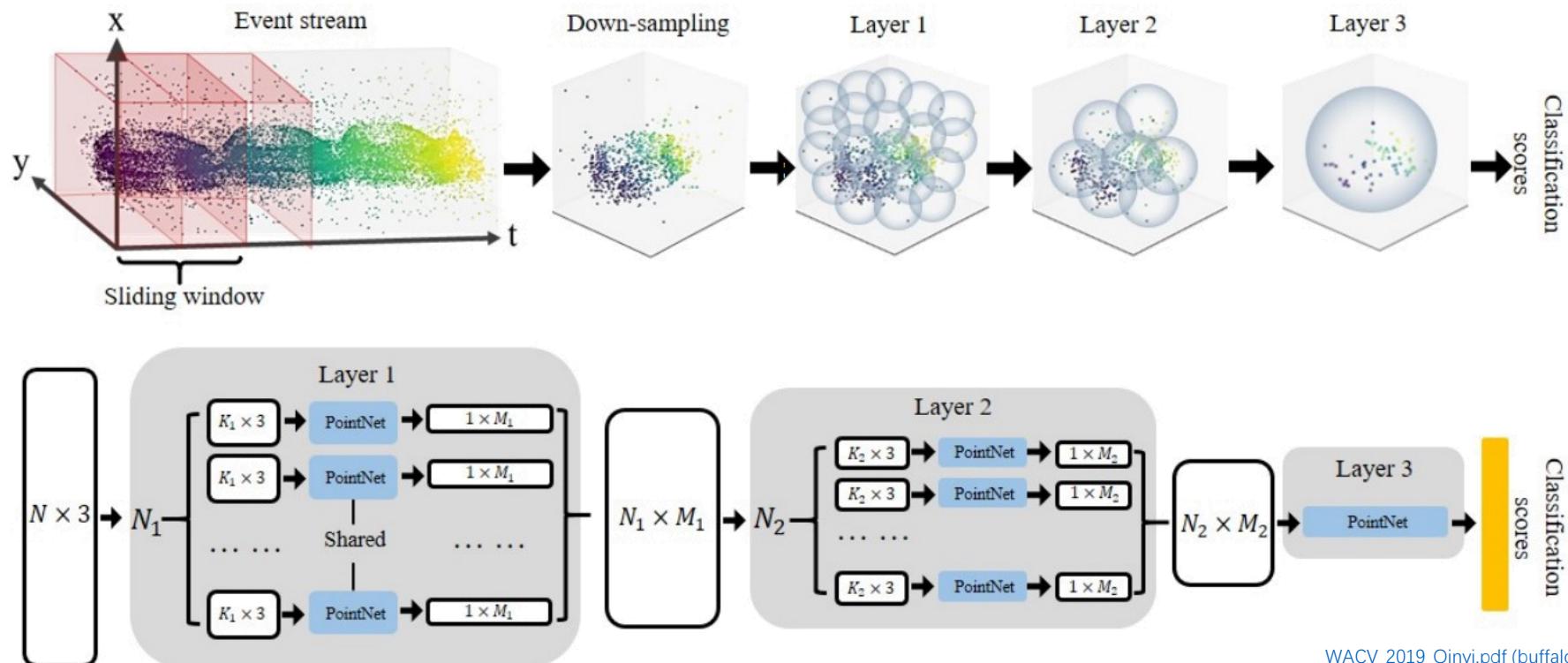
# Paper: Space-time Event Clouds for Gesture Recognition

- Left: the authors treat events as space-time event clouds in this paper.
- Right: the PointNet framework can be adapted to consume space-time event cloud for gesture recognition.



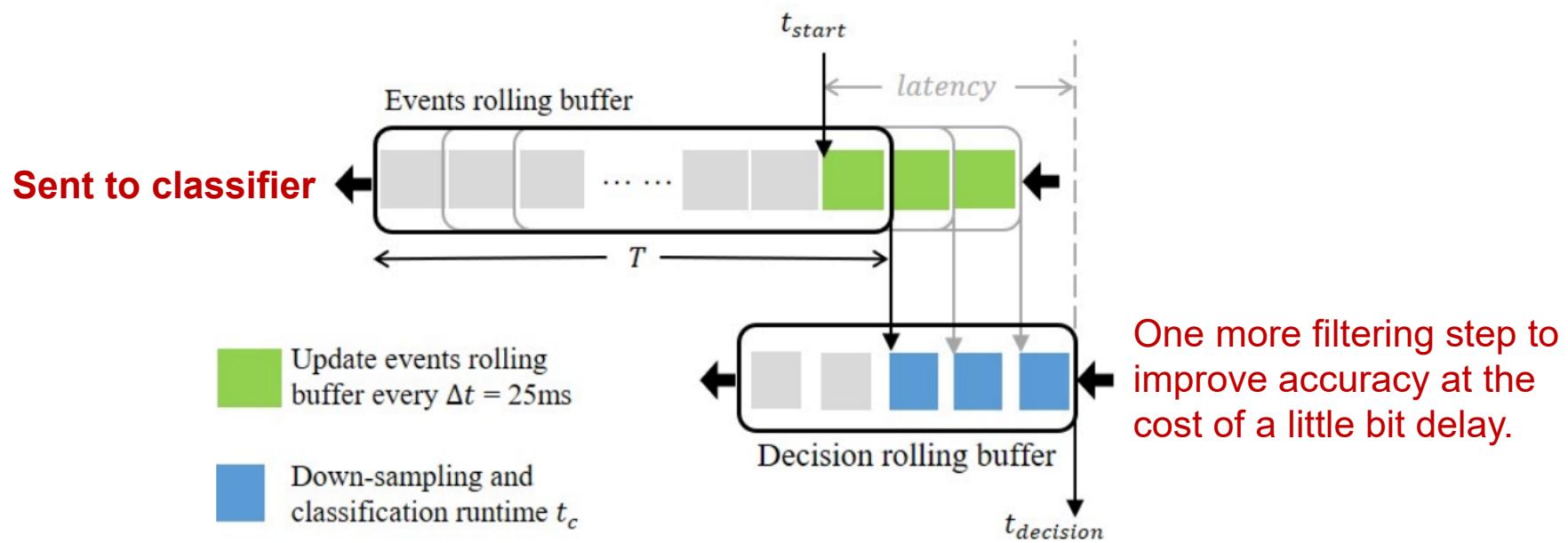
# Paper: Space-time Event Clouds for Gesture Recognition

- The gesture classification framework
  - Trained based on PointNet++, which applies PointNet recursively on a nested partitioning of the input set.
  - Fine-to-coarse, to extract local and global features!



# Paper: Space-time Event Clouds for Gesture Recognition

- The rolling buffer framework for real-time recognition
  - **Events rolling buffer** to store the event point clouds.
  - **Decision rolling buffer** to store the classification results.



# Paper: Space-time Event Clouds for Gesture Recognition

- Results
  - Tested on IBM DVS128 Gesture dataset
  - Results for recognizing 10 classes >95%
  - Small data scale and efficient computation, compared to CNN-based approach!

| Experiment<br>(model,classes,events) | Sliding window |              |           |
|--------------------------------------|----------------|--------------|-----------|
|                                      | $T=0.25s$      | $T=0.50s$    | $T=1.00s$ |
| PointNet,10,256                      | 87.85          | 89.63        | 88.54     |
| PointNet,10,512                      | 88.67          | <b>90.20</b> | 89.61     |
| PointNet,10,1024                     | 88.77          | 89.68        | 89.92     |
| PointNet++,10,256                    | 95.28          | 95.59        | 95.54     |
| PointNet++,10,512                    | 95.39          | <b>96.34</b> | 95.61     |
| PointNet++,10,1024                   | 94.93          | 95.89        | 95.97     |
| PointNet++,11,256                    | 91.92          | 93.38        | 93.61     |
| PointNet++,11,512                    | 92.23          | <b>94.10</b> | 93.83     |
| PointNet++,11,1024                   | 91.87          | 91.91        | 92.63     |

Why  $T=0.50s$  results are the best?

Why 1024-points results worse than 512-point results?