

Introduction to Visual Intelligence

Kuk-Jin Yoon

Visual Intelligence Lab.
Department of Mechanical Engineering, KAIST

What is *Artificial Intelligence*?

The intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans and other animals

from Wikipedia

*The theory and development of (computer) systems able to perform tasks normally requiring human intelligence, such as **visual perception**, **speech recognition**, **decision-making**, and **translation between languages***

from the Oxford dictionary

AI Fields for Different Tasks



Machine Learning
(Learn)



Computer Vision
(See)



Speech Recognition
(Hear)



Natural Language Processing (NLP)
(Communicate)



Expert Systems
(Think)



Motion Planning
(Move)

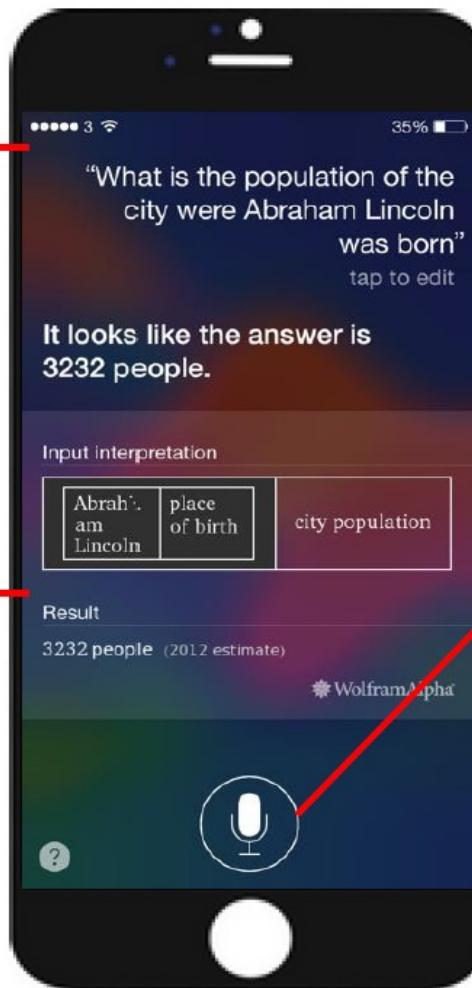
AI listens, thinks and communicates...



Natural language processing (NLP) focuses on human-computer interaction, enabling computers to derive meaning from human language input; and also generate natural language responses. Today, machines proficiently understand natural language syntax but face great challenge in interpreting sentiment (i.e. sarcasm, excitement).



Expert Systems emulate human expert decision-making abilities. It allows the computer to solve for complex problems by reasoning about knowledge, navigating if-then rules.



From the creators of Siri, Viv enables developers to create anything on top of its, conversational interface, making 'her' smarter.



Speech Recognition is the process of mapping audio speech data to textual sentences or key phrases. As humans can speak 150 words per minute on average, but can only type 40, speech recognition has great potential in computer efficiency.

As more voice usage data becomes available, speech recognition accuracy will get better and better. In 2010, accuracy for technology companies hovered around 70, and today sits between 95 and 99.

Sees, moves and learns...



Self-Driving Car Project



(Think)

Machine learning is training computers with datasets to recognise patterns, develop algorithms and self-improve. Machine Learning has been central to today's unprecedented momentum in AI, as it enables the progress of other AI branches.



(See)

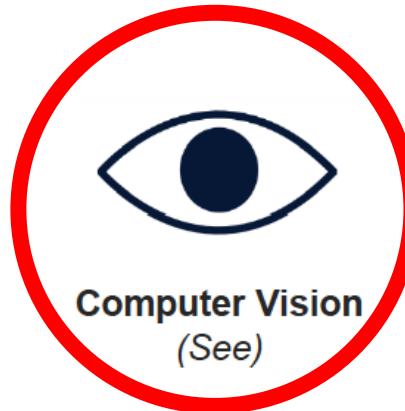
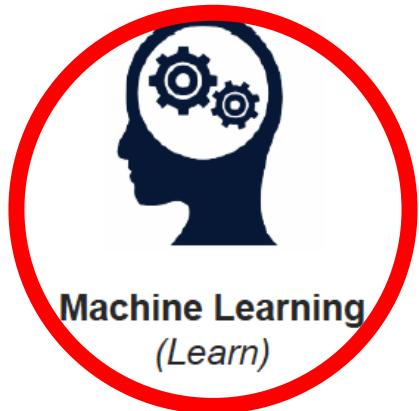
Computer vision is the ability to electronically perceive and understand image/video sources, extract meaningful information and take action. Up until now, image recognition has been driven by rules-based categorisation. Today, machines are fed data so they build their own vision.



(Move)

Motion Planning is the process of forming a strategy of action sequences to achieve a desired movement, typically for execution by intelligent agents, autonomous robots and unmanned vehicles. Today, we are at advanced levels of simple motion planning problems, such as 'move from A to B, while avoiding collision with any obstacles.'

Visual Intelligence in AI



Speech Recognition
(Hear)



Natural Language Processing (NLP)
(Communicate)



Expert Systems
(Think)



Motion Planning
(Move)

① What is Machine Learning?

A screenshot of a Google search results page for the query "machine learning". The top navigation bar shows the search term and various filter options like "이미지" (Images). Below the search bar, there are several image thumbnails arranged in a grid. These images include: 1) A network diagram with nodes labeled "Machine Learning". 2) A stylized human head with a circuit board pattern. 3) A brain filled with icons representing various technologies. 4) A graphic with the words "MACHINE LEARNING" and a neural network icon. 5) A diagram titled "A Standard Machine Learning Pipeline" showing a flow from data input to output. 6) A brain with a circuit board texture. 7) A graphic of gears labeled "Machine Learning" and "Practical Application". Other visible image thumbnails include comparisons between Deep Learning and Machine Learning, various brain-related infographics, and diagrams illustrating machine learning processes like clustering, classification, and regression.

Machine learning

From Wikipedia, the free encyclopedia

For the journal, see [Machine Learning \(journal\)](#).

"Statistical learning" redirects here. For statistical learning in linguistics, see statistical learning in language acquisition.

Machine learning (ML) is the **scientific study of algorithms and statistical models that computer systems use** to effectively perform a specific task without using explicit instructions, relying on patterns and inference instead. It is seen as a subset of **artificial intelligence**. Machine learning algorithms build a mathematical model of sample data, known explicitly programmed to perform **email filtering**, detection of new algorithm of specific instruction **statistics**, which focuses on methods that deliver results. Machine learning studies the question of how machines can learn from experience.

Machine learning and data mining

In its application across business problems, machine learning is also referred to as [predictive analytics](#).

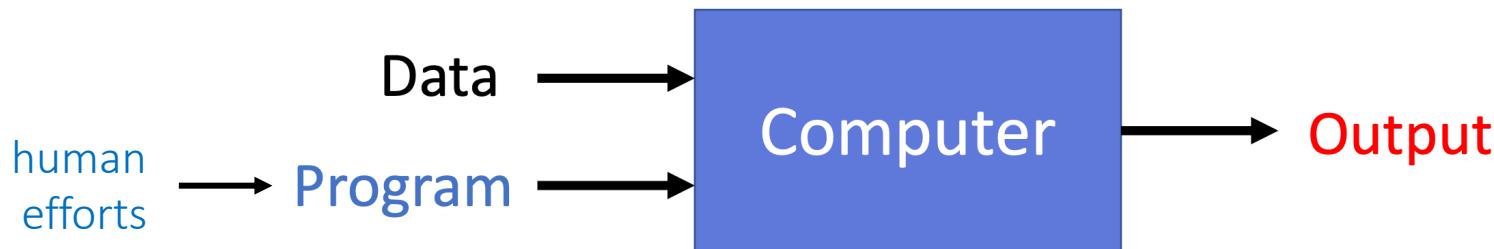
Contents [\[hide\]](#)

- 1 Overview of Machine Learning
 - 1.1 Machine learning tasks
 - 2 History and relationships to other fields
 - 2.1 Relation to data mining
 - 2.2 Relation to optimization

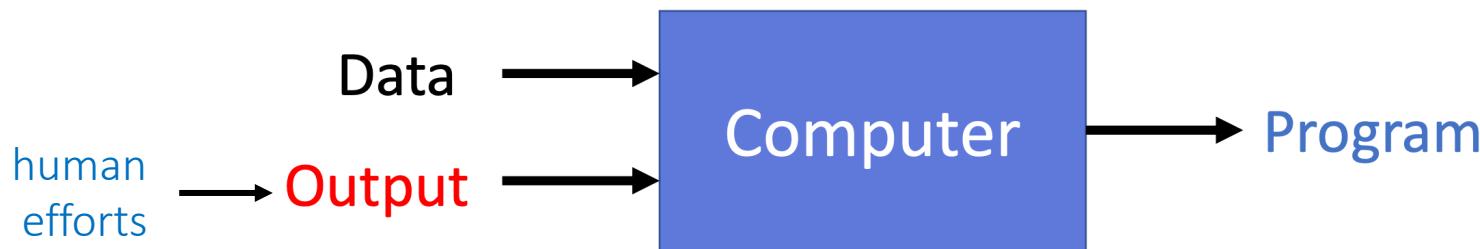
- | | |
|-----------------------------------|--------|
| Clustering | [show] |
| Dimensionality reduction | [show] |
| Structured prediction | [show] |
| Anomaly detection | [show] |
| Artificial neural networks | [show] |
| Reinforcement learning | [show] |
| Theory | [show] |
| Machine learning venues | [show] |

What Machine Learning is doing

Traditional Programming



Machine Learning



② What is Computer Vision?

Computer vision - Google 검색

https://www.google.com/search?q=Computer+vision&newwindow=1&source=lnms&tbo=isch&sa=X&ved=0ahUKEwjJh-WGffffAhUNzt4KHSkzClCQ_AUIDigB&biw=1280

Bookmarks Google 블라우저에서 가... KAIST 강의 준비 Signals and Syste... InCites KAIST 웹메일 시스템 기타 북마크

Google Computer vision

전체 이미지 동영상 뉴스 도서 더보기 설정 도구

컬렉션 세이프서치▼

동향 알고리즘 딥 러닝 시장 프로젝트 기술 ppt 인공 지능

8 cool new ways computer vision is cha... venturebeat.com

How Computer Vision Can Change the Automoti... medium.com

Computer Vision Lab | University of canterbury.ac.nz

Expand the Possibilities of Computer Vision software.intel.com

How To Build A Career in Computer Vision analyticsindiamag.com

Computer Vision Publications & Researcher... nvidia.com

What is Computer Vision? | Hayo hayo.io

Understanding Images: Computer Vision i... weareworldquant.com

All you need is attention -- Computer Visio... medium.com

Computer Vision for Quality Assurance | Catalys... catalysts.cc

Researchers improve patient safety with bedsid... engineering.stanford.edu

computer-vision-machine-learning - Alexandru ... alexvoica.com

Deep Dive Into Computer Vision - YouTube youtube.com

Learn Computer Vision - robotic.org

5 ways computer vision could impa... venturebeat.com

Computer Vision: Crash Course Computer Sc... youtube.com

Great Moments in Computer Vision thevisionary.com

What is Computer Vision - Post 5: A Vers... blippar.com

Computer Vision, Robotics and Machine Lear... surrey.ac.uk

Computer Vision jobs: Why you should consi... eu-recruit.com

computer-vision-techniques/README.md at ... github.com

MIT's computer vision could make robo... dnaIndia.com

Past, Present and Future of Computer Vi... aiforsecurity.com

Security + AI: 5 Best Computer V... aiforsecurity.com

Computer vision startup Tunicorn Tech... alltechasia.com

What is Computer Vision? | Hayo hayo.io

Computer Vision Takes Off - Dietrich College ... cmu.edu

Introduction to the Artificial Intelligence and ... slideshare.net

Creating Applications with the Intel® ... insidehpc.com

7 Steps to Understanding Computer Vision ... kdnuggets.com

DojoCV: Computer Vision With Python and Flask... codingdojo.com

The Most Exciting Applications of Compute... indatalabs.com

Comparing the Top Five Computer Vision ... goberol.com

Computer Vision at Princeton vision.princeton.edu

Computer Vision | Mitsubishi Electric Research... merrl.com

A Year in Computer Vision themranki.org

Computer Vision Vs. Image Processing: What ... fedtechmagazine.com

Computer vision and the AI boom ~ Hua... huawei.com

Computer Vision using Deep Learnin... trainings.analyticsvidhya.com

Detecting human facial expression by commo... interactivearchitecture.org

KAIST



WIKIPEDIA
The Free Encyclopedia

Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia
Wikipedia store

Interaction

Help
About Wikipedia
Community portal
Recent changes
Contact page

Tools

What links here
Related changes
Upload file

Article Talk

Read

Edit

View history

Search Wikipedia

Not logged in Talk Contributions Create account Log in

Computer vision

From Wikipedia, the free encyclopedia

Computer vision is an interdisciplinary scientific field that deals with how computers can be made to gain high-level understanding from digital images or videos. From the perspective of engineering, it seeks to automate tasks that the human visual system can do.^{[1][2]}

Computer vision tasks involve processing high-dimensional data to make decisions.^{[4][5][6][7]} Understanding descriptions of the world that can interface with other thought processes and elicit appropriate action. This image understanding can be seen as the disentangling of symbolic information from image data using models constructed with the aid of geometry, physics, statistics, and learning theory.^[8]

As a scientific discipline, computer vision is concerned with the theory behind artificial systems that extract information from images. The theory is drawn from a variety of disciplines, including mathematics, physics, computer science, neuroscience, and psychology.

Computer vision is concerned with the theory behind artificial systems that extract information from images.

Sub-domains of computer vision include scene reconstruction, event detection, video tracking, object recognition, 3D pose estimation, learning, indexing, motion estimation, and image restoration.^[6]

Connections to Other Disciplines for AI

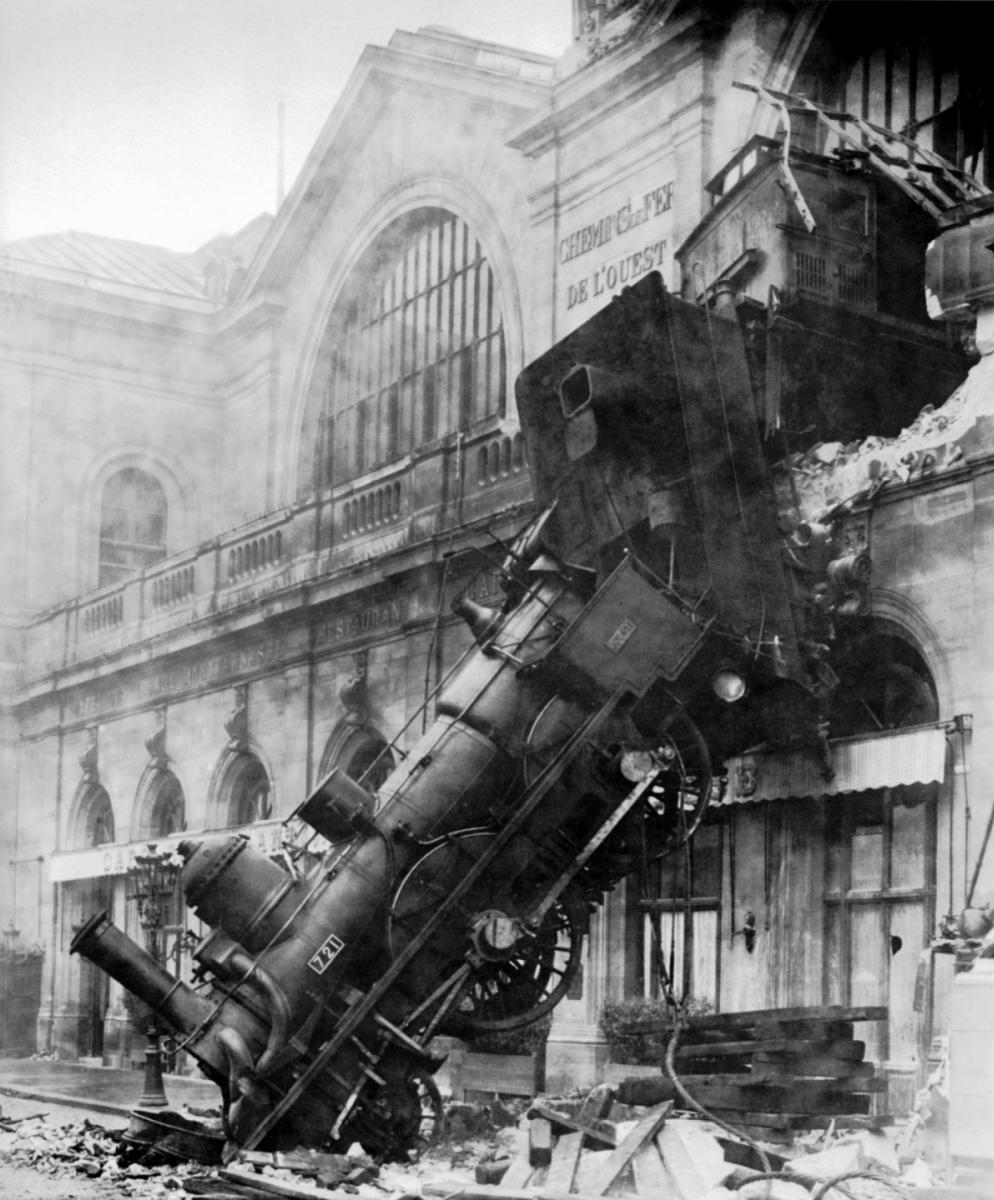


In Short, Visual Intelligence

- Takes visual data as an input
 - Visual data:
 - Image or video (video can be interpreted as a temporal sequence of images)
- and outputs meaningful semantic information such as
 - low-level (primitive) information
 - 3D structure, motion, etc.
 - high-level information
 - object identity (recognition), semantic segmentation, etc.

Image capturing





What does an image tell us?

We see a picture (or pictures) and perceive semantic information such as

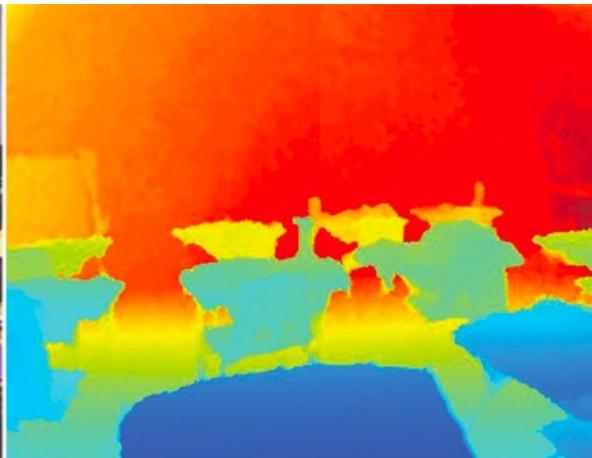
- color
- motion
- depth (or 3D)
- pattern
- object category
- context
- ...

We can even tell a long story from a single image.

Granville-Paris Express wreck on 22 October 1895
(image source: wikipedia
https://en.wikipedia.org/wiki/Montparnasse_derailment)

Visual Intelligence - What We Perceive from Images

- 3D perception
- Semantic object-level perception
- Color perception
- Motion perception
- Etc.



Visual Intelligence – 3D Perception

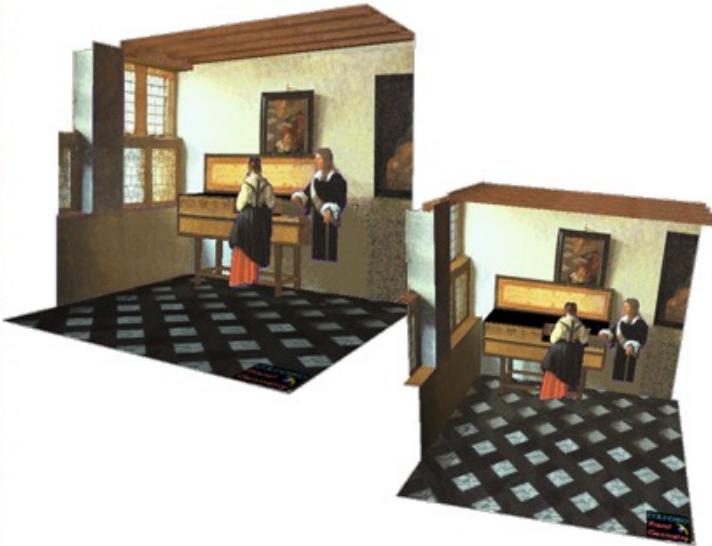
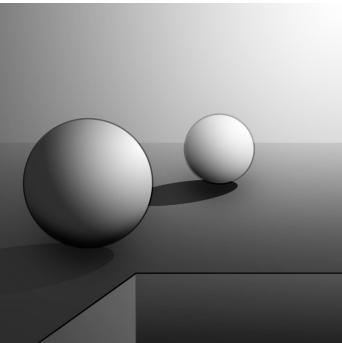


from a pattern or texture



from an object identity or shape

Visual Intelligence – 3D Perception



from multiple visual cues (shading, perspective)



(a) photo

(b) depth

Visual Intelligence – Semantic Object-level Perception

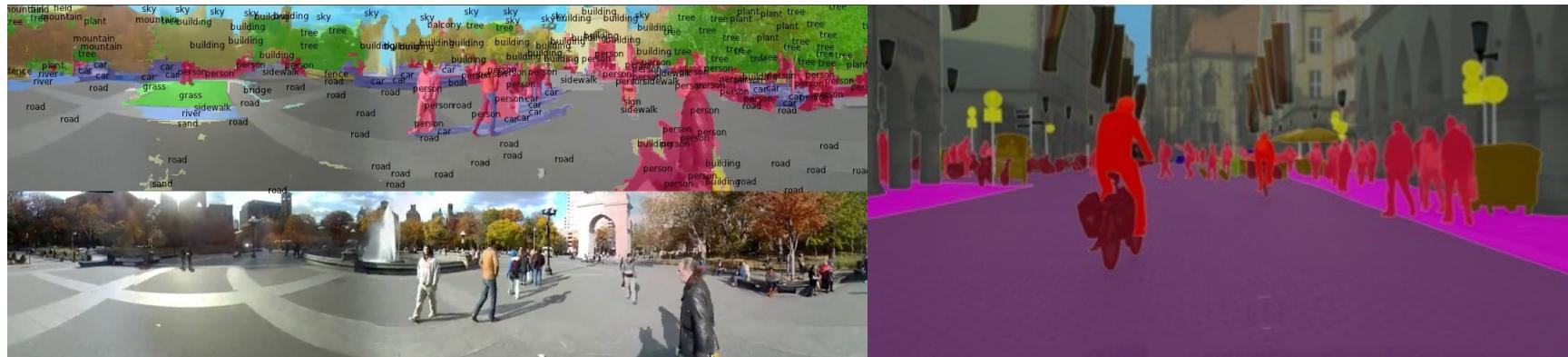


(a) photo

(c) instance

(d) class segmentation

Visual Intelligence – Scene Parsing/Semantic Segmentation

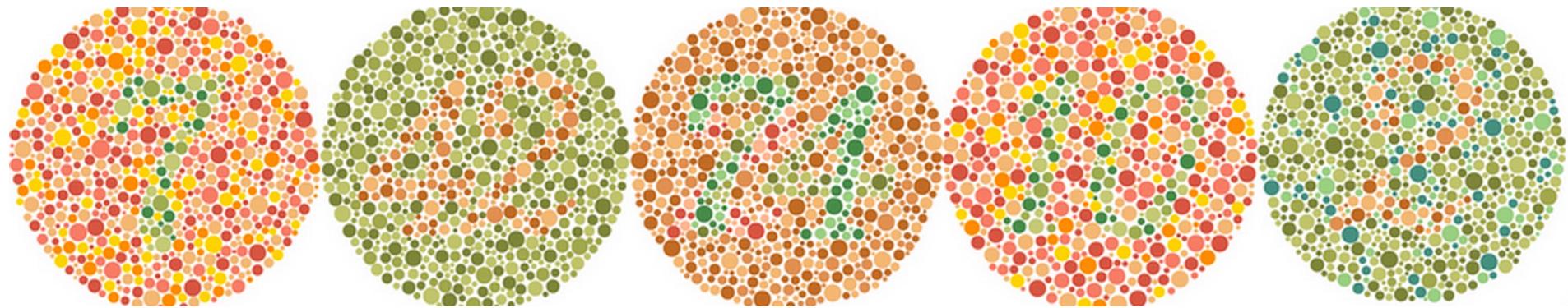


30.1 FPS

Visual Intelligence – Color Perception



Visual Intelligence – Color Perception



Visual Intelligence – Motion Perception

an open-source short film
by  **blender**™ Institute



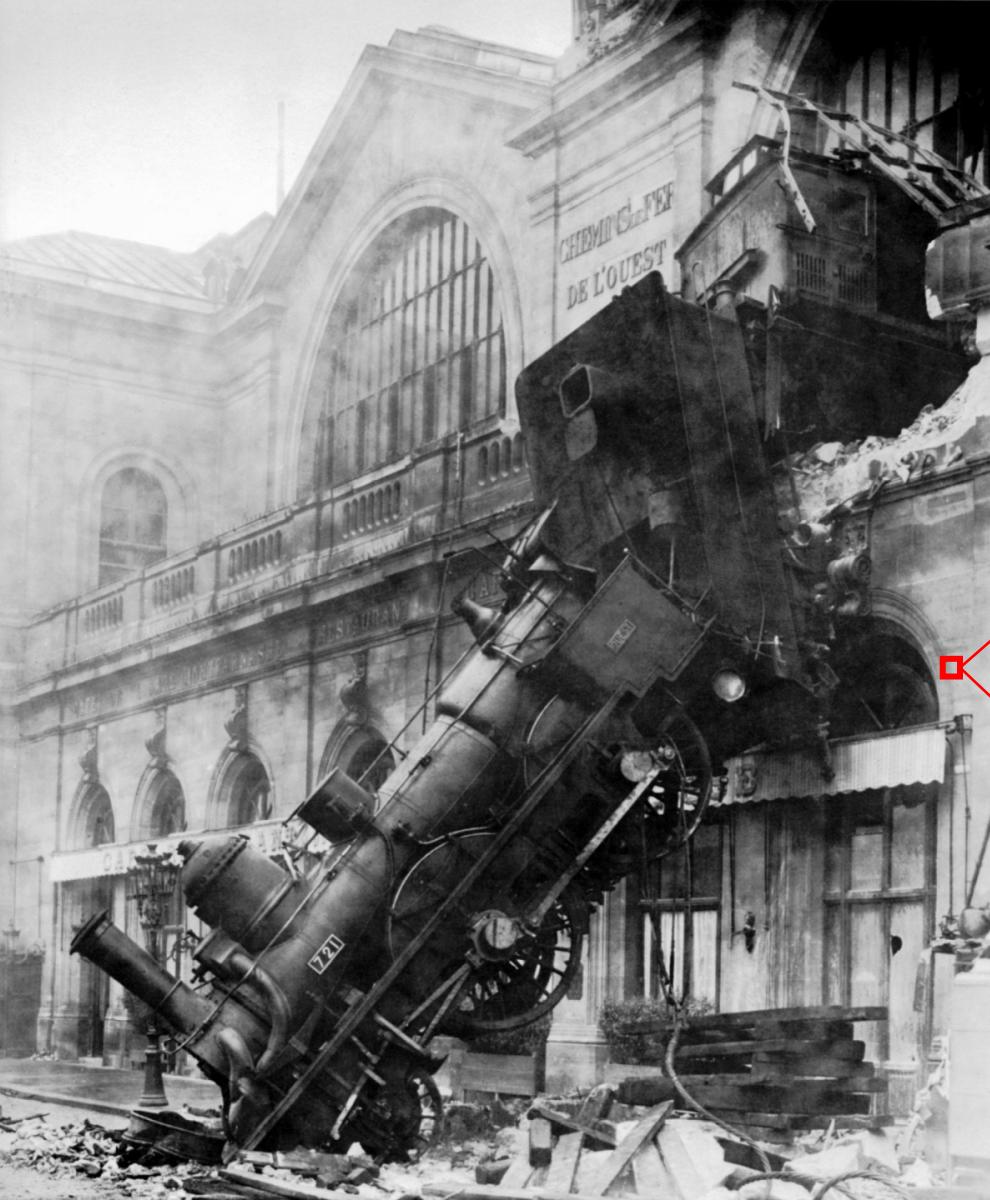
Butler et al, ECCV 2012

MPI Sintel Flow Dataset

A data set for the evaluation of optical flow derived from the open source 3D animated short film, Sintel.

Visual Intelligence – Wholistic Scene Understanding





What does an image tell **computers**?

Computers see (or get) numbers only.

| | | | | | | | |
|-----|-----|-----|-----|-----|-----|-----|-----|
| 183 | 160 | 94 | 153 | 194 | 163 | 132 | 165 |
| 183 | 153 | 116 | 176 | 187 | 166 | 130 | 169 |
| 179 | 168 | 171 | 182 | 179 | 170 | 131 | 167 |
| 177 | 177 | 179 | 177 | 179 | 165 | 131 | 167 |
| 178 | 178 | 179 | 176 | 182 | 164 | 130 | 171 |
| 179 | 180 | 180 | 179 | 183 | 169 | 132 | 169 |
| 179 | 179 | 180 | 182 | 183 | 170 | 129 | 173 |
| 180 | 179 | 181 | 179 | 181 | 170 | 130 | 169 |

Visual Intelligence is to bridge the gap between pixels and “meaning”.

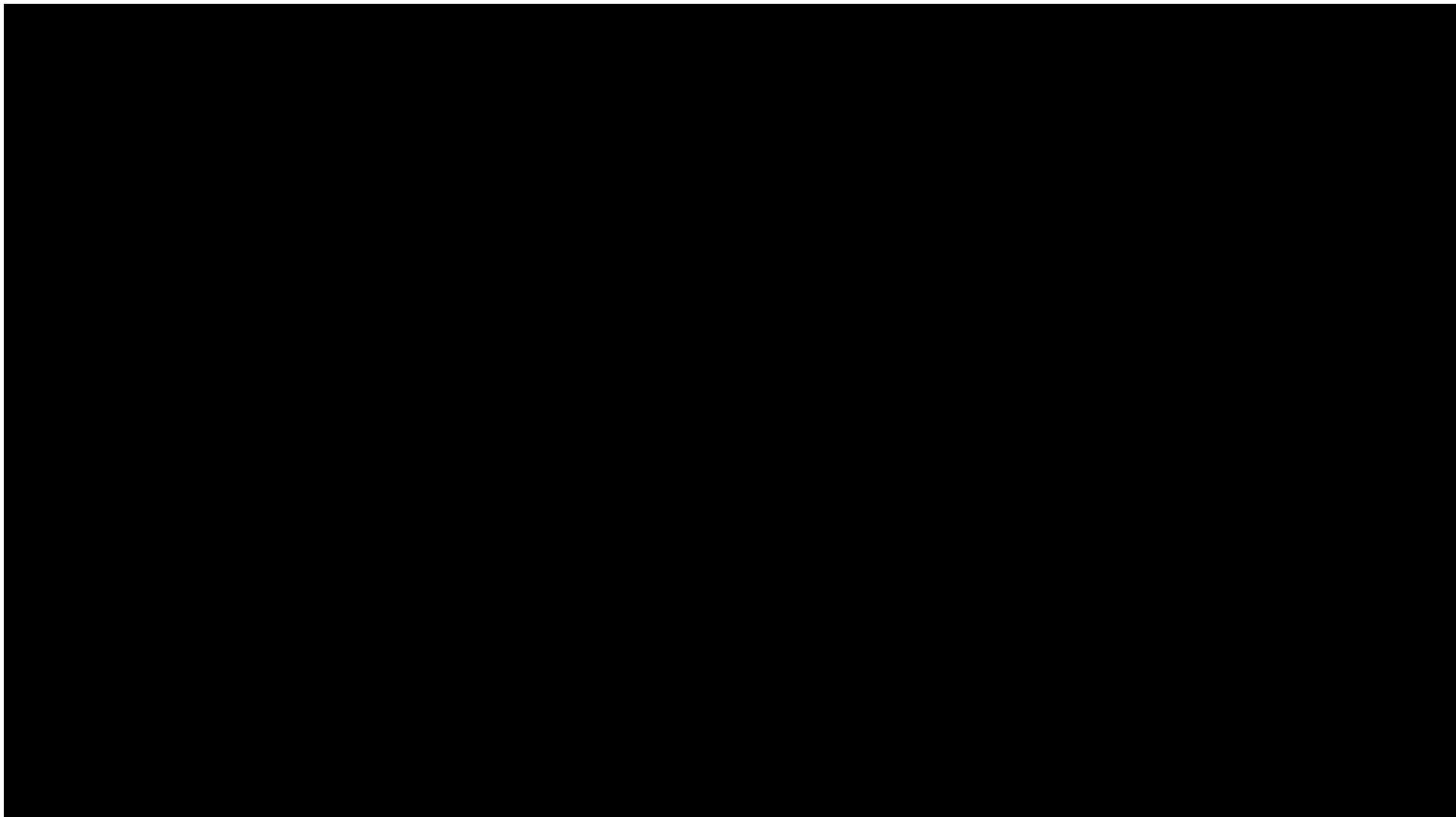
Granville-Paris Express wreck on 22 October 1895

(image source: wikipedia

(https://en.wikipedia.org/wiki/Montparnasse_derailment)

Ultimate Goal of Visual Intelligence Research

- Making computers perceive the images holistically as humans do



SOTA Papers on more than 1000 tasks

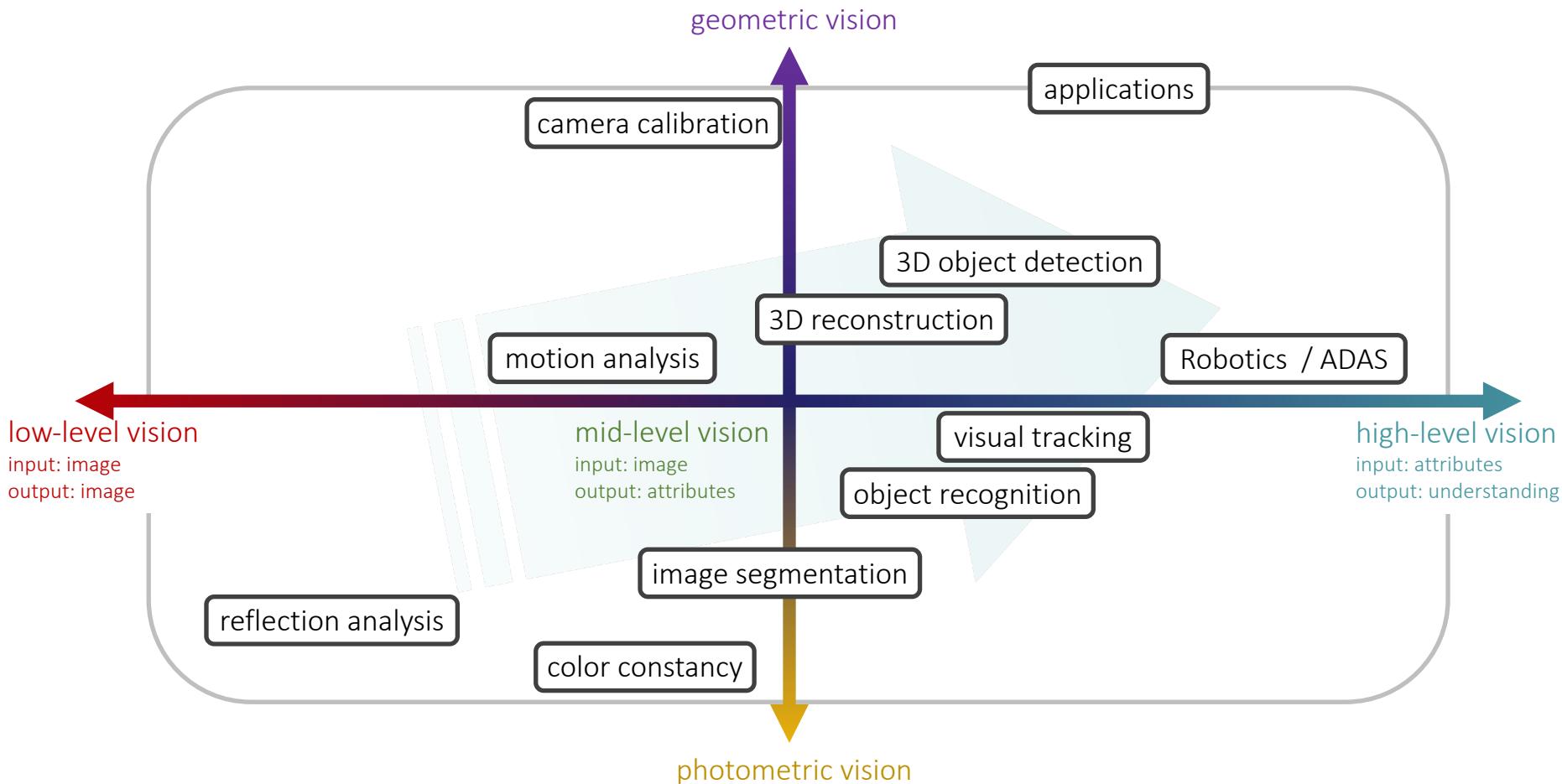
[Browse SoTA](#) > Computer Vision

Computer Vision

2960 benchmarks • 1002 tasks • 1974 datasets • 26695 papers with code

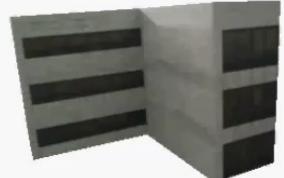
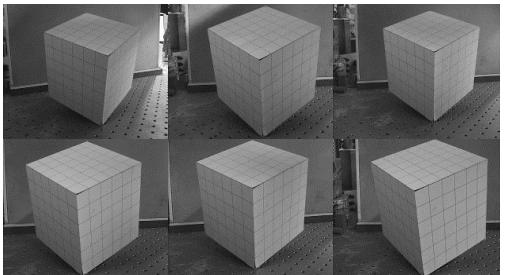
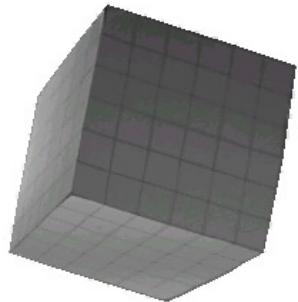
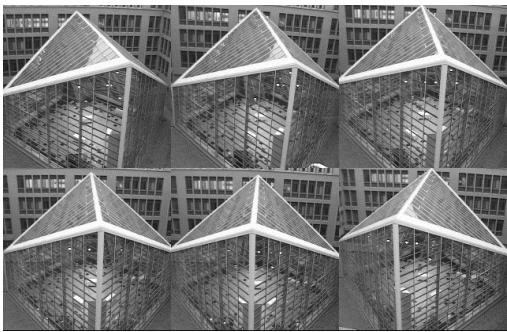
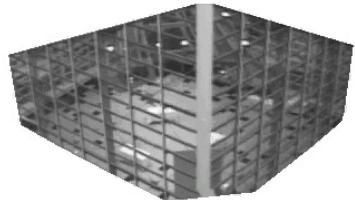
More topics and downloadable source codes are available at the website “Papers With Code: the latest in machine learning”(<https://paperswithcode.com/area/computer-vision>)

Classification of Research Topics



3D Reconstruction

Camera Self-calibration



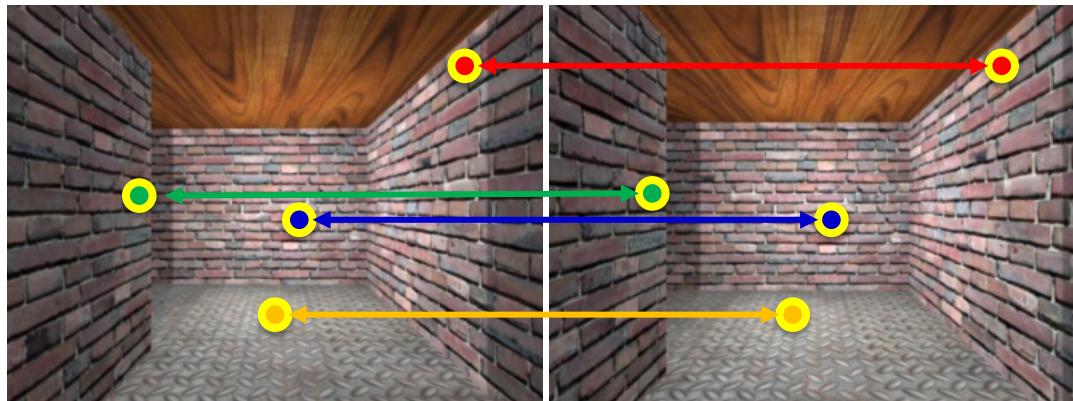
camera self-calibration (ICRA 2000)

3D Sensing – Stereo Vision

- Disparity
 - Informally, the difference between two pictures and allows us to gain a sense of depth
- Stereopsis
 - Ability to perceive depth from disparity
 - Gaining the sense of depth by fusing the images recorded by two or more cameras and exploiting the difference
- Stereo vision
 - Design algorithms that mimic stereopsis
 - can use two-view (binocular) or multi-view



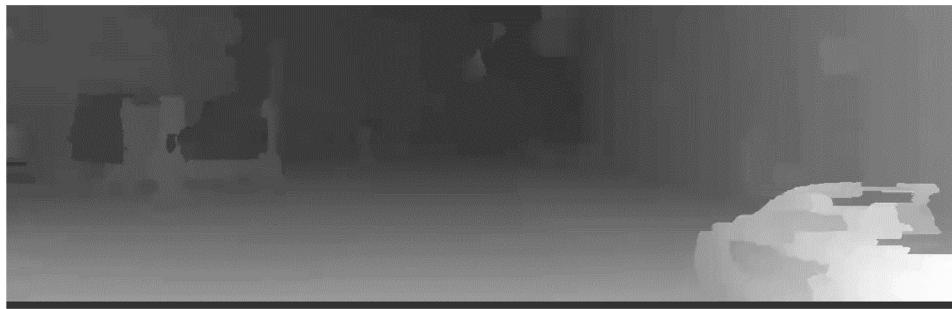
Stereo matching finds correspondences to obtain disparity, which is the difference in coordinates of correspondences.



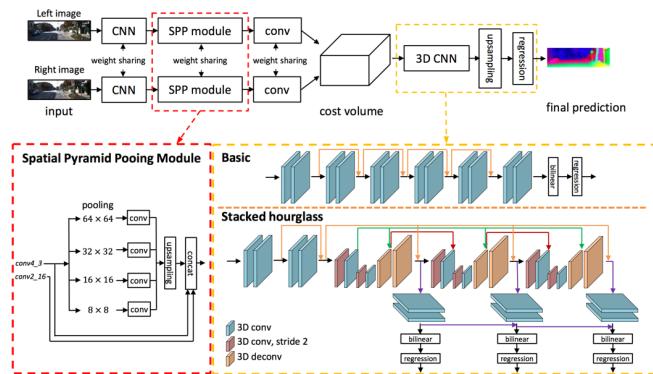
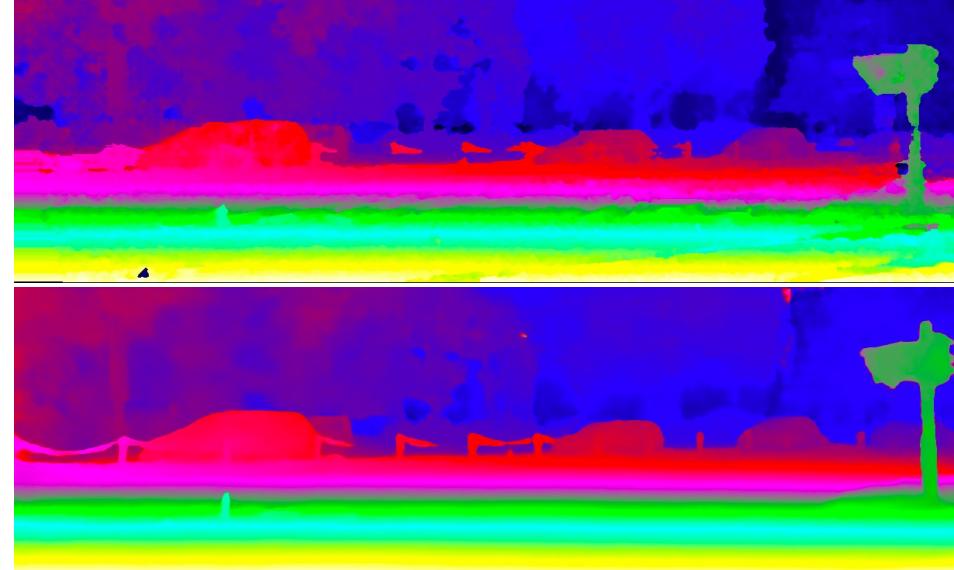
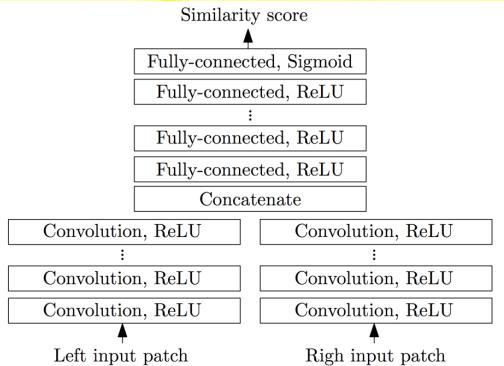
target image

disparity map

Spatiotemporal Stereo Matching with 3D Disparity Profiles



Deep-learning-based Stereo



Results using the KITTI Vision Benchmark Dataset (<http://www.cvlibs.net/datasets/kitti/>)

Stereo Matching in Real Applications

- Which direction should we go?
 - Considering all the difficulties that are possible to occur



Snowy night



Sun flare



Snow



Rain

Image Dehazing combined with Stereo Estimation

- Haze is a slight obscuration of the lower atmosphere, typically caused by fine suspended particles
- It is an annoying factor since it causes poor visibility and degradation of vision algorithms



Vehicle detection Ground Truth



Vehicle detection result of DPM

- **Dehazing** is a process of image enhancement by eliminating haze components in image

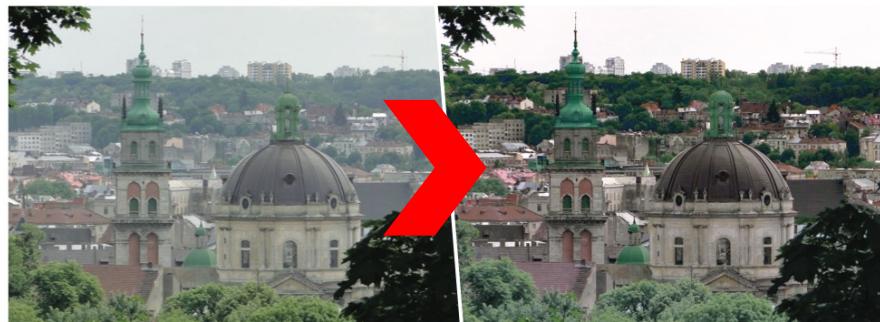
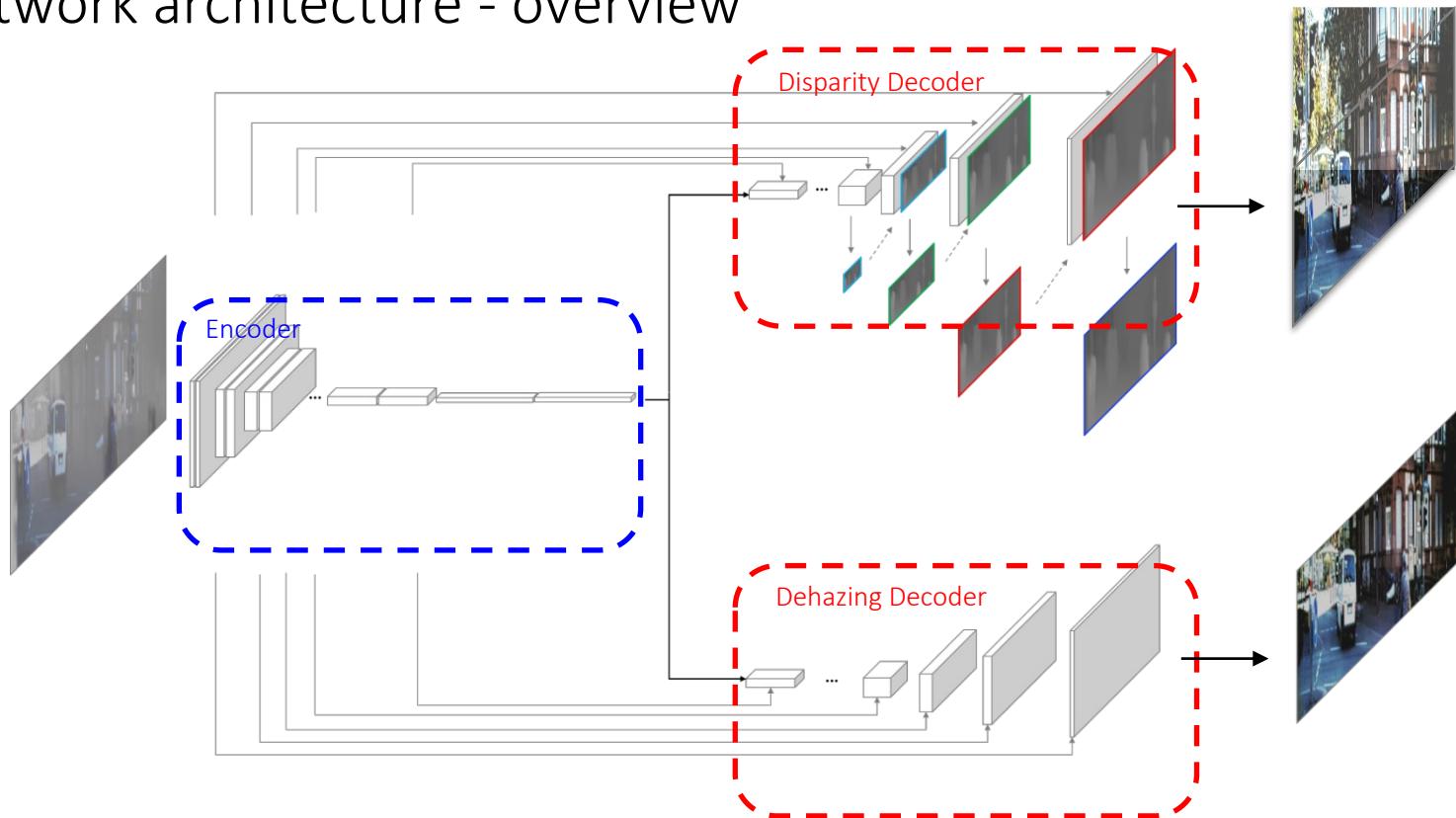


Image Dehazing combined with Stereo Estimation

[J.-Y. Na et al., CoView 2018]

- Network architecture - overview



Proposed Neural Network for Dehazing

Image Dehazing combined with Stereo Estimation

- Comparison of the proposed method with other methods
 - Test with the dataset created. (hazy KITTI)



Image Dehazing combined with Stereo Estimation

- Comparison of the proposed method with other methods
 - Test with the foggy Cityscapes (real + synthetic haze)



Hazy input



AOD-Net [9]



MS-CNN [8]



DCP [7]



Non-local [15]



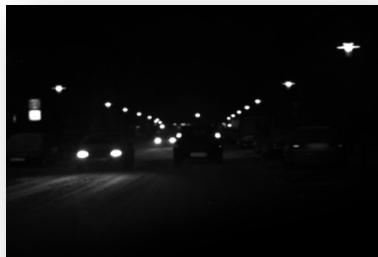
Ours(no connection)



Ours(disparity to dehazing)

Stereo Matching in Real Applications

- Which direction should we go?
 - Considering all the difficulties that are possible to occur
 - Remove rain drops, sun flares, photometric distortions, and etc.
 - The complexity of an algorithm increases enormously
 - Developing measures for estimating uncertainties (e.g., confidence measures) which can be used to identify wrong correspondences
 - The detection of wrong correspondences becomes more important as the **degree of ill-conditioning increases** (existing solutions fail)



Snowy night



Sun flare



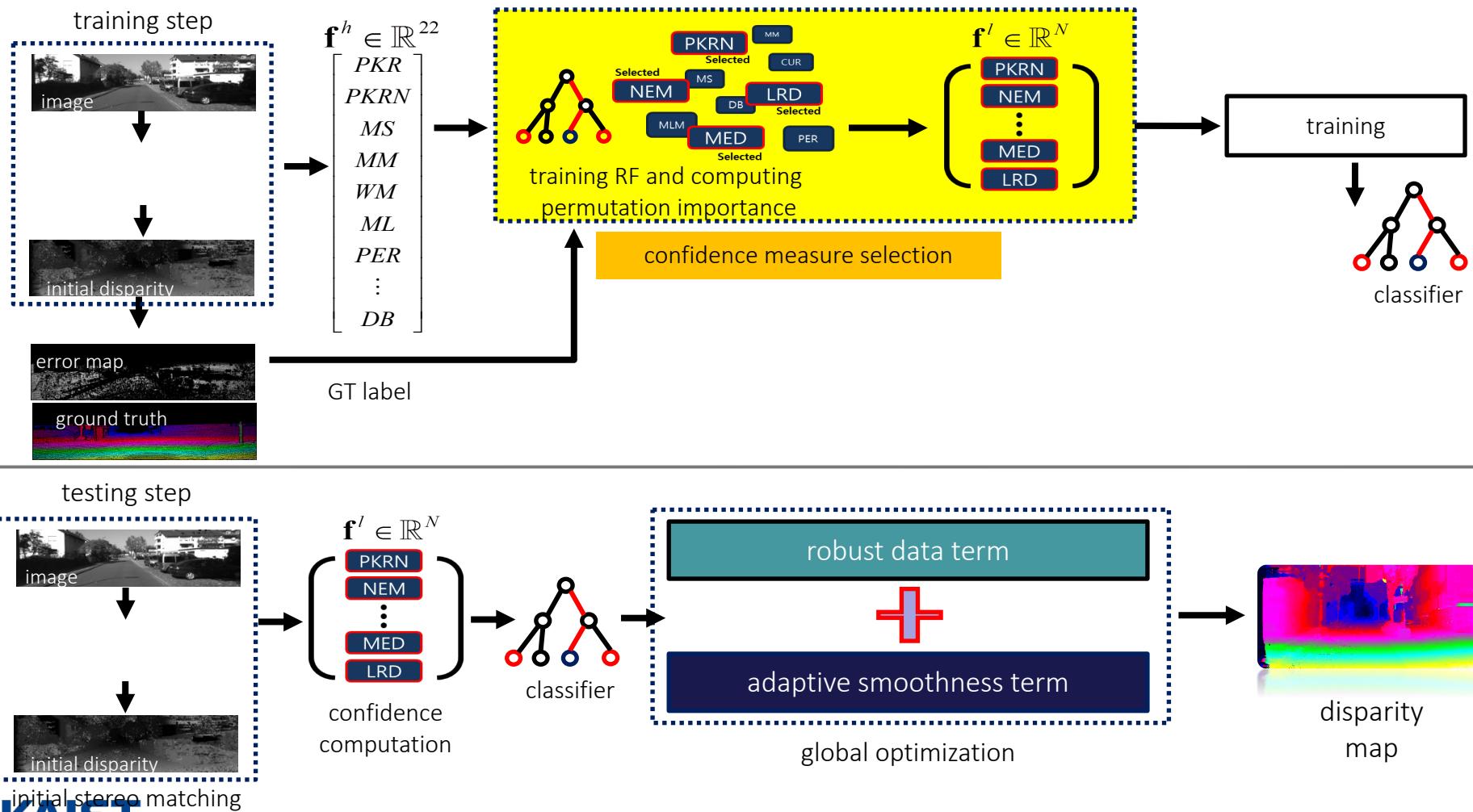
Snow



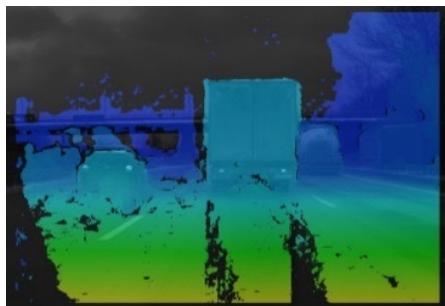
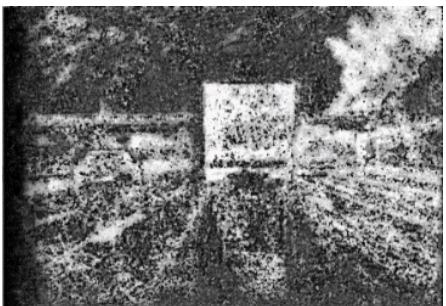
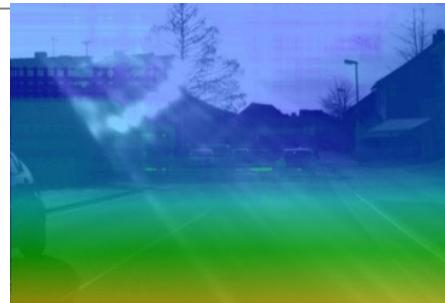
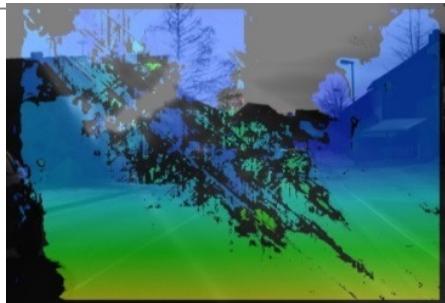
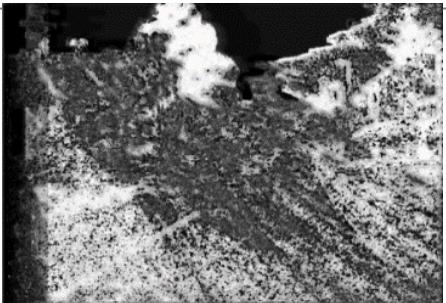
Rain

Stereo for Challenging Outdoor Scenes

[Park and Yoon, TPAMI 2019]
 [Park and Yoon, CVPR 2015]



Stereo for Challenging Outdoor Scenes

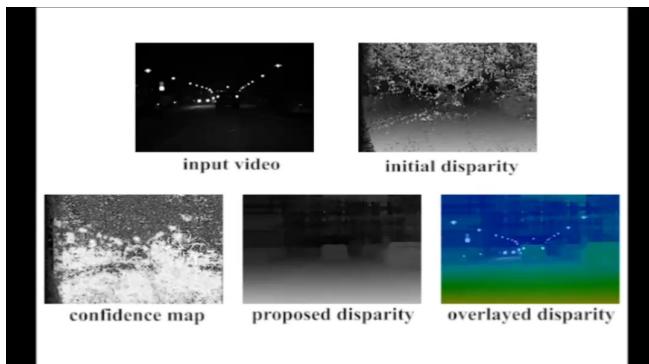


Input image

Predicted confidence

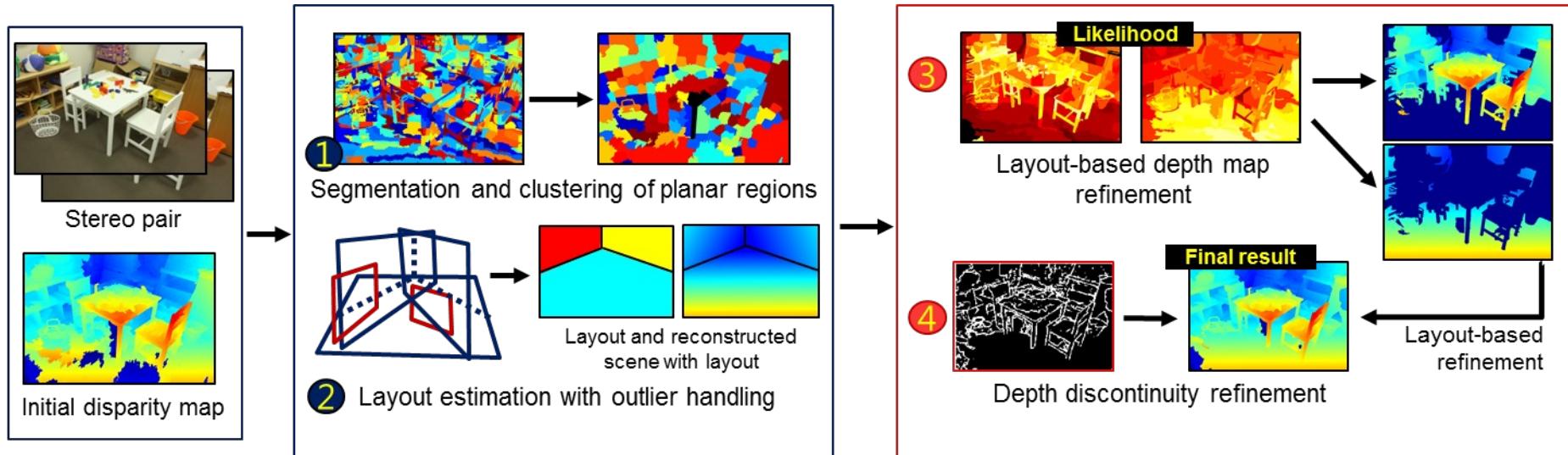
SGM

Our results

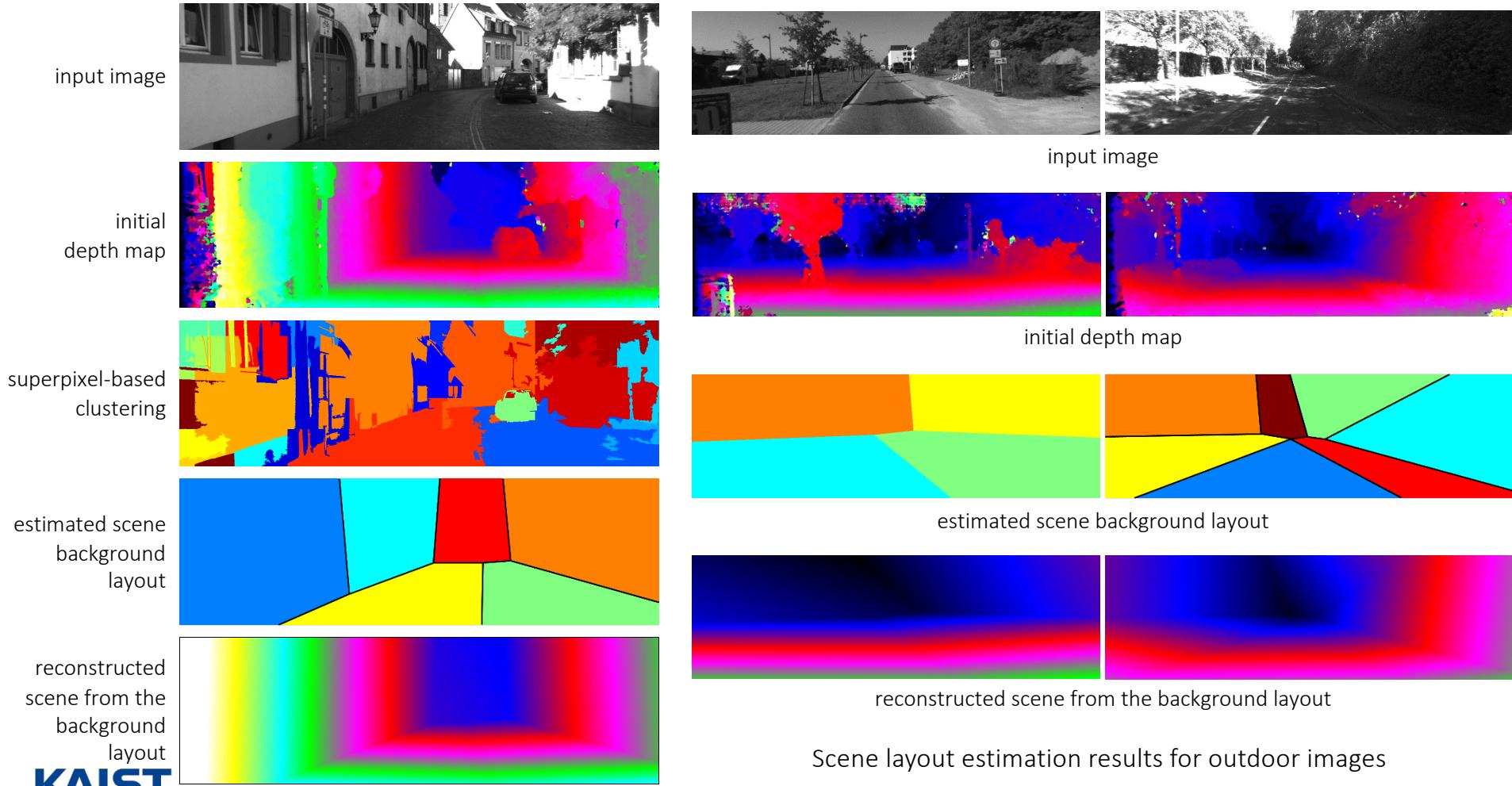


Stereo-based Scene Layout Estimation

[J.-K. Lee et al., ICCV 2017]
[Park and Yoon, CVIU 2019]



Stereo-based Scene Layout Estimation

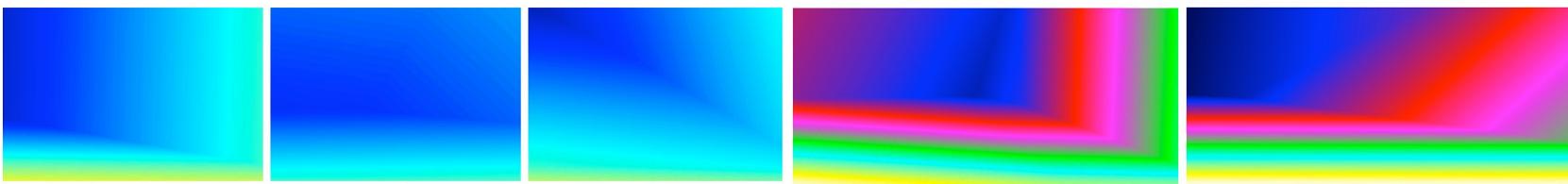


Layout-based Depth Refinement

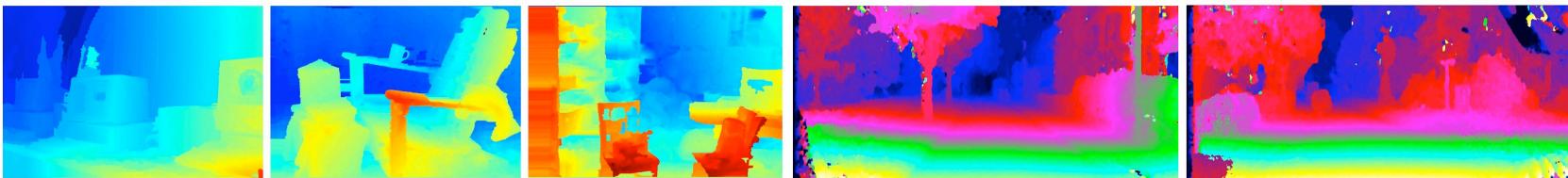
one of two
input images



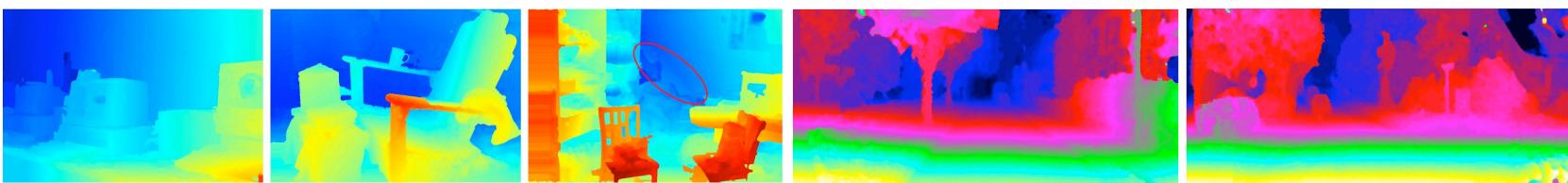
estimated scene
background
layout



initial
depth map

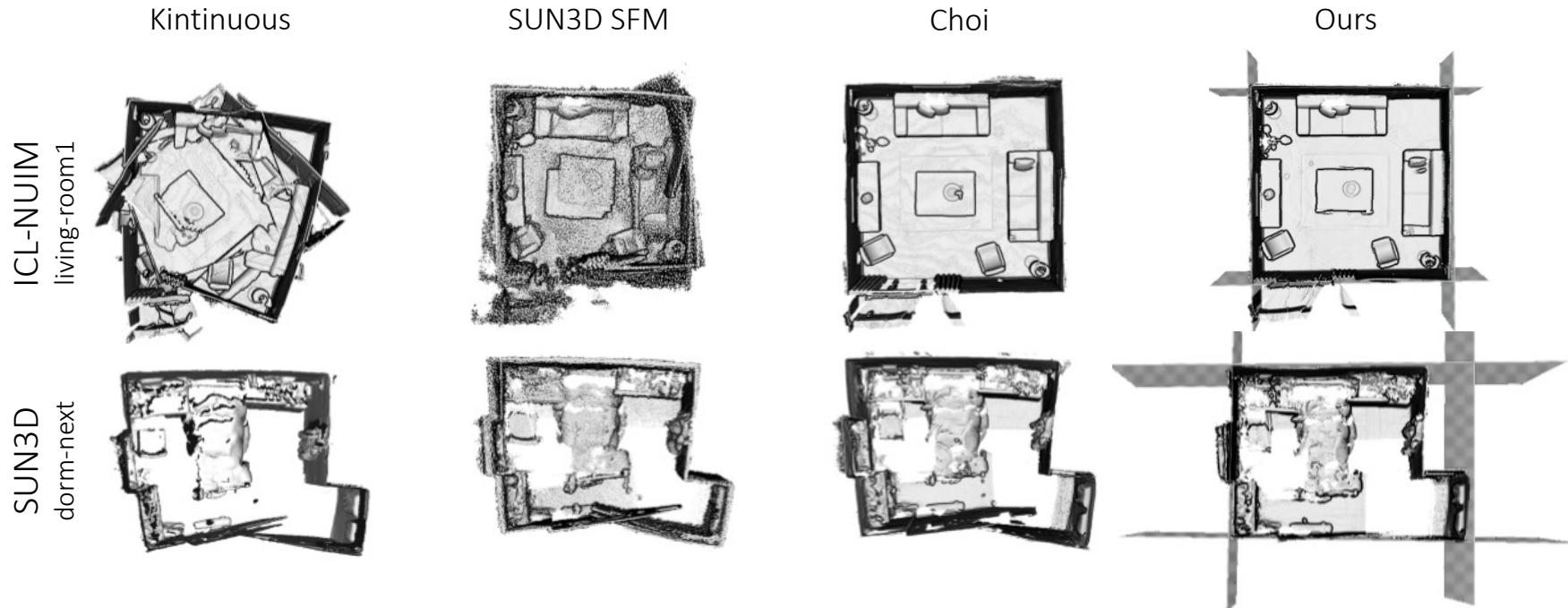


refined
depth map

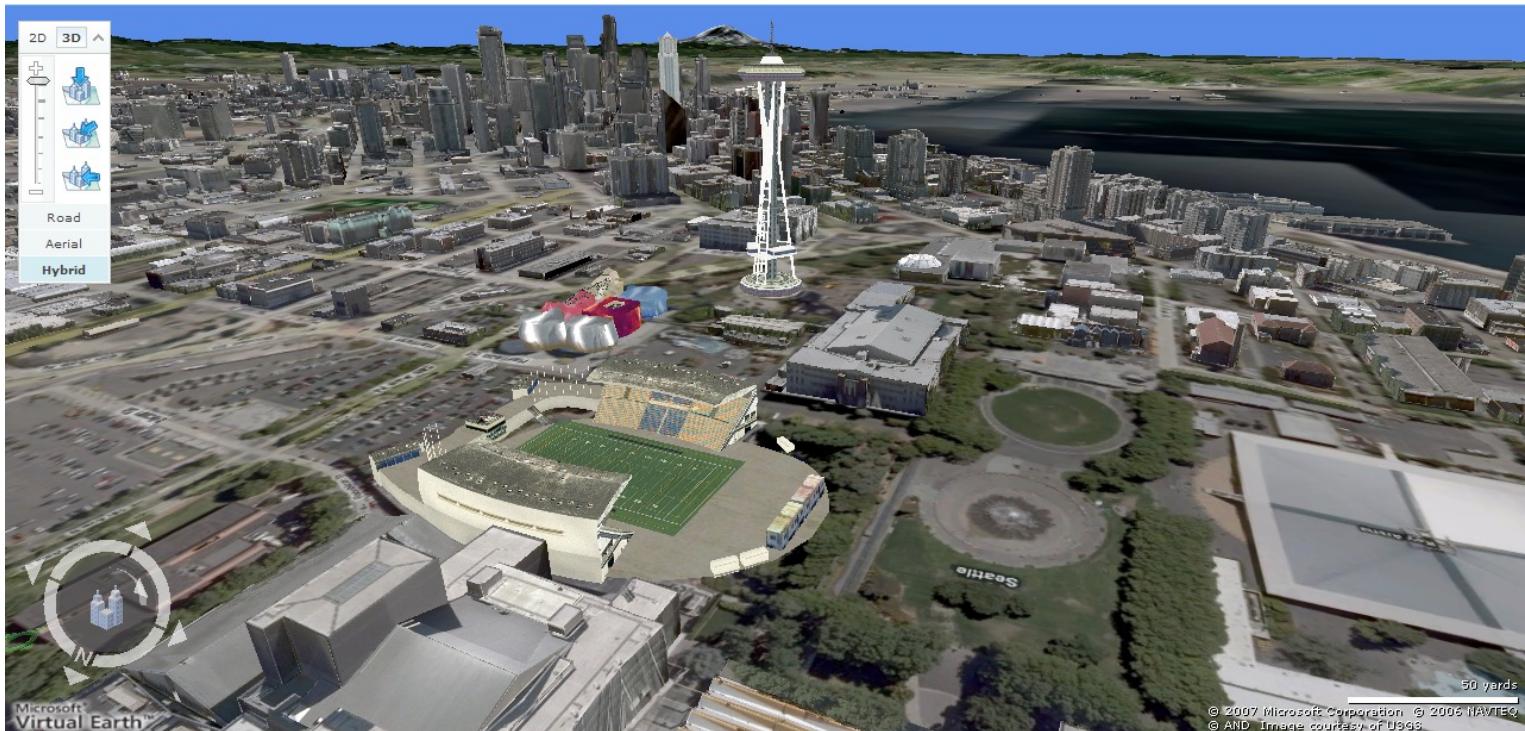


3D Indoor Reconstruction

- Qualitative evaluation in the ICL-NUIM and SUN3D datasets

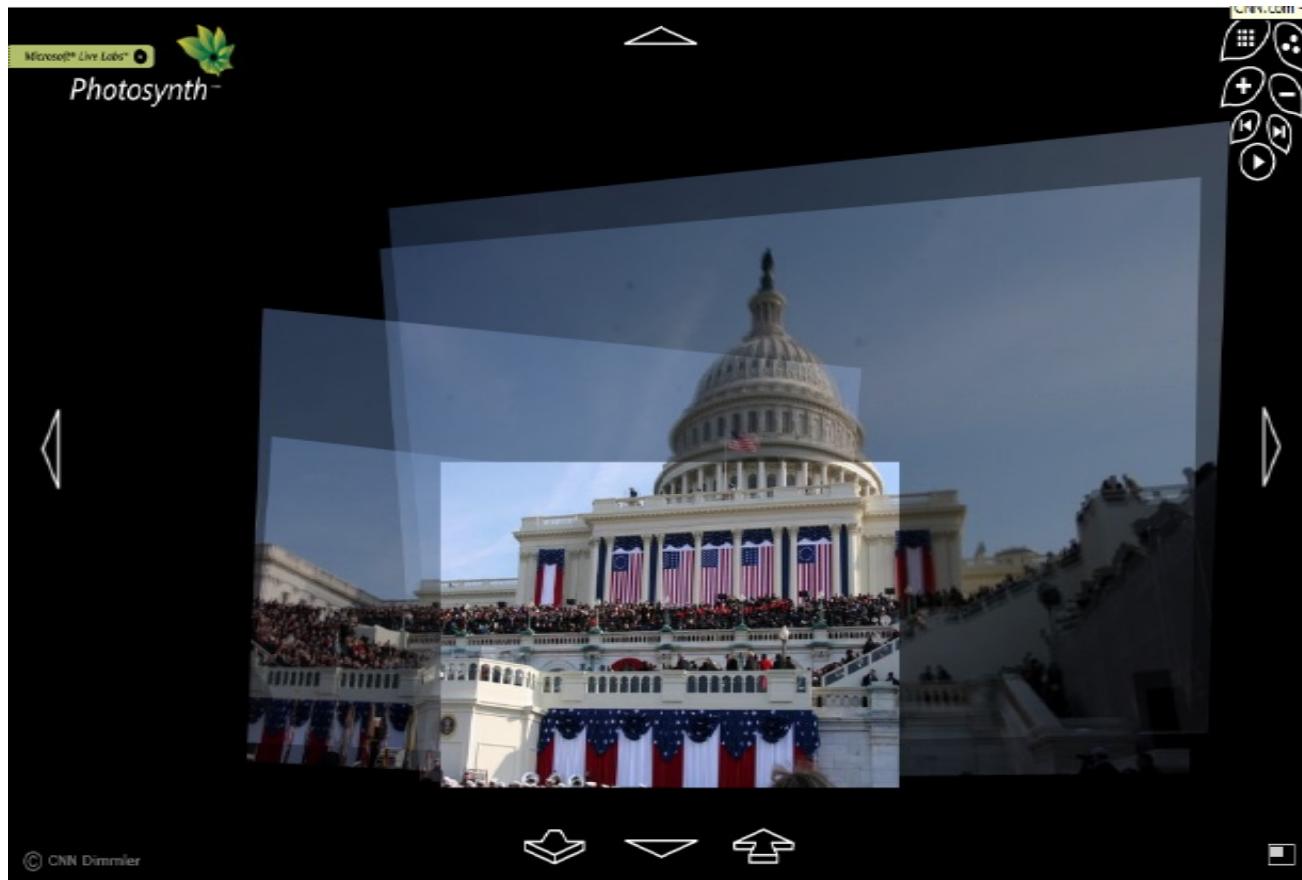


3D Urban Modeling



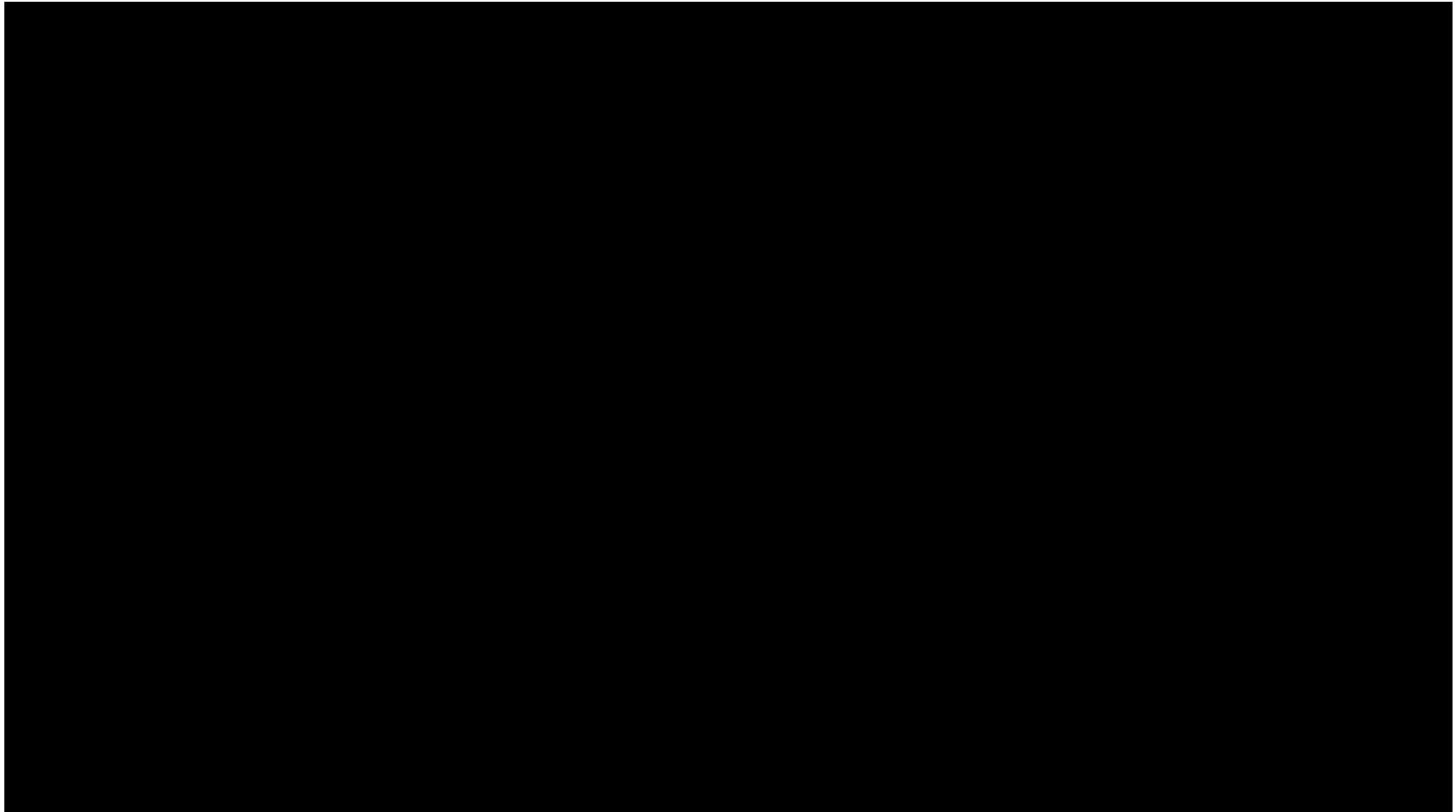
Bing maps, Google Streetview

3D Urban Modeling: Microsoft Photosynth



<http://labs.live.com/photosynth/>

3D Urban Modeling

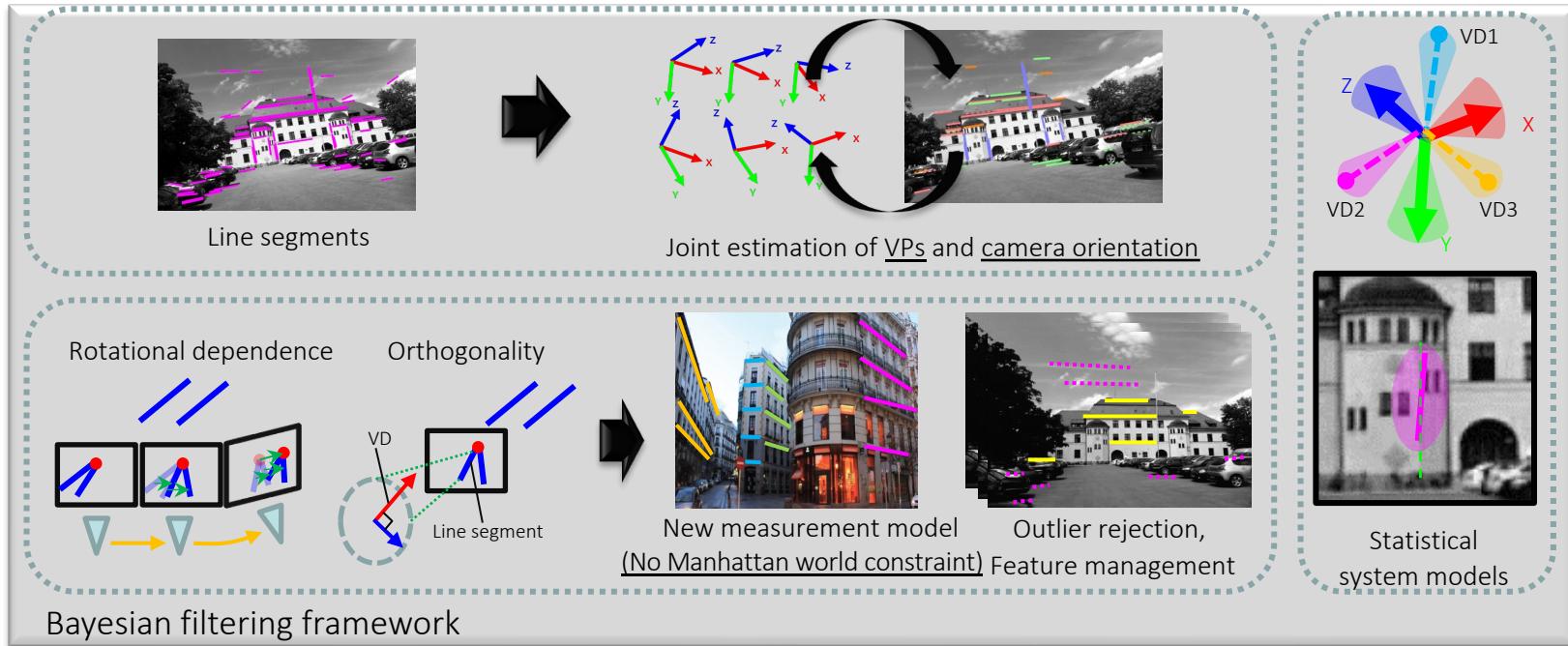


Motion

Joint Estimation of VPs and Camera Orientation

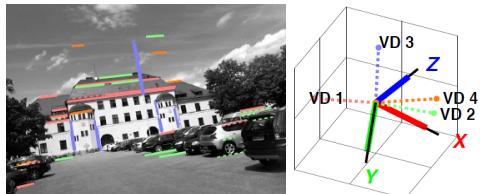
[Lee and Yoon, CVPR 2015]

[Lee and Yoon, IJCV 2019]



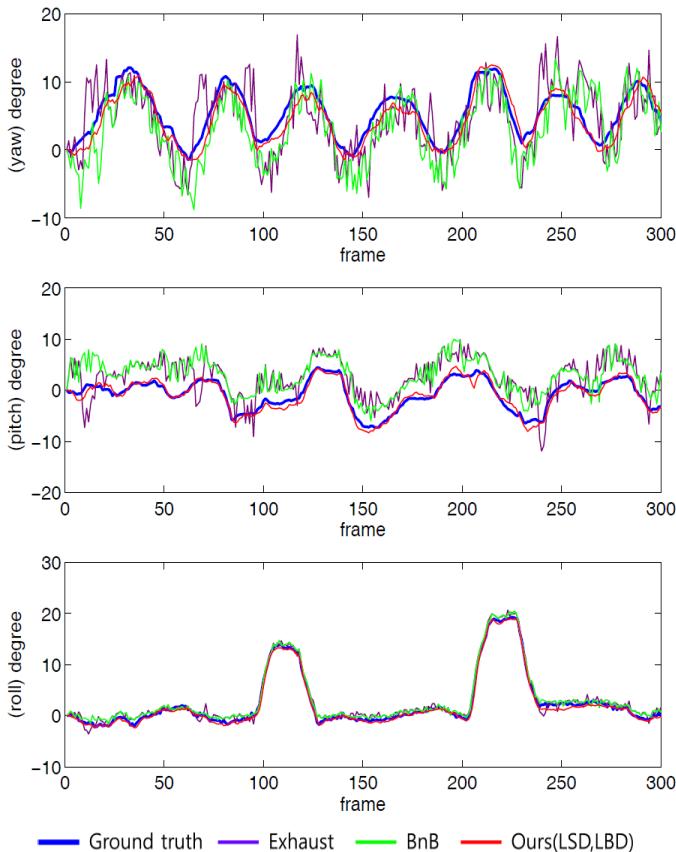
Estimate the camera orientation and VPs in sequential images:

1. Accurate
2. Robust to noisy or spurious line segments
3. No Manhattan world assumption
4. Fast

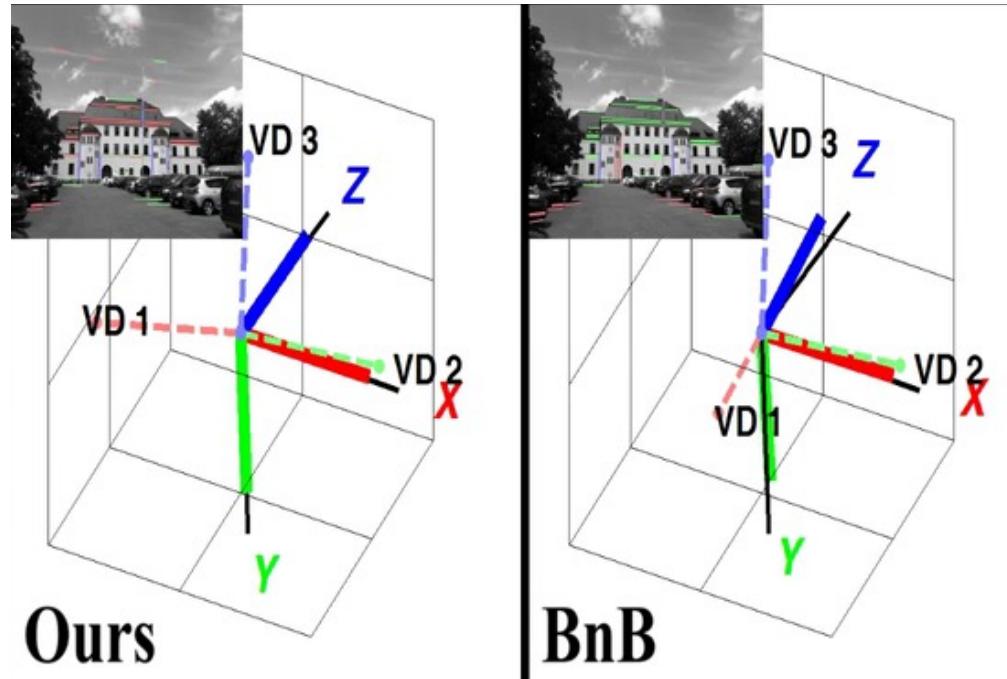


Joint Estimation of VPs and Camera Orientation

- ✓ Qualitative evaluation (in Metaio-100)
✓ Camera orientation estimates



- ✓ Experimental results of the proposed method and BnB[3]



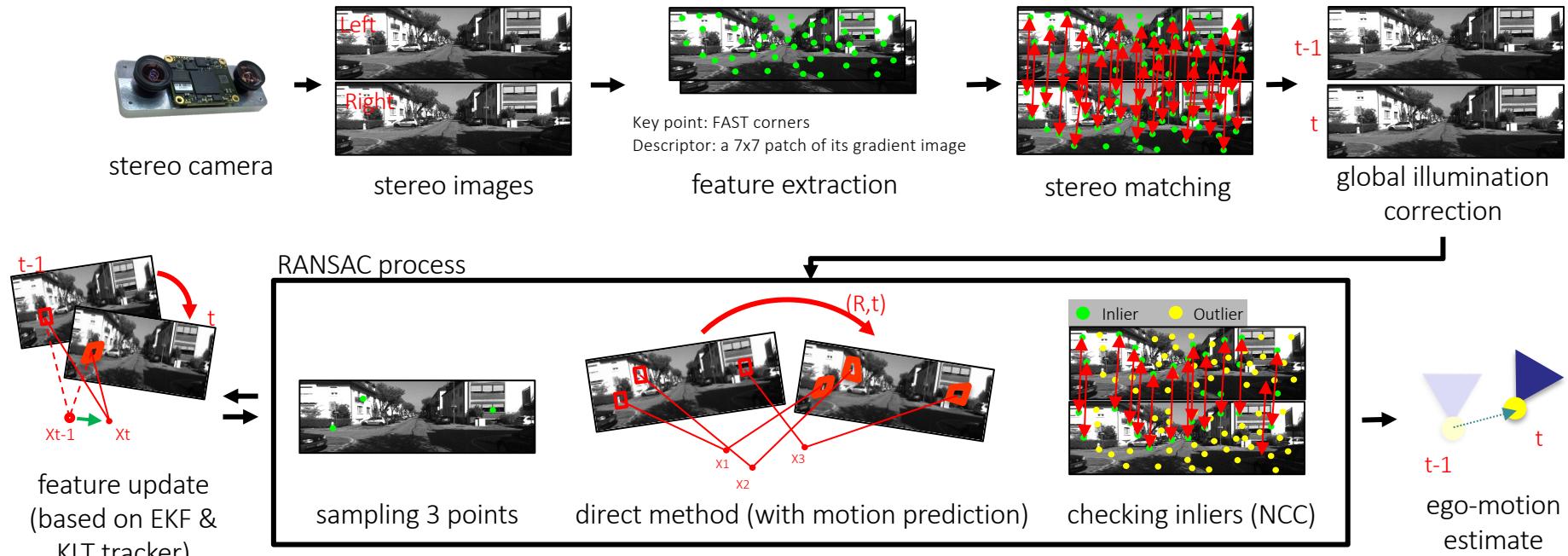
└ : Ground truth
pose

└ : Estimated pose

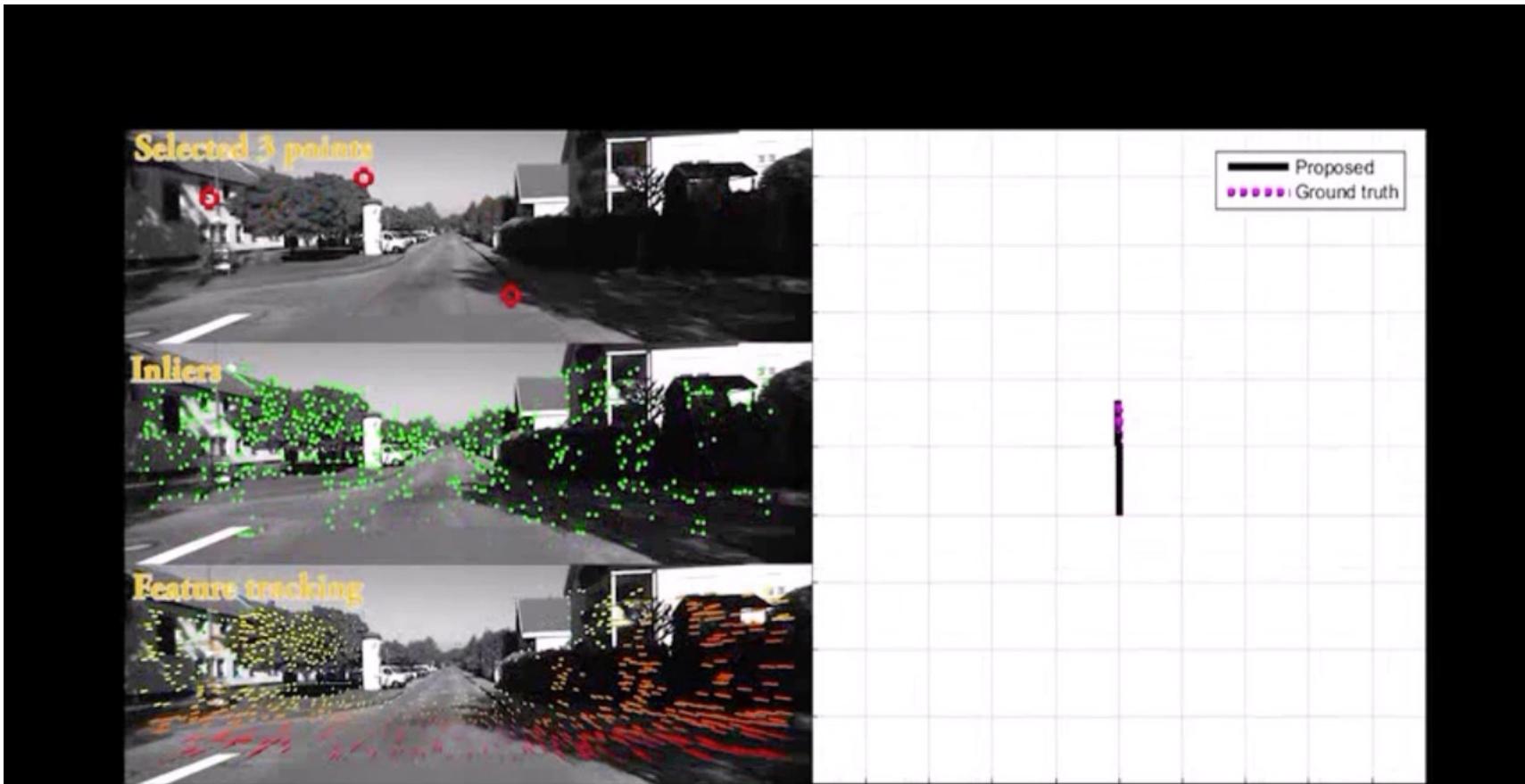
Three-point Direct Stereo Visual Odometry (SVO)

[Lee and Yoon, BMVC 2016]

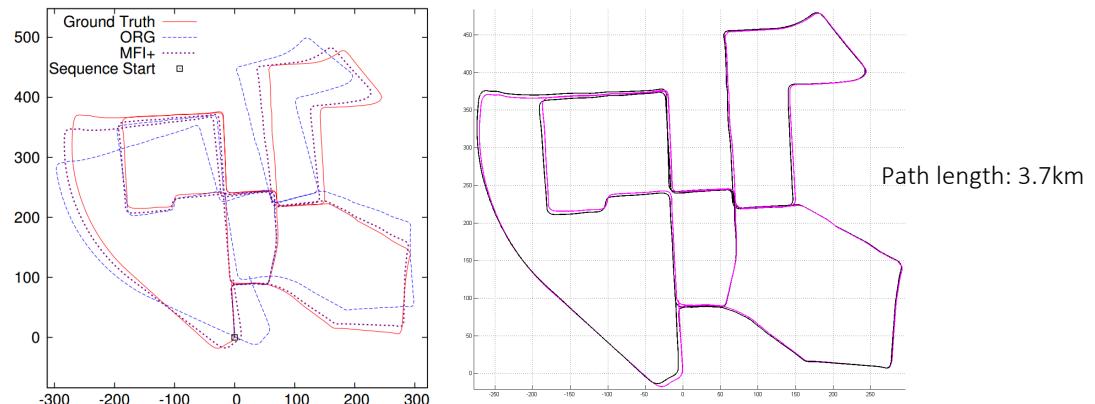
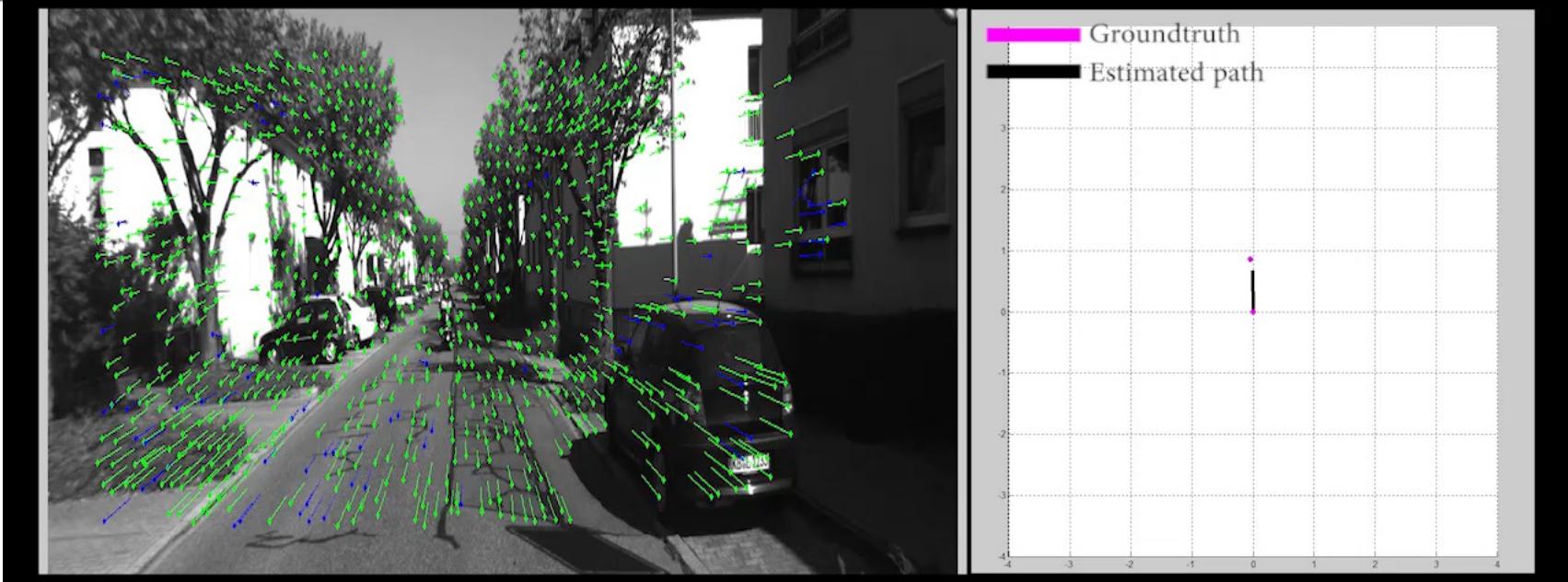
- Given sequential stereo images, finding the ego-motion of the camera
- Accurately estimating the ego-motion by using a minimum number of uncontaminated feature points



Three-point Direct Stereo Visual Odometry (SVO)



Three-point Direct Stereo Visual Odometry (SVO)

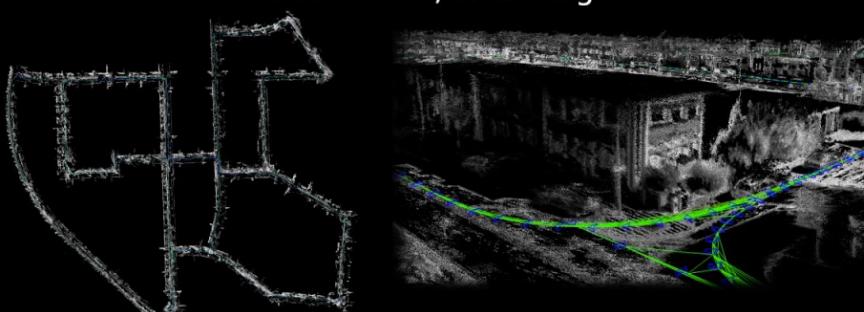


SLAM (Simultaneous Localization And Mapping)

- In robotic mapping and navigation, simultaneous localization and mapping (SLAM) is the computational problem of constructing or updating a map of an unknown environment while simultaneously keeping track of an agent's location within it.

**Large-Scale Direct SLAM
with Stereo Cameras**

Jakob Engel, Jörg Stückler, Daniel Cremers
IROS 2015, Hamburg

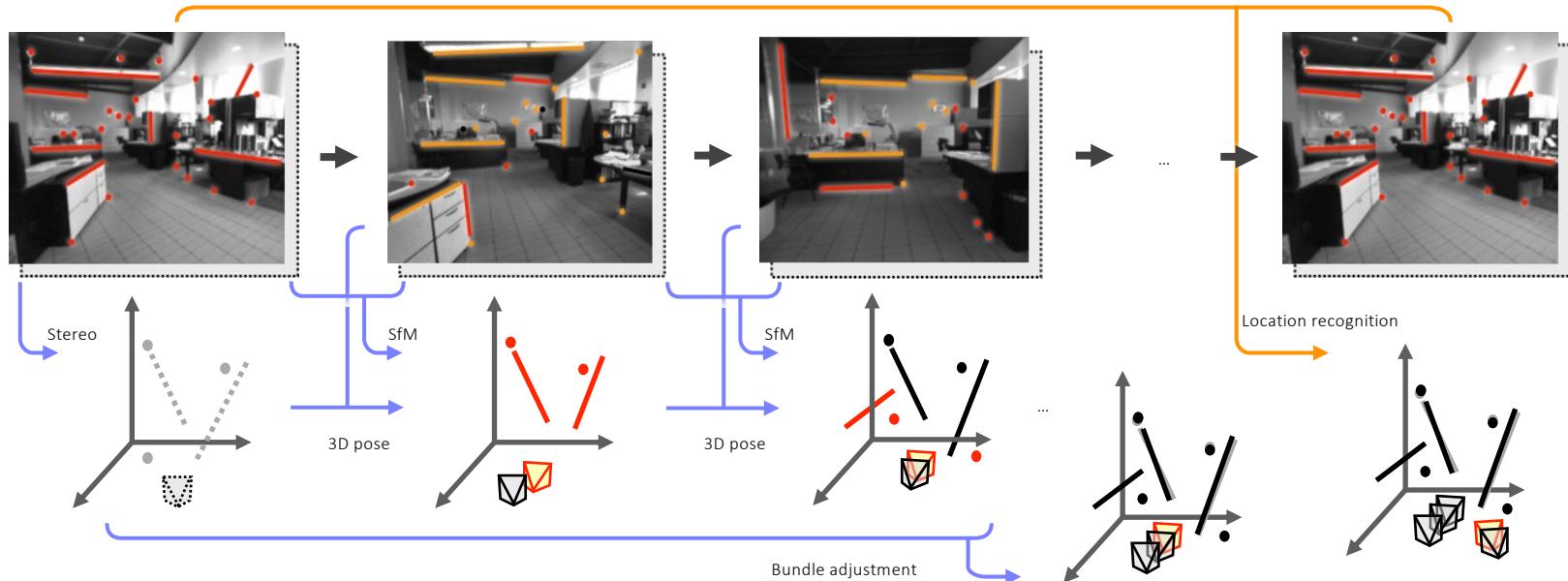


Computer Vision Group
Technical University Munich



Visual SLAM

- SLAM using a camera as the main sensor.
 - Monocular, stereo, and RGB+D sensors.
 - Visual features (points) vs. depth maps.
 - 3D sparse environment map vs. dense occupancy map.



Optical Flow Estimation

- FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks, CVPR 2017

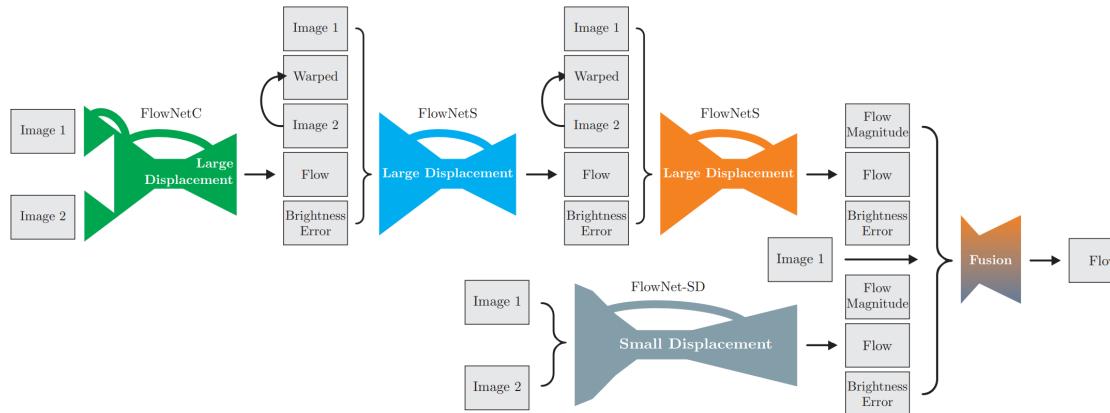


Figure 2. Schematic view of complete architecture: To compute large displacement optical flow we combine multiple FlowNets. Braces indicate concatenation of inputs. *Brightness Error* is the difference between the first image and the second image warped with the previously estimated flow. To optimally deal with small displacements, we introduce smaller strides in the beginning and convolutions between upconvolutions into the FlowNetS architecture. Finally we apply a small fusion network to provide the final estimate.



Figure 1. We present an extension of FlowNet. FlowNet 2.0 yields smooth flow fields, preserves fine motion details and runs at 8 to 140fps. The accuracy on this example is four times higher than with the original FlowNet.



Figure 6. Examples of flow fields from different methods estimated on Sintel. FlowNet2 performs similar to FlowFields and is able to extract fine details, while methods running at comparable speeds perform much worse (PCA-Flow and FlowNetS).

Loop-Net: Joint Unsupervised Disparity and Optical Flow Estimation of Stereo Videos

[T. Kim et al., RA-L 2020]

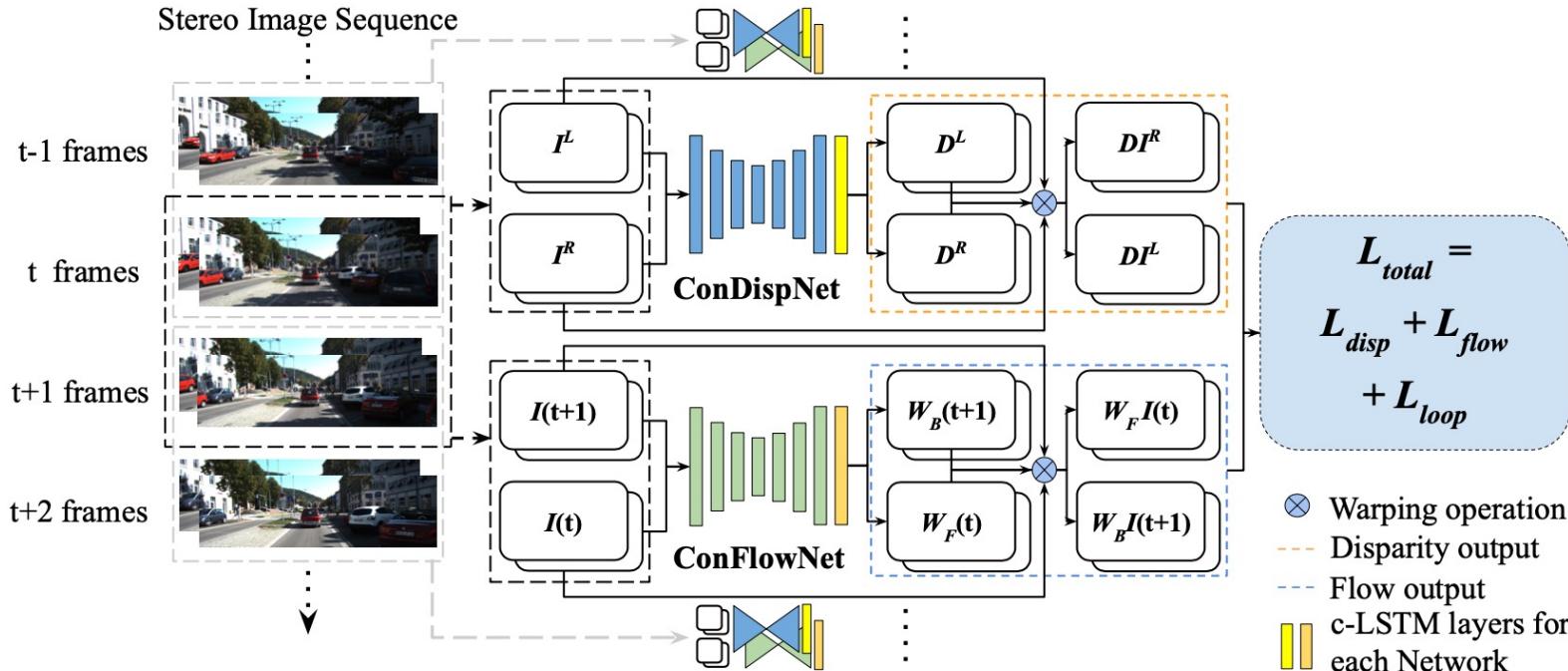


Figure 2: Overview of our framework. Our network consists of two subnetworks, ConDispNet and ConFlowNet, that estimate the disparity and the optical flow, respectively. Using the estimated dense correspondence map from each subnetwork, the reference image is warped to the target image. DI^s and $WI(t)$ denote the warped image corresponding to the target image I^s and $I(t)$, respectively.

Loop-Net: Joint Unsupervised Disparity and Optical Flow Estimation of Stereo Videos

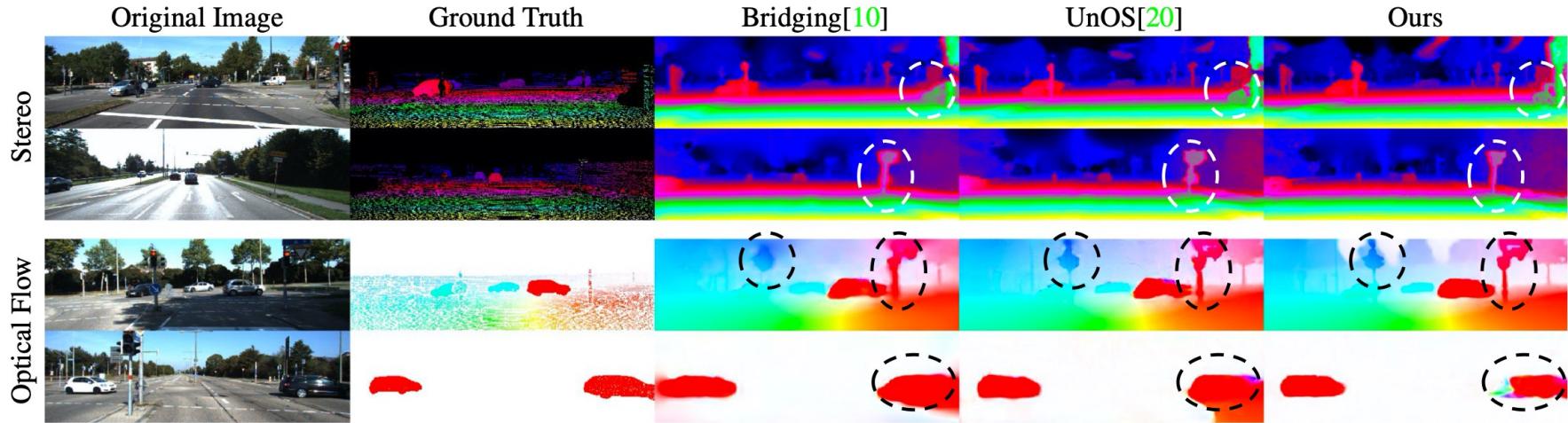


Table 4: Quantitative evaluation of the optical flow task on KITTI 2015 training/test dataset. “noc”, “occ” denote non-occlusion and occlusion regions, respectively. The boldface and underscore denote the best and second best performances, respectively.

| Method | Joint learning | Super vised | Train EPE | Train F1 | Train EPE noc | Test F1 |
|----------------------------|----------------|-------------|-------------|---------------|---------------|---------------|
| Flownet2 [6] | | ✓ | 10.06 | 30.37% | - | - |
| SpyNet [13] | | ✓ | 20.56 | 44.78% | - | - |
| UnFlow-C [12] | | | 8.80 | 28.94% | - | 29.46% |
| UnFlow-CSS [12] | | | 8.10 | 23.27% | - | - |
| GeoNet [25] | ✓ | | 10.81 | - | 8.05 | - |
| Wang <i>et al.</i> [21] | | | 8.88 | - | - | - |
| Jainai <i>et al.</i> [7] | | | 6.59 | - | 3.22 | 22.94% |
| DFnet [34] | ✓ | | 8.98 | 26.01% | - | 25.70% |
| CC [14] | ✓ | | 6.21 | 26.41% | - | - |
| CC-uft [14] | ✓ | | 5.66 | 20.93% | - | 25.27% |
| Bridging-R [10] | ✓ | | 7.02 | 27.34% | 4.26 | - |
| Bridging-P [10] | ✓ | | 6.66 | 21.50% | - | - |
| UnOS [20] | ✓ | | 5.58 | - | - | 18.00% |
| Ours (Flow only - SSIM) | | | 6.15 | 19.91% | 2.94 | - |
| (Flow only - ZNCC) | | | 6.08 | 18.63% | 2.78 | - |
| (Full - w/o c-LSTM) | | | 5.77 | 18.09% | 2.88 | - |
| (Full-1) | ✓ | | 5.53 | 18.00% | 2.63 | - |
| (Full-2) | ✓ | | 5.27 | 17.57% | 2.53 | 18.42% |

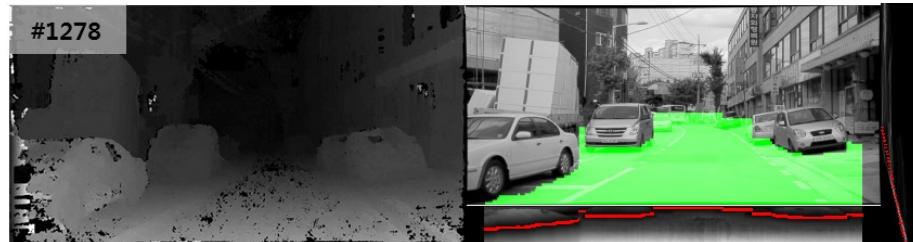
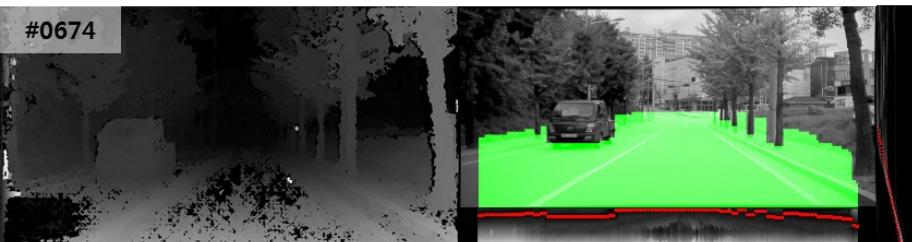
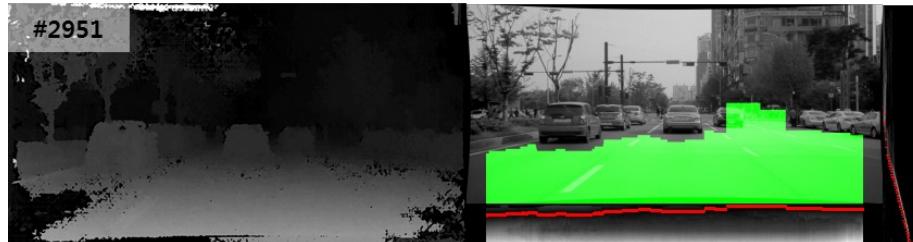
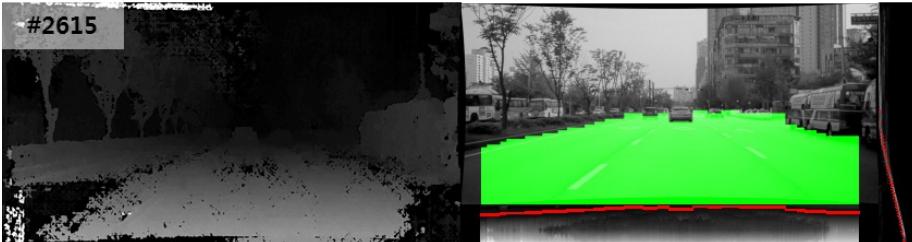
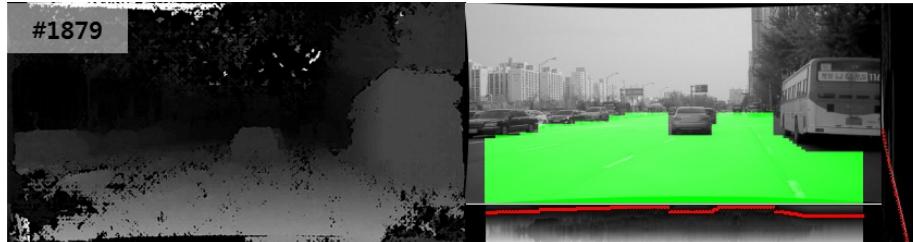
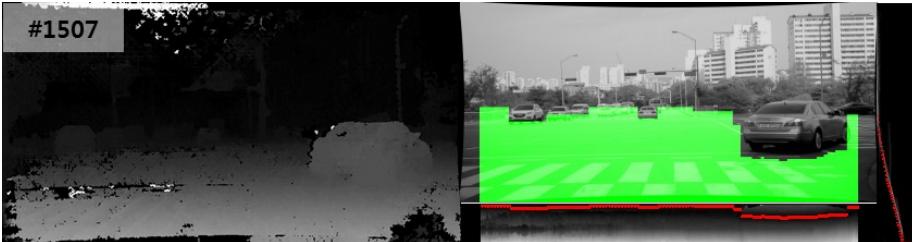
Table 5: Quantitative evaluation of spatiotemporal consistency on KITTI 2015 dataset. ξ is a set of LC(loop consistency). μ , med denote mean and median, respectively. Note that high accuracy does not guarantee consistent results.

| Method | Higher the better | | | Lower the better | |
|-----------------|-------------------|---------------|---------------|------------------|---------------|
| | $\xi < 2$ | $\xi < 3$ | $\xi < 4$ | $\mu(\xi)$ | $med(\xi)$ |
| Bridging-R [10] | 0.5114 | 0.5608 | 0.5964 | 4.9295 | 1.3659 |
| Bridging-P [10] | 0.5868 | 0.6213 | 0.6475 | 4.4213 | 0.9663 |
| UnOS [20] | 0.3412 | 0.3822 | 0.4184 | 16.6032 | 6.2568 |
| Ours(Full-2) | 0.7581 | 0.7822 | 0.8003 | 3.4095 | 0.7855 |

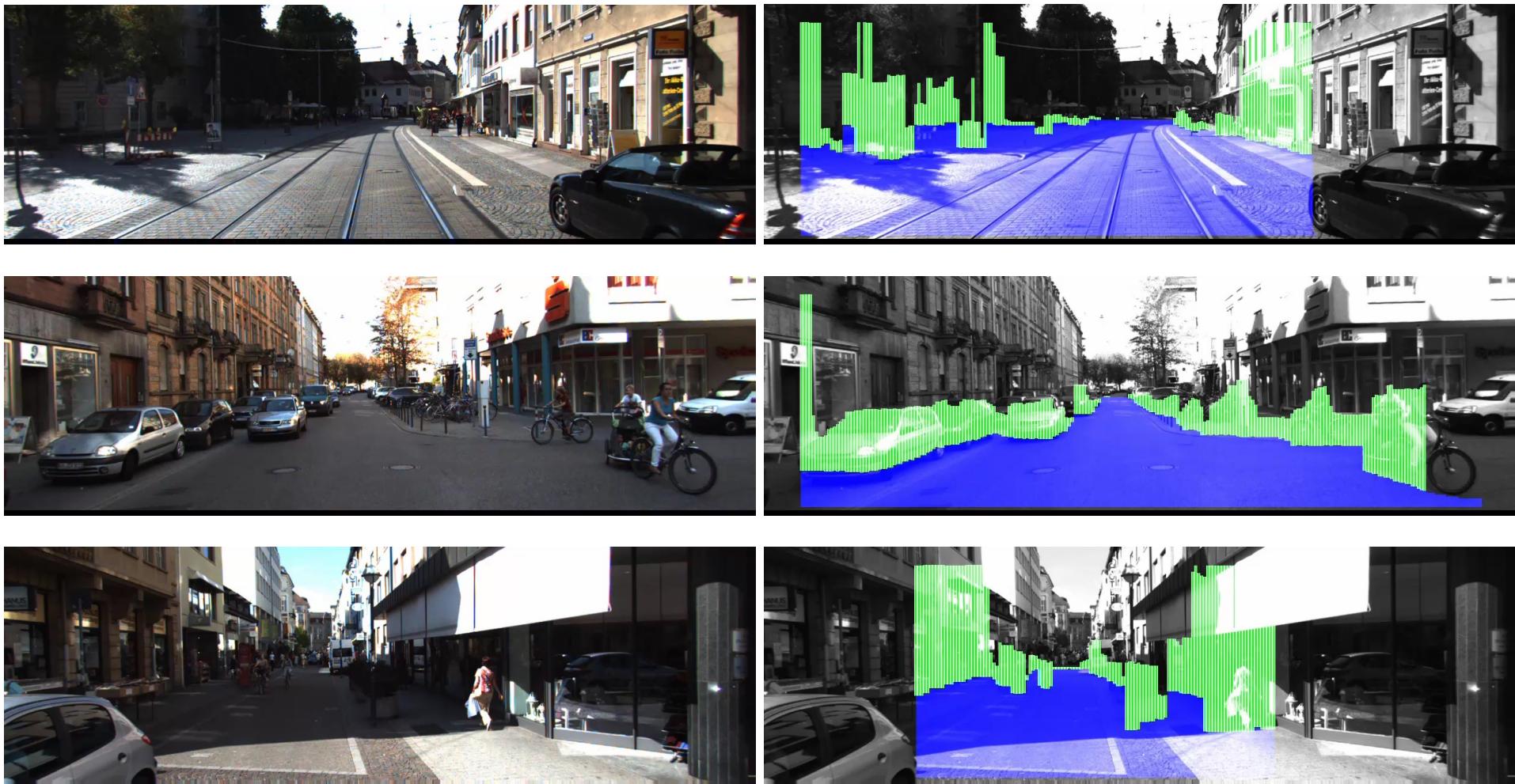
Scene

Ground Plane Detection

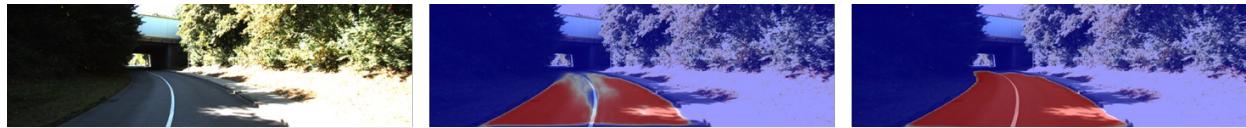
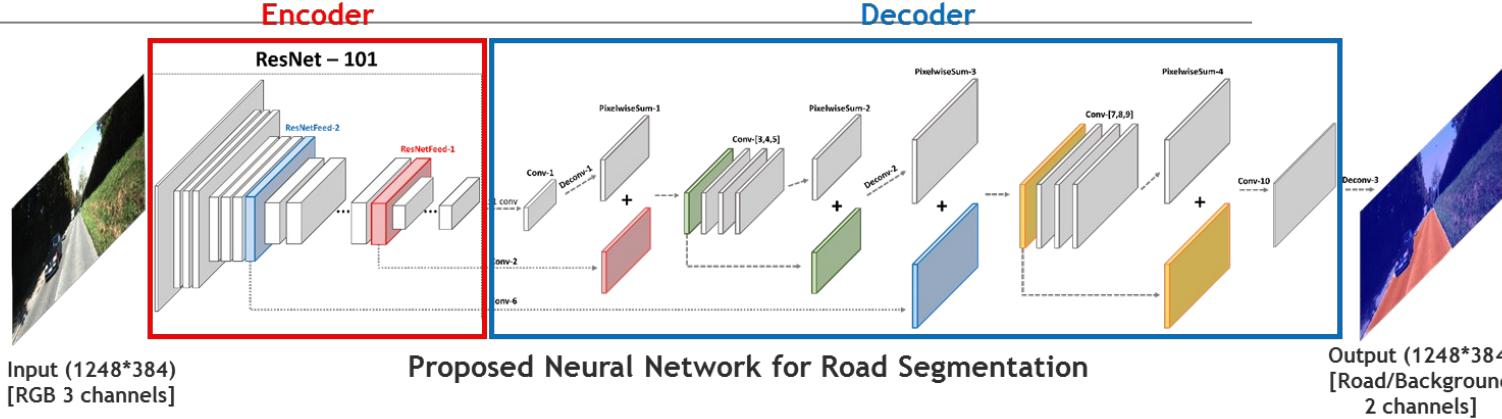
- Detection results with stereo disparity maps



Free Space & Obstacle Region Estimation



Road Segmentation using a Single Image based on Deep Learning



(a) Input

(b) MultiNet[13]

(c) Proposed method



(d) Input

(e) MultiNet[13]

(f) Proposed method



(g) Input

(h) MultiNet[13]

(i) Proposed method

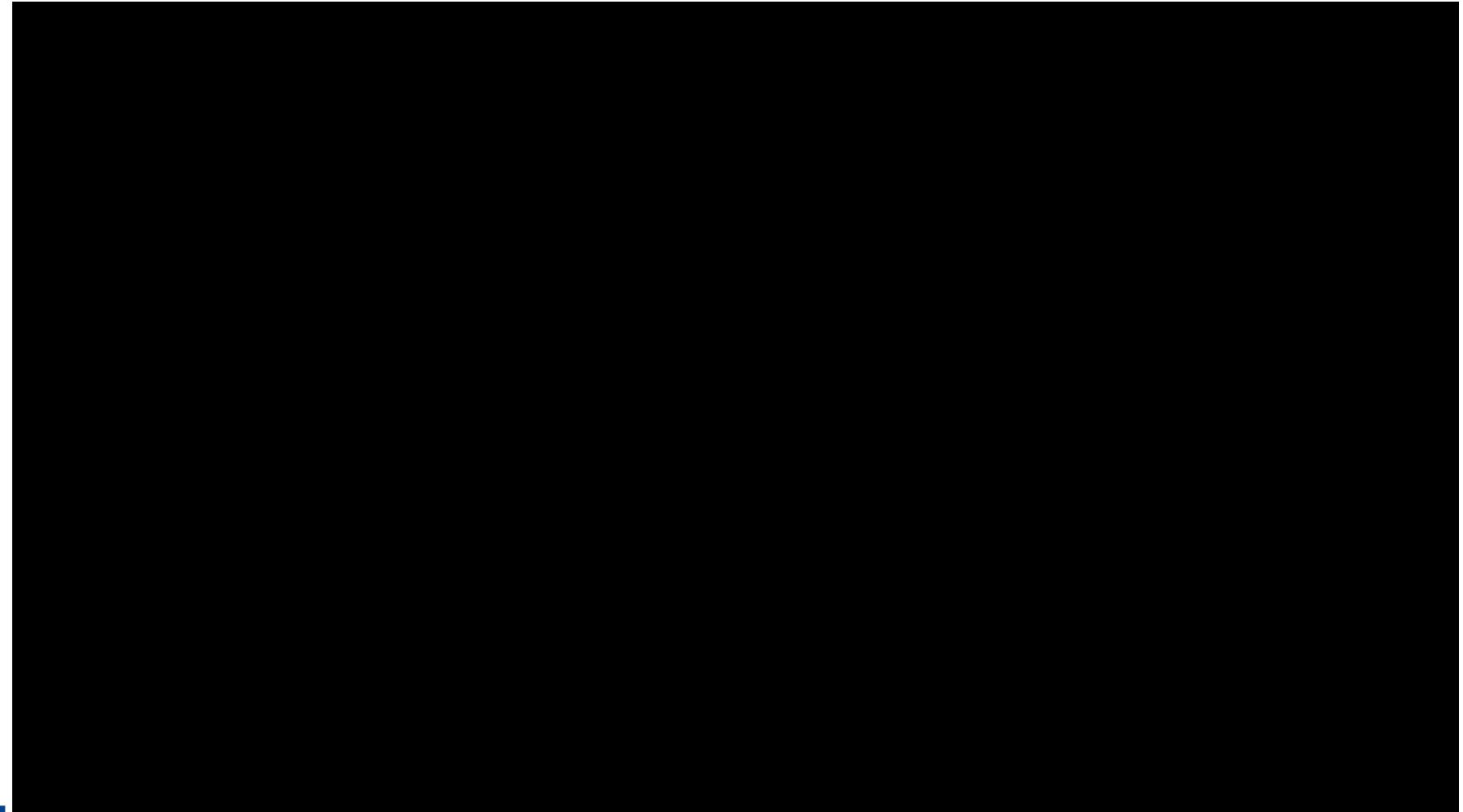


(j) Input

(k) MultiNet[13]

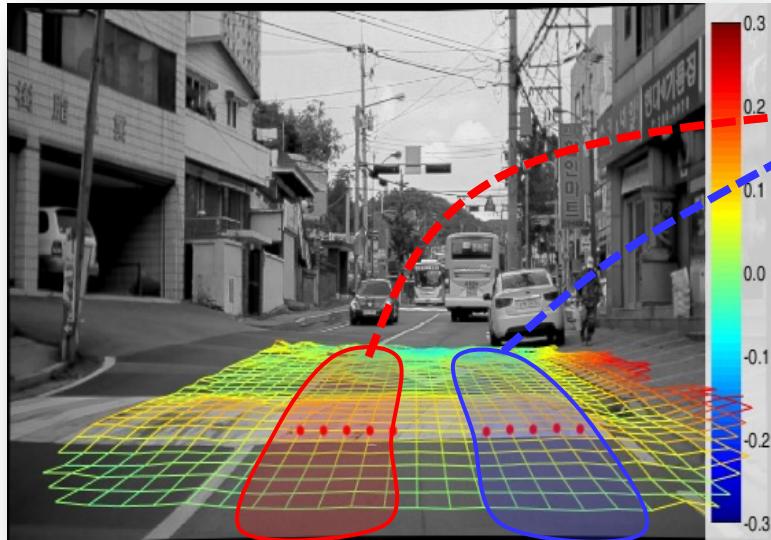
(l) Proposed method

Road Segmentation using a Single Image based on Deep Learning

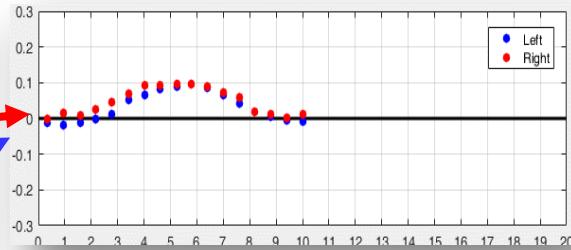


Road Surface Profile Estimation

- Road surface profile estimation
 - Estimating a road surface profile using depth images obtained from a stereo camera



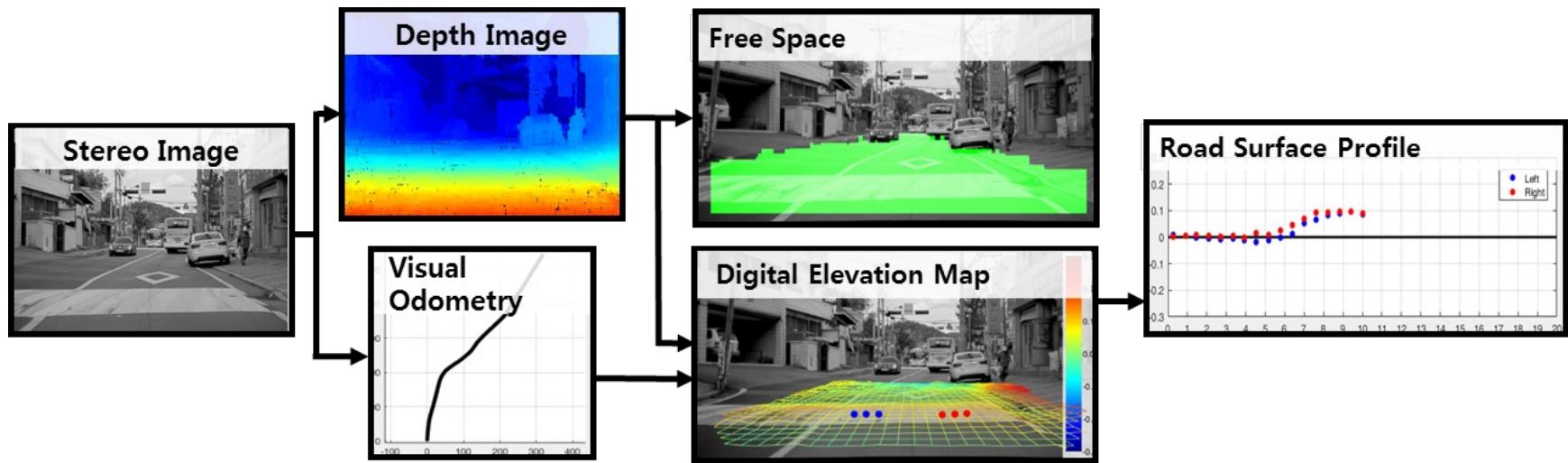
digital elevation map estimation
in a real-world driving scene



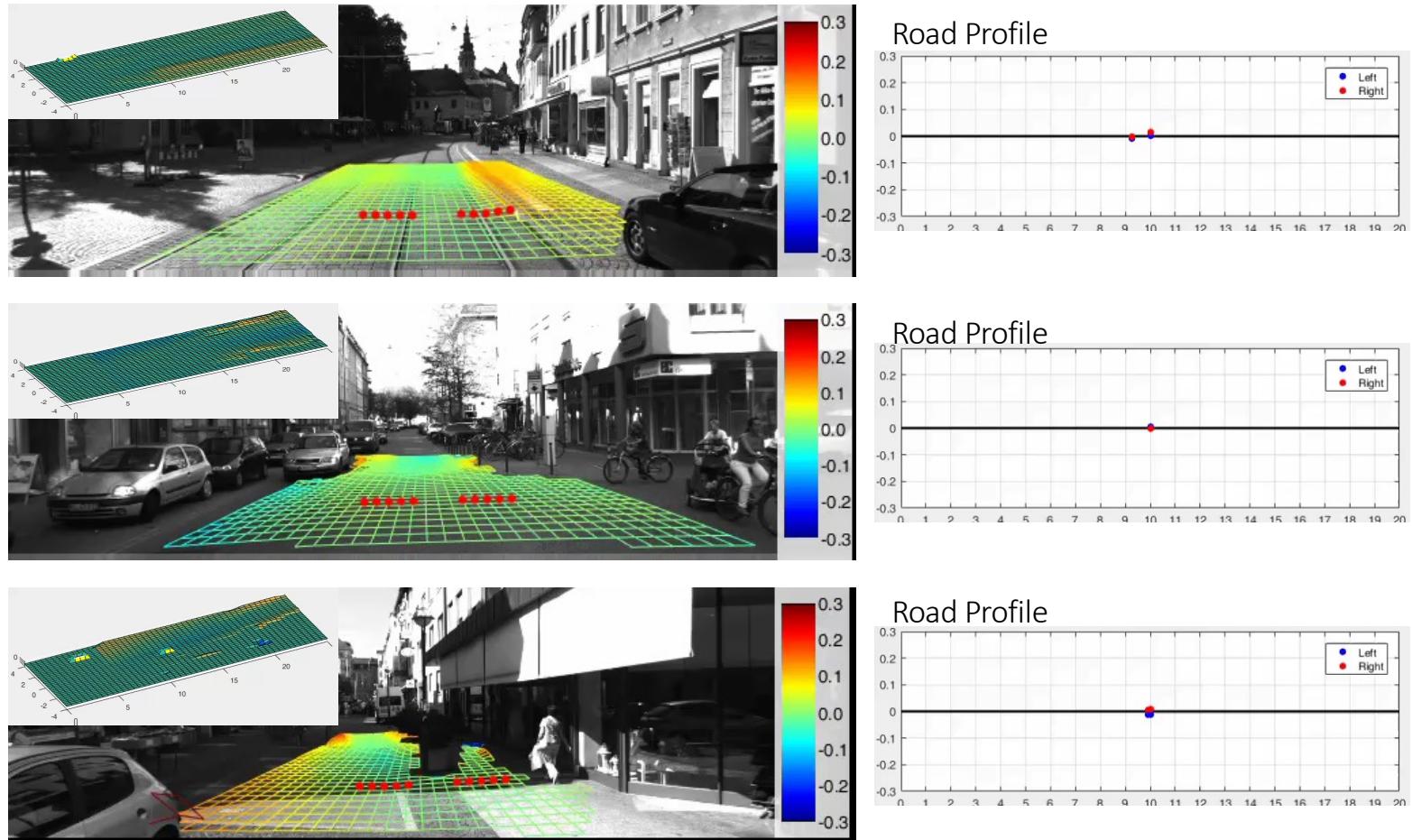
road surface profile
for front wheels' trajectories

Road Surface Profile Estimation

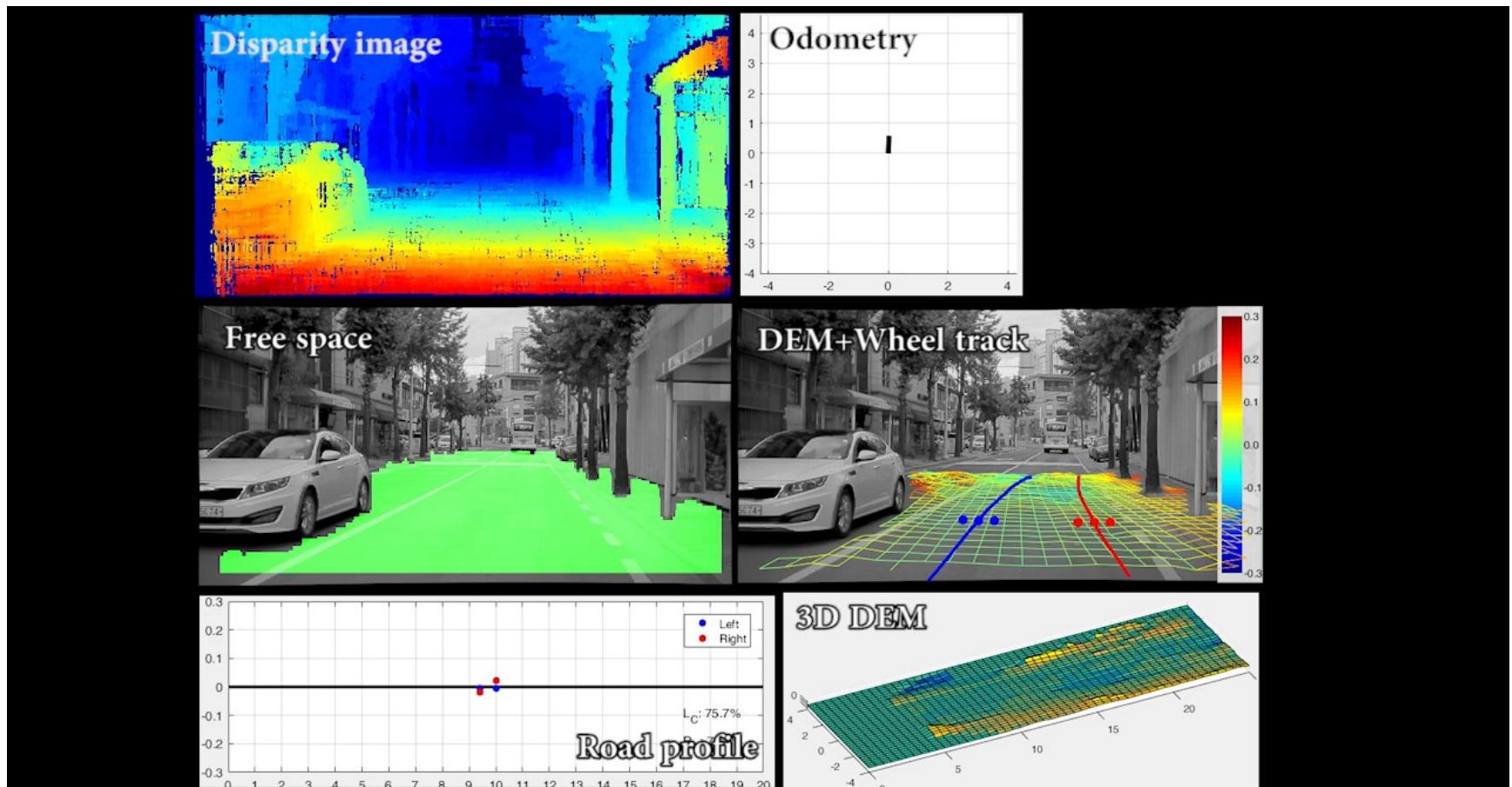
- Road surface profile estimation
 - Estimating a road surface profile using depth images obtained from a stereo camera



Ground Plane Detection and DEM Estimation



Road Surface Profile Estimation



Semantic Segmentation

Image



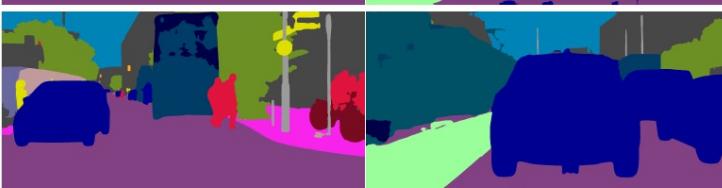
HRNetV2



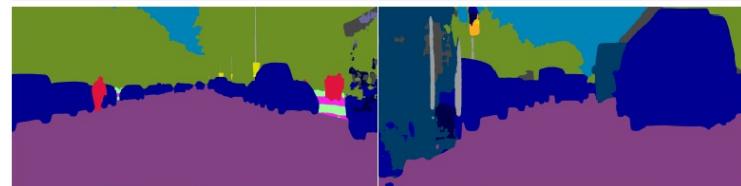
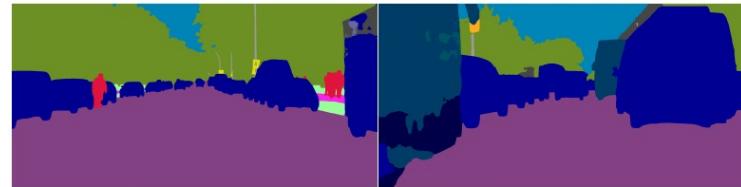
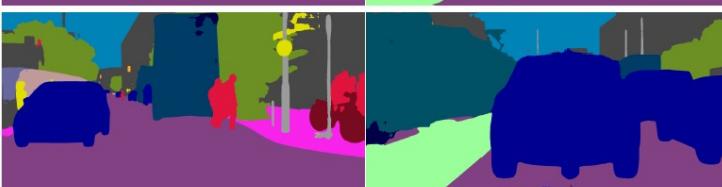
PSPNet



HRNetV2 \oplus PSPNet

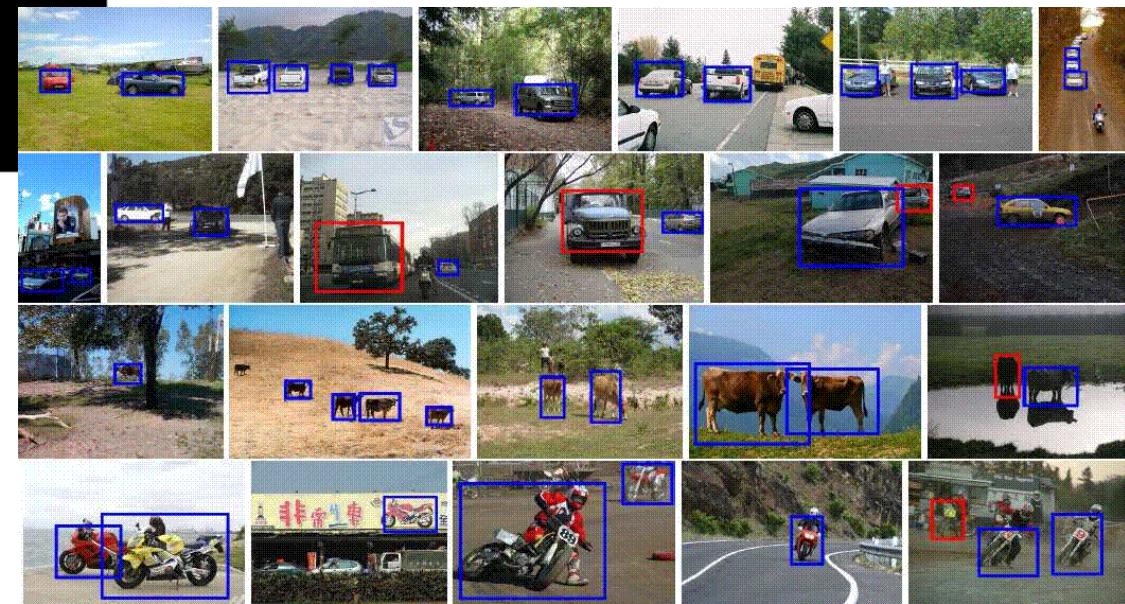
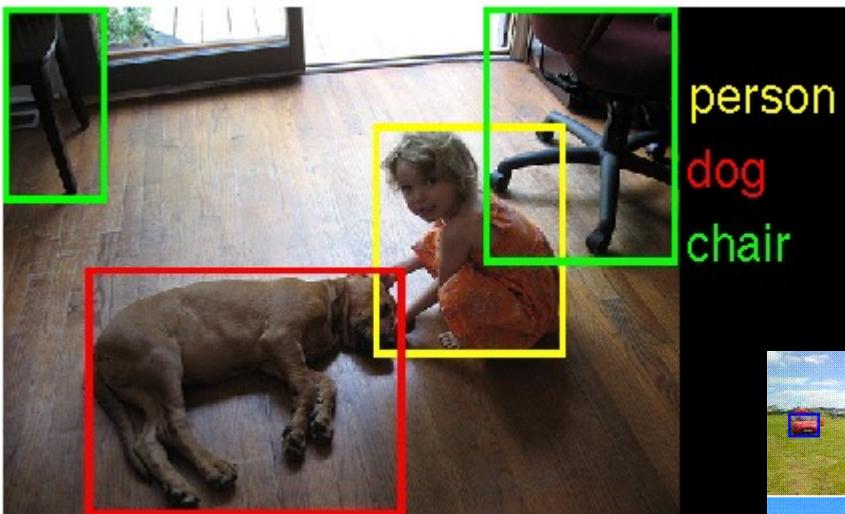


Proposed
 $F_{f+R}(\text{HRNetV2}, \text{PSPNet})$

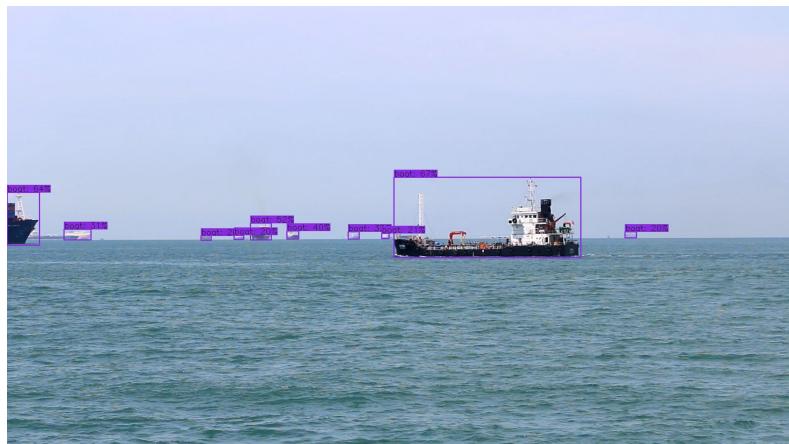
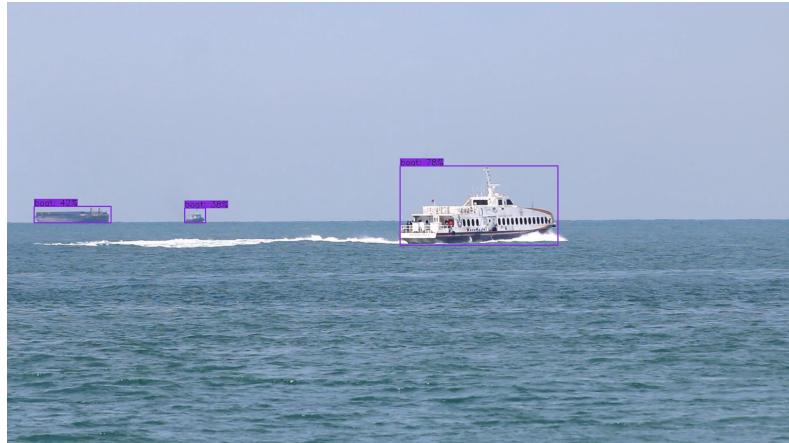


Object

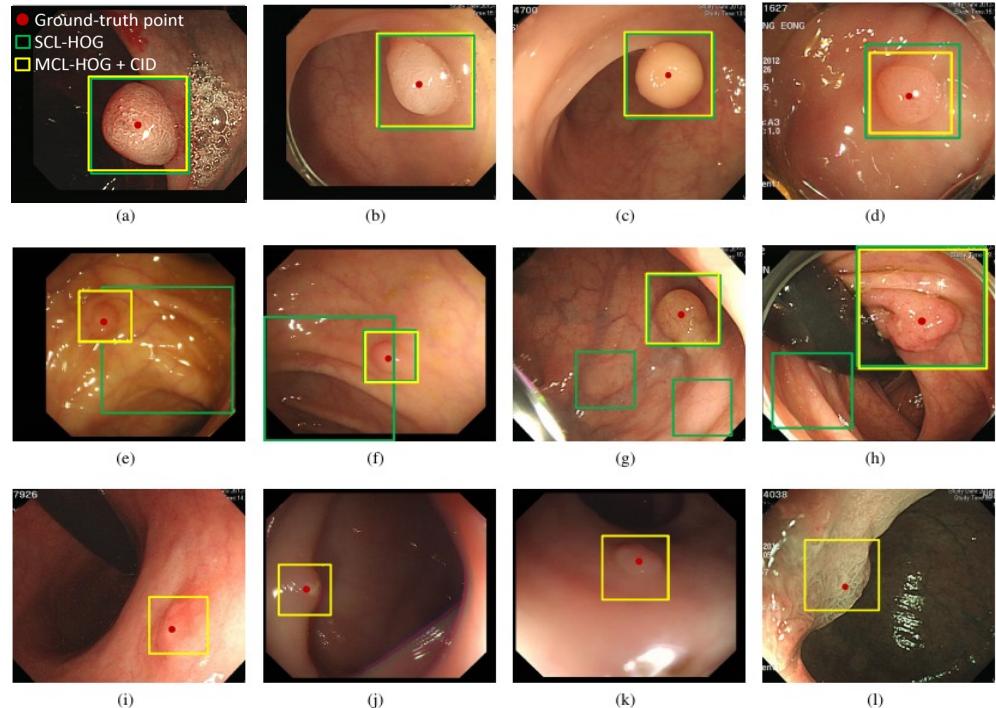
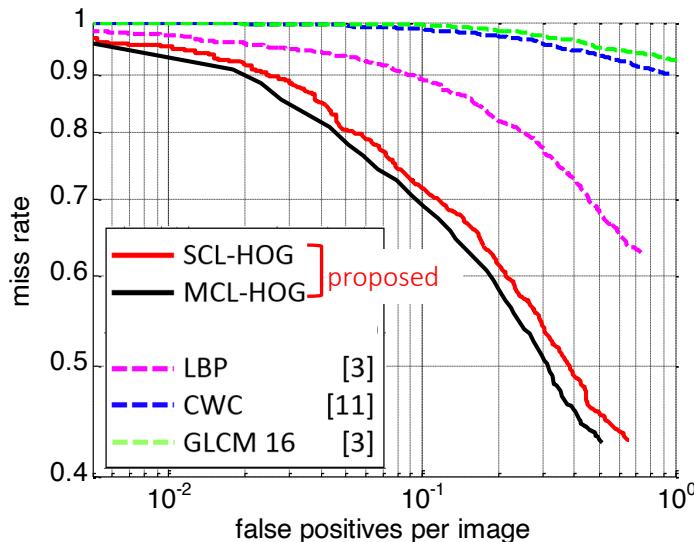
Object Detection/Recognition/Categorization



Object Detection (Ocean-Day)



Polyp Detection Result



Performance comparison between single classifier and multi-classifier

[3] Ameling, Stefan, et al. "Texture-based polyp detection in colonoscopy." Bildverarbeitung für die Medizin 2009. Springer Berlin Heidelberg, 2009. 346-350.

[11] Karkanis, Stavros A., et al. "Computer-aided tumor detection in endoscopic video using color wavelet features." Information Technology in Biomedicine, IEEE Transactions on 7.3 (2003): 141-152.

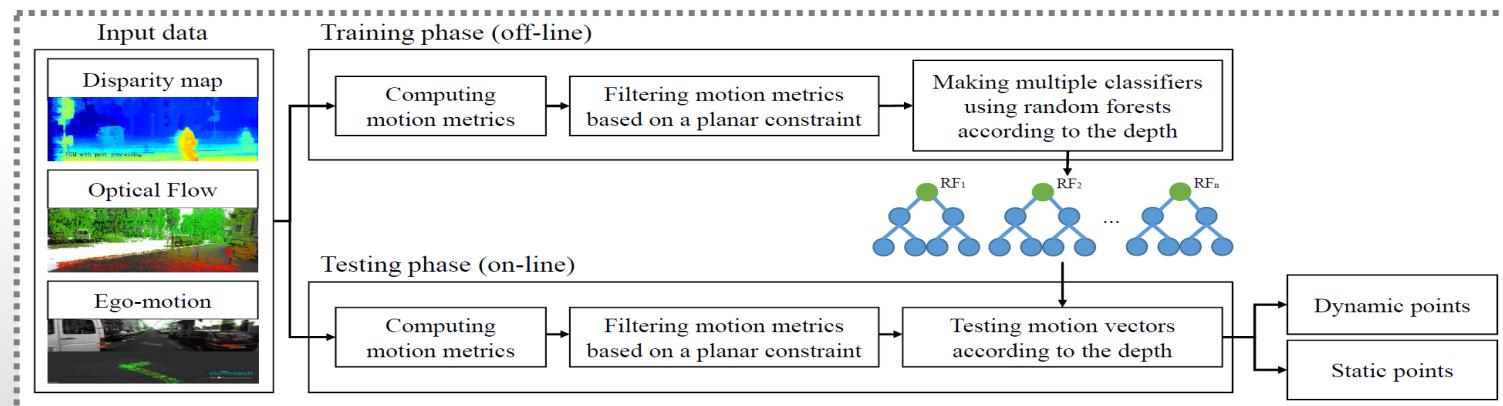
Multi-classifier Framework for Automatic Polyp Detection in Endoscopic Images

MOD (moving object detection) based on Stereo and VO

[Park et al., IV 2017]

Learning to Detect Dynamic Feature Point

- Robust detection using only **two frames** without scene-flow
- The first approach to detect dynamic points using a **learning technique**
 - Non-linear data classification
 - No ambiguous **thresholding** process required
- Noise reduction based on a foreground grouping



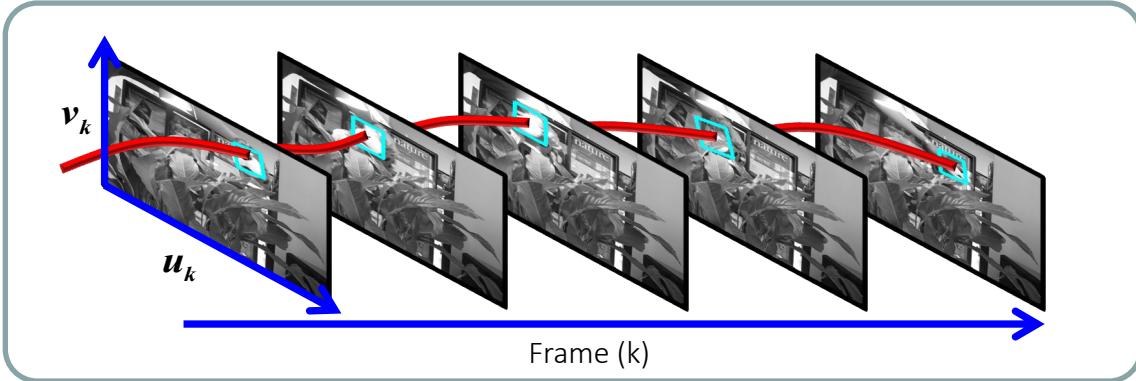
Overall procedures of the proposed dynamic point detection using a learning technique

MOD (moving object detection) based on Stereo and VO

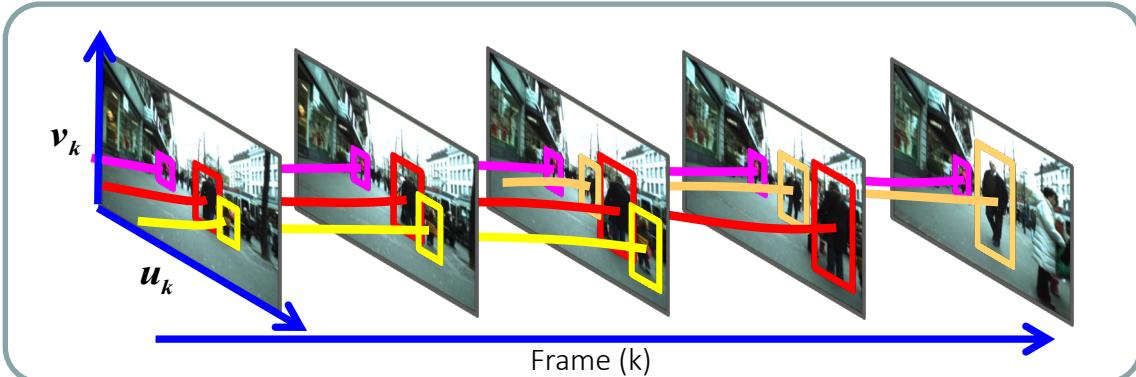
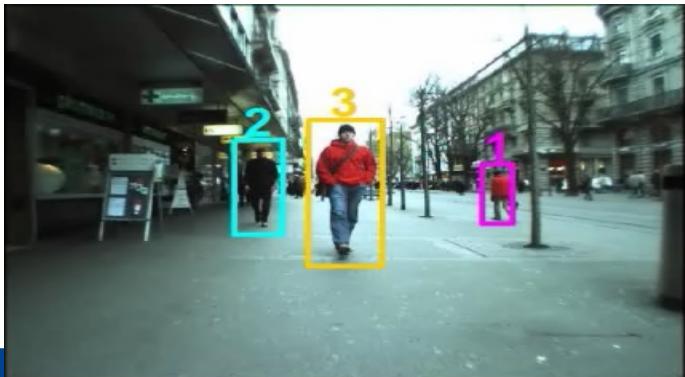
Is it possible to detect dynamic
points from noisy input data
using only **two** frames?

Object Tracking in Computer Vision

- Single Object Tracking
 - Tracking an object under various environmental conditions

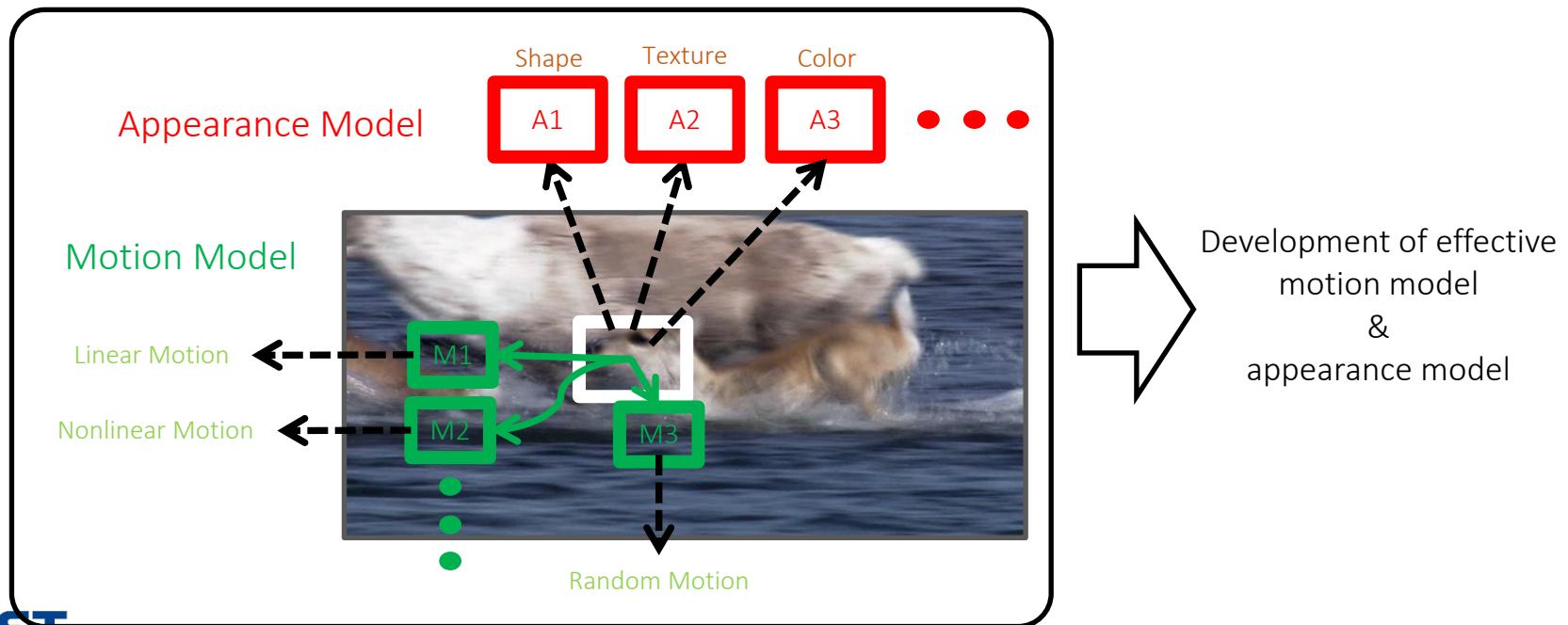


- ▶ Multi-Object Tracking
 - Constructing object trajectories while keeping their identities (label).



Important Cues in Object Tracking

- Appearance model
 - shape, texture, color, intensity,...
- Motion model
 - linear motion, nonlinear motion, random motion

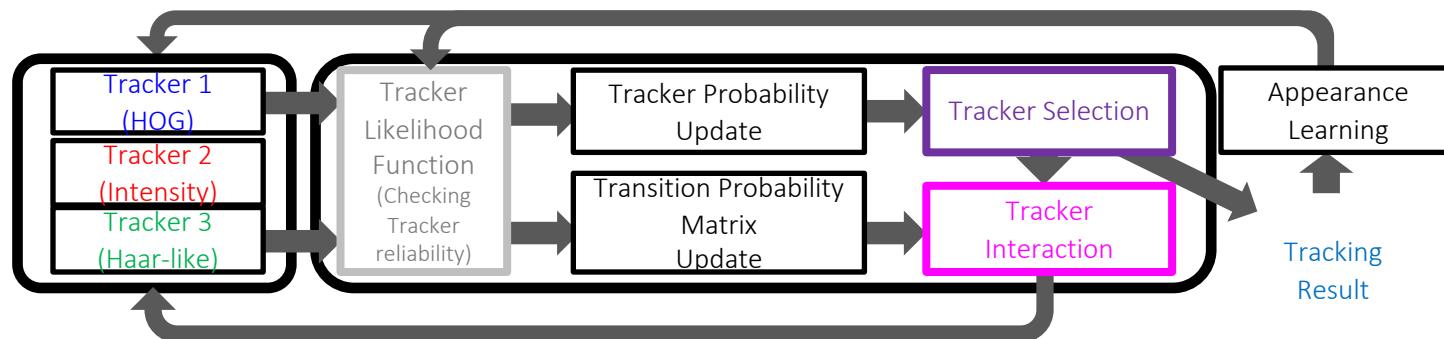


Single Object Tracking with Interaction of Multiple Appearances

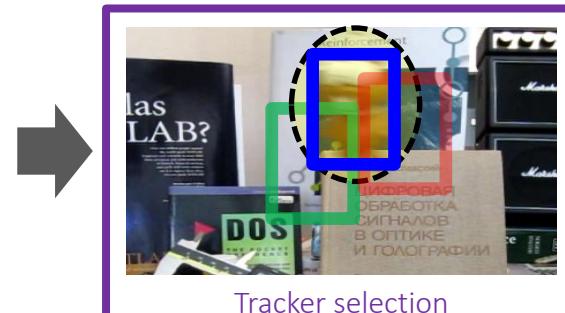
- Interacting Multiveiw Tracker
 - Multiple trackers with multiple features
 - The most reliable tracker is selected
 - Preventing unreliable trackers from drift via tracker interaction

[Yoon et al., PAMI 2016]

[Yoon et al., ECCV 2012]



Tracker Likelihood Function

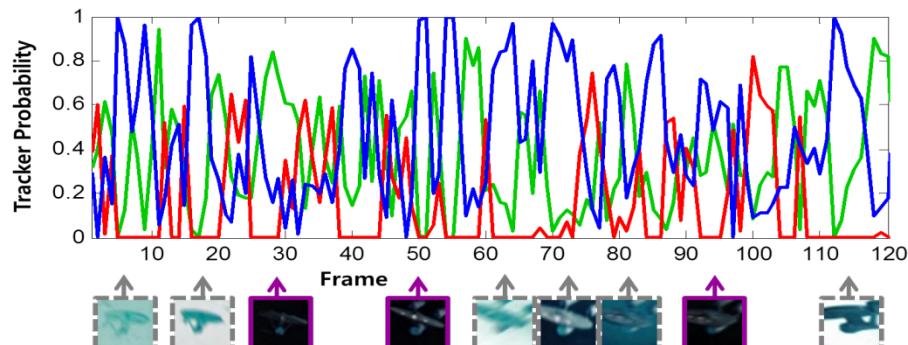
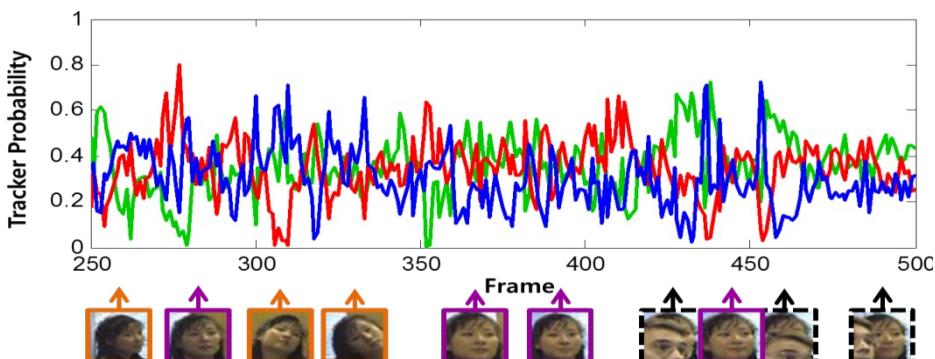
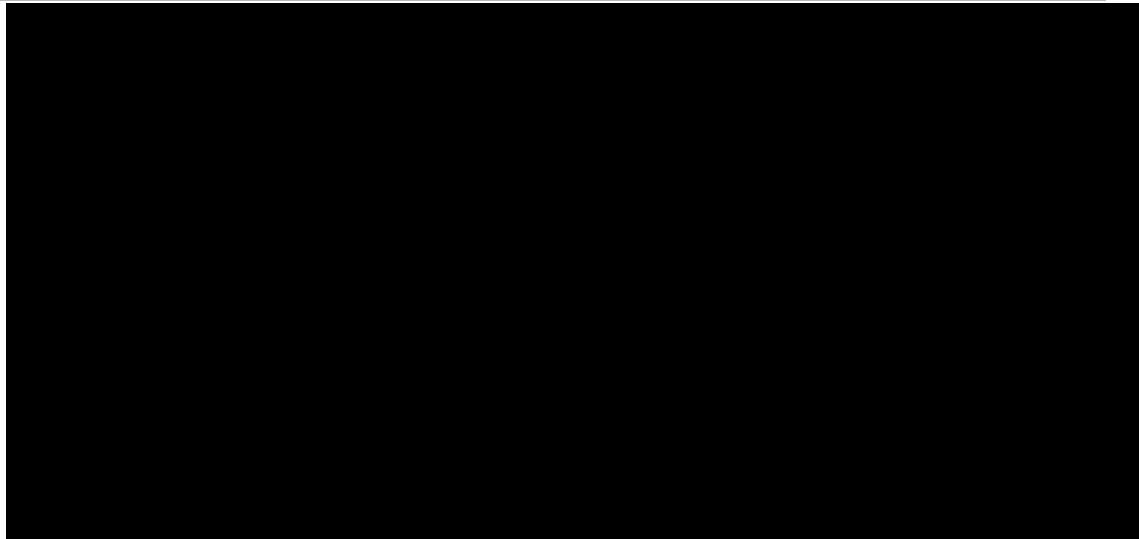


Tracker selection



After tracker interaction

Single Object Tracking with Interaction of Multiple Appearances



moderate variation



motion blur



occlusion



pose variation

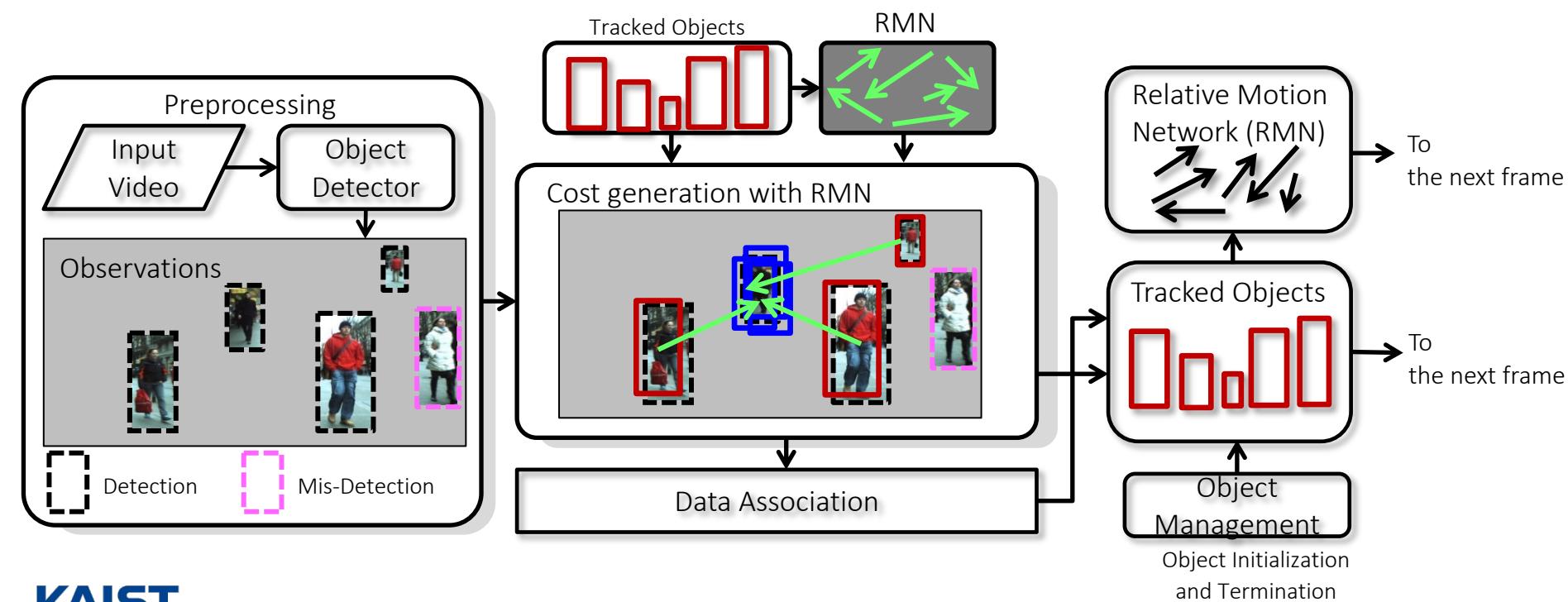


illumination change

Online Multi-object Tracking with Relative Motion

- Online MOT with Relative Motion Network
 - RMN (Relative Motion Network)
 - Data association with RMN
 - Robust relative motion update without a detection

[Yoon et al., IJCV 2019]
[Yoon et al., CVPR 2016]
[Yoon et al., WACV 2015]



Online Multi-object Tracking with Relative Motion

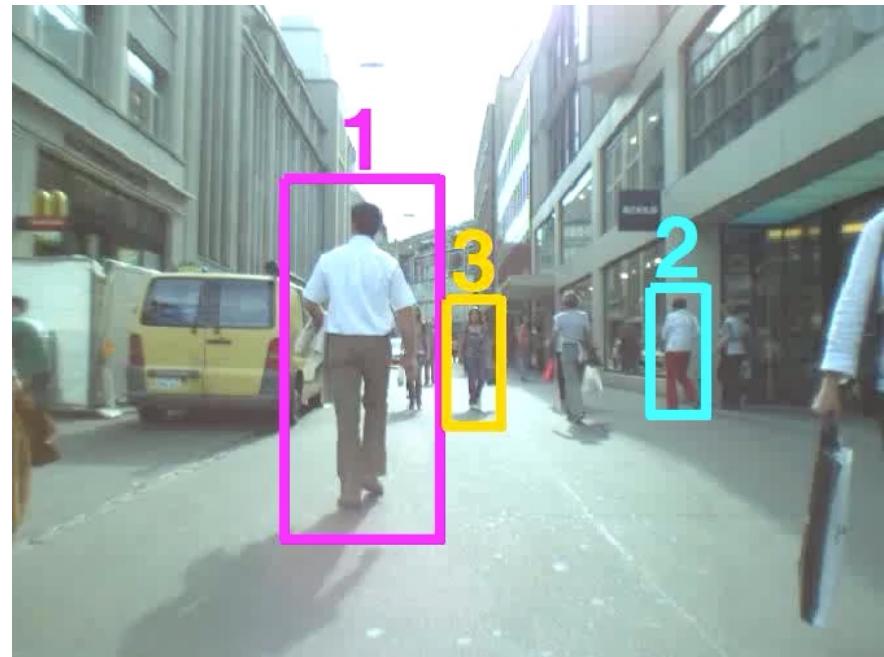
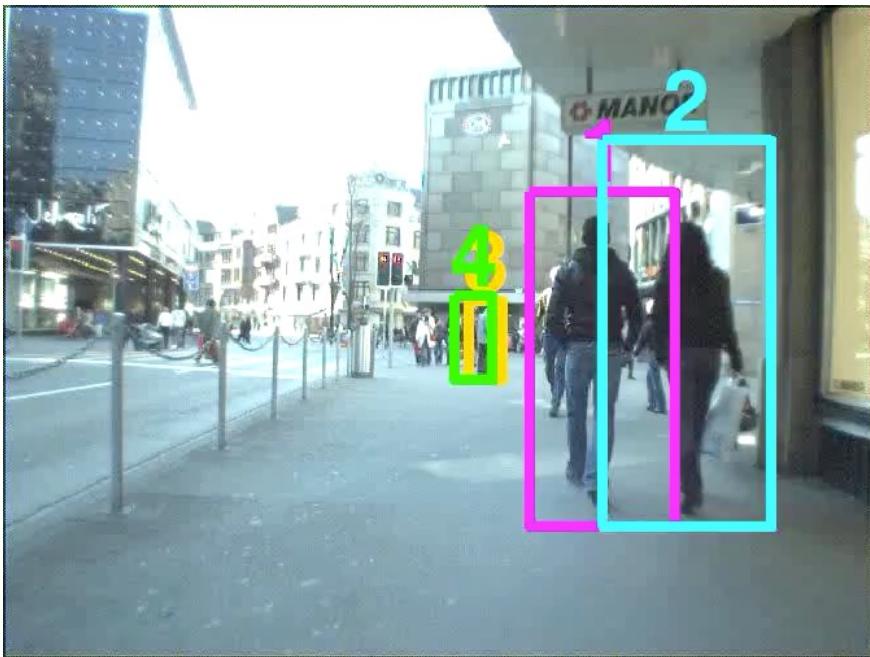


| (1) Marathon1 + Marathon2 + Bahnhof + Sunnyday + Jelmoli | | | | | | |
|--|--------------------|--------------------|------------------|--------------------|----------------------|---------------------|
| | MOTA(\uparrow) | MOTP(\uparrow) | MT(\uparrow) | ML(\downarrow) | Frag(\downarrow) | IDS(\downarrow) |
| RMOT | 76.0% | 78.9% | 58.6% | 16.0% | 39 | 13 |
| RMOT-A | 74.5% | 78.4% | 56.9% | 16.6% | 64 | 26 |
| SMOT | 69.5% | 78.7% | 51.9% | 17.1% | 112 | 34 |
| BPF[15] | 54.1% | 76.3% | 27.6% | 20.9% | 127 | 47 |
| GGA[17] | 63.9% | 76.0% | 44.2% | 17.1% | 233 | 57 |

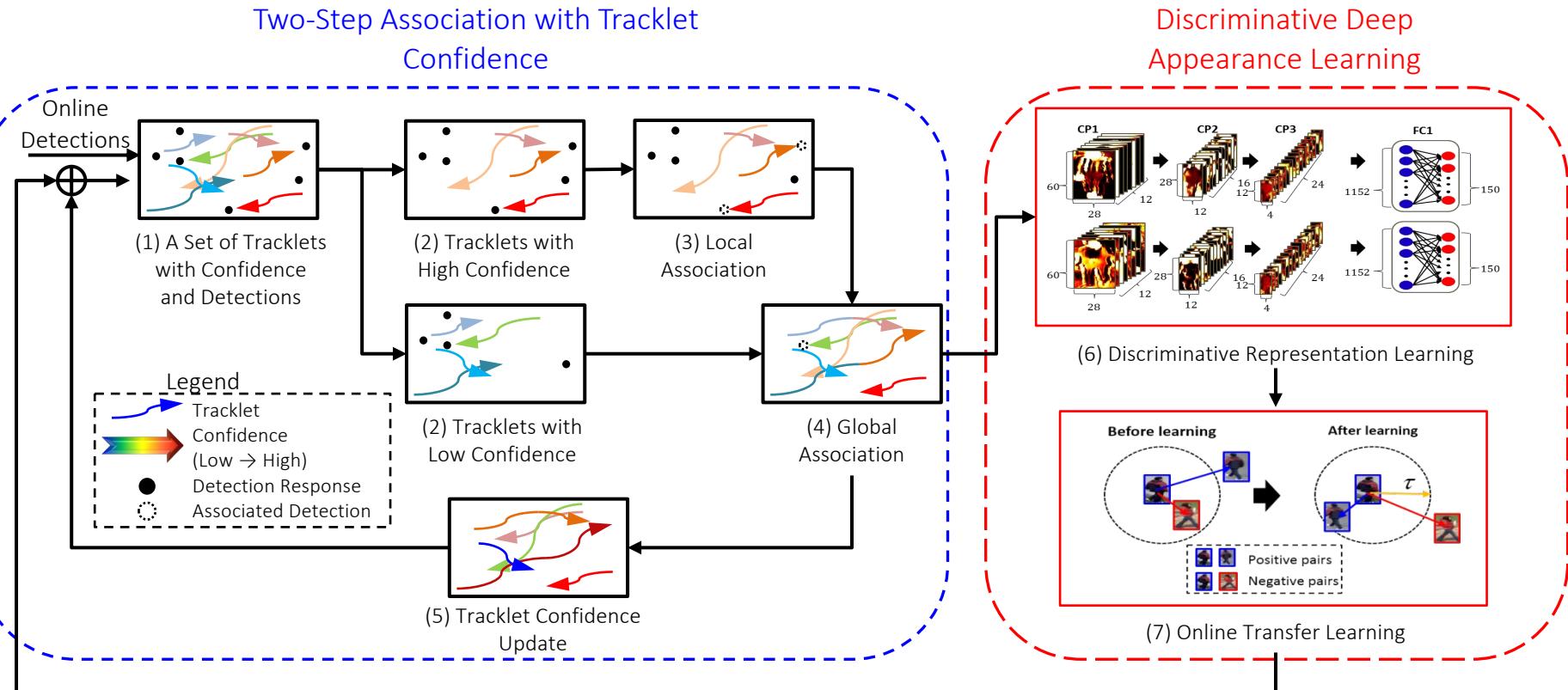
Object Detection & Online Multi-object Tracking with Relative Motion



Online Multi-object Tracking with Camera Motion



Robust Online Multi-Object Tracking based on Tracklet Confidence and Online Discriminative Appearance Learning

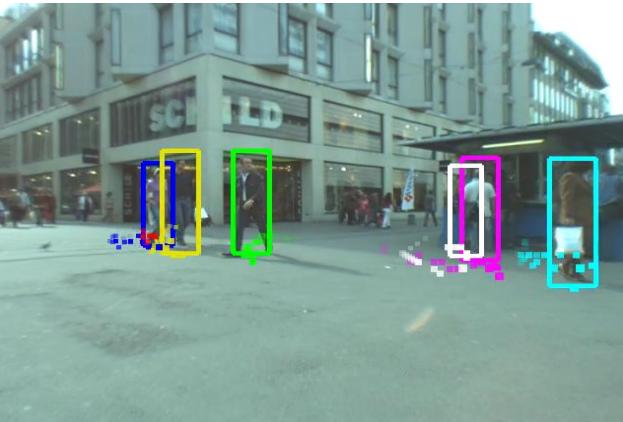


[Bae and Yoon, PAMI 2018]
[Bae and Yoon, CVPR 2014]
[Bae and Yoon, TIP 2014]

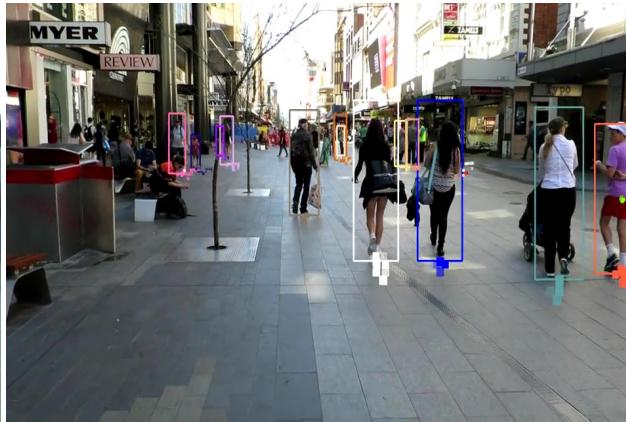
Results for 2015/2016 MOT Benchmark Challenge



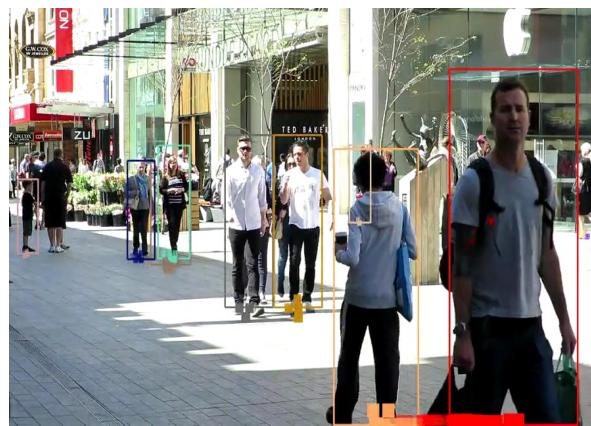
MOT16-03



MOT16-06



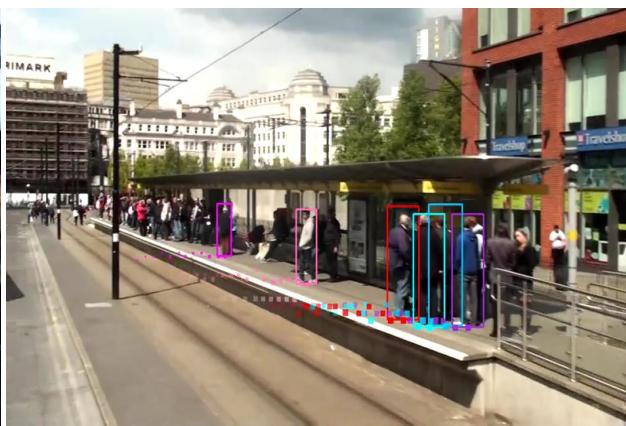
MOT16-07



MOT16-08

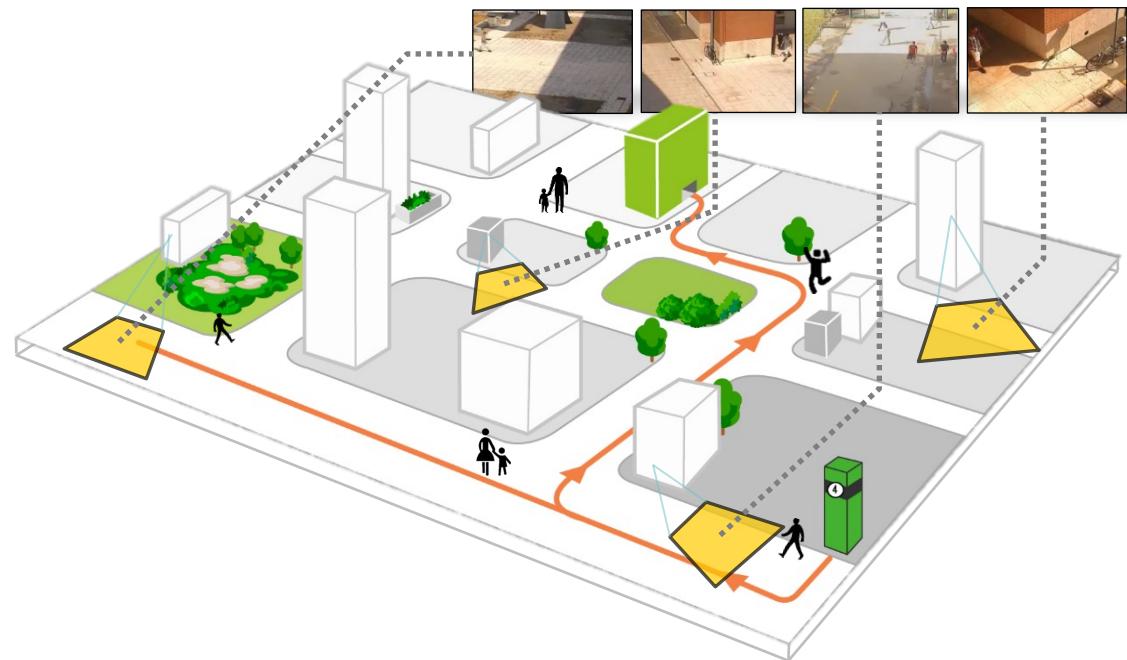
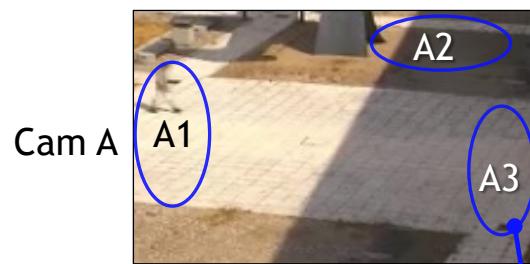


MOT16-12



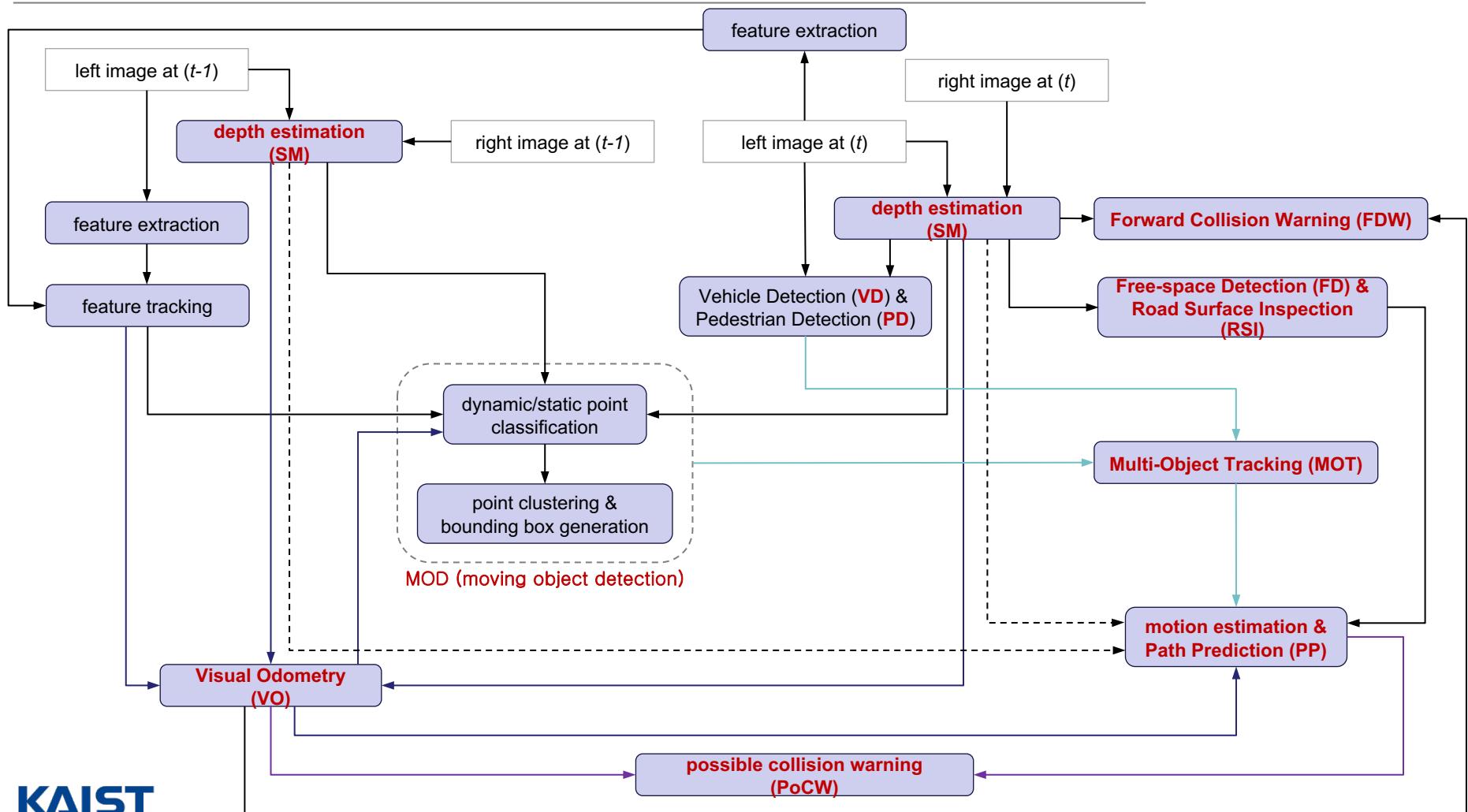
MOT16-14

Multiple People Tracking and Re-identification in Camera Networks



Integration

Stereo-camera-based ADAS

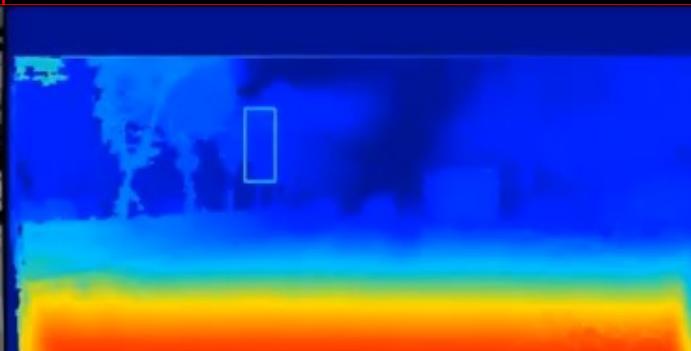


MOD (Moving Object Detection), MOT(Multi-Object Tracking), and PP (Path Prediction)

Tracking Bounding Box & Vehicle Distance



Detection Bounding Box & Depth Map



Feature Extraction & Clustering Result

Feature Tracking Result & Vehicle Motion

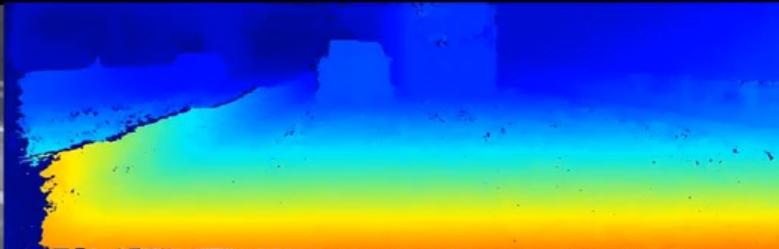
Top View

MOD (Moving Object Detection), MOT(Multi-Object Tracking), and PP (Path Prediction)

Tracking Bounding Box & Vehicle Distance



Detection Bounding Box & Depth Map



Feature Extraction & Clustering Result



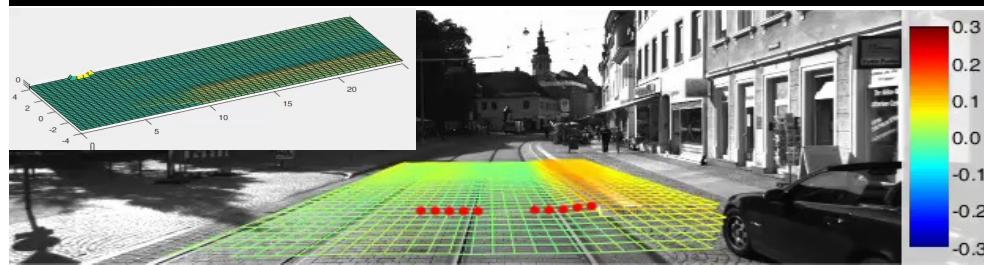
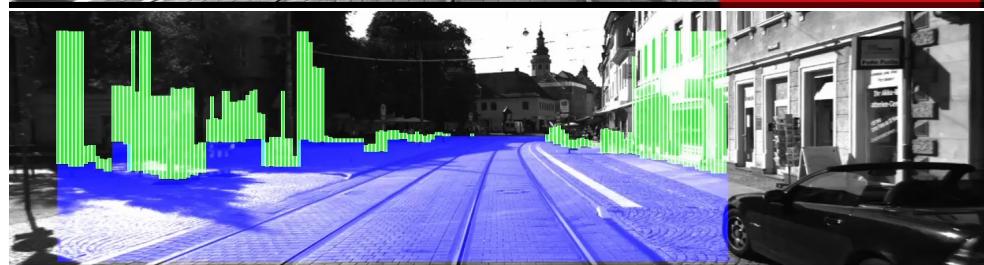
Feature Tracking Result & Vehicle Motion



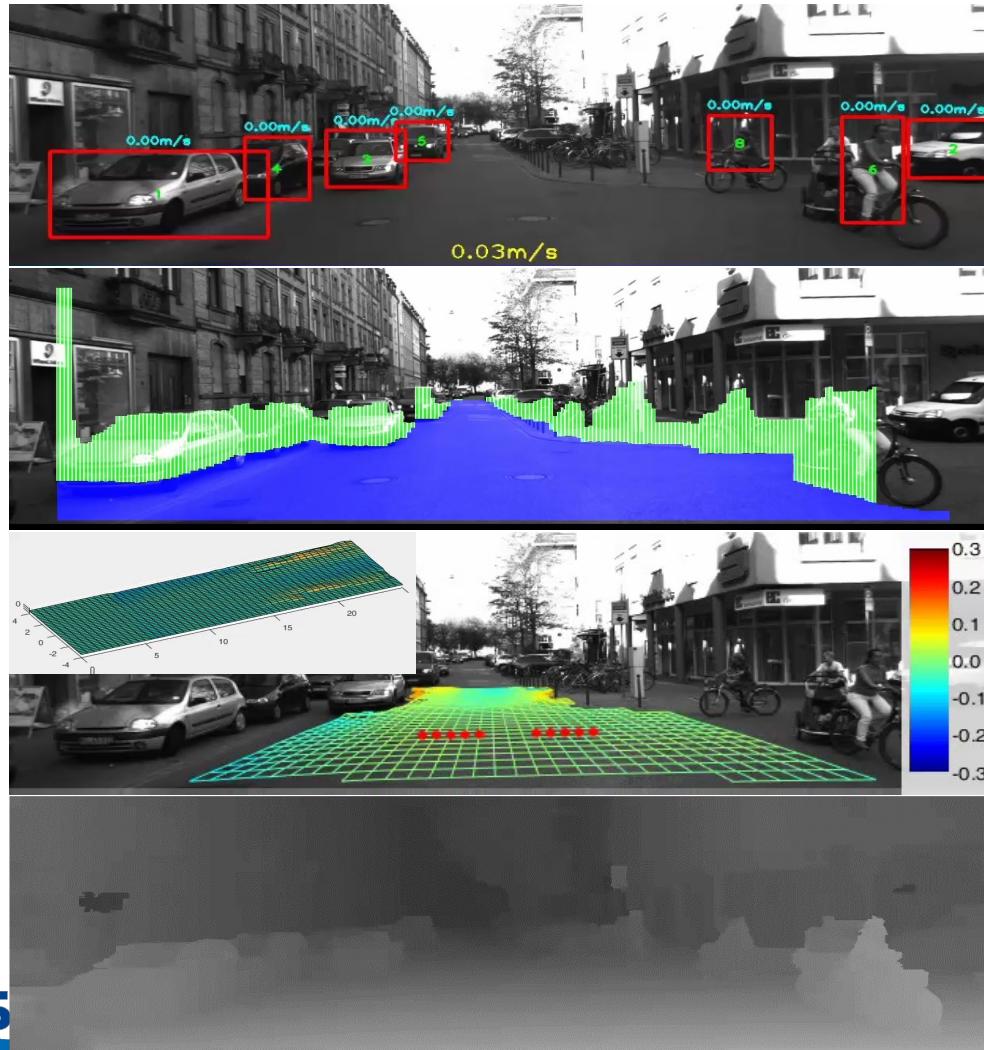
Top View



Integrated System (SM+FD+RSI+MOD+OD+VD+PD+RPE+FCW+MOT+VO+PP+PoCW)



Integrated System (SM+FD+RSI+MOD+OD+VD+PD+RPE+FCW+MOT+VO+PP+PoCW)



Automotive Safety

▶ manufacturer products

consumer products ◀ ▶

Our Vision. Your Safety.

The diagram shows a top-down view of a car. Three yellow lines extend from the front of the car, labeled 'forward looking camera'. One line extends from the back, labeled 'rear looking camera'. One line extends from the side, labeled 'side looking camera'.

› EyeQ Vision on a Chip

A close-up photograph of a square electronic chip with a green PCB underneath. The chip has 'mobileye' and 'EyeQ' printed on it.

> read more

› Vision Applications

An illustration of a person walking across a crosswalk. A yellow rectangular frame highlights the person, indicating the field of view of a camera.

Road, Vehicle, Pedestrian Protection and more

> read more

› AWS Advance Warning System

A circular display screen showing a small car icon and the number '0.8'.

> read more

News

- › Mobileye Advanced Technologies Power Volvo Cars World First Collision Warning With Auto Brake System
- › Volvo: New Collision Warning with Auto Brake Helps Prevent Rear-end

> all news

A photograph of a car's dashboard with a digital display showing a small car icon and the number '0.8'.

Events

- › Mobileye at Equip Auto, Paris, France
- › Mobileye at SEMA, Las Vegas, NV

> read more

- Mobileye: Vision systems in high-end BMW, GM, Volvo models

Waymo - Google Project for Autonomous Driving

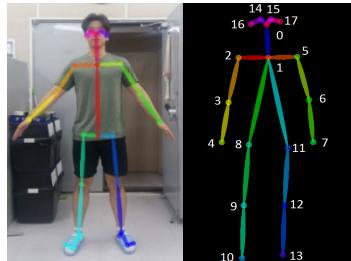


Human

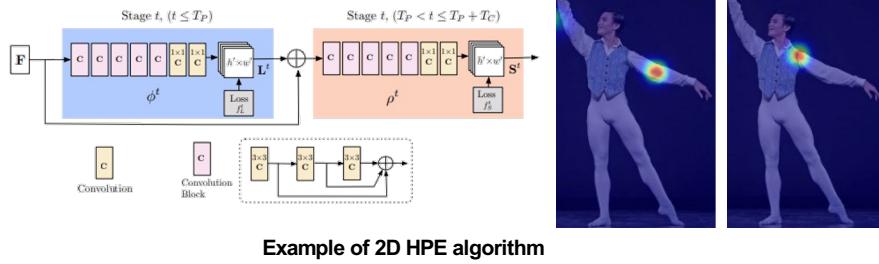
Face Detection



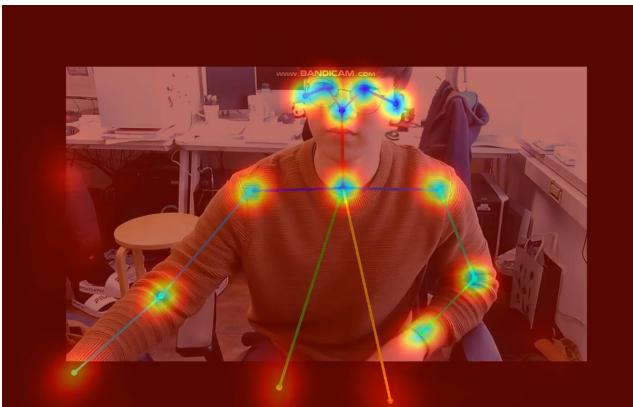
Image-based 2D/3D Human Pose Estimation



Skeleton model



Example of 2D HPE algorithm

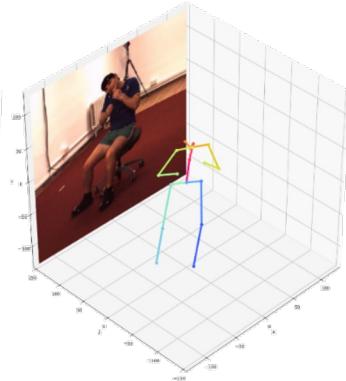
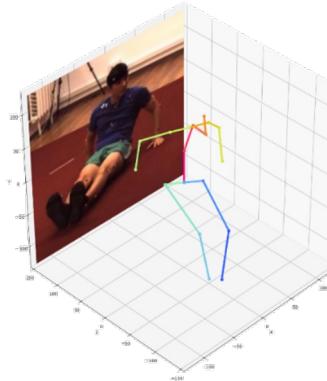
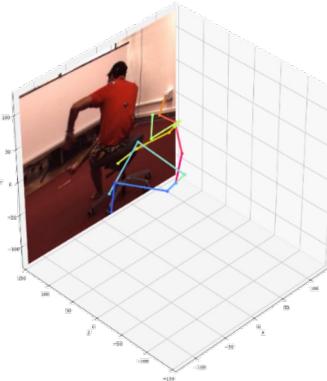
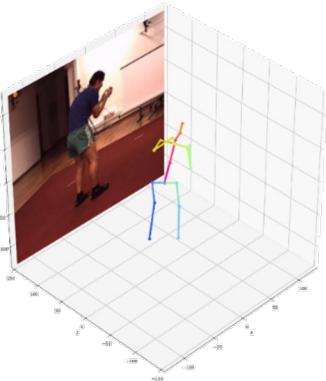
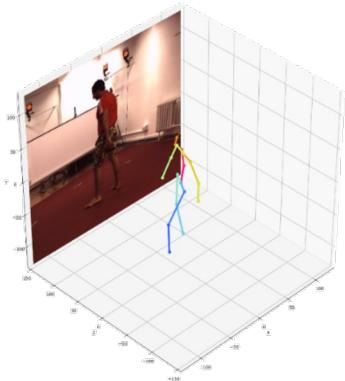


Web demo examples [6]

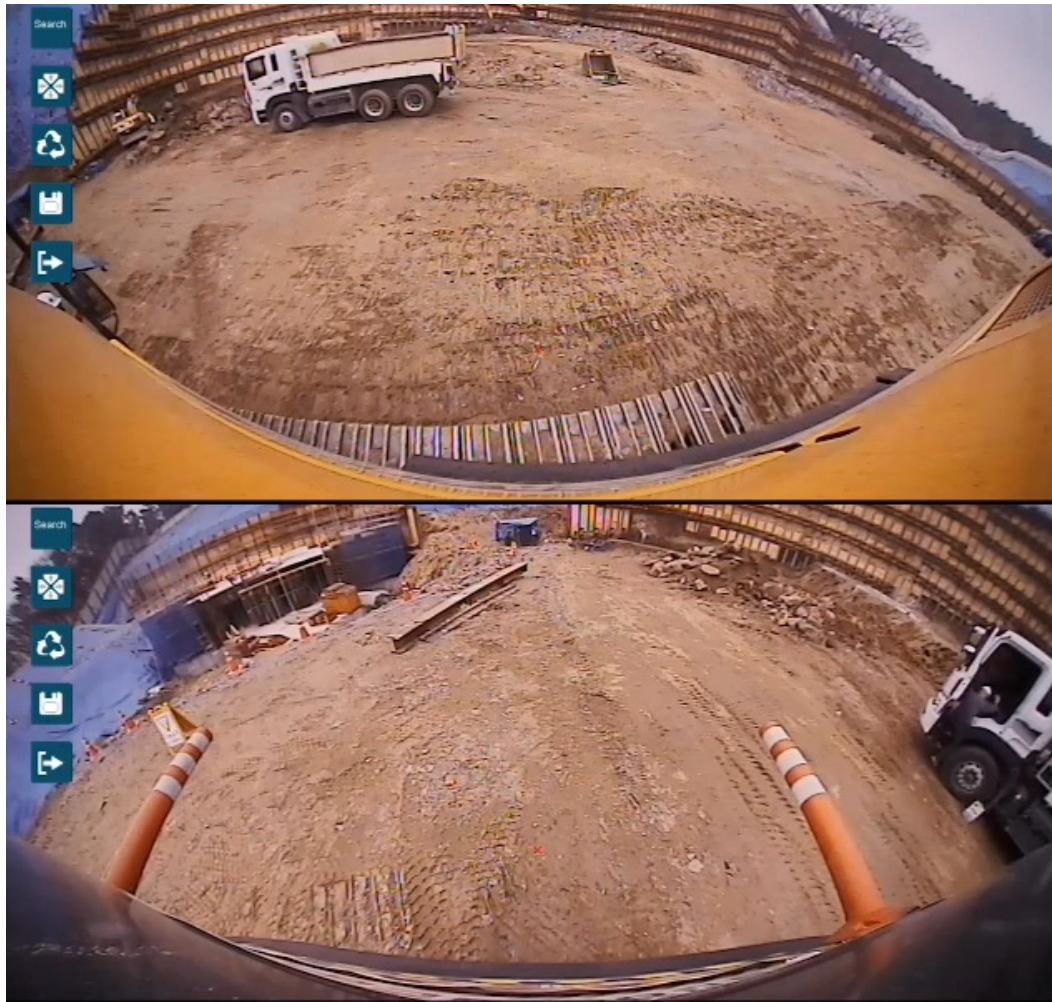


Picture demo examples [6]

Image-based 2D/3D Human Pose Estimation



Smart Scene Monitoring System



Smart Scene Monitoring System



HCI / HRI

Vision-based Interaction (and Games)



Nintendo Wii has camera-based IR tracking built in. See [Lee's work at CMU](#) on clever tricks on using it to create a [multi-touch display](#)!



Sony EyeToy

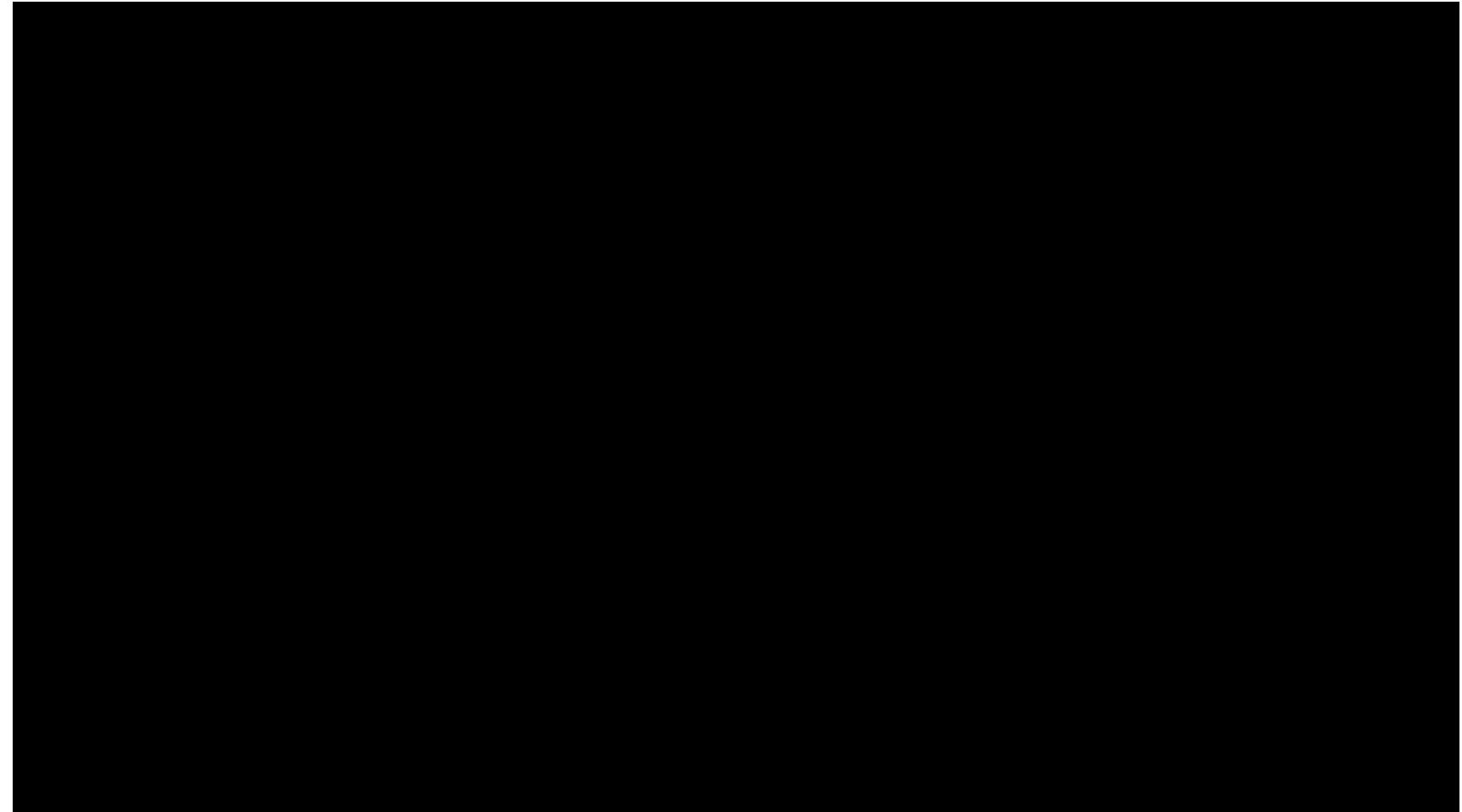


Assistive technologies



Kinect

Vision for Robotics with Depth Sensing

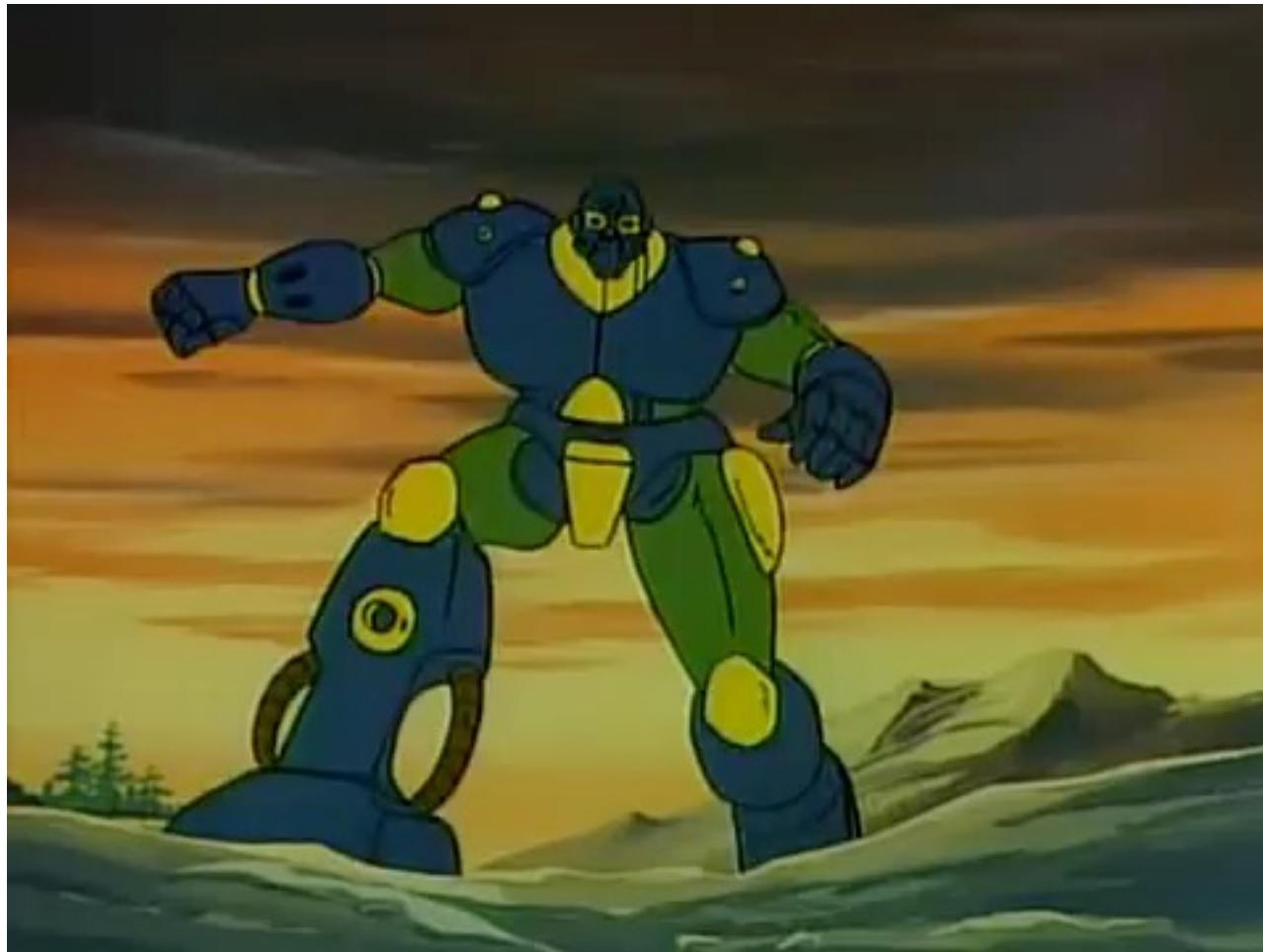


Visual Servoing

- Visual servoing, also known as vision-based robot control, is a technique which uses feedback information extracted from a vision sensor to control the motion of a robot.



Motion Capture for Robot Control



Motion Capture for Robot Control

- Human motion capture with cameras



Motion Capture for Robot Control

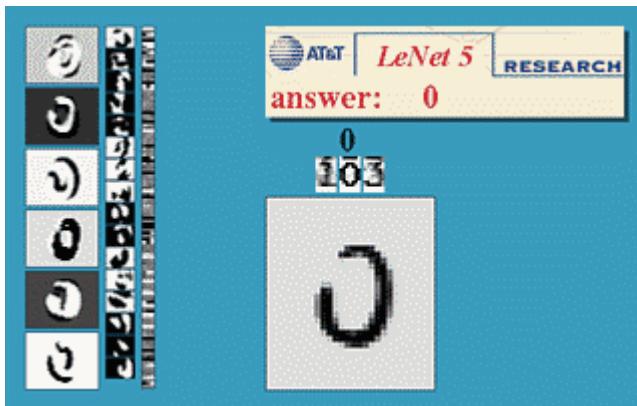
- Hand motion capture



Others

Optical Character Recognition (OCR)

- Technology to convert scanned docs to text
 - If you have a scanner, it probably came with OCR software



Digit recognition, AT&T labs



License plate readers

http://en.wikipedia.org/wiki/Automatic_number_plate_recognition

VQA (Visual Question Answering)

- Given an image and a free-form, natural language question about the image, the machine's task is to automatically produce a concise, accurate, free-form, natural language answer.

Who is wearing glasses?

man



woman



Where is the child sitting?

fridge



arms



Is the umbrella upside down?

yes



no



How many children are in the bed?

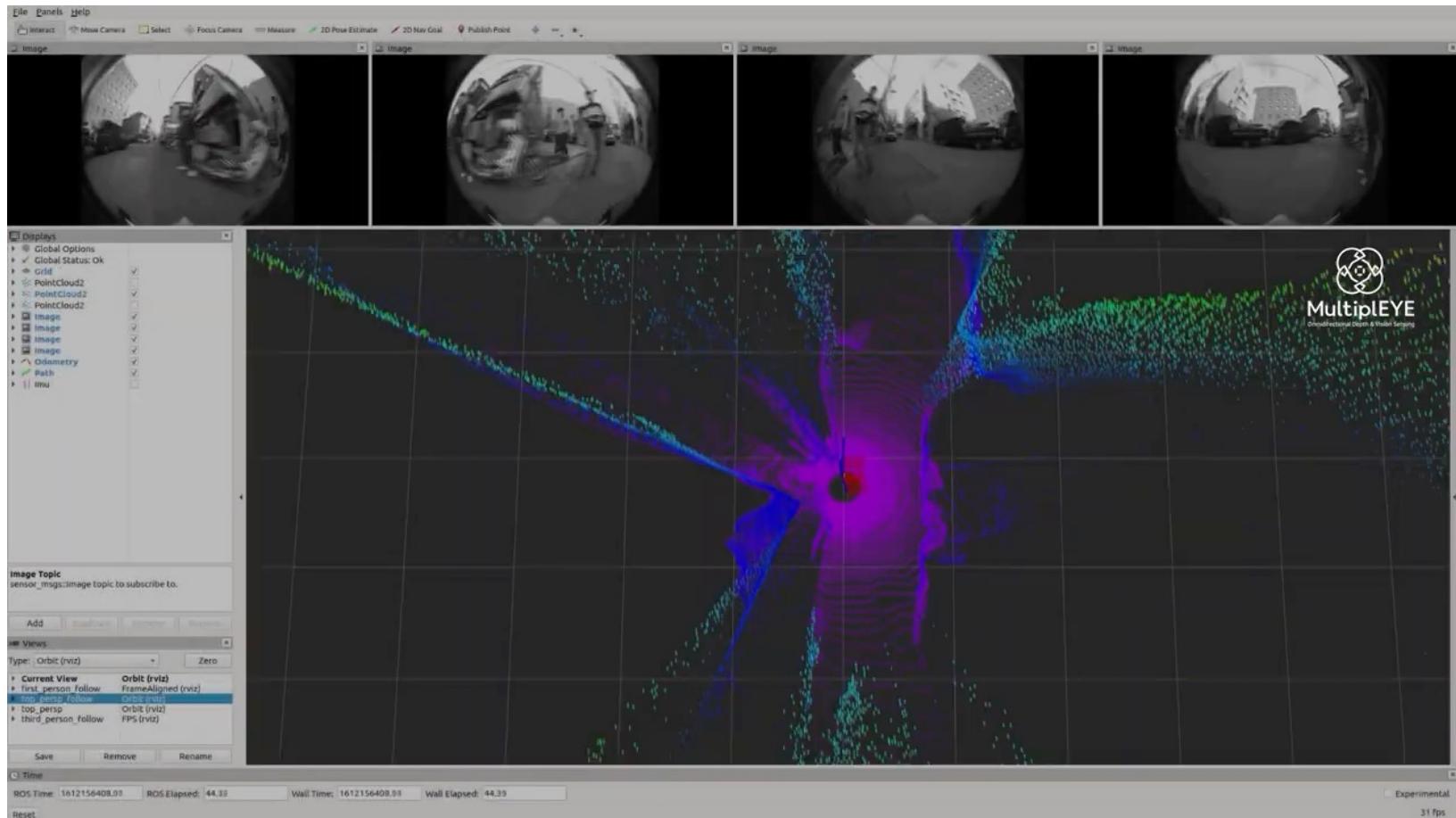
2



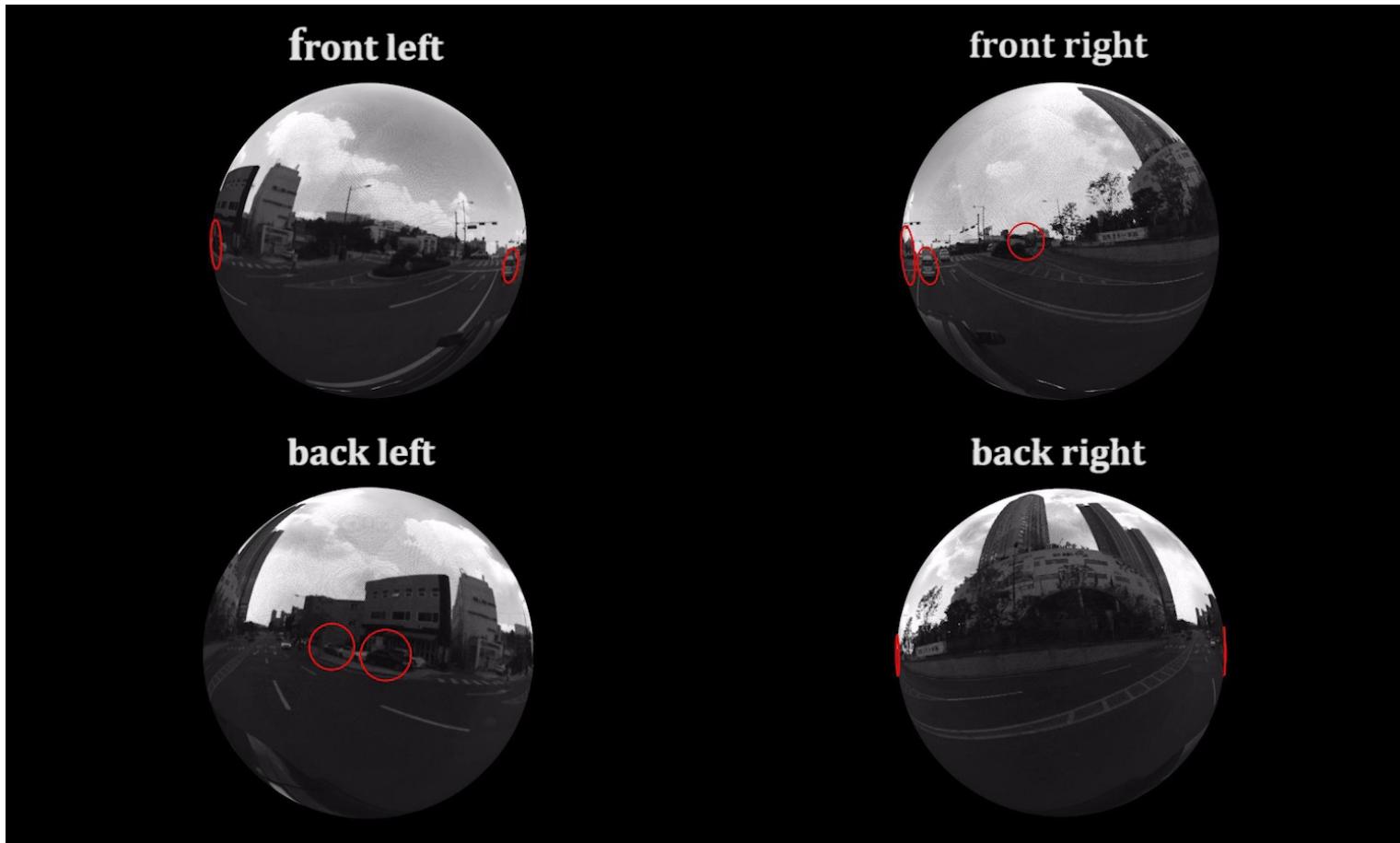
1



Multi-WFoV Camera-based ADAS

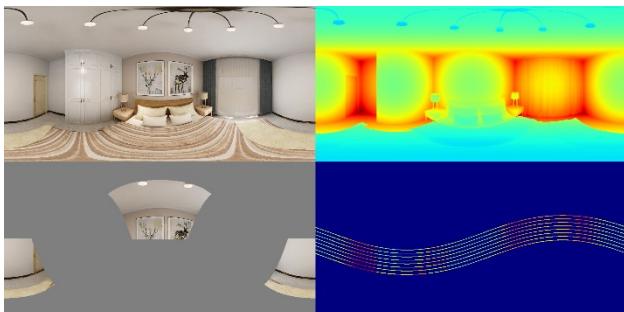


Multi-WFoV Camera-based ADAS

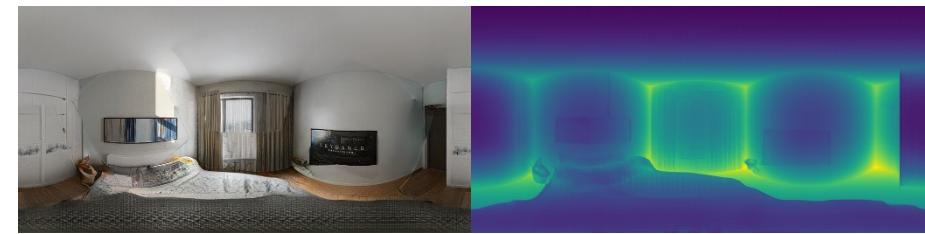
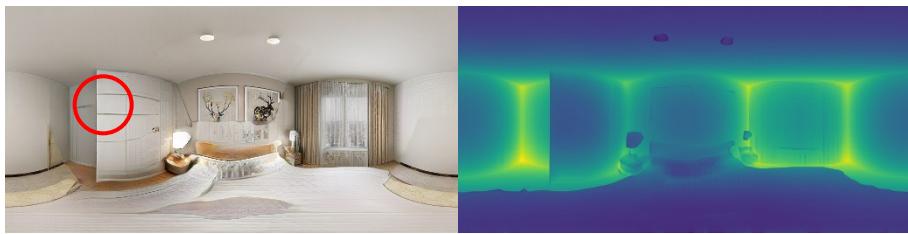
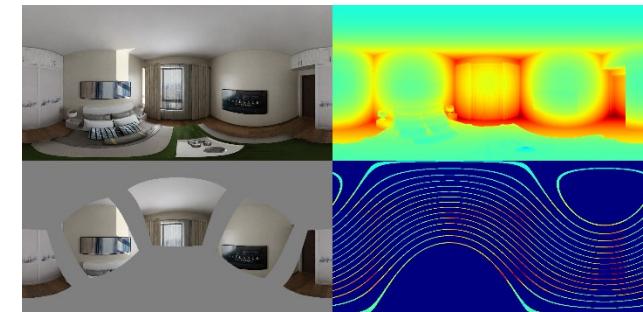


Bi-modal Indoor Panorama Synthesis

GT



Input



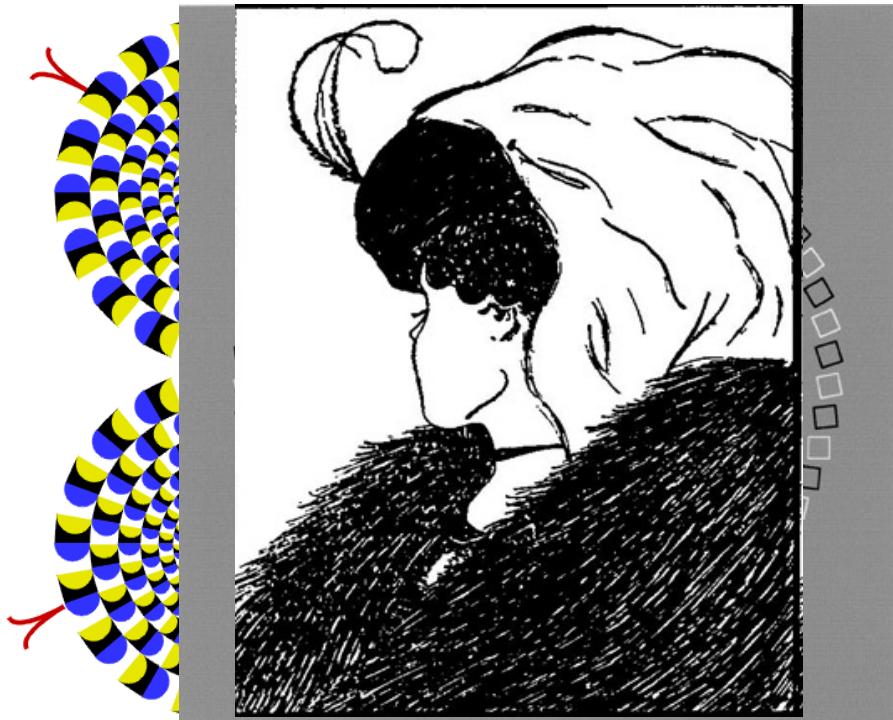
(6) Event-guided Deblurring of Videos



Why is it so hard to achieve human-level
visual intelligence?

Human vs. Computer

- Which one works better?



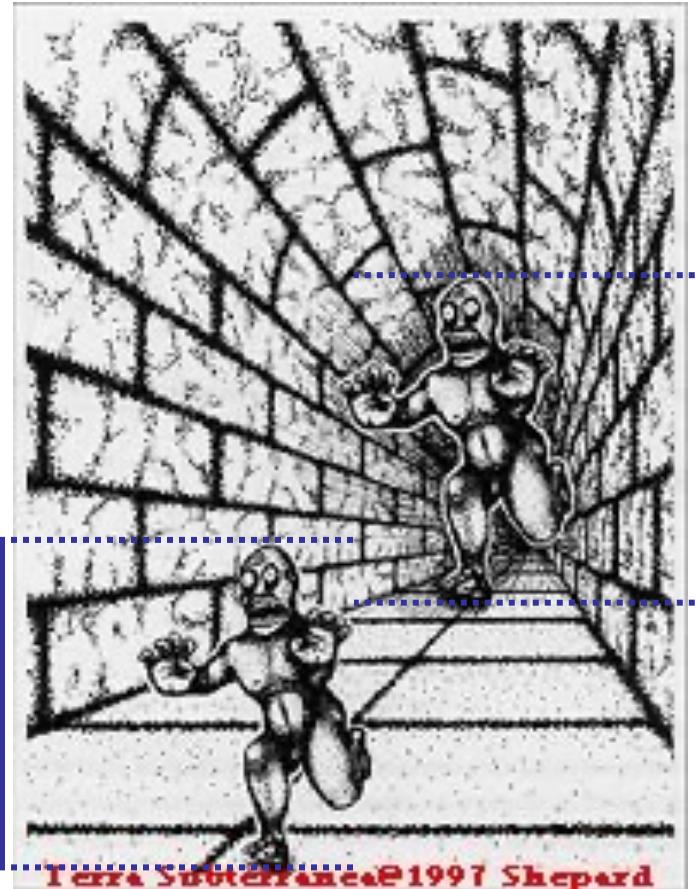
"Rotating snakes 3" : Rings of snakes appear to rotate.
(Copyright A.Kitaoka 2004)

Image: Projection of the World



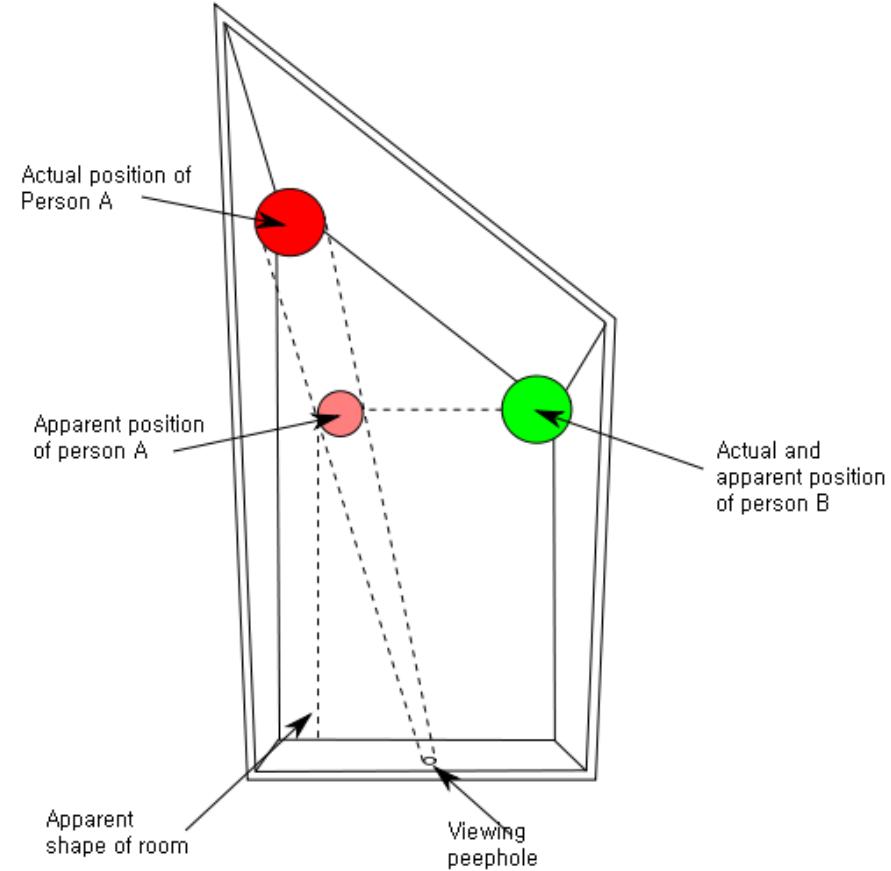
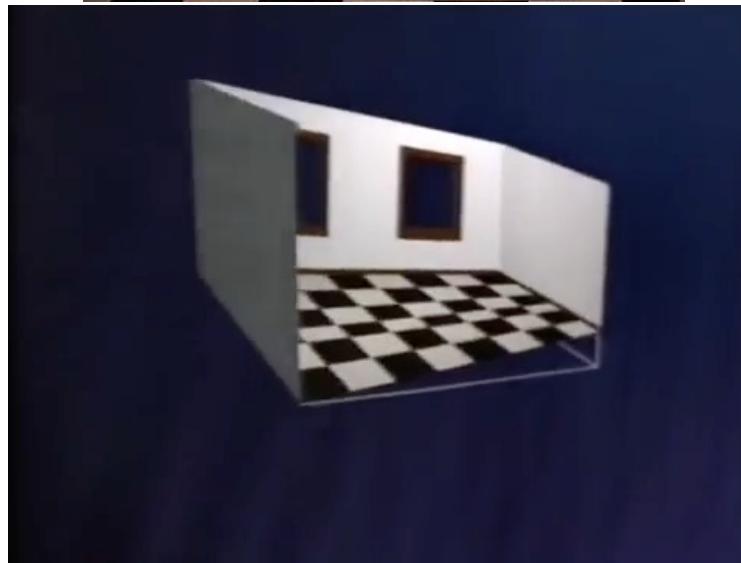
Figures from Seitz

Illusion with a Vanishing Point

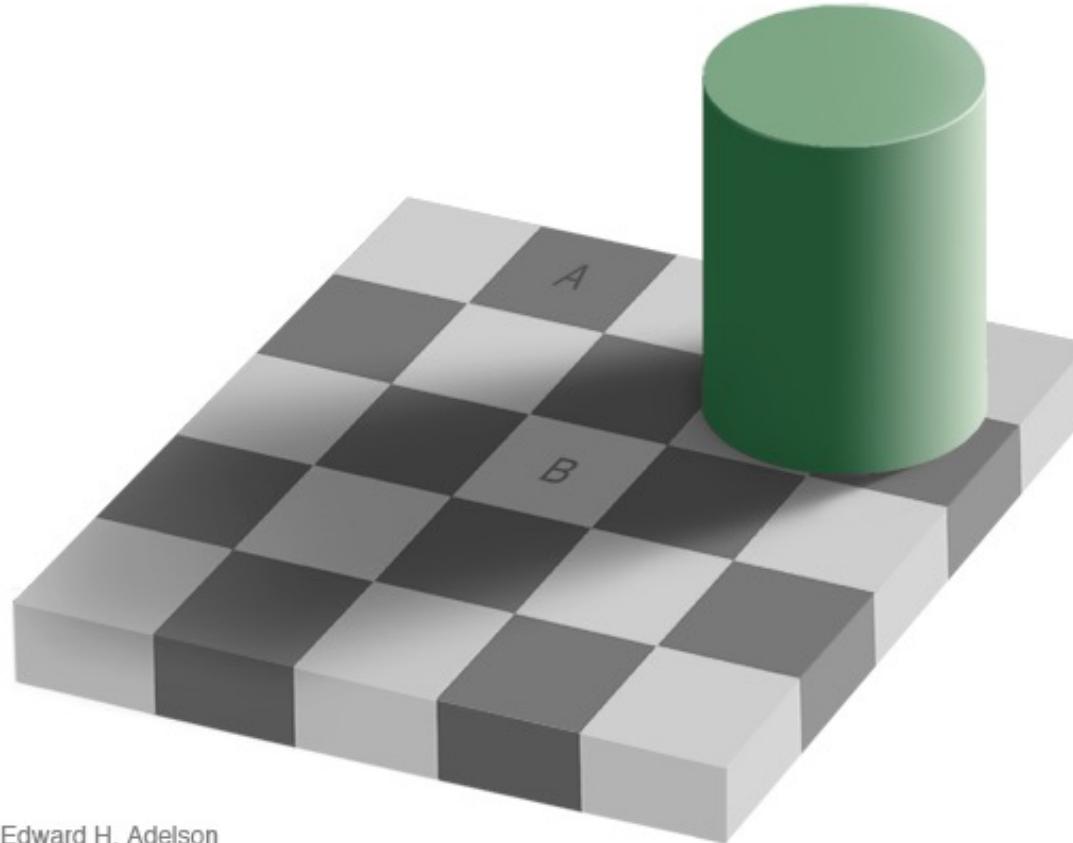


Terra Subterranea © 1997 Shepard

Ames Room



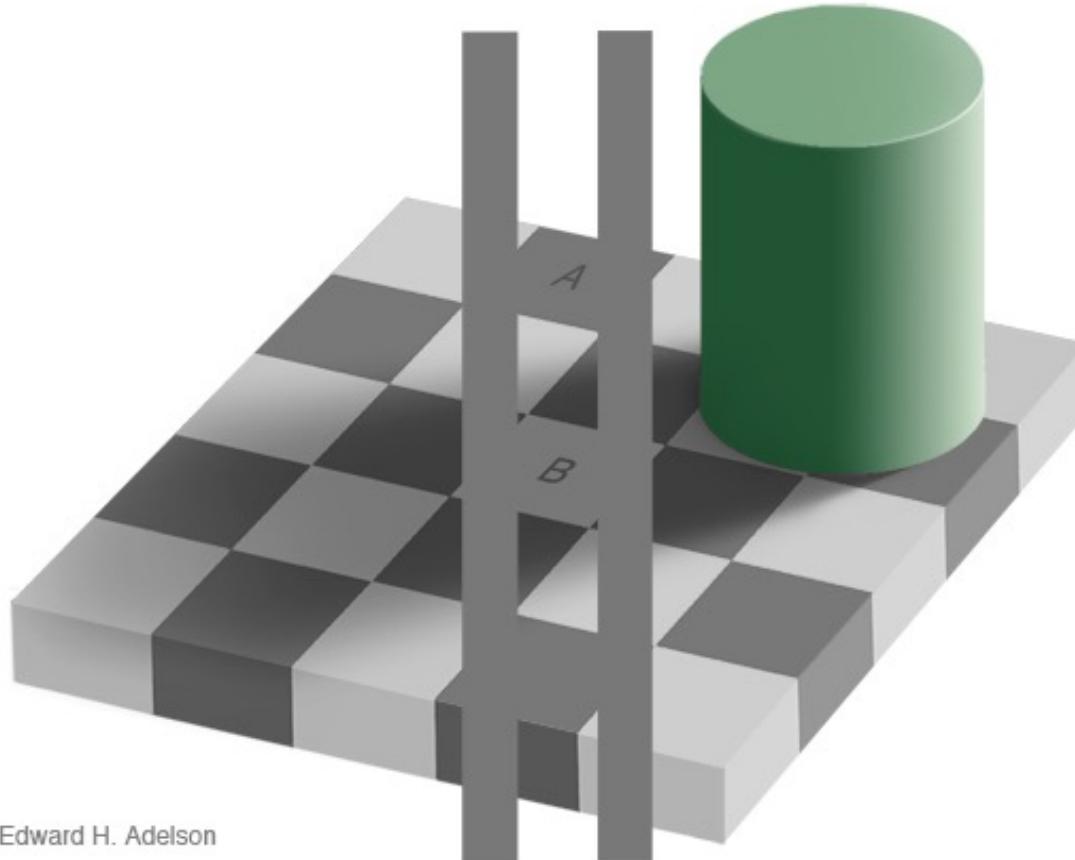
Lightness Constancy



Edward H. Adelson

http://web.mit.edu/persci/people/adelson/checkershadow_illusion.html

Lightness Constancy



Edward H. Adelson

http://web.mit.edu/persci/people/adelson/checkershadow_illusion.html

Color Perception

Green

Blue

Yellow

Purple

Orange

Red

White

Blue

Green

Green

Blue

Yellow

Purple

Orange

Red

White

Blue

Green

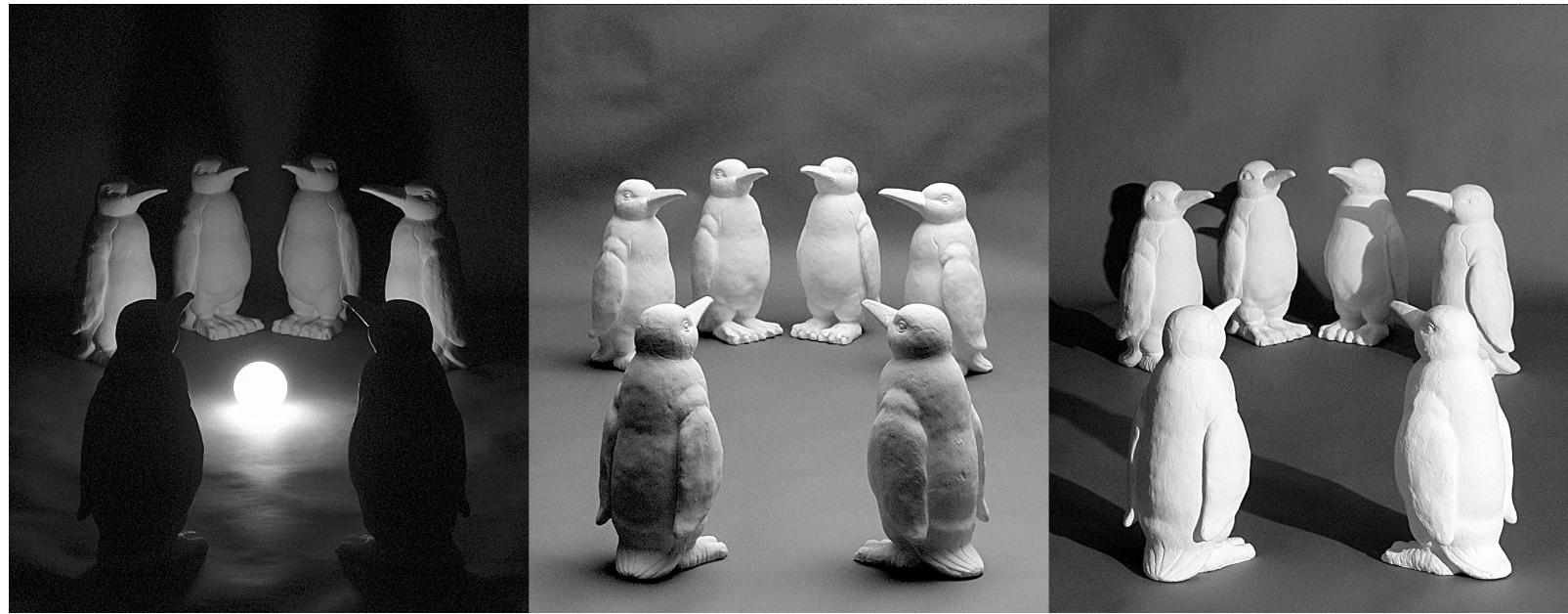
Challenges: Viewpoint Variation



Michelangelo 1475-1564

slide credit: Fei-Fei, Fergus & Torralba

Challenges: Illumination

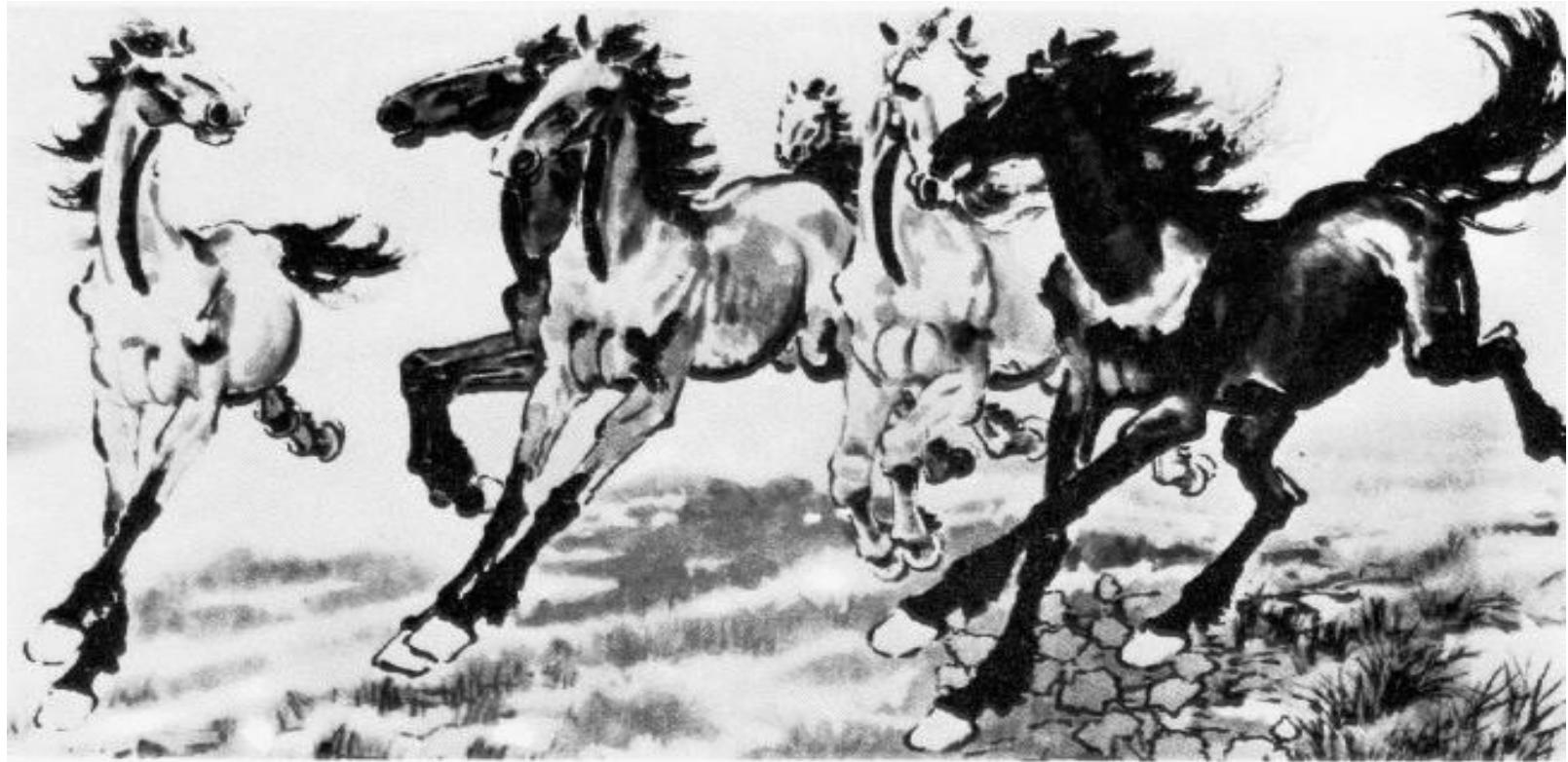


Challenges: Scale



slide credit: Fei-Fei, Fergus & Torralba

Challenges: Deformation



Xu, Beihong 1943

Challenges: Occlusion



Magritte, 1957

slide credit: Fei-Fei, Fergus & Torralba

Challenges: Background Clutter



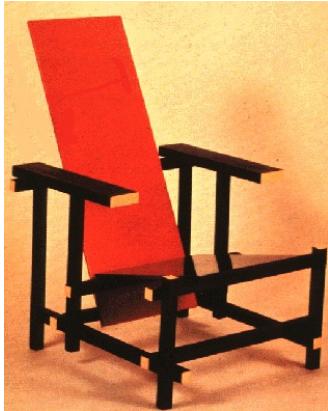
Emperor shrimp and commensal crab on a sea cucumber in Fiji
Photograph by Tim Laman

NATIONAL
GEOGRAPHIC

Challenges: Motion

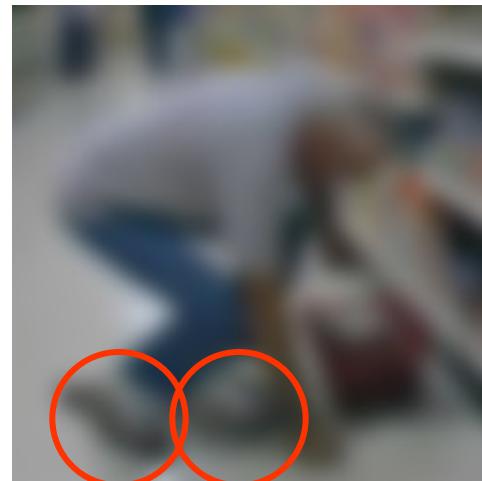
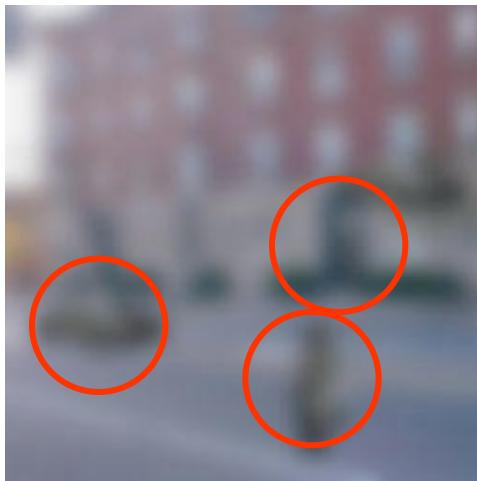
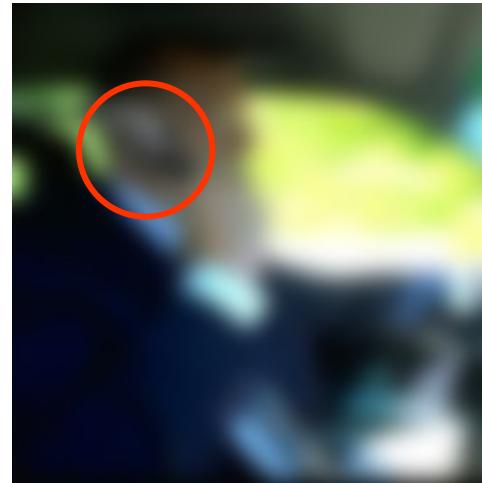
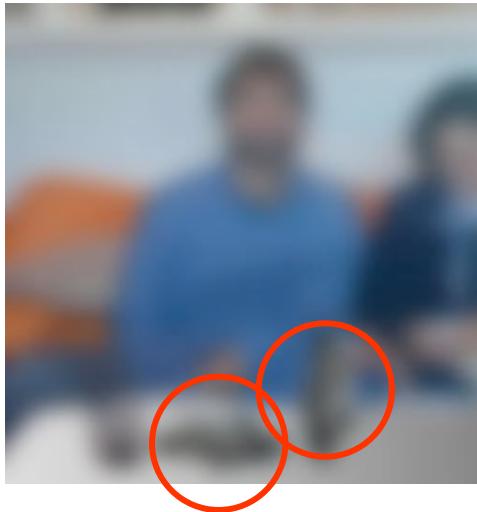


Challenges: Object Intra-class Variation



slide credit: Fei-Fei, Fergus & Torralba

Challenges: Local Ambiguity



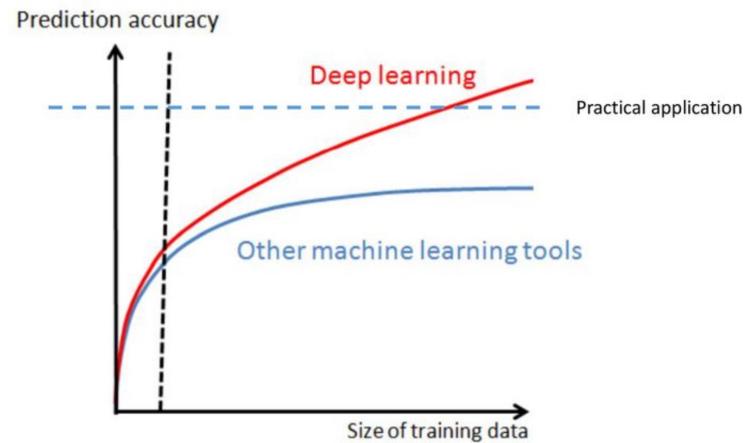
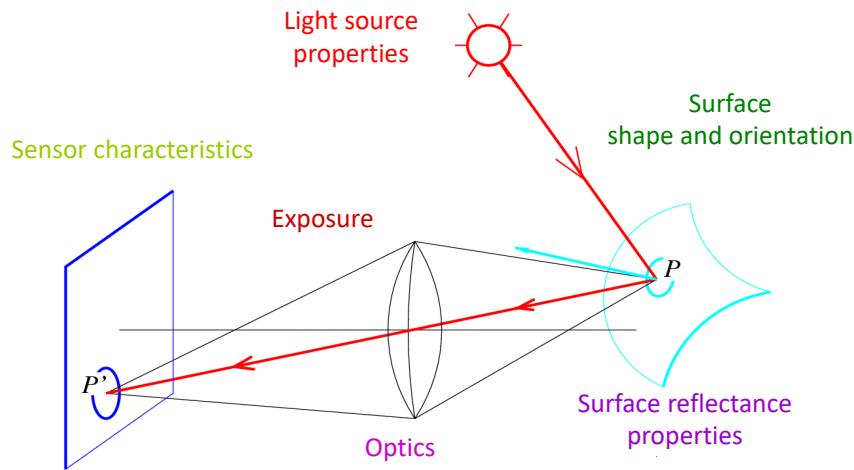
Bottom Line

- Perception is an inherently ambiguous problem
 - Many different 3D scenes could have given rise to a particular 2D picture



- Possible solutions
 - Bring in more constraints (more images)
 - Use prior knowledge about the structure of the world
 - Need a combination of different methods
 - **Learning-based approaches!!**

Conventional vs. Deep-learning-based Approaches in VI



- Conventional
 - Model each process (using math.)
 - Inversely estimate each property (3d shape, light, reflectance properties) based on the specific model
 - No need to collect training images

- Deep-learning based
 - Do not explicitly model the process
 - Need to prepare training images
 - Need to be generalized for the unseen data

Conventional vs. Deep-learning-based Approaches in VI

- Performance comparison
 - Classifications tasks (detection, stereo matching, semantic segmentation, etc.)
 - DL-based approaches \geq conventional approaches
 - Regression tasks (visual odometry, camera calibration, etc)
 - conventional approaches \geq DL-based approaches

Conventional vs. Deep-learning-based Approaches in VI

- *Why do we still need to learn conventional approaches and models?*
 - From conventional approaches, we can gather the **domain knowledge**.

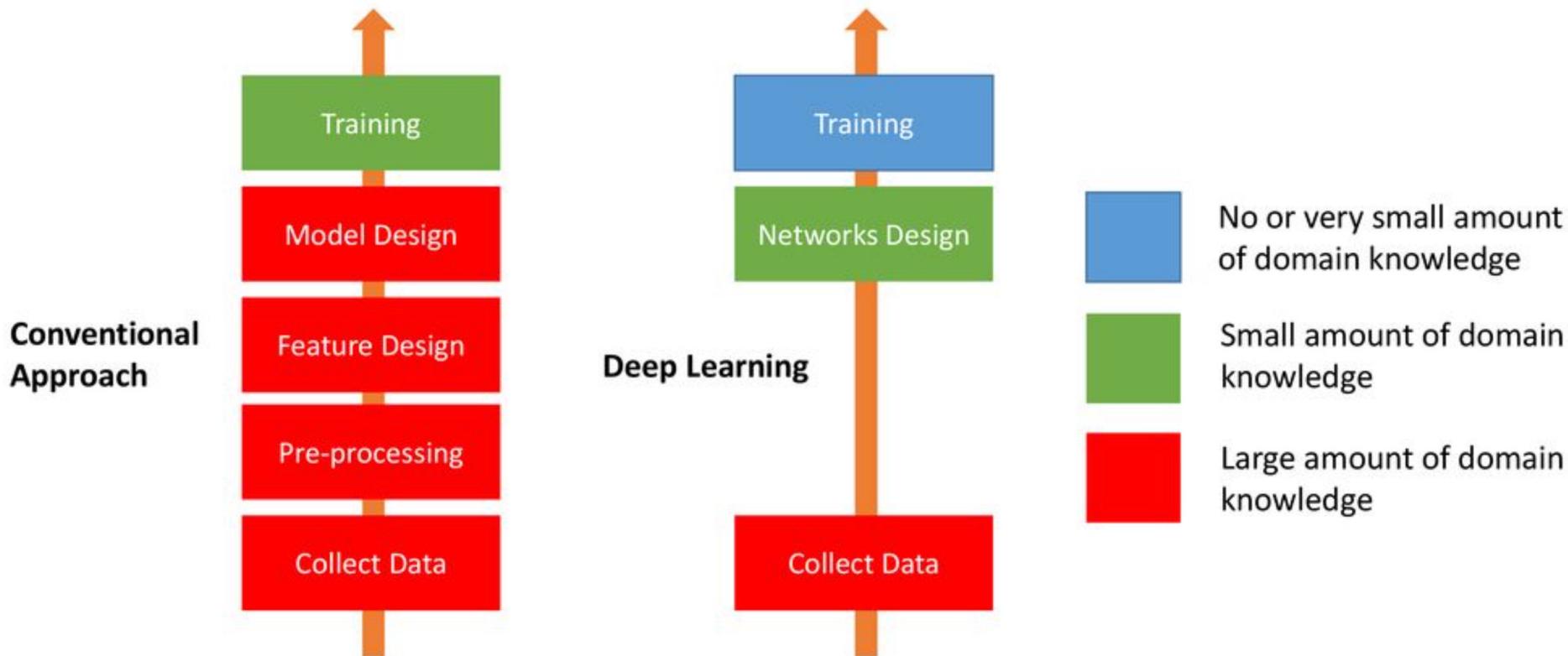
Domain knowledge is the knowledge about the environment in which the data is processed to reveal secrets of the data. In other words, the knowledge of the field that the data belongs to is known as Domain Knowledge.

- **Domain knowledge** is very important when designing DL networks.
 - What are the main challenges? What are the remaining problems?
- Without knowing the conventional models, hard to develop new DL networks that overcome the limitation of conventional approaches.

Domain Knowledge in Data Science

- Domain Knowledge makes you have the ability to anticipate which features are probably important for the problem at hand and which ones are not.
- Domain Knowledge not only affects feature engineering. It affects how you define success, what model you can use, how often you have to retrain it, how it can be deployed, and how you make sense of the results.
- Knowledge is indeed power and it's something you have to gain to take your data science or any skills to the next level.

Conventional vs. Deep-learning-based Approaches (cont.)



DL Network Development with Domain Knowledge

- DL network design w/o domain knowledge about the problem addressed
 - More trial-and-error required
 - Difficult to analyze the results
 - Hard to update the current network to improve the performance
- DL network design w/ domain knowledge about the problem addressed
 - Less trial-and-error
 - Easy to analyze the results
 - Possible to adopt existing networks designed for other tasks