

# Lecture 18

## Motion estimation

Deqing Sun (Google)

# Credits

Some slides, images, and videos from

- Dr. Ce Liu@Microsoft
- Dr. Huaizu Jiang@Northeastern
- Dr. David Fouhey@UMich
- Dr. Justin Johnson@UMich
- Dr. Svetlana Lazebnik@UIUC
- Dr. Shree K. Nayar@Columbia
- Dr. Jia Deng@Princeton
- Dr. Ming-Hsuan Yang@UC Merced
- Dr. Rick Szeliski's book
- Book by Antonio, Phillip, and Bill

Suggestions from Dr. Bill Freeman, Dr. Rick Szeliski, Dr. Noah Snavely and Dr. Junhwa Hur

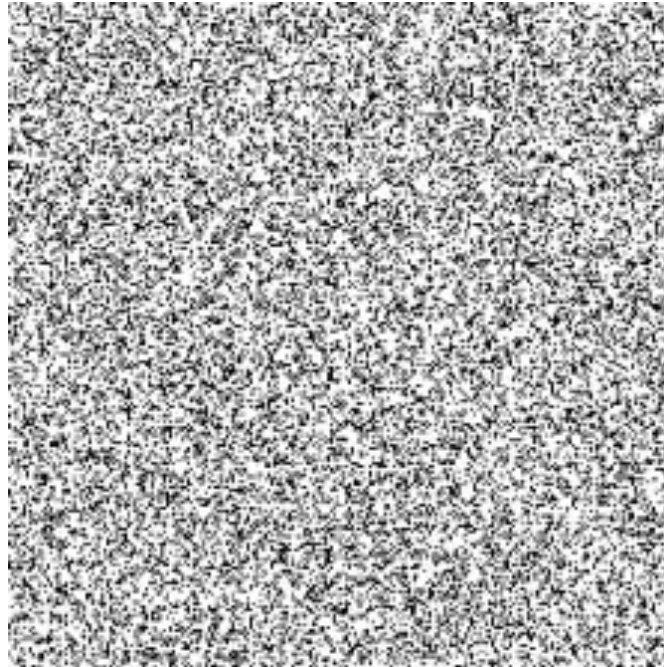
# We live in a dynamic world

Perceiving, understanding and predicting motion is an important part of our daily lives



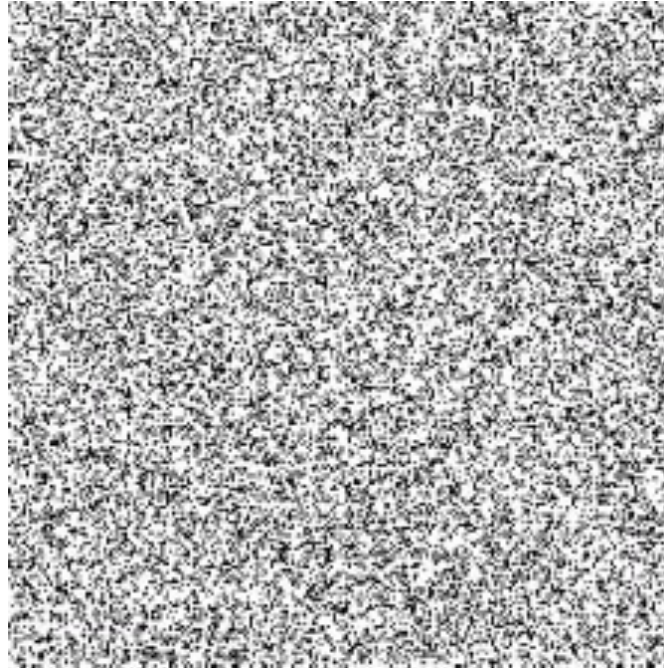
# Motion and perceptual organization

Sometimes motion is the only cue



# Motion and perceptual organization

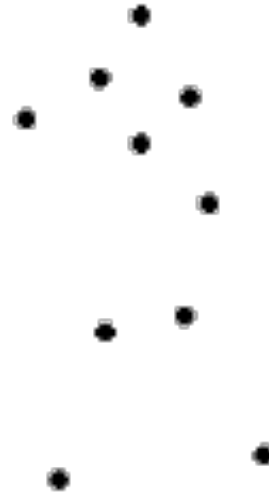
Sometimes motion is the only cue





# Motion and perceptual organization

Even impoverished data can create a strong percept



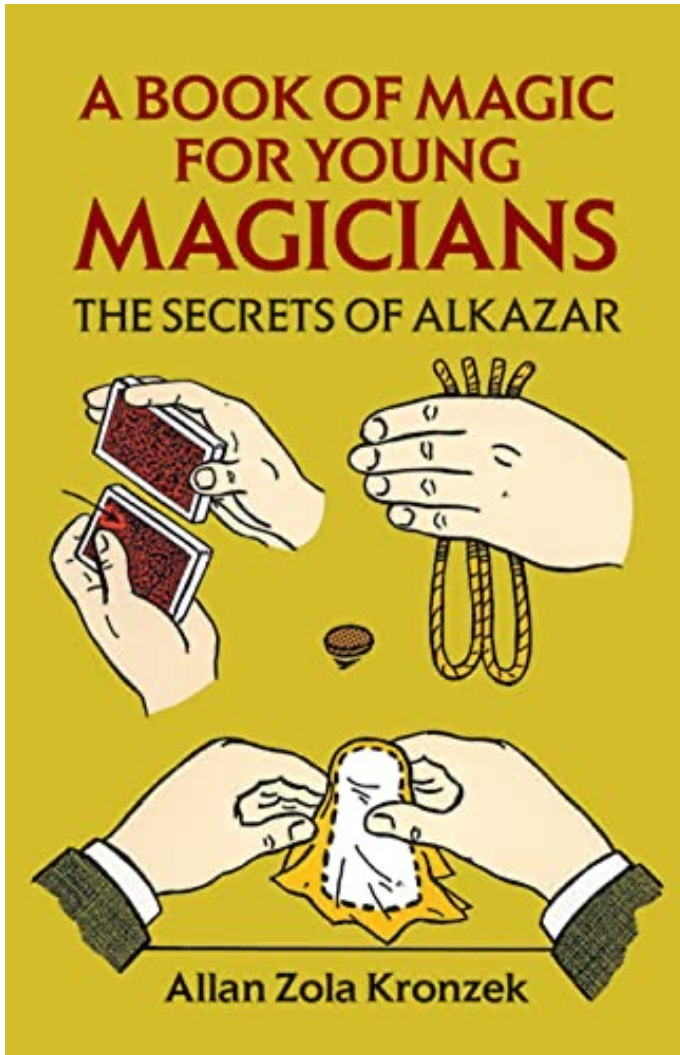
# Motion and perceptual organization

Even impoverished data can create a strong percept





# We pay attention to motion



Alkazar's principles of misdirection

The key to misdirection lies in learning to control **attention**.

Principle 1

The audience will pay attention to what **moves**. ...

What doesn't **move** ... doesn't attract attention.

...

# Content

- Classical approach
- Deep learning-based approach
- Applications: What is motion for?

# **Classical approach**

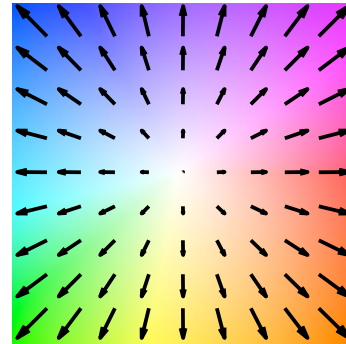
# Optical flow: 2D motion of every pixel



Input [Liu *et al.* CVPR'08]



Optical flow (2D motion vector)



Color key  
[Baker *et al.* IJCV'11]

# Fundamental assumption: Brightness constancy

[Horn & Schunck AI'81]

$$I_t(\mathbf{p}) \approx I_{t+1}(\mathbf{p} + \mathbf{w}_p)$$



First image ( $t$ )



Second image ( $t+1$ )

# Matching-based motion estimation

Similarity between a pixel in image 1 with pixels in image 2

$$\min_{\mathbf{w}_p} \left( I_t(\mathbf{p}) - I_{t+1}(\mathbf{p} + \mathbf{w}_p) \right)^2$$



Image 1



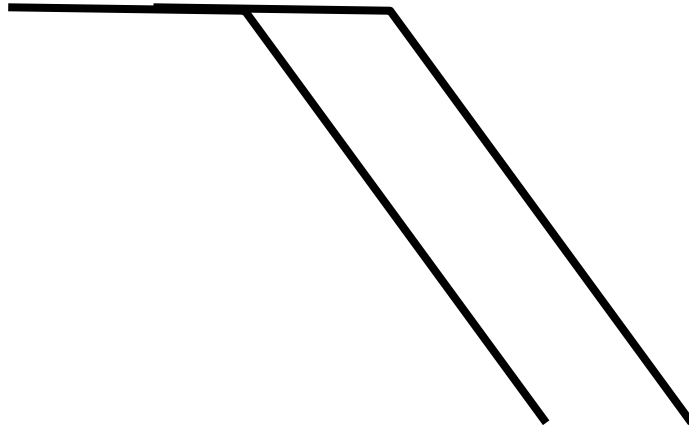
Image 2

# Comparing pixel colors



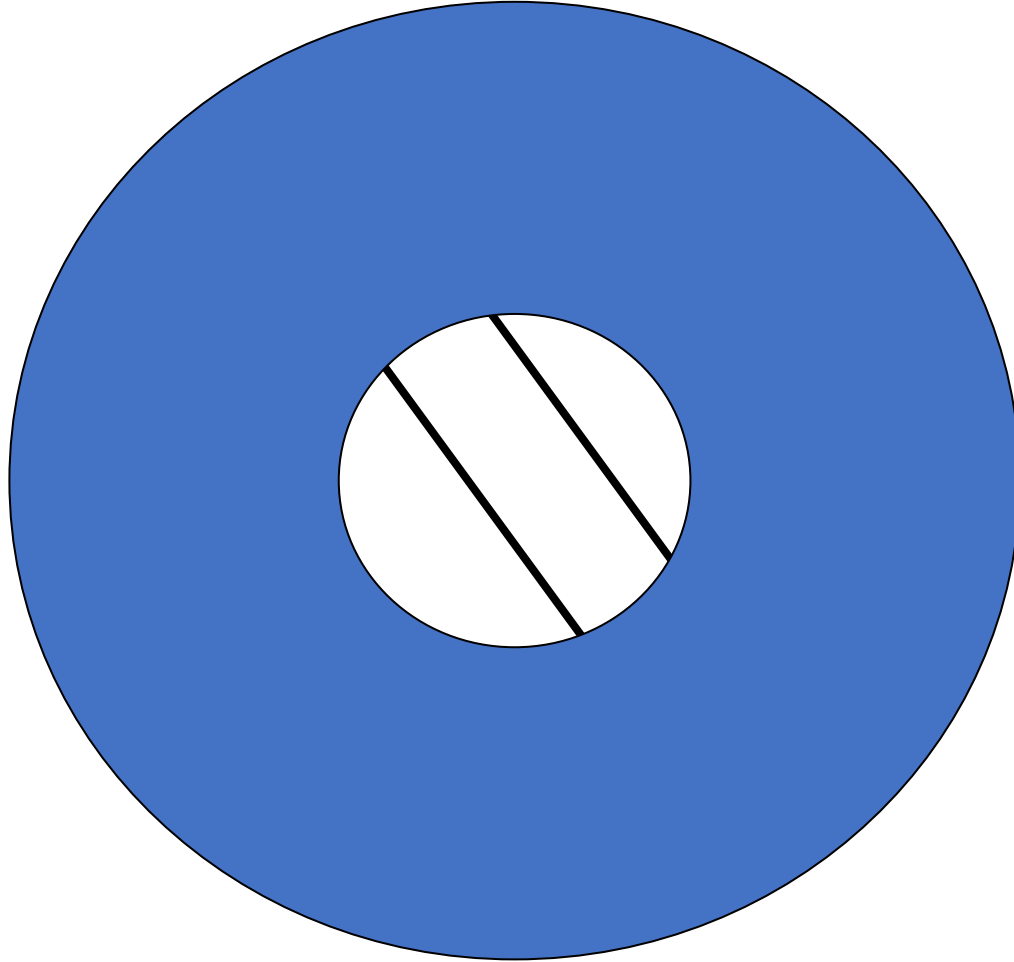
1x1

# Aperture problem

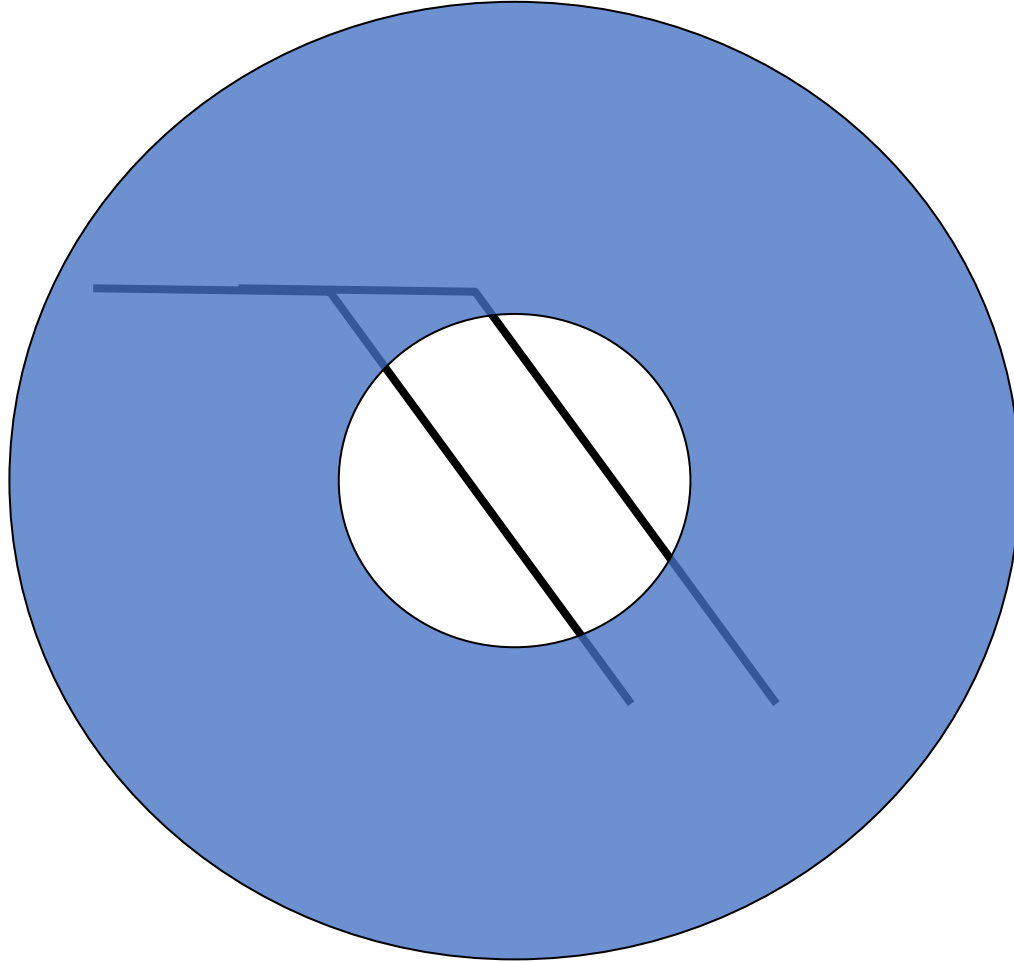




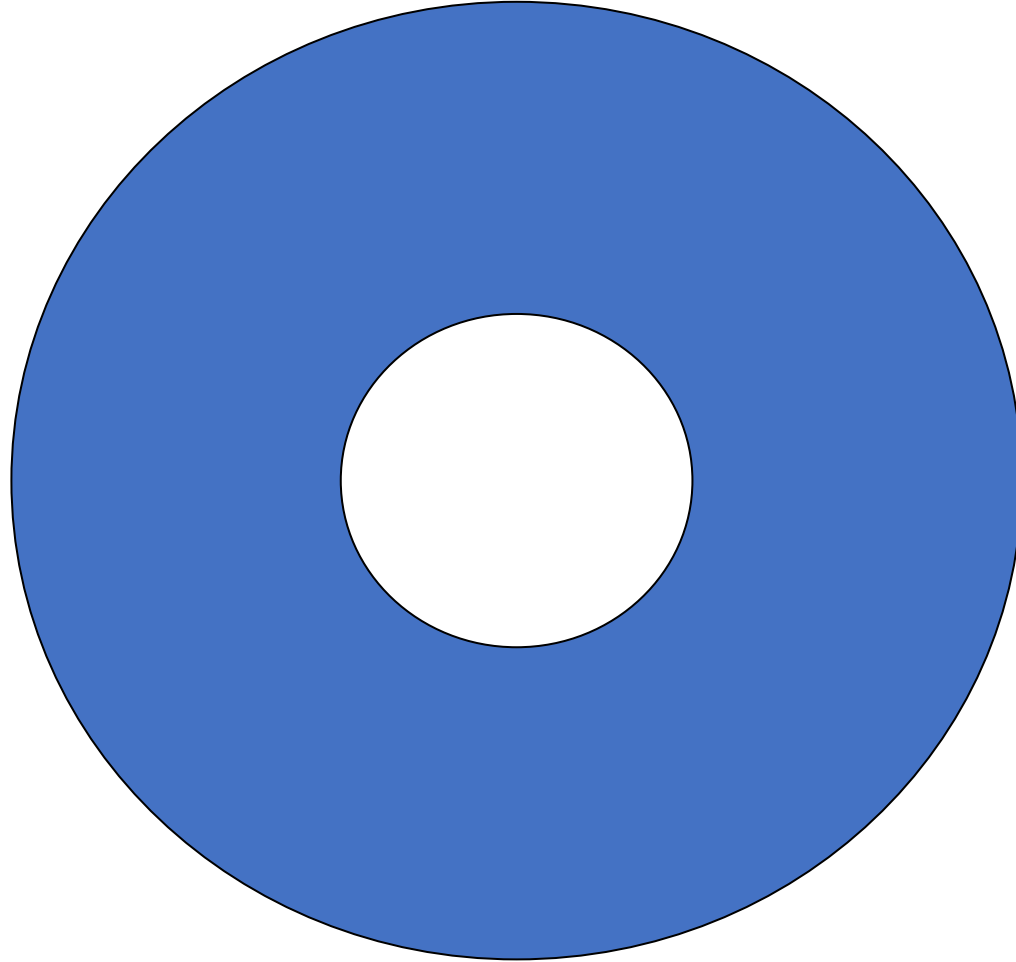
# Aperture problem



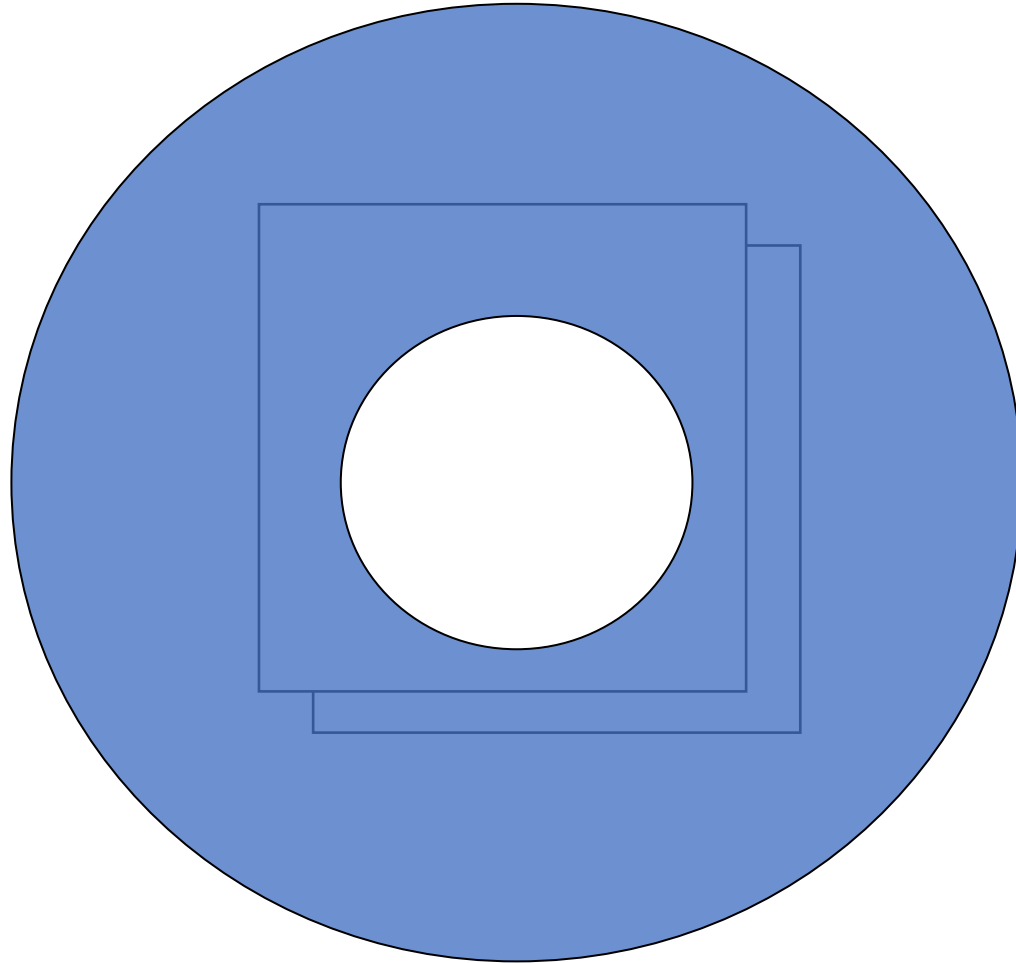
# Aperture problem



# Other invisible flow



# Other invisible flow



# What do you perceive?



Actual motion

# Matching-based motion estimation

Similarity between a pixel in image 1 with pixels in image 2

$$\min_{\mathbf{w}_p} \left( I_t(\mathbf{p}) - I_{t+1}(\mathbf{p} + \mathbf{w}_p) \right)^2$$



Image 1



Image 2

# Solving ambiguities: Lucas-Kanade

Similarity between a patch in image 1 with patches in image 2

- Pixels in a patch share the same (parametric) motion

$$\min_{\mathbf{w}_p} \sum_{\mathbf{q} \in N_p} \left( I_t(\mathbf{q}) - I_{t+1}(\mathbf{q} + \mathbf{w}_p) \right)^2$$



Image 1



Image 2

# Similarity between patches (cost volume)



Image 1



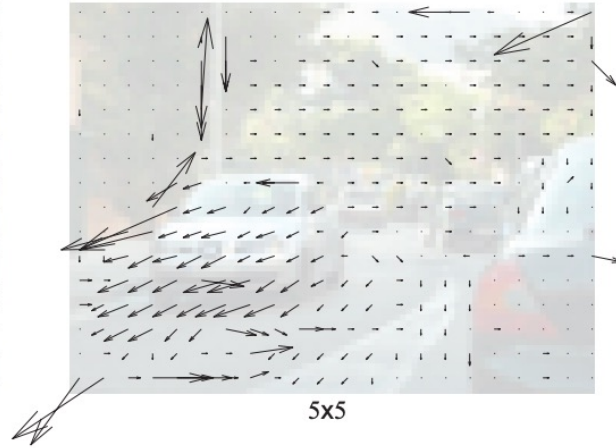
Cost volume (darker, more similar)



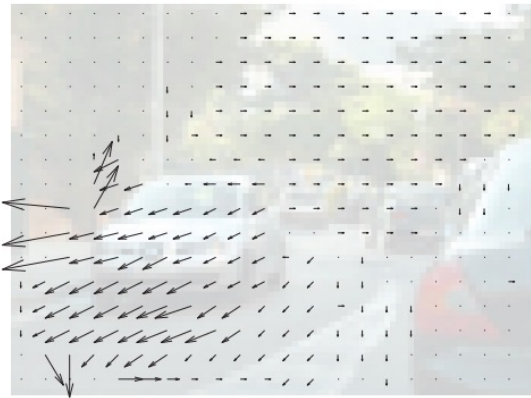
# Effect of patch size



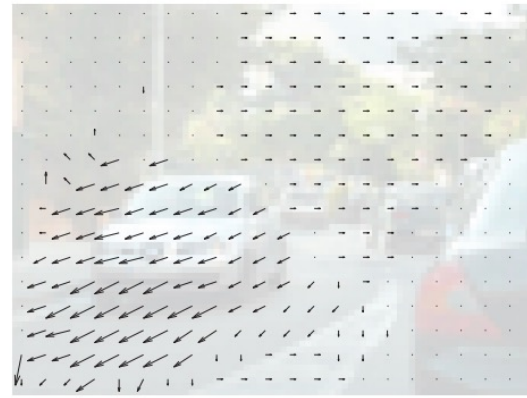
1x1



5x5



11x11



21x21

# Issue: Brute force is too expensive

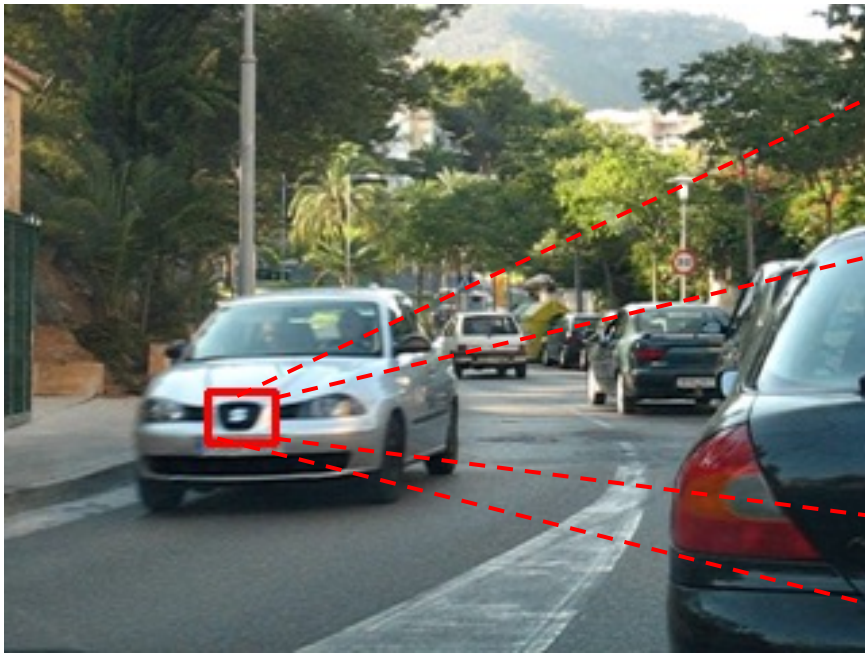


Image 1

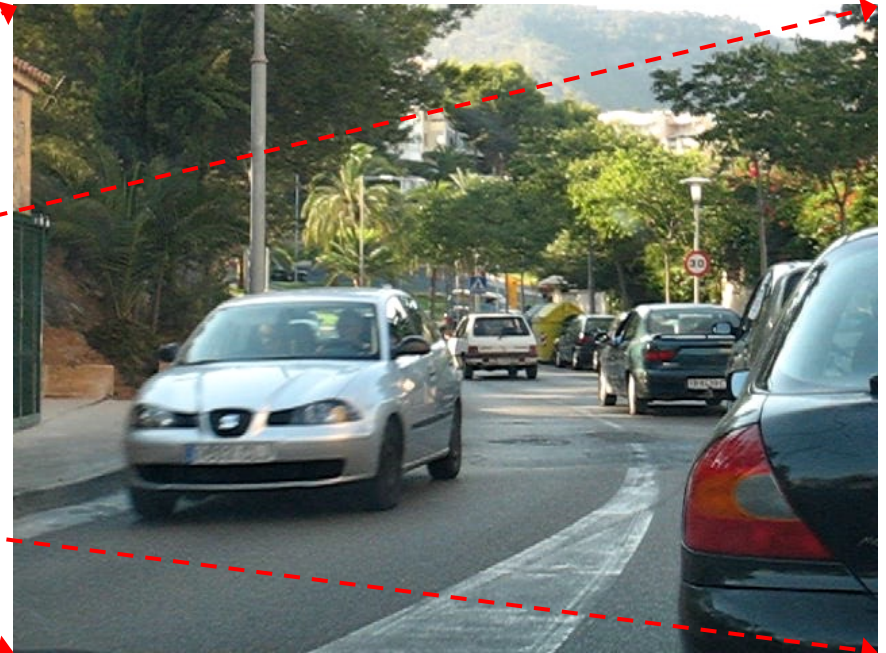
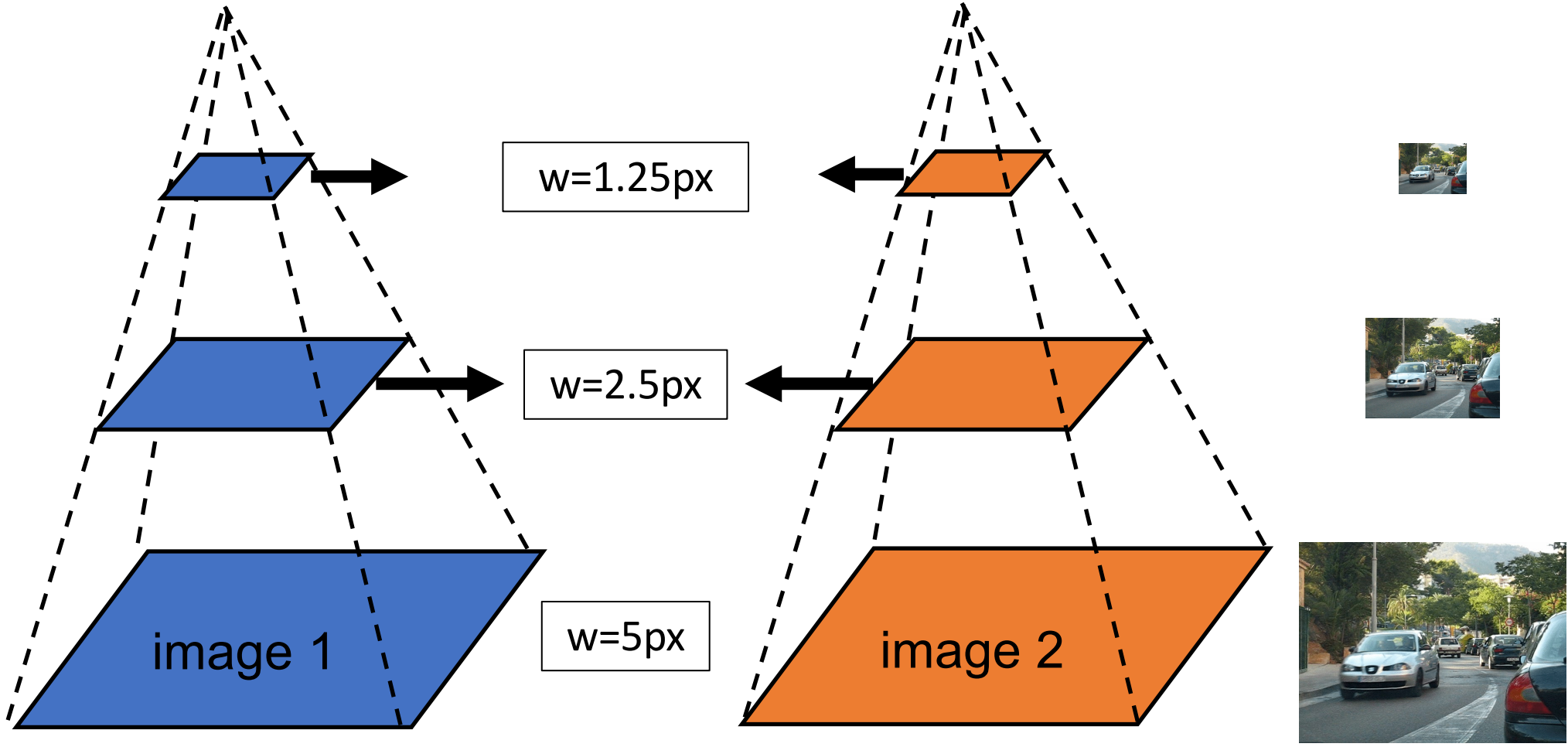


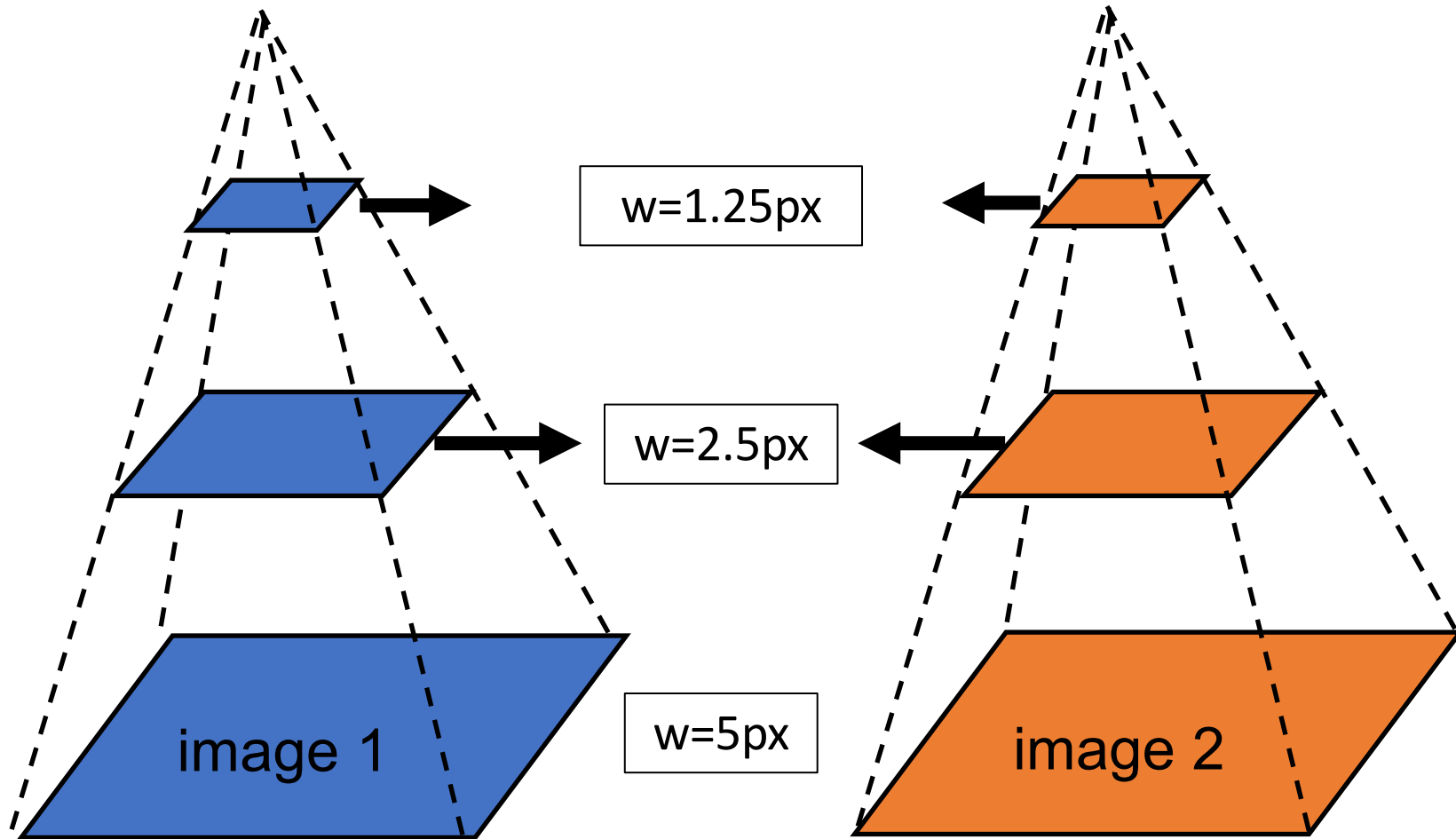
Image 2

# Coarse-to-fine iterative estimation



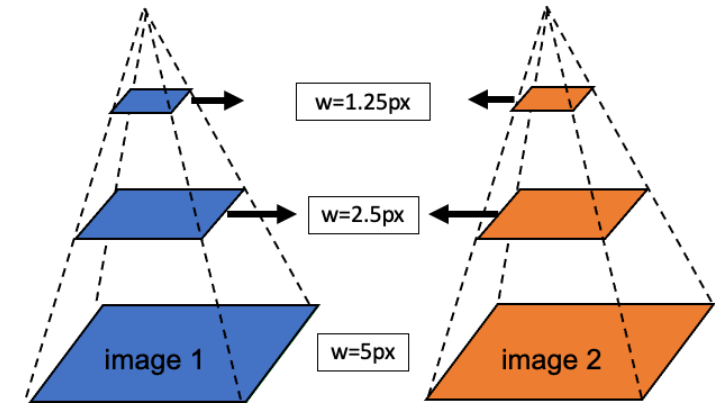
Start from top or bottom?

# Coarse-to-fine iterative estimation



How to use estimates from the upper level?

# Coarse-to-fine iterative estimation



Current estimate



$$I_{t+1}(\mathbf{q} + \mathbf{w}_p) = I_{t+1}(\mathbf{q} + \mathbf{w}_p^k + \delta \mathbf{w}_p)$$



Small increment

Warped image

$$I_w(\mathbf{q}) = I_{t+1}(\mathbf{q} + \mathbf{w}_p^k)$$

# Warping operation

$$I_w(\mathbf{q}) = I_{t+1}(\mathbf{q} + \mathbf{w}_p^k)$$

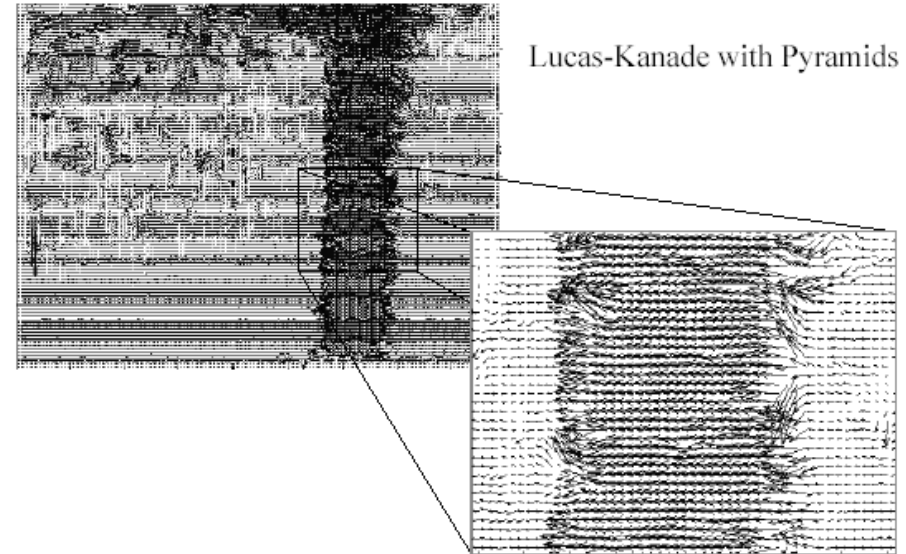
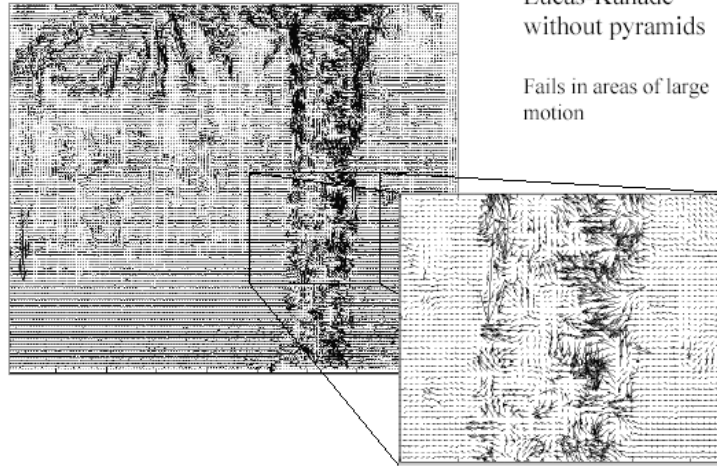


Input images  $t$  and  $t+1$



Image  $t$  and warped image

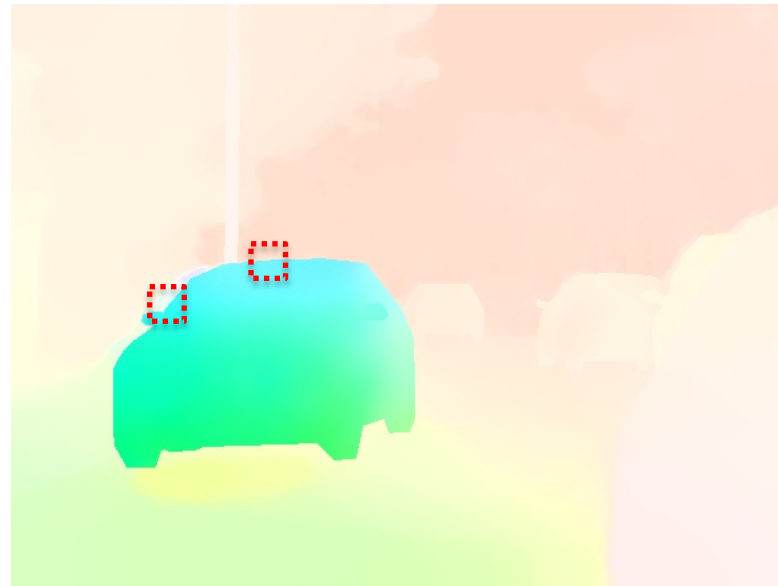
# Optical Flow Results



# Issue: Motion boundaries

Pixels in a patch share the same (parametric) motion

$$\min_{\mathbf{w}_p} \sum_{\mathbf{q} \in N_p} \left( I_t(\mathbf{q}) - I_{t+1}(\mathbf{q} + \mathbf{w}_p) \right)^2$$



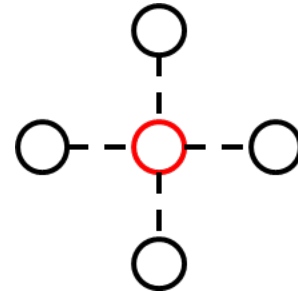


# Solving ambiguities: Horn & Schunck

Smoothness: neighboring pixels have similar motion

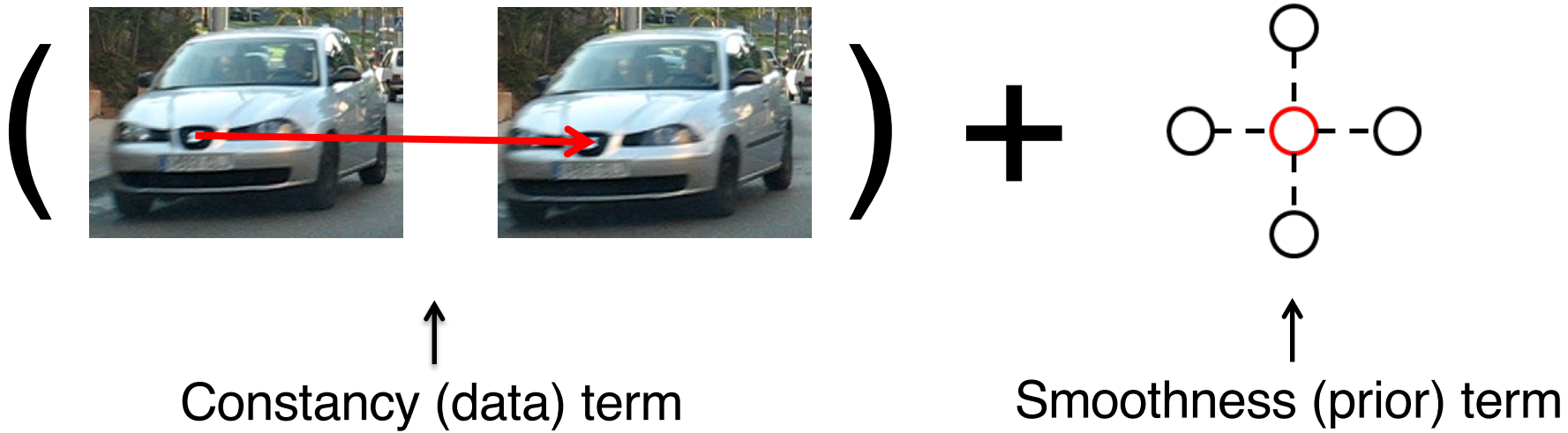


$$\mathbf{w}_p \approx \mathbf{w}_q, \mathbf{q} \in N_p$$



# Optimization/energy minimization

[Horn & Schunck AI'81]



$$E(\mathbf{w}) = \sum_{\mathbf{p}} |I_t(\mathbf{p}) - I_{t+1}(\mathbf{p} + \mathbf{w}_p)|^2 + \lambda \sum_{\mathbf{q} \in N_p} |\mathbf{w}_p - \mathbf{w}_q|^2$$

# Horn & Schunck



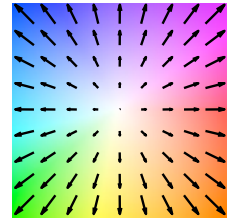
Input



Horn & Schunck



Ground truth

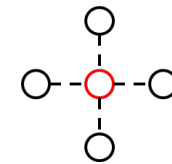


Color key  
[Baker *et al.* IJCV'11]

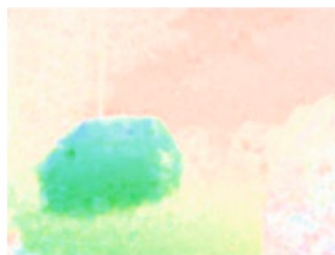
# Improving Horn & Schunck

[Sun *et al.* CVPR'10, IJCV'14]

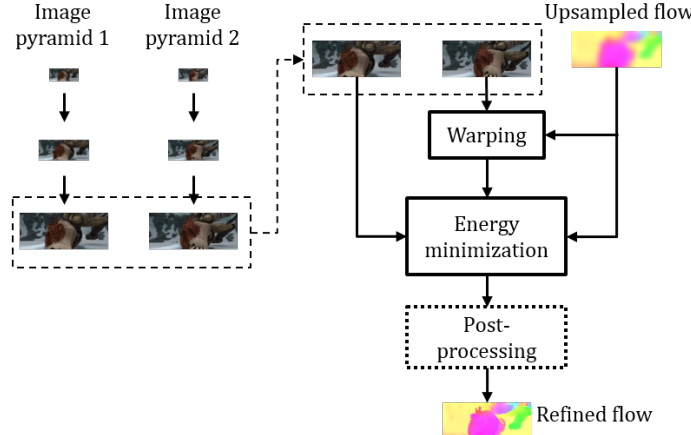
$$E(\mathbf{w}) = \sum_{\mathbf{p}} |I_t(\mathbf{p}) - I_{t+1}(\mathbf{p} + \mathbf{w}_{\mathbf{p}})|^2 + \lambda \sum_{\mathbf{q} \in N_{\mathbf{p}}} |\mathbf{w}_{\mathbf{p}} - \mathbf{w}_{\mathbf{q}}|^2$$



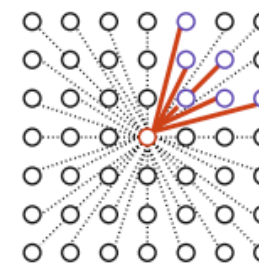
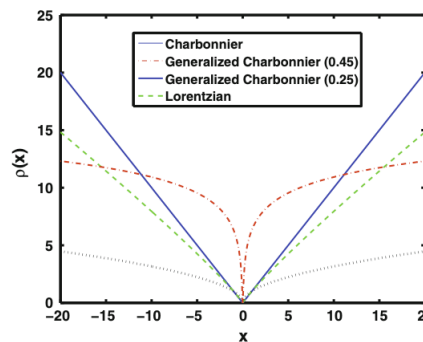
“Old”



Better optimization/  
implementation



Ground truth



# Improving Horn & Schunck

[Sun *et al.* CVPR'10, IJCV'14]



(a) "Old" HS [58]

(b) "New" HS

(c) Classic++

(d) Classic+NL

(e) Ground truth

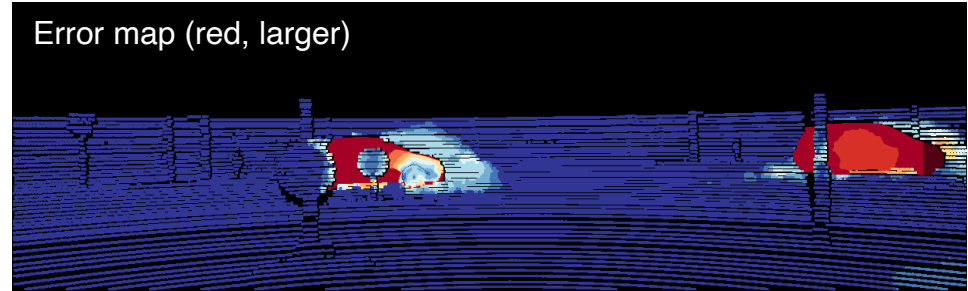
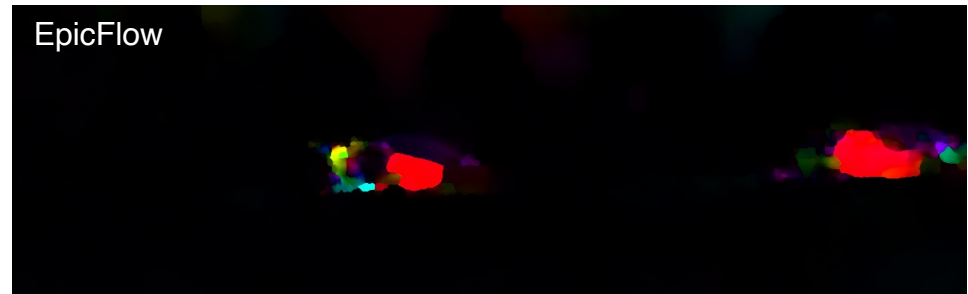
(f) First frame

Implementation    Robust function    Non-local  
smoothness term

# Challenges for classical methods

Large motion  
Motion blur  
Occlusions  
Lighting changes  
Noise...

Hard to modify  
objective function  
and even harder to  
optimize it



# Content

- Classical approach
  - Constancy assumption -> matching by comparison (cost volume)

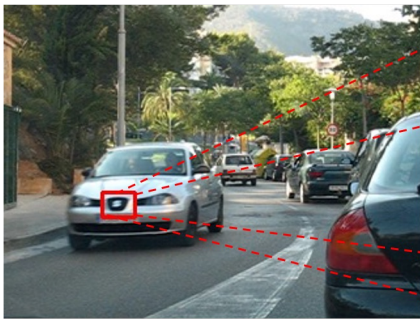


Image 1

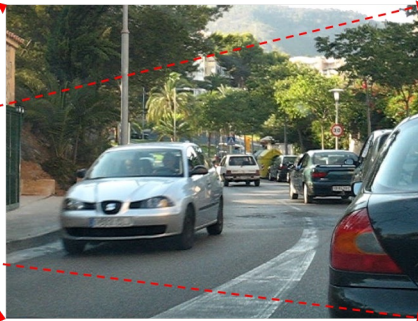


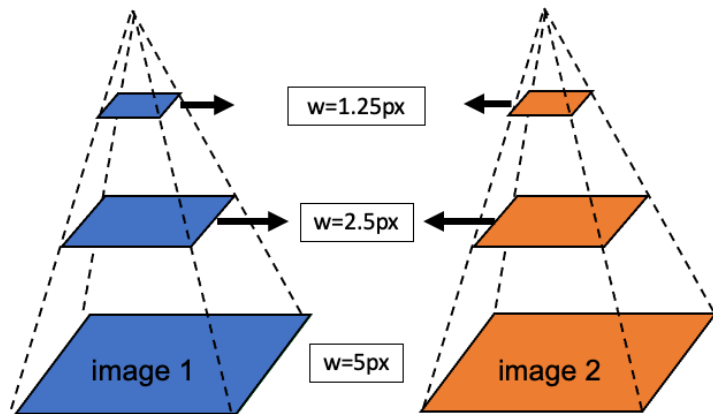
Image 2



Cost volume

# Content

- Classical approach
  - Constancy assumption -> matching by comparison (cost volume)
  - Coarse-to-fine, warping-based iterative estimation

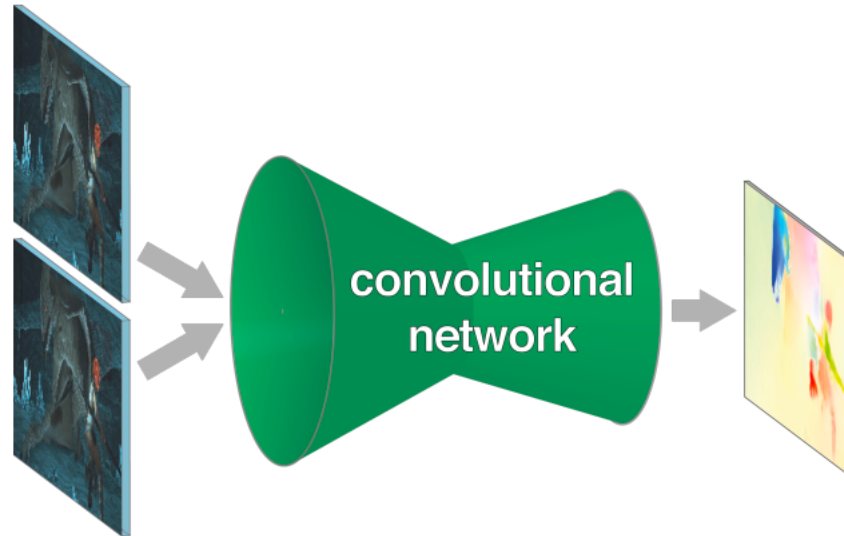




# **Deep learning-based approach**

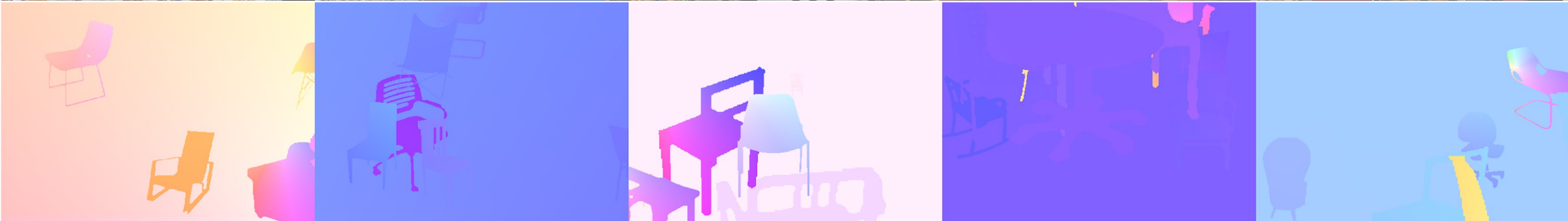
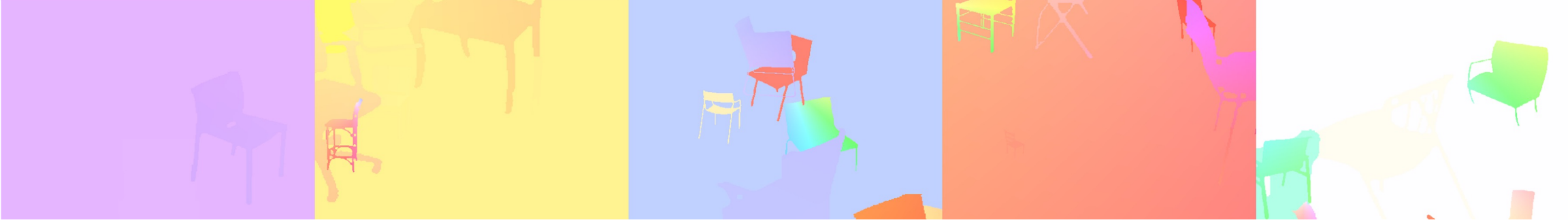
# Supervised optical flow

[Dosovitskiy *et al.* ICCV'15]



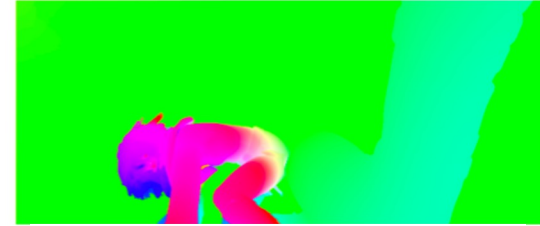
What is the training/test data?

# Training Data: FlyingChairs

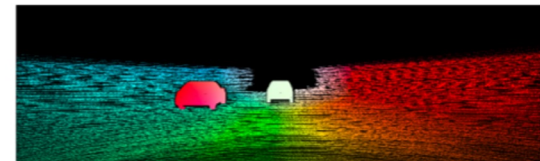
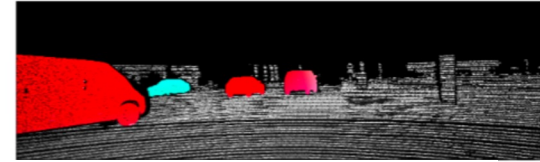


# Two widely-used benchmarks for optical flow

Sintel (Blender movie)

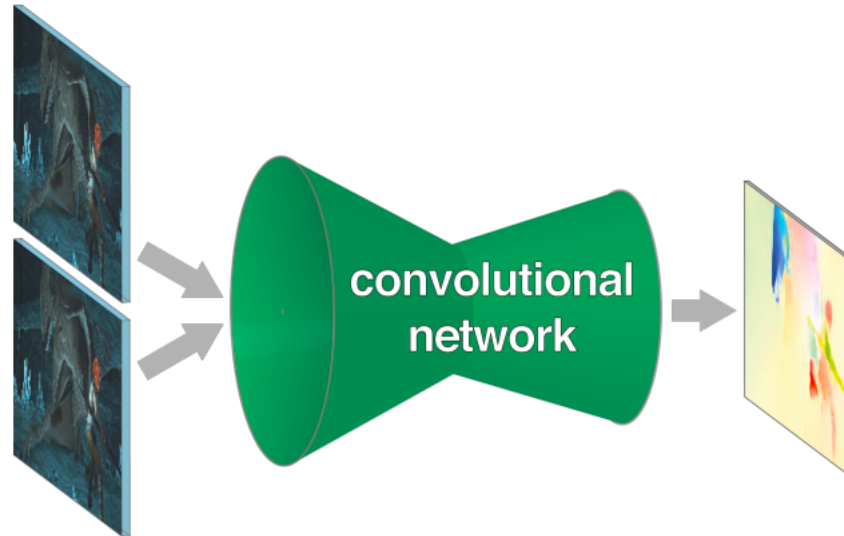


KITTI (driving)



# Supervised optical flow

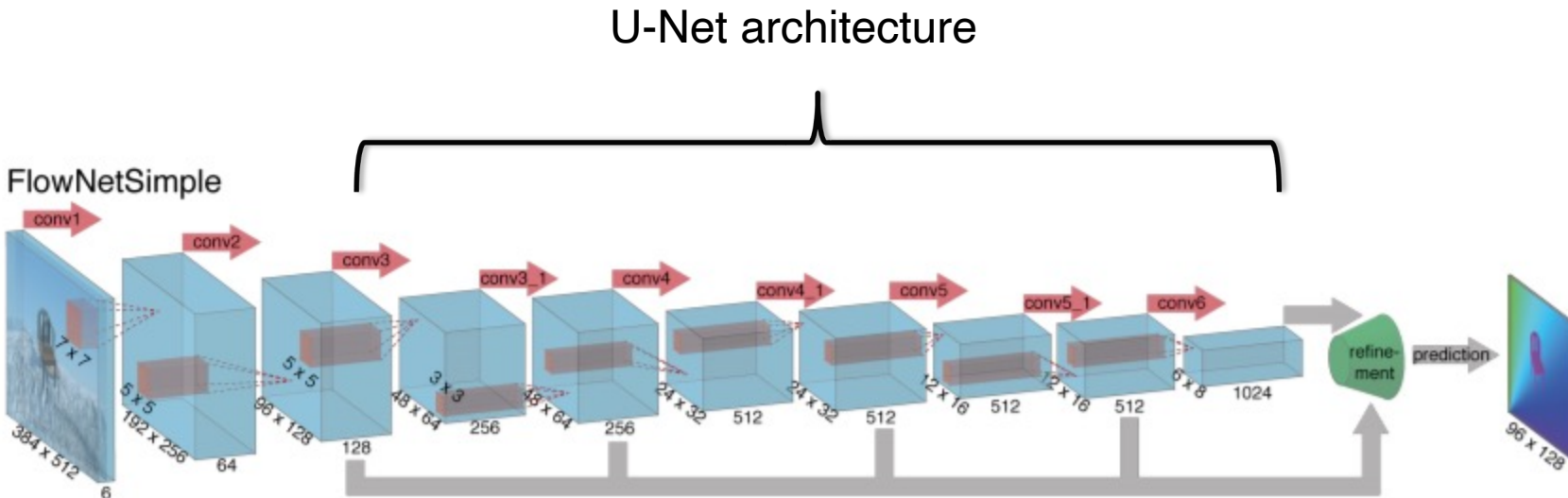
[Dosovitskiy *et al.* ICCV'15]



What is the network/architecture?

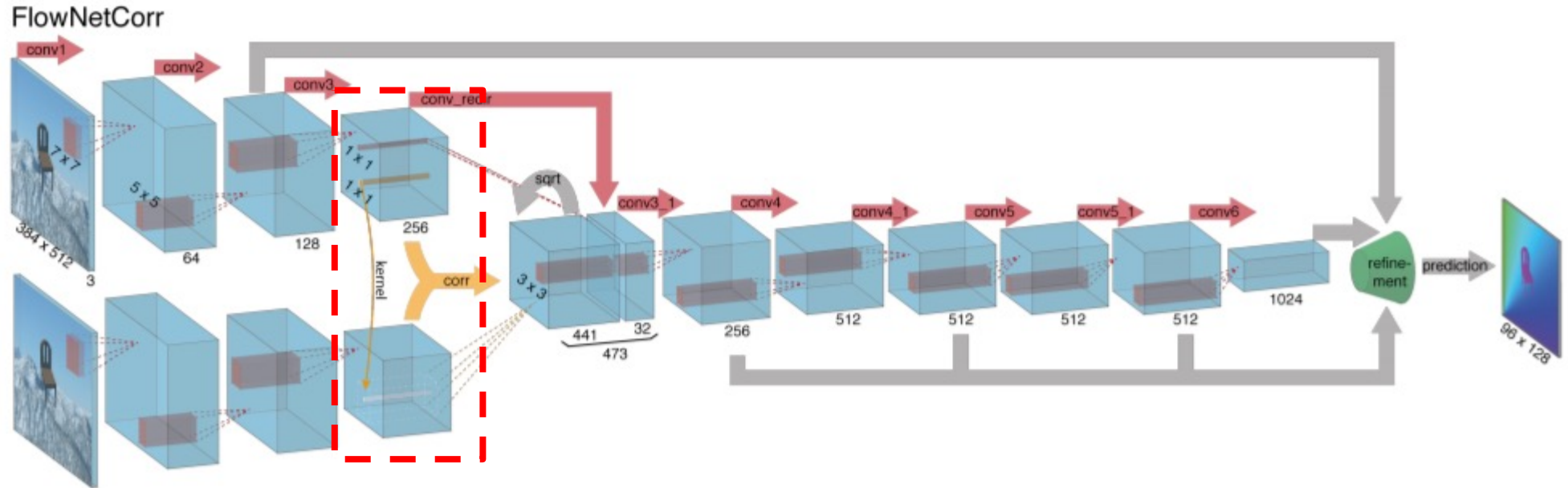
# FlowNetS(imple): Mapping from images to flow

[Dosovitskiy *et al.* ICCV'15]



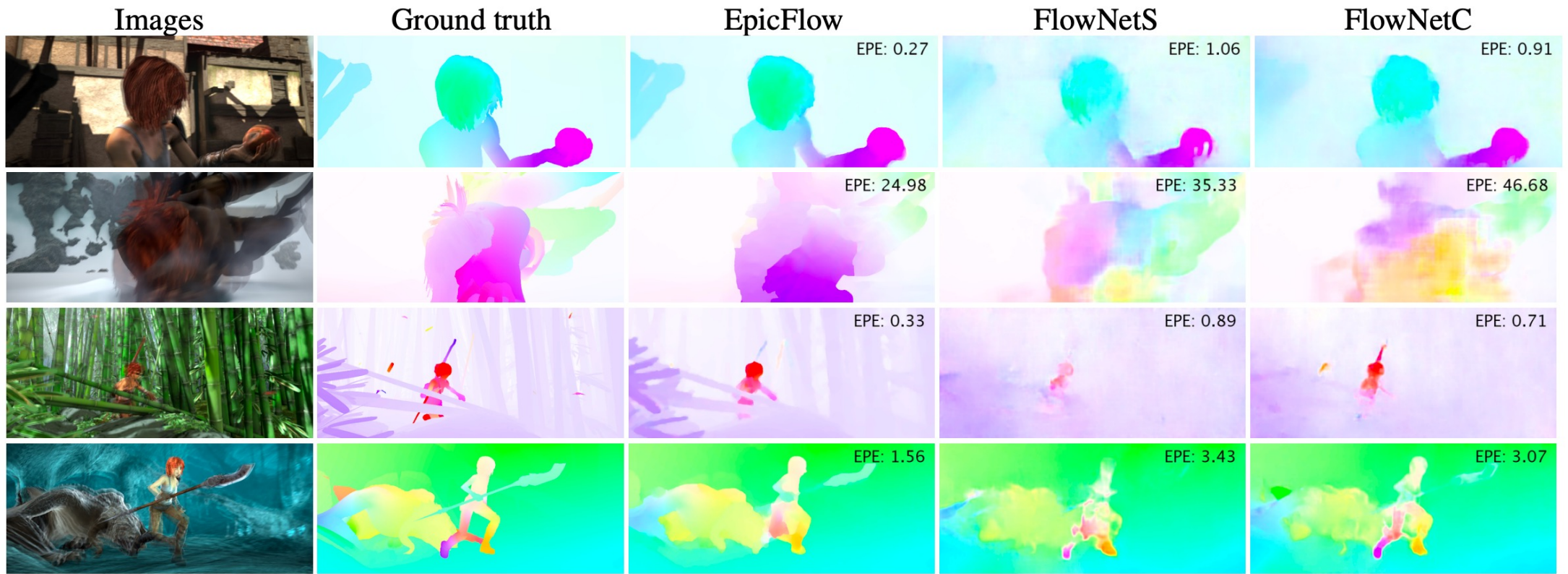
# FlowNetC(orrelation): Compare features

[Dosovitskiy *et al.* ICCV'15]



# Promising but behind contemporary state of the art

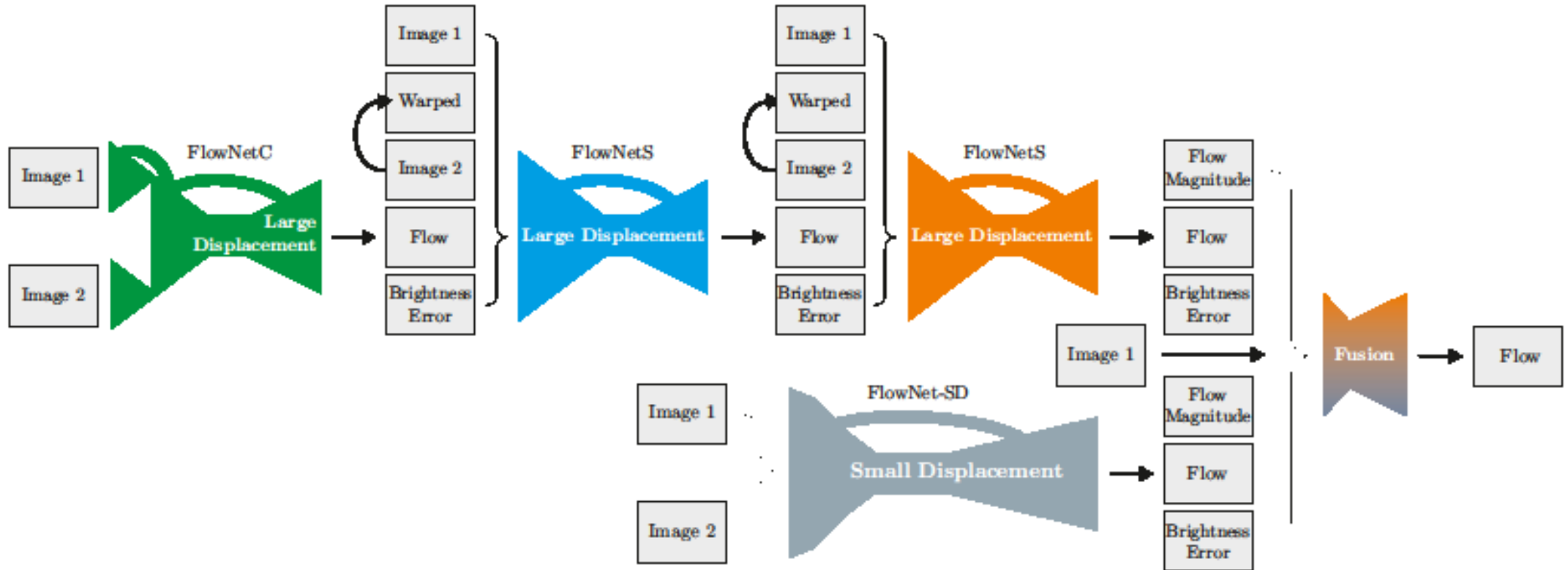
[Dosovitskiy *et al.* ICCV'15]



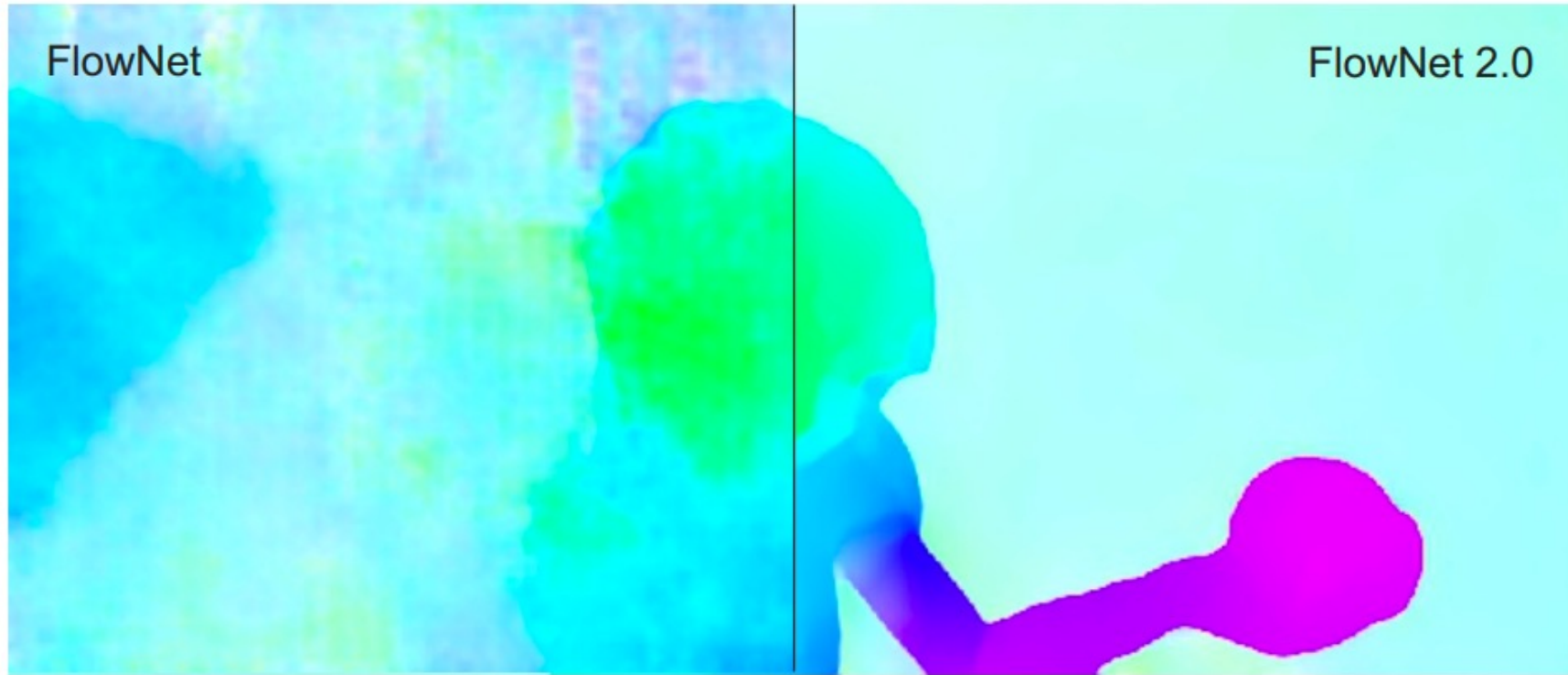


# FlowNet2: Scaling up by stacking up FlowNetS/C

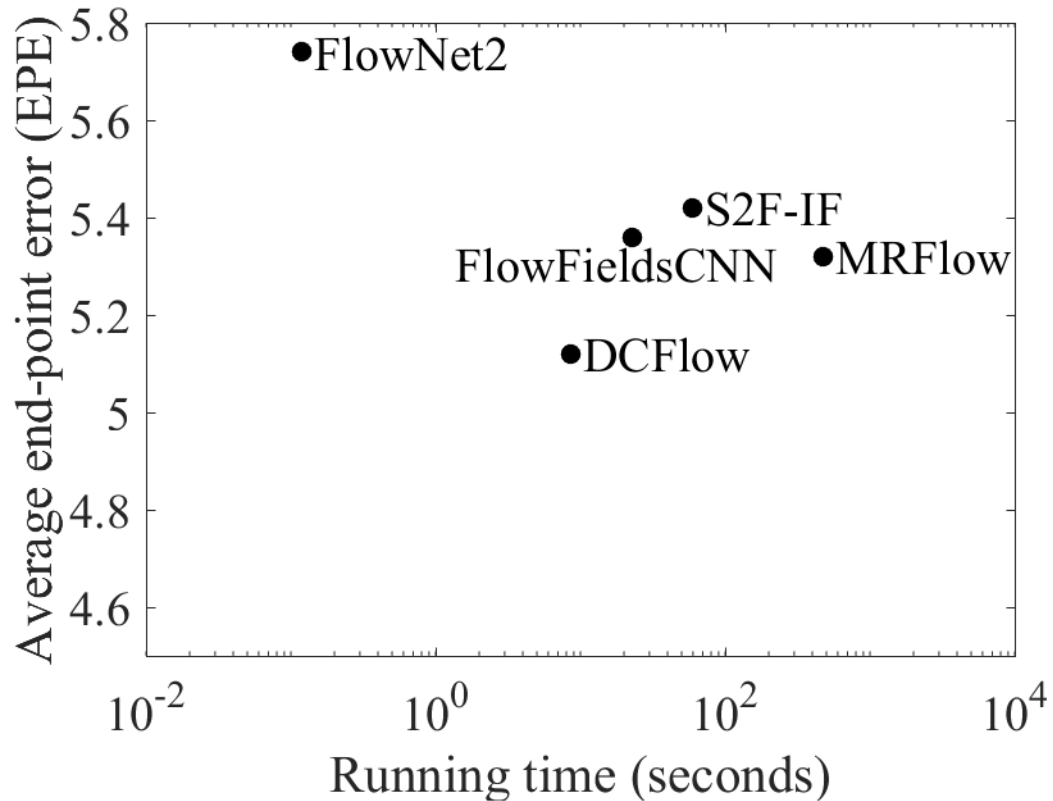
[Ilg et al. CVPR'17]



# Significant improvement

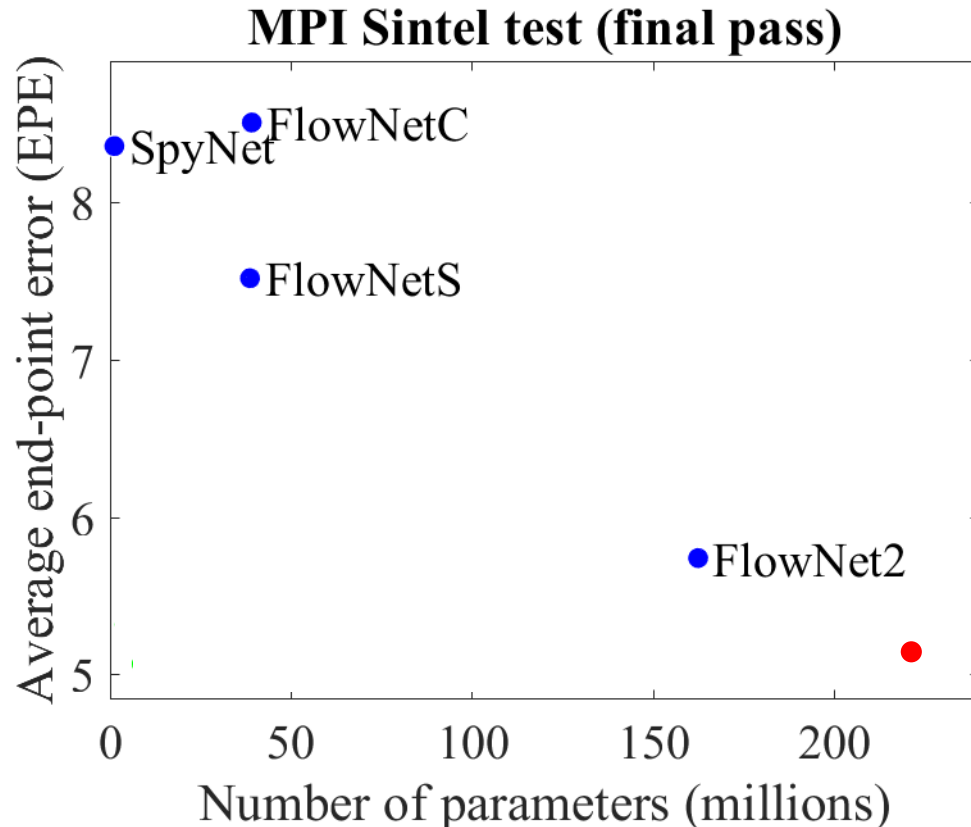


# Trade-off between accuracy and running time



FlowNet2: Ilg *et al.* CVPR'17  
S2F-IF: Yang & Soatto CVPR'17  
FlowFieldsCNN: Bailer *et al.* CVPR'17  
MRFlow: Wulff *et al.* CVPR'17  
DCFlow: Xu *et al.* CVPR'17

# Trade-off between accuracy and size for CNN methods



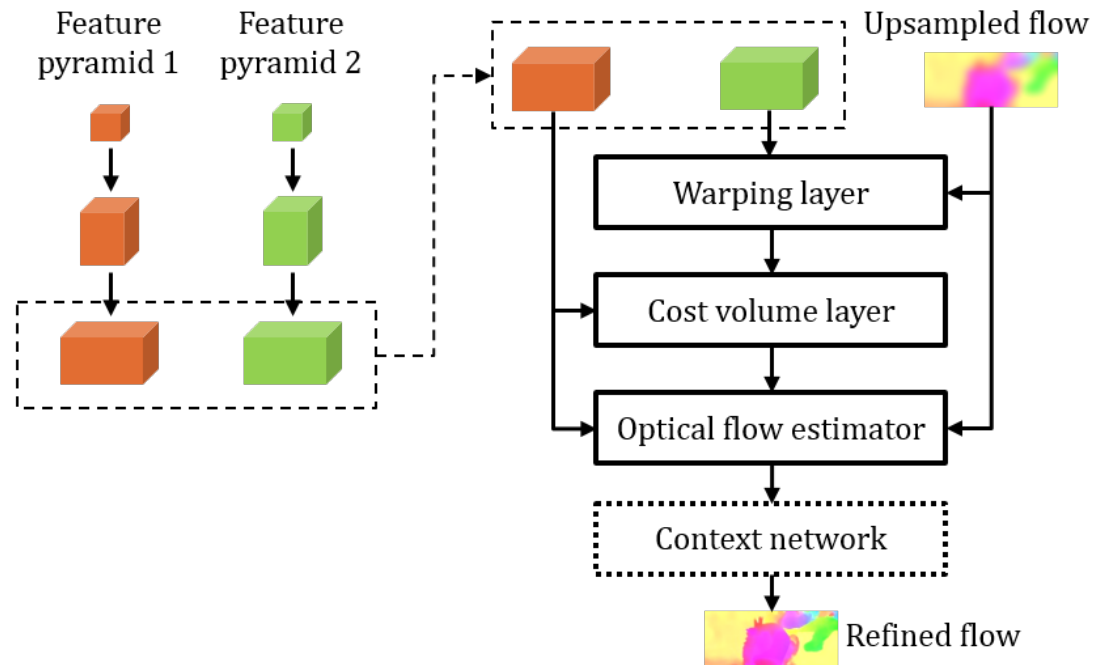
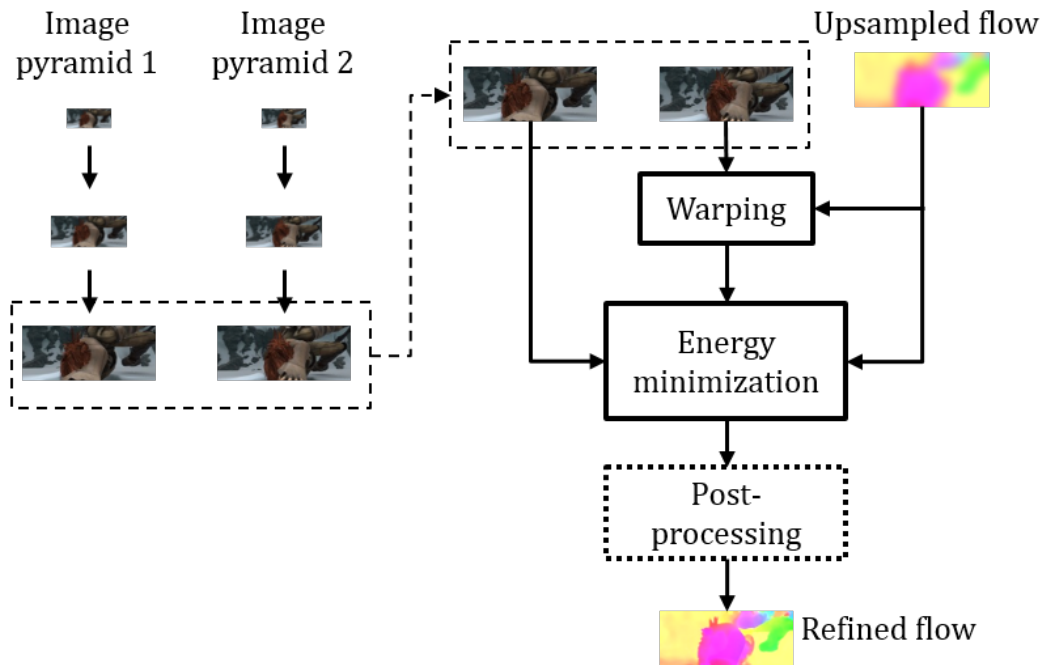
FlowNetS/C: Dosovitskiy *et al.* ICCV'15  
FlowNet2: Ilg *et al.* CVPR'17  
SpyNet: Ranjan & Black CVPR'17

PWC-Net: Sun *et al.* CVPR'18

**Pyramid, warping, & cost volume,**  
not PricewaterhouseCoopers 😊

Even bigger model?

# Inspired by classical approach



# Pyramid of learnable features

(width, height, channel number)

(1024, 512, 3)

(512, 256, 16)

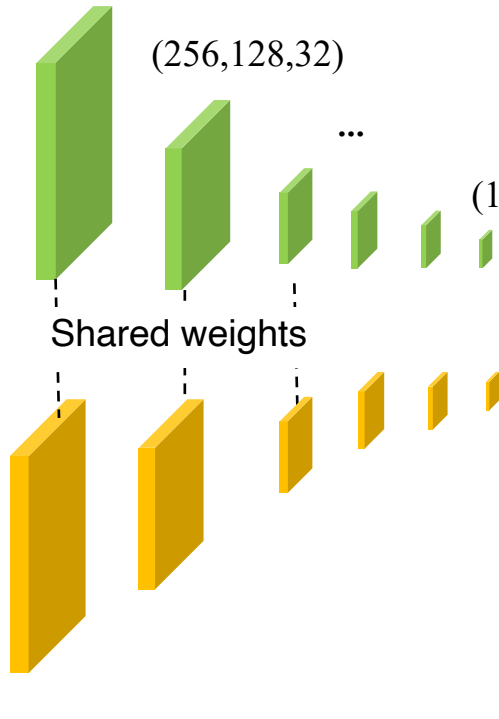
(256, 128, 32)

...

(16, 8, 192)

Large receptive field

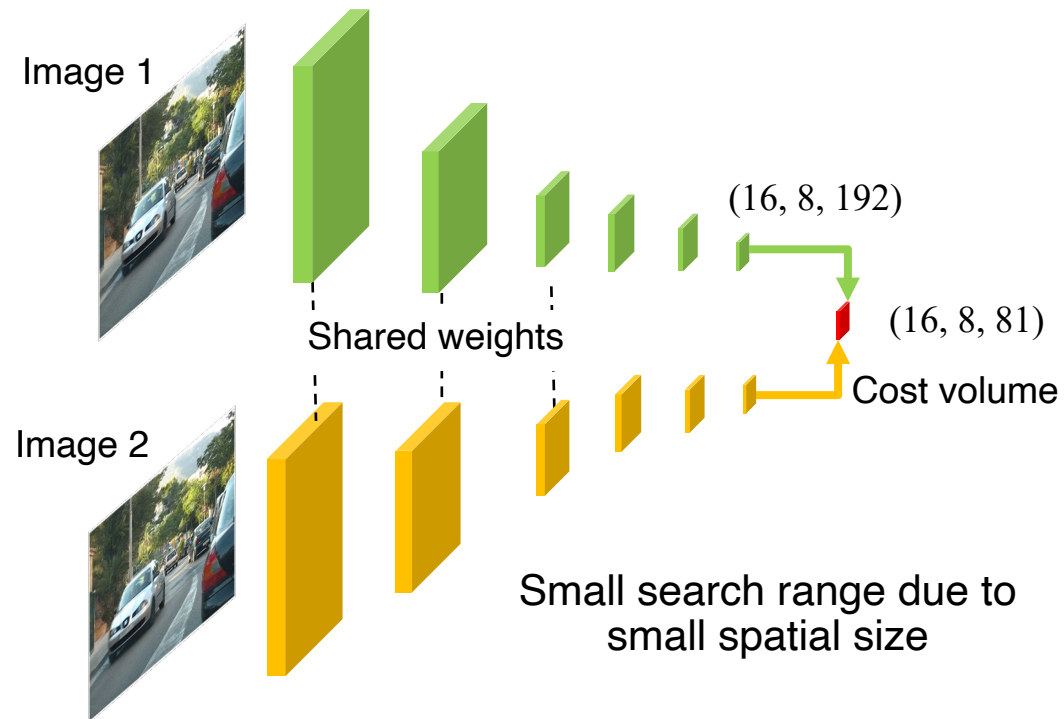
Shared weights



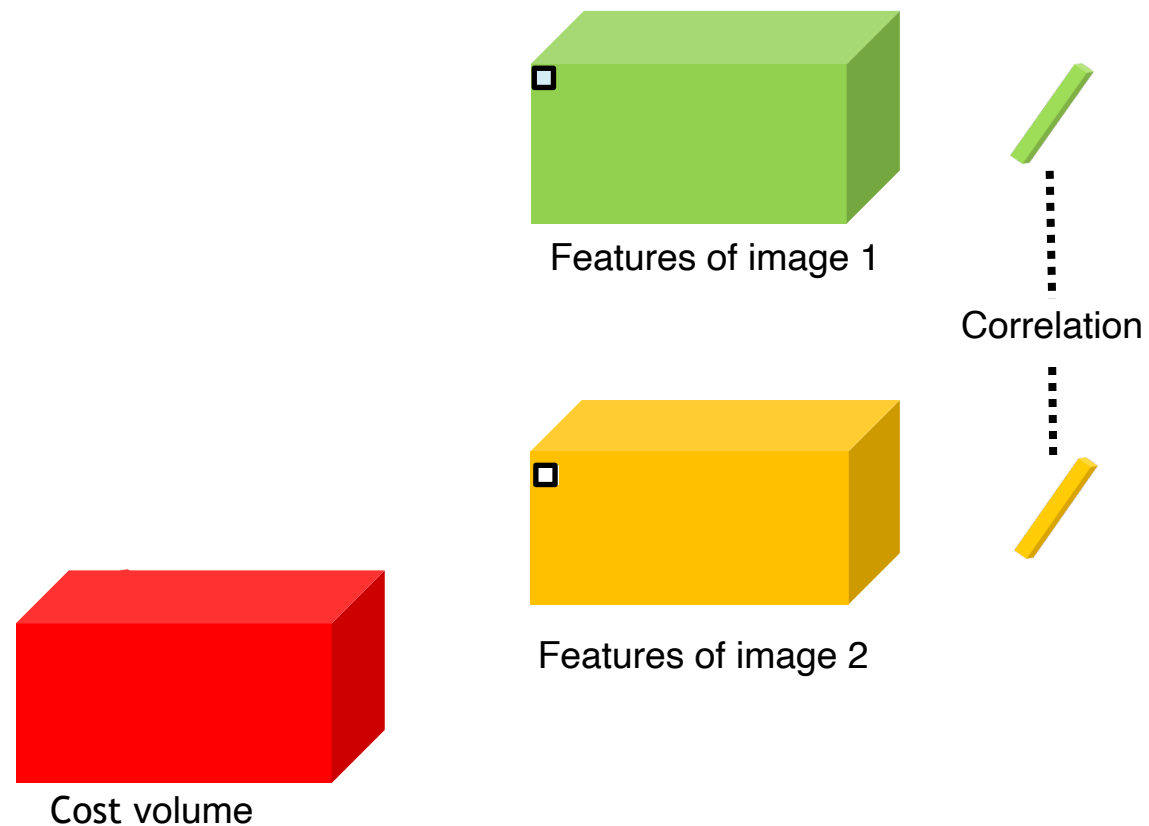
# Compute cost volume by correlation

(width, height, channel number)

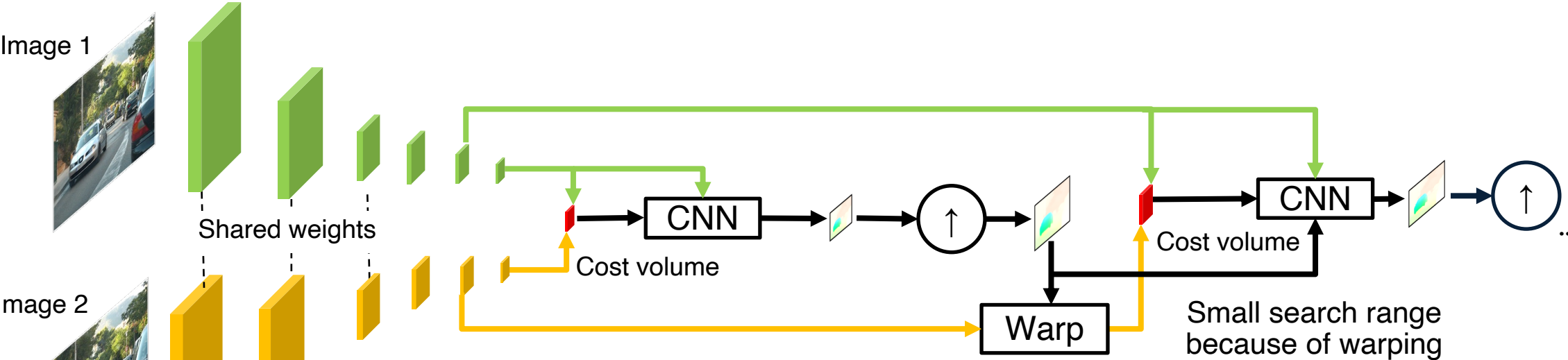
(1024, 512, 3)



[Dosovitskiy *et al.* FlowNet ICCV'15]



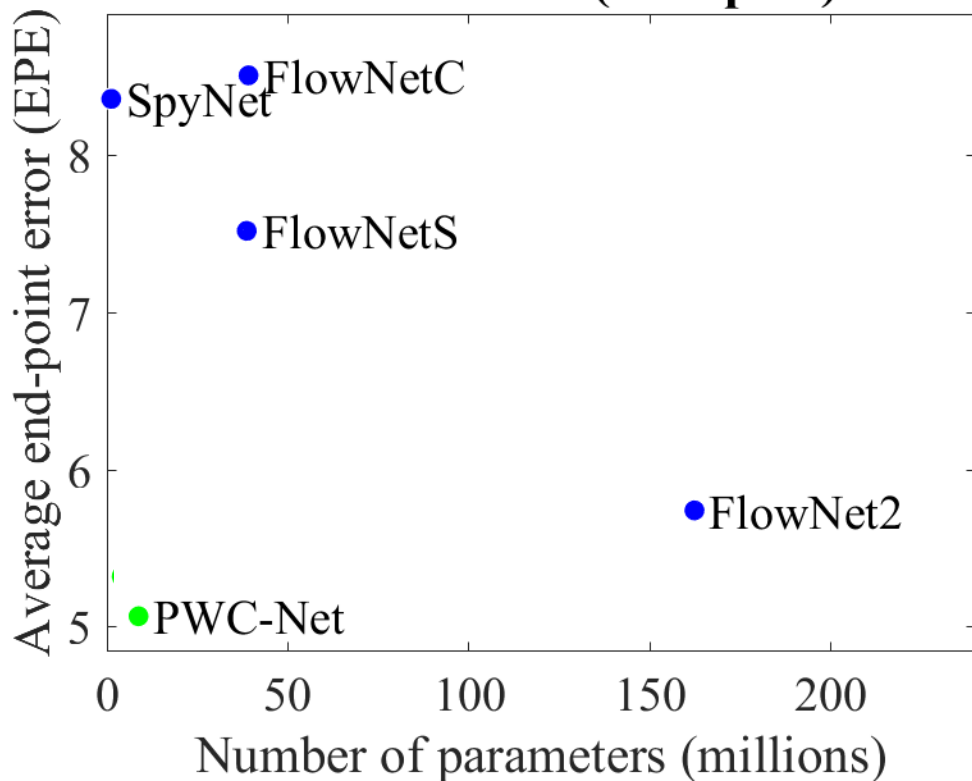
# Mapping cost volume to optical flow






# Architectures matter

MPI Sintel test (final pass)



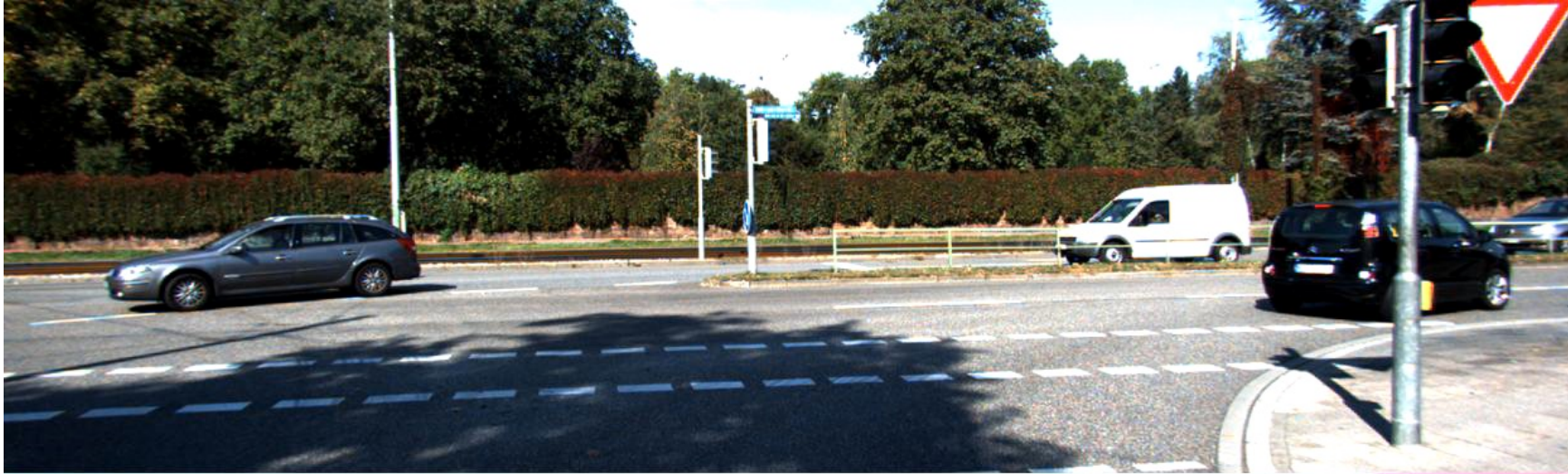
Flow Leaderboard

Final results of the ROB 2018 Challenge. New submissions will be accepted after CVPR 2018.

 Method	Middlebury <small>(Detailed subrankings)</small>	KITTI <small>(Detailed subrankings)</small>	MPI Sintel <small>(Detailed subrankings)</small>	HD1K <small>(Detailed subrankings)</small>
<b>1</b> PWC-Net_ROB	2	2	2	1
<small>PWC-Net: CNNs for Optical Flow Using Pyramid, Warping, and Cost Volume [Project page] - Submitted by Deqing Sun (NVIDIA)</small>				
2 ProFlow_ROB	1	5	1	3
<small>Submitted by Daniel Maurer (University of Stuttgart)</small>				
3 LFNNet_ROB	6	1	5	4
<small>Submitted by Anonymous</small>				
4 AugFNG_ROB	8	3	3	2
<small>Submitted by Anonymous</small>				
4 FF++_ROB	3	4	4	5
<small>FlowFields++: Accurate Optical Flow Correspondences Meet Robust Interpolation [Project page] - Submitted by René Schuster (DFKI)</small>				
6 DMF_ROB	4	7	6	7
<small>DeepFlow: Large displacement optical flow with deep matching [Project page] - Submitted by Alexander Brock (HCI)</small>				
6 ResPWCR_ROB	5	6	7	6
<small>Submitted by Anonymous</small>				
8 WOLF_ROB	7	8	8	8
<small>Submitted by Anonymous</small>				

# Visual results on KITTI video sequence

Tensorflow  
code



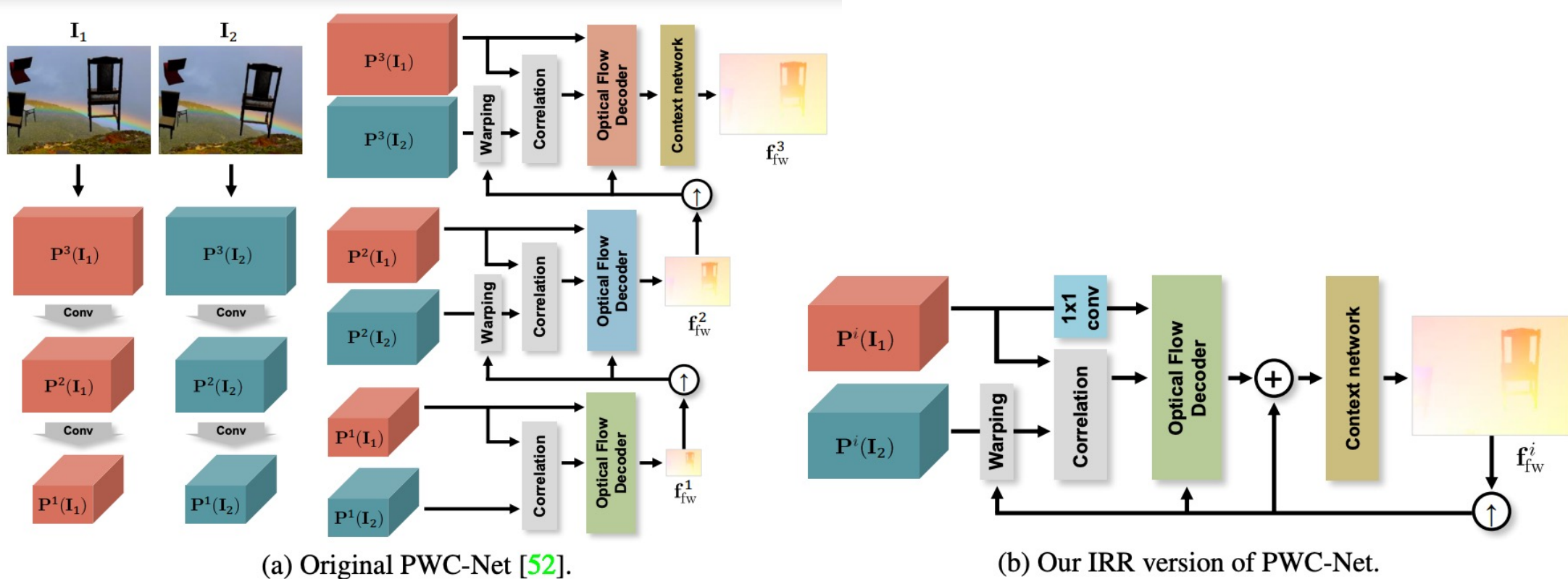
Caffe &  
PyTorch  
code



35 fps for  
Sintel  
(1024x448)  
resolution on  
NVIDIA Pascal  
TitanX

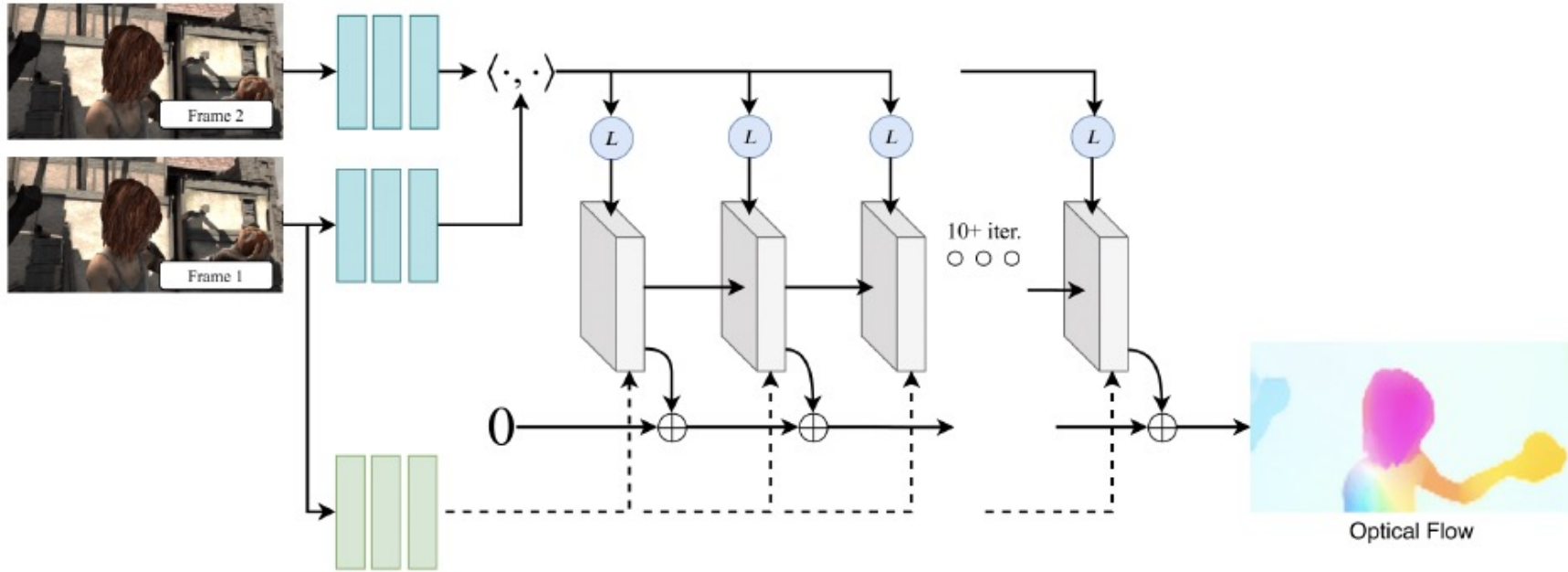
# Improvement: Iterative Residual Refinement (IRR)

[Hur and Roth CVPR 2019]



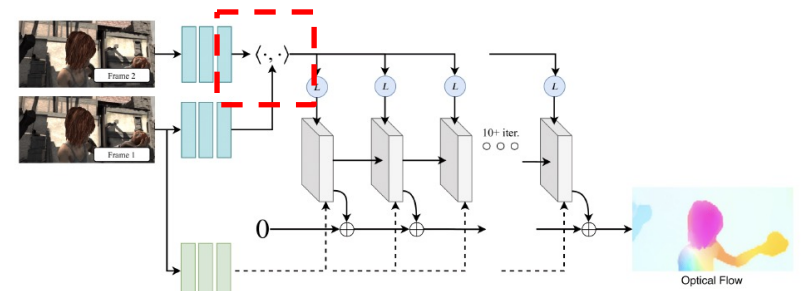
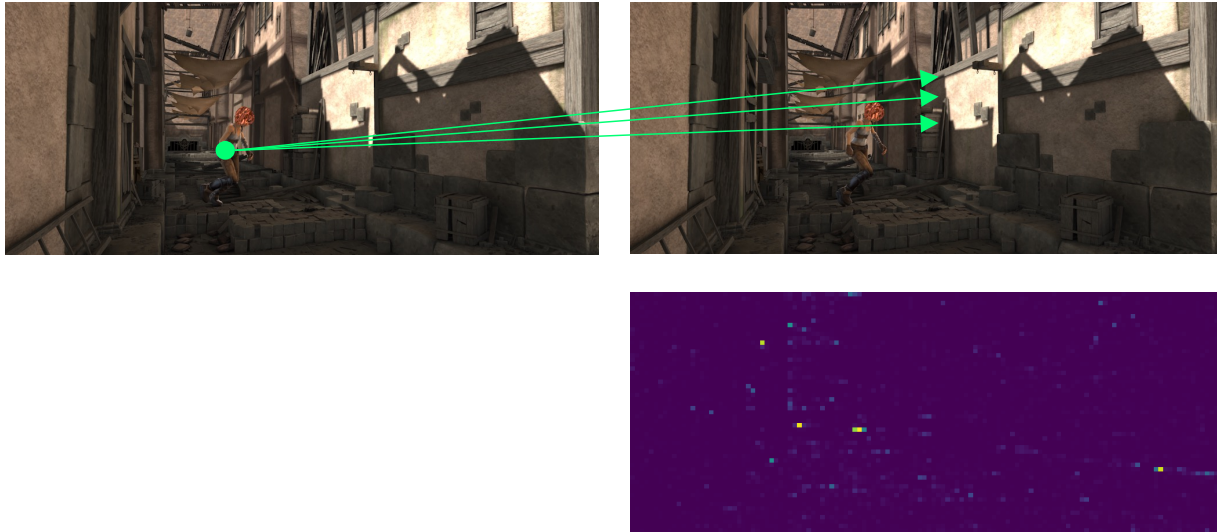
# RAFT: Recurrent All-pairs Field Transforms

[Teed and Deng ECCV 2020 **Best paper**]



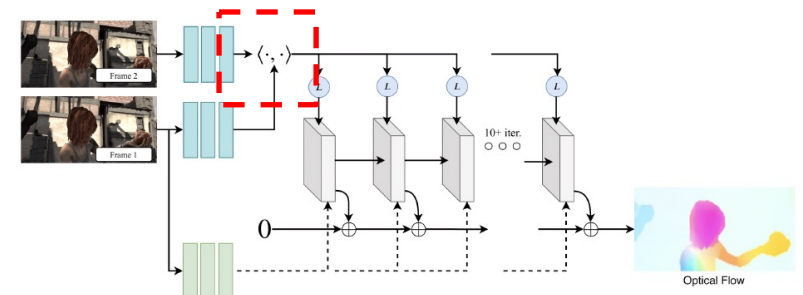
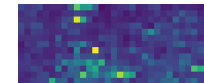
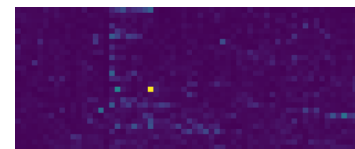
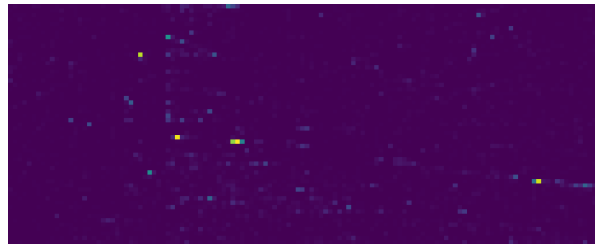
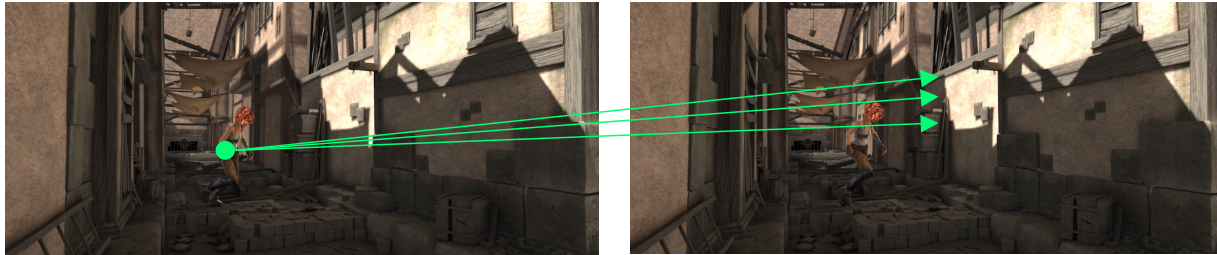
# All-pairs visual similarity (cost volume)

Inner product/correlation between features



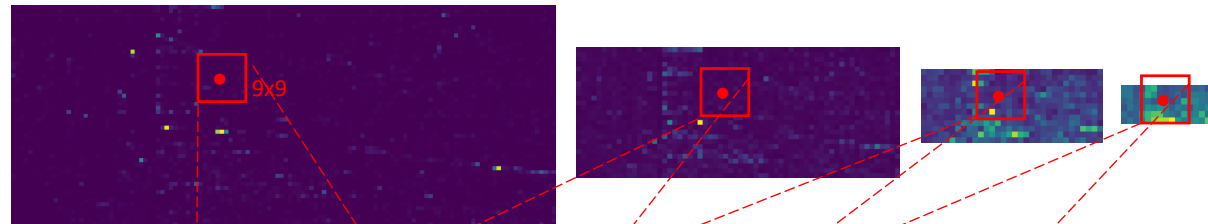
# Cost volume pyramid

## Spatial pooling

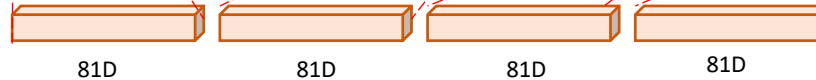


# Look up cost volume

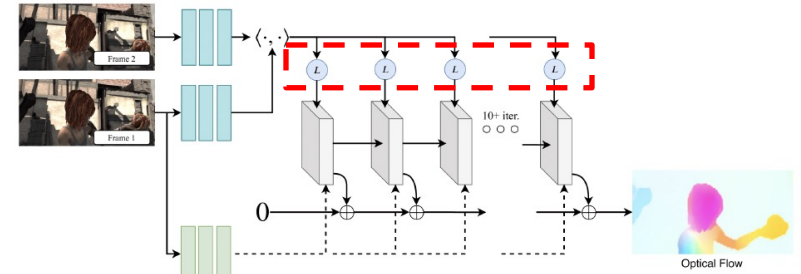
Retrieve using current motion estimates



retrieved feature vector:

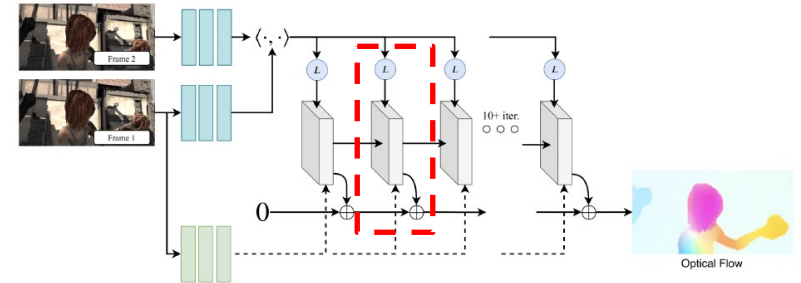
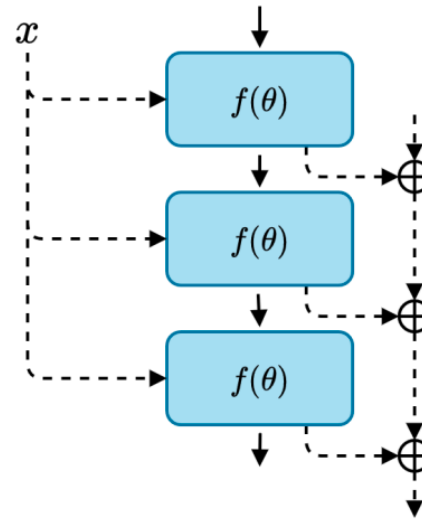
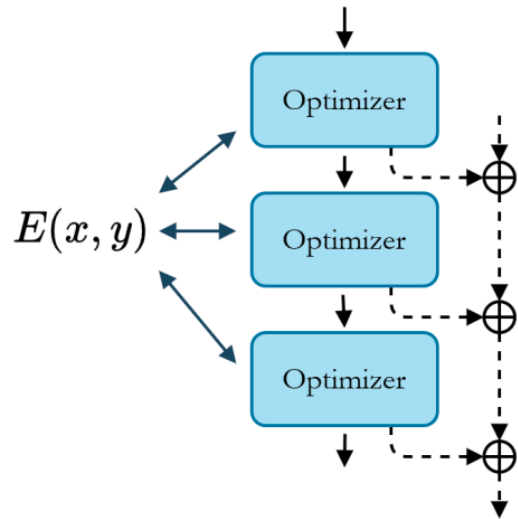


cues on how good the current flow is and where are better similarities



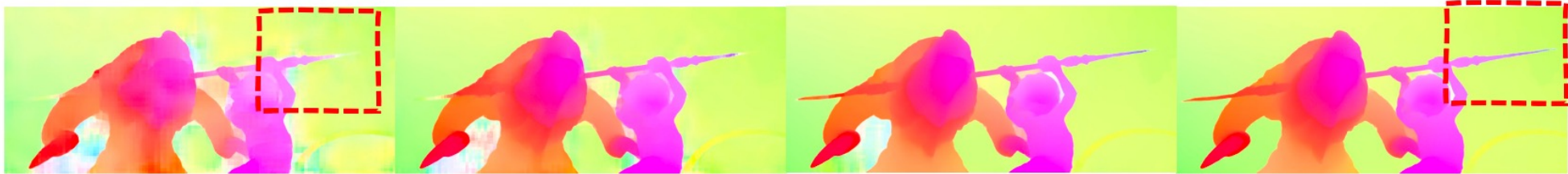
# Recurrent update

Like classical optimization algorithms





# Recurrent update



1 Iteration

2 Iterations

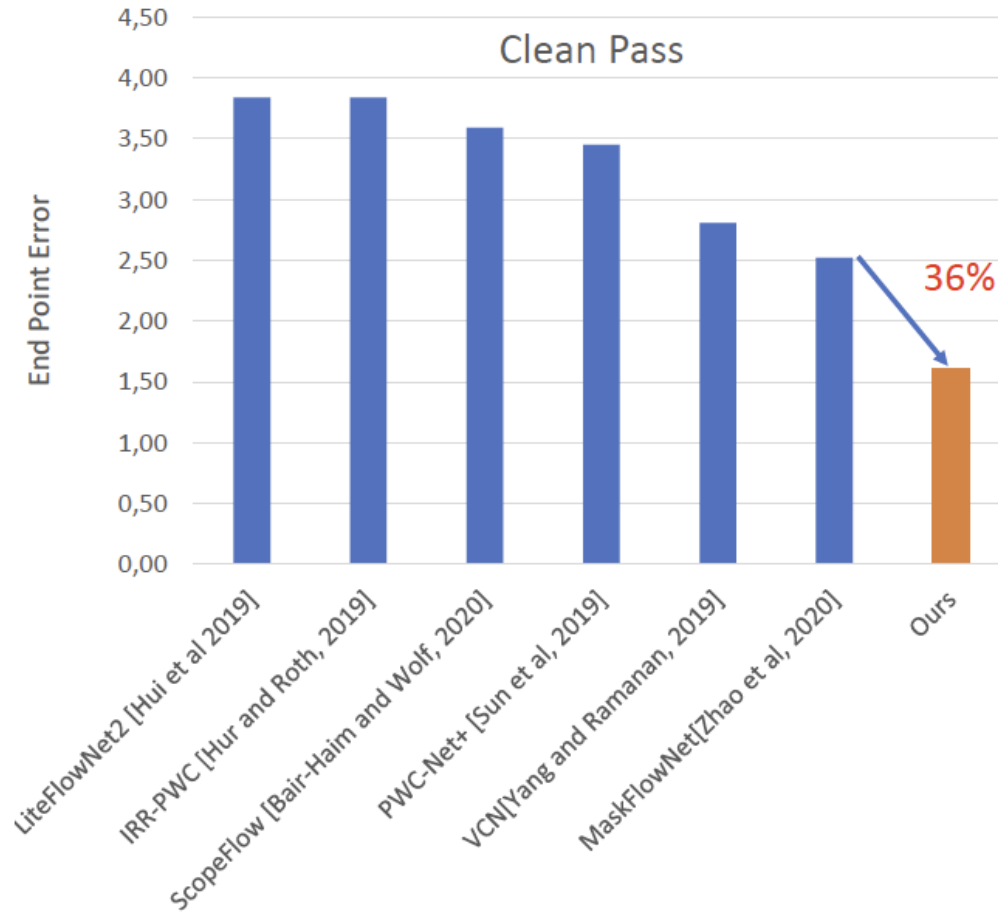
5 Iterations

32 Iterations

# Significant improvement over prior art

[Teed and Deng ECCV 2020 **Best paper**]

## Sintel Results

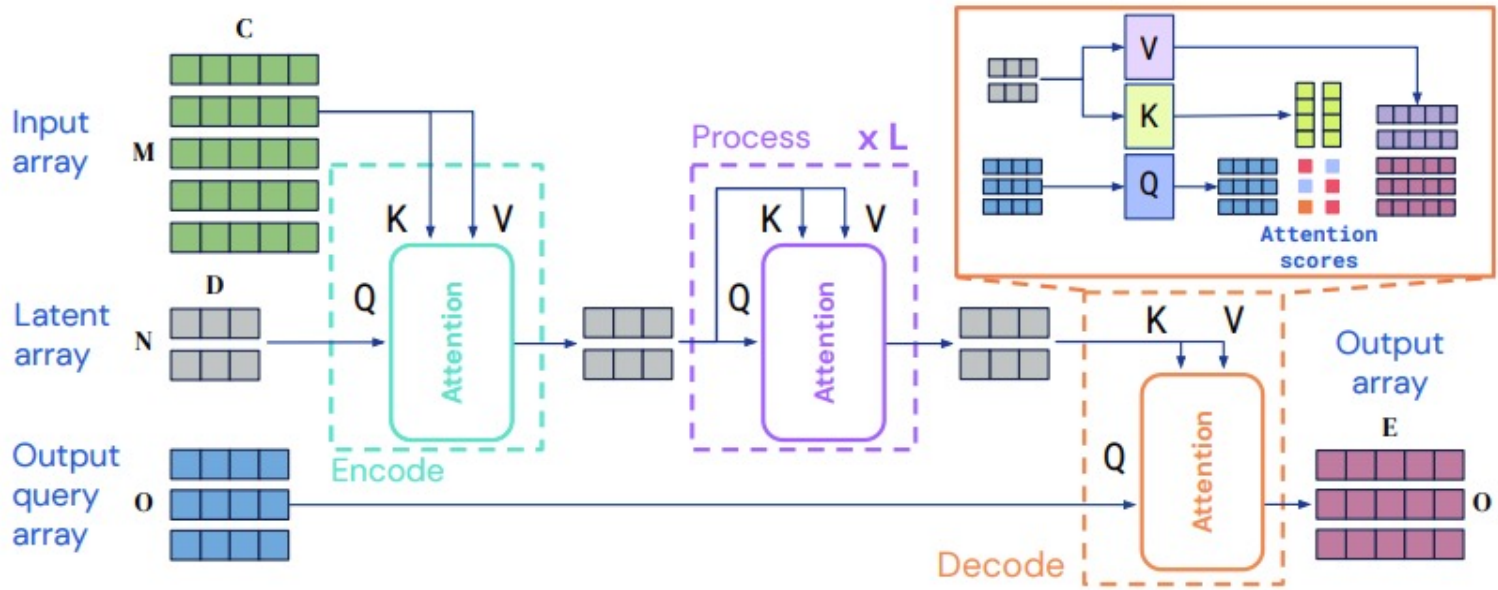
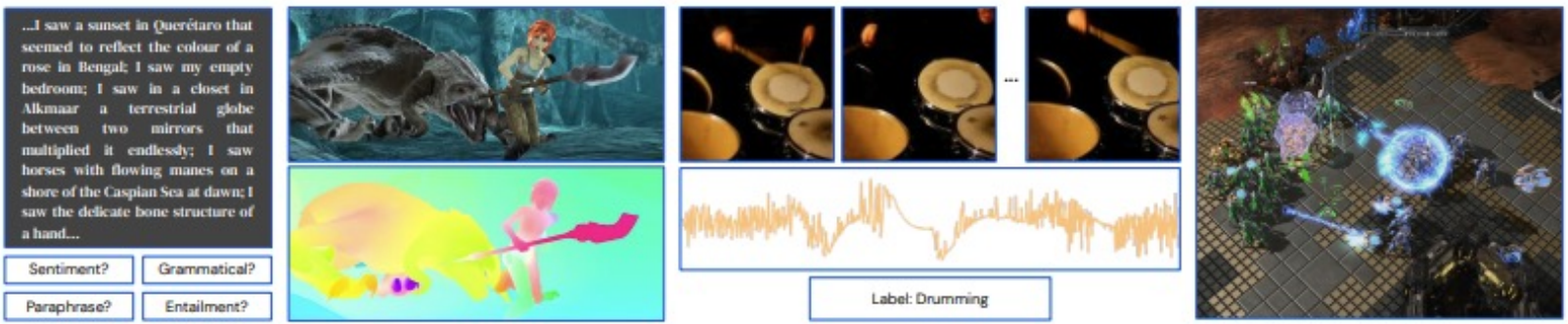


# Visual results on Davis (real-world)

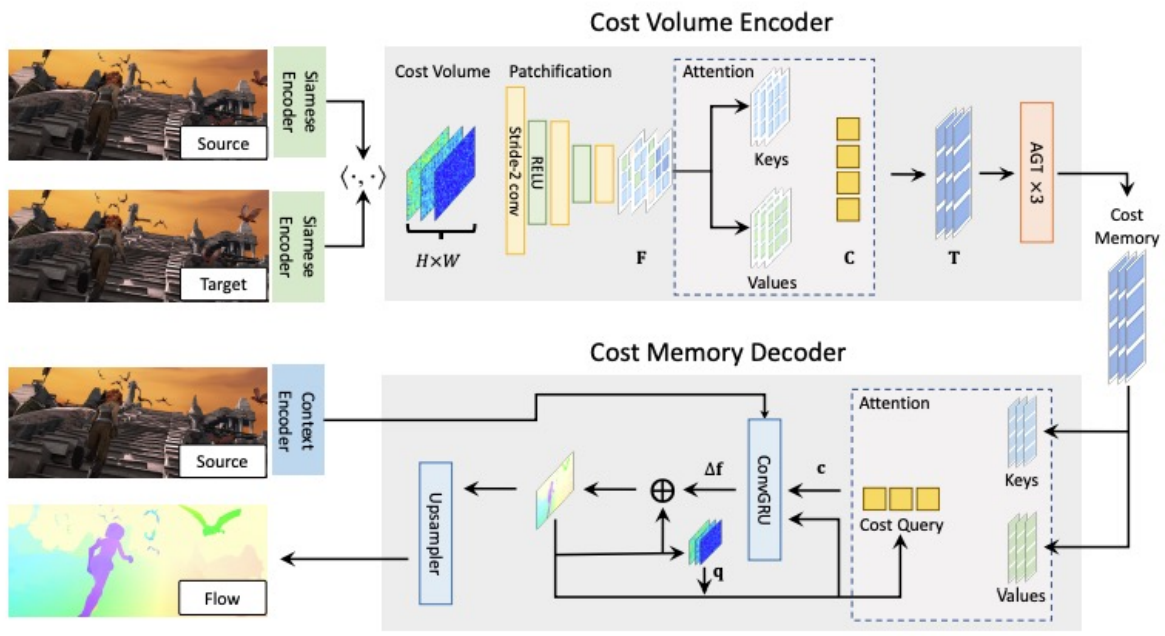


# Recent development: Attention/transformer

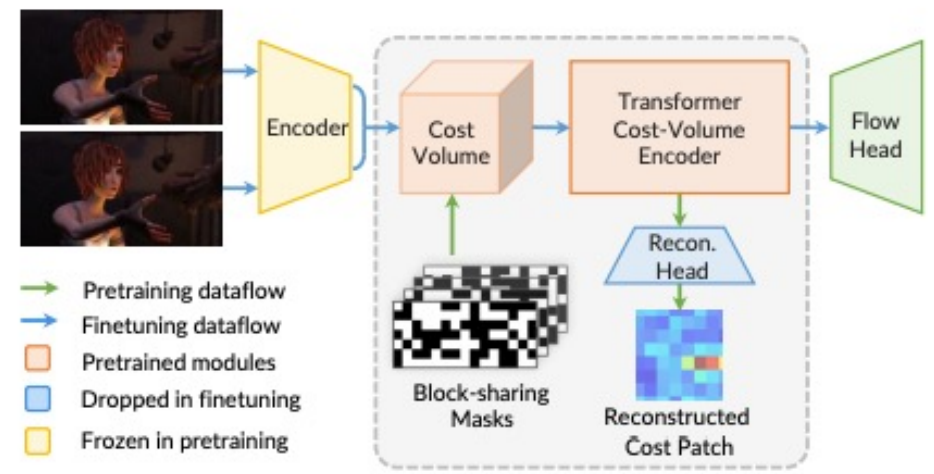
Perceiver IO: A General Architecture for Structured Inputs & Outputs. ICLR '22



# Flowformer: A transformer architecture for optical flow



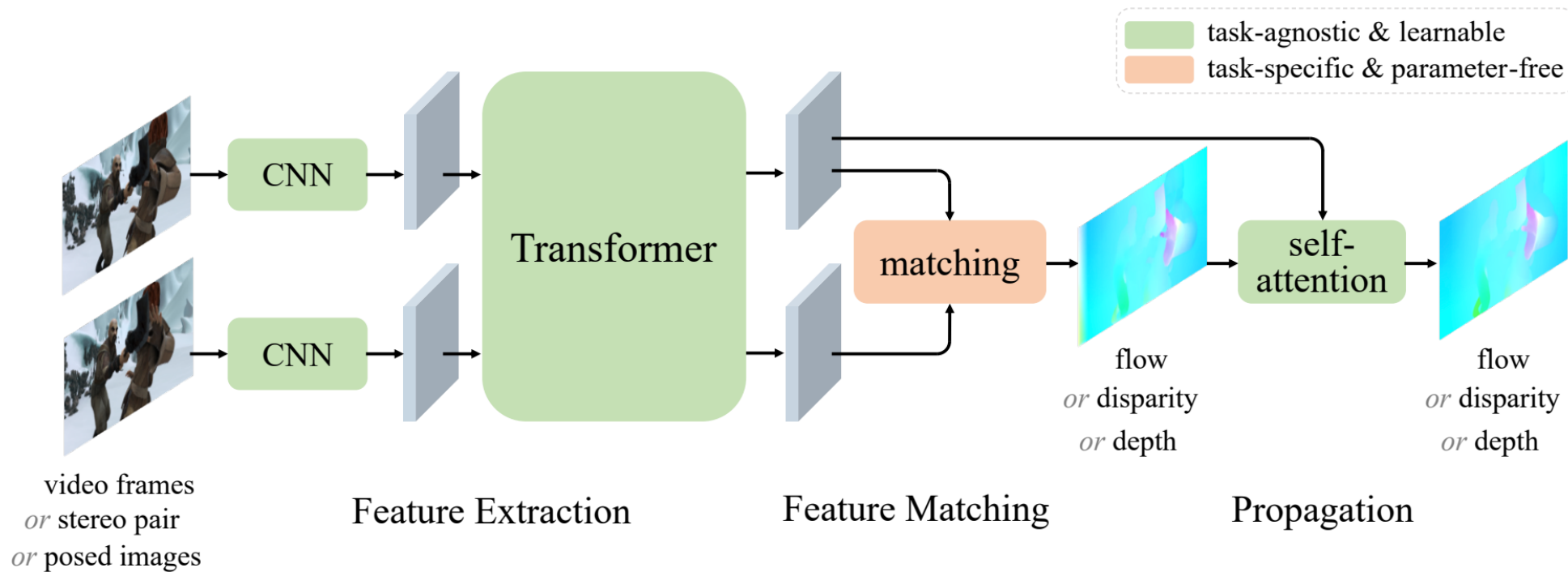
Flowformer [Huang *et al.* ECCV '22]



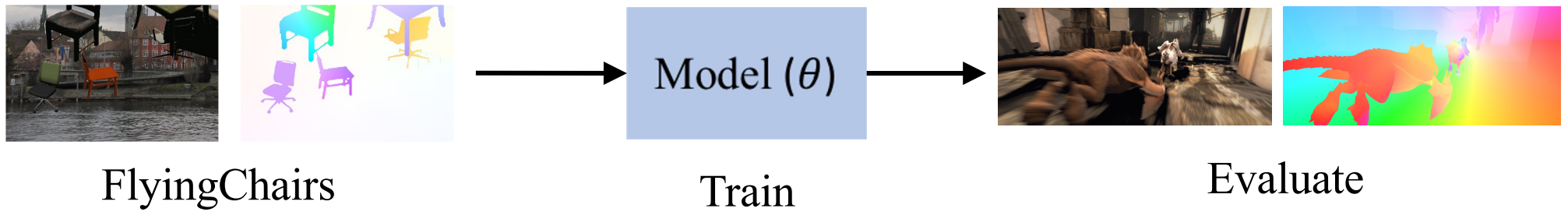
Flowformer++ [Shi *et al.* arXiv '23]

# GM-Flow: Unifying flow, stereo and depth estimation

[Xu *et al.* arXiv '23]



# Is architecture all we need?

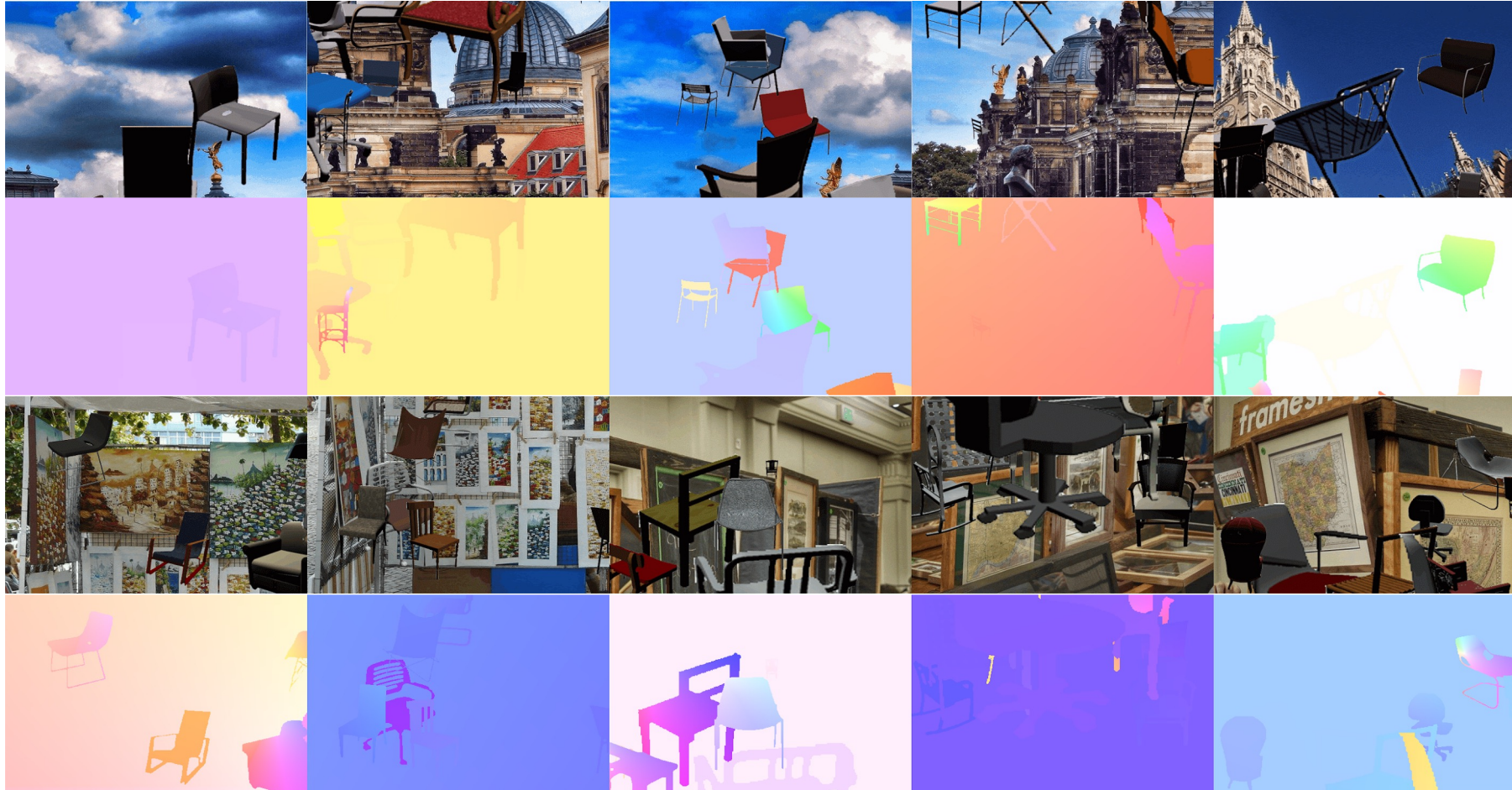


# Learning Data



# “FlyingChairs” manually designed in 2015

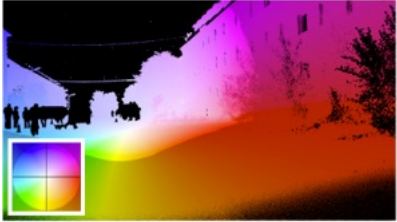
[Dosovitskiy *et al.* ICCV'15]



# “FlyingChairs” more effective than later datasets



“FlyingThings3D” [Mayer et al. 2016]



HD1K [Kondermann et al. 2016]



Multi-human [Ranjan et al. 2020]



Virtual KITTI [Gaidon et al. 2016]



Playing for benchmark [Richter et al. 2017]



Refresh [Lv et al. 2018]

...

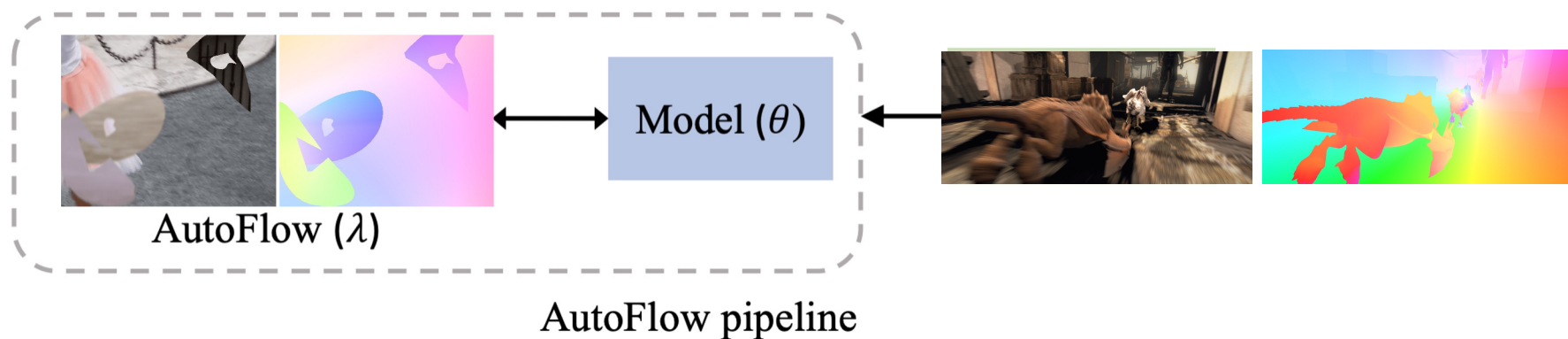
# Many important and interesting questions unanswered

- How realistic should the rendering be?
- Why does FlyingChairs work so well?
- Should we carefully match the motion statistics of Sintel?
- Are thin structures/fine motion details of FlyingChairs critical?

...

# What is the objective for rendering training data?

Optimize the performance of a network on a target dataset



Jointly render data and train model

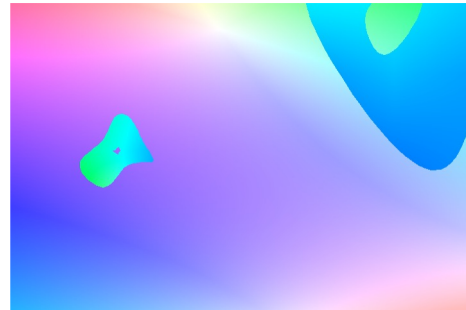
# How simple can the rendering be?

Start from the simplest rendering pipeline: 2D layered model

Images



Motion

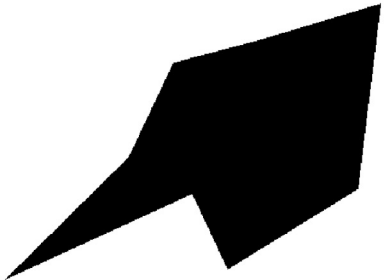


Background

+ 1 foreground objects

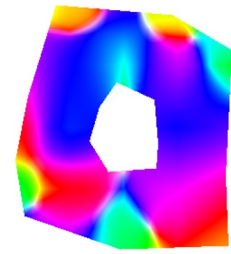
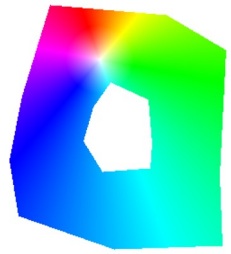
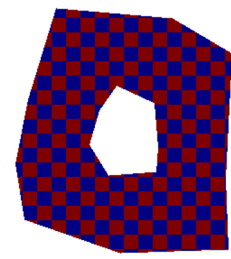
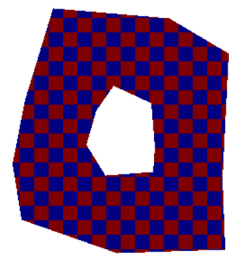
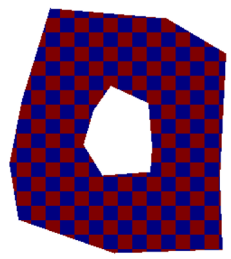
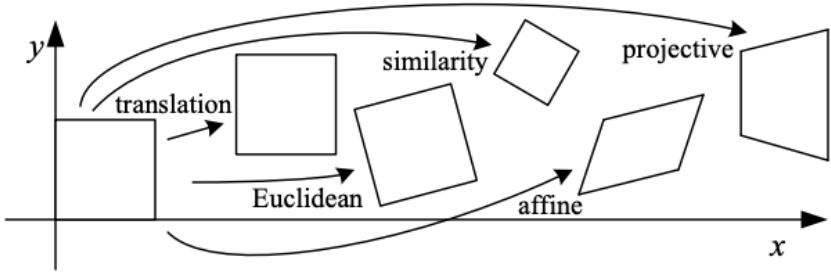
+ 2 foreground objects

# Modeling foreground shapes



Random polygons

# Modeling motion



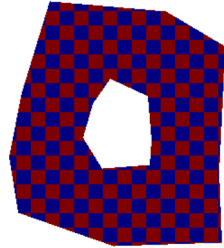
Affine

Perspective

Bilinear grid warp

# Visual effects

Motion blur



Fog

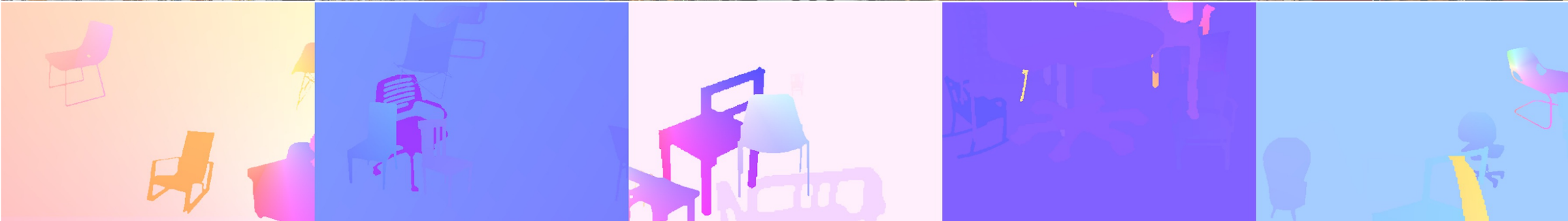
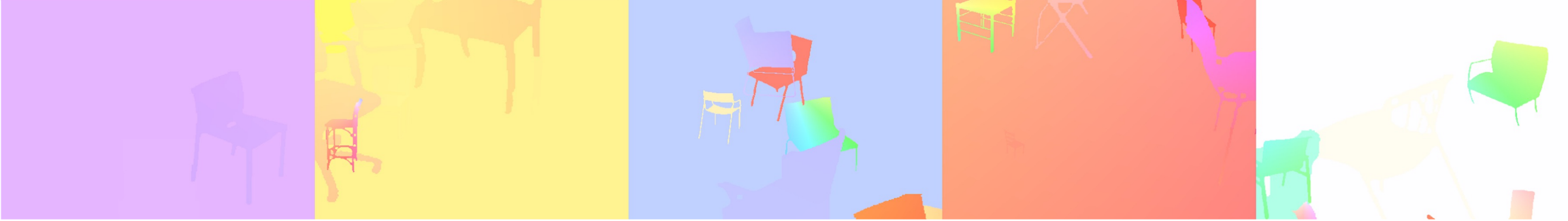




# AutoFlow



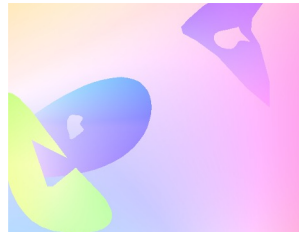
# FlyingChairs



# Results of pre-training (on training set)

Avg. end-point error (AEPE) ↓

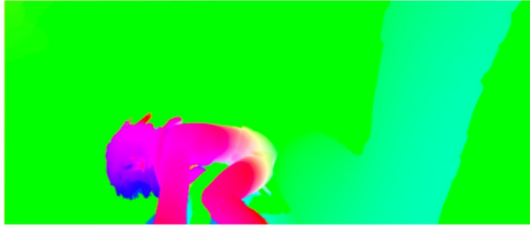
Model	Dataset	Sintel.clean	Sintel.final	KITTI
PWC-Net	FlyingChairs	3.27	4.42	11.43
	Chairs → Things	2.39	3.90	9.81
	AutoFlow	<b>2.17</b>	<b>2.91</b>	<b>5.76</b>
RAFT	FlyingChairs	2.27	3.76	7.63
	Chairs → Things	<b>1.68</b>	2.80	5.92
	AutoFlow	1.95	<b>2.57</b>	<b>4.23</b>



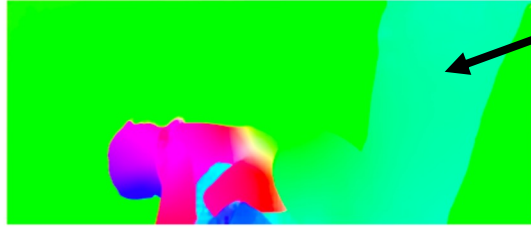
# AutoFlow vs. FlyingChairs



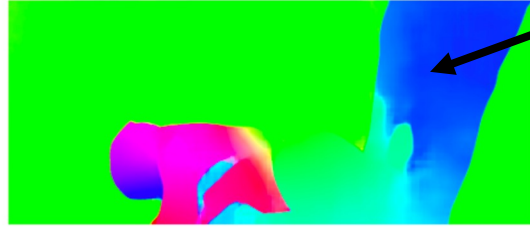
Frame 29 of final Ambush\_4



Ground truth



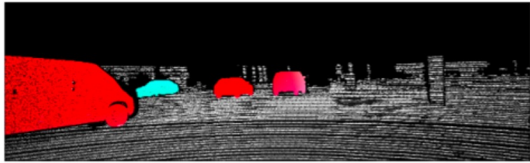
AutoFlow



FlyingChairs



Frame 000094 of KITTI training



Ground truth



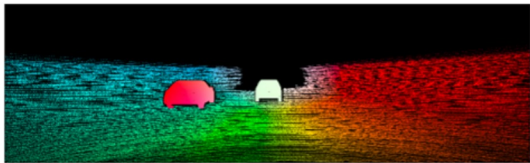
AutoFlow



FlyingChairs



Frame 000104 of KITTI training



Ground truth

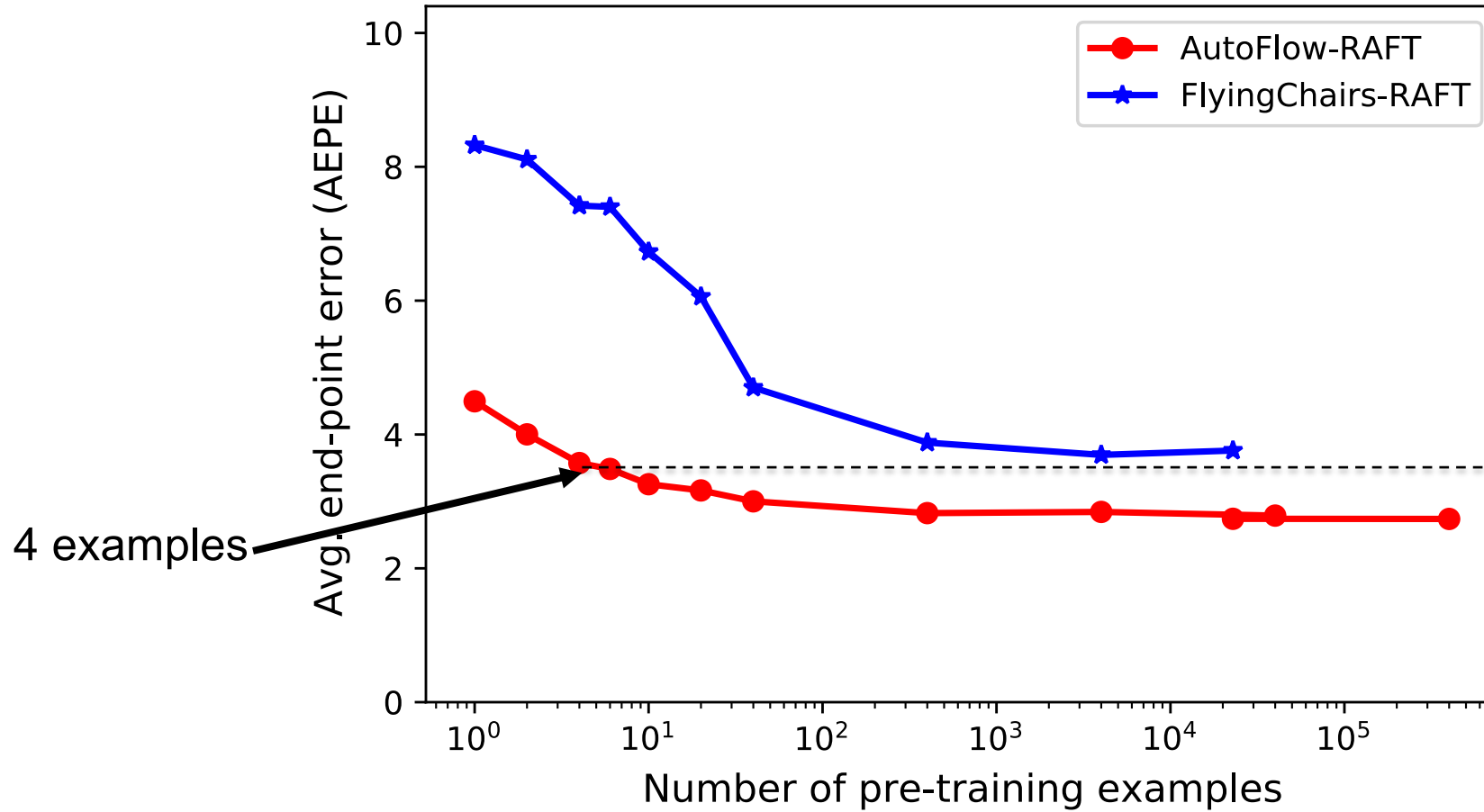


AutoFlow



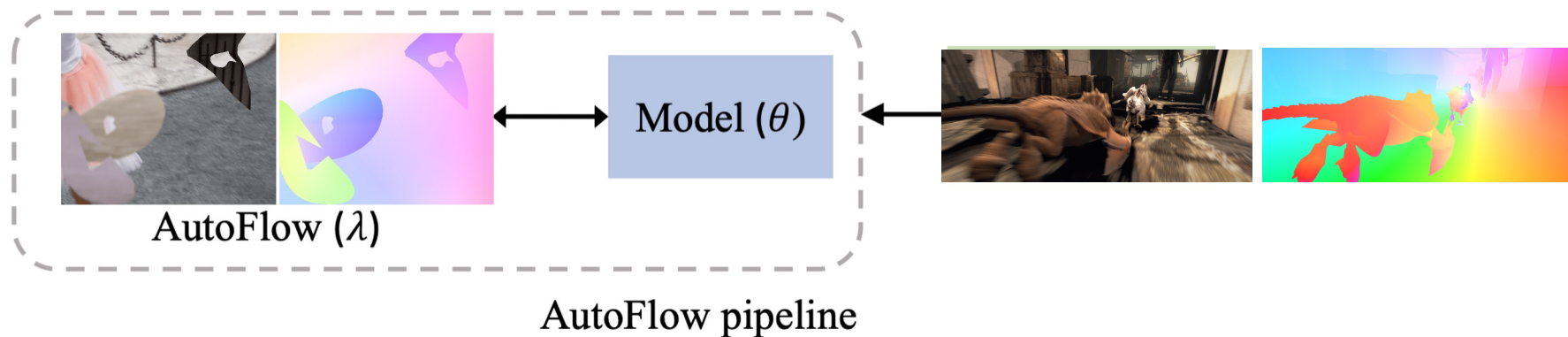
FlyingChairs

# Number of training examples



# Are architectures and data all we need?

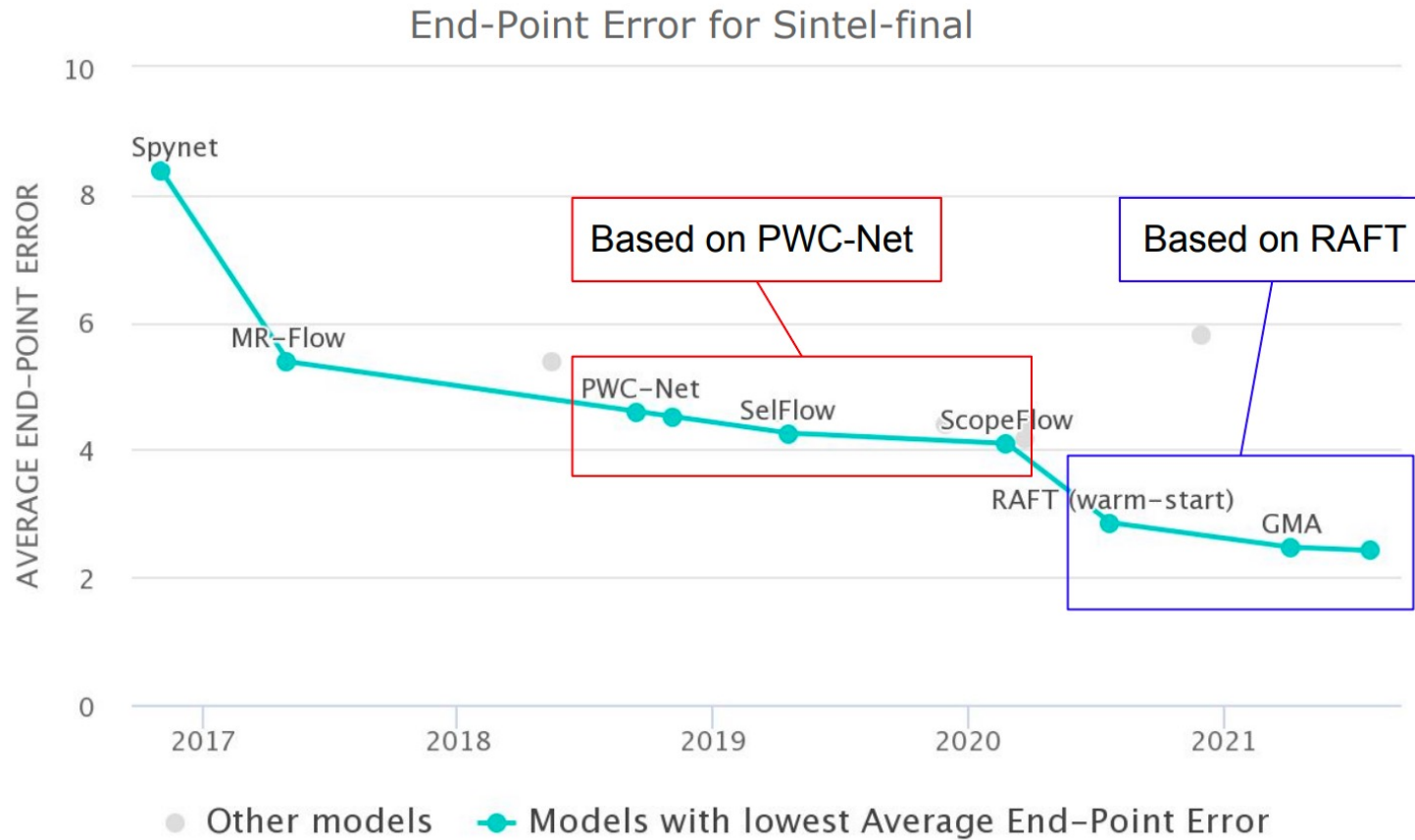
Optimize the performance of a network on a target dataset



Jointly render data and train model

**The devils are in the **training** details**

# Rapid progress on optical flow architectures



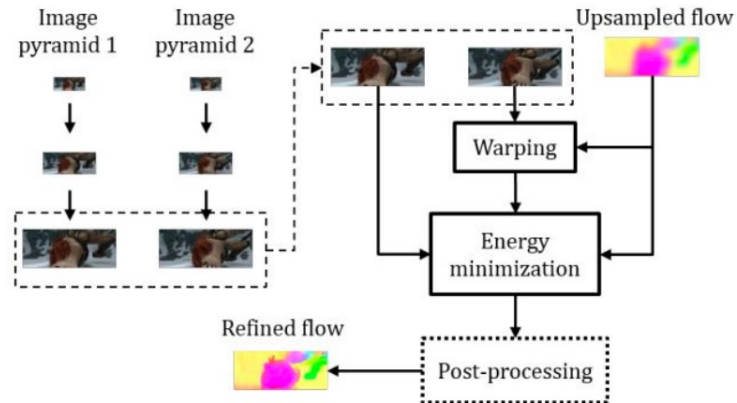


# Differences in architecture

## PWC-Net (2018)

Inspired by classical optical flow:

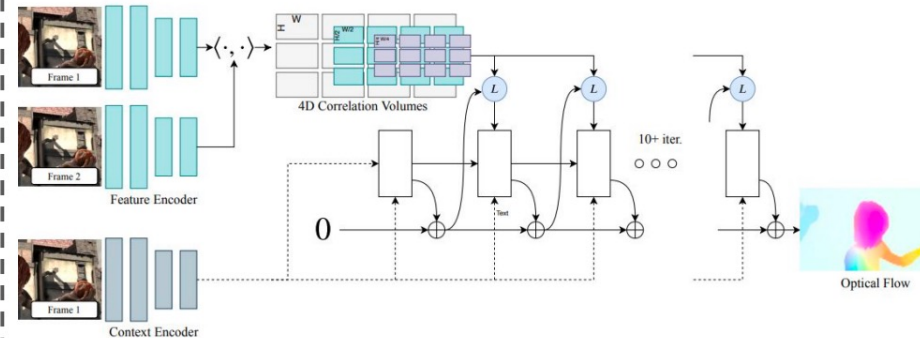
- Image pyramid
- Cost volume warping
- Multiscale loss function



## RAFT (2020)

Introduced new network elements:

- Produces flow at a single level
- Multi-scale all-pairs cost volume
- Recurrent update operator
- Upsample module



# Differences in training techniques

Lots of differences in training schemes

Training Details	PWC-Net (2018)	RAFT (2020)
Optimizer	Adam	AdamW
Learning rate schedule	Multi-Step LR	One Cycle LR
Gradient clipping	No	Yes
Training time	Moderate	Low
Augmentation	Color + Spatial	Different Color + Spatial

# Imbalanced focus on architecture/modeling

Optical flow work on Modeling	Optical flow work on Training
<i>Learning to estimate hidden motions with global motion aggregation. ICCV21</i>	
<i>High-resolution optical flow from 1d attention and correlation. ICCV21</i>	
<i>Separable flow: Learning motion cost volumes for optical flow estimation. ICCV21</i>	???
<i>Learning optical flow with adaptive graph reasoning. arxiv22</i>	
<i>Csflow: Learning optical flow via cross strip correlation for autonomous driving. arxiv22</i>	
<i>Pyramid recurrent all-pairs field transforms for optical flow estimation in robust vision challenge. arxiv22</i>	

# Disentangling architecture and training for optical flow

[Sun, Herrmann, *et al.* ECCV '22]

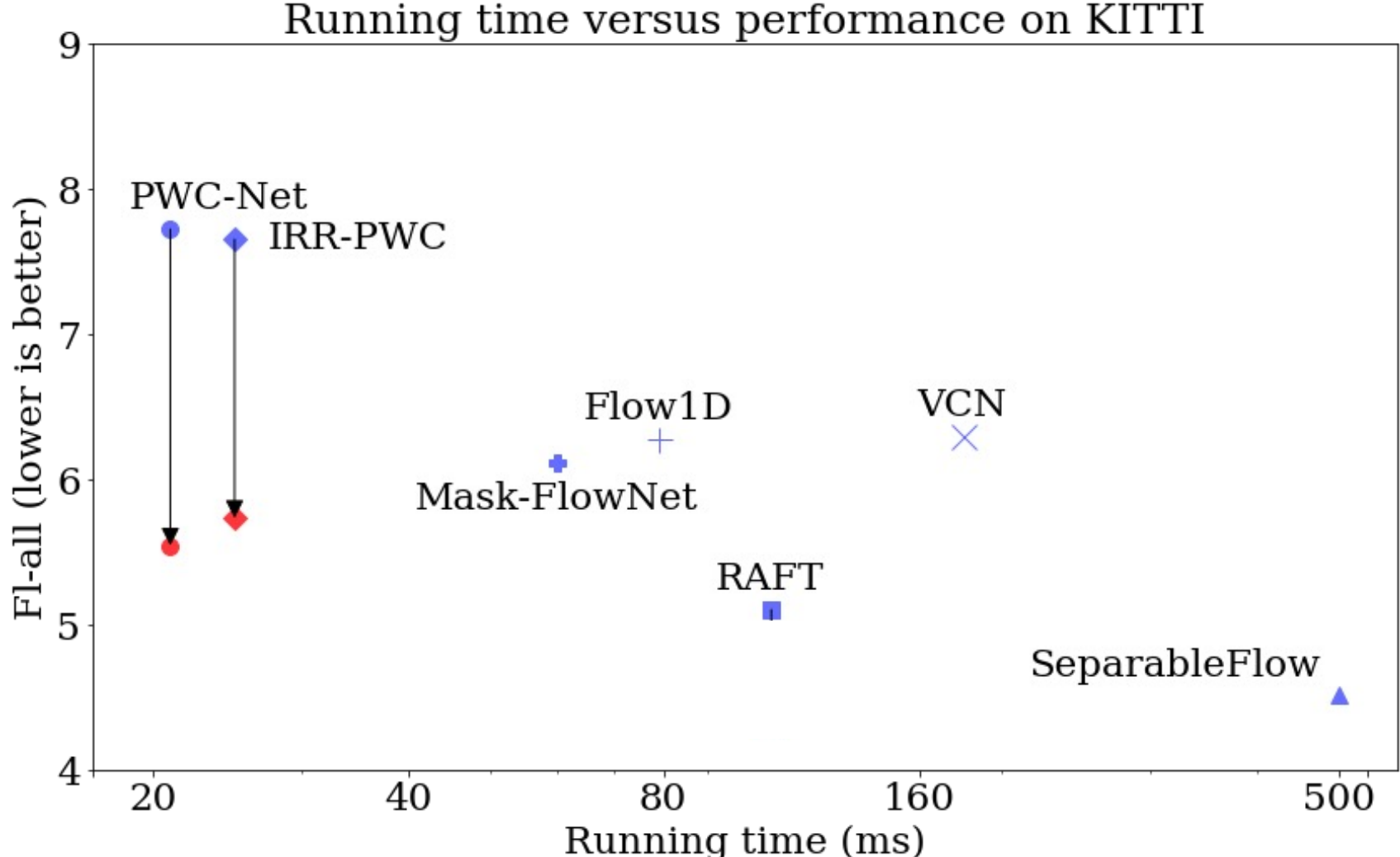
## ***Our goals***

- Understand the effect of modern training on performance and improve it.

## ***Our approach***

- We apply a modern training scheme to 3 prominent models
  - PWC-Net (2018)
  - IRR-PWC (2019)
  - RAFT (2020)
- We perform a thorough ablation study on pre-training and fine-tuning

# Better training significantly improves performance



# “Old” vs. “new” PWC-Net

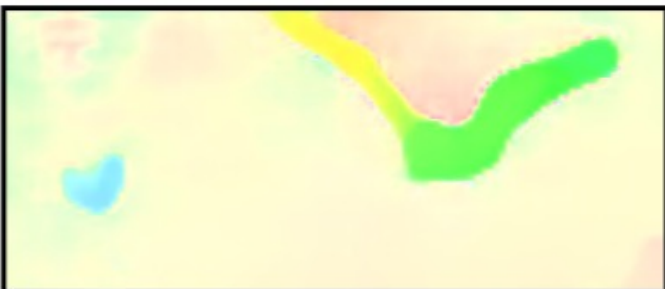
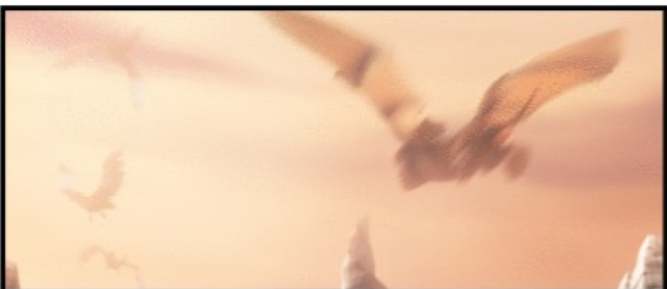
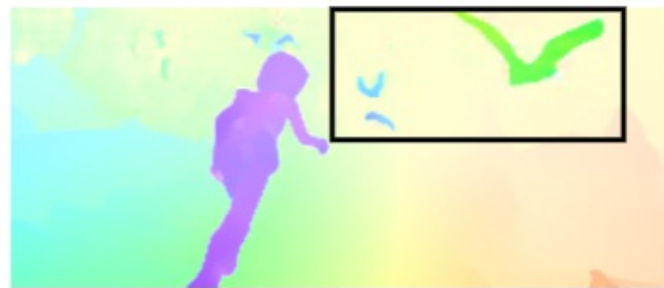
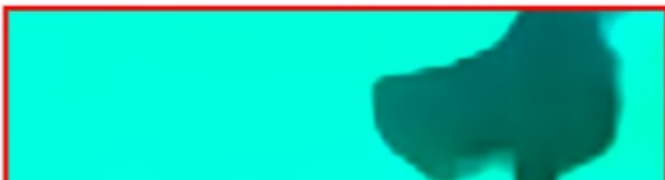
First frame



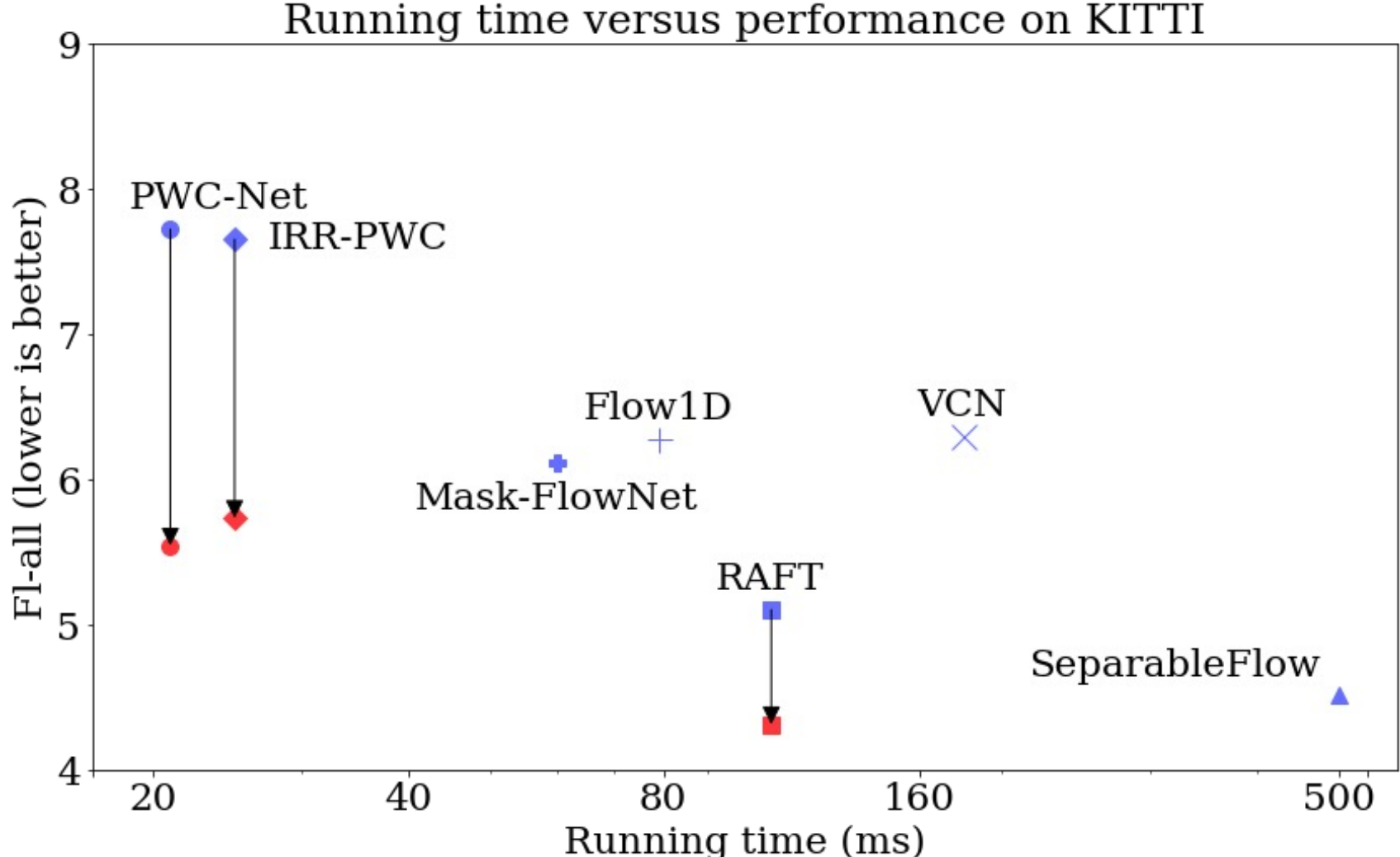
PWC-Net (Orig.)



PWC-Net (Ours)



# Better training significantly improves performance



# Compared with state of the art (April 2023)

	Method	Setting	Code	FI-bg	FI-fg	FI-all	Density	Runtime	Environment	Compare
1	<a href="#">CamLiRAFT</a>		<a href="#">code</a>	2.08 %	7.37 %	2.96 %	100.00 %	1 s	GPU @ 2.5 Ghz (Python + C/C++)	<input type="checkbox"/>
H. Liu, T. Lu, Y. Xu, J. Liu and L. Wang: <a href="#">Learning Optical Flow and Scene Flow with Bidirectional Camera-LiDAR Fusion</a> . arXiv preprint arXiv:2303.12017 2023.										
2	<a href="#">CamLiFlow</a>		<a href="#">code</a>	2.31 %	7.04 %	3.10 %	100.00 %	1.2 s	GPU @ 2.5 Ghz (Python + C/C++)	<input type="checkbox"/>
H. Liu, T. Lu, Y. Xu, J. Liu, W. Li and L. Chen: <a href="#">CamLiFlow: Bidirectional Camera-LiDAR Fusion for Joint Optical Flow and Scene Flow Estimation</a> . CVPR 2022.										
3	<a href="#">CamLiRAFT-NR</a>		<a href="#">code</a>	2.76 %	6.78 %	3.43 %	100.00 %	1 s	GPU @ 2.5 Ghz (Python + C/C++)	<input type="checkbox"/>
H. Liu, T. Lu, Y. Xu, J. Liu and L. Wang: <a href="#">Learning Optical Flow and Scene Flow with Bidirectional Camera-LiDAR Fusion</a> . arXiv preprint arXiv:2303.12017 2023.										
4	<a href="#">M-FUSE</a>		<a href="#">code</a>	2.56 %	7.47 %	3.46 %	100.00 %	1.3 s	GPU	<input type="checkbox"/>
L. Mehl, A. Jahedi, J. Schmalfluss and A. Bruhn: <a href="#">M-FUSE: Multi-frame Joint Optical Flow Estimation</a> . Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) 2023.										
5	<a href="#">RigidMask+ISF</a>		<a href="#">code</a>	2.63 %	7.85 %	3.50 %	100.00 %	3.3 s	GPU @ 2.5 Ghz (Python)	<input type="checkbox"/>
G. Yang and D. Ramanan: <a href="#">Learning to Segment Rigid Motions from Two Frames</a> . CVPR 2021.										
6	<a href="#">ScaleRAFT3D</a>			2.37 %	9.26 %	3.51 %	100.00 %	1 s	1 core @ 2.5 Ghz (C/C++)	<input type="checkbox"/>
7	<a href="#">TPCV+RAFT3D</a>			2.48 %	10.19 %	3.76 %	100.00 %	0.2 s	1 core @ 2.5 Ghz (C/C++)	<input type="checkbox"/>
8	<a href="#">RAFT3D+mscv</a>			2.48 %	10.23 %	3.77 %	100.00 %	0.2 s	1 core @ 2.5 Ghz (C/C++)	<input type="checkbox"/>
9	<a href="#">RAFT-it+ RVC</a>		<a href="#">code</a>	3.62 %	5.33 %	3.90 %	100.00 %	0.14 s	1 core @ 2.5 Ghz (Python)	<input type="checkbox"/>
D. Sun, C. Herrmann, F. Reda, M. Rubinstein, D. Fleet and W. Freeman: <a href="#">Disentangling Architecture and Training for Optical Flow</a> . ECCV 2022.										
10	<a href="#">RAFT-OCTC</a>			3.72 %	5.39 %	4.00 %	100.00 %	0.2 s	GPU @ 2.5 Ghz (Python)	<input type="checkbox"/>
J. Jeong, J. Lin, F. Porikli and N. Kwak: <a href="#">Imposing Consistency for Optical Flow Estimation (Qualcomm AI Research)</a> . CVPR 2022.										
11	<a href="#">RCA-Flow</a>			3.67 %	6.25 %	4.10 %	100.00 %	0.16 s	1 core @ 2.5 Ghz (Python)	<input type="checkbox"/>
12	<a href="#">SF2SE3</a>		<a href="#">code</a>	3.17 %	8.79 %	4.11 %	100.00 %	2.7 s	GPU @ >3.5 Ghz (Python)	<input type="checkbox"/>
L. Sommer, P. Schröppel and T. Brox: <a href="#">SF2SE3: Clustering Scene Flow into SE (3)-Motions via Proposal and Selection</a> . DAGM German Conference on Pattern Recognition 2022.										
13	<a href="#">RAFT-CF-CE-PL3</a>			3.80 %	5.65 %	4.11 %	100.00 %	0.05 s	GPU @ 2.5 Ghz (Python)	<input type="checkbox"/>

Using stereo input





# What matters?

We iteratively build our training and show that each step improves Sintel.F

- *Dataset matters*  
FlyingChairs → AutoFlow  
for IRR 5.09 → 4.11 [-19%]  
for RAFT 4.03 → 3.36 [-17%]
- *Gradient clipping matters*  
NoGC → YesGC  
for IRR 4.11 → 3.29 [-20%]  
for RAFT 3.36 → 3.20 [-5%]
- *Training schedule matters*  
Piecewise → OneCycle  
for IRR 3.29 → 2.93 [-11%]  
for RAFT 3.20 → 2.75 [-14%]
- *Training iterations matter*  
standard → 4-10 times more  
for IRR 2.93 → 2.76 [-6%]  
for RAFT 2.75 → 2.41 [-12%]

# Tradeoff between accuracy and speed/memory

	Inference Time (msec)↓			Peak Memory (GB)↓		
	1024×448	Full HD	4K	1024 × 448	Full HD	4K
PWC-Net	20.61	28.77	63.31	1.478	2.886	7.610
IRR-PWC	24.71	33.67	57.59	1.435	2.902	8.578
RAFT	107.38	499.63	n/a	2.551	9.673	OOM

**Table 6.** Inference time and memory usage for 1024×448, Full HD (1920×1080) and 4K (3840×2160) frame sizes, averaged over 100 runs on an NVIDIA V100 GPU.

# Running on Full HD and 4K images

Inputs



IRR-PWC



PWC-Net and IRR-PWC can run on 4K images, which cause RAFT to OOM



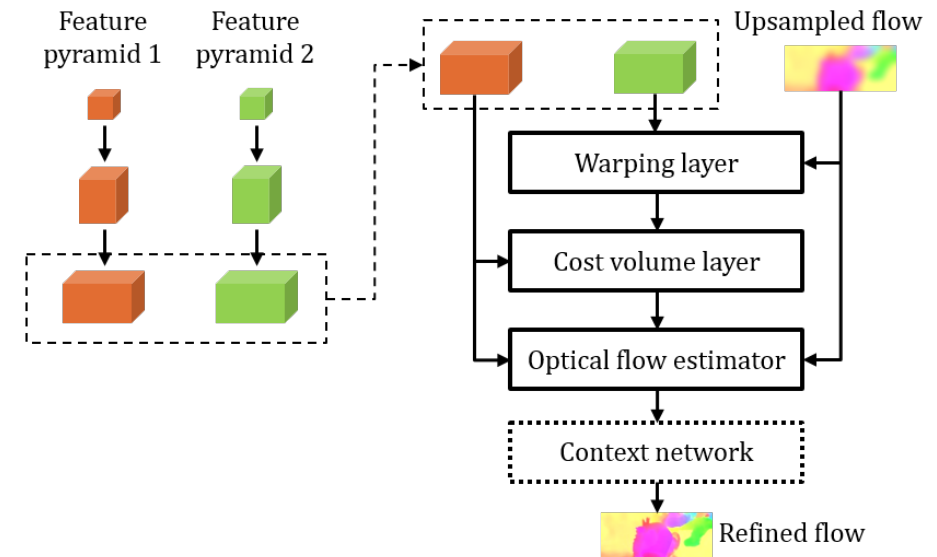
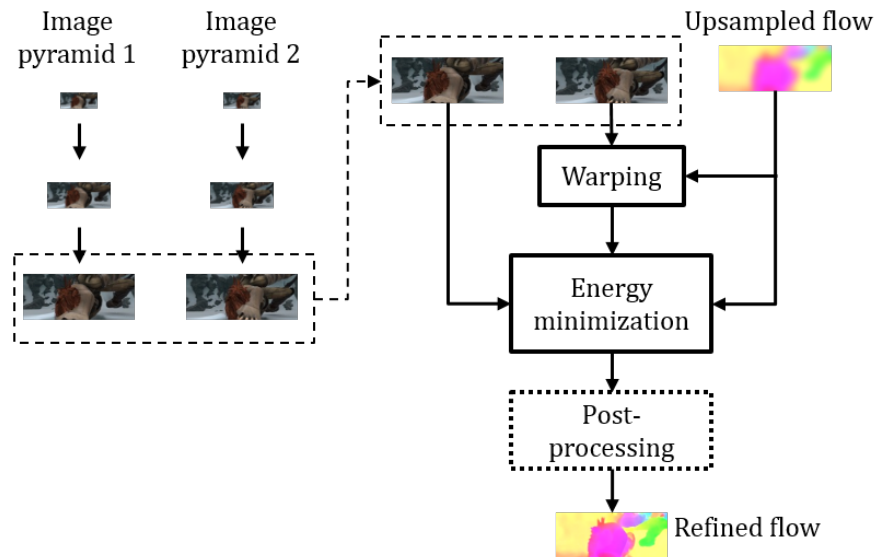
IRR-PWC on 4K input



RAFT on full HD input + x2 upsampling

# Content

- Deep learning-based approach
  - Designing architecture (using domain knowledge)



# Content

- Deep learning-based approach
  - Designing architecture (using domain knowledge)
  - Learning data (matters)



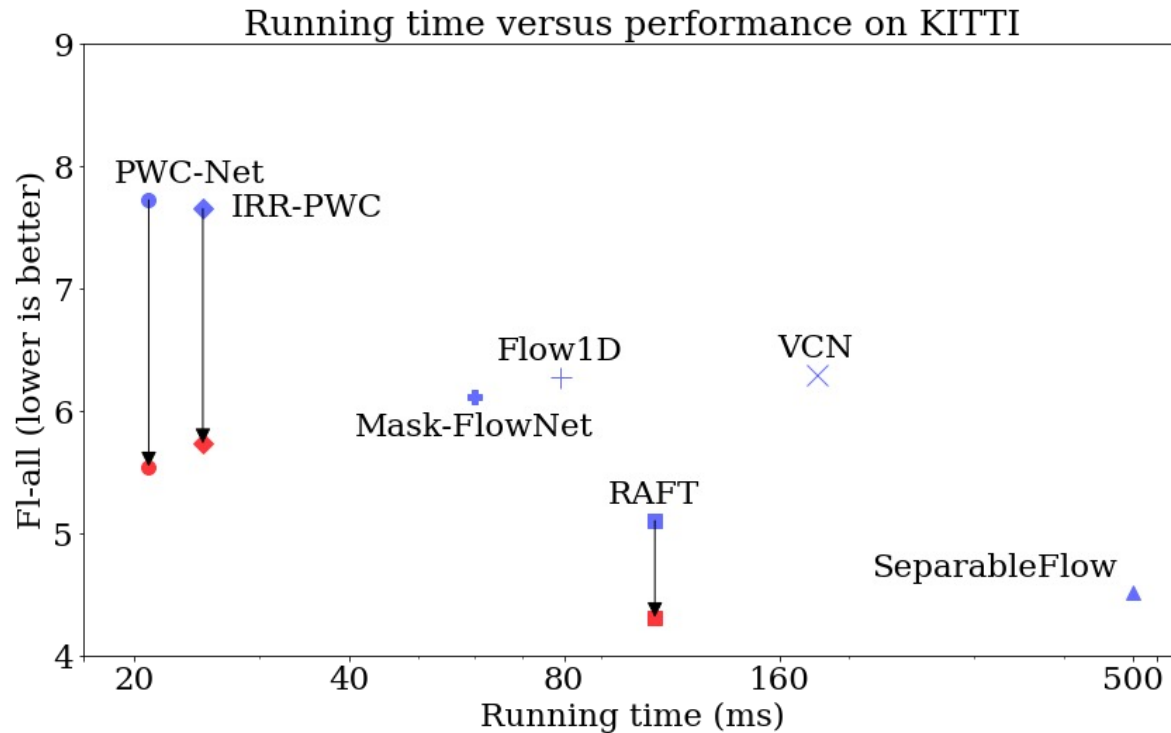
FlyingChairs: Manually designed



AutoFlow: Joint data generation and network training

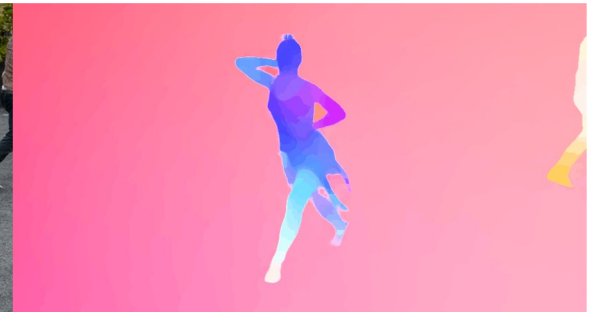
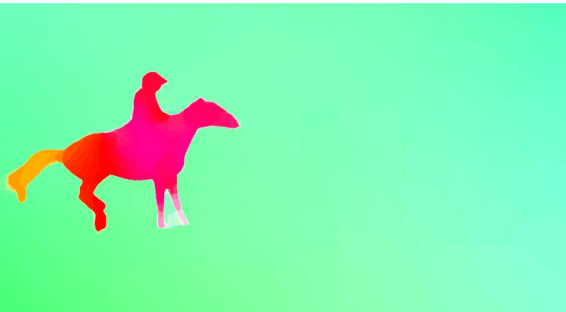
# Content

- Deep learning-based approach
  - Designing architecture (using domain knowledge)
  - Learning data (matters)
  - Evaluating architectures fairly (trade-off in accuracy and speed/memory)





# Results on real-world videos



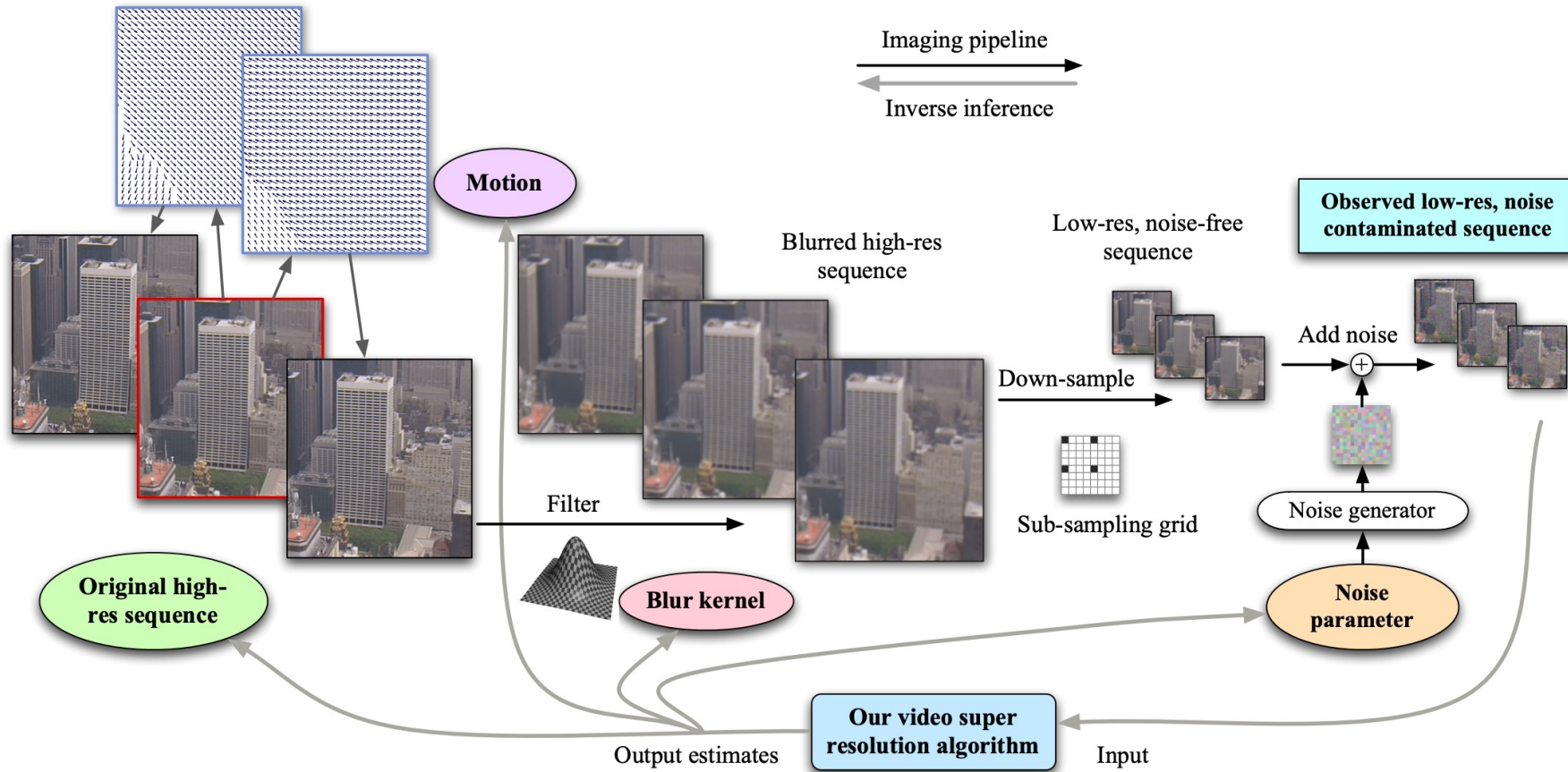
**What is motion for?**

# Video super-resolution

[Liu & Sun CVPR 2011, TPAMI 2014]



# Video super-resolution



# Video frame interpolation



Super SloMo [Jiang, Sun, *et al.* CVPR 2018]  
Incorporated into **NVIDIA NGX** SDK for the Turing GPU.

# Idea: frame synthesis using optical flow

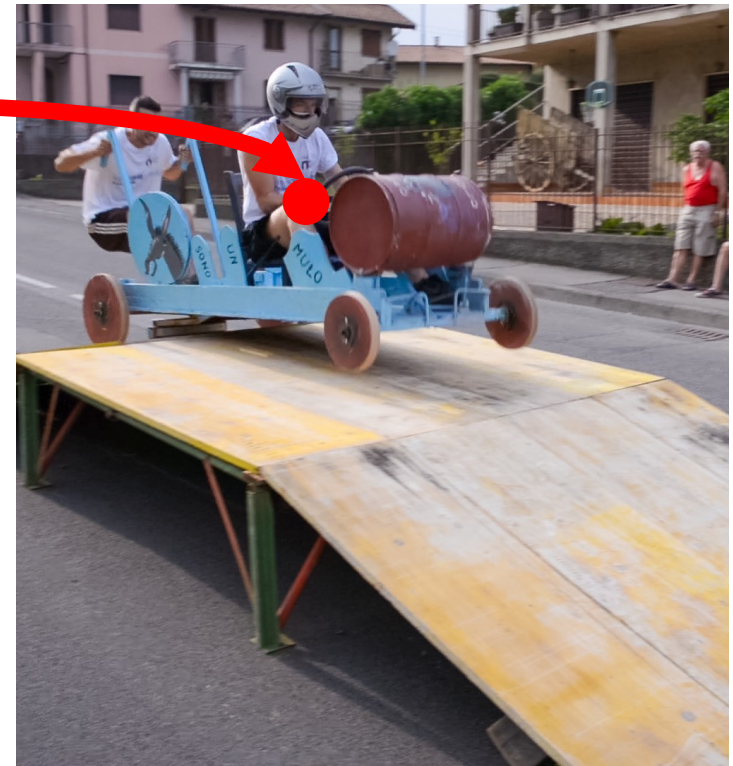
$T = 0$



$T = t \in (0,1)$

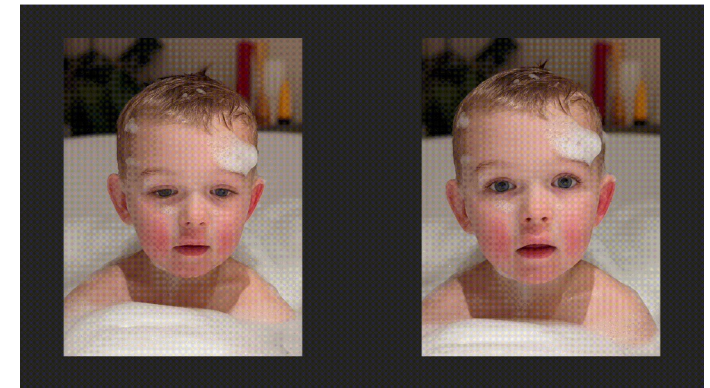
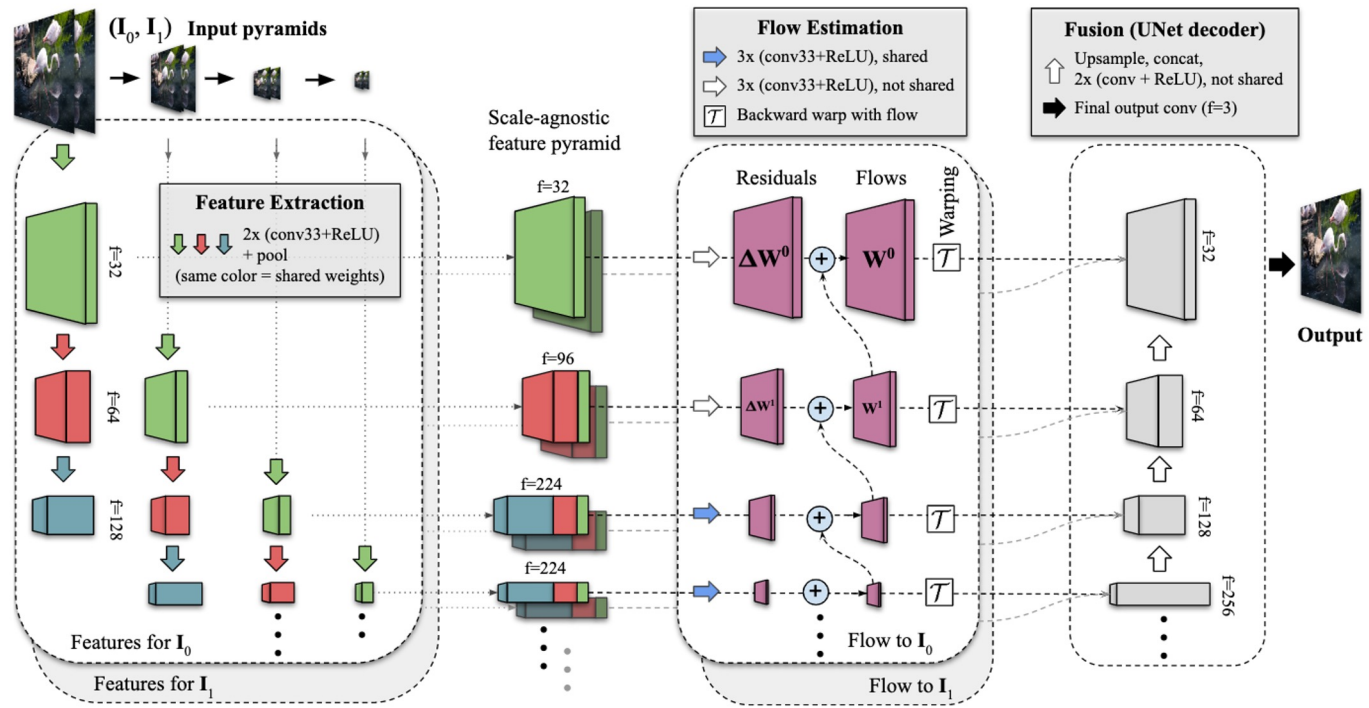


$T = 1$

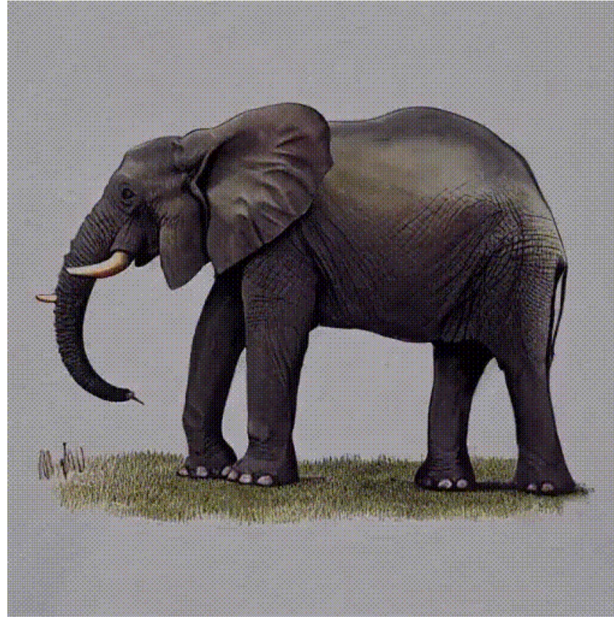
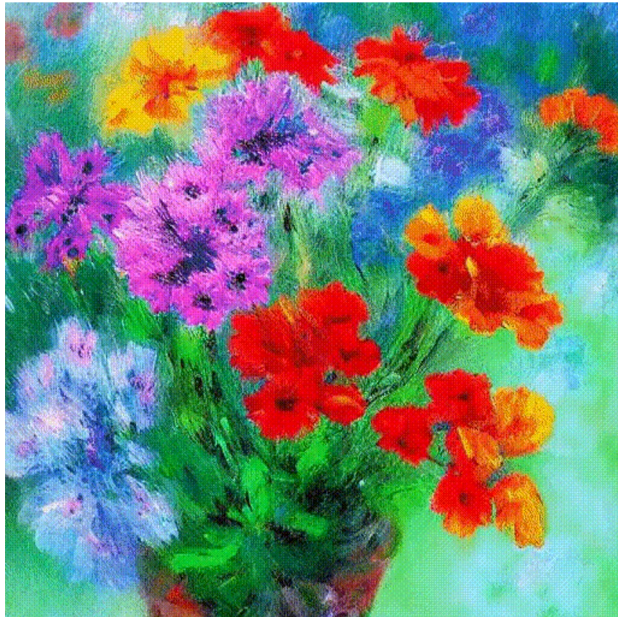


# FILM: Frame Interpolation for Large Motion

[Reda, *et al.* ECCV 2022]



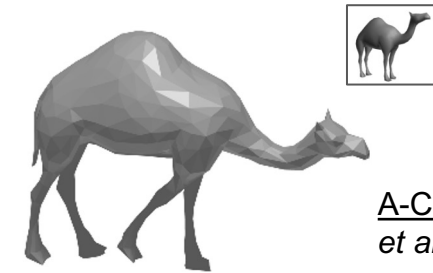
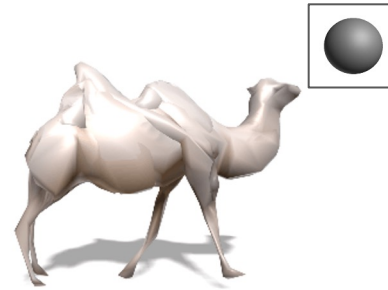
# Stable diffusion + FILM (from twitter)



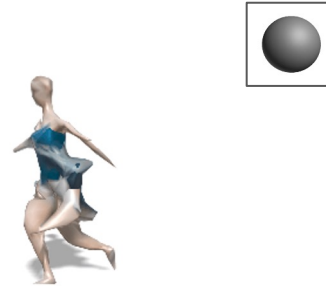


# LASR: Learning Articulated Shape Reconstruction from a monocular video

[Yang, *et al.* CVPR 2021]



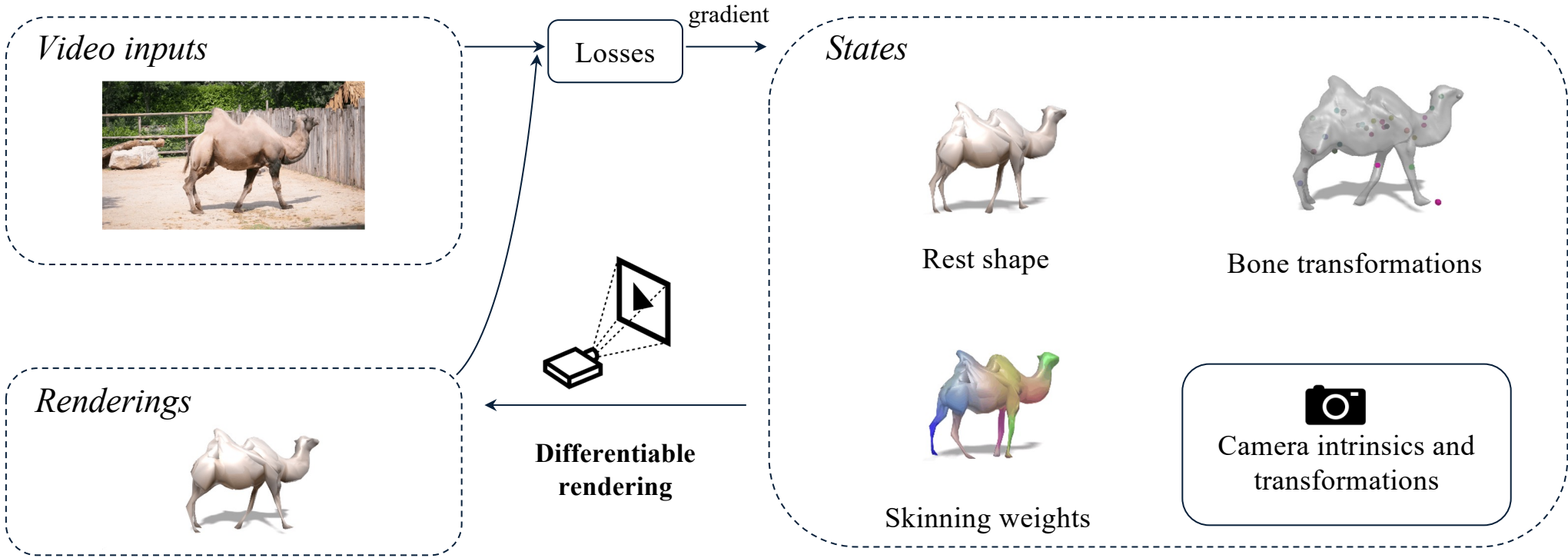
A-CSM: Kulkarni  
*et al.* CVPR 2020



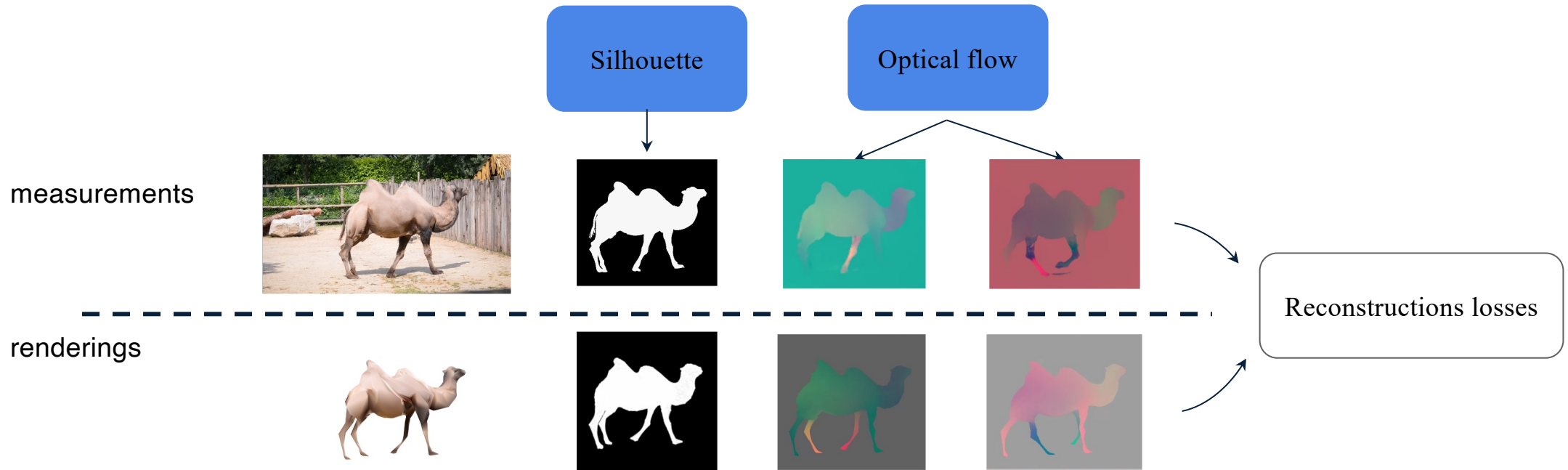
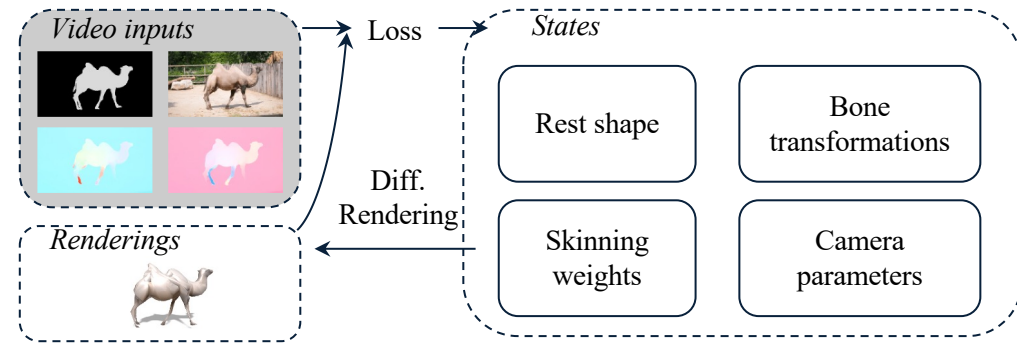
VIBE: Kocabas *et al.* CVPR 2020

Challenge: Solving non-rigid 3D shape from 2D measurements **without template or category prior** is highly *under-constrained*

# Approach: Analysis-by-synthesis



# Supervision from silhouette, flow and pixels



# Reconstructions on more real videos



Input

Reconstructions

Input

Reconstructions

Input

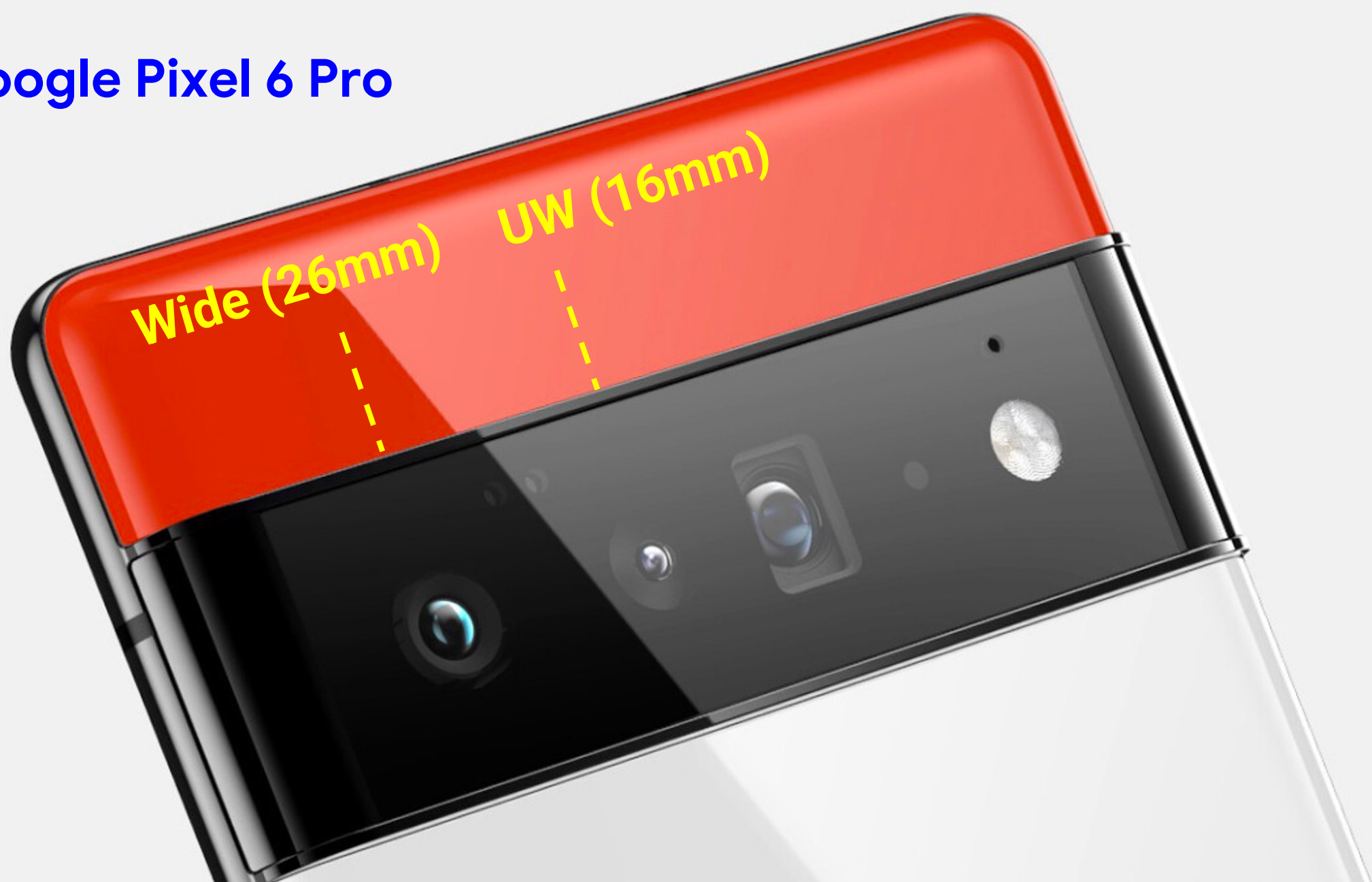
Reconstructions

# Face Unblur for Pixel 6

[Lai *et al.* SIGGRAPH 2022]



# Google Pixel 6 Pro



# Key Idea: Wide + Ultrawide Dual Camera Fusion



Wide (1/120 s)

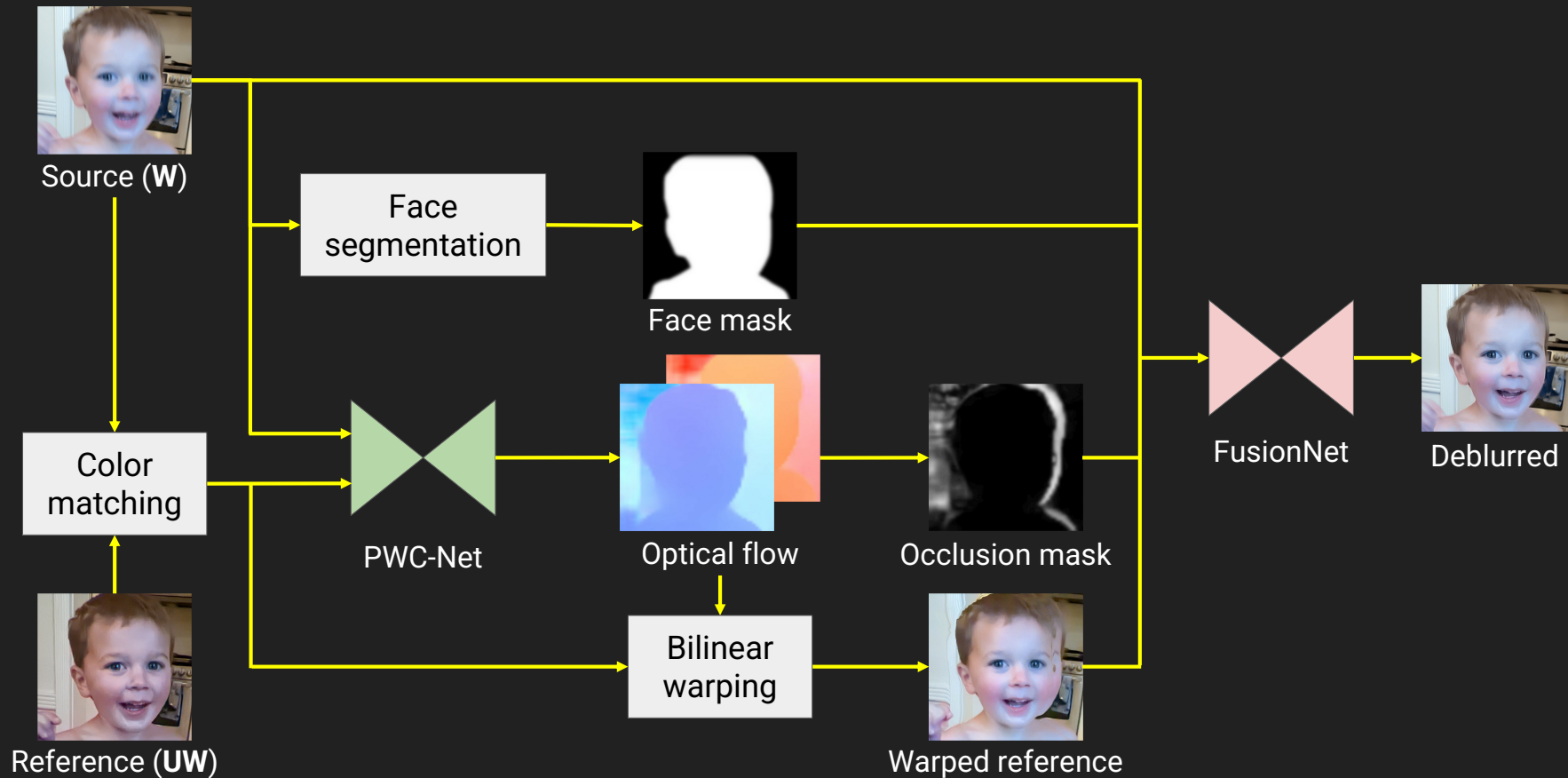


UW (1/480 s)



# Face Unblur

## Alignment and Fusion Algorithm

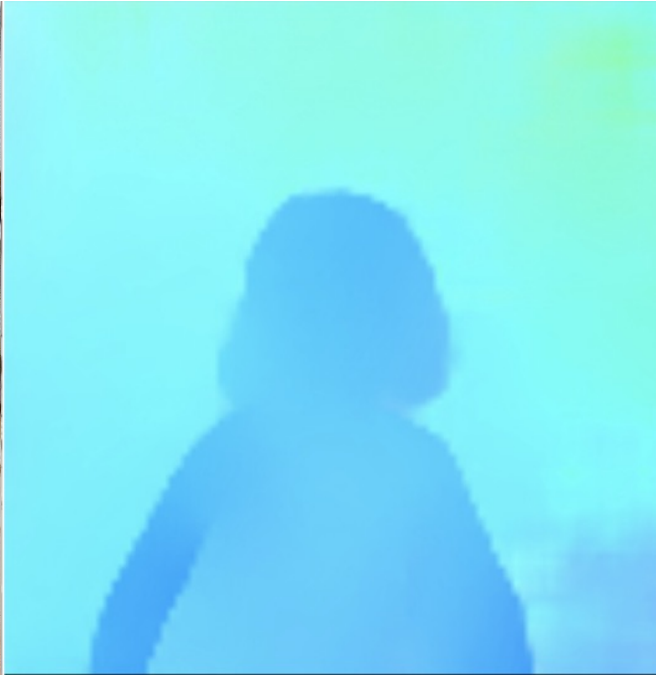




# Real-time optical flow on Pixel 6



Input image pairs

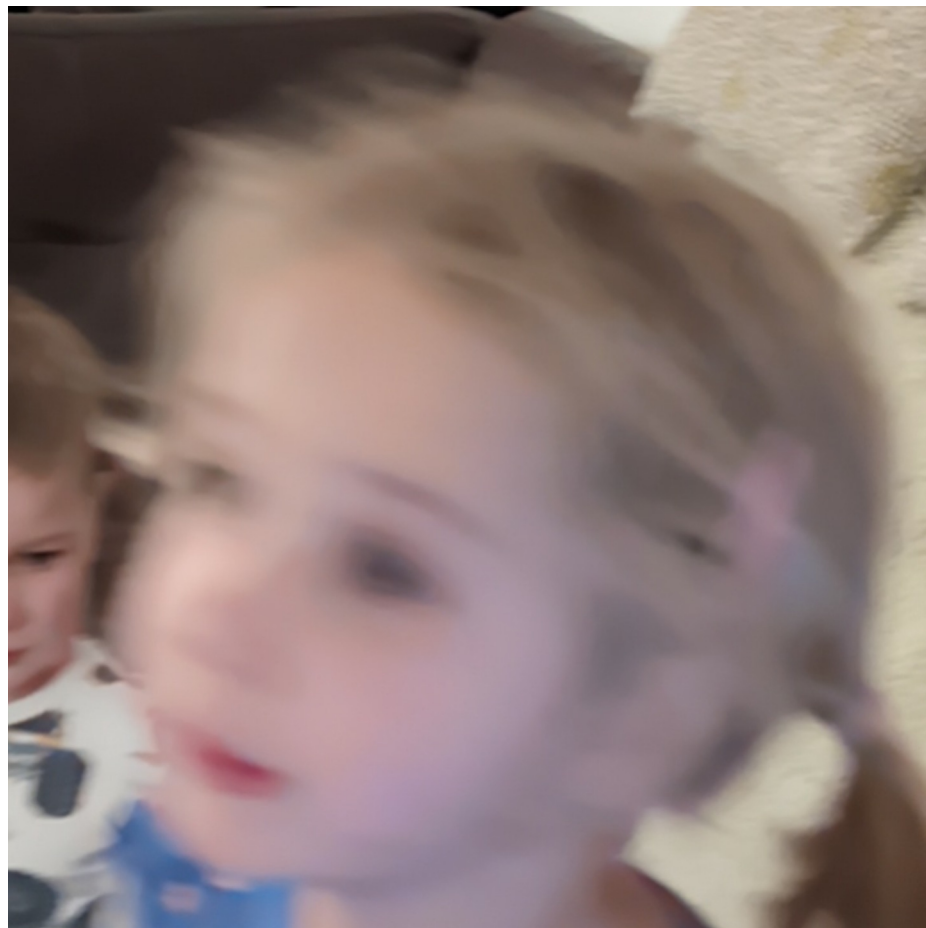


Flow by unoptimized model  
(>9000ms, 2GB memory)



Flow by optimized model  
(~**13ms**, **34MB** memory)

Input: Kids Standing Up



Exposure time = 33 ms, iso = 1024

## Our Deblurred Result



Exposure time = 33 ms, iso = 1024

# Input: Dynamic Motion



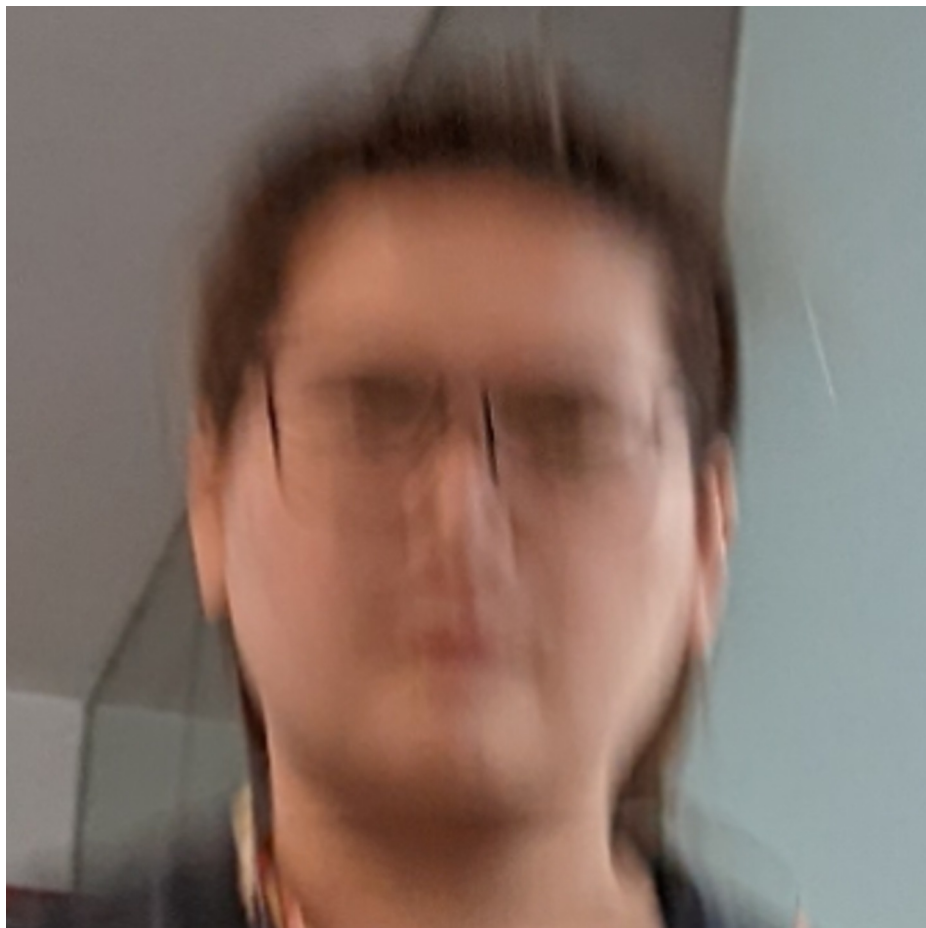
Exposure time = 8.6 ms, iso = 43

## Our Deblurred Result



Exposure time = 8.6 ms, iso = 43

Input: Walking



Exposure time = 8.3 ms, iso = 125

# Our Deblurred Result



Exposure time = 8.3 ms, iso = 125

**What we haven't covered**



# Multiple motions



# Fine details



846314168

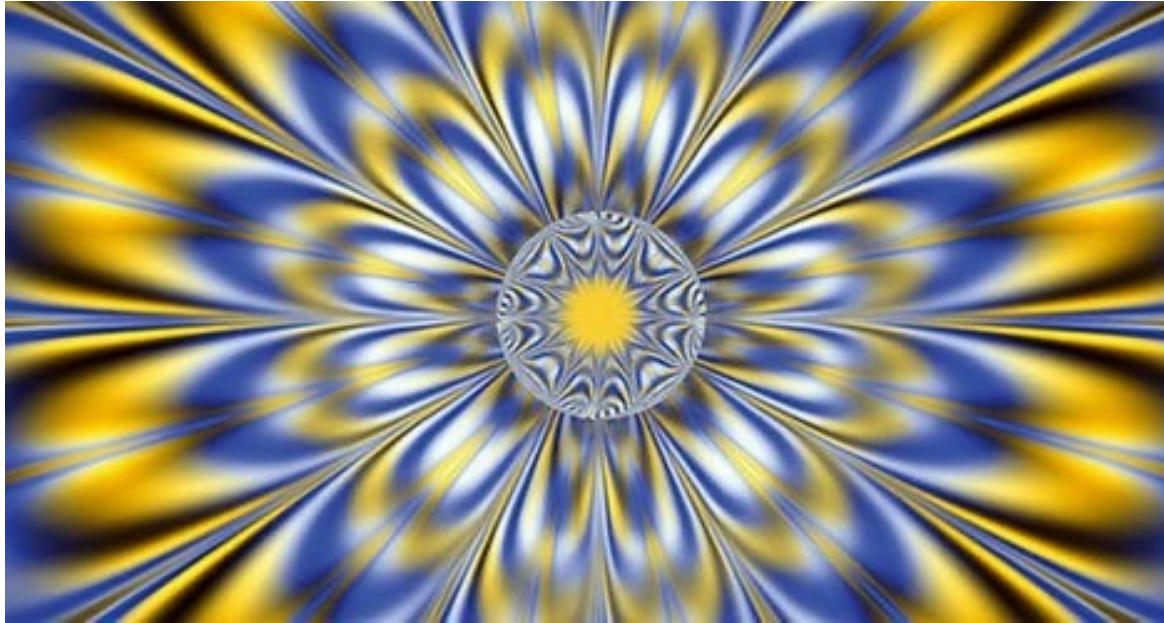
# Even harder



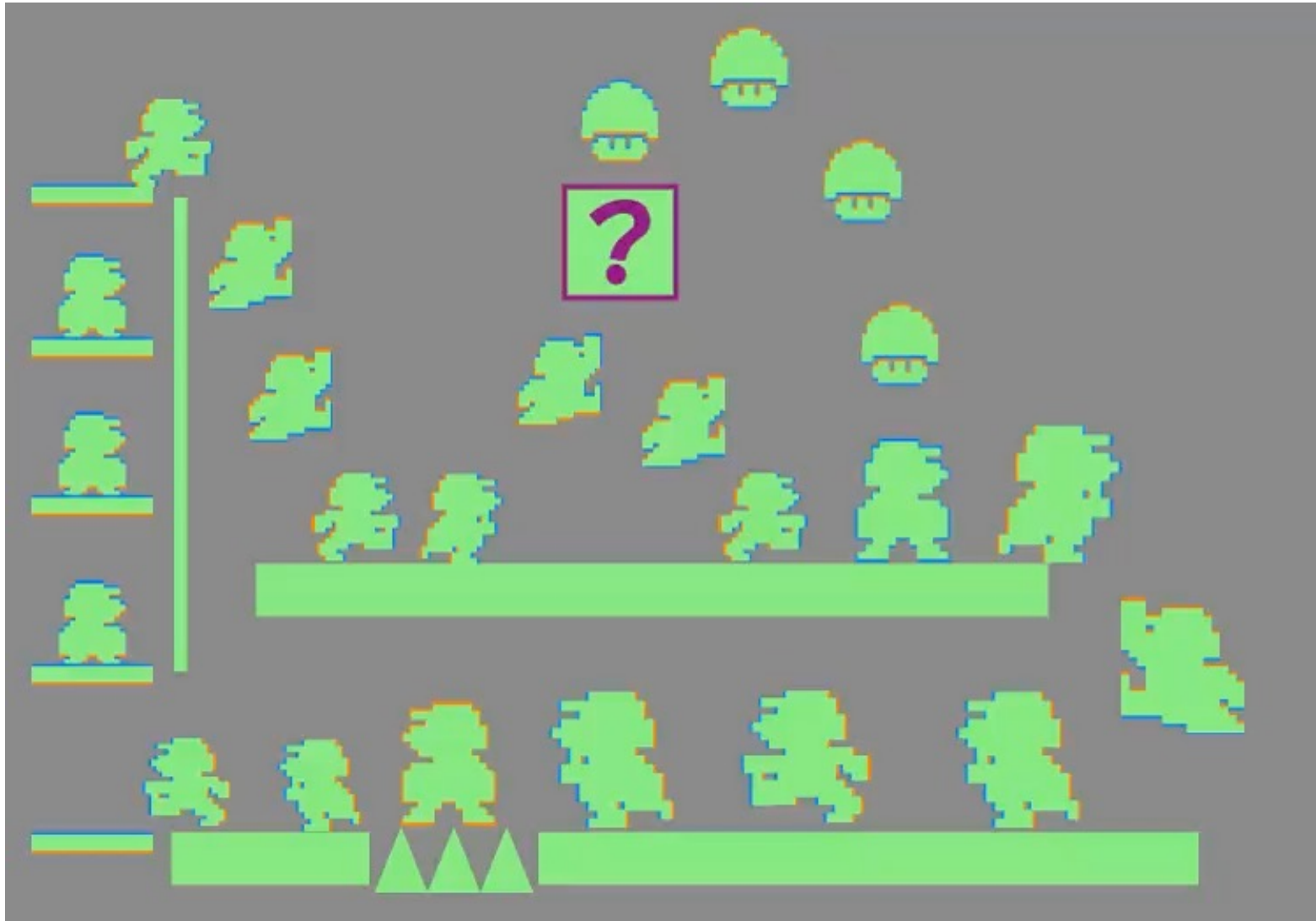
# How to obtain ground truth for real-world videos?



# How do humans perceive motion?



# What do you perceive?

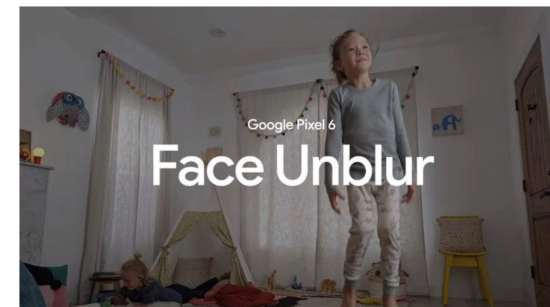
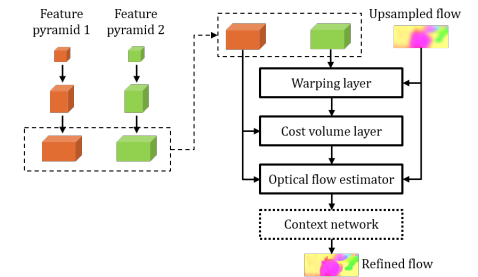
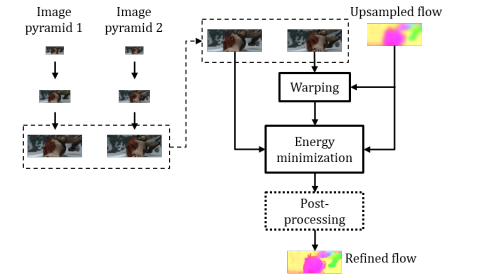


# A single Mario



# Content

- Classical approach
  - Constancy assumption -> matching by comparison (cost volume)
  - Coarse-to-fine, warping-based iterative estimation
- Deep learning-based approach
  - Designing architecture (using domain knowledge)
  - Learning data (matters)
  - Evaluating architectures fairly (trade-off in accuracy and speed/memory)
- Applications: What is motion for?
  - Super-resolution, frame interpolation, articulated 3D reconstruction ...
  - Face Unblur (real-time dense accurate flow on mobile device)





# A “biased” reading list

- Dr. Rick Szeliski’s book (2<sup>nd</sup> edition) chapter 9 on motion estimation
- Chapters 40-43 of book by Antonio, Phillip and Bill
- Horn & Schunck, Lucas & Kanade, Secrets of optical flow
- FlowNet, PWC-Net, IRR-PWC, RAFT, Perceiver IO, GM-Flow, FlowFormer(++), AutoFlow, Disentangling architecture and training
- Bayesian VSR, Super SloMo, FILM, LASR, Face Unblur

# Thank you!

Deqing Sun  
deqingsun@google.com

