

Optical Flow

- Introduction

2016.08.08

Hyeongseok Kim

hskim@capp.snu.ac.kr



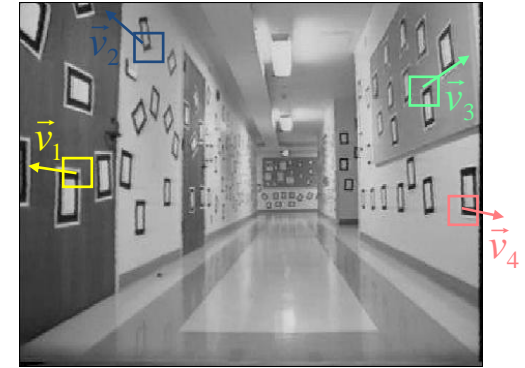
Optical Flow



$I(t), \{p_i\}$



$I(t+1)$



Velocity vectors $\{\vec{v}_i\}$

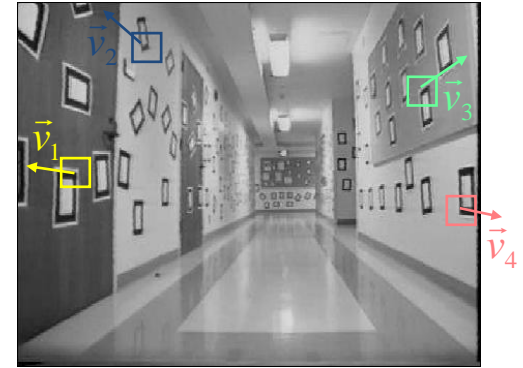
Optical Flow



$I(t), \{p_i\}$



$I(t+1)$



Velocity vectors $\{\vec{v}_i\}$

→ Get 'Motion' information

Optical Flow



$I(t), \{p_i\}$



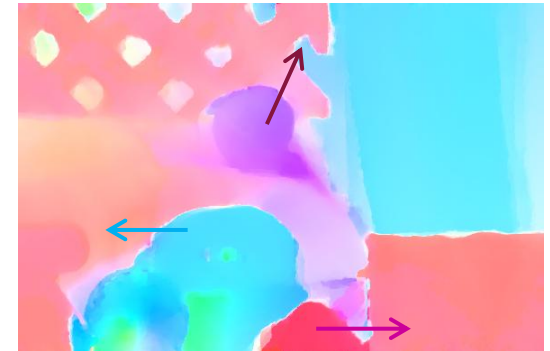
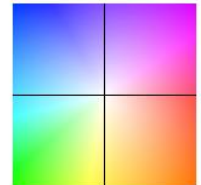
$I(t+1)$



Velocity vectors $\{\vec{v}_i\}$

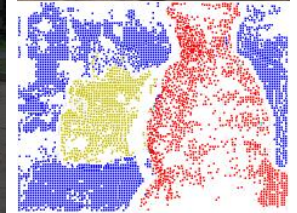
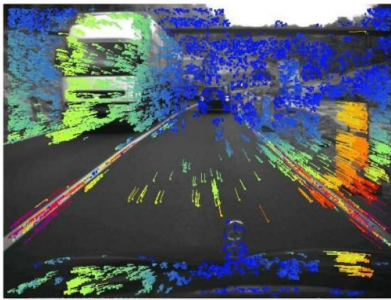
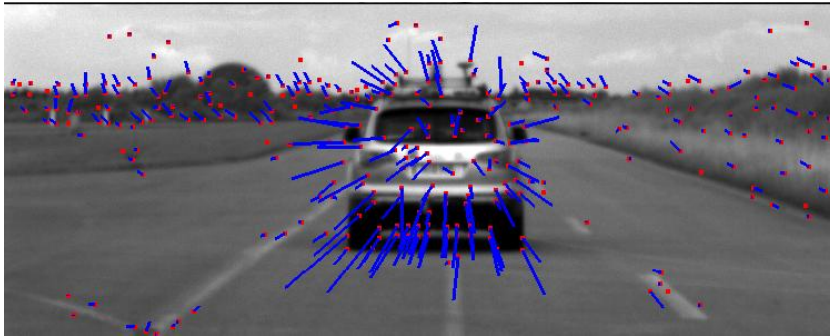
→ Get 'Motion' information

optical flow color encoding scheme



Optical Flow

→ Get 'Motion' information →



Optical Flow

Local method

Brightness constancy

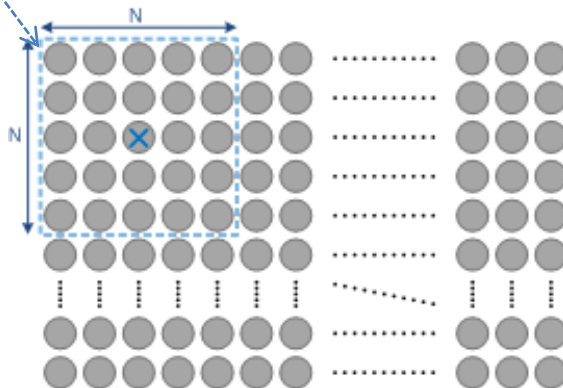
Spatial Coherence

$$I(x(t), y(t), t) = \text{constant}$$

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

$$\begin{bmatrix} \sum w I_{x_i} I_{x_i} & \sum w I_{x_i} I_{y_i} \\ \sum w I_{x_i} I_{y_i} & \sum w I_{y_i} I_{y_i} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum w I_{x_i} I_{t_i} \\ \sum w I_{y_i} I_{t_i} \end{bmatrix}$$

ex) LK (Lucas-Kanade)



Global method

minimize Energy function

$$E(f) = \sum_{(p,q) \in \mathcal{N}} V(f_p, f_q) + \sum_{p \in \mathcal{P}} D_p(f_p)$$

ex) Horn-Schunck, TV-L1

Optical Flow

Local method

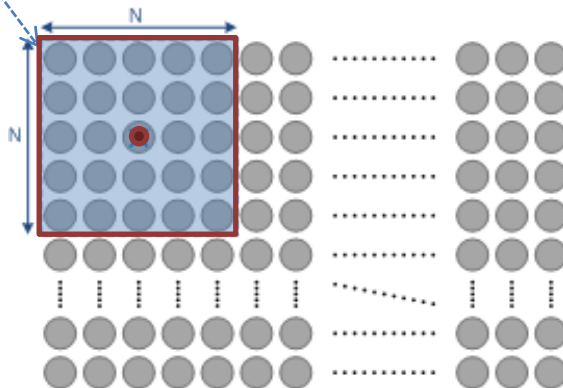
Brightness constancy
Spatial Coherence

$$I(x(t), y(t), t) = \text{constant}$$

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

$$\begin{bmatrix} \sum w I_{x_i} I_{x_i} & \sum w I_{x_i} I_{y_i} \\ \sum w I_{x_i} I_{y_i} & \sum w I_{y_i} I_{y_i} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum w I_{x_i} I_{t_i} \\ \sum w I_{y_i} I_{t_i} \end{bmatrix}$$

ex) LK (Lucas-Kanade)



Global method

minimize Energy function

$$E(f) = \sum_{(p,q) \in \mathcal{N}} V(f_p, f_q) + \sum_{p \in \mathcal{P}} D_p(f_p)$$

ex) Horn-Schunck, TV-L1

Optical Flow

Local method

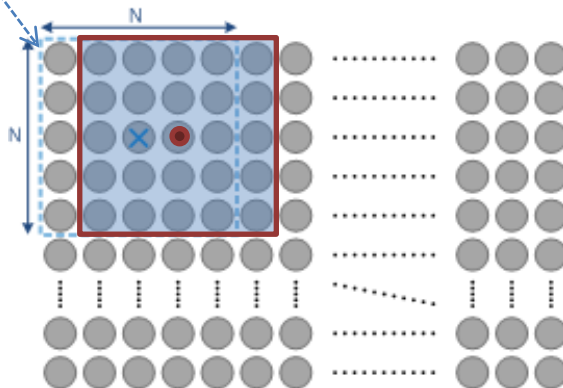
Brightness constancy
Spatial Coherence

$$I(x(t), y(t), t) = \text{constant}$$

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

$$\begin{bmatrix} \sum w I_{x_i} I_{x_i} & \sum w I_{x_i} I_{y_i} \\ \sum w I_{x_i} I_{y_i} & \sum w I_{y_i} I_{y_i} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum w I_{x_i} I_{t_i} \\ \sum w I_{y_i} I_{t_i} \end{bmatrix}$$

ex) LK (Lucas-Kanade)



Global method

minimize Energy function

$$E(f) = \sum_{(p,q) \in \mathcal{N}} V(f_p, f_q) + \sum_{p \in \mathcal{P}} D_p(f_p)$$

ex) Horn-Schunck, TV-L1

Optical Flow

Local method

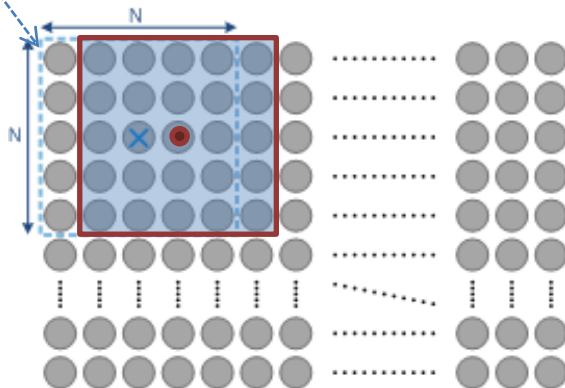
Brightness constancy
Spatial Coherence

$$I(x(t), y(t), t) = \text{constant}$$

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

$$\begin{bmatrix} \sum w I_{x_i} I_{x_i} & \sum w I_{x_i} I_{y_i} \\ \sum w I_{x_i} I_{y_i} & \sum w I_{y_i} I_{y_i} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum w I_{x_i} I_{t_i} \\ \sum w I_{y_i} I_{t_i} \end{bmatrix}$$

ex) LK (Lucas-Kanade)

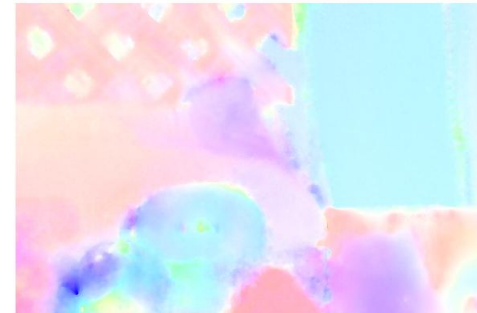


Global method

minimize Energy function

$$E(f) = \sum_{(p,q) \in \mathcal{N}} V(f_p, f_q) + \sum_{p \in \mathcal{P}} D_p(f_p)$$

ex) Horn-Schunck, TV-L1



Optical Flow

Local method

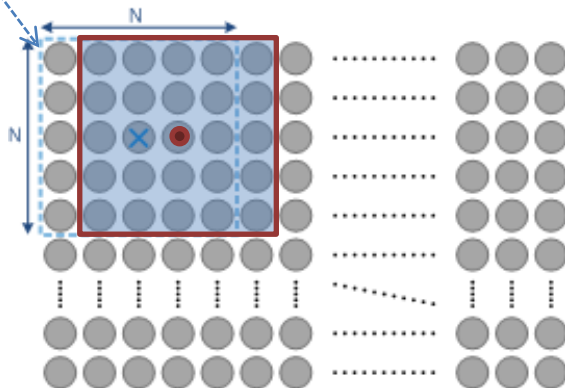
Brightness constancy
Spatial Coherence

$$I(x(t), y(t), t) = \text{constant}$$

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

$$\begin{bmatrix} \sum w I_{x_i} I_{x_i} & \sum w I_{x_i} I_{y_i} \\ \sum w I_{x_i} I_{y_i} & \sum w I_{y_i} I_{y_i} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum w I_{x_i} I_{t_i} \\ \sum w I_{y_i} I_{t_i} \end{bmatrix}$$

ex) LK (Lucas-Kanade)



Global method

minimize Energy function

$$E(f) = \sum_{(p,q) \in \mathcal{N}} V(f_p, f_q) + \sum_{p \in \mathcal{P}} D_p(f_p)$$

ex) Horn-Schunck, TV-L1



Optical Flow

Local method

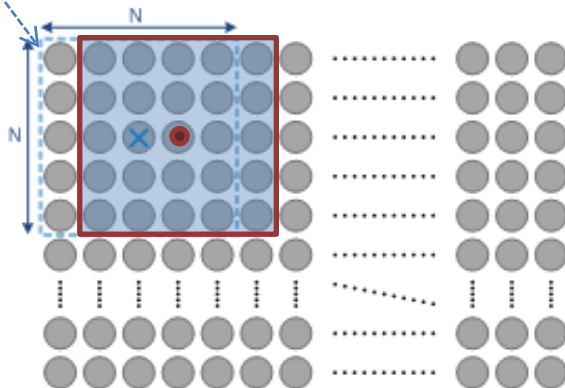
Brightness constancy
Spatial Coherence

$$I(x(t), y(t), t) = \text{constant}$$

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

$$\begin{bmatrix} \sum w I_{x_i} I_{x_i} & \sum w I_{x_i} I_{y_i} \\ \sum w I_{x_i} I_{y_i} & \sum w I_{y_i} I_{y_i} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum w I_{x_i} I_{t_i} \\ \sum w I_{y_i} I_{t_i} \end{bmatrix}$$

ex) LK (Lucas-Kanade)



Global method

minimize Energy function

$$E(f) = \sum_{(p,q) \in \mathcal{N}} V(f_p, f_q) + \sum_{p \in \mathcal{P}} D_p(f_p)$$

ex) Horn-Schunck, TV-L1



Optical Flow

Local method

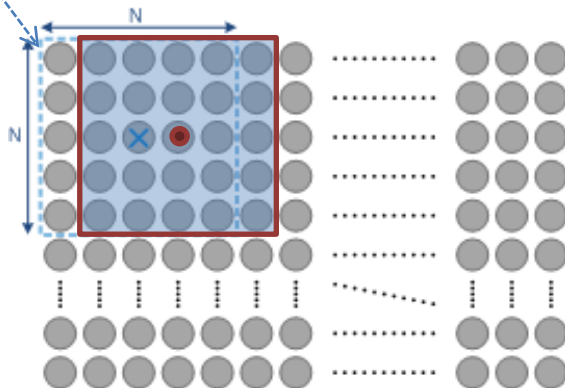
Brightness constancy
Spatial Coherence

$$I(x(t), y(t), t) = \text{constant}$$

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

$$\begin{bmatrix} \sum w I_{x_i} I_{x_i} & \sum w I_{x_i} I_{y_i} \\ \sum w I_{x_i} I_{y_i} & \sum w I_{y_i} I_{y_i} \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum w I_{x_i} I_{t_i} \\ \sum w I_{y_i} I_{t_i} \end{bmatrix}$$

ex) LK (Lucas-Kanade)



Global method

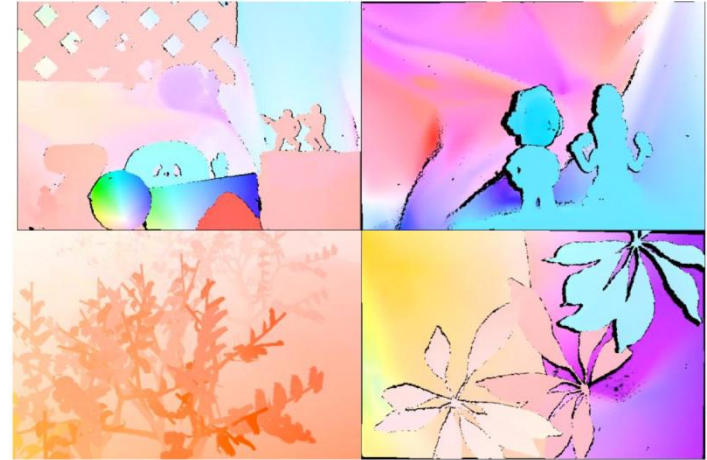
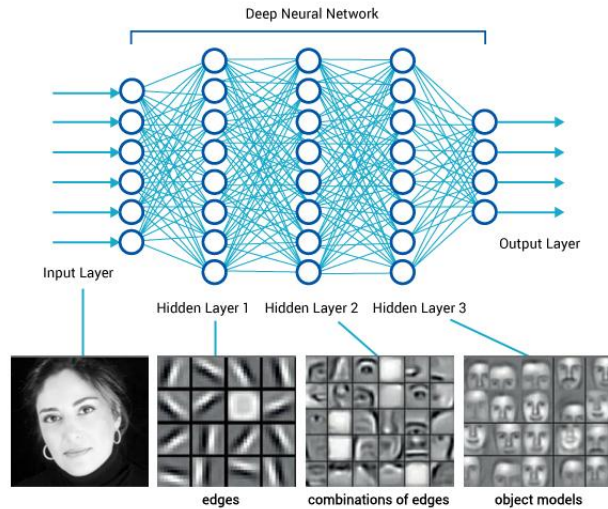
minimize Energy function

$$E(f) = \sum_{(p,q) \in \mathcal{N}} V(f_p, f_q) + \sum_{p \in \mathcal{P}} D_p(f_p)$$

ex) Horn-Schunck, TV-L1



Deep Learning for Optical Flow



Per-pixel prediction task

- like semantic segmentation, depth prediction, keypoint prediction and edge detection

From a pair of images

Ex. DeepFlow(@ICCV 2013), EpicFlow(@CVPR 2015),
FlowNet(@ICCV 2015), PathBatch(@CVPR 2016), ...

FlowNet: Learning Optical Flow with Convolutional Networks

Philip Fischer¹ Alexey Dosovitskiy¹ Eddy Ilg¹ Philip Hager¹ Case Huxley¹ Vladislav Golkov²
¹University of Freiburg ²Technical University of Munich
 {fischer, dosovitskiy, ilg, ilg, hager, huxley, golkov}@informatik.uni-freiburg.de {dosovitskiy, golkov}@informatik.uni-muenchen.de
 Patrick van der Smagt¹ Daniel Cremers¹ Thomas Brox¹
 Technical University of Munich Technical University of Munich University of Freiburg
 smagt@informatik.uni-freiburg.de cremers@informatik.uni-muenchen.de brox@informatik.uni-freiburg.de

Abstract

Convolutional neural networks (CNNs) have recently been very successful in a variety of computer vision tasks, especially in those related to recognition. Optical flow estimation has not been getting the same level of success so far. In this paper we construct approximate CNNs which are capable of solving the optical flow estimation problem in a supervised learning task. We propose and compare two architectures, a generic architecture and another one including a layer that correlates feature vectors at different image locations.

Since existing ground truth datasets are not sufficiently large to train a CNN, we propose a synthetic FlowNet dataset. We show that networks trained on this synthetic data still generalize very well to existing datasets such as Sintel and KITTI, achieving competitive accuracy at frame rates of 3 to 30 fps.

1. Introduction

Convolutional neural networks have become the method of choice in many fields of computer vision. They are usually applied to classification [15, 21], but recently pre-trained networks also allow for pixel-level prediction like semantic segmentation [20] or depth estimation from single images [18]. In this paper, we propose training CNNs end-to-end to learn predicting the optical flow field from a pair of images.

While optical flow estimation needs precise per pixel information, it also requires finding correspondences between two image regions. This involves not only detecting image features but also learning to match them at different locations in the two images. In this respect, optical

¹Supported by the Deutsche Forschungsgemeinschaft (DFG) under the Collaborative Program SFB/TR 8.



Figure 1. We propose neural networks which learn to estimate optical flow, being trained end-to-end. The information is then quickly incorporated in a completely part of the network and then output as an optical flow.

flow estimation fundamentally differs from previous applications of CNNs.

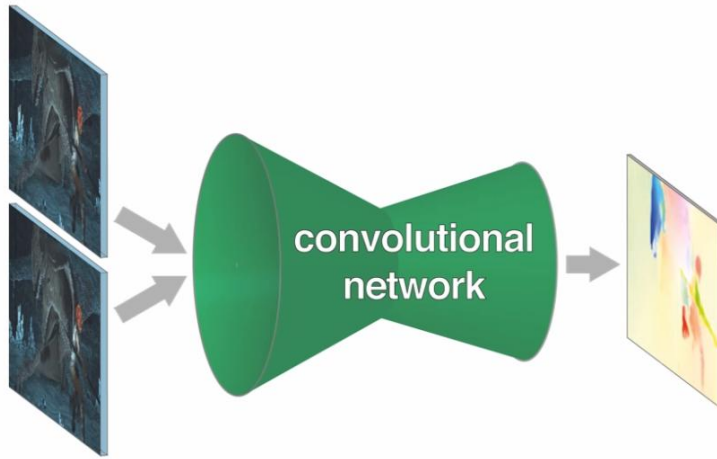
Since it was not clear whether this task could be solved with a standard CNN architecture, we deliberately designed an architecture with a correlation layer that explicitly provides matching capabilities. This architecture is based end-to-end. The idea is to exploit the ability of convolutional networks to learn strong features at multiple levels of scale and abstraction and to couple them tightly to the actual correspondences based on flow features. The layers on top of the correlation layer learn how to produce flow from these matches. Surprisingly, training the network this way is not necessary and even the core network can learn to predict optical flow with competitive accuracy.

Training such a network to predict precise optical flow requires a sufficiently large training set. Although data augmentation does help, the existing optical flow datasets are still too small to train a network on per pixel scale of the an-

FlowNet

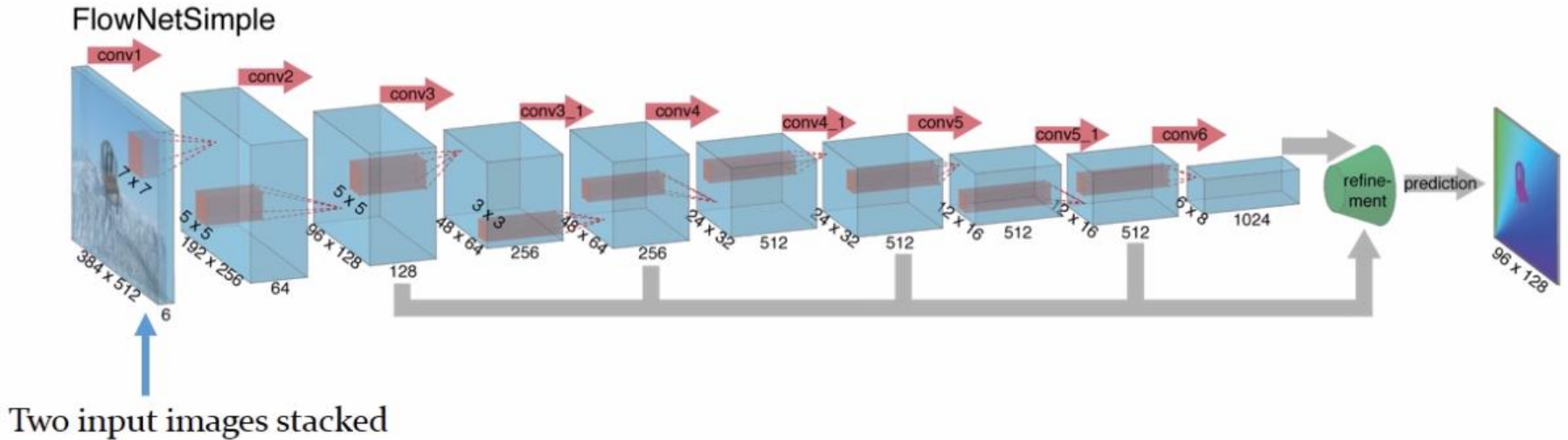
: Learning Optical Flow with Convolutional Networks

@ 2015 ICCV

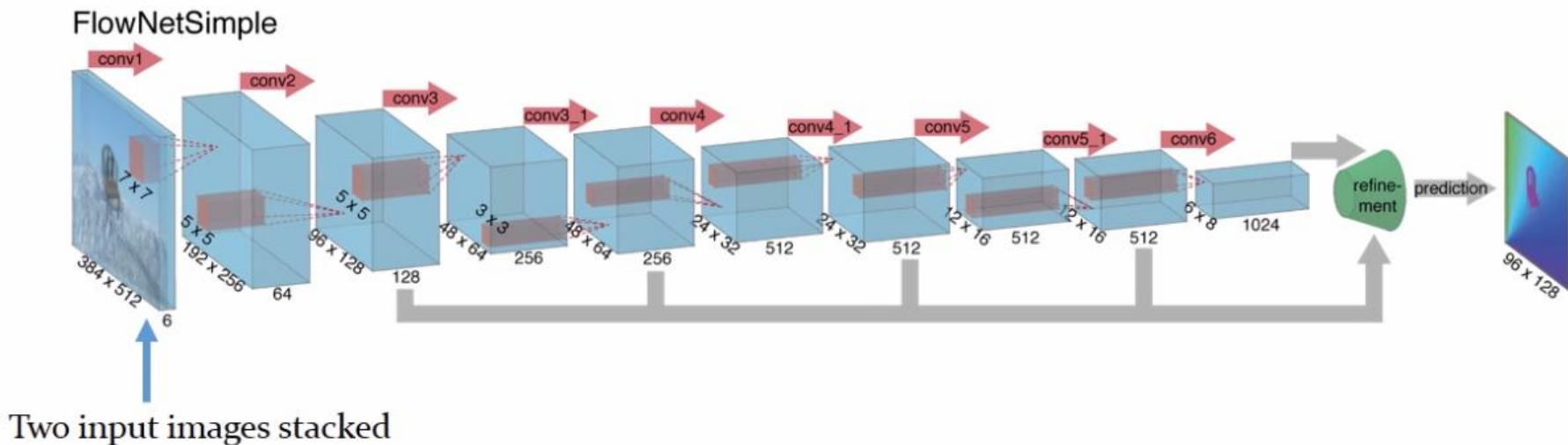


First spatially compressed in a contractive part of the network and then refined in an expanding part.

FlowNet



Stack both input images together and feed them through a network
- a rather generic architecture

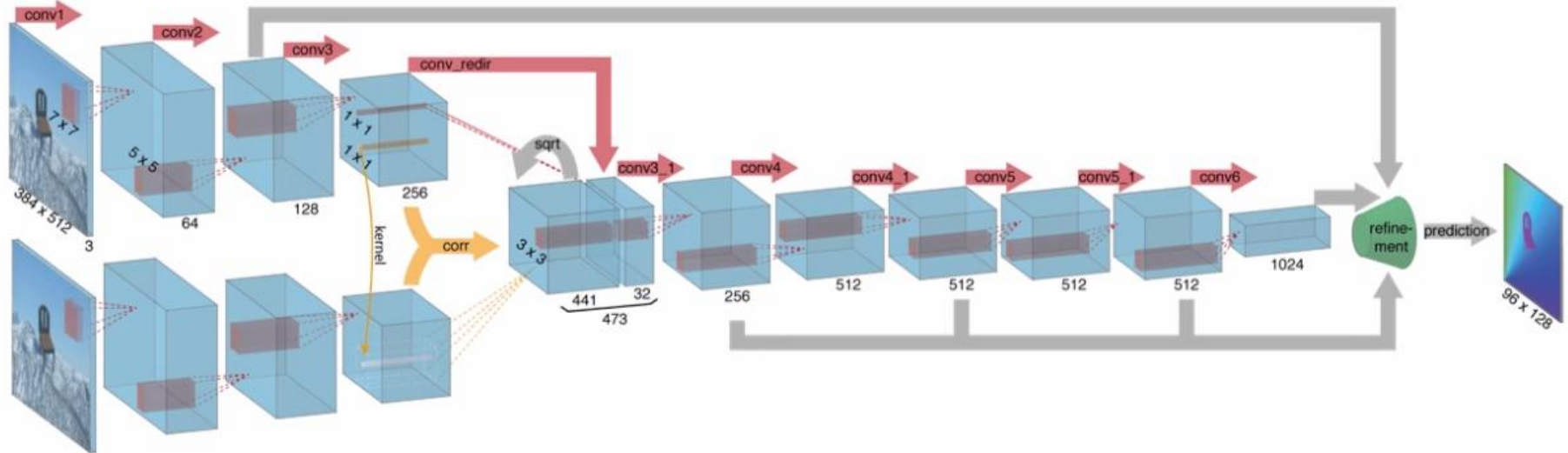


Stack both input images together and feed them through a network
- a rather generic architecture

If this network is large enough, it could learn to predict optical flow.
However, we can never be sure that a **local gradient optimization** like gradient descent can get the network to this point.

→ Therefore, it could be beneficial to hand-design an architecture
which is less generic, but may perform better with the given data and optimization
techniques. → “FlowNetCorr”

FlowNetCorr

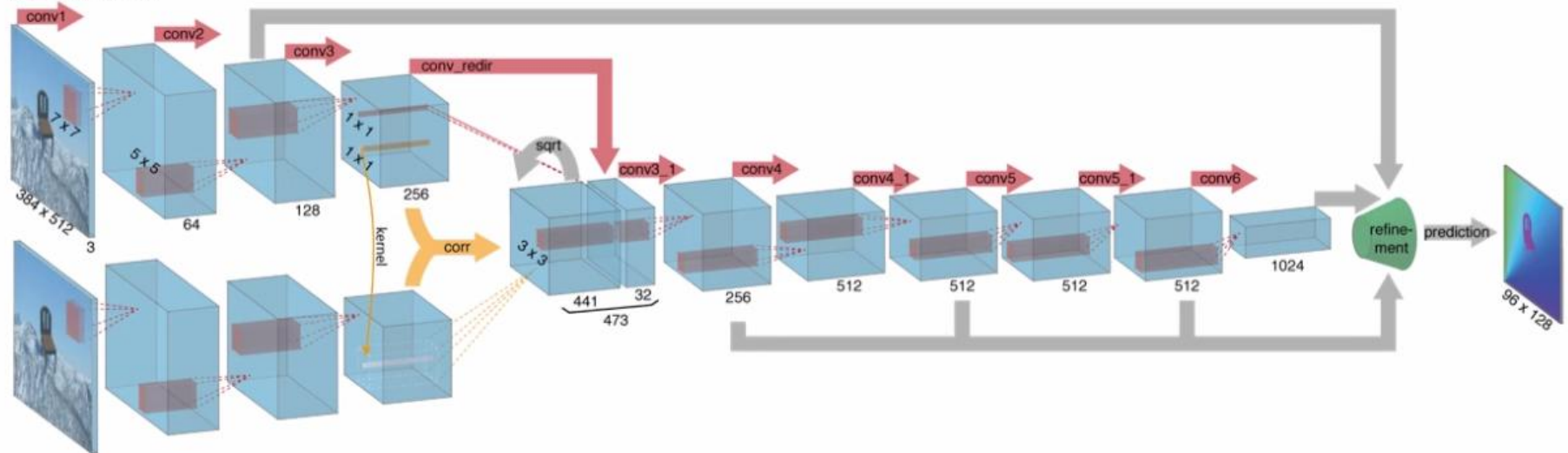


Two separate steps.

First produce meaningful representations of the two images separately
And then combine them on a higher level (by correlation layer)

→ Roughly resembles the standard matching approach.
when one first extracts features from patches of both images
and then compares those feature vectors.

FlowNetCorr



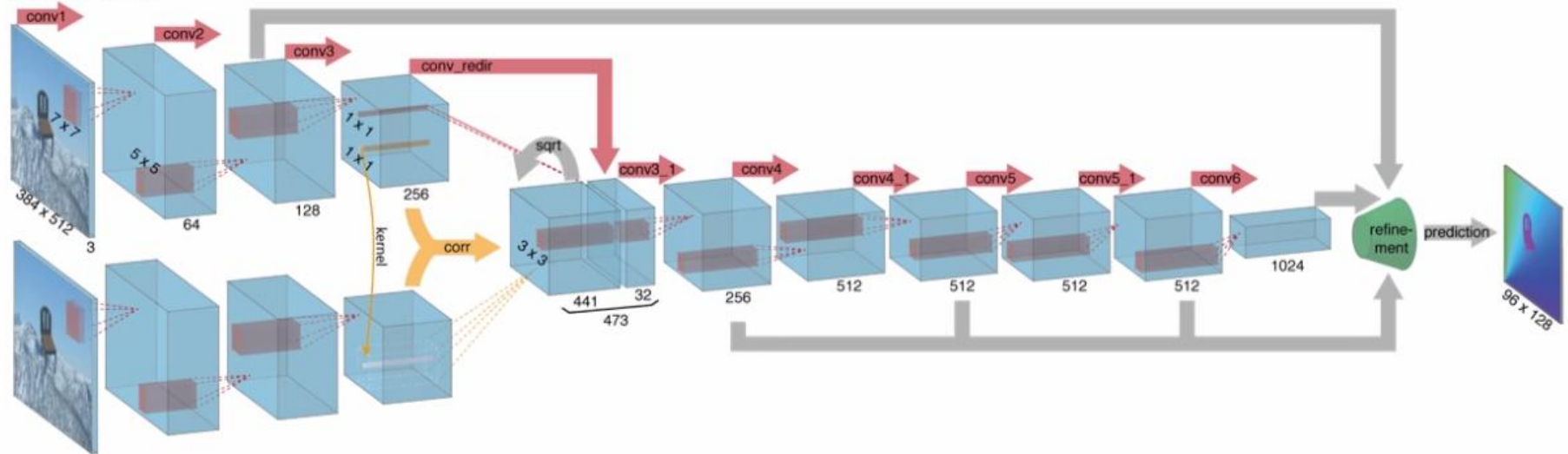
‘Correlation layer’ helps the network find correspondences.

Given two multi-channel feature maps
(with w , h , and c being their width, height, and number of channels),

‘correlation’ of two patches centered at \mathbf{x}_1 in the first map and \mathbf{x}_2 in the second map :

$$c(\mathbf{x}_1, \mathbf{x}_2) = \sum_{\mathbf{o} \in [-k, k] \times [-k, k]} \langle \mathbf{f}_1(\mathbf{x}_1 + \mathbf{o}), \mathbf{f}_2(\mathbf{x}_2 + \mathbf{o}) \rangle$$

FlowNetCorr



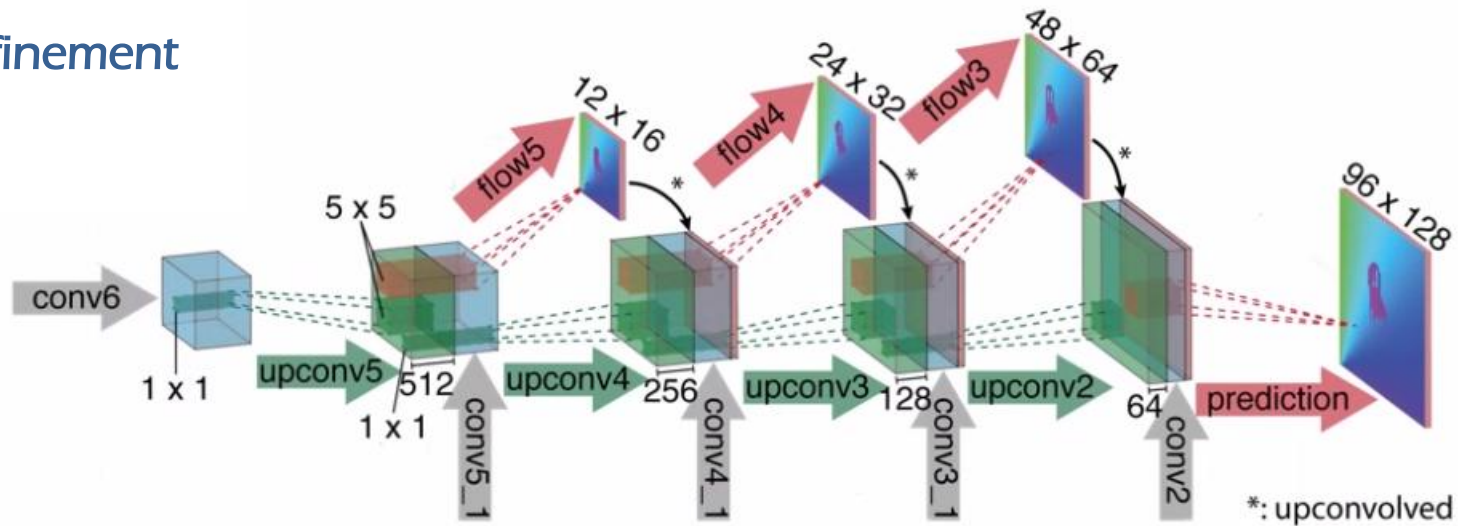
'Correlation layer'

Comparing all patch combinations : $w^2 \cdot h^2$ computations \rightarrow inefficient (too large)

Maximum displacement d , : $w \times h \times D^2$ computations(outputs), $\bar{D} := 2d + 1$

- limit the maximum displacement for comparisons
- strides s_1 and s_2 to quantize \mathbf{x}_1 globally and to quantize \mathbf{x}_2 within the neighborhood centered around .

Refinement

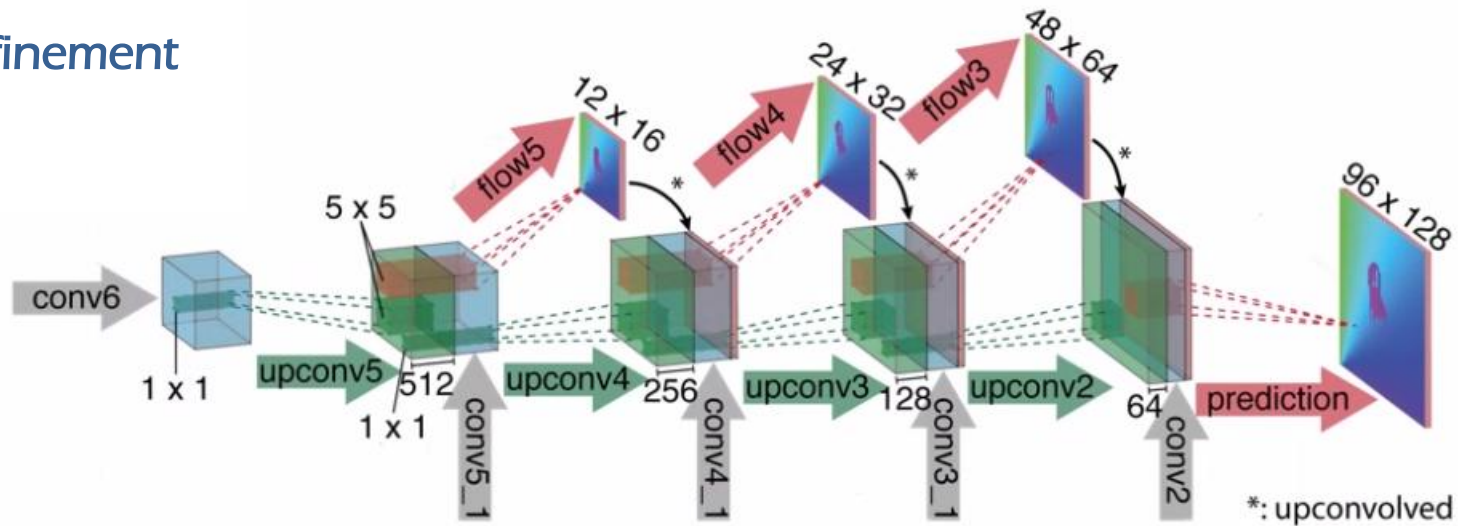


To provide dense per-pixel predictions we need to refine the coarse pooled representation.

Apply the 'upconvolution' to feature maps, and concatenate it with corresponding feature maps from the 'contractive' part of the network and an upsampled coarser flow prediction (if available).

→ Preserve both high-level information passed from coarser feature maps and fine local information provided in lower layer feature maps.

Refinement



Repeat 4 times

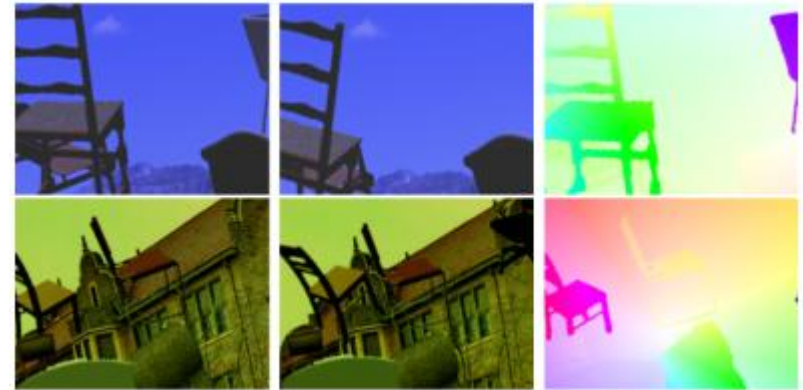
- resulting in a predicted flow for which the resolution is still 4 times smaller than the input.
- further refinement does not significantly improve the results, compared to computationally less expensive bilinear upsampling to full image resolution.
- Alternative scheme : bilinear upsampling \rightarrow variational approach

Training Data

	Frame pairs	Frames with ground truth	Ground truth density per frame
Middlebury	72	8	100%
KITTI	194	194	~50%
Sintel	1,041	1,041	100%
Flying Chairs	22,872	22,872	100%



Generated image pair



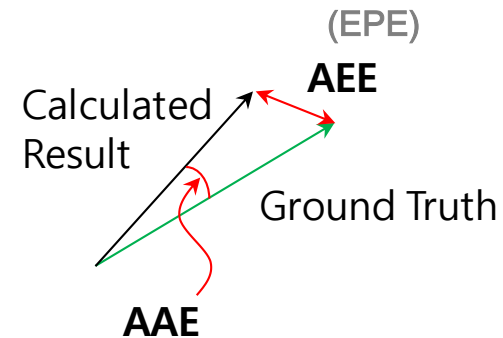
Augmented image pair

- geometric transformation
: translation, rotation, scaling
- Gaussian noise
- brightness, contrast, gamma, color

Results

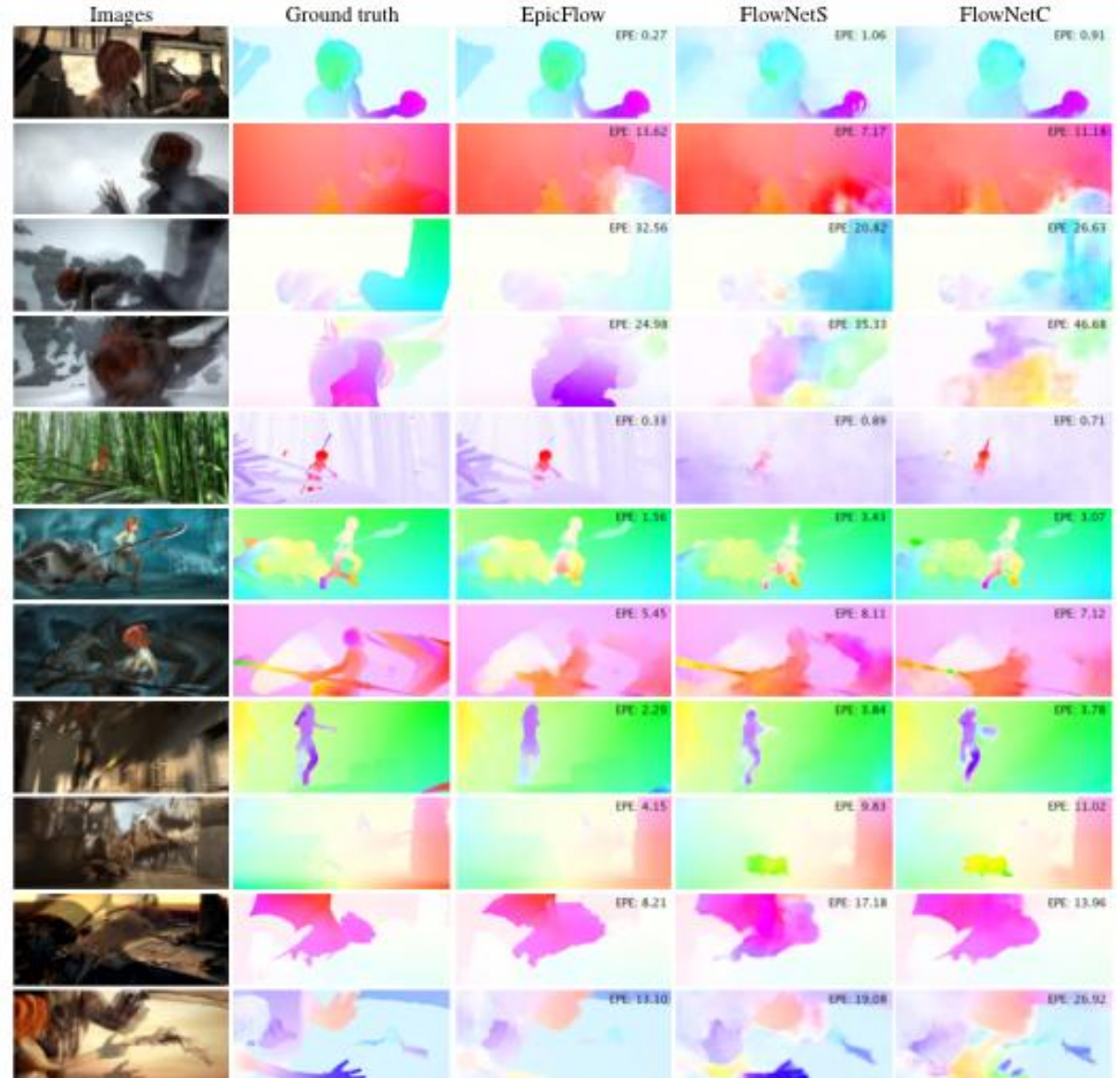
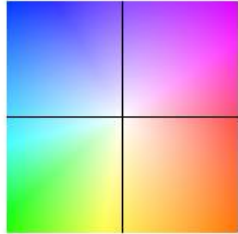
Method	Sintel Clean		Sintel Final		KITTI		Middlebury train		Middlebury test		Chairs test	Time (sec)	
	train	test	train	test	train	test	AEE	AAE	AEE	AAE		CPU	GPU
EpicFlow [30]	2.40	4.12	3.70	6.29	3.47	3.8	0.31	3.24	0.39	3.55	2.94	16	-
DeepFlow [35]	3.31	5.38	4.56	7.21	4.58	5.8	0.21	3.04	0.42	4.22	3.53	17	-
EPPM [3]	-	6.49	-	8.38	-	9.2	-	-	0.33	3.36	-	-	0.2
LDOF [6]	4.29	7.56	6.42	9.12	13.73	12.4	0.45	4.97	0.56	4.55	3.47	65	2.5
FlowNetS	4.50	7.42	5.45	8.43	8.26	-	1.09	13.28	-	-	2.71	-	0.08
FlowNetS+v	3.66	6.45	4.76	7.67	6.50	-	0.33	3.87	-	-	2.86	-	1.05
FlowNetS+ft	(3.66)	6.96	(4.44)	7.76	7.52	9.1	0.98	15.20	-	-	3.04	-	0.08
FlowNetS+ft+v	(2.97)	6.16	(4.07)	7.22	6.07	7.6	0.32	3.84	0.47	4.58	3.03	-	1.05
FlowNetC	4.31	7.28	5.87	8.81	9.35	-	1.15	15.64	-	-	2.19	-	0.15
FlowNetC+v	3.57	6.27	5.25	8.01	7.45	-	0.34	3.92	-	-	2.61	-	1.12
FlowNetC+ft	(3.78)	6.85	(5.28)	8.51	8.79	-	0.93	12.33	-	-	2.27	-	0.15
FlowNetC+ft+v	(3.20)	6.08	(4.83)	7.88	7.31	-	0.33	3.81	0.50	4.52	2.67	-	1.12

AAE and AEE with several datasets



Results

optical flow color encoding scheme



with Sintel dataset

https://youtu.be/k_wkDLJ8IJE

Thank you