

# Your Objective Is Wrong: Rethink Unsupervised learning of Optical Flow

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

*The ABSTRACT is to be in fully-justified italicized text, at the top of the left-hand column, below the author and affiliation information. Use the word “Abstract” as the title, in 12-point Times, boldface type, centered relative to the column, initially capitalized. The abstract is to be in 10-point, single-spaced type. Leave two blank lines after the Abstract, then begin the main text. Look at previous CVPR abstracts to get a feel for style and length.*

## 1. Introduction

### 1.1. Related Work

### 1.2. Novel Contribution

We extend FlowNet [?] in this work, in summary our contributions are three folds. First, we proposed to combine traditional layered approach for optical flow estimation with deep learning. The proposed approach does not require pre-segmentation of images, instead, the separation of layers is automatically done during training the network. Second, a soft-masks module is proposed. This soft-masks module implements a channel-wise maxout operation among masks. As a result, the estimated optical flow will be separated to layers. each of which will contain optical flow that is estimated using a quadratic function. Third, we extend the FlowNet by adding the proposed soft-mask module in output layers, the resulting network is trained to compare with both supervised and unsupervised optical flow estimation approaches using neural networks. The empirical results show that the proposed network structure achieves comparable or lower error in each experimental group.

## 2. Methodology

The proposed approach and corresponding analysis will be introduced in this section.

### 2.1. Annotation

Given a pair of images  $I_a, I_b \in \mathbb{R}^{H \times W \times C}$  as input, where  $H, W$  and  $C$  are height, width and channels of the input images. The proposed approach is going to estimate an optical flow field  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{H \times W}$ , where  $\mathbf{u}$  and  $\mathbf{v}$  are the horizontal and vertical components of the optical flow field to be estimated that transform image from  $I_a$  to  $I_b$ . The original formulation of optical flow estimation is proposed by Horn and Schunck in [?]. In classical formulation, an objective function is composed with a combination of a data term which makes a local constancy assumption of some image property and a spatial term that models how the flow is expected to vary across images. We write the classical optical flow objective function as:

$$E(\mathbf{u}, \mathbf{v}) = \sum_i^H \sum_j^W (I_1(i+u_{ij}, j+v_{ij}) - I_0(i, j))^2 + \varphi(\mathbf{u}, \mathbf{v}) \quad (1)$$

where  $\varphi(\mathbf{u}, \mathbf{v})$  is a regularization term that constrains smoothness of optical flow.

Nowadays, the above objective is still being used by many optical flow estimation using deep neural network based on unsupervised training framework [?][?][?]. We also use above objective when training our network and comparing with results of unsupervised methods. Experiments results are presented in Section 3.

### 2.2. Soft-masks module

FlowNet [?] is the first work that uses deep convolutional neural network for flow estimation. The network architecture used by FlowNet is very similar to classical structure of auto-encoder, where optical flows are generated using deconvolution at each scale level of the image pyramid. To refine estimation of the flows, shortcuts are built to connect layers of corresponding level in encoder and decoder. Let's take a look at a single computation of convolution, and for simplicity, let's assume  $f$  represents both horizontal and vertical components of an output flow. Given  $x \in \mathbb{R}^{s \times s \times c}$ , representing a volume feature vector in the input feature volume, inputted to output layer, where  $s$  is kernel size and

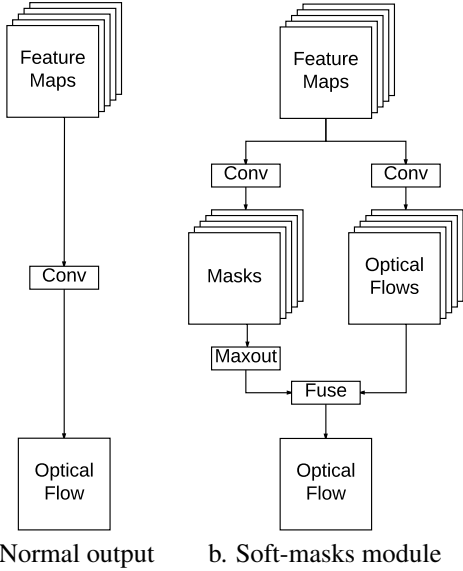


Figure 1. An illustration of the structure of the proposed soft-masks module compared with normal linear optical flow output.

$c$  is number of channels. FlowNet employs a linear activation to compute optical flow:

$$f = x^T W + b \quad (2)$$

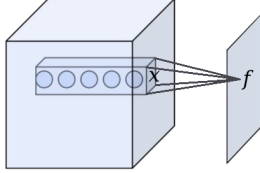


Figure 2. An illustration of annotation used in convolution.

Given actual optical flow field a nonlinear and piece-wise smooth representation of motions contained in images, using linear function to fit the flow field is less accurate. We introduce a soft-masks module in this paper. The proposed module can be used to replace the linear output of optical flow estimation. We will show that by using this module, we are able to separate optical flow field to multiple layers. Flow estimation in each layer is smooth and can be estimated using a quadratic function, which results in a more accurate and flexible optical flow estimation.

An illustration of soft-masks module is shown in Figure 1. The essential part of the soft-masks module is its dual-branch structure which contains mask branch and optical flow branch. The same input feature maps represented as a set of volume feature vectors,  $\mathbf{x} = \{x\}, x \in \mathbb{R}^{s \times s \times c}$  thus are imported to both branches. The most significant contribution of this work is to separate one optical flow field to multiple layers. For a separation of  $k$  layers.  $k$  masks will be generated in mask branch as illustrated in

Figure 1, which requires  $k$  convolutional filters labeled as  $\{W_n^m, b_n^m\}_{n=1}^k$  being used in mask branch. Correspondingly, the same number of filters are used in optical flow branch which are labeled as  $\{W_n^f, b_n^f\}_{n=1}^k$ . Then the mask and intermediate optical flow could be computed as following:

$$m_n = x^T W_n^m + b_n^m \quad \text{for } n = 1 \dots k \quad (3)$$

$$f_n = x^T W_n^f + b_n^f \quad \text{for } n = 1 \dots k \quad (4)$$

Thus, give  $k$  filters used in both branches, we will obtain  $k$  corresponding pairs of mask and intermediate optical flow. Basically, by using  $k$  filters in optical flow branch and generating  $k$  intermediate optical flow, we assume each filter will work independently and being active only to a single type or a few types of object motions. Correspondingly, filters in mask branch are expected to have some behaviors that each of the generated masks which indeed are active maps should be high active for certain types of motions in some region and low active for other types of motions in other regions. This inspires us to use a maxout operation to extract mask entries with maximal activation along channel axis. So, after maxout operation, for each mask  $m_n, n = 1 \dots k$ , all entries will be zero-out except entries whose activation value are maximal in the some region among all masks. We denote masks after maxout as  $\{m'_n\}_{n=1}^k$ . Thus, there is no intersection among masks in  $\{m'_n\}_{n=1}^k$  and the union of all  $m'_n, n = 1 \dots k$  has activation in full region. Using  $m_n(i, j)$  to denote the activation at  $i, j$  location of mask  $m_n$ , the maxout can be represented as following:

$$m'_n(i, j) = \begin{cases} m_n(i, j), & \text{if } m_n(i, j) = \max_{p=1 \dots k} (m_p(i, j)) \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

## 2.3. Target Image Estimation Network

## 2.4. Flow Estimation Network

## 3. Empirical Evaluation

### 3.1. Datasets

### 3.2. Training Details

### 3.3. Results

## 4. Discussion and Summary