Devon: Deformable Volume Network for Learning Optical Flow

Yao Lu, Jack Valmadre, Heng Wang, Juho Kannala, Mehrtash Harandi, Philip Torr

ANU & Oxford & Facebook & Aalto

November 25, 2023

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Optical Flow

Pixel-wise motion between two images







◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへぐ

Learning Optical Flow



Why?

- Learn good features for matching
- Scene reasoning (e.g. occlusion, segmentation, semantics)

Standard CNNs only have small receptive field

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Standard solution: multi-resolution + warping

Multi-resolution + Warping

Coarse-to-fine optical flow estimation



◆□ > ◆□ > ◆三 > ◆三 > 三 - のへで

Multi-resolution + Warping

Problems:

small things move fast

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

warping

Small Things Move Fast

- At low resolution, small things disappear.
- At high resolution, the motion is too large to be covered.

The Problem of Warping



(a) First image



(b) Second image



(c) Ground truth flow



<□> <@> < => < => < => < => < => <</p>

The Problem of Warping

$$\tilde{J}(x) = J(x + F(x))$$

$$F(\mathbf{p}_1) = 0, \quad F(\mathbf{p}_2) = \mathbf{p}_1 - \mathbf{p}_2,$$
 (1)

we have

$$\widetilde{J}(\mathbf{p}_1) = J(\mathbf{p}_1 + F(\mathbf{p}_1))$$
(2)

$$= J(\mathbf{p}_1 + 0) = J(\mathbf{p}_1), \qquad (3)$$

$$\widetilde{J}(\mathbf{p}_2) = J(\mathbf{p}_2 + F(\mathbf{p}_2)) \tag{4}$$

$$= J(\mathbf{p}_2 + \mathbf{p}_1 - \mathbf{p}_2) = J(\mathbf{p}_1).$$
 (5)

$$\widetilde{J}(\mathbf{p}_1) = \widetilde{J}(\mathbf{p}_2) = J(\mathbf{p}_1)$$

Solution

Key idea: instead of deforming the images (downsample + warping), we deform the correspondences.



Devon



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

- No multi-resolution
- No warping
- Only 1M parameters

Results

