



# Volumetric Correspondence Networks for Optical Flow

# Code available:

github.com/gengshan-y/VCN



# Gengshan Yang<sup>1</sup>, Deva Ramanan<sup>1,2</sup>

Carnegie Mellon University<sup>1</sup>, Argo AI<sup>2</sup>

## Introduction

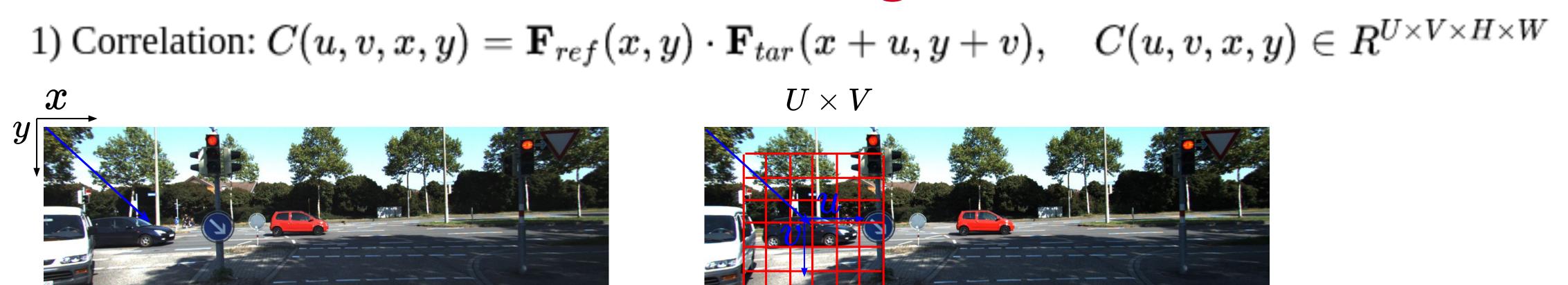
### Optical flow / visual correspondences





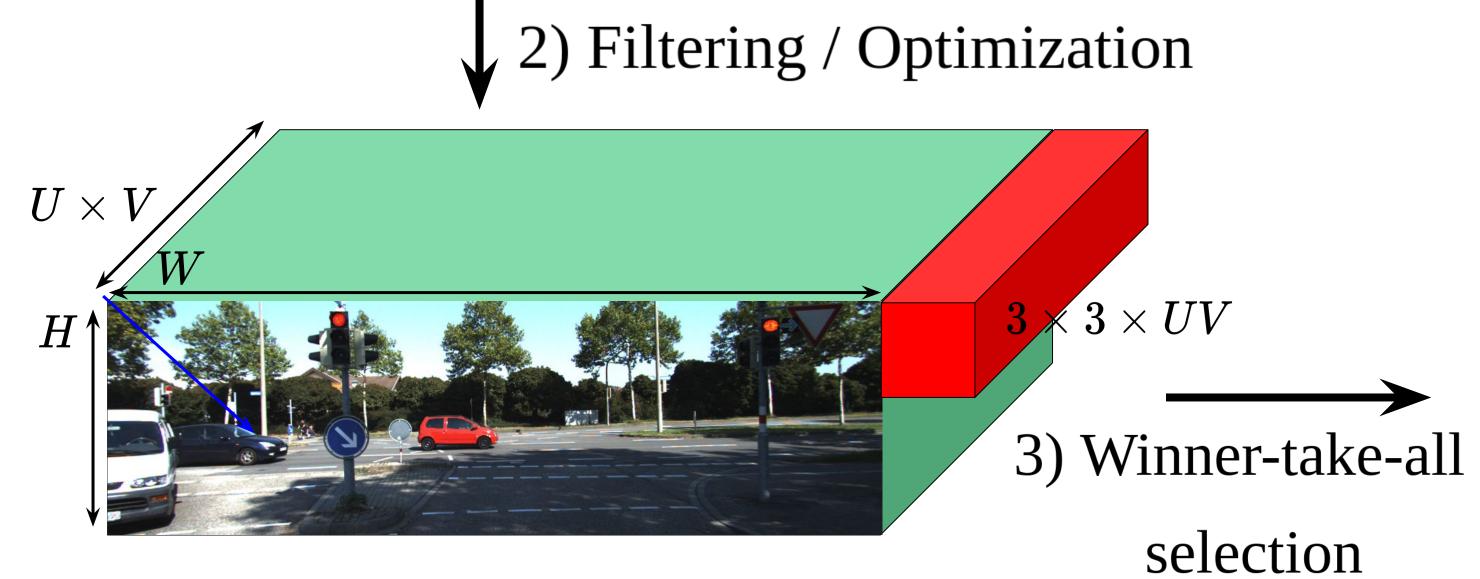
Overlaid reference and target image

#### Related work: cost volume filtering



Reference image

Target image





Dense 2D displacement fields

#### Contributions

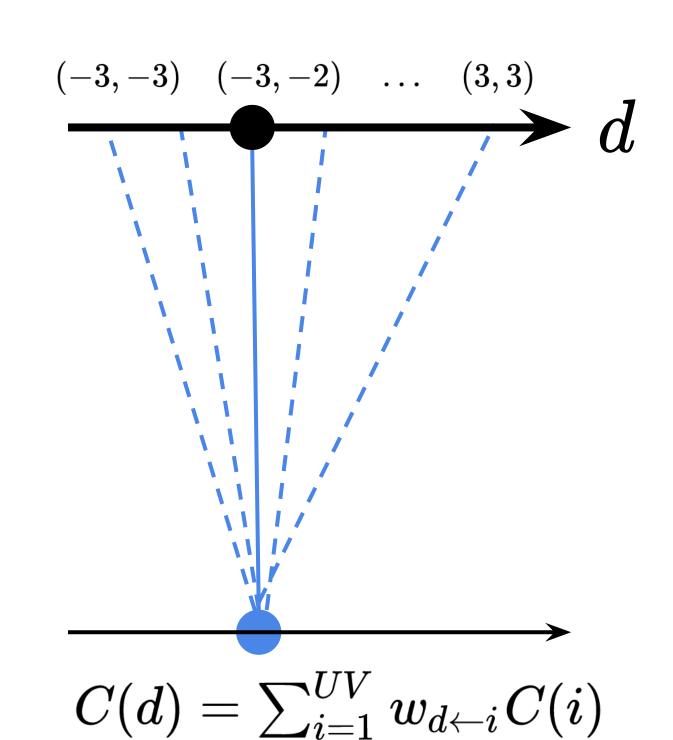
- ·Efficient higher-dimension (4D) cost volume processing
- Separable volumetric filters: reduce computation and parameters
- Multi-channel cost volumes: capture multiple dimensions of pixel similarity
- Adaptive cost volumes allow networks to generalize across tasks: train and test a single network for both flow+stereo
- · Results
- SOTA accuracy on optical flow benchmarks
- Training converges in up to 10X fewer iterations than prior art

# Approach

### Key idea: 4D filters

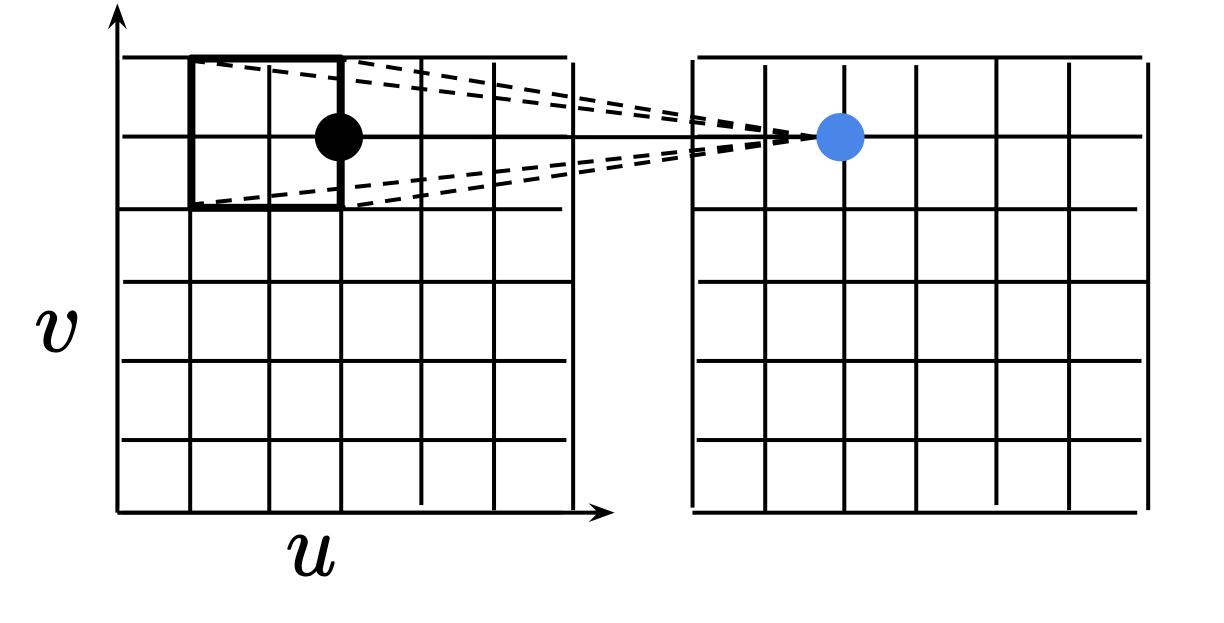
#### Prior method

- 1) reshape to  $C(d,x,y) \in R^{UV imes H imes W}$
- 2) multi-channel 2D convs with  $W \in R^{UV imes 3 imes 3}$



#### Ours

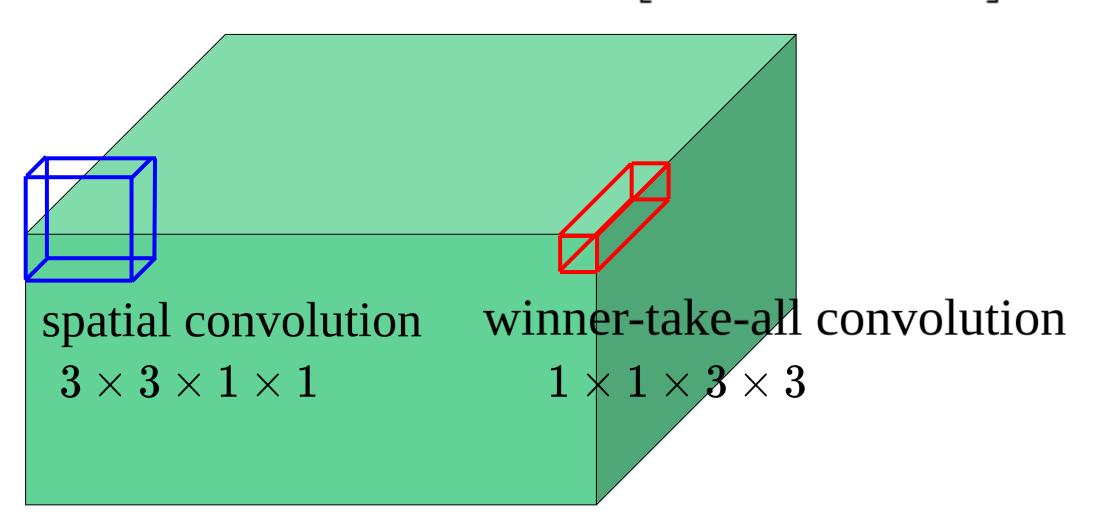
- 1) keep  $C(u,v,x,y) \in R^{U imes V imes H imes W}$
- 2) 4D convs with  $W \in R^{3 \times 3 \times 3 \times 3}$



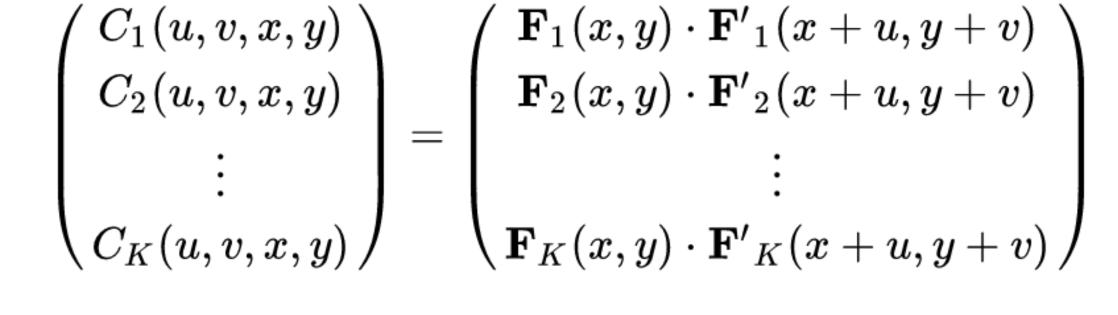
 $C(u,v) = \sum_{i=1}^{3} \sum_{j=1}^{3} w_{(u,v) \leftarrow (i,j)} C(i,j)$ 

### Separable filters

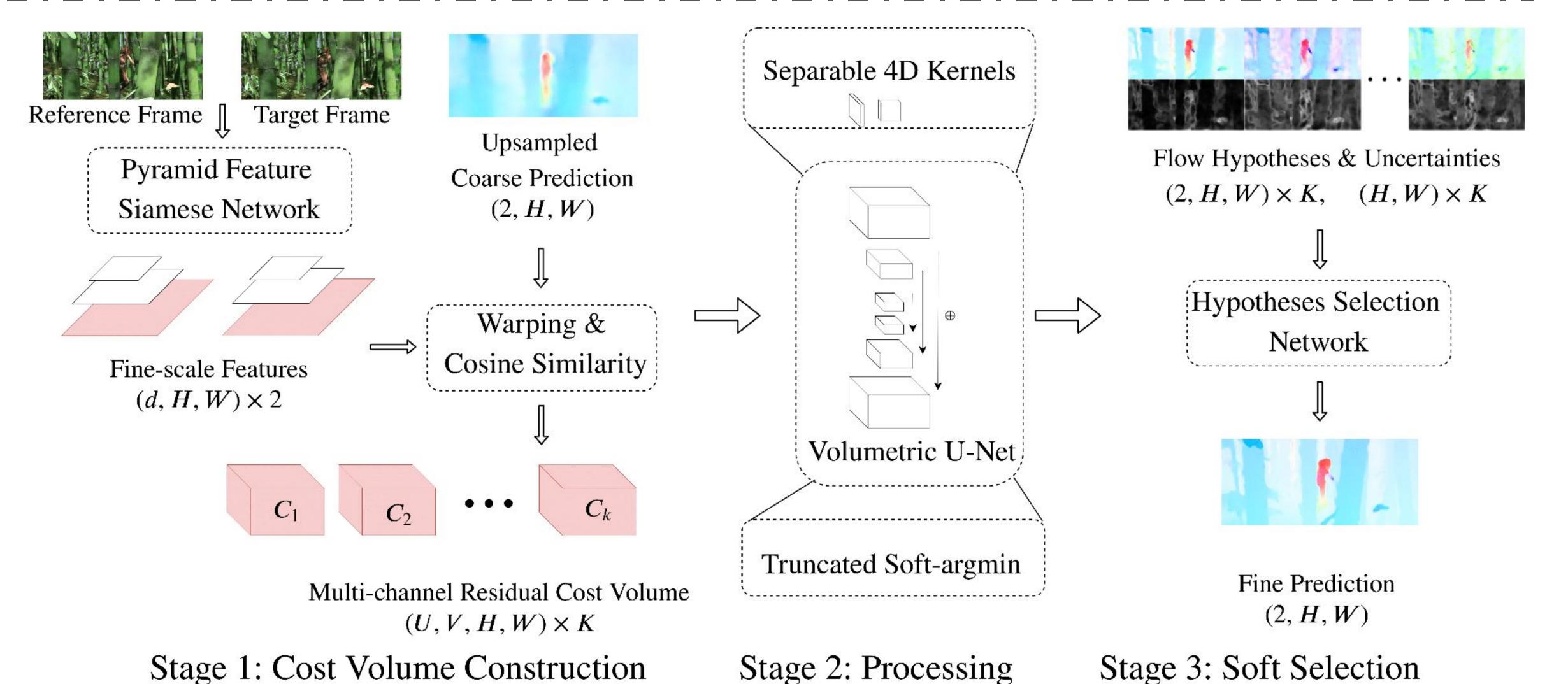
### $K(\mathbf{u}, \mathbf{x}) * C(\mathbf{u}, \mathbf{x}) = K_{WTA}(\mathbf{u}) * \left| K_S(\mathbf{x}) * C(\mathbf{u}, \mathbf{x}) \right|$



### Multi-channel 4D cost volumes







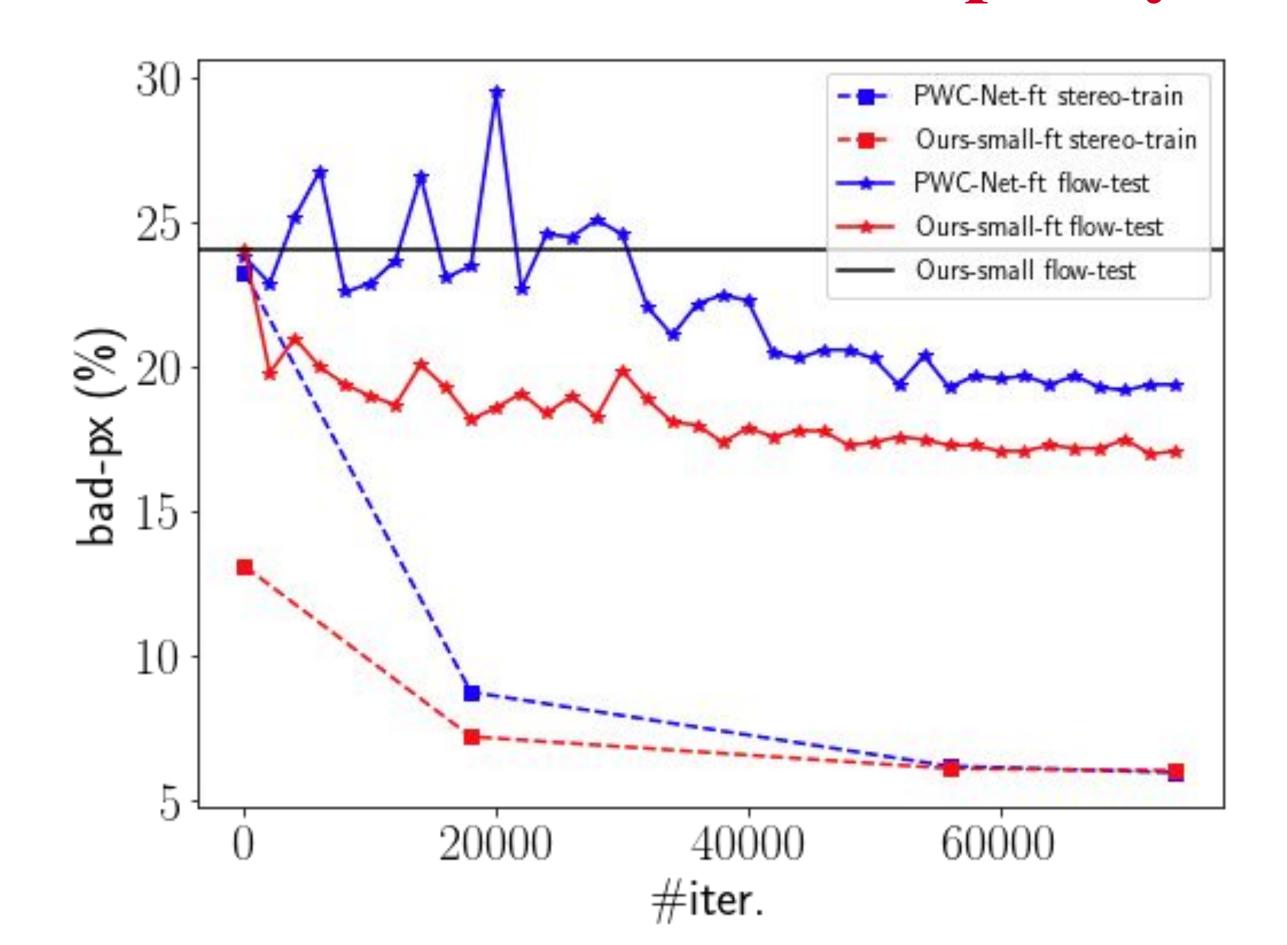
# Experiments

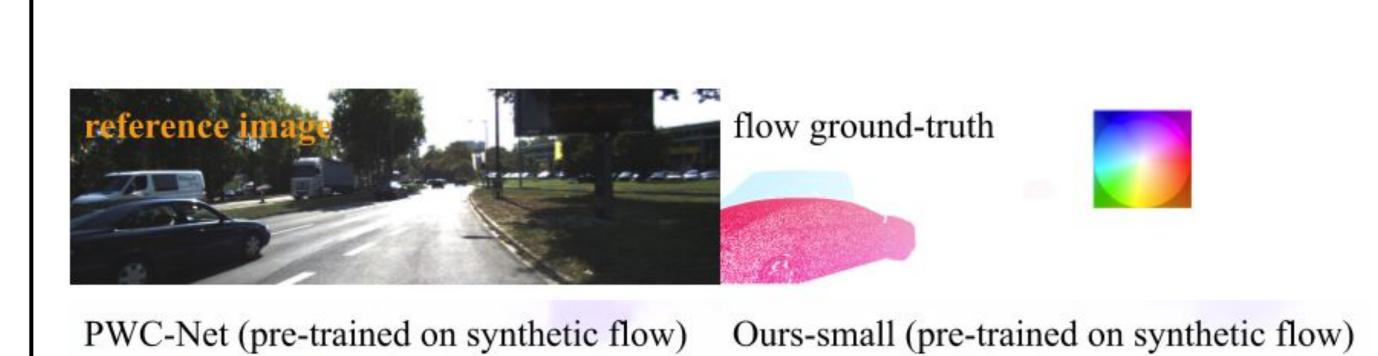
#### Benchmarks

Train	Method	K-15-train		K-15-test		S-train (epe) $\downarrow$		S-test (epe) $\downarrow$		Param.	GFlops	Iter.
dataset		Fl-epe↓	Fl-all↓	Fl-all↓	$\mathrm{D}1\text{-}\mathrm{all}^\dagger\downarrow$	Clean	Final	Clean	Final		СТЮрь	1001.
Pre-train: Chairs and Things	FlowNet2	10.08	30.0	_	-	2.02	3.54	3.96	6.02	162.5M	368.3	7100K
	PWC-Net	10.35	33.7	-	23.30	2.55	3.93	10 <del>-10</del>	-	8.8M	101.6	1700K
	$\mathrm{HD}^{\wedge}\mathrm{3F}$	13.17	24.0	<del></del>	-	3.84	8.77	· <del>-</del>	1 <del>17</del> 8	39.6M	174.8	: <del></del> -
	Ours-small	9.43	33.4	<del></del>	13.12	2.45	3.63	10 <del>70</del>	1 <del>77</del> 4	5.6M	41.0	220K
	Ours-full	8.36	25.1	: <del>-</del> :%	8.73	2.21	3.62	227	1 <del></del>	6.2M	101.7	220K
$\begin{array}{c} \text{Fine-tune:} \\ \text{K(ITTI)} \\ \text{or} \\ \text{S(intel)} \end{array}$	FlowNet2	(2.30)	(8.6)	11.48	-	(1.45)	(2.01)	4.16	5.74	162.5M	368.3	+500K
	PWC-Net+	(1.50)	(5.3)	7.72	9.17	(1.71)	(2.34)	3.45	4.60	8.8M	101.6	+750K
	LiteFlowNet2	(1.47)	(4.8)	7.74	-	(1.30)	(1.62)	3.45	4.90			: <del></del> -
	IRR-PWC	(1.63)	(5.3)	$7.65_{3}$	-	(1.92)	(2.51)	3.84	4.58	6.4M		+750K
	$\mathrm{HD}^{\wedge}\mathrm{3F}$	(1.31)	(4.1)	$6.55_{2}$	-	(1.87)	(1.17)	4.79	4.67	39.6M	174.8	-
	Ours-small	(1.41)	(5.5)	7.74	6.10	(1.84)	(2.44)	3.26	4.73	5.6M	41.0	+140K
	Ours-full	(1.16)	(4.1)	$6.30_{1}$	4.67	(1.66)	(2.24)	$2.81_{1}$	$4.40_{1}$	6.2M	101.7	+140K

\* Metrics: Fl-epe is the average end-point (L2) error of optical flow vectors; Fl-all and D1-all are percentage of flow/stereo predictions with error less than some threshold.

#### Generalize: 1D to 2D disparity

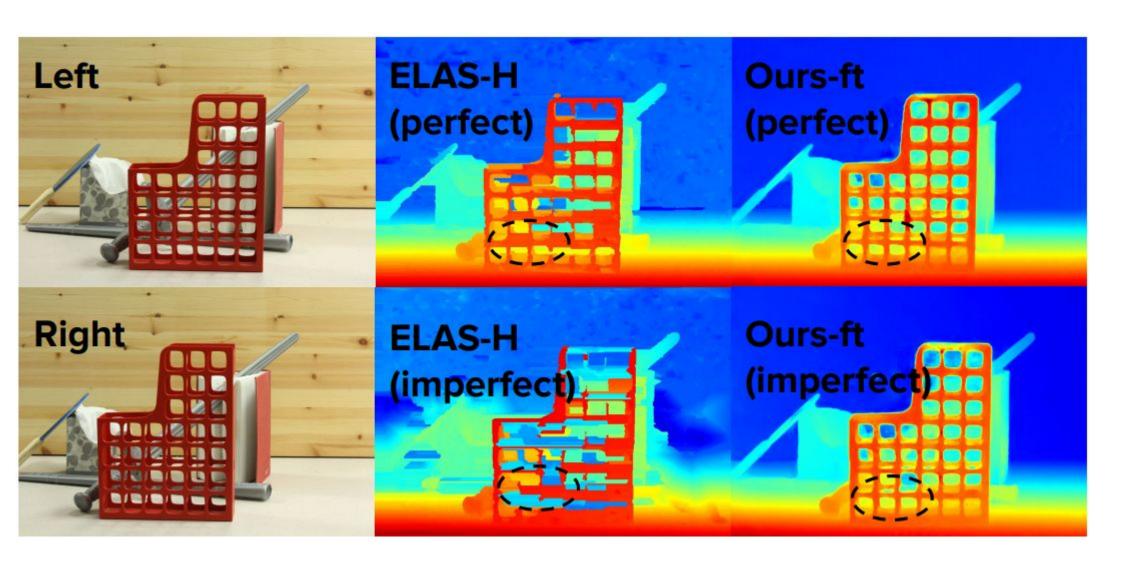




PWC-Net-ft (fine-tuned on real stereo)	Ours-small-ft (fine-tuned on real stereo)

## Application: stereo matching with imperfect rectification

Method	avge	inc.(%)	
	perfect	imperfect	
SGBM2	14.51	15.89	9.5
ELAS	9.89	11.79	19.2
PWC-Net	9.41	9.92	5.4
Ours	9.03	8.79	-2.7



### Ablation

	Method	EPE $(px) \downarrow$	GFlops	# Params.
(2)	DenseNet-2D	2.64	25.5	8.2
	Ours-full-4D	2.30	52.5	1.83
	Ours-sep-4D	2.31	23.4	1.78
	Ours-final	1.73	28.5	2.94
<u> </u>	Ours-K = 1	2.05	27.8	2.94