

# Convolutional Hough Matching Networks

Juhong Min Minsu Cho  
POSTECH CSE & GSAI

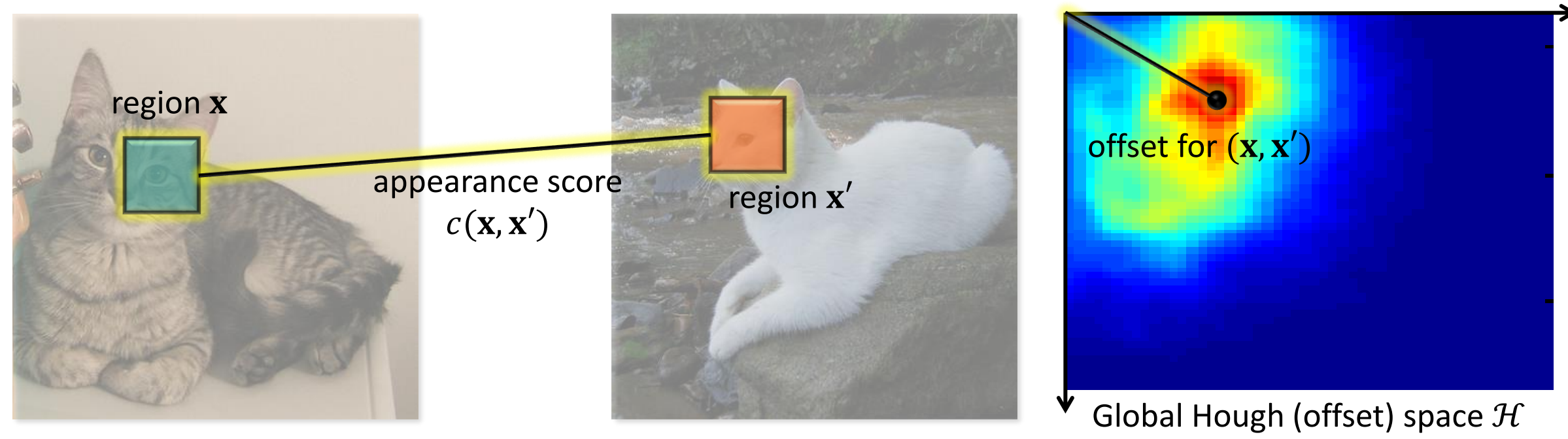


## Problem Introduction & Contribution

- Semantic visual correspondence:**  
Matching images depicting different instances of the same object class  
Visual correspondence “in the wild,” generalizing to other matching tasks  
Core component for many vision tasks, *e.g.*, tracking, retrieval, etc.
- Our contributions:**  
Introduce a **Hough transform perspective** on convolutional matching  
Develop trainable CHM layer with **semi-isotropic high-dimensional** kernel  
Propose CHMNet with a **small number of interpretable parameters**  
**SOTA** on three standard benchmarks of semantic correspondence

## Hough Matching (HM)

Hough matching is the algorithm of *Cho et al.* (CVPR 2015) which reweights appearance similarity by Hough ‘voting’ to enforce geometric consistency



$$v(\mathbf{h}) = \sum_{(\mathbf{x}, \mathbf{x}') \in \mathcal{X} \times \mathcal{X}'} c(\mathbf{x}, \mathbf{x}') k_{\text{iso}}(\|\mathbf{x}' - \mathbf{x} - \mathbf{h}\|_g)$$

Appearance score, *e.g.*, cosine similarity

Kernel that assigns a voting score according to how close  $(\mathbf{x}' - \mathbf{x})$  is to  $\mathbf{h}$

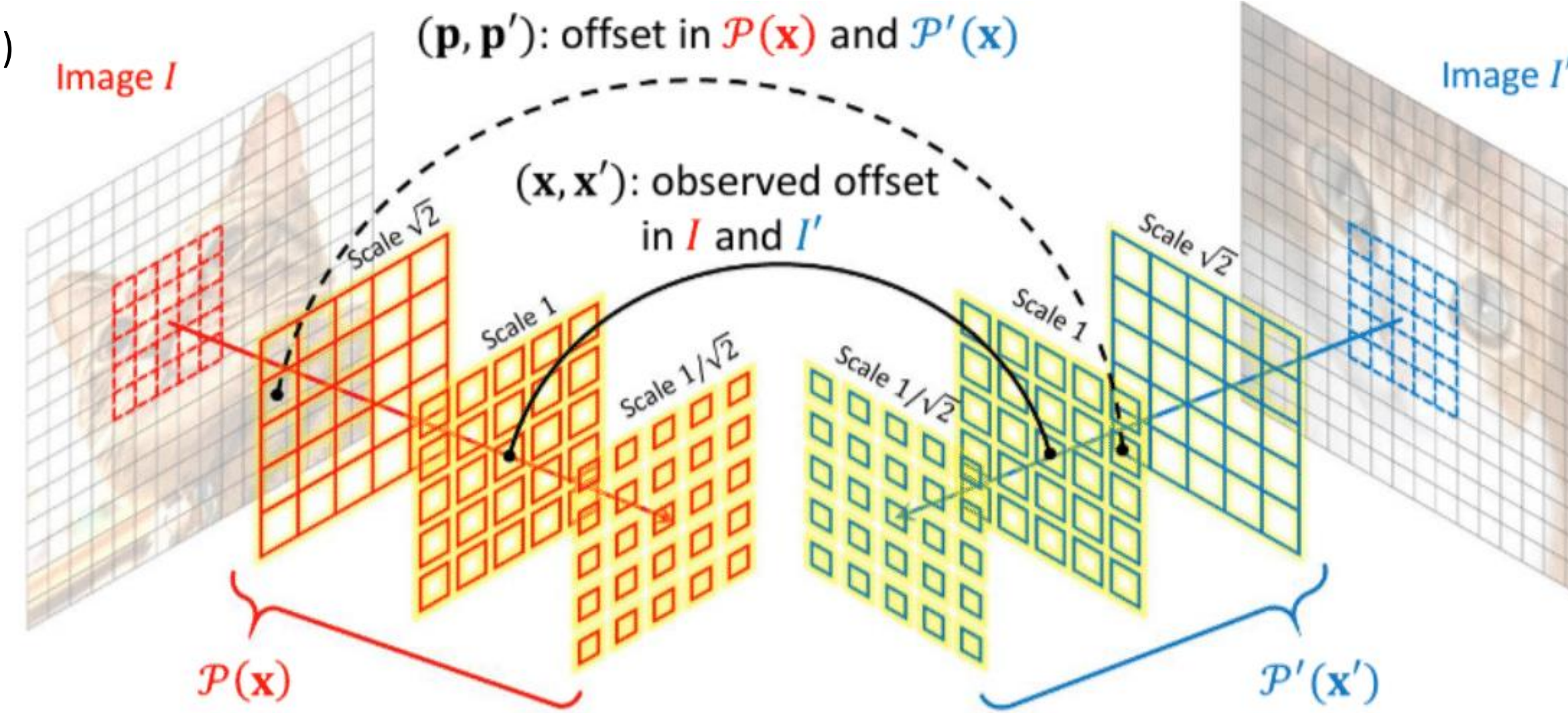
Distance between observed offset  $(\mathbf{x}' - \mathbf{x})$  and the given offset  $\mathbf{h}$  in the Hough space

- Limitation:** The Hough space is shared for all candidate matches so it cannot capture the reliability of a specific candidate matches, thus being less accurate and weak to background clutters.
- Solution:** to create a **local & individual** voting space for each match, *i.e.*, convolutionalization of the Hough matching algorithm.

## Convolutional Hough Matching (CHM)

(global & shared) (local & individual)  
**HM**  $\rightarrow$  **CHM**

- We introduce local windows,  $\mathcal{P}(\mathbf{x})$  and  $\mathcal{P}'(\mathbf{x}')$ , around regions  $\mathbf{x}$  and  $\mathbf{x}'$ .
- The local voting space is now dedicated to  $(\mathbf{x}, \mathbf{x}')$ .
- Let  $k(\mathbf{z}, \mathbf{z}')$  represent kernel value corresponding to two positions,  $\mathbf{z}$  and  $\mathbf{z}'$ .

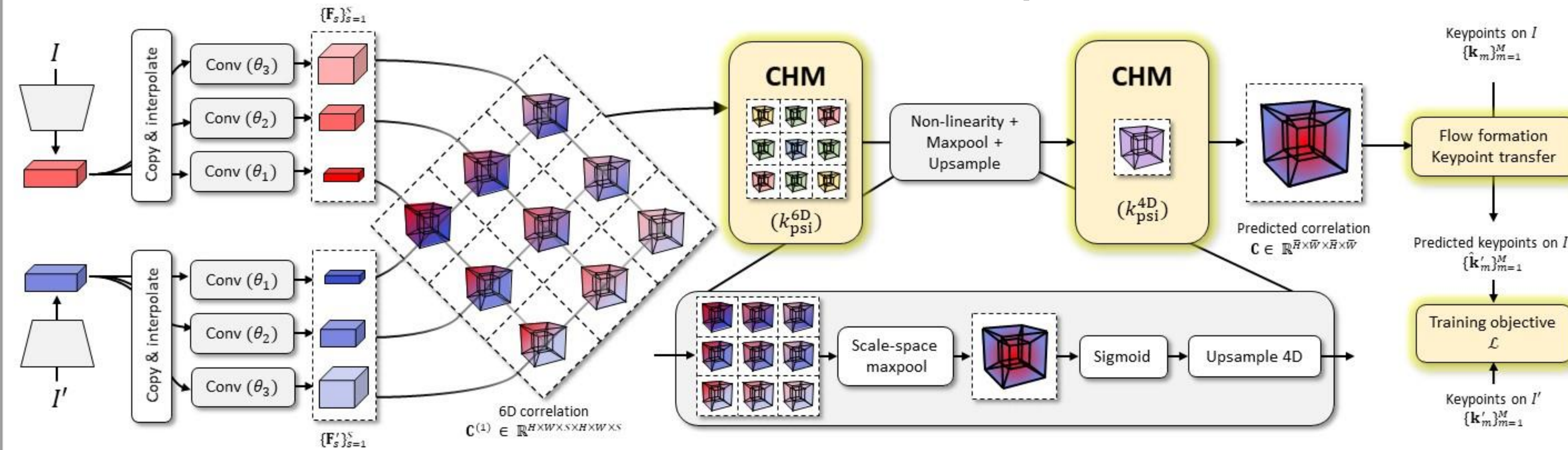


- Then we have,  $c_{\text{CHM}}(\mathbf{x}, \mathbf{x}') = \sum_{(\mathbf{p}, \mathbf{p}') \in \mathcal{P}(\mathbf{x}) \times \mathcal{P}'(\mathbf{x}')} c(\mathbf{p}, \mathbf{p}') k(\mathbf{p} - \mathbf{x}, \mathbf{p}' - \mathbf{x}')$  (We refer the readers to [the paper](#) for a complete derivation of CHM)

**High-dimensional convolution on correlation tensor as local & individual Hough matching**

We further relax isotropy and propose position-sensitive isotropic kernel  $k_{\text{psi}}$  which shares parameters whose triplets  $(\|\mathbf{p}' - \mathbf{p}\|_g, \|\mathbf{p} - \mathbf{x}\|_g, \|\mathbf{p}' - \mathbf{x}'\|_g)$  are the same.

- Advantage of CHM over existing 4D convolutions on correlation:**
  - Generalizability:** voting space can be extended to higher dim. beyond 4D, *e.g.*, 6D (translation & scale)
  - Scalability:** channel size of 1 & parameter sharing  $\rightarrow$  small number of parameters
  - Interpretability:** a single kernel for each layer eases kernel visualization
  - Performance:** state-of-the-art performance on three standard benchmark datasets
- Convolutional Hough matching networks** (The CHM part with  $k_{\text{psi}}$  has 275 parameters):

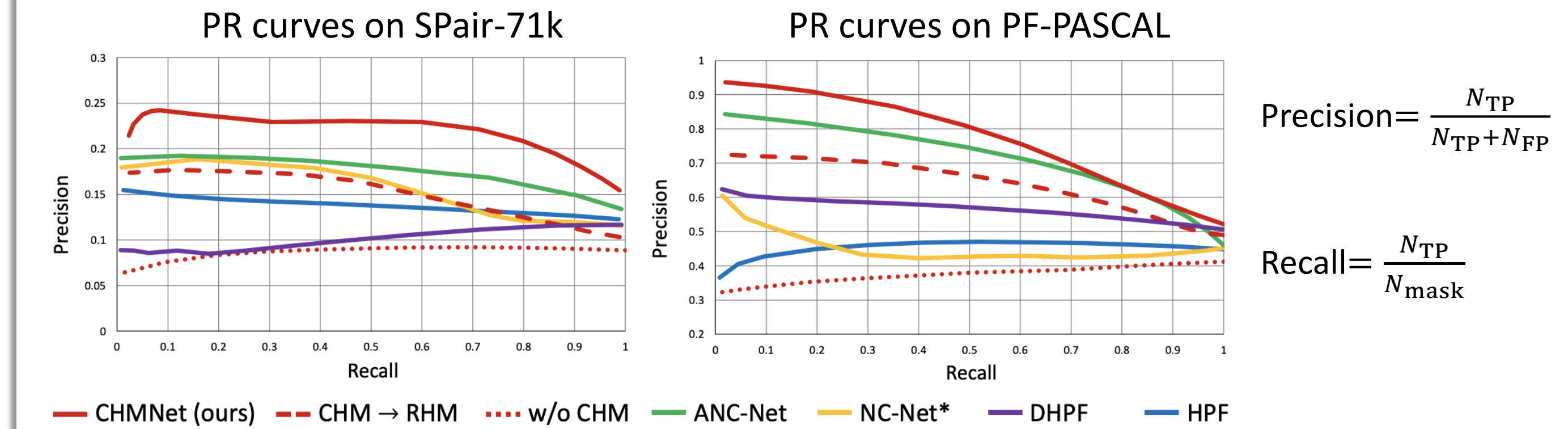


- Flow formation & keypoint transfer:** soft-argmax with soft sampler
- Objective:** minimizes L2-distance between predicted and GT keypoints:  $\mathcal{L} = \frac{1}{M} \sum_{m=1}^M \|\hat{\mathbf{k}}'_m - \mathbf{k}'_m\|$

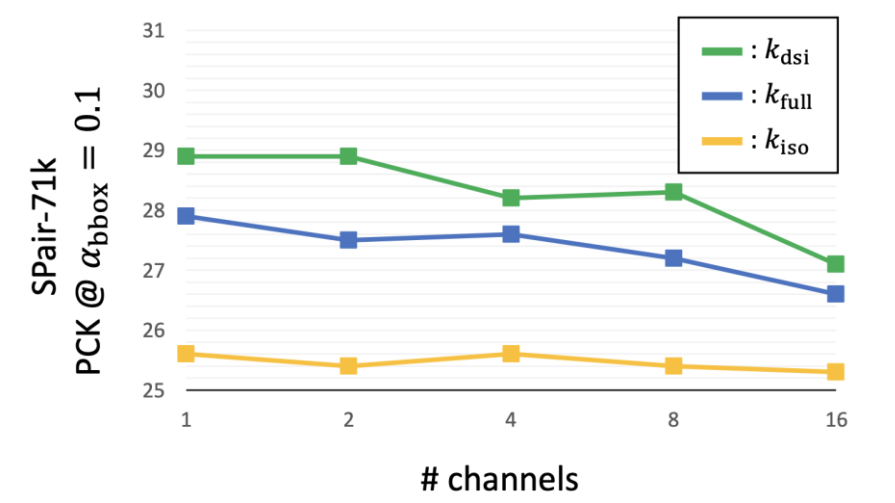
## Experimental Results and Analyses

- Evaluation results on semantic correspondence benchmarks:**

Methods	SPair-71k		PF-PASCAL		PF-WILLOW		uses nD conv?	FLOPs (G)	time (ms)	memory (GB)
	PCK @ $\alpha_{\text{bbox}}$ 0.1 (F)	0.1 (T)	PCK @ $\alpha_{\text{img}}$ 0.05	0.1	PCK @ $\alpha_{\text{bbox}}$ 0.05	0.1				
UCN <sub>res101</sub> (NeurIPS'16)	-	17.7	-	75.1	-	-	<b>X</b>	-	-	-
HPF <sub>res101</sub> (ICCV'19)	28.2	-	60.1	84.8	45.9	74.4	<b>X</b>	-	63	-
SCOT <sub>res101</sub> (CVPR'20)	35.6	-	63.1	85.4	47.8	76.0	<b>X</b>	6.2	151	4.6
DHPF <sub>res101</sub> (ECCV'20)	<u>37.3</u>	27.4	<u>75.7</u>	<u>90.7</u>	49.5	77.6	<b>X</b>	<b>2.0</b>	<u>58</u>	1.6
NC-Net <sub>res101</sub> (NeurIPS'18)	-	-	-	81.9	-	-	4D	44.9	222	1.2
DCC-Net <sub>res101</sub> (ICCV'19)	-	-	-	83.7	-	-	4D	47.1	567	2.7
ANC-Net <sub>res101</sub> (CVPR'20)	-	<u>28.7</u>	-	86.1	-	-	4D	44.9	216	<b>0.9</b>
CHMNet <sub>res101</sub> (ours)	<b>46.3</b>	<b>30.1</b>	<b>80.1</b>	<b>91.6</b>	<b>52.7</b>	<u>79.4</u>	6D	19.6	<b>54<sup>†</sup></b> (248)	1.6



- Channel size experiments:** High-dimensional convolution on a correlation tensor is to learn a *reliable voting strategy* rather than to capture diverse patterns in the correlation tensor.



- Qualitative results & learned kernel ( $k_{\text{psi}}$ ) visualization:**

