# Recurrent Transformer Networks for Semantic Correspondence

**Seungryong Kim**[1], Stepthen Lin[2], Sangryul Jeon[1], Dongbo Min[3], Kwanghoon Sohn[1]

Dec. 05, 2018

1) YONSEI UNIVERSITY   2) Microsoft Research   3) 이화여자대학교 EWHA WOMANS UNIVERSITY
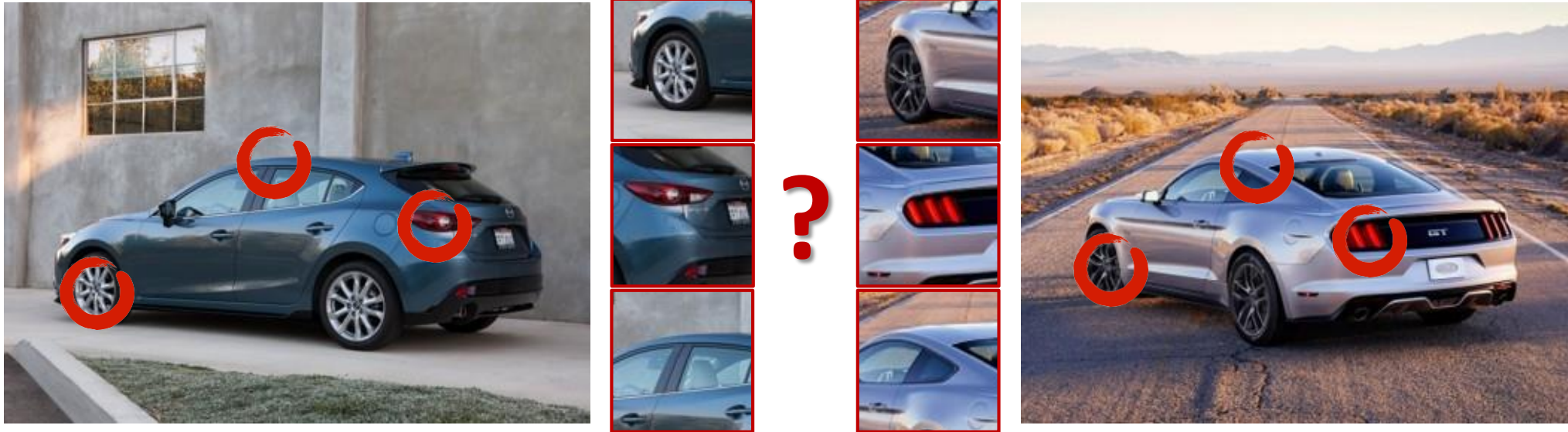
# Introduction

**Semantic correspondence**



- Establishing ***dense correspondences*** between ***semantically similar images***, i.e., different instances within the same object or scene categories
- *For example, the wheels of two different cars, the body of people or animals*

# Introduction

## Challenges in semantic correspondence



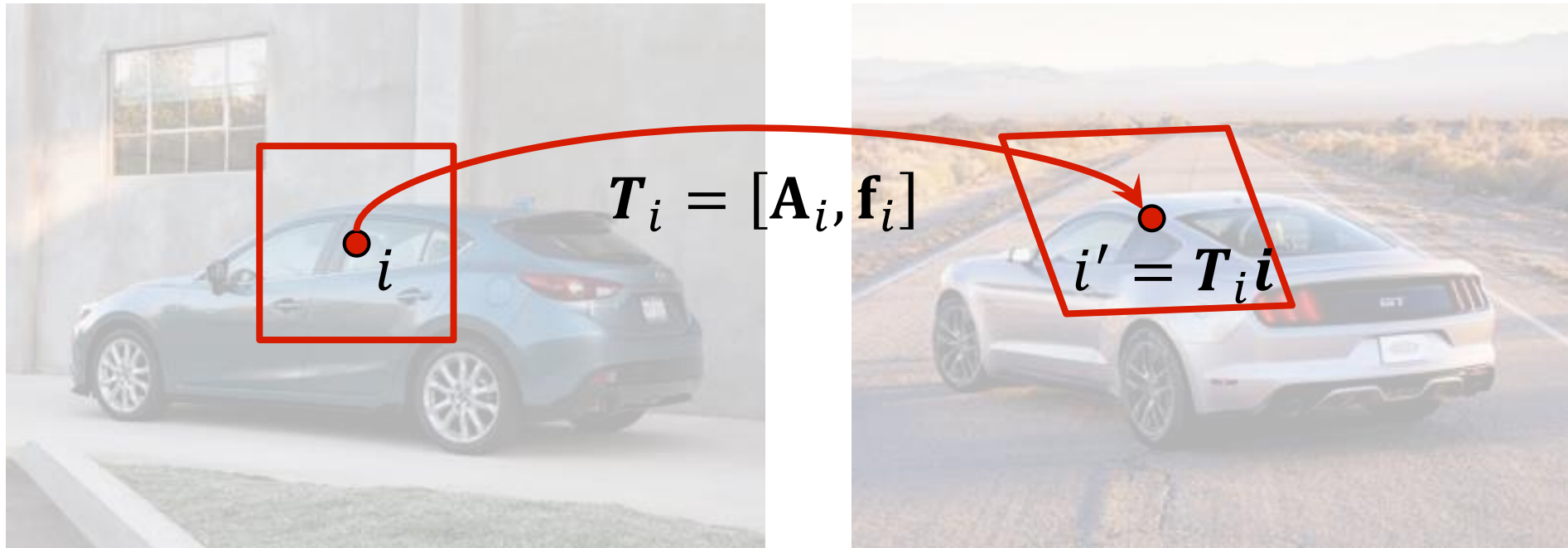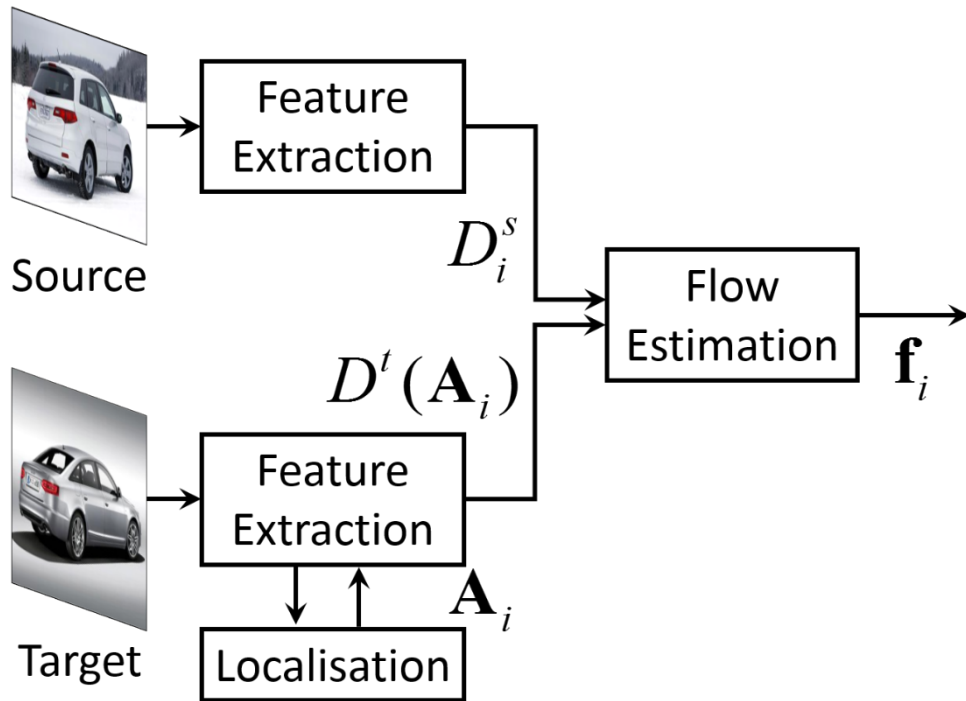| Photometric Deformations | Geometric Deformations | Lack of Supervision |
|---|---|---|
| • *Intra-class appearance and attribute variations*<br>• *Etc.* | • *Different viewpoint or baseline*<br>• *Non-rigid shape deformations*<br>• *Etc.* | • *Labor-intensive of annotation*<br>• *Degraded by subjectivity*<br>• *Etc.* |

# Problem Formulation

**Objective: locally-varying affine transformation fields**



$$T_i = [\mathbf{A}_i, \mathbf{f}_i]$$

$i$

$i' = T_i i$

$\rightarrow$ *How to estimate* <u>*dense affine transformation fields*</u> *between semantically similar images?*

# Related Works

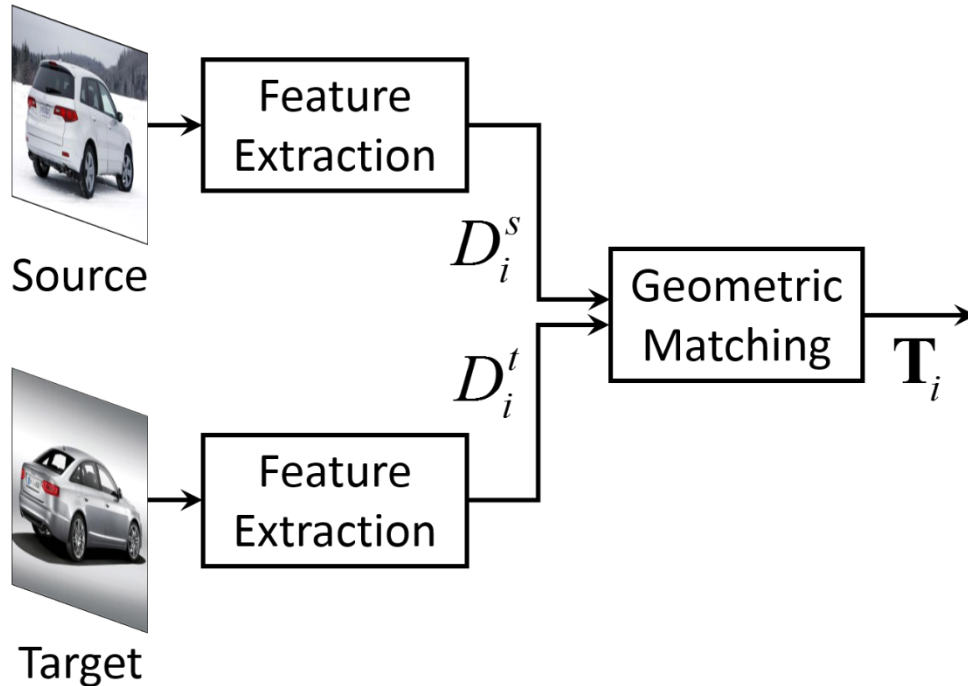## Methods for geometric invariance in the feature extraction



- **UCN** [Choy *et al.*, NIPS'16]
- **CAT-FCSS** [Kim *et al.*, TPAMI'18]
- Etc.

✓ Spatial Transformer Networks (STNs)-based methods [Jaderberg et al., NIPS'15]

✓ $\mathbf{A}_i$ is learned wo/$\mathbf{A}_i^*$

✗ But, $\mathbf{f}_i$ is learned w/$\mathbf{f}_i^*$
✗ Geometric inference based on only source or target image

# Related Works

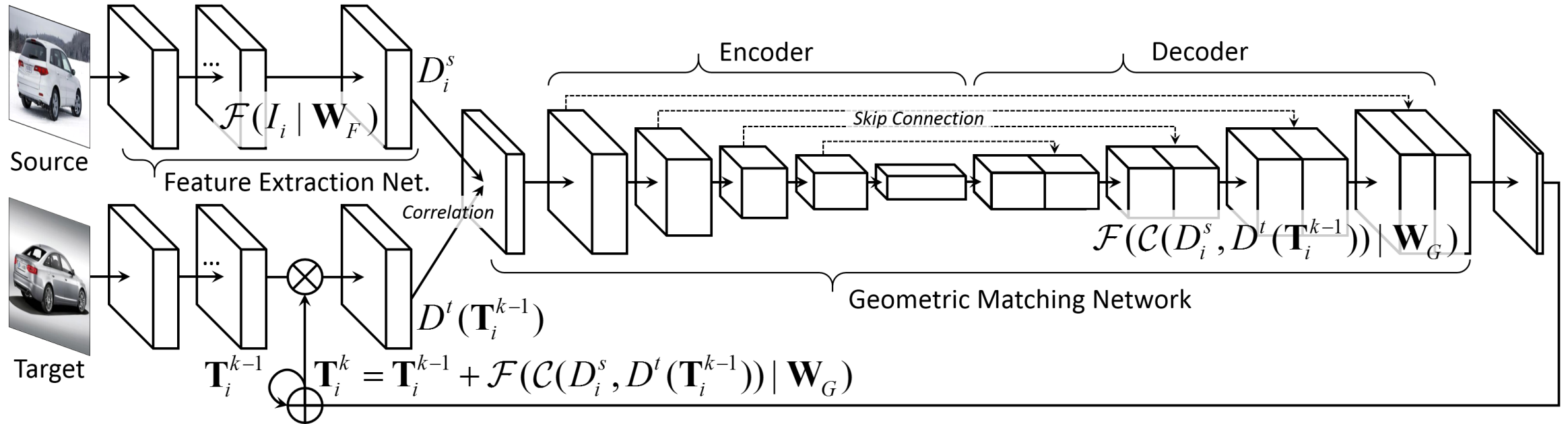## Methods for geometric invariance in the regularization



- **GMat.** [Rocco *et al.*, CVPR'17]
- **GMat. w/Inl.** [Rocco *et al.*, CVPR'18]
- Etc.

✓ $\mathbf{T}_i$ is learned wo/$\mathbf{T}_i^*$ using self- or meta-supervision

✓ Geometric Inference using source/target images

✗ Globally-varying geometric Inference only

✗ only fixed, untransformed versions of the features
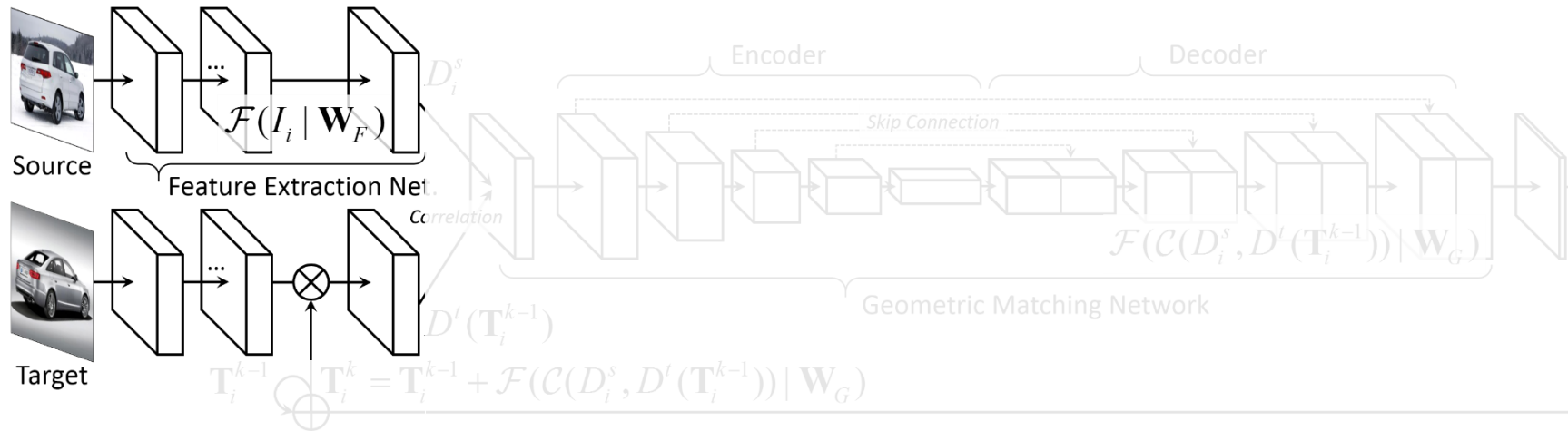
# Recurrent Transformer Networks (RTNs)

## Networks configuration



- Weaves the advantages of **STN-based methods** and **geometric matching methods** by recursively estimating geometric transformation residuals using geometry-aligned feature activations
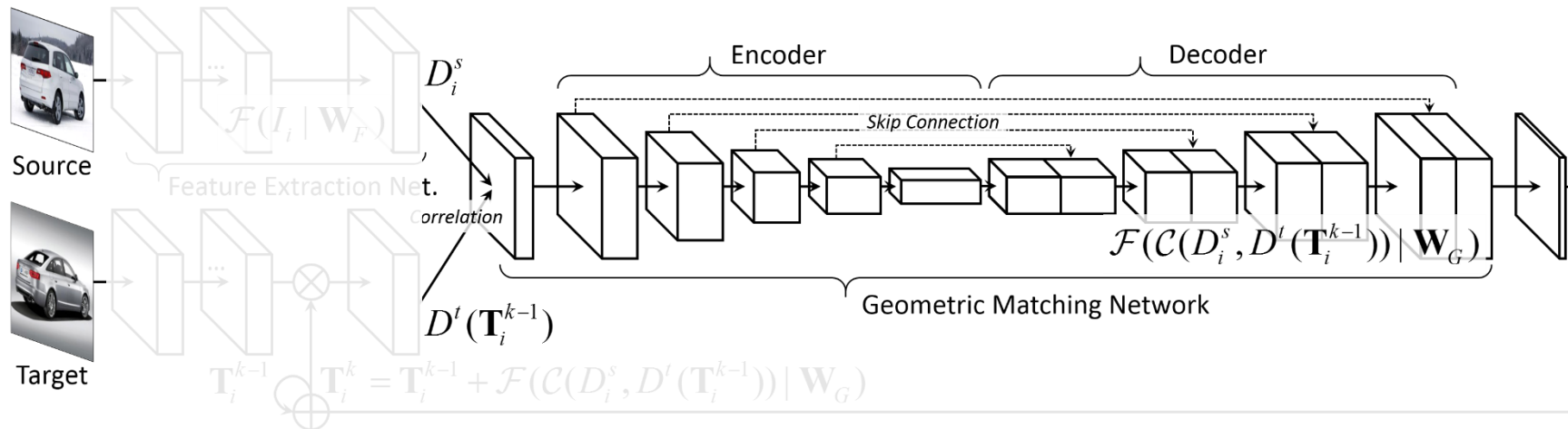
# Recurrent Transformer Networks (RTNs)

## Feature Extraction Networks



- To extract features $D^s$ and $D^t$, input images $I^s$ and $I^t$ are passed through convolution networks with parameters $\mathbf{W}_F$ such that $D_i = F(I|\mathbf{W}_F)$ using CAT-FCSS, VGGNet (conv4-4), ResNet (conv4-23)
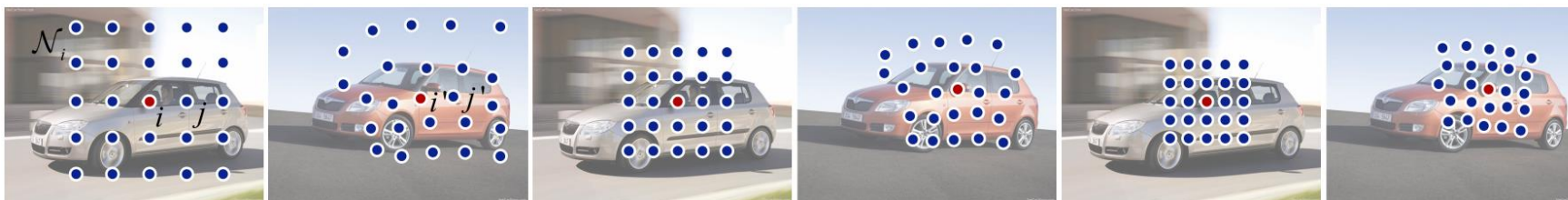
# Recurrent Transformer Networks (RTNs)

## Recurrent Geometric Matching Networks



- *Constraint correlation volume construction*

$$C(D_i^s, D^t(\mathbf{T}_j)) = < D_i^s, D^t(\mathbf{T}_j) > / \sqrt{< D_i^s, D^t(\mathbf{T}_j) >^2}$$
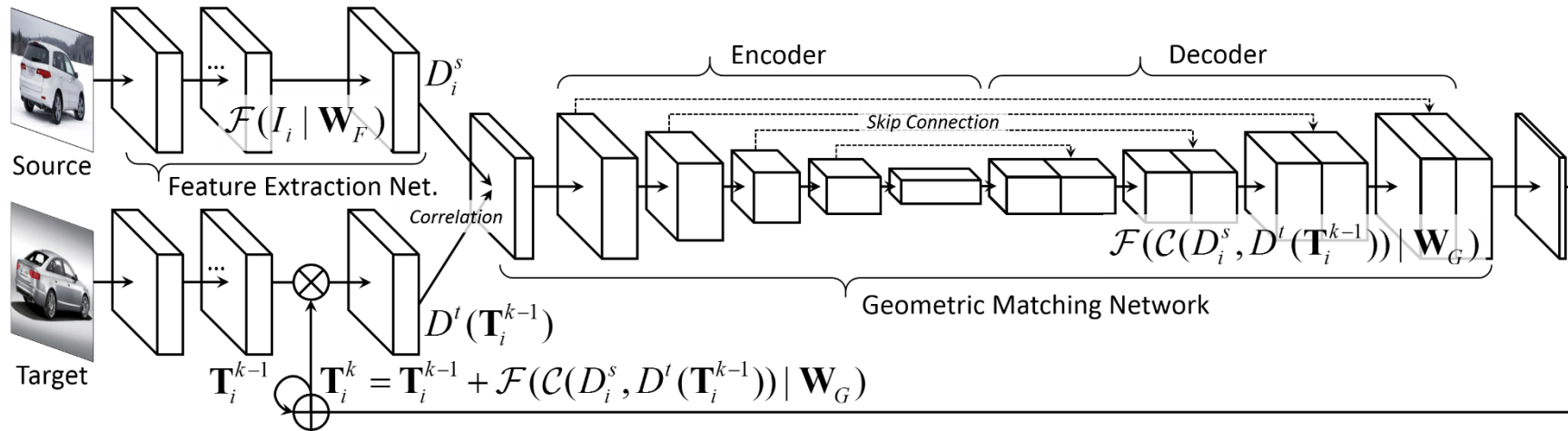


Search window 4          Search window 2          Search window 1

# Recurrent Transformer Networks (RTNs)

## Recurrent Geometric Matching Networks



- *Recurrent geometric inference*

$$\mathbf{T}_i^k - \mathbf{T}_i^{k-1} = F(C(D_i^s, D^t(\mathbf{T}_i^{k-1}))|\mathbf{W}_G)$$



| Source | Target | Iter. 1 | Iter. 2 | Iter. 3 | Iter. 4 |

# Recurrent Transformer Networks (RTNs)

## Weakly-supervised Learning

- **Intuition**: matching score between the source $D^s$ at each pixel $i$ and the target $D^t(T_i)$ should be maximized while keeping the scores of other candidates low!

- **Loss Function**:

$$L\left(D_i^s, D^t(T)\right) = -\sum_{j \in M_i} p_j^* \log(p(D_i^s, D^t(T_j)))$$

where the function $p(D_i^s, D^t(T_j))$ is a Softmax probability

$$p(D_i^s, D^t(T_j)) = \frac{\exp(C(D_i^s, D^t(T_j)))}{\sum_{l \in M_i} \exp(C(D_i^s, D^t(T_j)))}$$

where $p_j^*$ denotes a class label defined as 1 if $j = i$, 0 otherwise

# Experimental Results

## Results on the TSS Benchmark



| Source images | Target images | SCNet [Han *et al.*, ICCV'17] | GMat. w/Inl. [Rocco *et al.*, CVPR'18] | **RTNs** |

# Experimental Results

## Results on the PF-PASCAL Benchmark



| Source images | Target images | SCNet [Han *et al.*, ICCV'17] | GMat. w/Inl. [Rocco *et al.*, CVPR'18] | **RTNs** |

# Experimental Results

## Results on the PF-PASCAL Benchmark



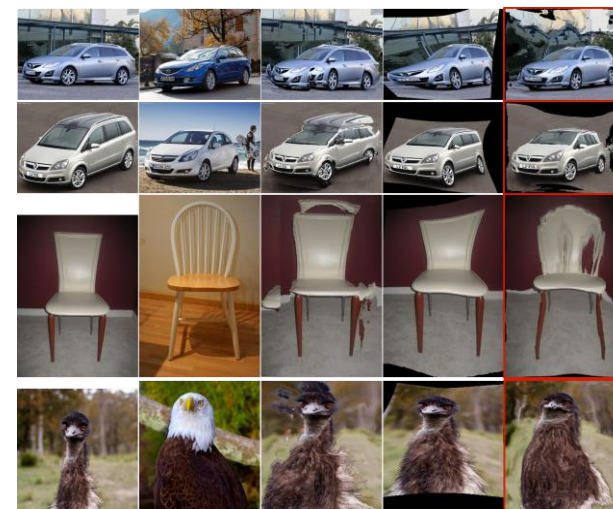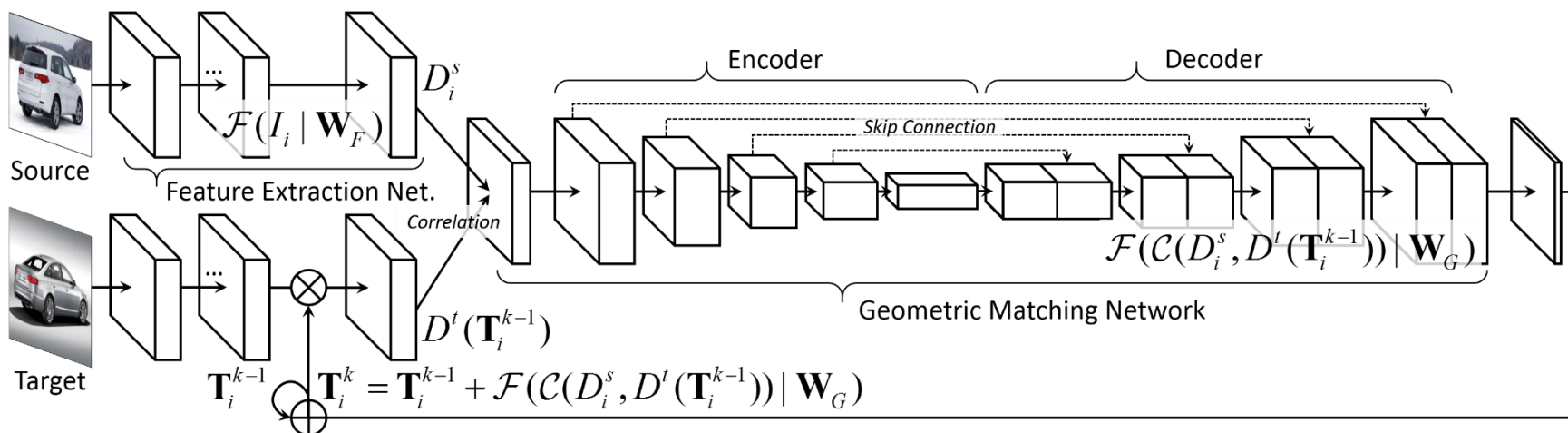| Source images | Target images | SCNet [Han *et al.*, ICCV'17] | GMat. w/Inl. [Rocco *et al.*, CVPR'18] | **RTNs** |

# Concluding Remarks

- RTNs learn to infer **locally-varying geometric fields** for semantic correspondence in an end-to-end and weakly-supervised fashion

- The key idea is to utilize and iteratively refine **the transformations and convolutional activations via matching** between the image pair

- A technique is presented for **weakly-supervised training** of RTNs

# Thank you!

## See you at 210 & 230 AB #119

Seungryong Kim, Ph.D.
Digital Image Media Lab.
Yonsei University, Seoul, Korea
Tel: +82-2-2123-2879
E-mail: srkim89@yonsei.ac.kr
Homepage: http://diml.yonsei.ac.kr/~srkim/