# Unconstrained Fashion Landmark Detection via Hierarchical Recurrent Transformer Networks

Sijie Yan, Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, Xiaoou Tang

The Chinese University of Hong Kong
香港中文大学

ACMmultimedia25

## Motivation



**(a) Constrained Fashion Landmark Detection**

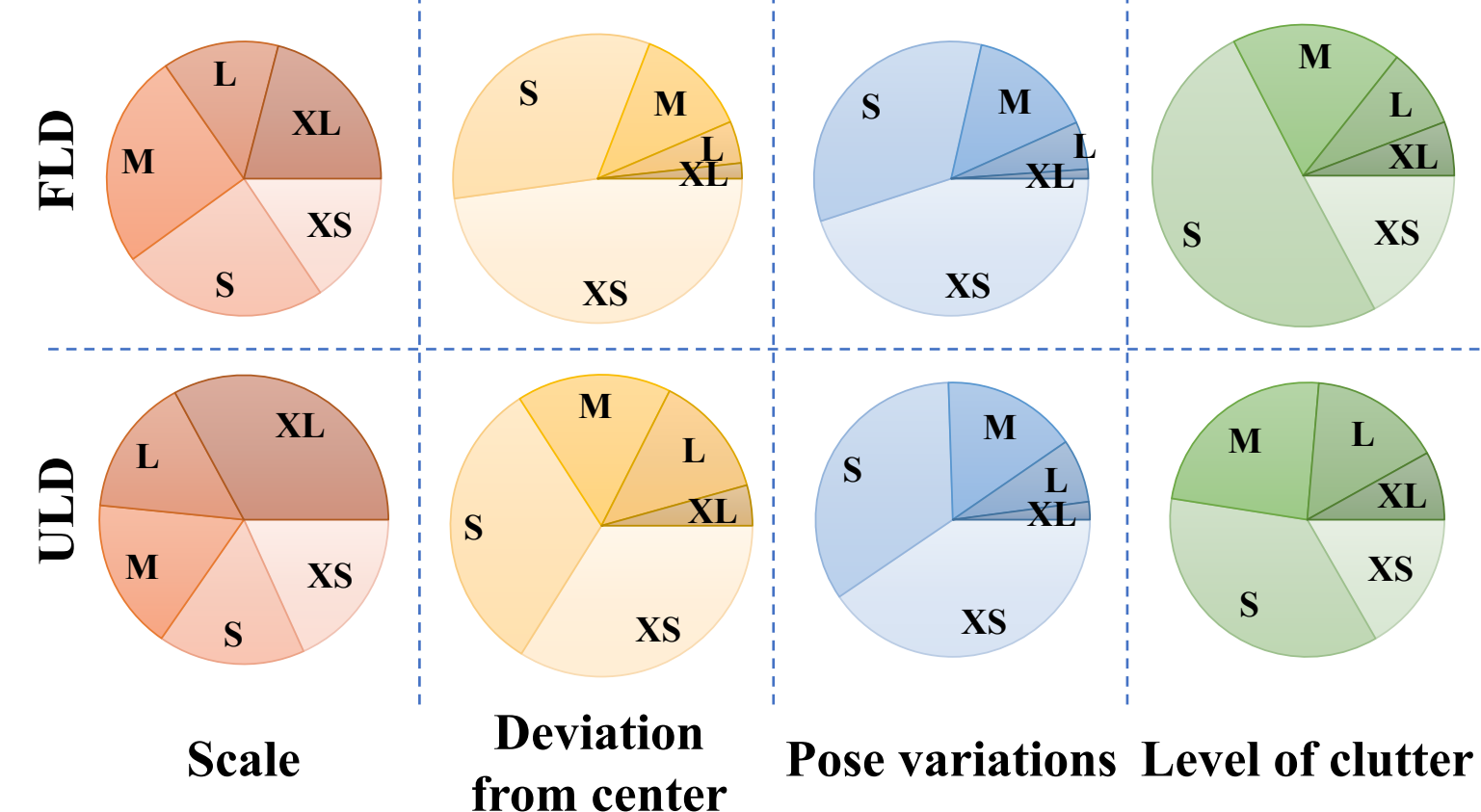**(b) Unconstrained Fashion Landmark Detection**

**Problem**

How to detect fashion landmarks without bounding boxes of clothes?

**Difficulty**

Background clutters, human poses, and scales variations due to:
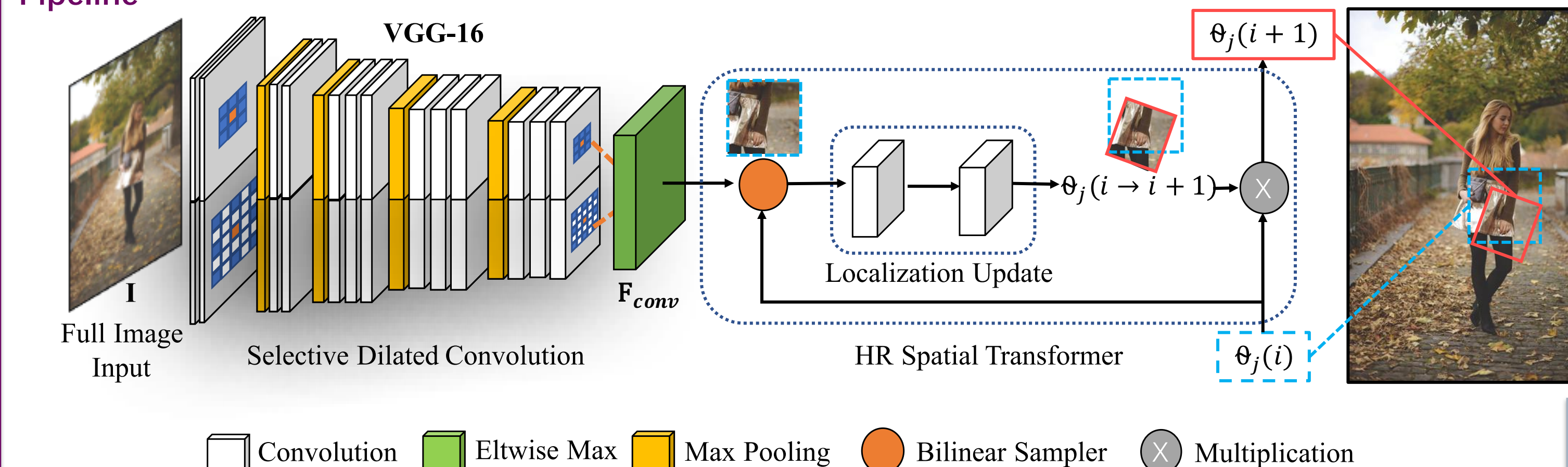- deformable objects,
- complex application scenarios.

## Dataset



Scale | Deviation from center | Pose variations | Level of clutter

**ULD v.s. FLD**

- 30K images with comprehensive fashion landmark annotations.
- collected from fashion blogs, forums and online shop.

## Approach

### Pipeline



**I** Full Image Input — **VGG-16** — Selective Dilated Convolution — $F_{conv}$ — HR Spatial Transformer — Localization Update — $\Theta_j(i \rightarrow i+1)$ — $\Theta_j(i+1)$ — $\Theta_j(i)$

□ Convolution  ▨ Eltwise Max  ▨ Max Pooling  ● Bilinear Sampler  ⊗ Multiplication

(a) $F_{conv}$ — STN — $\Theta$ — $F_{trans}$ — $\widehat{l_j'}$ — $\widehat{l_j} = \Theta \cdot \widehat{l_j'}$ — Original Coordinates

(b) $F_{conv}$ — Step 1 — HR-ST — $\Theta(1)$ — HR-ST — $\Theta_j(2)$ — HR-ST — $\Theta_j(3)$ — $\widehat{l_j} = \Theta_j(3) \cdot \widehat{l_j'}$ — $\widehat{l_j'}$ — $\widehat{l_j}$ — Original Coordinates

Step 1  Step 2  Step 3  Relative Coordinates

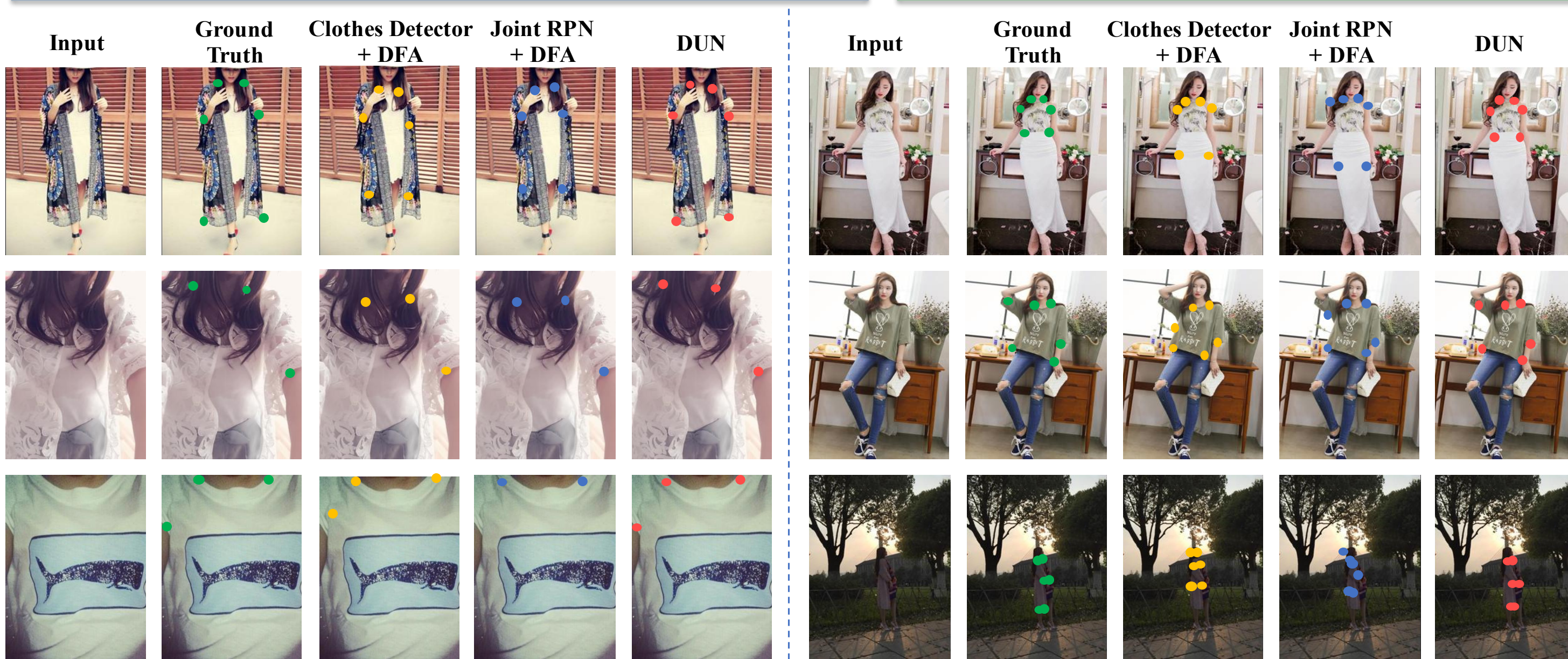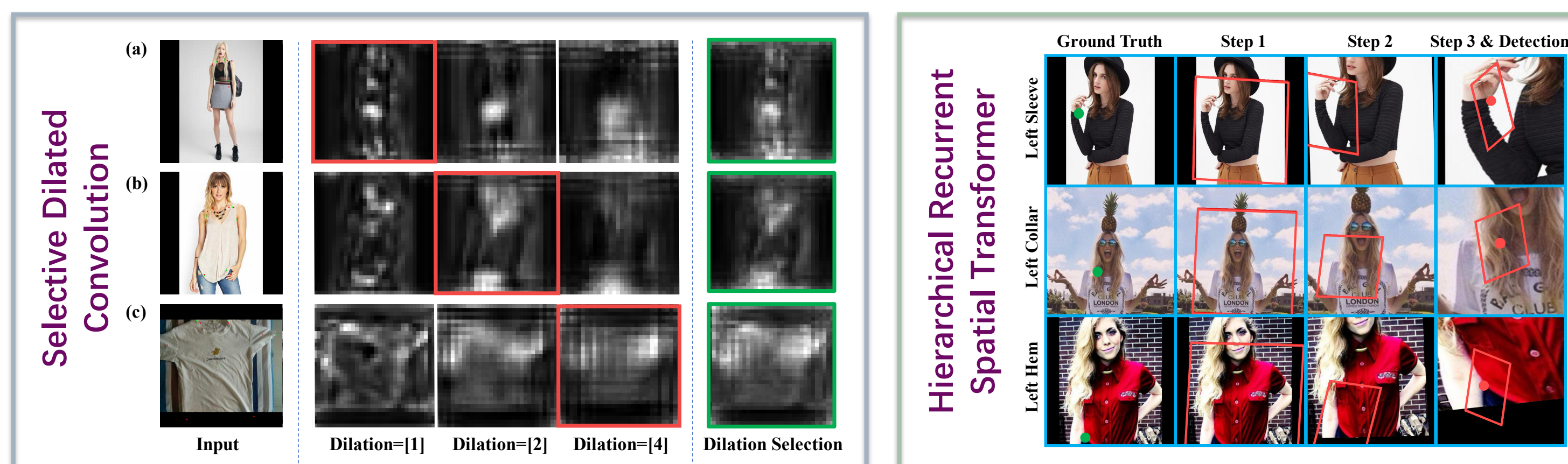▮ Inner Product  ⊗ Multiplication

### Selective Dilated Convolution

- Share weights
- element-wise max

$$F_{conv} = \max_s F_{conv5 \star 2^s}$$

- $F_{conv}$: Final convolutional response
- $s$: scale. $F_{conv5 \star 2^s}$: convolution
- $\star$: expanded sampling.

### Hierarchical Recurrent Spatial Transformer

- Recover coordinates $\quad \widehat{l_j} = \Theta_{global} \cdot \widehat{l_j'}$
- $\widehat{l_j'}$: relative landmark coordinates
- $\widehat{l_j}$: original coordinates.
- $\Theta_{global}$: geometric transformation for clothes
- Hierarchical model $\quad \Theta_j = \Theta_{global} \cdot \Theta_j$
- $\Theta_j$: geometric transformation for each landmark
- $\Theta_j(i)$: local refinement transformation for recurrent step $i$
- Recurrent update

$$\Theta_j(i) \leftarrow \Theta_j(i-1) \cdot \Theta_j(i-1 \rightarrow i)$$

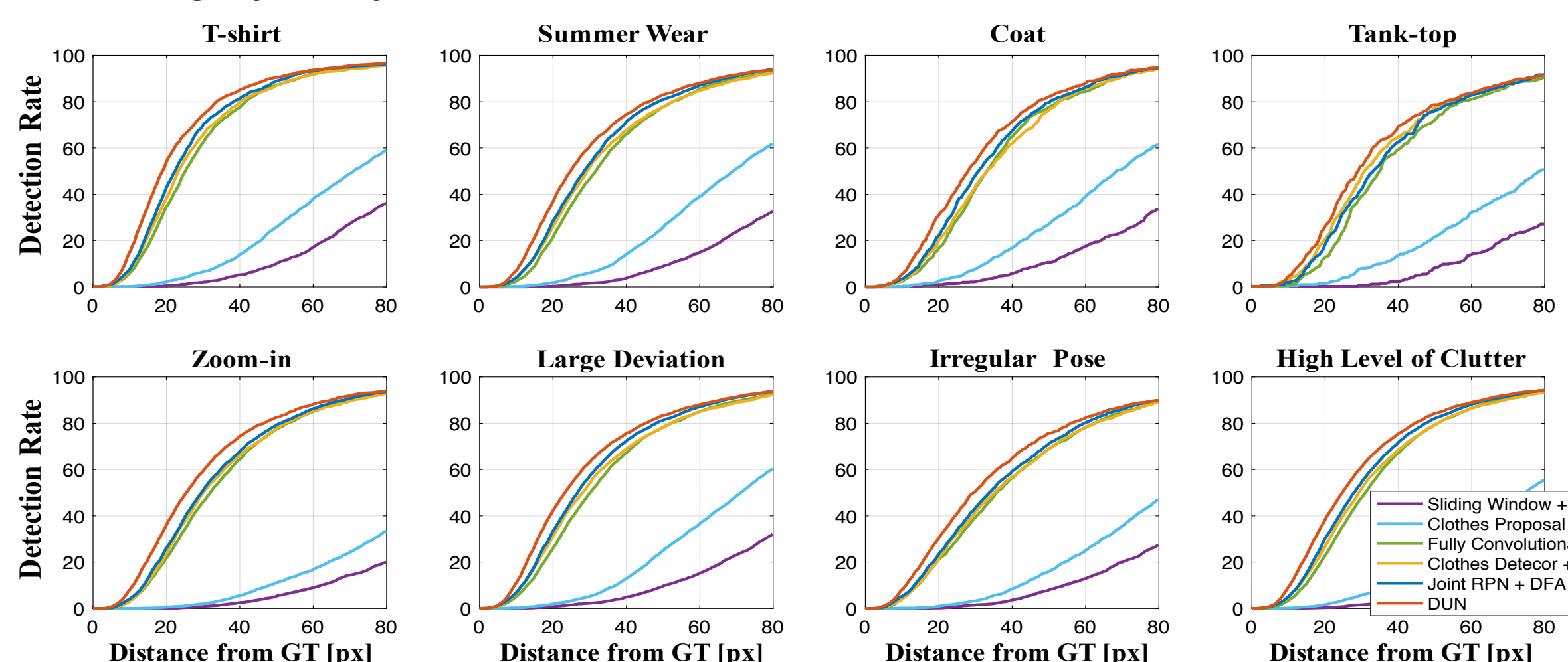- $\Theta_j(i-1 \rightarrow i)$: refinement transformation

## Visual Results



**Selective Dilated Convolution**

Input | Dilation=[1] | Dilation=[2] | Dilation=[4] | Dilation Selection

**Hierarchical Recurrent Spatial Transformer**

Ground Truth | Step 1 | Step 2 | Step 3 & Detection

Input | Ground Truth | Clothes Detector + DFA | Joint RPN + DFA | DUN

Input | Ground Truth | Clothes Detector + DFA | Joint RPN + DFA | DUN

## Experiments

### Comparison

| | # VGGs | # bbox anno. | end-to-end | # inference pass | speed (fps) | det. rate (%) |
|---|---|---|---|---|---|---|
| Sliding Window + DFA [18] | 1 | ✗ | ✗ | 17 | 3.2 | 2.7 |
| Clothes Proposal + DFA [18] | 1 | ✗ | ✗ | 100 | 0.5 | 9.7 |
| Clothes Detector + DFA [18] | 2 | 16K | ✗ | 1 | 5.0 | 63.1 |
| Joint RPN [20] + DFA [18] | 2 | 16K | ✓ | 1 | 3.9 | 66.0 |
| **Deep LAndmark Network** | 1 | ✗ | ✓ | 1 | 5.2 | 73.8 |

### Per-landmark Analysis

| | L. Collar | R. Collar | L. Sleeve | R. Sleeve | L. Hem | R. Hem | Mean |
|---|---|---|---|---|---|---|---|
| Fully Convolutional DFA | 75.4% | 75.7% | 52.1% | 52.7% | 61.2% | 61.6% | 60.8% |
| Clothes Detector + DFA | 76.3% | 76.1% | 56.3% | 57.6% | 61.7% | 61.6% | 63.1% |
| Joint RPN + DFA | 79.5% | 79.8% | 55.0% | 57.7% | 65.4% | 66.6% | 66.0% |
| DUN | **83.3%** | **83.7%** | **64.6%** | **66.7%** | **71.7%** | **72.4%** | **73.8%** |

### Per-category Analysis



T-shirt | Summer Wear | Coat | Tank-top

Zoom-in | Large Deviation | Irregular Pose | High Level of Clutter

### Ablation Study

| | | | | | | DUN |
|---|---|---|---|---|---|---|
| Fully Convolutional DFA? | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Spatial Transformer? | | ✓ | ✓ | ✓ | ✓ | |
| Selective Dilated Convolutions? | | | ✓ | ✓ | ✓ | |
| HR Spatial Transformer? | | | | ✓ | ✓ | |
| Scale Regularization? | | | | | ✓ | |
| detection rate (%) | 56.9 | 62.8 | 64.8 | 71.2 | **73.8** | |

### Selective Dilated Convolution



Dilation=[1] | Dilation=[1, 2] | Dilation=[1, 2, 3] | Dilation=[1, 2, 4, 8]

Maximum Selection | Average Selection | Small Scale | Medium Scale | Large Scale

(a) (b)

### Generalization of DLAN



(a) # Training Samples — Fully Convolutional DFA / Clothes Detector + DFA / Joint RPN + DFA / DLAN

(b) # Training Samples — Joint RPN + DFA(Shop) / Joint RPN + DFA(User) / DLAN(Shop) / DLAN(User)