



Microsoft



CoCosNet v2: Full-Resolution Correspondence Learning for Image Translation

Xingran Zhou¹, Bo Zhang², Ting Zhang², Pan Zhang⁴, Jianmin Bao²,
Dong Chen², Zhongfei Zhang³, Fang Wen²

1 Zhejiang University 2 Microsoft Research Asia 3 Binghamton University 4 USTC



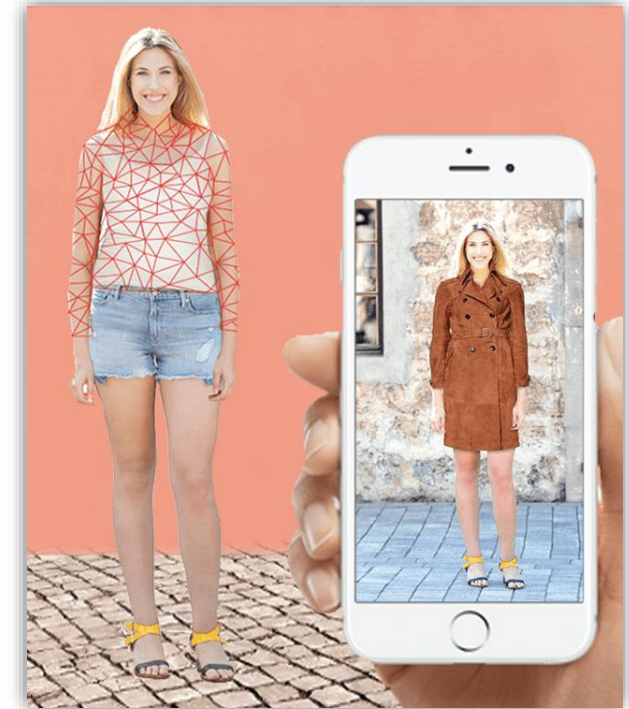
Image-to-image translation



old photos restoration



semantic editing



virtual try-on



Exemplar-based translation



exemplar



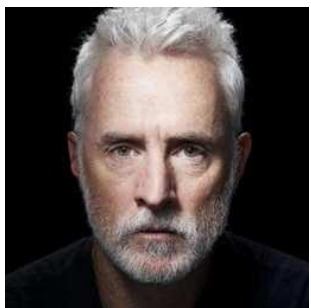
segmentation



synthesis

□ Pros

- Flexible user control
- Improved generation quality



exemplar



edge



synthesis



Exemplar-based translation



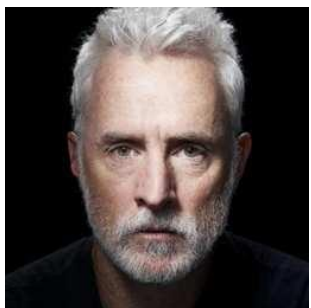
exemplar



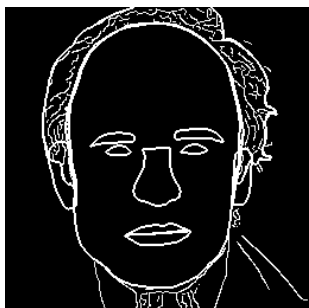
segmentation



synthesis



exemplar



edge



synthesis

□ Pros

- Flexible user control
- Improved generation quality

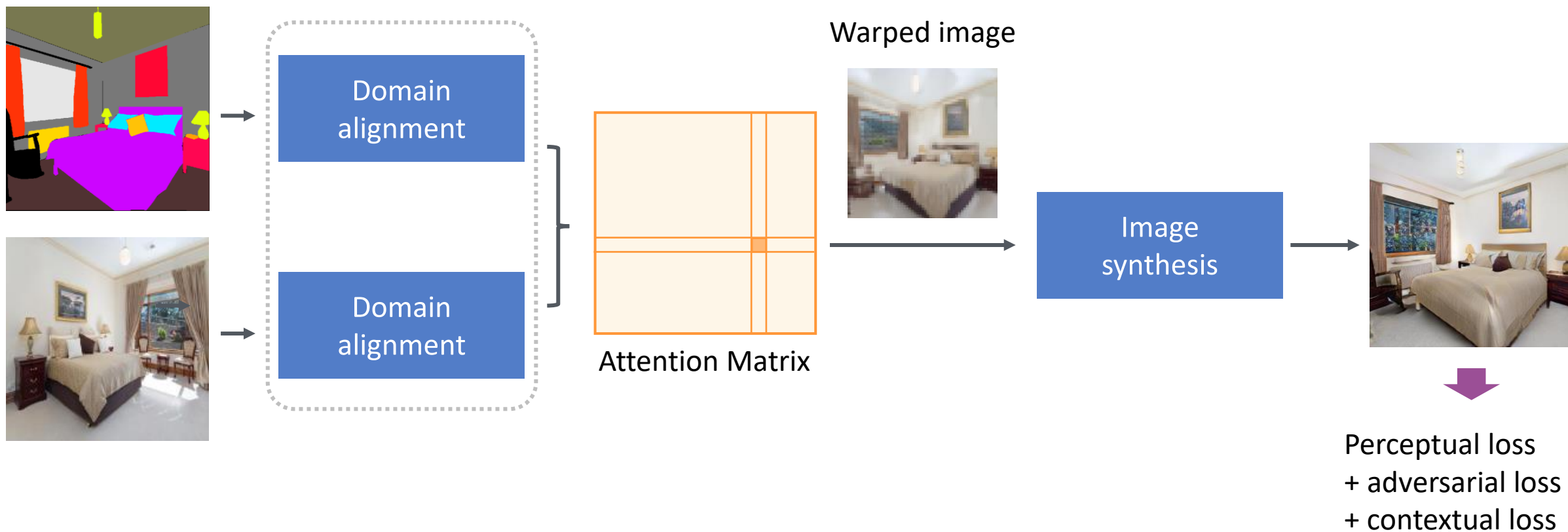
□ Cons

- Significant **artifacts** for complex scenes
- Lack of **fine-grained style** controllability
- Lack of using **fine textures** from exemplar



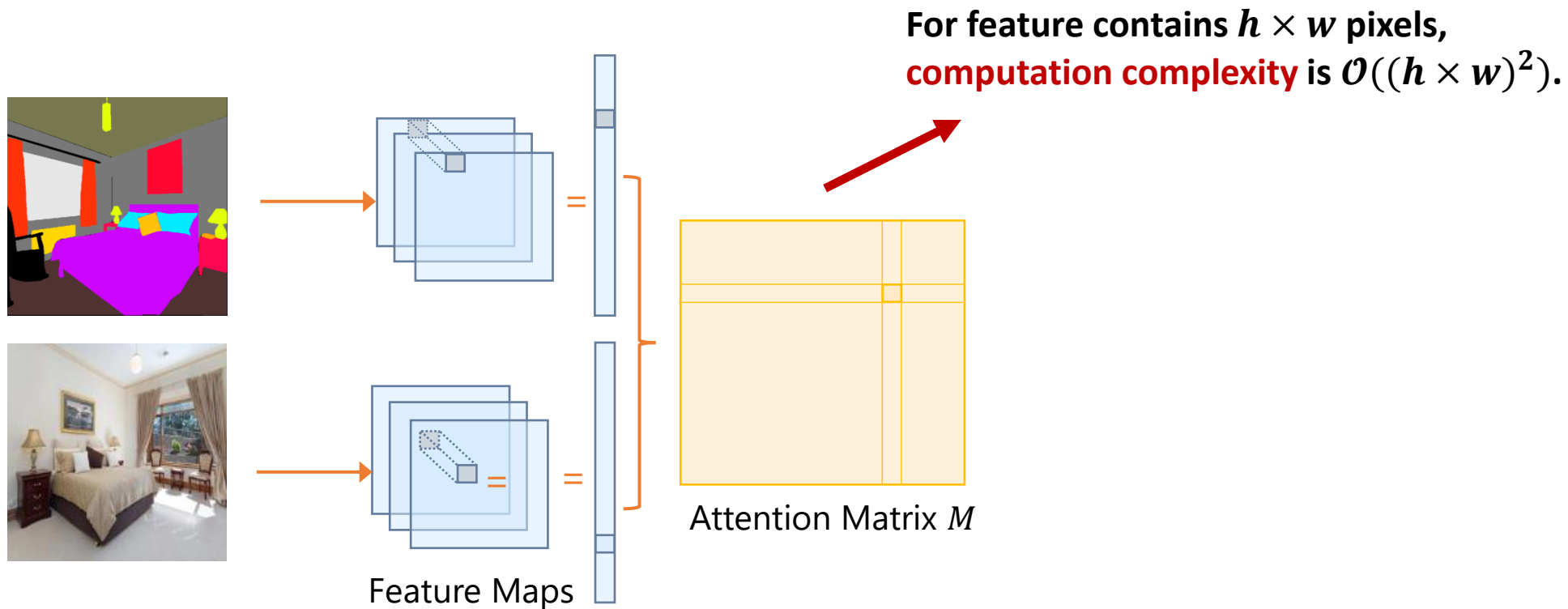
CoCosNet

Cross-domain correspondence learning with image synthesis



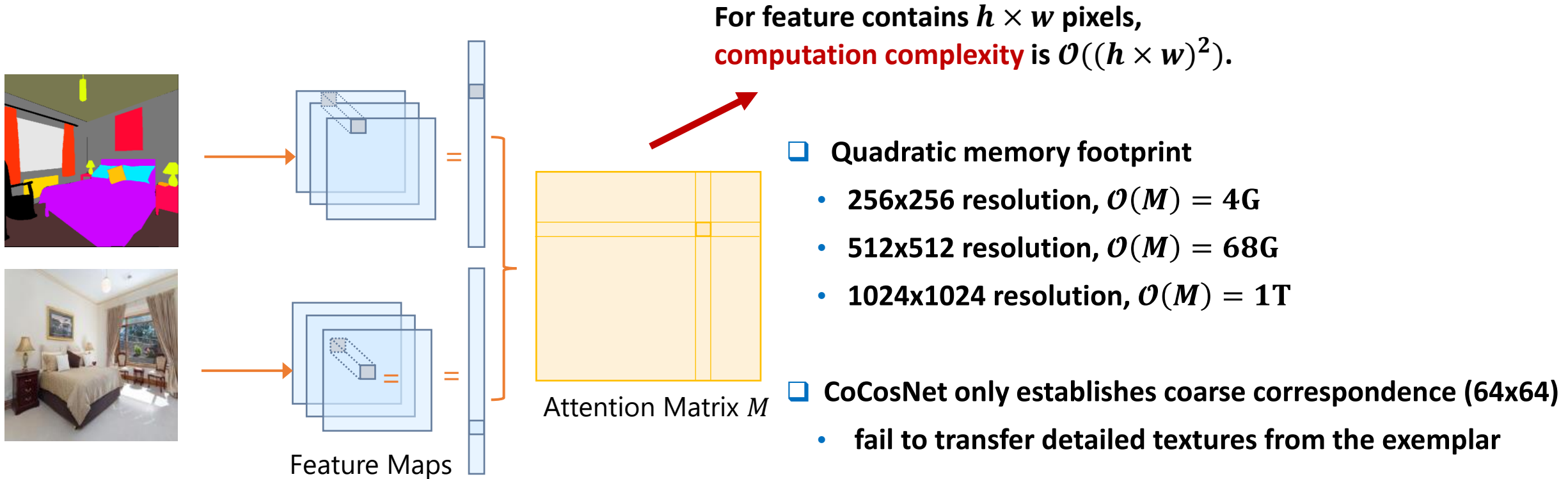
CoCosNet

Full resolution attention is computationally **inhibitive**



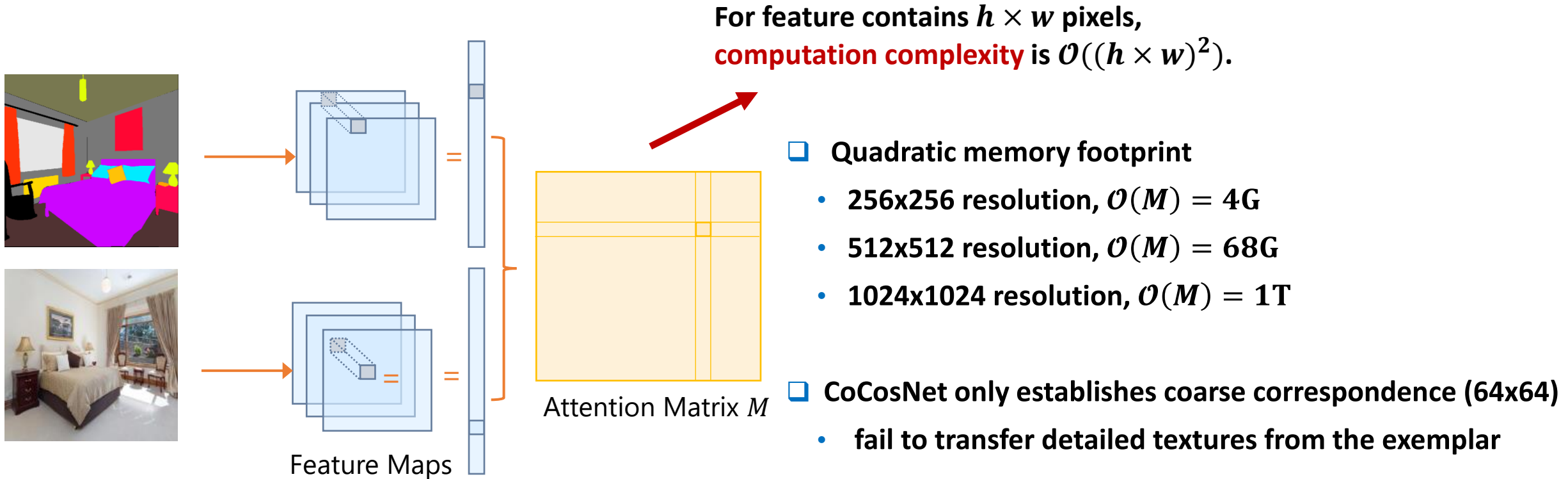
CoCosNet

Full resolution attention is computationally **inhibitive**



CoCosNet

Full resolution attention is computationally **inhibitive**



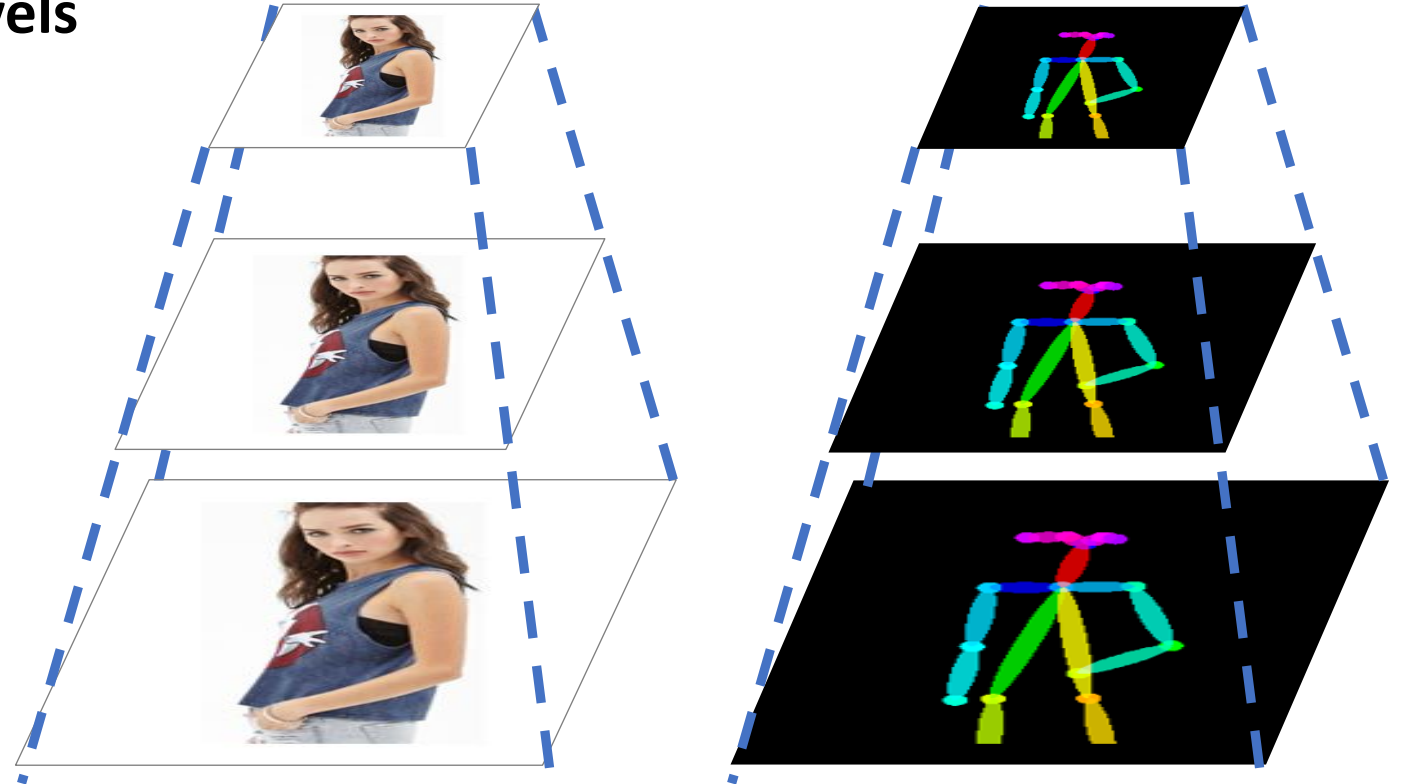
How can we compute the correspondence on high resolution?



CoCosNet v2

❑ Coarse-to-fine strategy

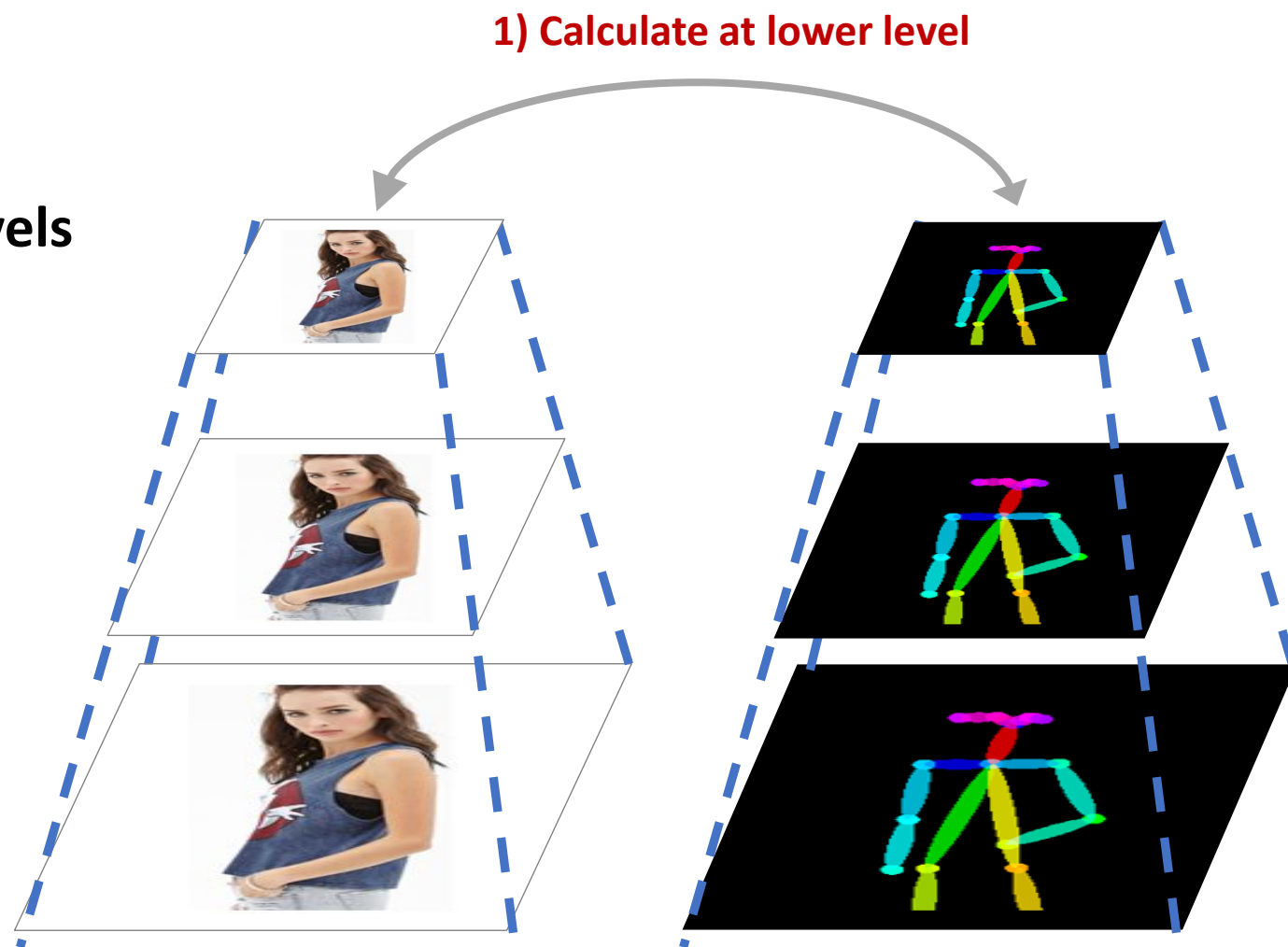
- Coarse level guides the finer levels



CoCosNet v2

❑ Coarse-to-fine strategy

- Coarse level guides the finer levels



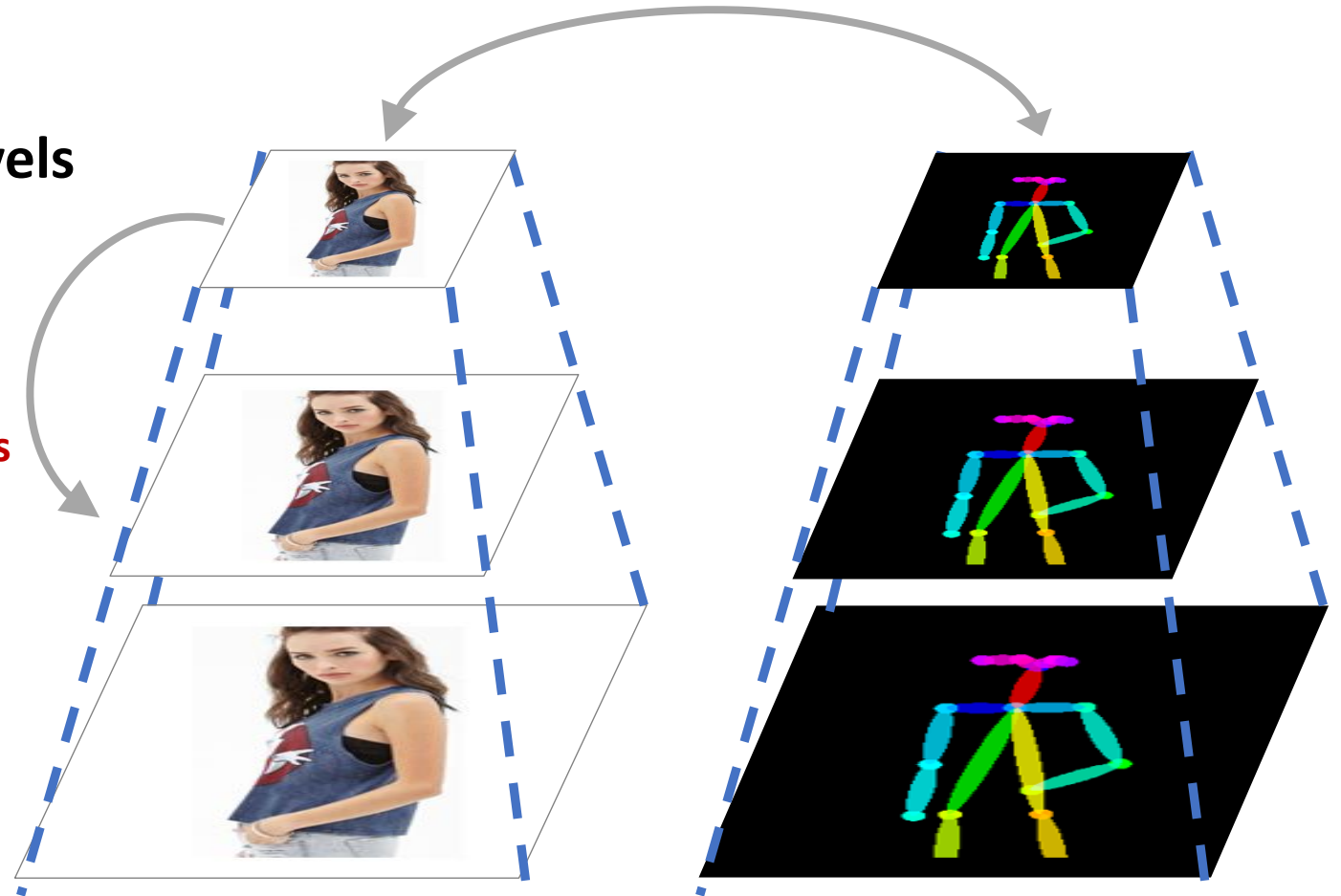
CoCosNet v2

❑ Coarse-to-fine strategy

- Coarse level guides the finer levels

2) Initialization for finer levels

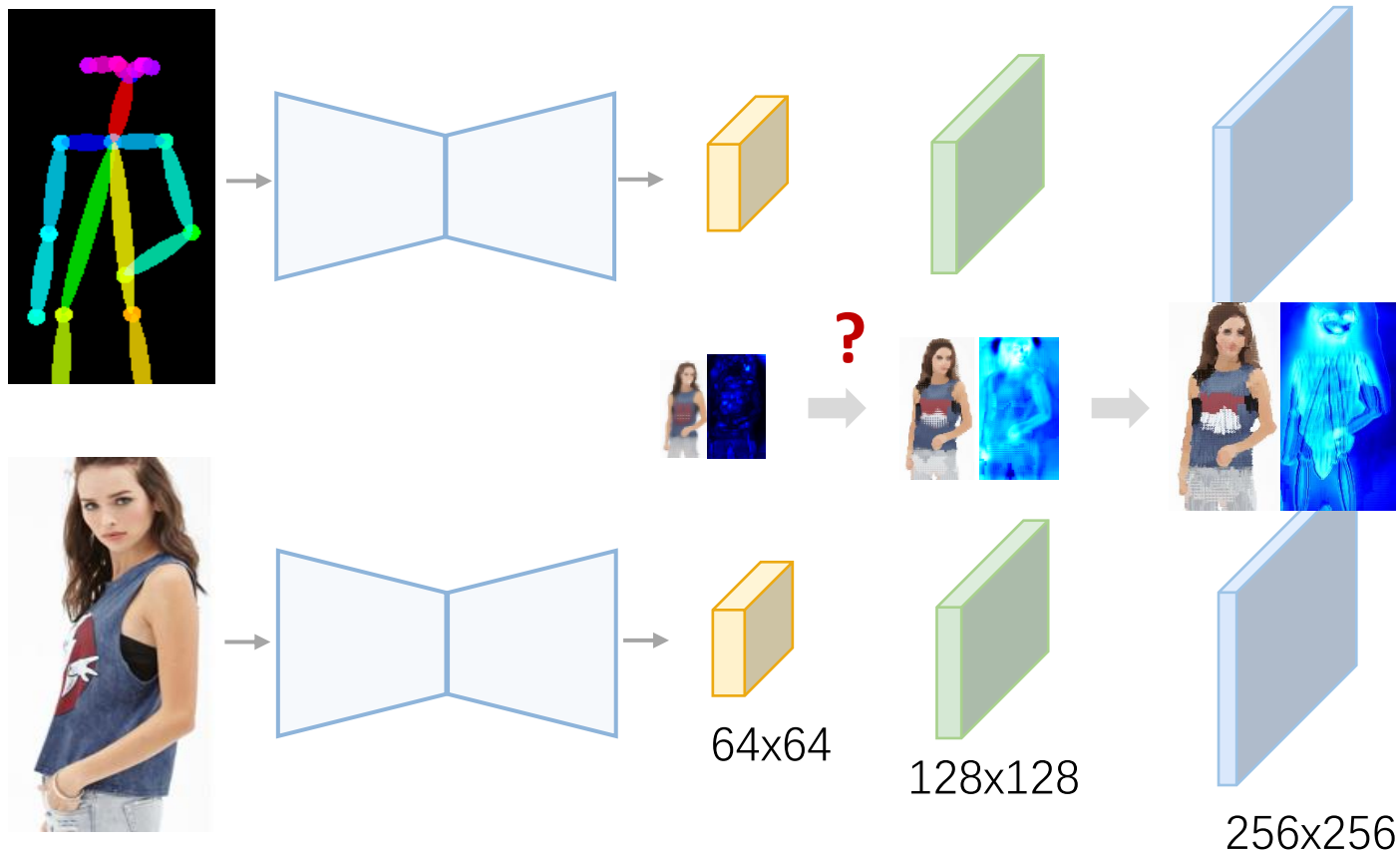
1) Calculate at lower level



CoCosNet v2

Coarse-to-fine strategy

- Coarse level guides the finer levels



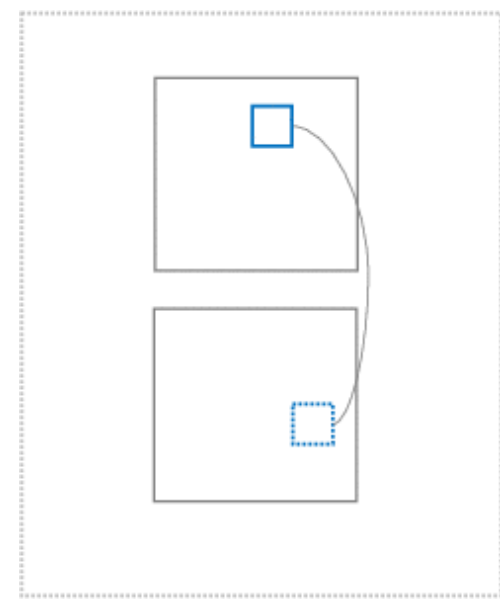
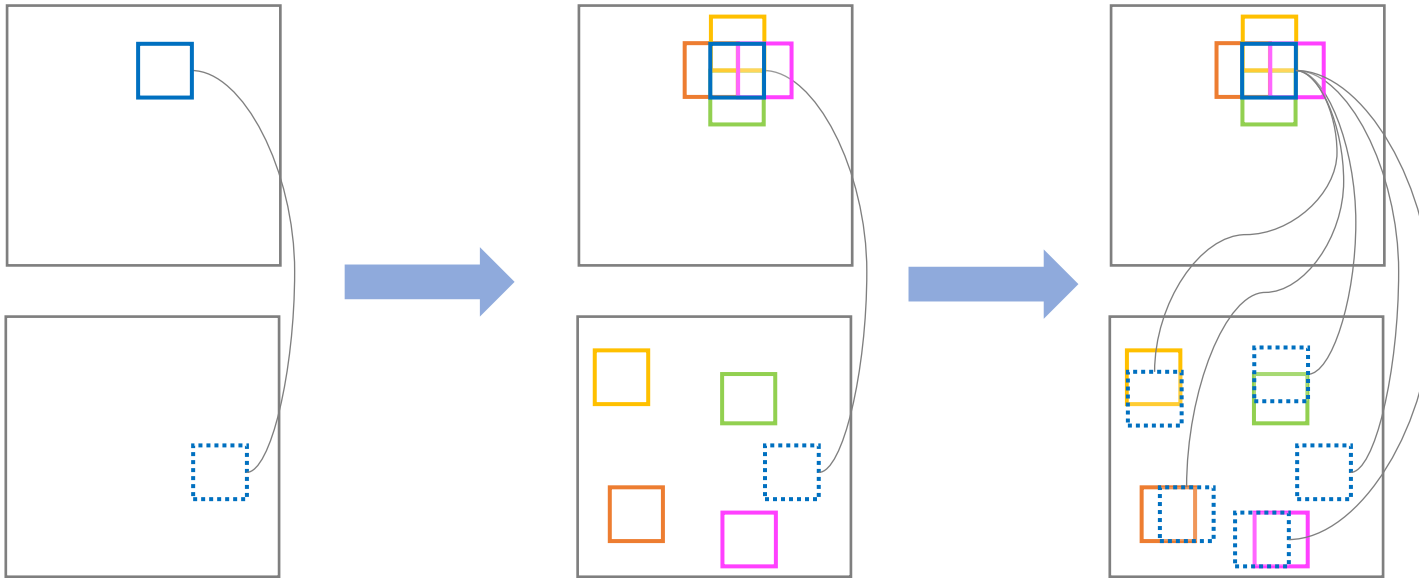
How can we make use of the initialization from the lower level?



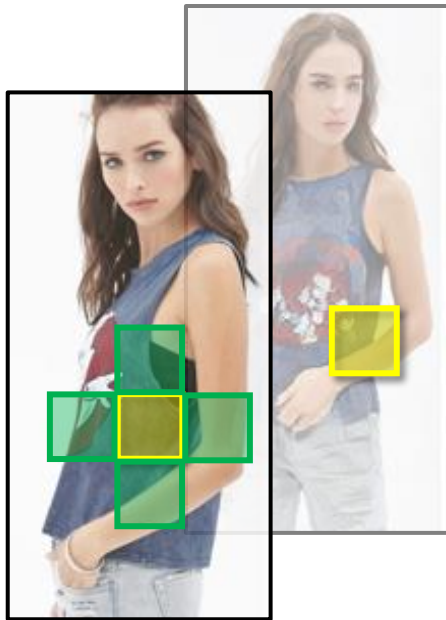
PatchMatch

Global ➔ Local : PatchMatch

PatchMatch searches from the **neighborhood** rather than searching **globally**.



Differentiable PatchMatch

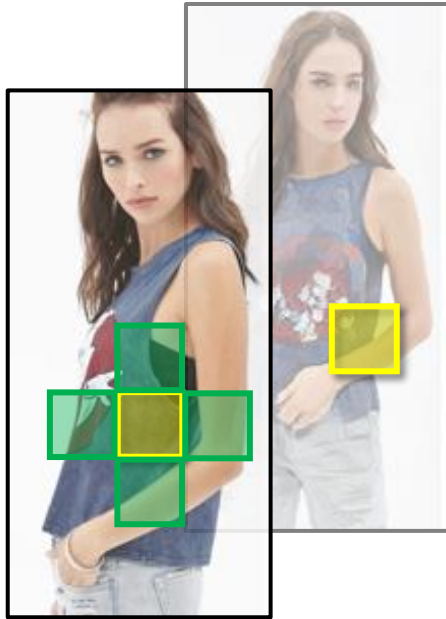


- Our supervision comes from the image warping, *i. e.*,

$$y_B(H(p)) \quad \text{where, } H(p) = \arg \min_q ||f_x(p) - f_y(q)||$$



Differentiable PatchMatch



- Our supervision comes from the image warping, *i. e.*,

$$y_B(H(p)) \quad \text{where, } H(p) = \arg \min_q ||f_x(p) - f_y(q)||$$

- To make it differentiable, we consider the K possible matchings. Now the warping becomes:

$$w^{y \rightarrow x}(p) = \sum_{k=1}^K \text{softmax}\{\underbrace{S(p; k)}_{\text{confidence}} \underbrace{y_B(H(p; k))}_{k^{th} \text{ matching}}\}$$

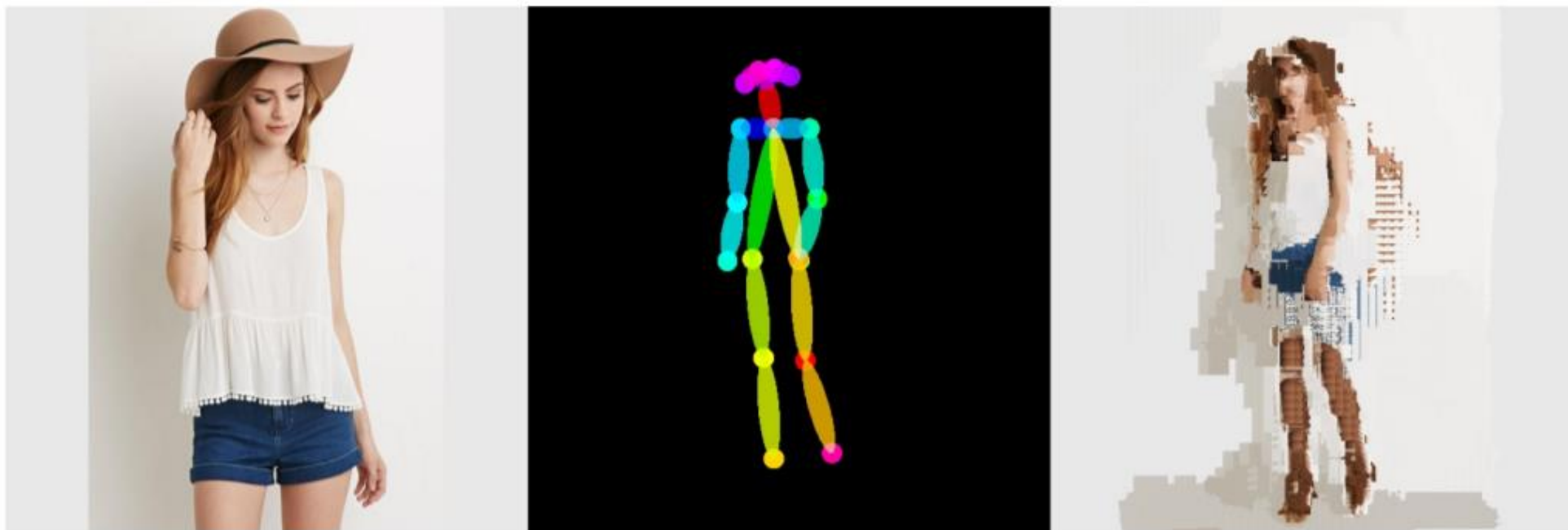
confidence

k^{th} matching



Differentiable PatchMatch

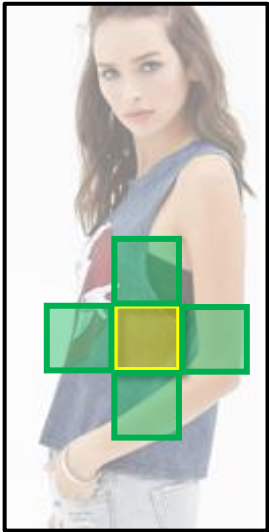
- ❑ PatchMatch implicitly assume **local smoothness**, which is true to natural images
- ❑ However, this is violated because deep features are not well-trained at the beginning



warped image
(**only** PatchMatch propagation)



ConvGRU-assisted PatchMatch



- ❑ PatchMatch only considers the **adjacent** patches
- ❑ Conv makes the propagation consider the **distant patches**
- ❑ The gradient can be propagated to more locations

Receptive field of ConvGRU-assist propagation



CoCosNet v2

Warped images via different **variants** of our method.

CoCosNet v2 produces the **most faithful** warped image.



only PatchMatch propagation



only ConvGRU



PatchMatch propagation with Conv



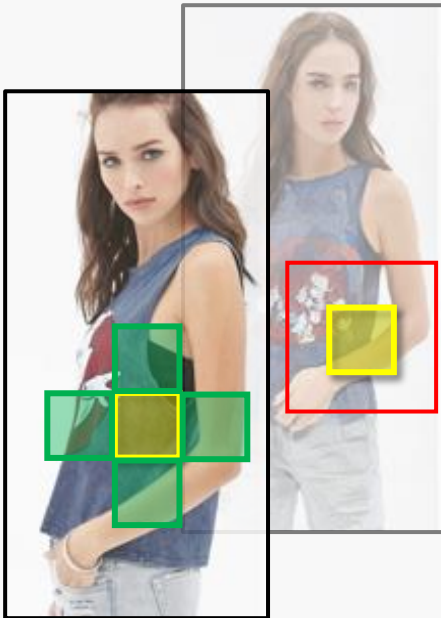
CoCosNet v2



CoCosNet v2

ConvGRU-assisted Patch Match

Propagation from neighborhood

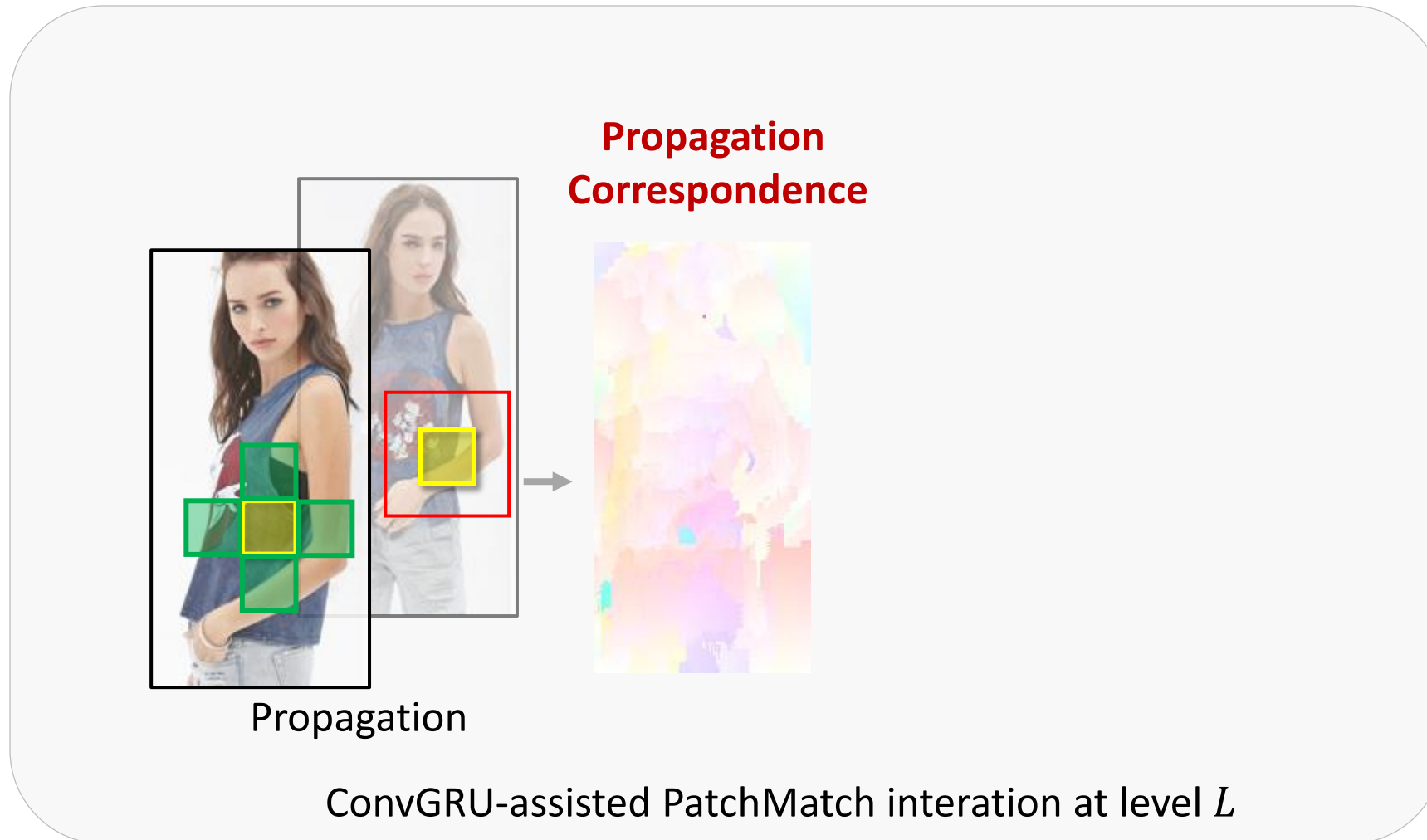


ConvGRU-assisted PatchMatch interaction at level L



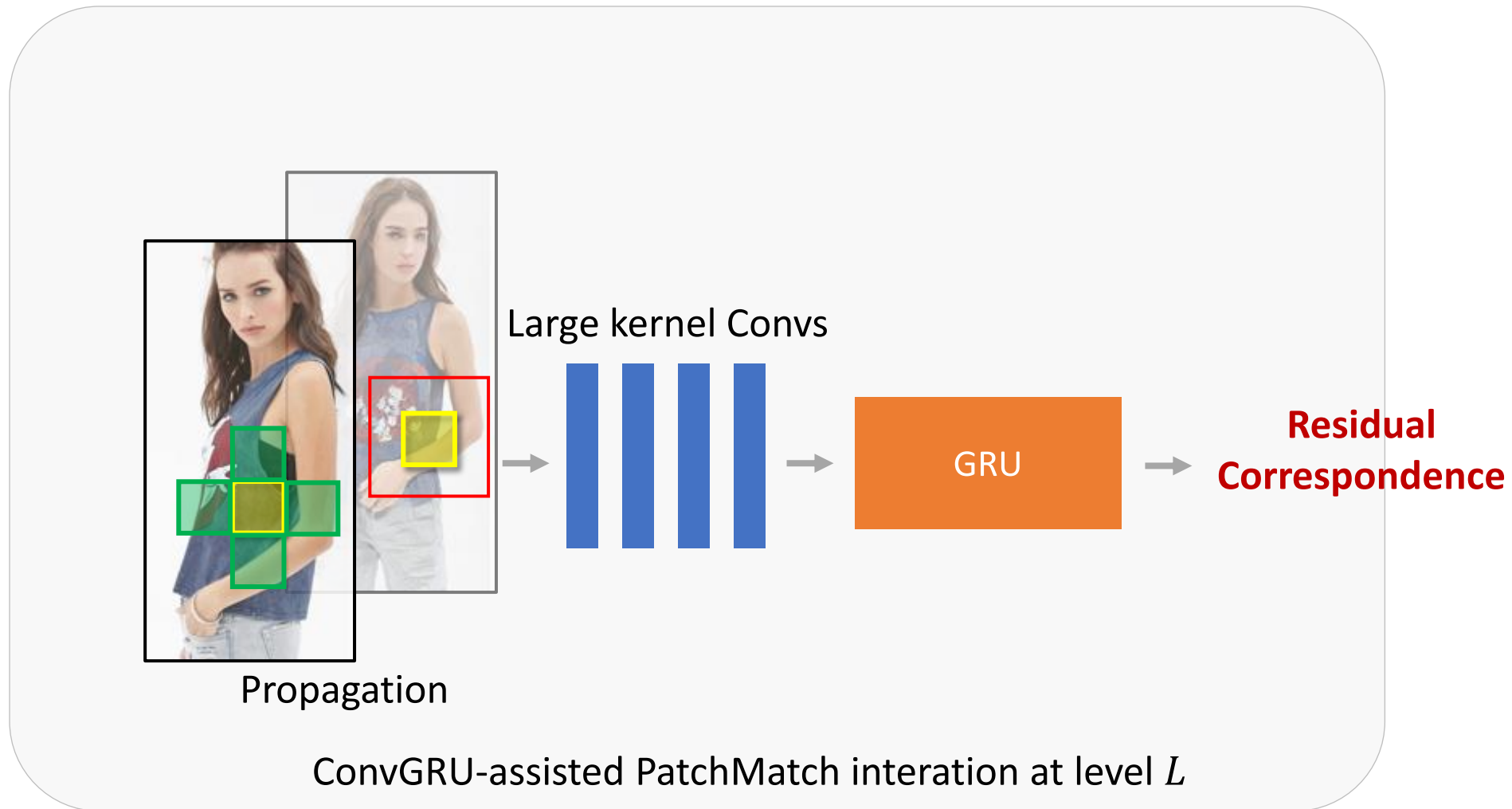
CoCosNet v2

ConvGRU-assisted Patch Match



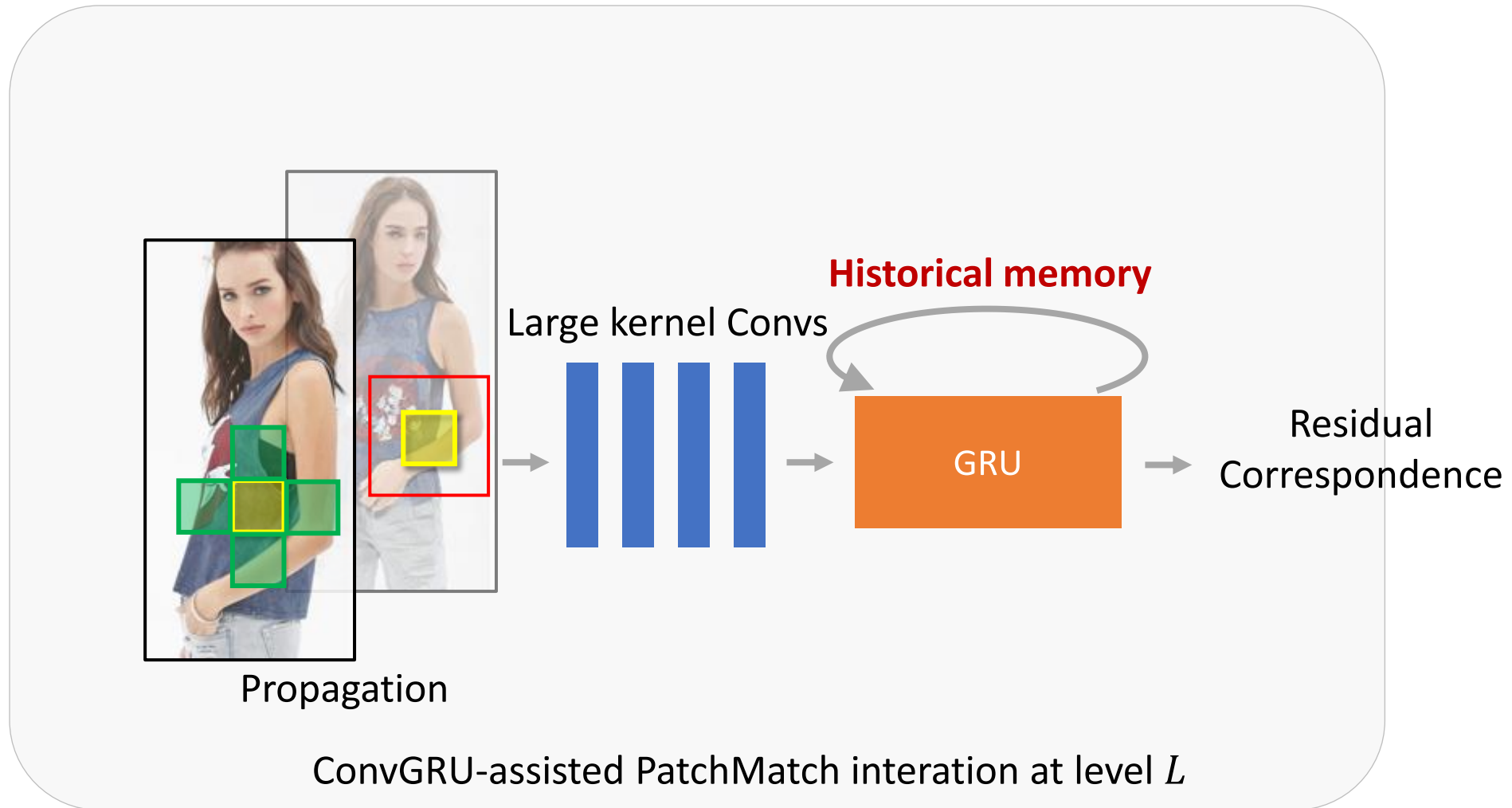
CoCosNet v2

ConvGRU-assisted Patch Match



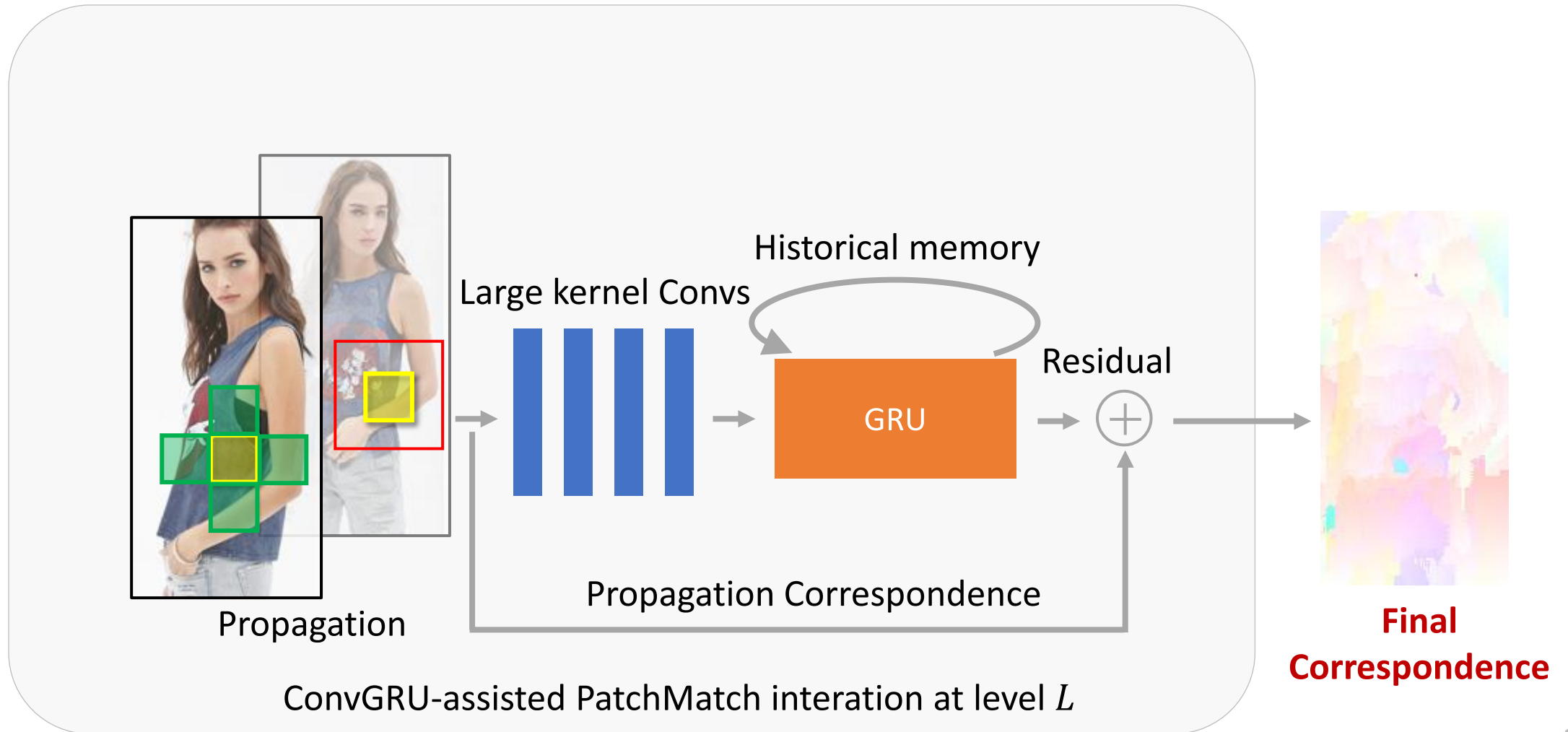
CoCosNet v2

ConvGRU-assisted Patch Match



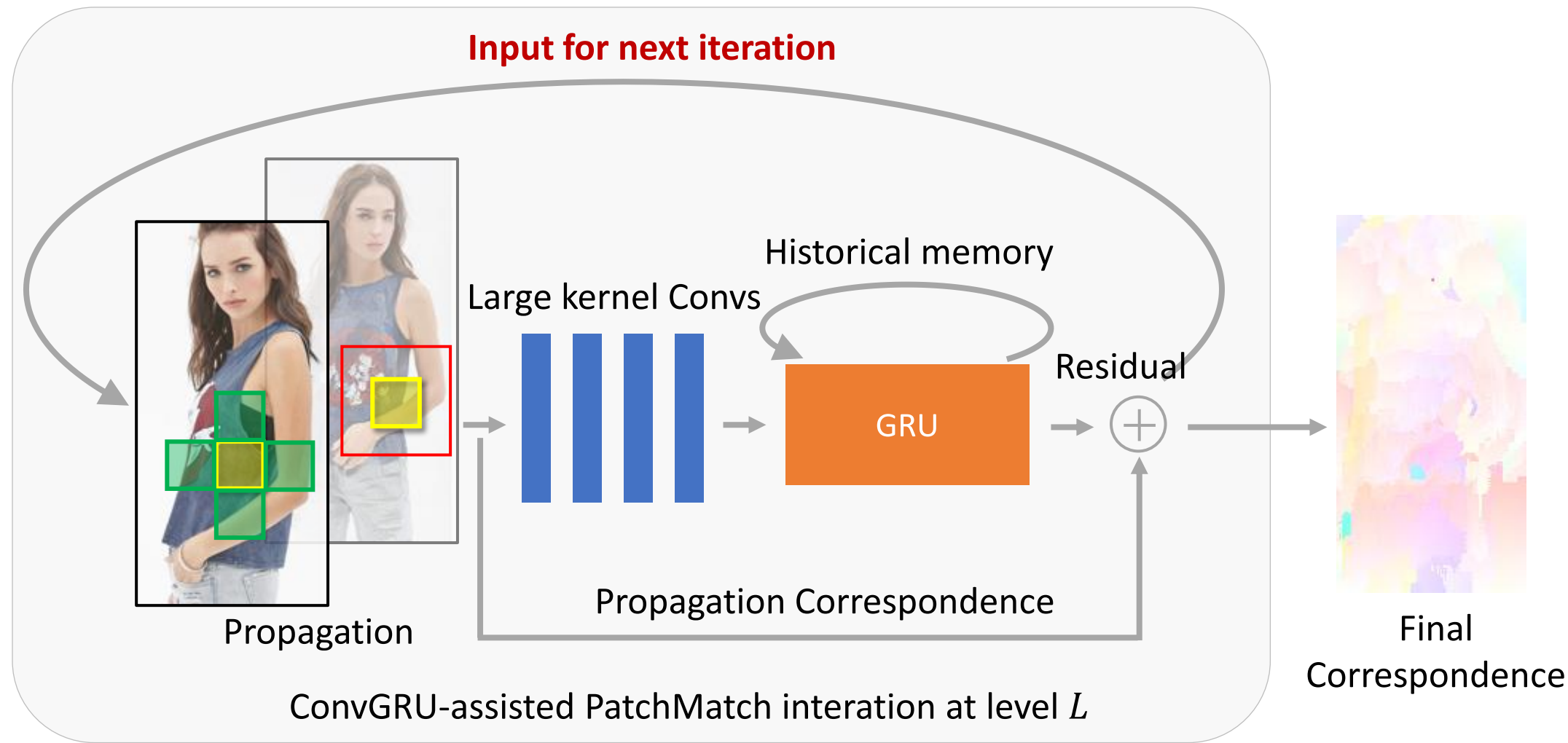
CoCosNet v2

ConvGRU-assisted Patch Match



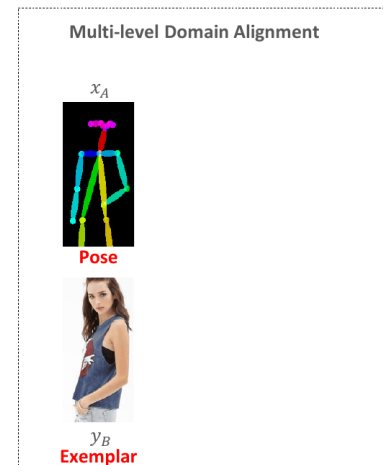
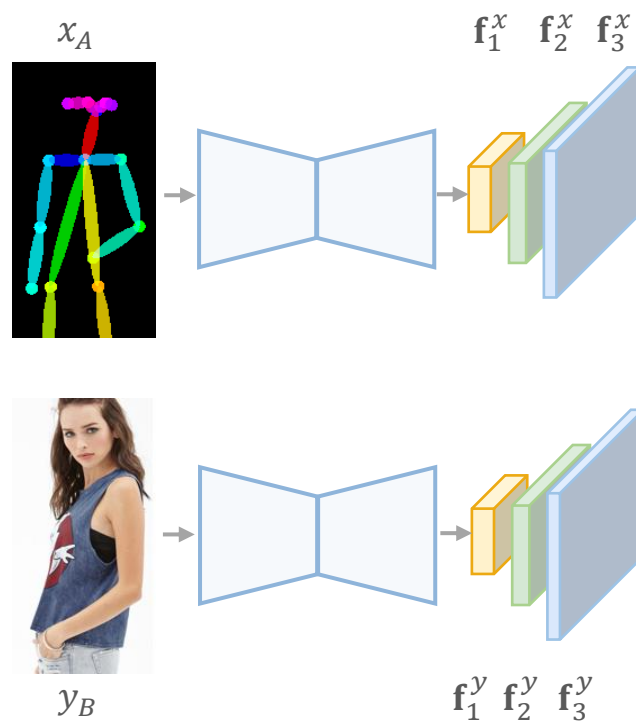
CoCosNet v2

ConvGRU-assisted Patch Match



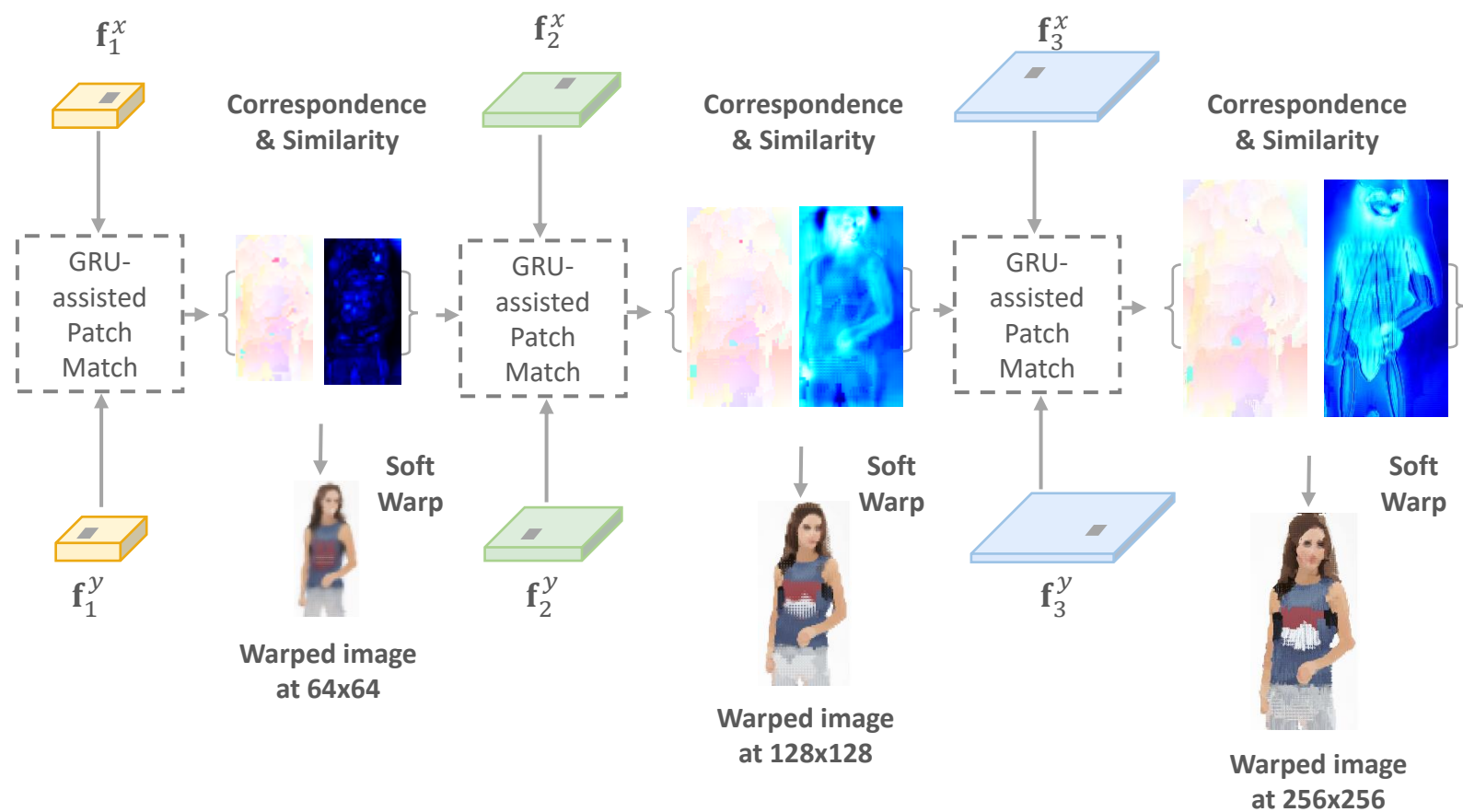
CoCosNet v2

Multi-level Domain Alignment Feature Extraction



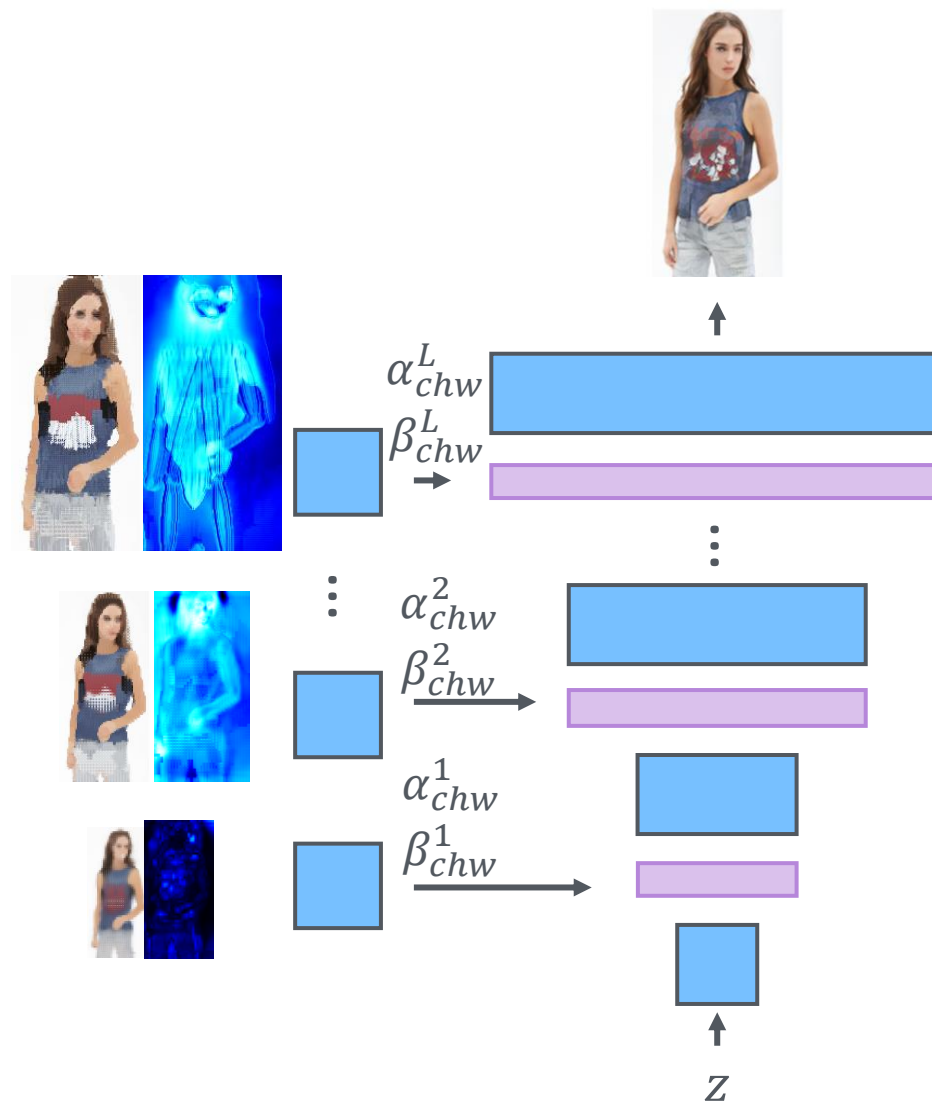
CoCosNet v2

Hierarchical ConvGRU-assisted Patch Match

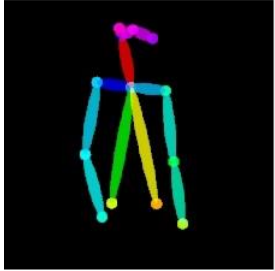


CoCosNet v2

Translation Network



Pose-to-body



Pose



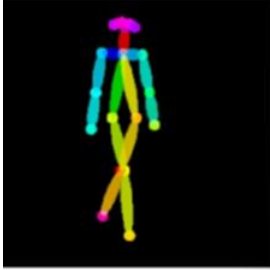
Exemplar

Synthesis

Deepfashion dataset (512x512 resolution)



Pose-to-body



Pose

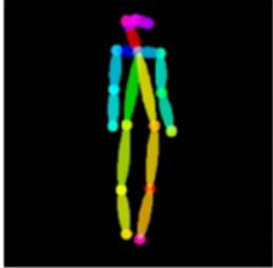


Exemplar

Synthesis



Pose-to-body



Pose



Exemplar

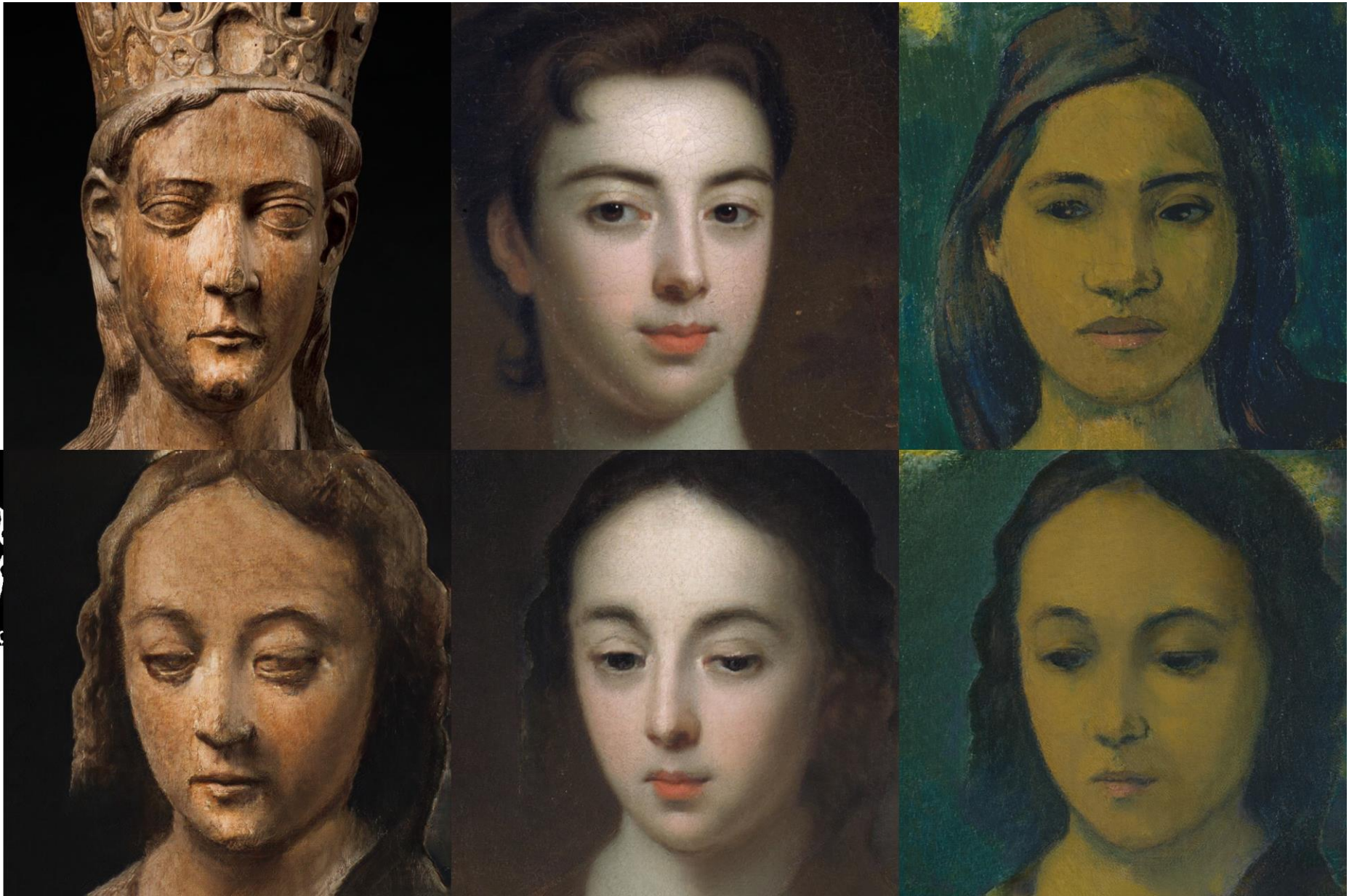
Synthesis



Edge-to-face



Edge



Exemplar

Synthesis

MetFaces dataset (1024x1024 resolution)



Edge-to-face

Exemplar

Edge

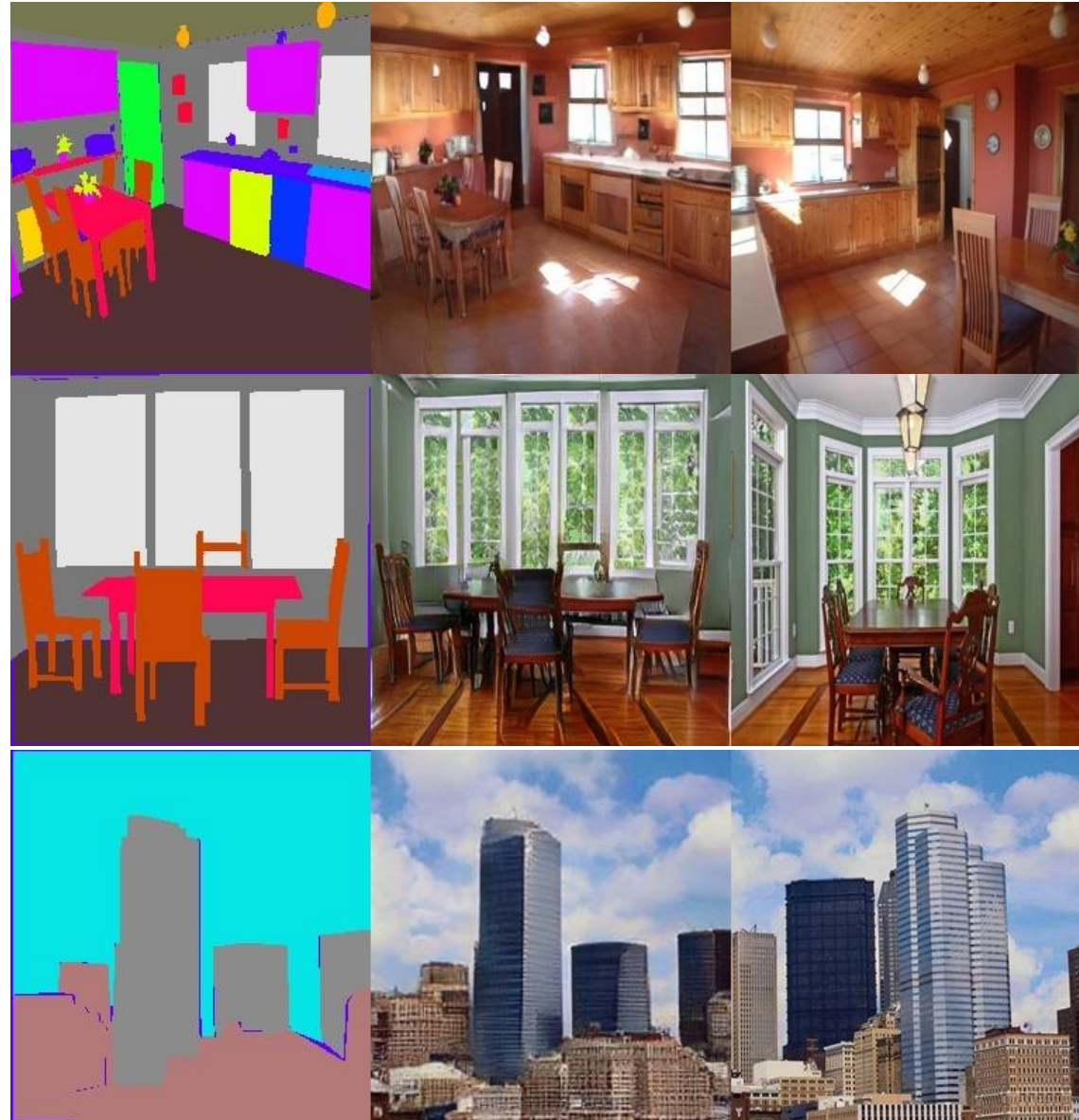


Synthesis

MetFaces dataset (1024x1024 resolution)



Mask-to-image



Segmentation

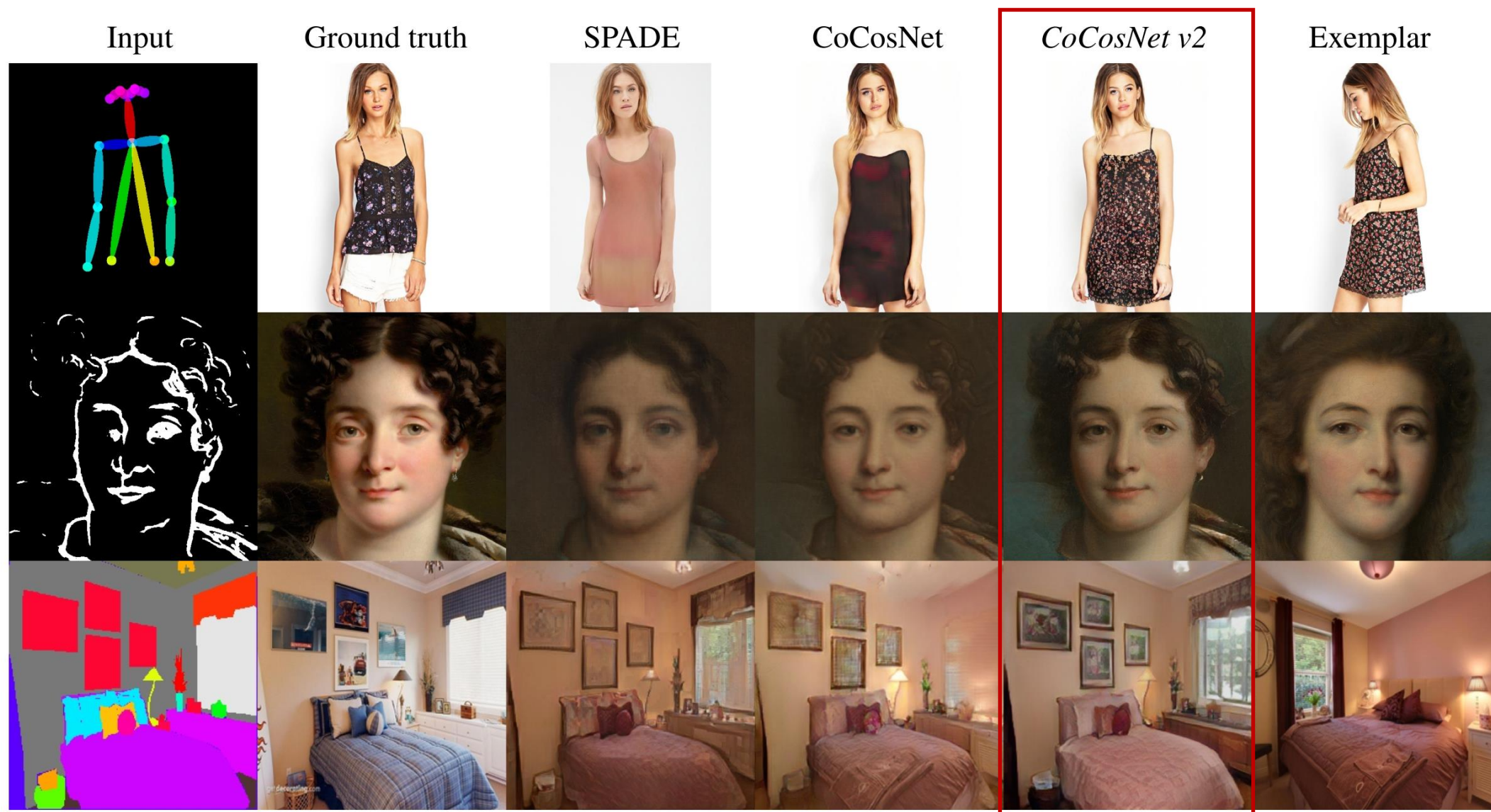
Exemplar

Synthesis

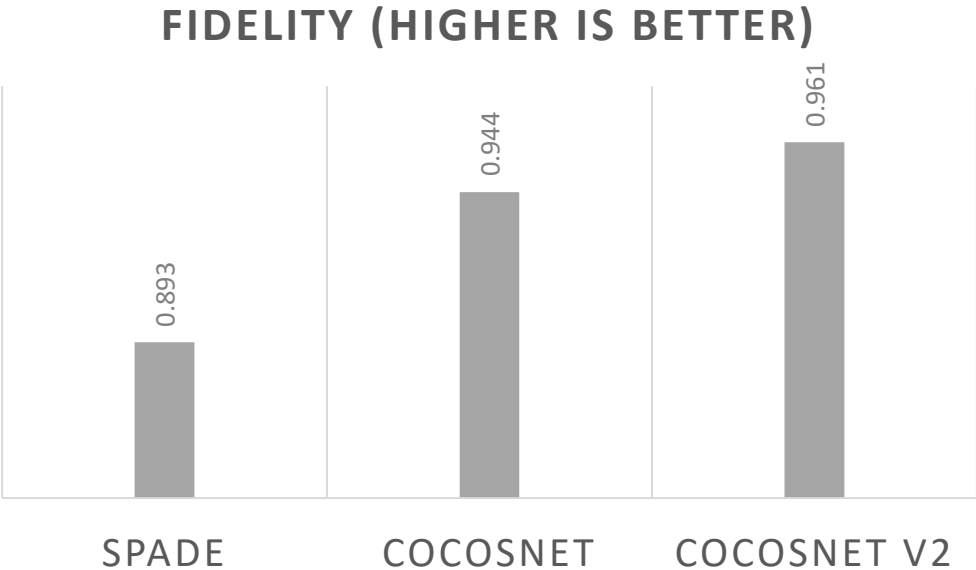
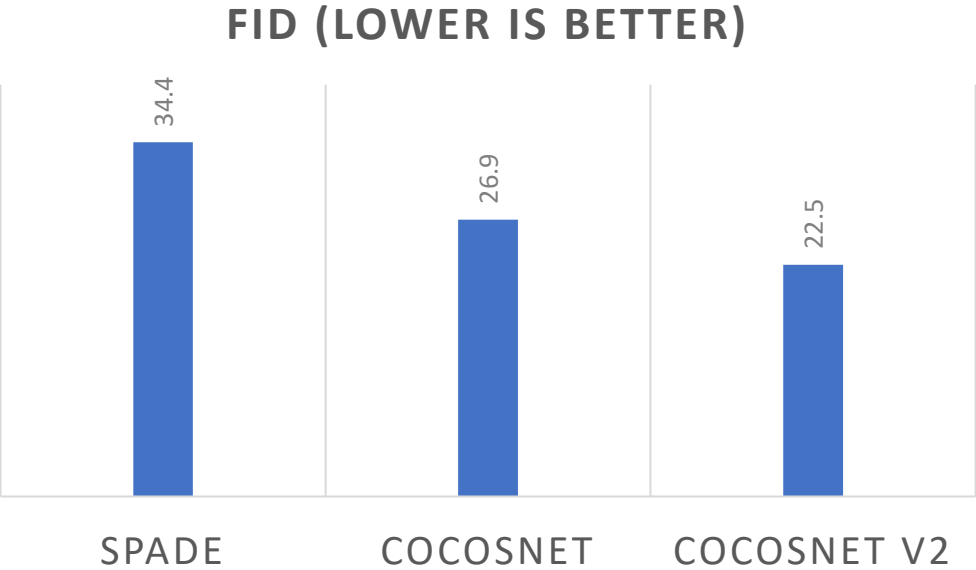
ADE20k dataset



Quantitative comparison



Quantitative comparison



Application

Exemplar



Photos of real person



Synthesis

Synthesis

Oil portrait application



Application

Exemplar



Photos of real person



Synthesis

Oil portrait application



Thank you !

