

# Learning Co-segmentation by Segment Swapping for Retrieval and Discovery

Xi Shen<sup>1</sup> Alexei A. Efros<sup>2</sup> Armand Joulin<sup>3</sup> Mathieu Aubry<sup>1</sup>

<sup>1</sup>LIGM (UMR 8049), École des Ponts ParisTech <sup>2</sup>UC Berkeley <sup>3</sup>Meta AI

## 1. Introduction

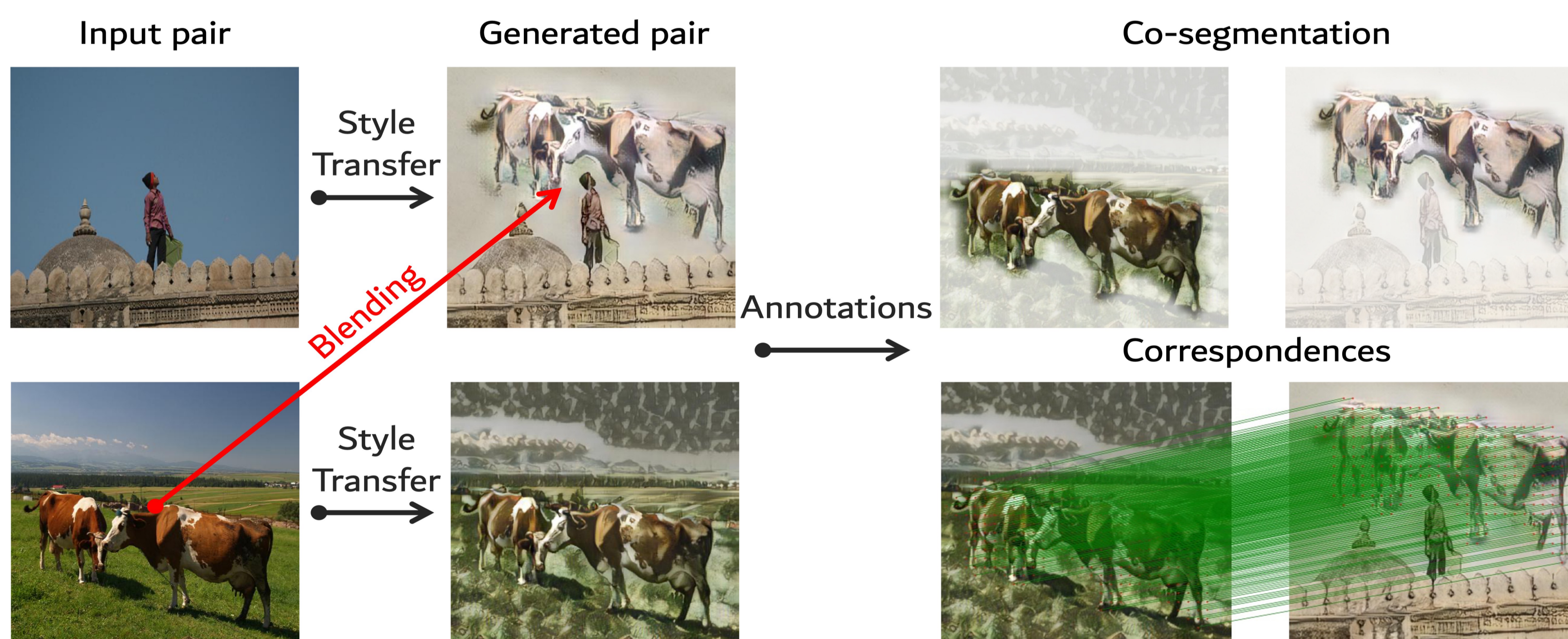
### Motivation



### Challenges

- No training data available.

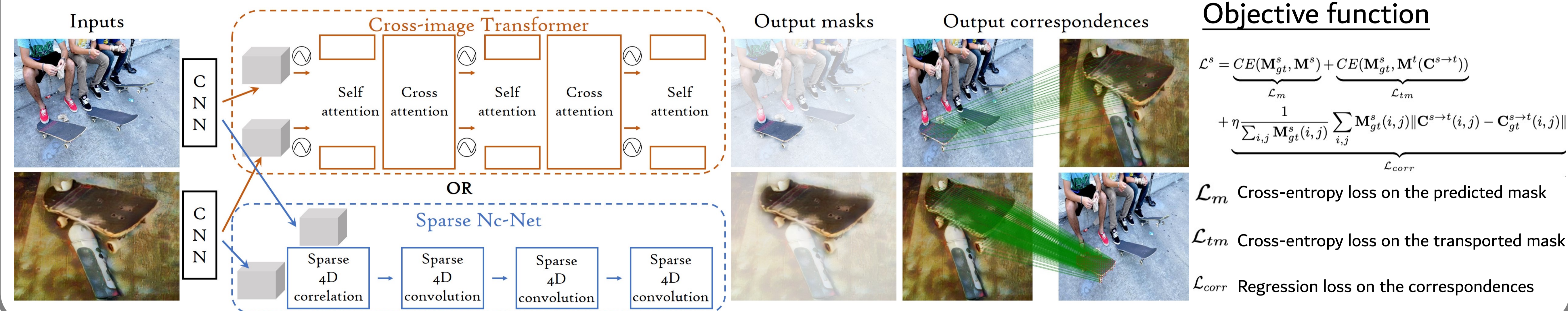
### Key idea : training with synthetic pairs



### Main contributions

- A method of generating synthetic pairs to learn co-segmentation
- A transformer-based architecture for co-segmentation producing competitive performances on art detail retrieval and place recognition
- Spectral clustering on a correspondence graph for discovery in image collections.

## 2. Learning co-segmentation



## 3. Experiments

Query → Top-3 retrieved images

Feat. + Methods	Retrieval	mAP Det.(IoU > 0.3)
Shen et al. [54] + cos [54]	75.5	75.3
Shen et al. [54] + discovery [54]	76.6	76.4
MocoV2 [9] + cos [54]	79.0	78.7
MocoV2 [9] + discovery [54]	80.8	79.6
<b>Ours + Unsupervised segments</b>		
Transformer	81.8	79.4
Sparse-Ncnet	82.8	73.4
<b>Ours + COCO segments [36]</b>		
Transformer	<b>84.4</b>	<b>81.8</b>
Sparse-Ncnet	83.3	73.7

Score between a pair of images

$$\mathcal{S}(I^s, I^t) = \sum_{i,j} \underbrace{M_{joint}^s(i,j)}_{Mask} \underbrace{\cos(\mathbf{F}^s(i,j), \mathbf{F}^t(C^{s \rightarrow t}(i,j)))}_{Feat. similarity}$$

Where  $M_{joint}^s$  is the product of the source and the transported target mask:

$$M_{joint}^s(i,j) = M^s(C^{s \rightarrow t}(i,j))M^t(i,j)$$

$\mathbf{F}^s(i,j)$  Feature at coordinate  $(i,j)$  in the source  
 $\mathbf{F}^t(C^{s \rightarrow t}(i,j))$  Feature at warped coordinate  $(i,j)$  in the target

Method	ECCSD [57]			DUTS [74]			DUT-OMRON [77]		
	max $F_{\beta}$	IoU	Acc.	max $F_{\beta}$	IoU	Acc.	max $F_{\beta}$	IoU	Acc.
HS [76]	0.673	0.598	0.847	0.504	0.369	0.826	0.561	0.433	0.843
wCsr [50]	0.684	0.517	0.862	0.522	0.392	0.835	0.541	0.416	0.838
WSC [36]	0.683	0.498	0.852	0.528	0.384	0.862	0.523	0.387	0.865
DeepUSPS [43]	0.584	0.440	0.795	0.425	0.305	0.773	0.414	0.305	0.779
BigBiGAN [73]	0.782	0.672	0.899	0.608	0.498	0.878	0.549	0.453	0.856
E-BigBiGAN [73]	0.797	0.684	0.906	0.624	0.511	0.882	0.563	0.464	0.860
LOST [58]	0.758	0.654	0.895	0.611	0.518	0.871	0.473	0.410	0.797
LOST [58] + Bilateral Solver [5] (Ours)	<b>0.837</b>	<b>0.723</b>	<b>0.916</b>	<b>0.697</b>	<b>0.572</b>	<b>0.887</b>	<b>0.578</b>	<b>0.489</b>	<b>0.818</b>

Validation of unsupervised segments on unsupervised saliency detection benchmarks : ECCSD, DUTS, DUT-OMRON.

## 4. Object discovery

Discovered clusters on Bruegel dataset.

Horse, Airplane, Car

Co-segmentation results in Internet dataset for the Horse, Airplane and Car categories.

**Correspondences graph**

Nodes : correspondences

$$e_{i,j} = \frac{1}{2} m_i m_j \exp\left(-\frac{\|x_i^s - x_j^t\|}{\sigma}\right) \left( \exp\left(-\frac{\|x_i^s - C^{t \rightarrow s}(x_j^t)\|}{\sigma}\right) + \exp\left(-\frac{\|x_j^t - C^{s \rightarrow t}(x_i^s)\|}{\sigma}\right) \right)$$

Method	Airplane		Car		Horse		Avg	
	$\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	$\mathcal{J}$
DOCS [17]	0.946	0.664	0.940	0.83	0.914	0.65	0.933	0.70
Sun et al. [59]	0.886	0.36	0.870	0.73	0.876	0.55	0.877	0.55
Joulin et al. [28]	0.493	0.15	0.587	0.37	0.638	0.30	0.572	0.27
Kim et al. [10]	0.802	0.08	0.889	0.00	0.751	0.06	0.754	0.05
Rubinstein et al. [51]	0.880	0.56	0.854	0.64	0.828	0.52	0.827	0.43
Chen et al. [1]	0.902	0.40	0.876	0.65	0.893	0.58	0.896	0.54
Quan et al. [45]	0.910	0.56	0.885	0.67	0.893	0.58	0.896	0.60
Hsu et al. [19]	0.777	0.33	0.621	0.43	0.738	0.20	0.712	0.32
Chang et al. [7]	0.728	0.27	0.759	0.36	0.797	0.26	0.761	0.33
Lee et al. [3]	0.528	0.36	0.647	0.42	0.701	0.30	0.625	0.39
Jerripothula et al. [26]	0.905	0.61	0.880	0.71	0.883	0.61	0.889	0.64
Jerripothula et al. [25]	0.816	0.48	0.847	0.69	0.813	0.50	0.826	0.56
Hsu et al. [22]	0.956	0.66	0.914	0.79	0.876	0.59	0.909	0.68
Chen et al. [10]	<b>0.941</b>	<b>0.65</b>	<b>0.940</b>	<b>0.82</b>	<b>0.922</b>	<b>0.63</b>	<b>0.935</b>	<b>0.70</b>
Ours + Unsupervised segments								
transformer	<b>0.941</b>	<b>0.66</b>	<b>0.919</b>	<b>0.79</b>	<b>0.887</b>	<b>0.57</b>	<b>0.916</b>	<b>0.67</b>
Nc-Net	0.682	0.19	0.791	0.36	0.774	0.27	0.749	0.34
Ours + COCO segments [36]								
transformer	<b>0.941</b>	<b>0.67</b>	<b>0.928</b>	<b>0.82</b>	<b>0.916</b>	<b>0.66</b>	<b>0.928</b>	<b>0.70</b>
Nc-Net	0.655	0.23	0.857	0.61	0.873	0.43	0.795	0.42

\* uses learnt keypoint detector Superpoint [14]

## 5. Training pairs



Project page (code, more experimental and visual results): <http://imagine.enpc.fr/~shenx/SegSwap/>