



# Joint Inference in Image Databases via Dense Correspondence

Michael Rubinstein

MIT CSAIL

(while interning at Microsoft Research)

# My work

- Throughout the year (and my PhD thesis): **Temporal Video Analysis and Visualization**



**Pulse signal amplified**



**Breathing motions amplified**

- This short talk: my work during the summers (MSR 2011, 2012)
  - Inference in large, weakly-annotated image databases

# Videos vs. Image Datasets

- Goal: we want to infer properties of pixels/regions
  - Semantics, layers, geometry (depth), motion, ...
- Recent advances allow us to **treat a set of images like videos!**
  - Correspondence between adjacent frames in videos: optical flow, layer models, tracking, ...
  - Correspondence between similar images in databases: Feature Matching, graph matching, Spatial Pyramid Matching (SPM), SIFT flow, ...



# Image Correspondence is Challenging...

Query

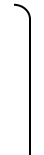
Best match



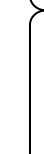
Multiple objects; no global transform



Changing perspective, occlusions



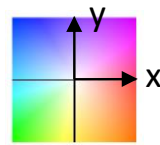
Intra-class variation



Background clutter



# ...but Good Solutions Exist



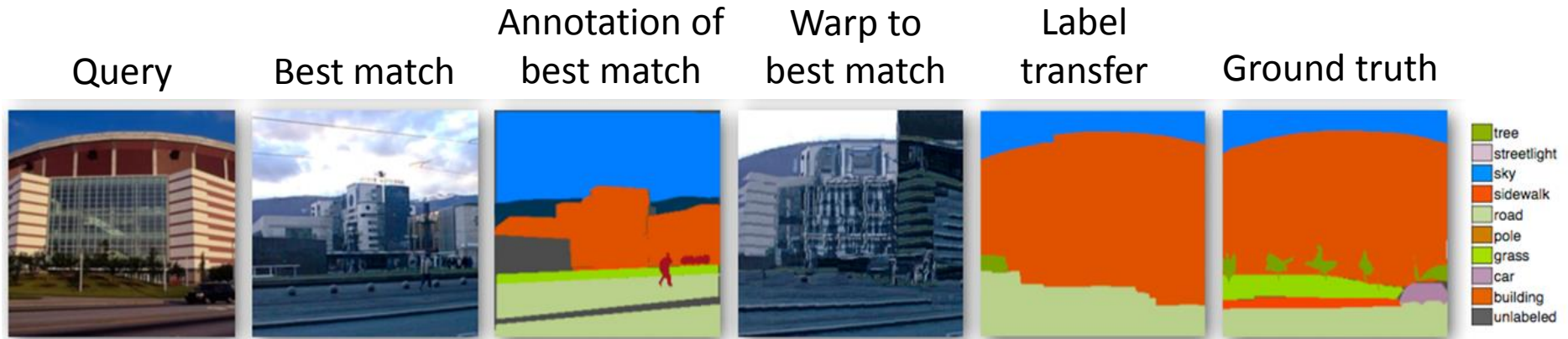
SIFT Flow [Liu et al. TPAMI 2011]

Query

Best match

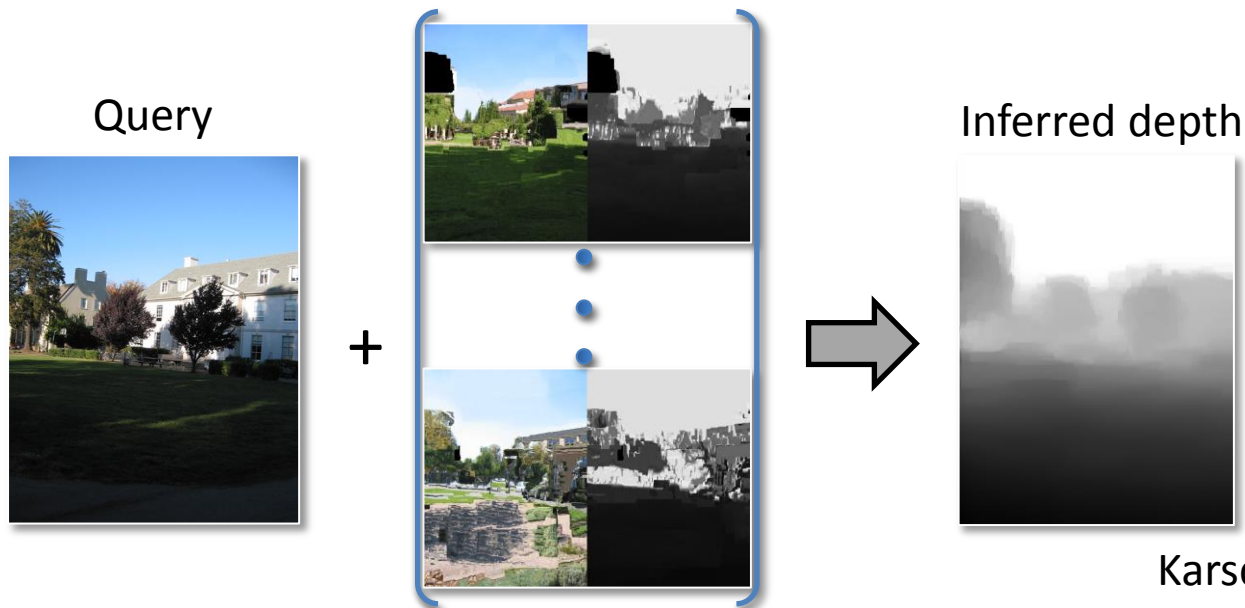


# Correspondence-driven Approaches to Computer Vision



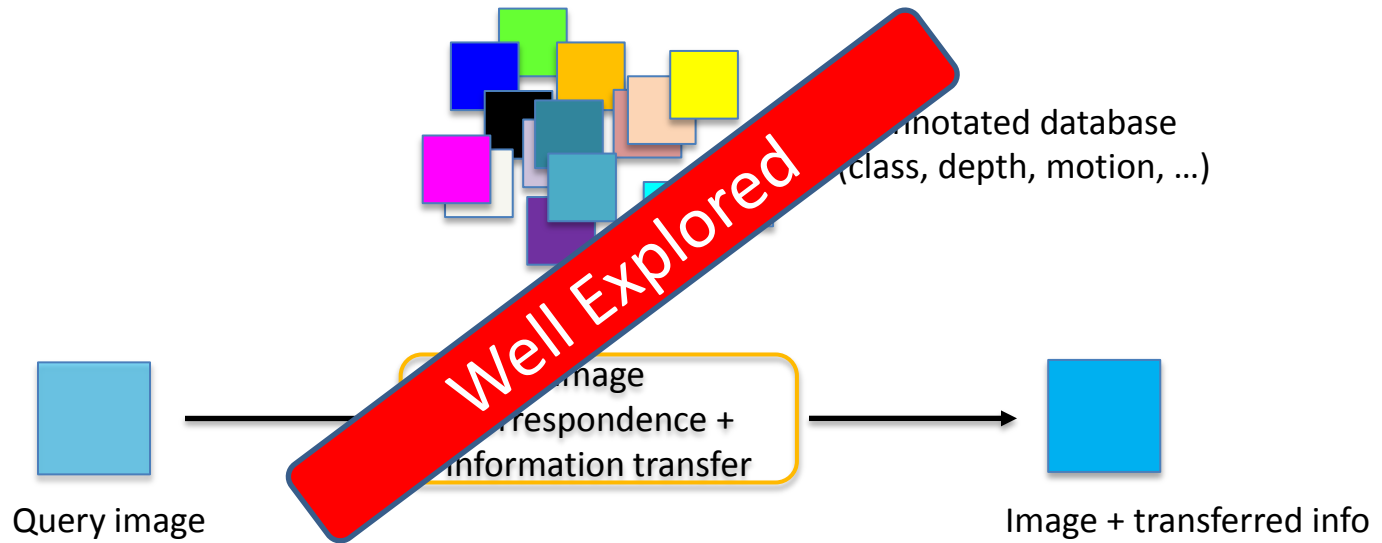
Liu et al. TPAMI'11

Warped candidates and depths



Karsch et al. ECCV'12

# How to *densely* label new images?



# Big Visual Data

Pixel labels usually unavailable!

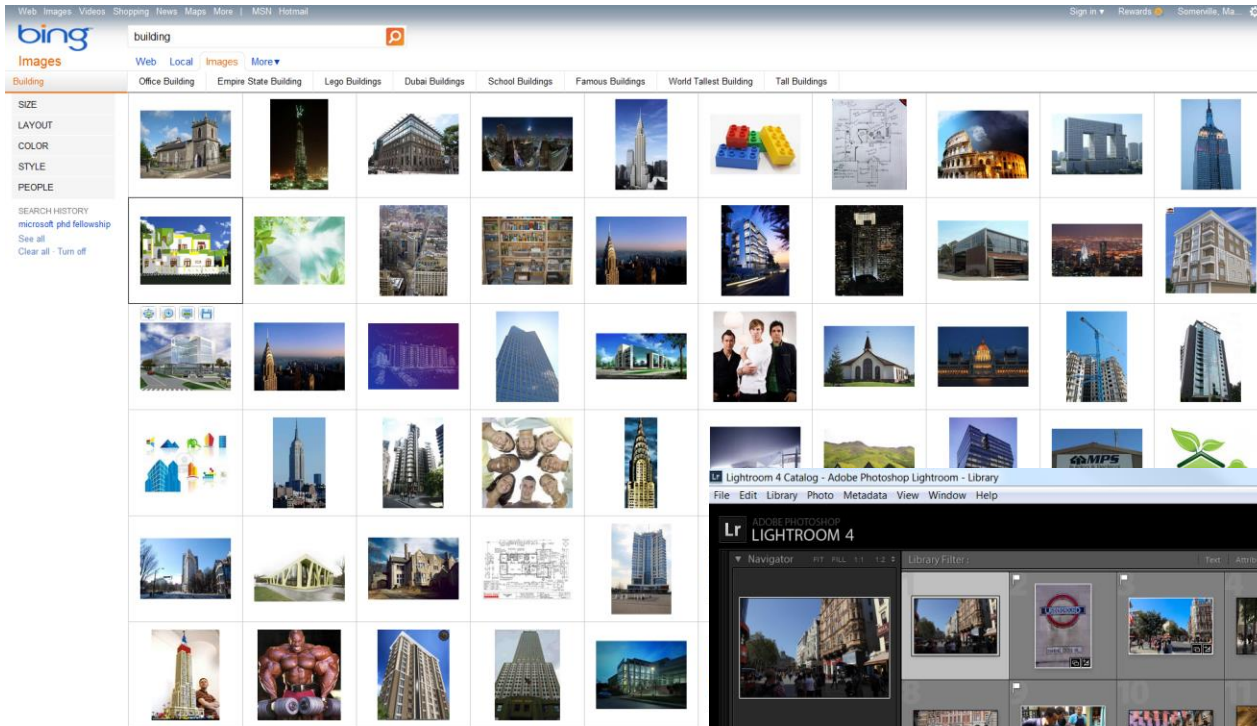
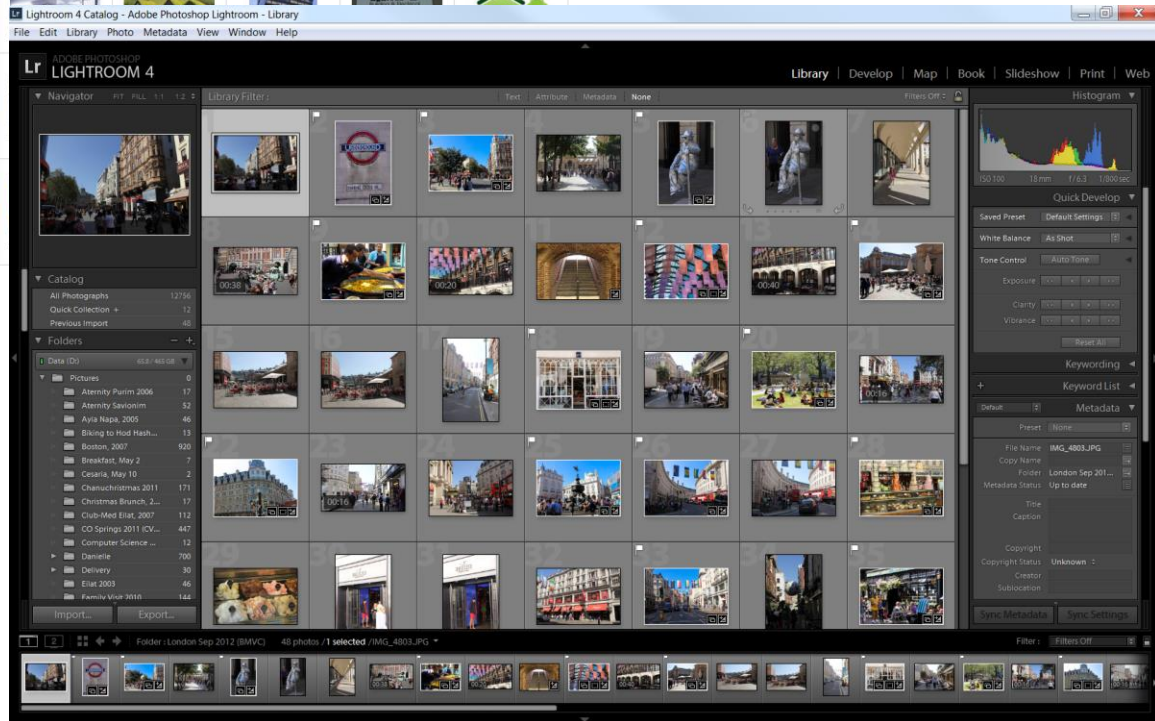


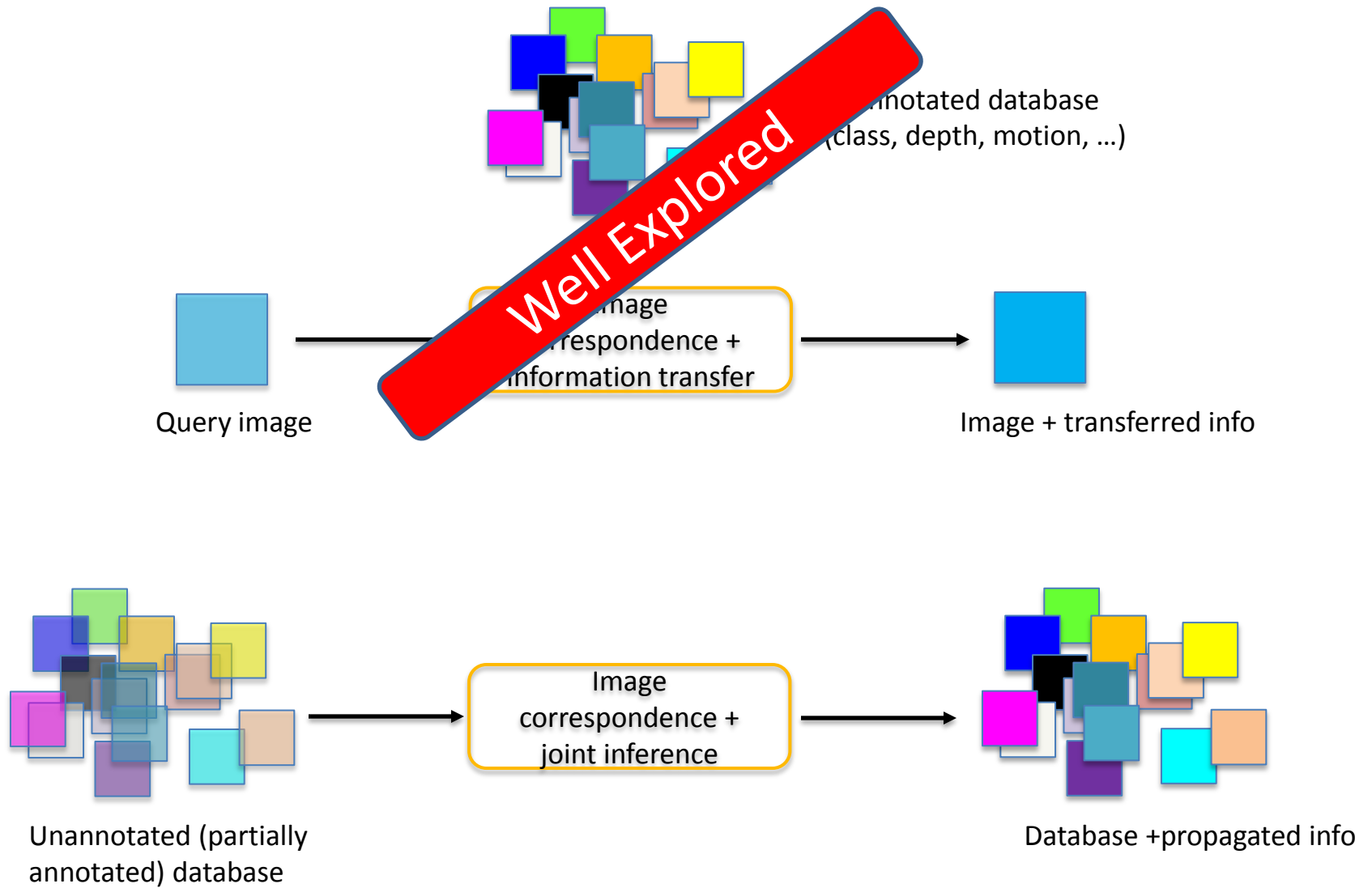
Photo collections



Internet



# How to *densely* label new images?

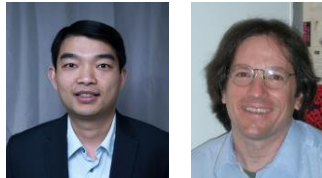


# Joint Inference for Image Databases

- Weakly supervised

## Annotation Propagation in Large Image Databases via Dense Image Correspondence (ECCV 2012)

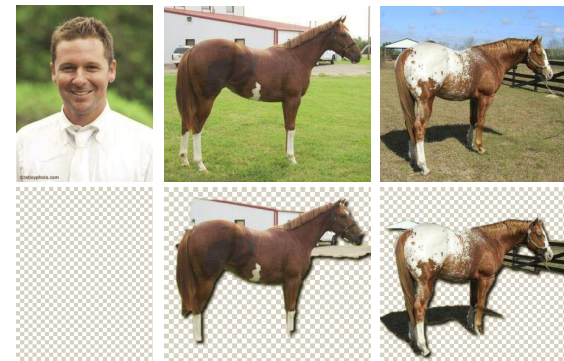
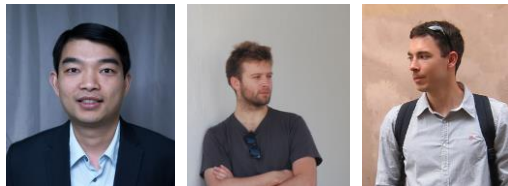
With Ce Liu, William T. Freeman



- Unsupervised

## Unsupervised Joint **Object Discovery and Segmentation** in Internet Images (CVPR 2013)

With Ce Liu, Armand Joulin, Johannes Kopf



# Annotation Propagation

**Input:** A large database of images where only some are tagged and very few (possibly none) are densely labeled



tree, sky, river  
mountain



sky, mountain  
tree



sidewalk, road, car  
building, tree, sky



























sky, river  
building, bridge

# Annotation Propagation

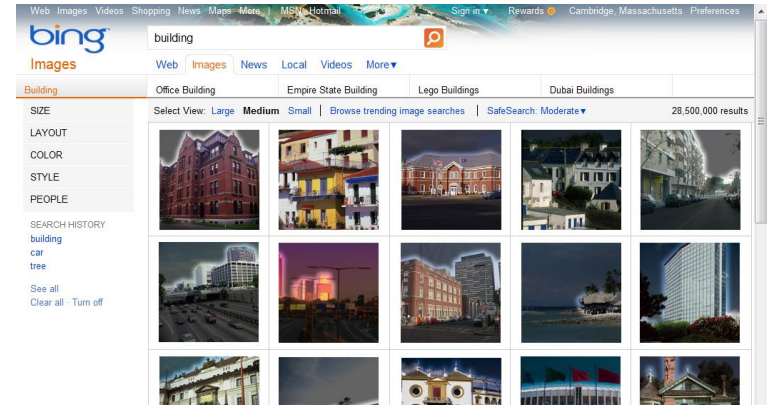
**Output:** The same database with all the pixels labeled and all the images tagged

<p>tree, sky, river mountain</p>	<p>tree, sky, road car, building</p>	<p>tree, sky, plant grass</p>	<p>mountain, field, building, sky, tree</p>
<p>sky, mountain tree</p>	<p>tree, sky, sidewalk road, car, building</p>	<p>sidewalk, road, car building, tree, sky</p>	<p>tree, sky, car building</p>
<p>tree, staircase, sky road, plant, door sidewalk, car, building</p>	<p>tree, sky, river person, mountain</p>	<p>sky, building tree</p>	<p>sky, river building, bridge</p>

-  window
-  tree
-  streetlight
-  staircase
-  sky
-  sign
-  sidewalk
-  road
-  river
-  pole
-  plant
-  person
-  mountain
-  grass
-  field
-  fence
-  door
-  crosswalk
-  car
-  bus
-  building
-  bridge
-  balcony
-  awning

# Dense pixel/region labeling is important

- Enhanced image search



- Constructing training sets for detectors/classifiers



PASCAL 2012

- Image editing
  - User edit propagation



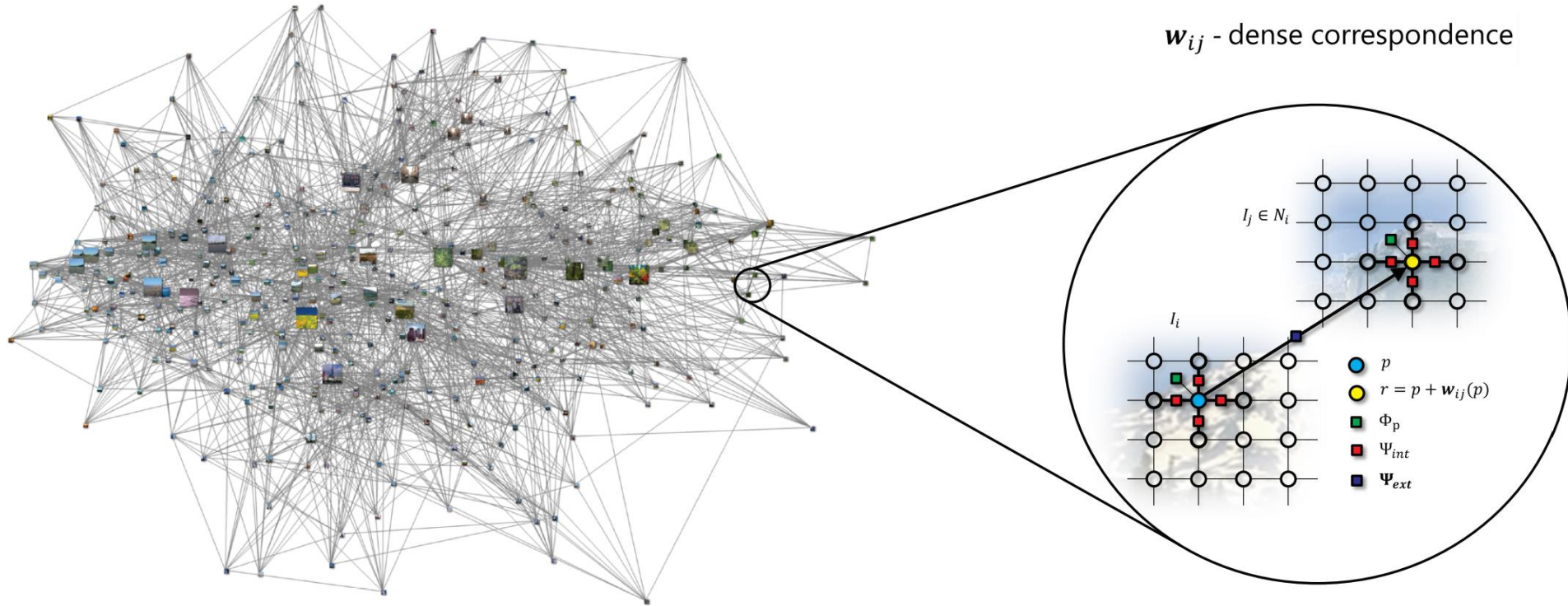
HaCohen et al. 2013

# Pixel-wise image graph

$P(\text{word} \mid I(\mathbf{p}))$  – using machine learning

$$E(\mathbf{C}) = \sum_{i=1}^N \sum_{\mathbf{p} \in \Lambda_i} \left[ \underbrace{\Phi_{\mathbf{p}}(c_i(\mathbf{p}))}_{\text{Local evidence}} + \underbrace{\sum_{\mathbf{q} \in N_{\mathbf{p}}} \Psi_{\text{int}}(c_i(\mathbf{p}), c_i(\mathbf{q}))}_{\text{Intra-image regularization}} + \underbrace{\sum_{j \in N_i} \Psi_{\text{ext}}(c_i(\mathbf{p}), c_j(\mathbf{p} + \mathbf{w}_{ij}(\mathbf{p})))}_{\text{Inter-image regularization}} \right]$$

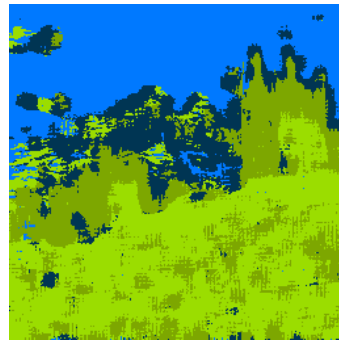
$\mathbf{w}_{ij}$  - dense correspondence



# Inference Results



Input image



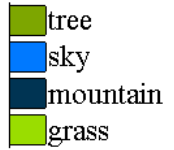
MAP appearance



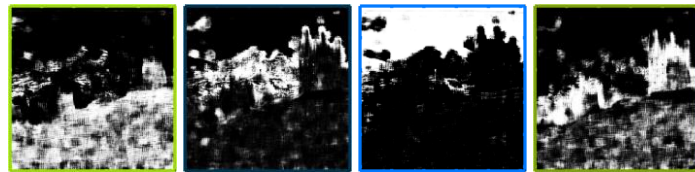
+ Intra-image reg.



+ Inter-image reg.



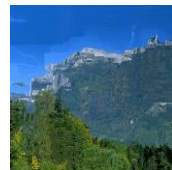
Input image local evidence →



Neighbors



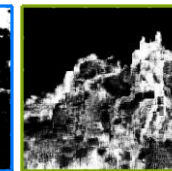
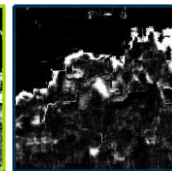
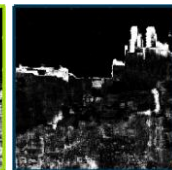
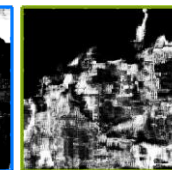
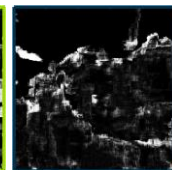
Dense corr.



Neighbors warped



Neighbors local evidence warped



# Optimization

1. Initialize from given tags and labels
2. Repeat until convergence
  - 2.1. Update appearance model parameters  $\Theta$  and compute local evidences  $P(c_i(\mathbf{p}); \Theta)$
  - 2.2. For each image, repeat until convergence
    - 2.2.1. Intra-image message passing
    - 2.2.2. Update color model  $\mathbf{h}_{i,l}^c$
  - 2.3. Inter-image message passing
  - 2.4. Compute MAP label estimate  $C^* = \arg \min_C E(C)$
  - 2.5. Update spatial prior  $\mathbf{h}_l^s$  and co-occurrence matrix  $\mathbf{h}^o$  from  $C^*$
3. Output  $C^*$

Appearance Modeling



Propagation

- Coordinate descent, iterating between estimating the appearance model (learning) and tag propagation (inference)
- Lots of engineering, but nothing revolutionary
  - Partition message passing into intra- and inter-image updates
  - Intra-image message passing on separate cores
  - Parallel inter-image message passing



# From stronger local evidence to weaker local evidence



Input image



Local evidence + intra-image reg.

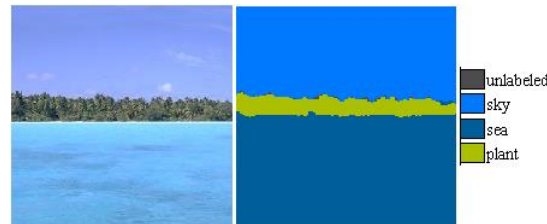


+ Inter-image reg.

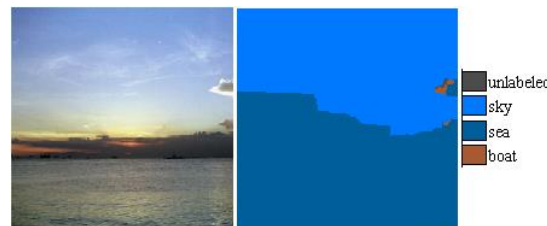
sky  
sea  
mountain



sky  
sea  
mountain  
building



unlabeled  
sky  
sea  
plant



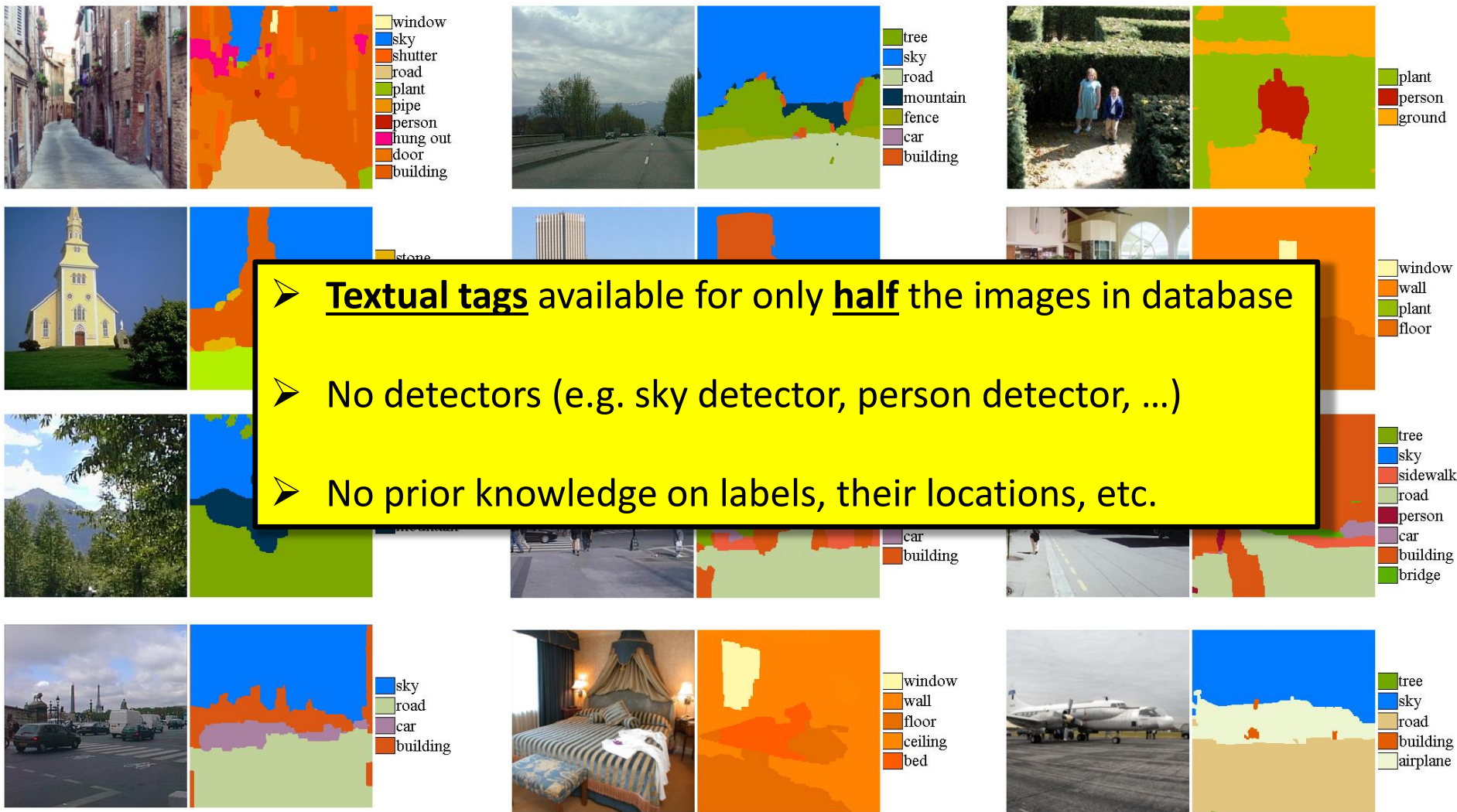
unlabeled  
sky  
sea  
boat

neighbors



warp

# Results on SUN Dataset



SUN dataset [Xiao et al. 2010] - 9556 images, 522 labels

# Joint Inference for Image Databases

- Weakly supervised

Annotation Propagation in Large Image Databases via Dense Image Correspondence (ECCV 2012)

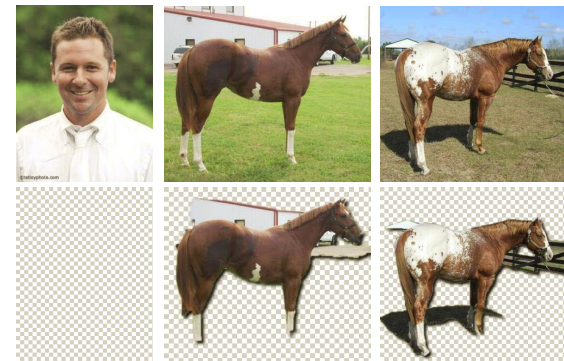
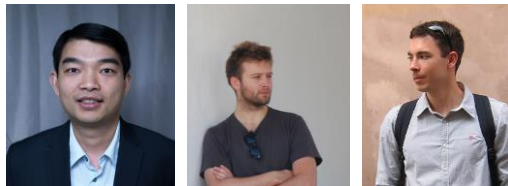
With Ce Liu, William T. Freeman



- Unsupervised

Unsupervised Joint **Object Discovery and Segmentation** in Internet Images (CVPR 2013)

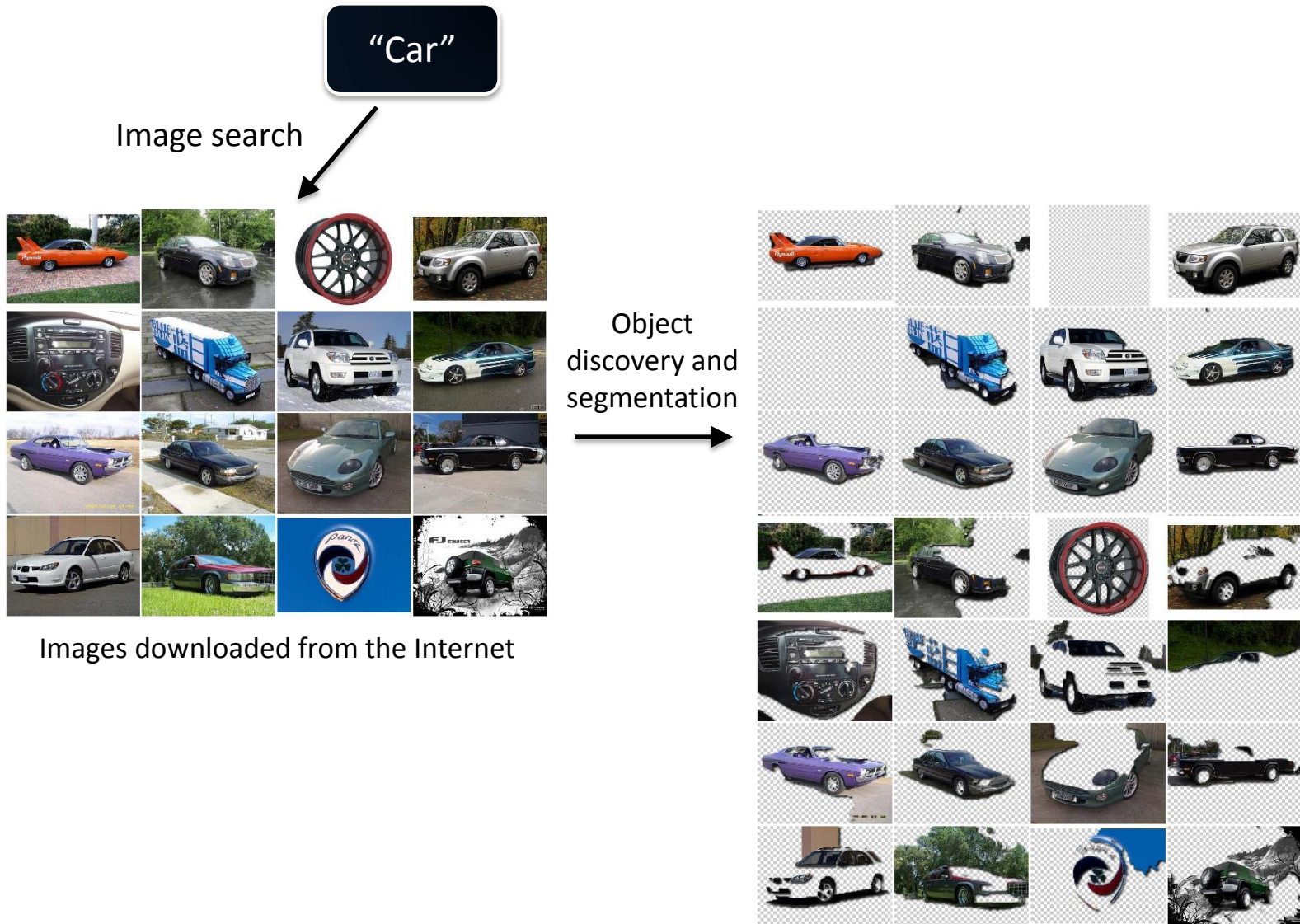
With Ce Liu, Armand Joulin, Johannes Kopf



# Object discovery and Co-segmentation

- **Input**: A set of images containing some “common object”
- **Output**: Every pixel in the dataset marked as belonging or not belonging to the “common object”
- **No additional information on the images or the object class**

# Object discovery and Co-segmentation



State of the art co-segmentation [Joulin et al. CVPR 2012]

# Benchmark “plane” Dataset (MSRC)



4\_7\_s.bmp



4\_8\_s.bmp



4\_9\_s.bmp



4\_10\_s.bmp



4\_11\_s.bmp



4\_12\_s.bmp



4\_13\_s.bmp



4\_14\_s.bmp



4\_15\_s.bmp



4\_16\_s.bmp



4\_17\_s.bmp



4\_18\_s.bmp



4\_19\_s.bmp



4\_20\_s.bmp



4\_21\_s.bmp



4\_22\_s.bmp



4\_23\_s.bmp



4\_24\_s.bmp



4\_25\_s.bmp



4\_26\_s.bmp



4\_27\_s.bmp



4\_28\_s.bmp



4\_29\_s.bmp

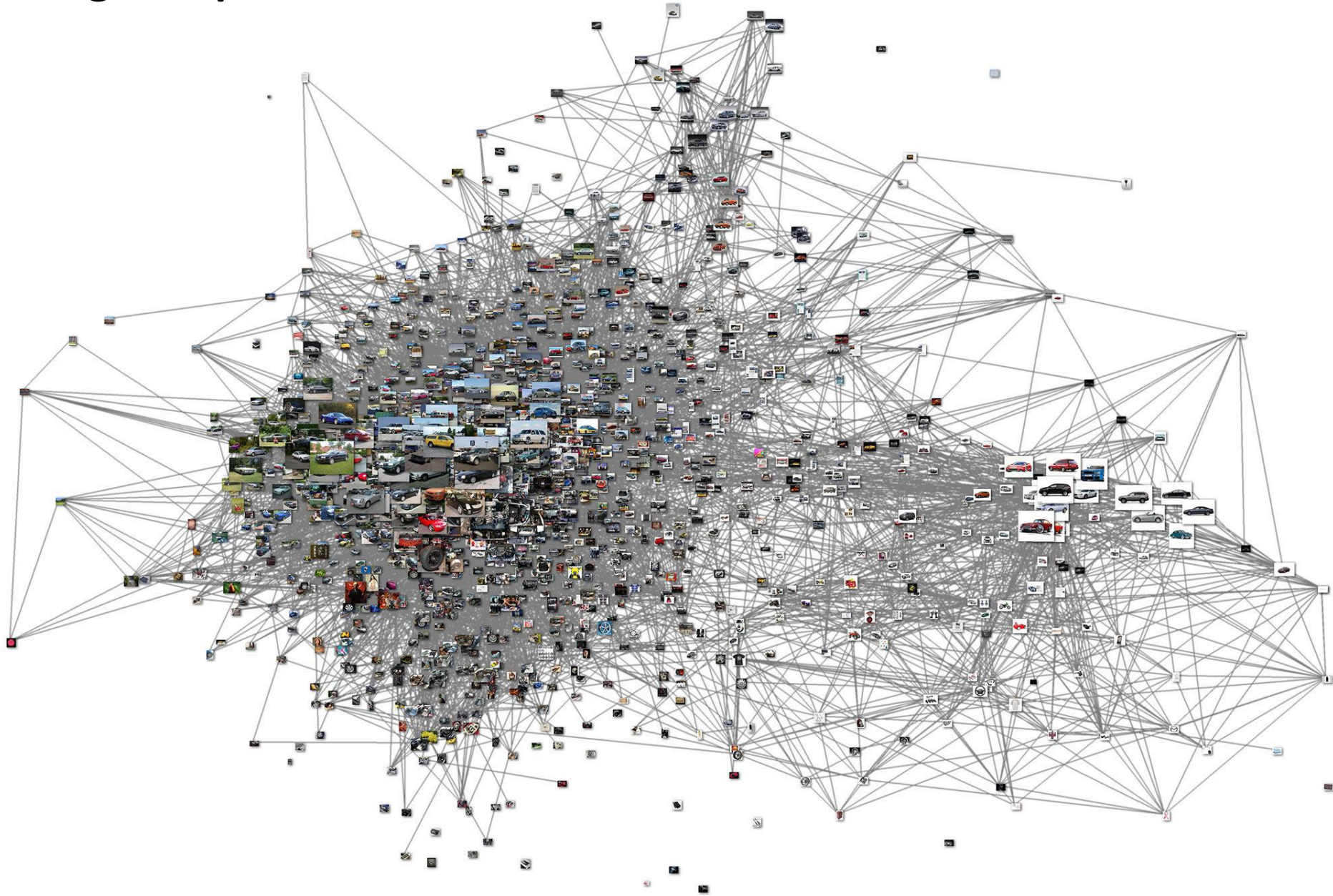


4\_30\_s.bmp

# Real-world "plane" Dataset (Internet Search)



# Image Graph





# Basic Idea

- Pixels (features) belonging to the common object should be:

**1. *Salient*** - Dissimilar to other pixels (features) in their image

Captured by image *saliency measures*

**2. *Sparse*** - Similar to other pixels (features) in other images (with respect to smooth transformations)

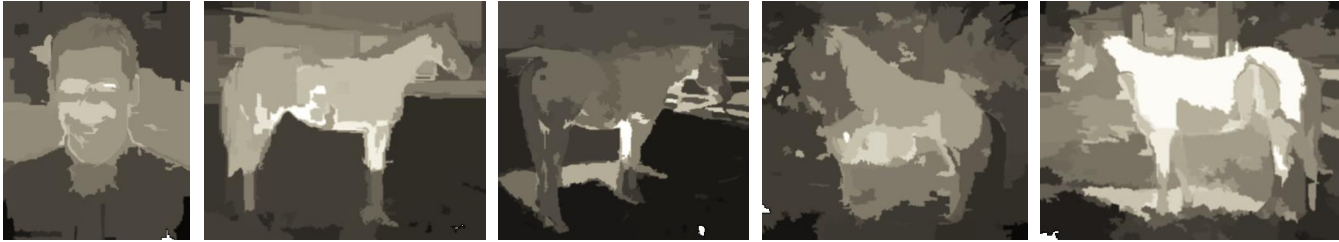
Captured by (dense) **image correspondence**

# One of these things is not like the others

Source



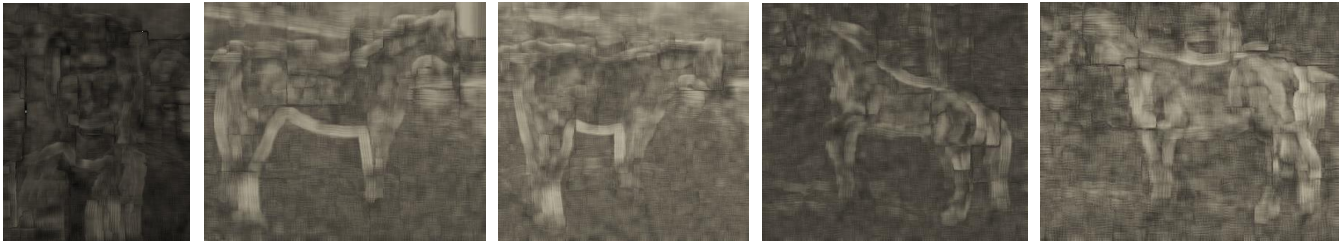
Saliency



Warped neighbor



Matching Score



Segmentation



# One of these things is not like the others

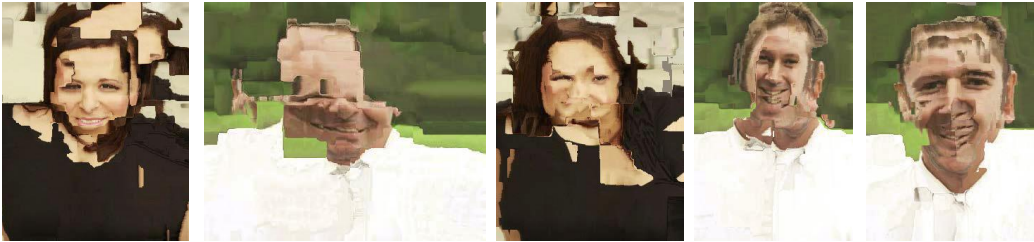
Source



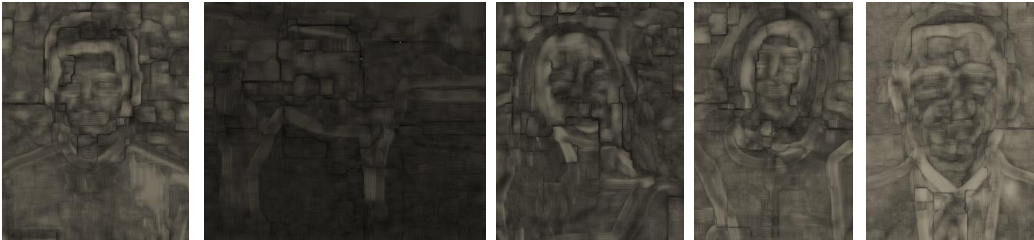
Saliency



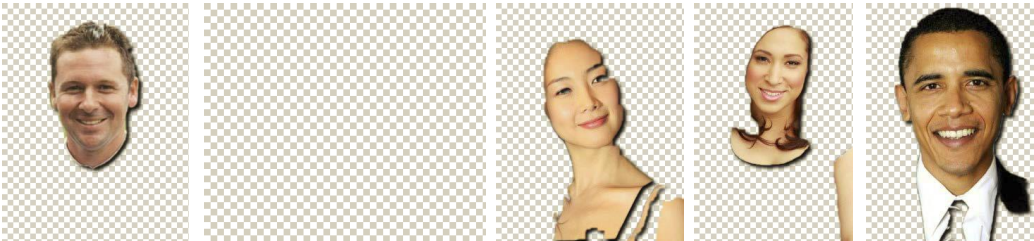
Warped neighbor



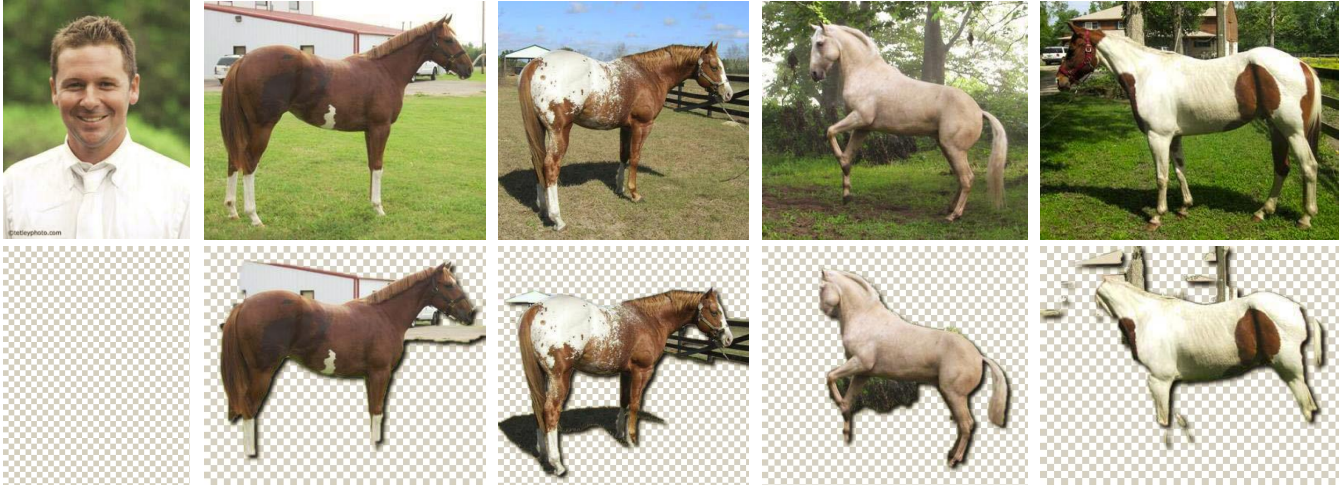
Matching Score



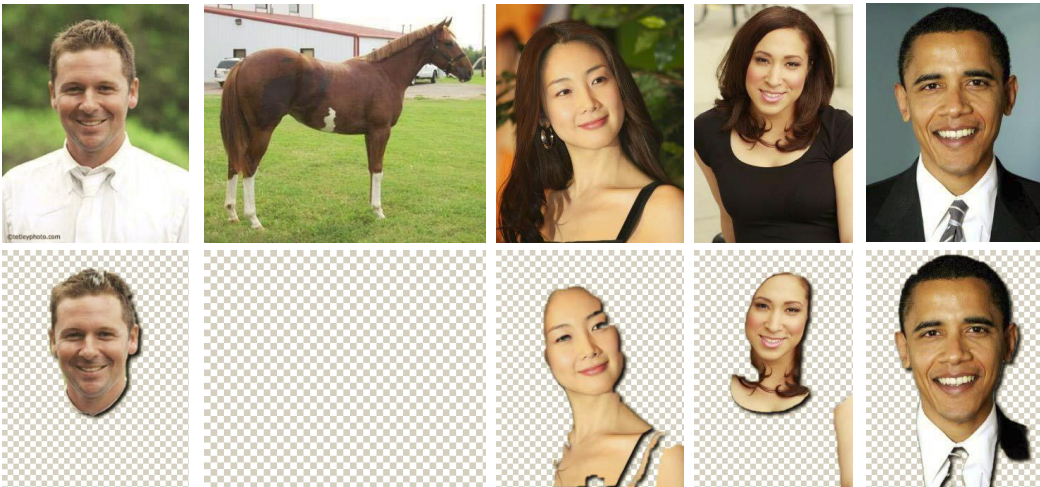
Segmentation



# One of these things is not like the others

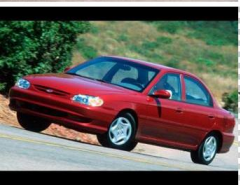
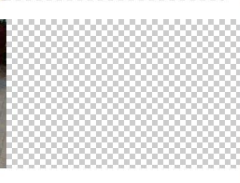
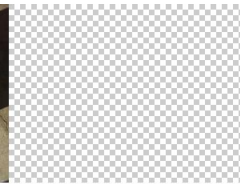


Horse

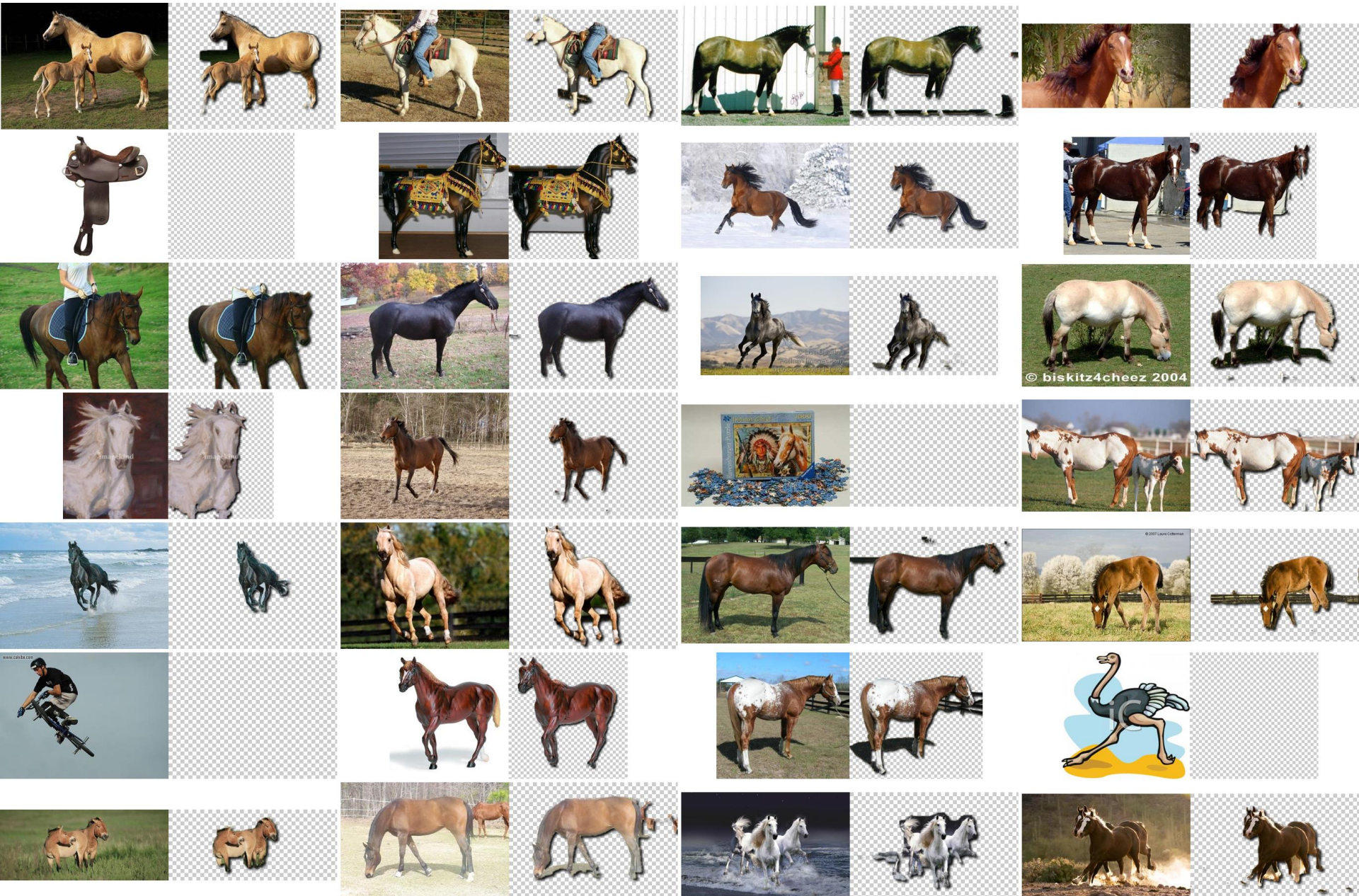


Face

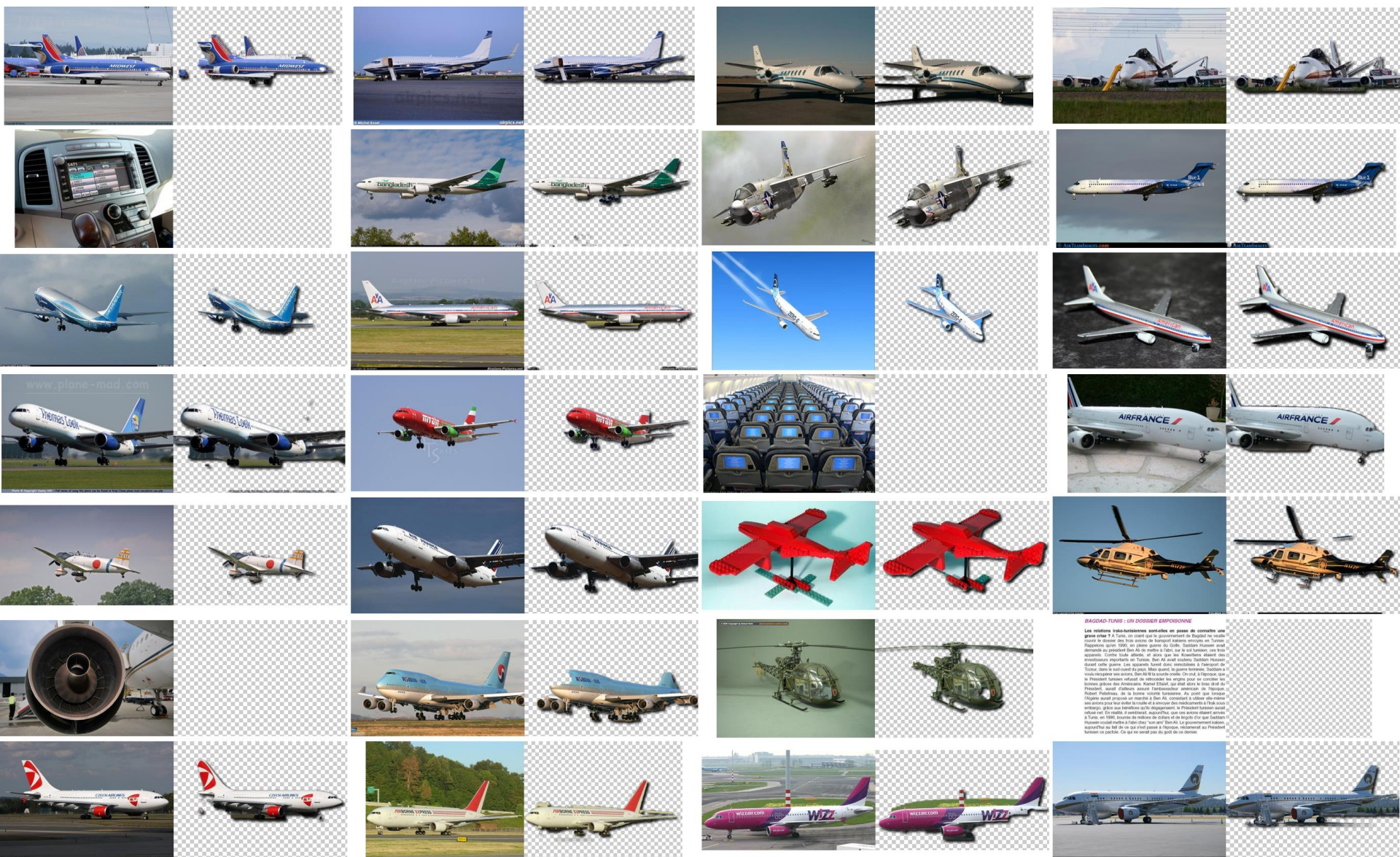
# Car (4,347 images, 11% noise)



# Horse (6,381 images, 7% noise)

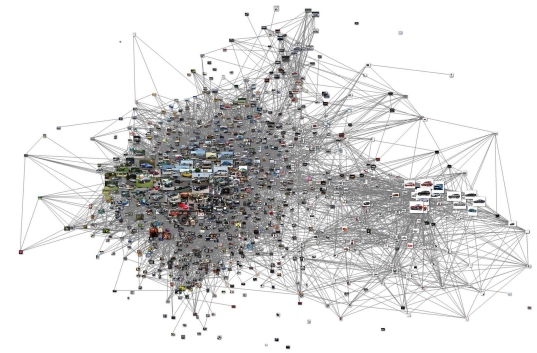


# Airplane (4,542 images, 18% noise)



# Conclusion

- Labels in big visual data are often unavailable/noisy
- Dense image correspondence (SIFT flow, and others) useful to capture structure, resolve visual ambiguity
  - Becoming a mature technology
- Joint inference for weakly-labeled image databases
  - Annotation Propagation: partial tags + very few (possibly none) pixel labels
  - Object discovery and segmentation: only assuming some underlying “common object”







**Thank you!**

Michael Rubinstein

MIT CSAIL