Image Completion / Inpainting

- Goal: Remove object(s) from digital photographs, and then fill the hole with information extracted from the surrounding area, **preserving image structures**
- Issues:
 - Underconstrained problem
 - Requires inference about regions, boundaries, textures
 - Data-driven but needs scene semantics
 - Result should "look real" without seams or perceptual artifacts







Photoshop CS5 "Content-Aware Fill"



1

Fundamental Issue: Texture vs. Structure





"Onion Skin" Filling Order



Fill Order



- In what order should we fill the pixels?
 - choose pixels that have more neighbors already filled
 - choose pixels that are continuations of lines/curves/edges so that linear structures are propagated

Criminisi's Approach

- Combine the strengths of texture synthesis and inpainting approaches
 - Use a texture synthesis algorithm to quickly fill patches with similar texture regions
 - Fill patches near linear structures first, thereby finding structure/texture boundaries first, and then filling in the texture inside regions
- Result
 - Preserve linear structures
 - Avoid blurring

Image Completion by Example-Based Inpainting

• A. Criminisi, P. Perez, and K. Toyama, CVPR 2003



Best-First Filling Algorithm

- Extract the initial fill-front, δΩ, and compute priority values, *P*(**p**), for all pixels **p** in δΩ
- · Repeat until no more unfilled pixels remain:
 - 1. Find target pixel **p** and patch $\Psi_{\mathbf{p}}$ with highest "priority," $P(\mathbf{p})$
 - 2. Find source patches, $\Psi_{\bm{q}'} \; \Psi_{\bm{q}''}$, that most closely match the target patch
 - 3. Paste most similar source patch into target location
 - 4. Update fill front and priority values





Derivatives with Convolution	
For 2D function, $f(x, y)$, the partial derivative is:	
$\frac{\partial f(x, y)}{\partial x} = \lim_{\varepsilon \to 0} \frac{f(x + \varepsilon, y) - f(x, y)}{\varepsilon}$	
For discrete data, we can approximate using <i>finite</i> differences:	
$\frac{\partial f(x,y)}{\partial x} \approx \frac{f(x+1,y) - f(x,y)}{1}$	
To implement using convolution, what would be the associated filter?	















http://research.microsoft.com/en-us/projects/i3l/patchworks.aspx

Results: Filling Large Regions

12% of image area filled



Original



Object Cut

Texture Synthesis

Criminisi



Results: Filling Large Regions





10% of image area filled



Inpainting Examples



Inpainting Examples



Image Completion with Structure Propagation

J. Sun, L. Yuan, J. Jia, and H. Shum SIGGRAPH 2005

Image Completion with Structure Propagation

- The method of Criminisi *et al.* does **not** guarantee continuity of salient structures such as curves or junctions
- Missing structure is hard to recover automatically, but can be easily specified manually



Example Result



Algorithm

- 1. User *interactively* draws a few curves on the image that extend from given structures into hole, specifying most salient missing structures and partitioning image into regions
- 2. Synthesize patches along the curves, defining structure boundaries
- 3. Fill in remaining holes using patch-based texture synthesis with samples from same segmented regions











Scene Completion Using Millions of Photographs

James Hays and Alexei A. Efros SIGGRAPH 2007

Slides by J. Hays and A. Efros























Data

<u>2.3 Million</u> unique images from Flickr groups and keyword searches



landscapes, travel, city, outdoors, vacation, etc.





Challenges

- Computational
 - How to efficiently search millions or billions of source images?
- Finding semantically appropriate image fragments
 - Generically (e.g., car) and specifically (red BMW SUV)
- Differences in appearance
 - Lighting, color, shadows, etc.

Finding a Good Fill Patch

- 1. Find a set of candidate images that are "semantically similar scenes" to the input image
 - Don't assume image tags or structured datasets (e.g., ImageNet)
- 2. Find most similar patch from each candidate image to fill hole





How to Find Semantically-Similar Scenes? The "Gist" of a Scene





If this is a street, this must be a pedestrian

Gist: abstract meaning of scene; low-dim., global image rep

- Obtained within 150 ms (Biederman, 1981, Thorpe S. et.al 1996)
- · Obtained without attention (Oliva & Schyns, 1997, Wolfe, J.M. 1998)
- Possibly derived via statistics of low-level structures (e.g. Swain & Ballard, 1991)

Image Representations

- The space of all possible images is huge!
 - $-32 \times 32 \times 8$ bit image $\Rightarrow 10^{7400}$ possible images
 - In 100 years a person only sees about 10¹¹ images
- But the number of "natural" / "real" images is much, much smaller and there are semantic clusters
- How many images are needed to be able to find a similar image to match any given natural input image?

What is the Gist? A Scene-based Image Descriptor

What features are sufficient to represent a scene?

- Statistics of local, low-level features
- Color histograms
- Oriented band-pass filters



80 Million Tiny Images A. Torralba, R. Fergus, and W. T. Freeman



Dataset available at: groups.csail.mit.edu/vision/TinyImages









Learning a Hierarchy of Feature Extractors

• Each layer of hierarchy extracts features from output of previous layer

.....

- All the way from pixels \rightarrow classifier / features
- Layers have (nearly) the same structure















Step 2: Context Matching to Find Best Patch



Query image template = band within 80 pixels of hole boundary
SSD match with all translations and 3 scales of database image







Result Ranking

Assign each of the 200 results a score that is the sum of:



Scene matching distance



Context matching distance (color + texture)



Graph-cut cost





















































































Image Completion: Summary

- The key challenge is propagating the image "structure"
- Approaches:
 - For simple enough problems, the heuristic of extending existing edges (Criminisi *et al.*) is sufficient
 - For more complex problems, the user can provide the high-level structure information (Sun *et al.*)
 - Alternatively, high-level information can be obtained from a large database of images (Hays and Efros)