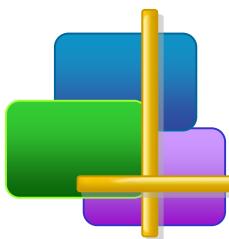


(I'm Fun) Digital Image Fundamentals



Week 6: Segmentation

Thuong Nguyen Canh

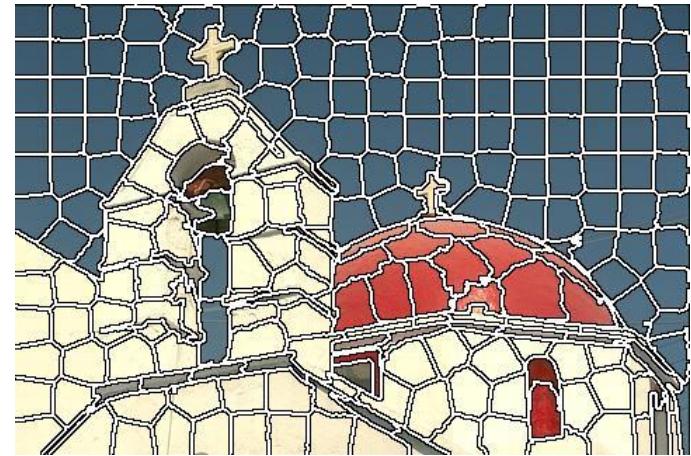
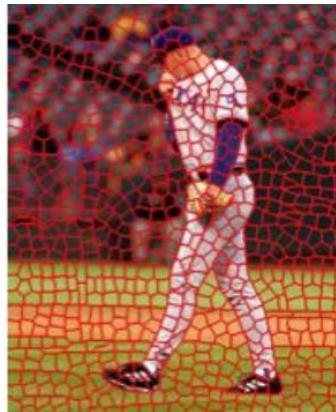
Institute for Datability Science, Osaka University

November 2019

Image Segmentation?

- Group pixels into “meaningful” regions that share some similar properties

“superpixels”



[X. Ren and J. Malik 2003, Mori et al. 2005]

- Segment images into meaningful objects

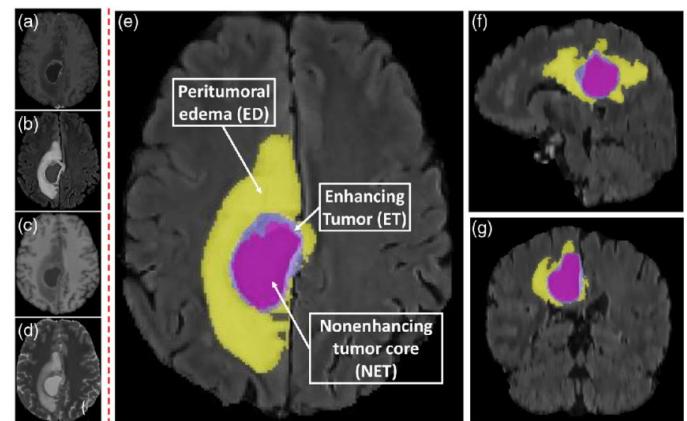


Image segmentation is the process of assigning a label to every pixels

Semantic Segmentation

- Label every pixel: recognize the class of every pixel
- Do not differentiate instances



Mottaghi et al, “[The role of context for object detection and semantic segmentation in the wild](#)”, CVPR 2014

Instance Segmentation

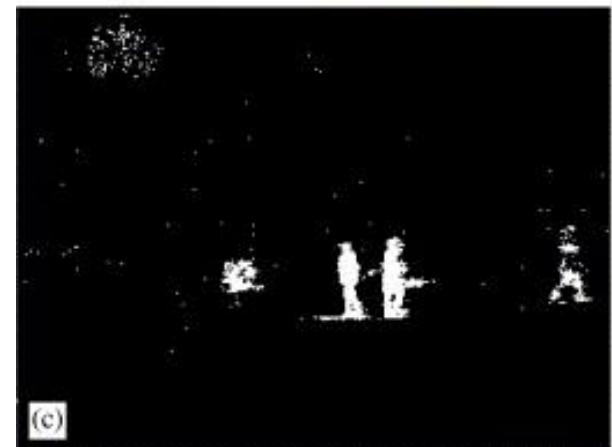
- Detect instances, categorize and label every pixel
- Labels are class-aware and instance-aware



Arnab,Torr “[Pixelwise instance segmentation with a dynamically instantiated network](#)”, CVPR 2017

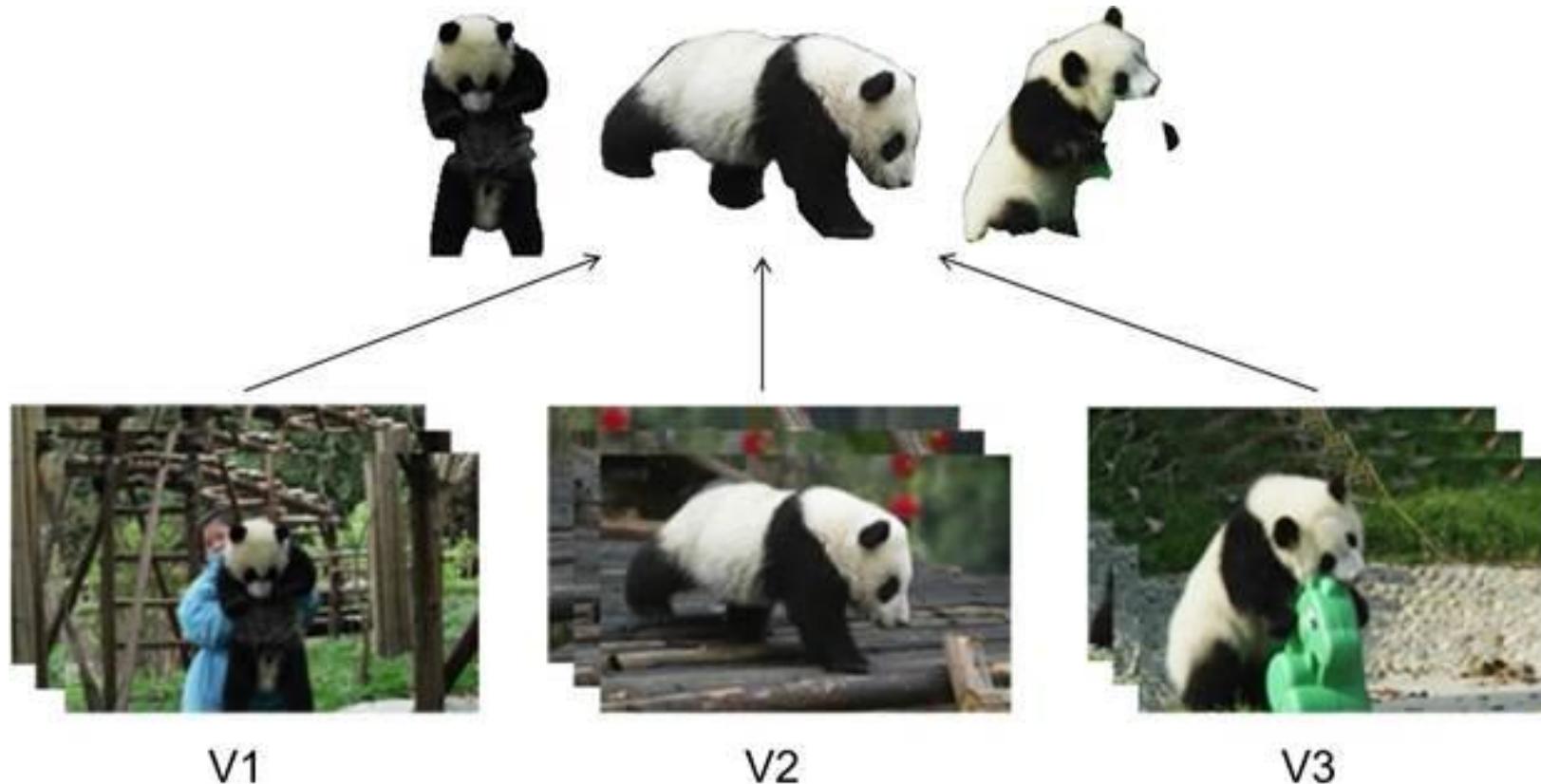
Figure-ground Segmentation

- An binary semantic segmentation



Co-segmentation

- Segment the common objects



Applications



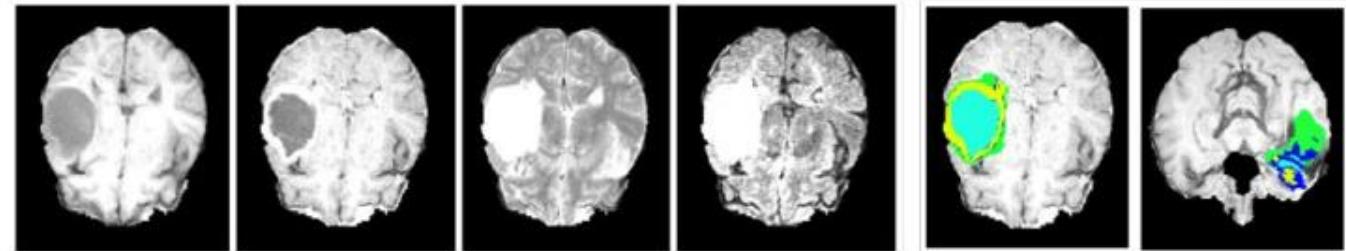
Image editing and composition (Xu, 2016)

Robotics

Autonomous driving
(cordts, 2016)



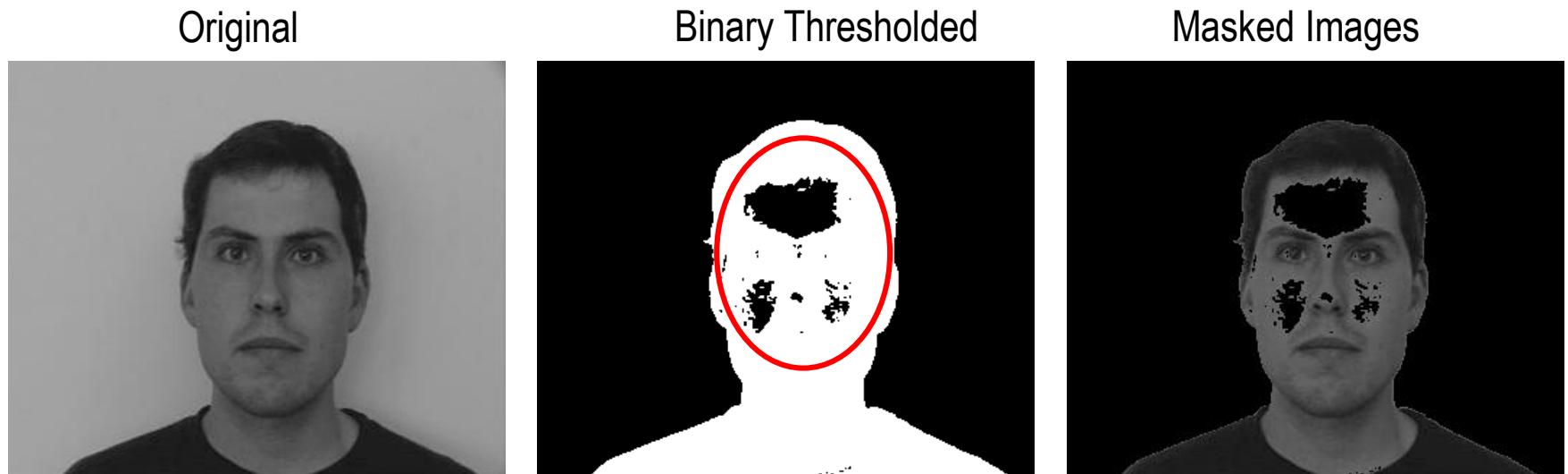
Medical image analysis
(Casamitjana, 2017)



Segmentation Method

- ➊ Segmentation as “clustering”
 - ➌ Clustering: Group similar data points and represent them with a single token
 - ➌ Key challenges
 - ➍ What make two points/images/patches similar?
 - ➍ How to compute overall grouping?
- ➋ Thresholding based
- ⌁ Edge based
- ⌂ Region based
- ⌃ Energy based
- ⌄ Learning based

Gray-Level Thresholding



$$f[x, y] \in \mathbb{N}^{286 \times 334}$$

Range 0 ~ 255

$$m[x, y] \in \mathbb{B}^{286 \times 334}$$

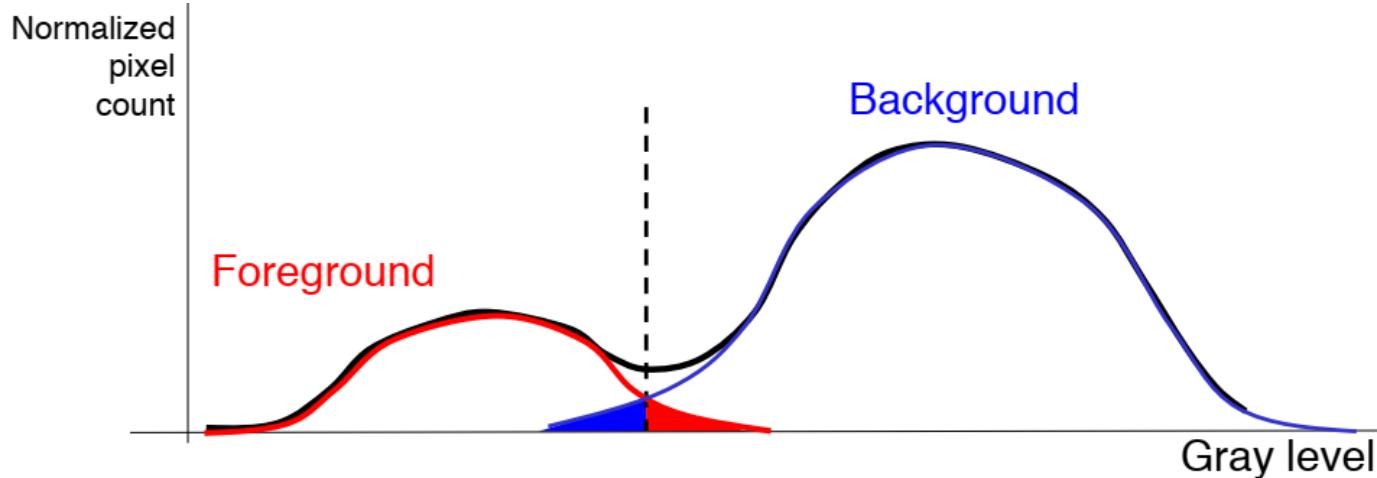
Fixed Threshold 105

$$f[x, y] \cdot m[x, y]$$

Pixel-wise multiplication

How can holes be filled?

How to choose the threshold?



- Unsupervised thresholding
 - Find threshold T that minimizes within-class variance

$$\sigma_{within}^2(T) = \frac{N_{Fgnd}(T)}{N} \sigma_{Fgnd}^2(T) + \frac{N_{Bgrnd}(T)}{N} \sigma_{Bgrnd}^2(T)$$

$$\sigma_{between}^2(T) = \sigma^2 - \sigma_{within}^2(T)$$

Maximize between-class variance

$$\begin{aligned} &= \left(\frac{1}{N} \sum_{x,y} f^2[x,y] - \mu^2 \right) - \frac{N_{Fgnd}}{N} \left(\frac{1}{N_{Fgnd}} \sum_{x,y \in Fgnd} f^2[x,y] - \mu_{Fgnd}^2 \right) - \frac{N_{Bgrnd}}{N} \left(\frac{1}{N_{Bgrnd}} \sum_{x,y \in Bgrnd} f^2[x,y] - \mu_{Bgrnd}^2 \right) \\ &= -\mu^2 + \frac{N_{Fgnd}}{N} \mu_{Fgnd}^2 + \frac{N_{Bgrnd}}{N} \mu_{Bgrnd}^2 = \frac{N_{Fgnd}}{N} (\mu_{Fgnd} - \mu)^2 + \frac{N_{Bgrnd}}{N} (\mu_{Bgrnd} - \mu)^2 \\ &= \frac{N_{Fgnd}(T) \cdot N_{Bgrnd}(T)}{N^2} (\mu_{Fgnd}(T) - \mu_{Bgrnd}(T))^2 \end{aligned}$$

[Otsu, 1979]

Unsupervised thresholding

- Algorithm: Search for threshold T to maximize

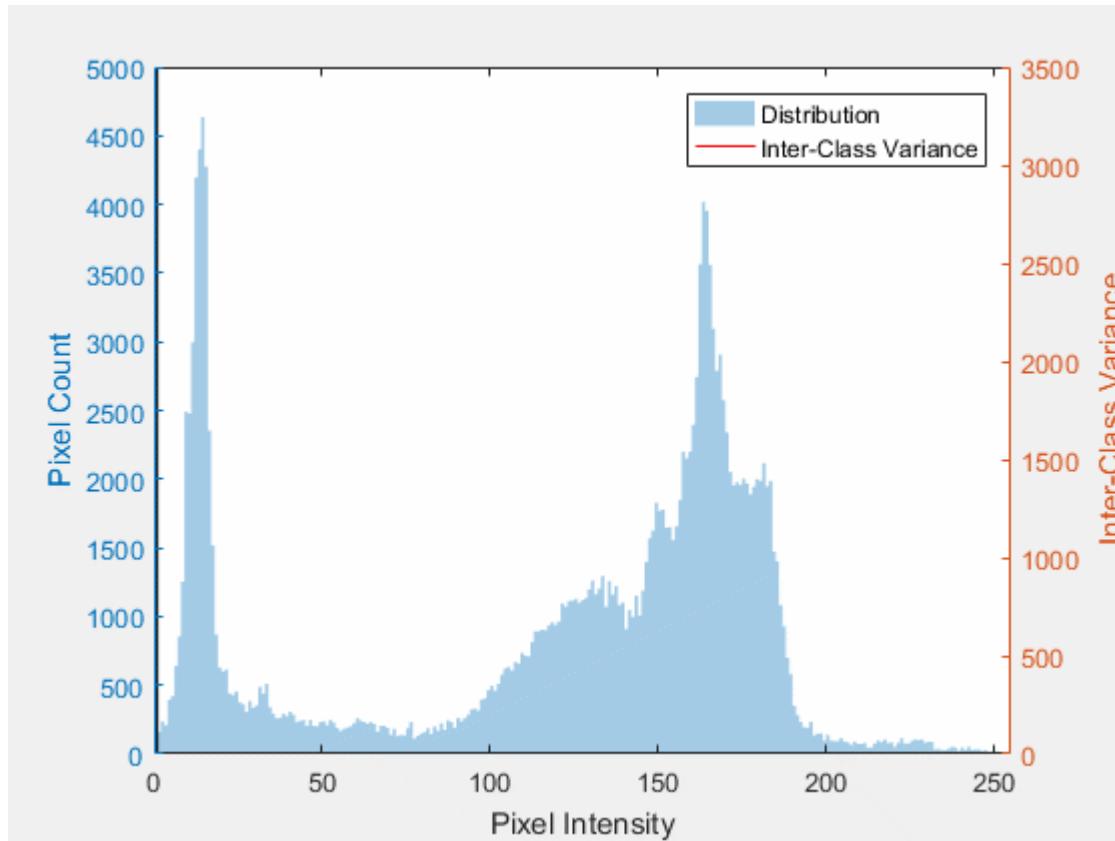
$$\sigma_{\text{between}}^2(T) = \frac{N_{\text{Fgnd}}(T) \cdot N_{\text{Bgrnd}}(T)}{N^2} (\mu_{\text{Fgnd}}(T) - \mu_{\text{Bgrnd}}(T))^2$$

- Useful recursion for sweeping T across histogram:

$$N_{\text{Fgnd}}(T+1) = N_{\text{Fgnd}}(T) + n_T$$
$$N_{\text{Bgrnd}}(T+1) = N_{\text{Bgrnd}}(T) - n_T$$
$$\mu_{\text{Fgnd}}(T+1) = \frac{\mu_{\text{Fgnd}}(T)N_{\text{Fgnd}}(T) + n_T T}{N_{\text{Fgnd}}(T+1)}$$
$$\mu_{\text{Bgrnd}}(T+1) = \frac{\mu_{\text{Bgrnd}}(T)N_{\text{Bgrnd}}(T) - n_T T}{N_{\text{Fgnd}}(T+1)}$$

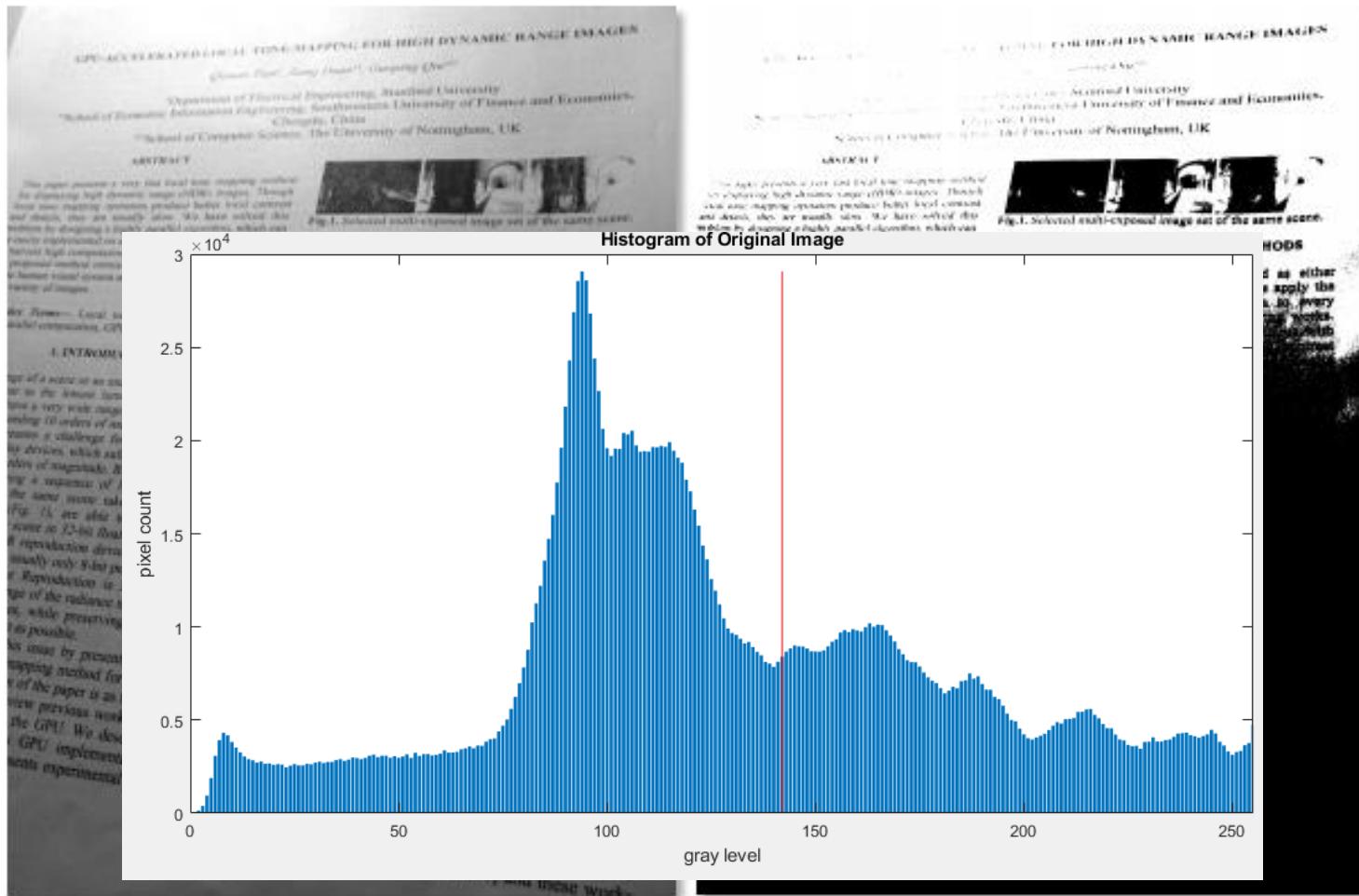
[Otsu, 1979]

Unsupervised thresholding



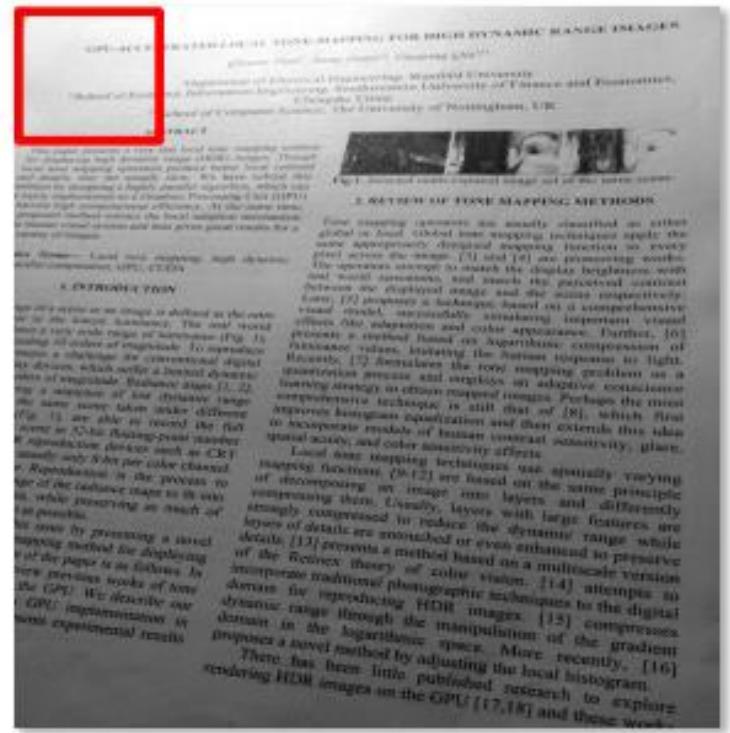
Locally Adaptive Thresholding

- A single thresholding value might not work



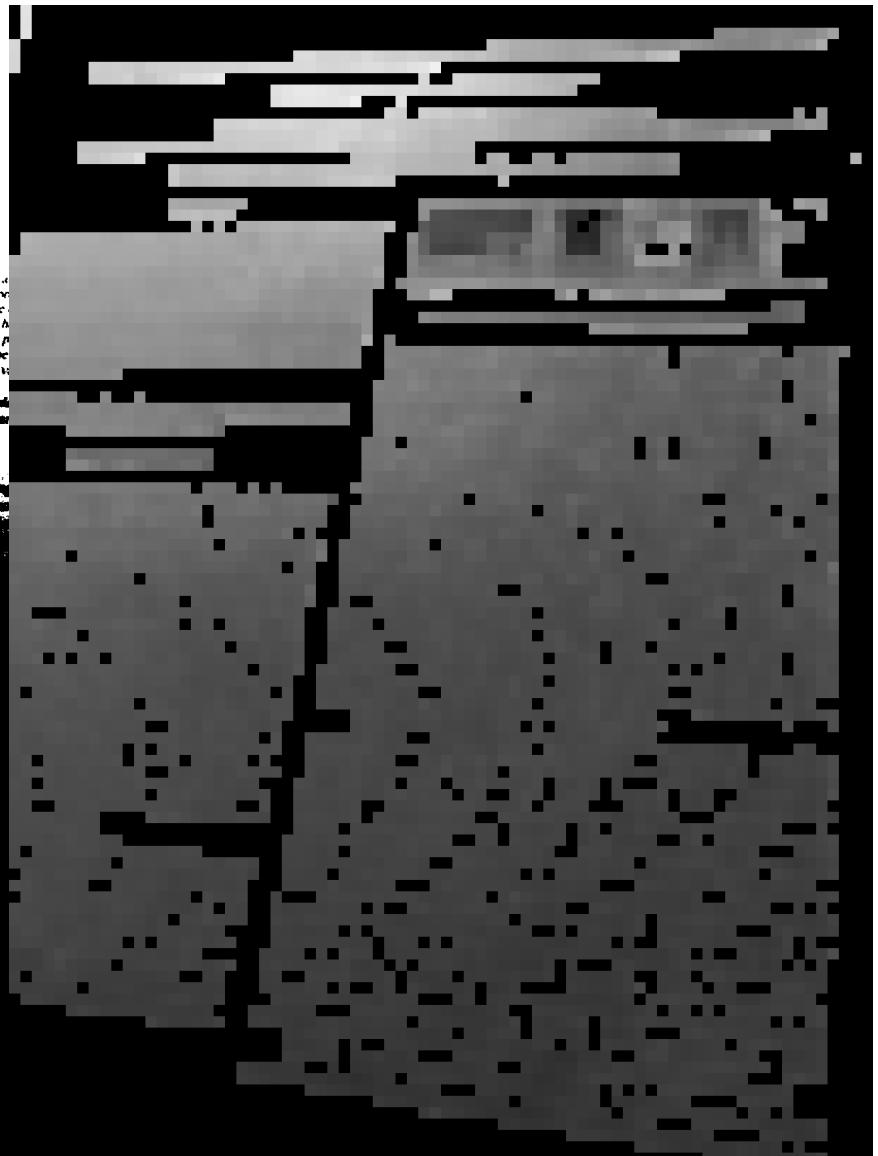
Locally Adaptive Thresholding

- Slide a window over the image
- For each window position, decide whether to perform thresholding
 - Thresholding should not be performed in uniform areas
 - Use variance or other suitable criterion
- Non-uniform areas: apply Otsu's method (based on local histogram)
- Uniform areas: classify the entire area as foreground or background based on mean value



Local tone mapping techniques use spatially varying mapping functions [9–12] or based on the same principle of decomposing an image into layers and differently compressing them. Usually, layers with large features are strongly compressed to reduce the dynamic range while layers of details are smoothed or even enhanced to preserve details. [13] presents a method based on a multiscale version of the Retinex theory of color vision. [14] attempts to incorporate traditional photographic techniques to the digital domain for reproducing HDR images. [15] compresses the dynamic range through the manipulation of the gradient domain in the logarithmic space. More recently, [16] proposes a novel method by adjusting the local histogram, rendering HDR images on the GPU [17, 18] and these were

Locally Adaptive Thresholding



GPU-ACCELERATED LOCAL TONE-MAPPING FOR HIGH DYNAMIC RANGE IMAGES

Qiyuan Tian^a, Jiaxin Liang^b, Guoping Ouyang^c

^aDepartment of Electrical Engineering, Stanford University
^bSchool of Economic Information Engineering, Southwest Jiaotong University
^cSchool of Computer Science, The University of Nottingham, UK

ABSTRACT

This paper presents a very fast local tone mapping method for displaying high dynamic range (HDR) images. Though local tone mapping operators produce better local contrast and details, they are usually slow. We have solved this problem by designing a highly parallel algorithm, which can be easily implemented on a Graphics Processing Unit (GPU) to harvest high computational efficiency. At the same time, proposed method mimics the local adaption mechanism in human visual system and thus gives good results for a variety of images.

Keywords— Local tone mapping, high dynamic range computation, GPU, CUDA

I. INTRODUCTION

Age of a scene or an image is defined as the ratio of the highest to the lowest luminance. The real world have a very wide range of luminance (Fig. 1), exceeding 10 orders of magnitude. To reproduce it presents a challenge for conventional digital display devices, which suffer a limited dynamic range of magnitude. Radiance maps [1, 2] store a sequence of low dynamic range images of the same scene taken under different lighting conditions [3]. They are able to record the full range of the radiance maps to fit into 32-bit floating-point number representation of reproduction devices such as CRT monitors and computer displays. Reproduction is the process to convert the radiance maps to fit into 8-bit per color channel reproduction devices such as CRT monitors and computer displays. While preserving as much of the original scene as possible.

In this paper, we propose a novel locally adaptive tone mapping method for displaying HDR images. The rest of the paper is as follows. In Section II, we review previous works of tone mapping methods. In Section III, we describe our GPU implementation in detail. Section IV presents experimental results and conclusions. Finally, Section V concludes the paper.

2. REVIEW OF TONE MAPPING METHODS

Tone mapping operators are usually classified as either global or local. Global tone mapping techniques apply the same appropriately designed mapping function to every pixel across the image. [3] and [4] are pioneering works. The operators attempt to match the display brightness with real-world sensations, and match the perceived contrast between the displayed image and the scene respectively. Later, [5] proposes a technique based on a comprehensive visual model successfully simulating important visual effects like adaptation and color appearance. Further, [6] presents a method based on logarithmic compression of luminance values, imitating the human response to light. Recently, [7] formulates the tone mapping problem as a quantization process and employs an adaptive consciousness learning strategy to obtain mapped images. Perhaps the most comprehensive technique is still that of [8], which first improves histogram equalization and then extends this idea to incorporate models of human contrast sensitivity, glare, spatial acuity, and color sensitivity effects.

Local tone mapping techniques use spatially varying mapping functions. [9-12] are based on the same principle of decomposing an image into layers and differently compressing them. Usually, layers with large features are strongly compressed to reduce the dynamic range while layers of details are untouched or even enhanced to preserve details. [13] presents a method based on a multiscale version of the Retinex theory of color vision. [14] attempts to incorporate traditional photographic techniques to the digital domain for reproducing HDR images. [15] compresses the dynamic range through the manipulation of the gradient domain in the logarithmic space. More recently, [16] proposes a novel method by adjusting the local histogram. There has been little published research to explore rendering HDR images on the GPU [17,18] and these works

Fig. 1. Selected multi-exposed image set of the same scene.

Maximally stable extremal regions

- Extremal region: any connected region in an image with all pixel values above (or below) a threshold

Region that differ in properties, such as brightness or color, compared to surrounding regions. Almost uniform intensity

- Observations:
 - Nested extremal regions result when the threshold is successively raised or lowered.
 - The nested extremal regions form a “component tree.”
- Key idea: choose thresholds θ such that the resulting bright (or dark) extremal regions are nearly constant when these thresholds are perturbed by $+/-\Delta$

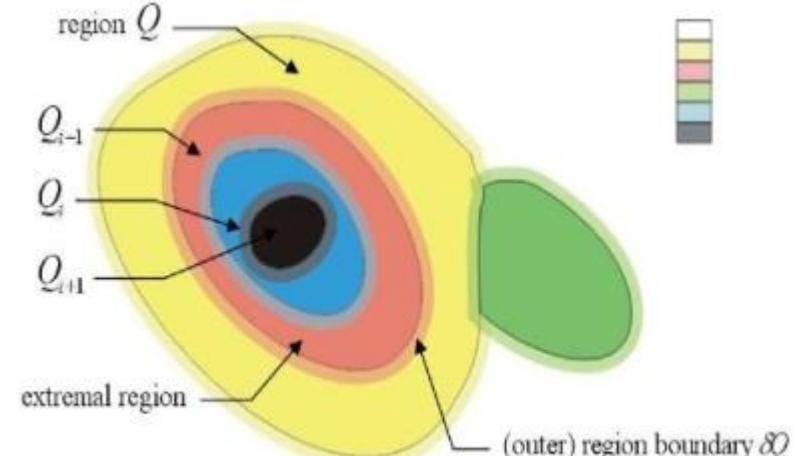
“maximally stable” extremal regions (MSER)

stable connected component of some gray-level sets of the image .

[Matas, Chum, Urba, Pajdla, 2002]

MSER Processing

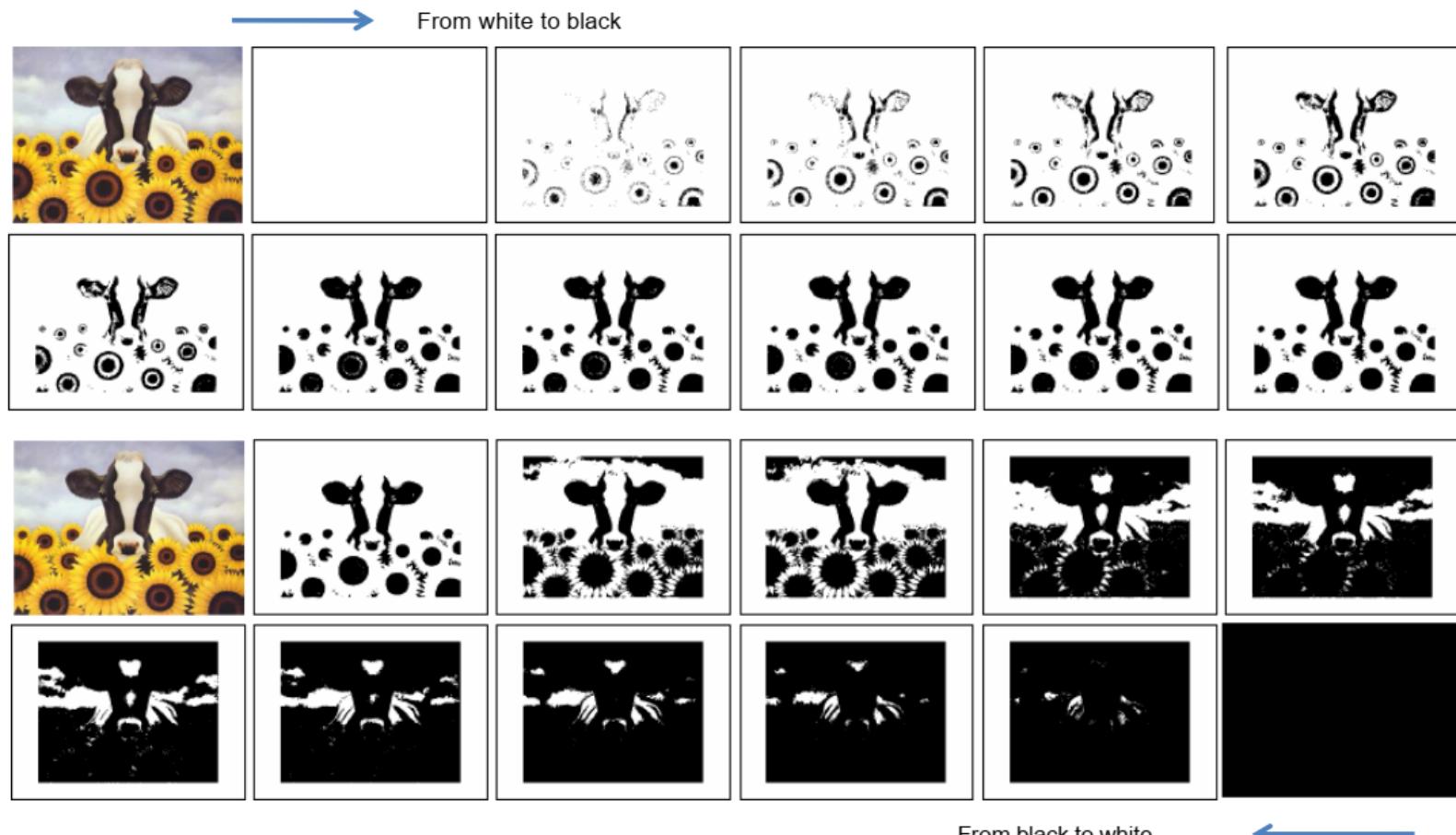
- Step by step
 - Sweep threshold: from black to white
 - Extract connected component (extremal regions)
 - Find a threshold make extremal region “maximally stable”
 - Approximate a region with an ellipse (optional)
 - Keep those regions descriptors as features
- Extremal region can be rejected if
 - It is too big/mall
 - Too unstable
 - Too similar to its parent MSER



$$Q_{i^*} : i^* = \arg \min_i |Q_{i+\Delta} \setminus Q_{i-\Delta}| / |Q_i|$$

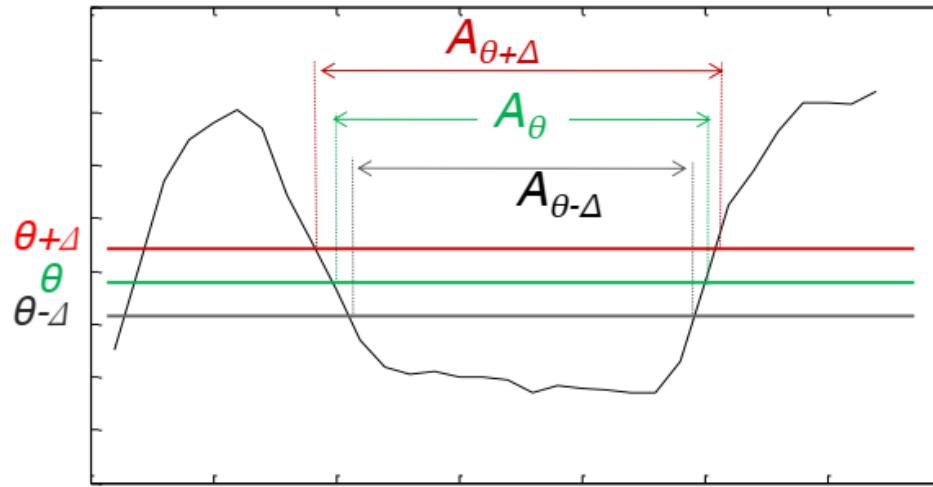
Sweeping Image threshold

- Apply a series of thresholds – one for each grayscale level.
- Threshold the image at each level to create a series of black and white images.
- One extreme will be all white, the other all black. In between, blobs grow and merge



MSER Processing

- Control how stable it is the region via Δ



$$\text{Local minimum of } \left| \frac{A_{\theta-\Delta} - A_{\theta+\Delta}}{A_\theta} \right| \rightarrow \text{MSER}$$

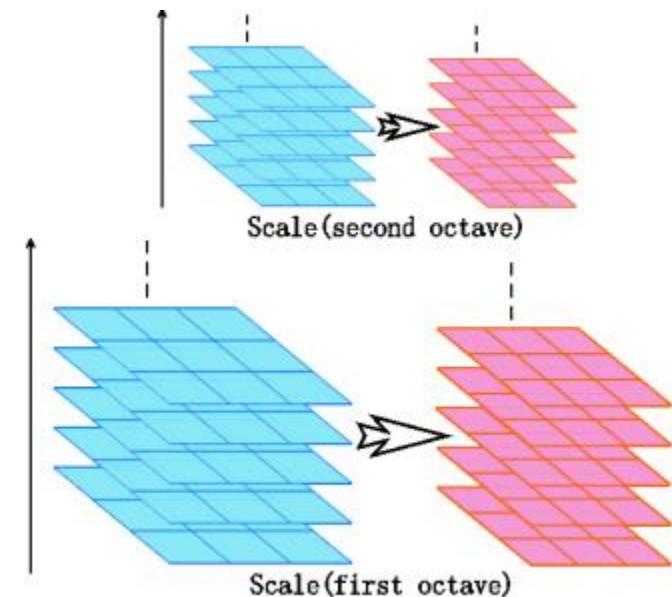
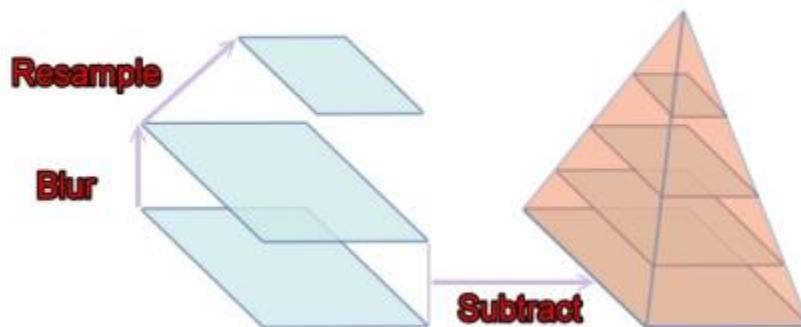
[Matas, Chum, Urba, Pajdla, 2002]

A_θ the area of the extremal region

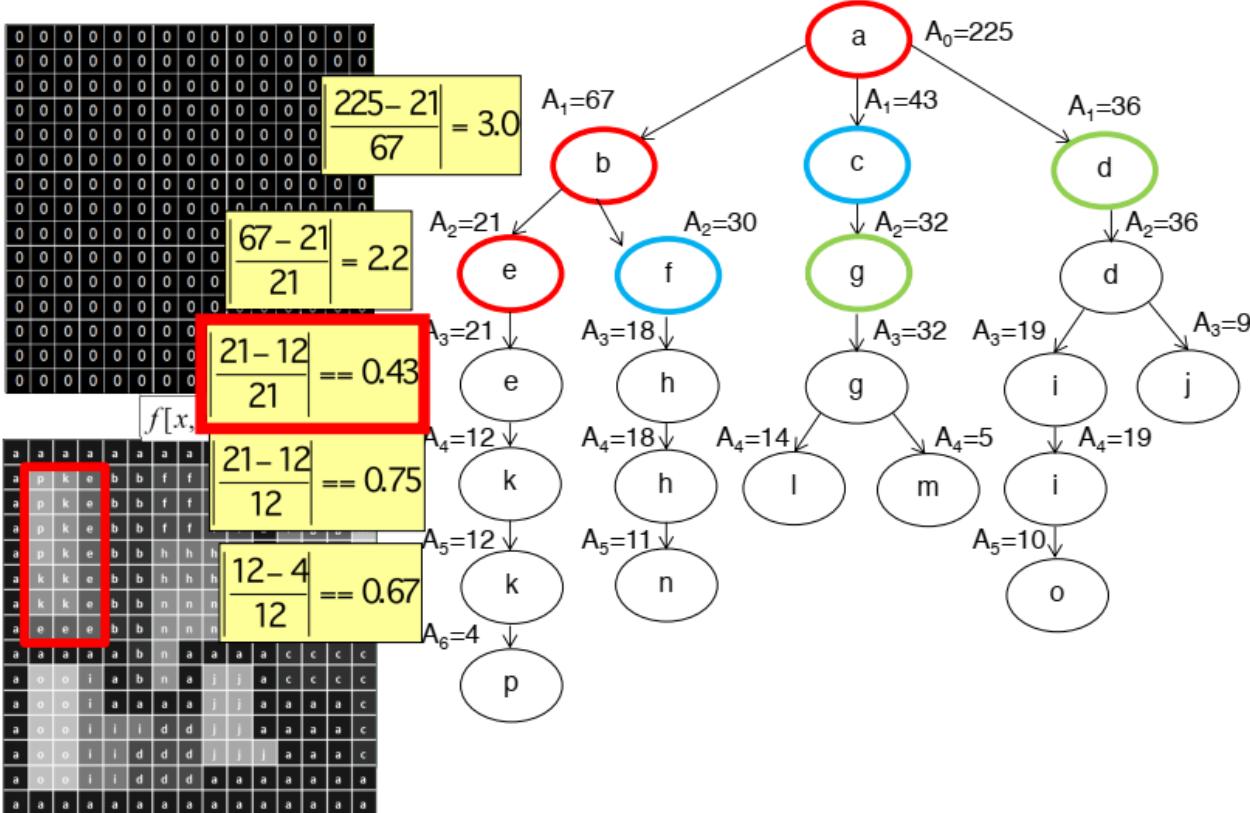
$A_{\theta-\Delta}, A_{\theta+\Delta}$ is the difference in the area of two regions with $\pm\Delta$ level

Multi-resolution MSER

- Step1 : Construct a scale pyramid with one octave between scales
- Step2 : Detect MSERs separately at each resolution
- Step3 : Duplicate MSERs are removed by eliminating fine scale MSERs with similar locations and sizes as MSERs detected at the next coarser scale



Component tree of an image



Local minima of sequence

$$\left| \frac{A_{\theta-\Delta} - A_{\theta+\Delta}}{A_\theta} \right|$$

$\theta = \Delta, \Delta + 1, \dots \rightarrow \text{MSERs}$

MSER Example



Dark MSERs, $\Delta=15$



Original image



Bright MSERs, $\Delta=15$



Dark MSERs, $\Delta=15$



Original image



Bright MSERs, $\Delta=15$

Level sets of an image

1	1	1	1	1	1	1	1	1	1	5	4	4	8	
1	7	6	4	2	2	3	3	3	3	1	5	4	4	8
1	7	6	4	2	2	3	3	3	3	1	5	4	4	8
1	7	6	4	2	2	3	3	3	3	1	5	4	4	8
1	7	6	4	2	2	5	5	5	5	1	5	4	4	8
1	6	6	4	2	2	5	5	5	6	1	5	4	4	4
1	6	6	4	2	2	6	6	6	6	1	5	5	5	5
1	4	4	4	2	2	6	6	6	6	1	5	5	5	5
1	1	1	1	1	2	6	1	1	1	1	2	2	2	2
1	8	8	5	1	2	6	1	7	7	1	2	2	2	2
1	8	8	5	1	1	1	1	7	7	1	1	1	1	2
1	8	8	5	5	5	3	3	7	7	1	1	1	1	2
1	8	8	5	5	3	3	3	7	7	7	1	1	1	2
1	8	8	5	5	3	3	3	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

$f[x, y]$

Image

0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0

$f[x, y] > 8$

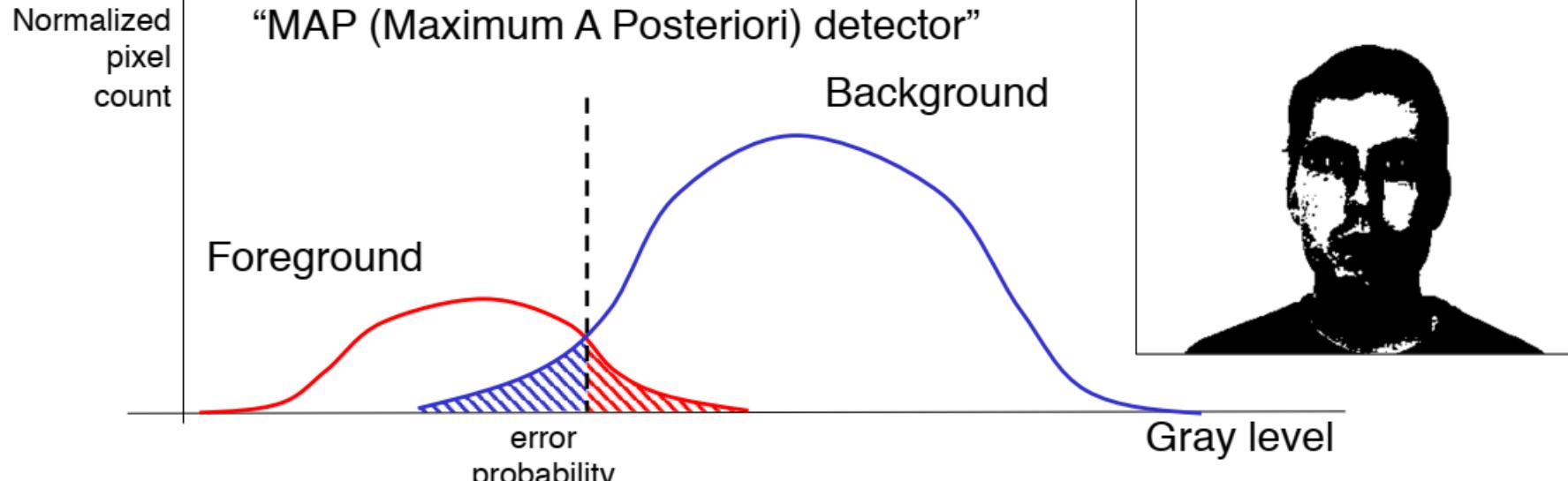
Level Set

MSER Properties

- MSER perform well on
 - images containing homogeneous regions with distinctive boundaries.
 - small regions
- MSER doesn't work well with images with any motion blur
- Multi-resolution MSER provides better
 - Robustness to large scale changes and blurred images
 - Improves matching performance over large scale changes & for blurred images
- Good repeatability
- Affine invariant
- A smart implementation makes it one of the fastest region detectors

<http://www.vlfeat.org/overview/mser.html>

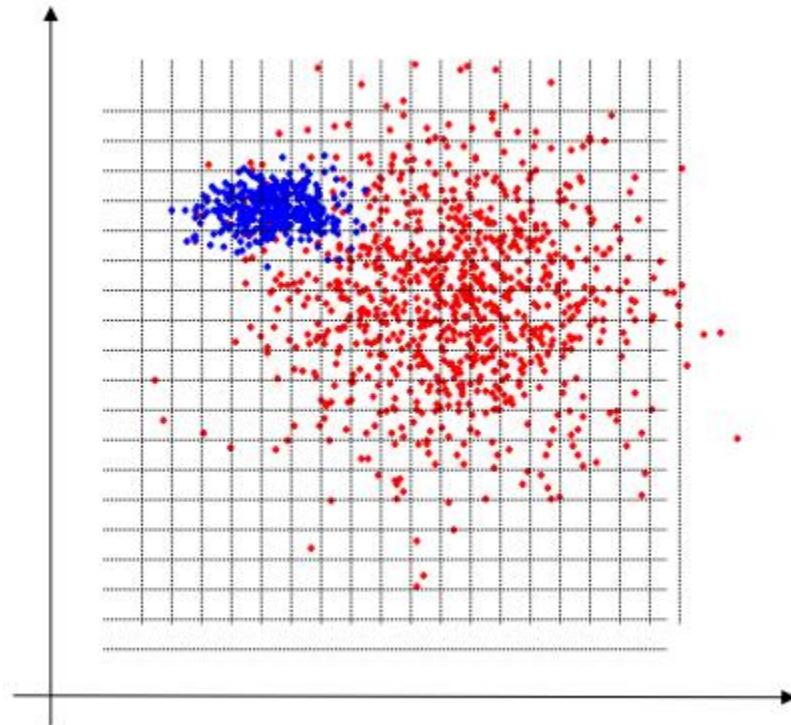
Supervised Thresholding



If errors $BG \rightarrow FG$ and $FG \rightarrow BG$ are associated with different costs:
“Bayes minimum risk detector” is optimal.

Multidimensional MAP detector

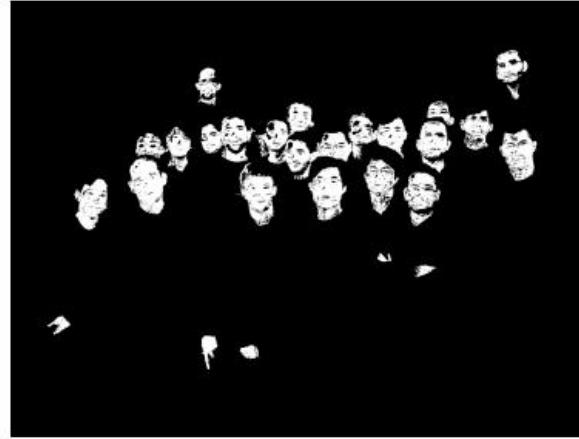
- Training
 - Provide labelled set of training data
 - Subdivide n-dimensional space into small bins
 - Count frequency of occurrence for each bin and class in training set, label bin with most probable class
 - (Propagate class labels to empty bins)
- For test data: identify bin, look up the most probable class



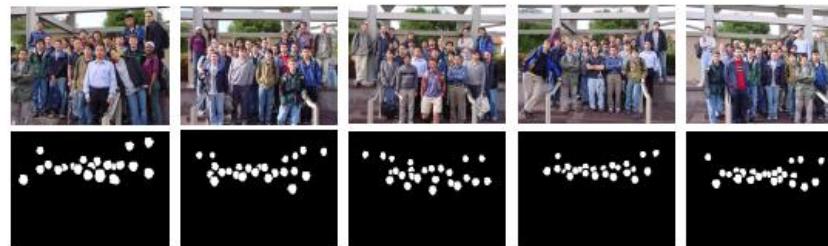
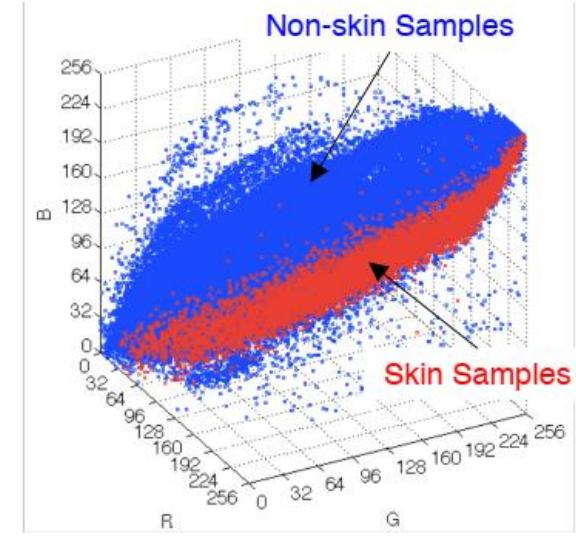
MAP detector in RGB-space



Original image



Skin color detector



Five training images

