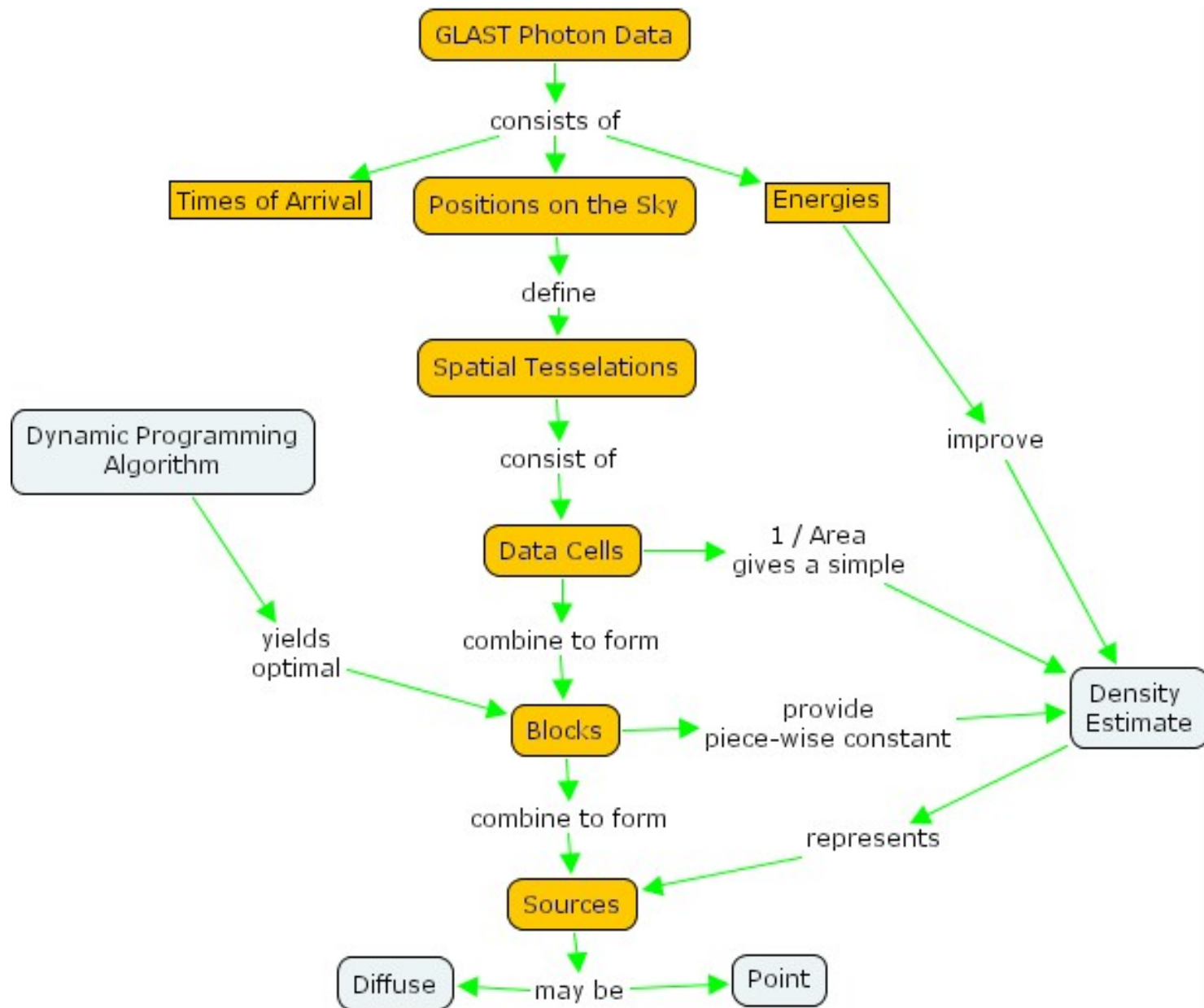


# Bootstrap Segmentation Analysis and Expectation Maximization to Detect and Characterize Sources

Jeffrey.D.Scargle@nasa.gov

Space Science Division  
NASA Ames Research Center

Diffuse Emission & LAT Source Catalog  
SLAC Workshop: May 23, 2005



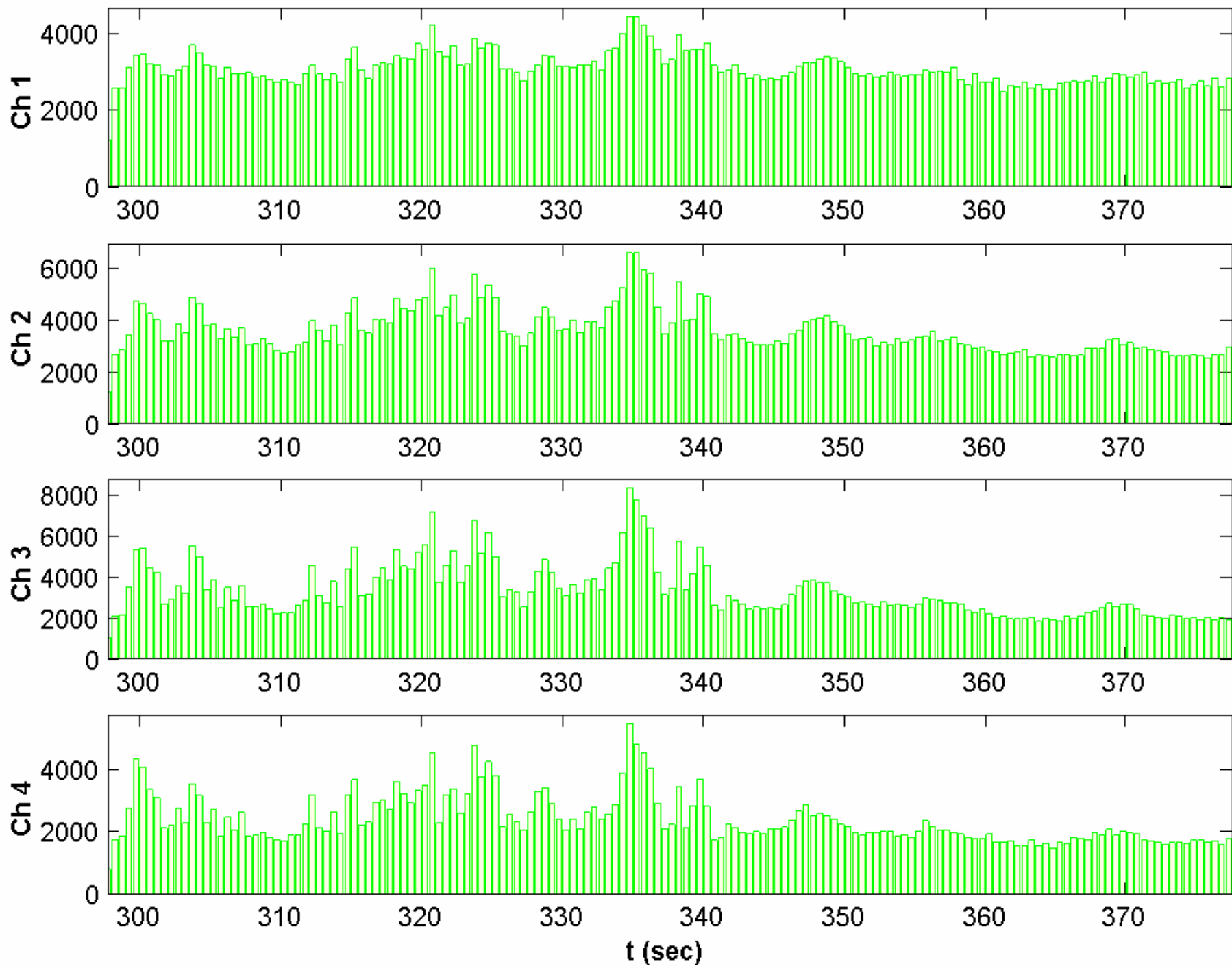
## Problems (Solutions)

- **Algorithm provides no Error Estimate**
  - **Bootstrap Errors**
  - **Block Posterior Probabilities**
- **Point Spread: Overlapping Sources**
  - **Expectation Maximization (EM)**
- **Point Spread: Function of Energy**
  - **Maximum Likelihood Models?**

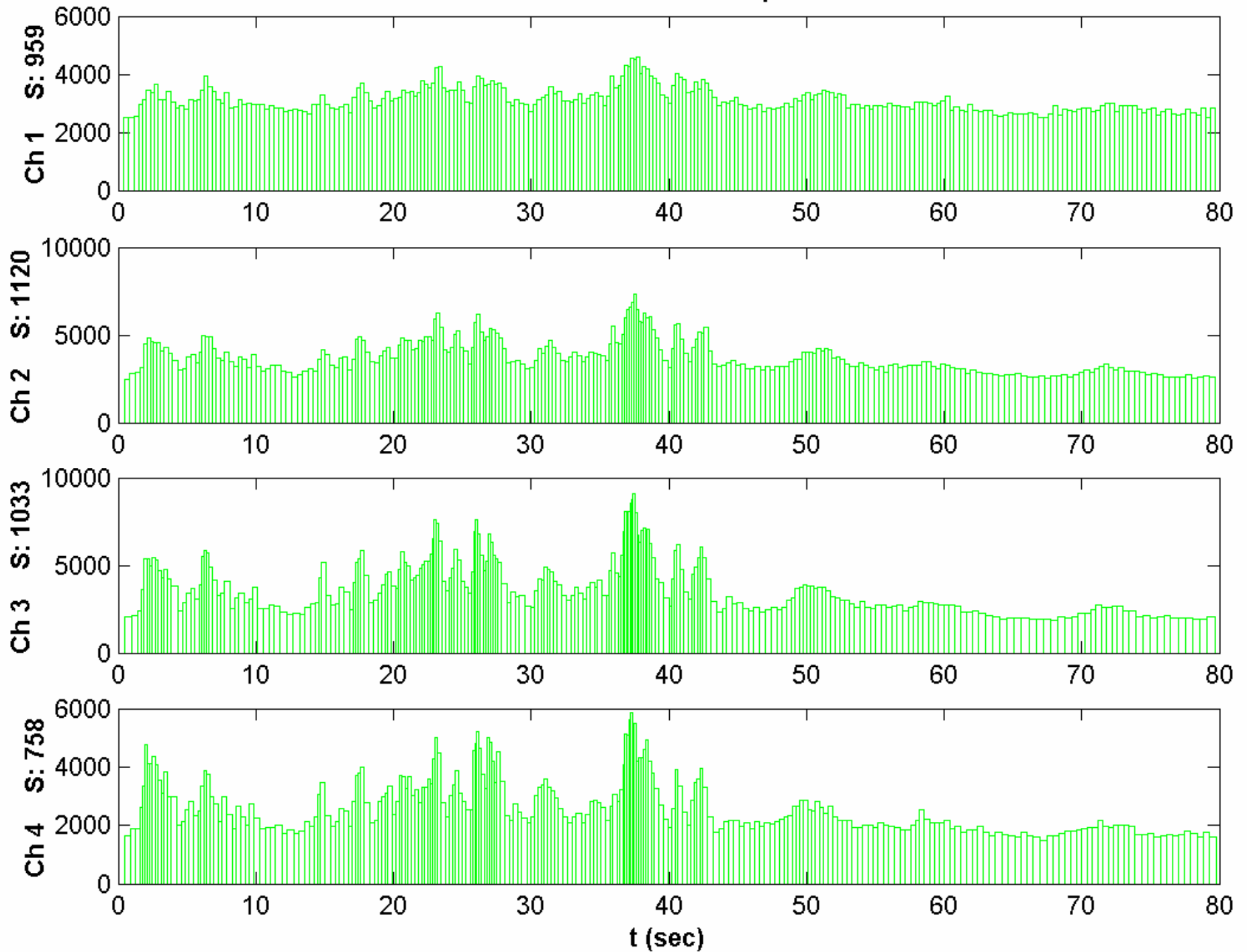
# One Dimensional<sup>1</sup> Example: Swift GRB Data

1. But everything applies to 2D and higher!

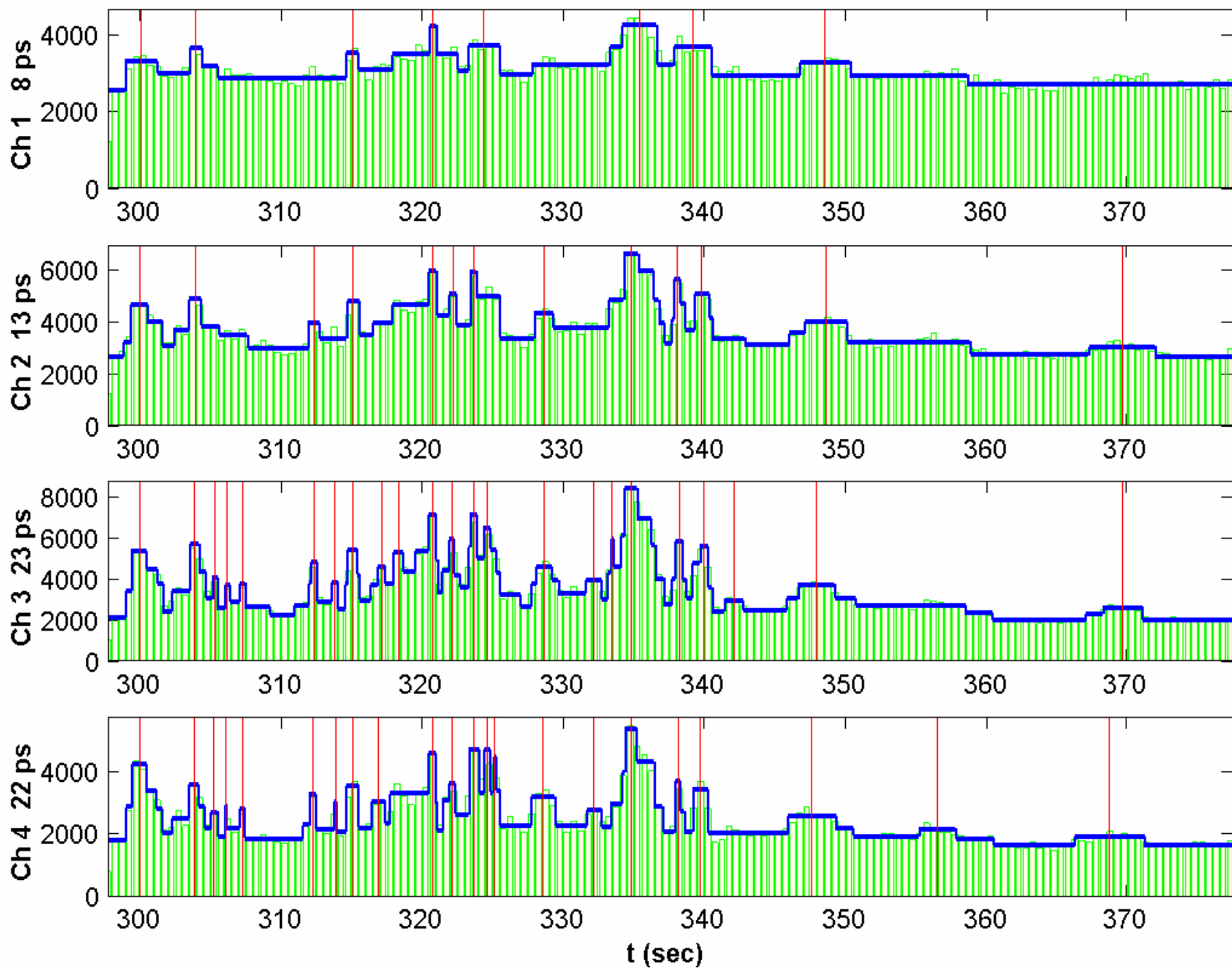
Swift Burst - 12/23/2004 Bin size: 0.5 (sec)



### Swift Burst - 12/23/2004 256 percentile bins



### Swift Burst - 12/23/2004 TTS factor 32



# Expectation Maximization (EM)

Initialize: Find a good guess at the mixture model:

- How Many Sources?
- Locations of the Sources
- Source Parameters (size, spectra, ...)

Iterate:

1. From the model: Divide the data into pieces that are relevant to each Source separately
2. Re-determine the Source Locations and Parameters by fitting to the data pieces in 1
3. Repeat as needed



# The “Maximization” Step with Point Data

Maximize the Unbinned log-Likelihood:

$$\text{Log(L)} = \sum_i \log[ \hat{x}_a(t_i) + b ] - \int w(t) [ \hat{x}_a(t) + b ]$$

Where the model for each source (pulse) is:

$$X(t) = [ \hat{x}_a(t) + b ]$$

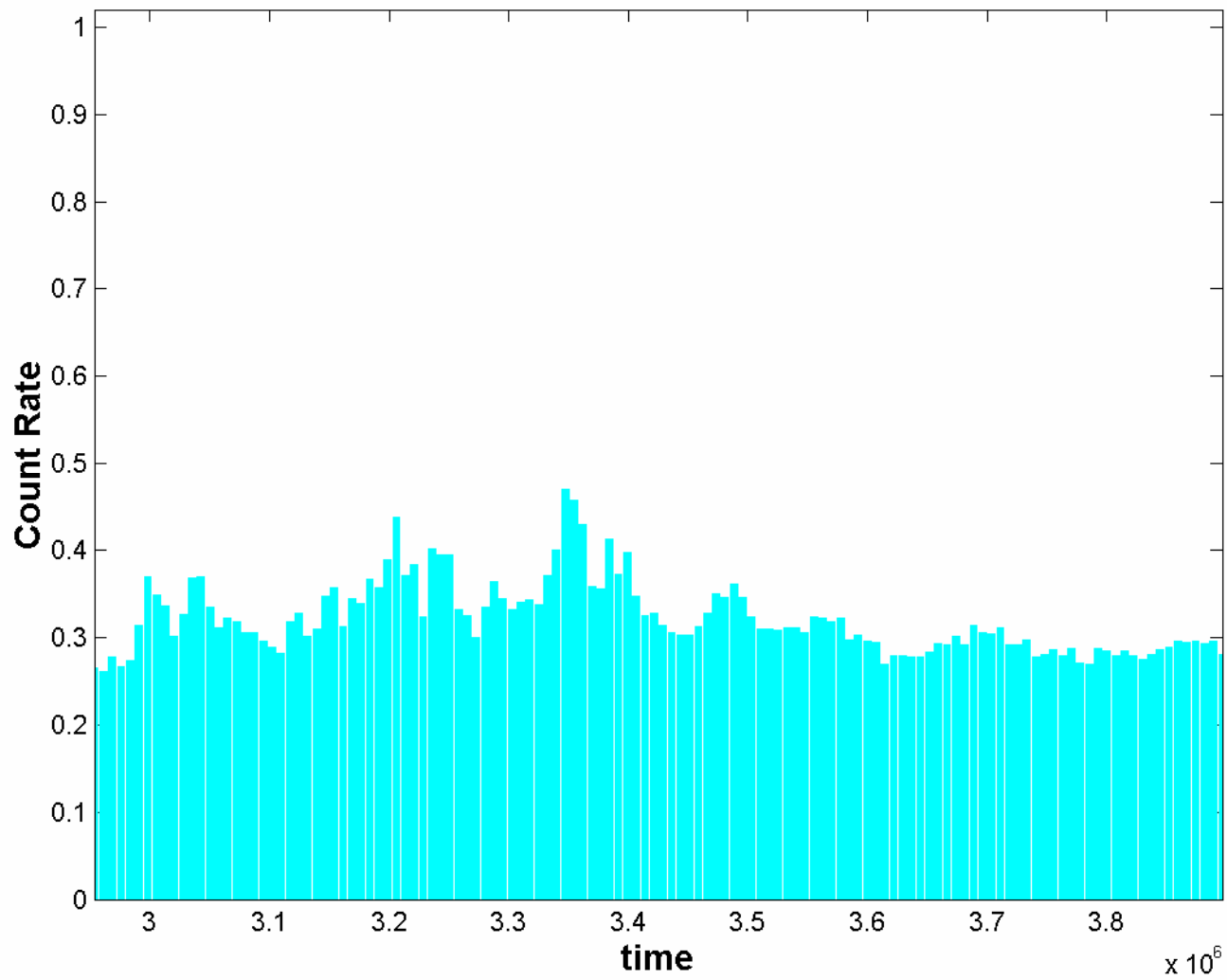
a = source parameters (location, size, ... )

b = local background constant

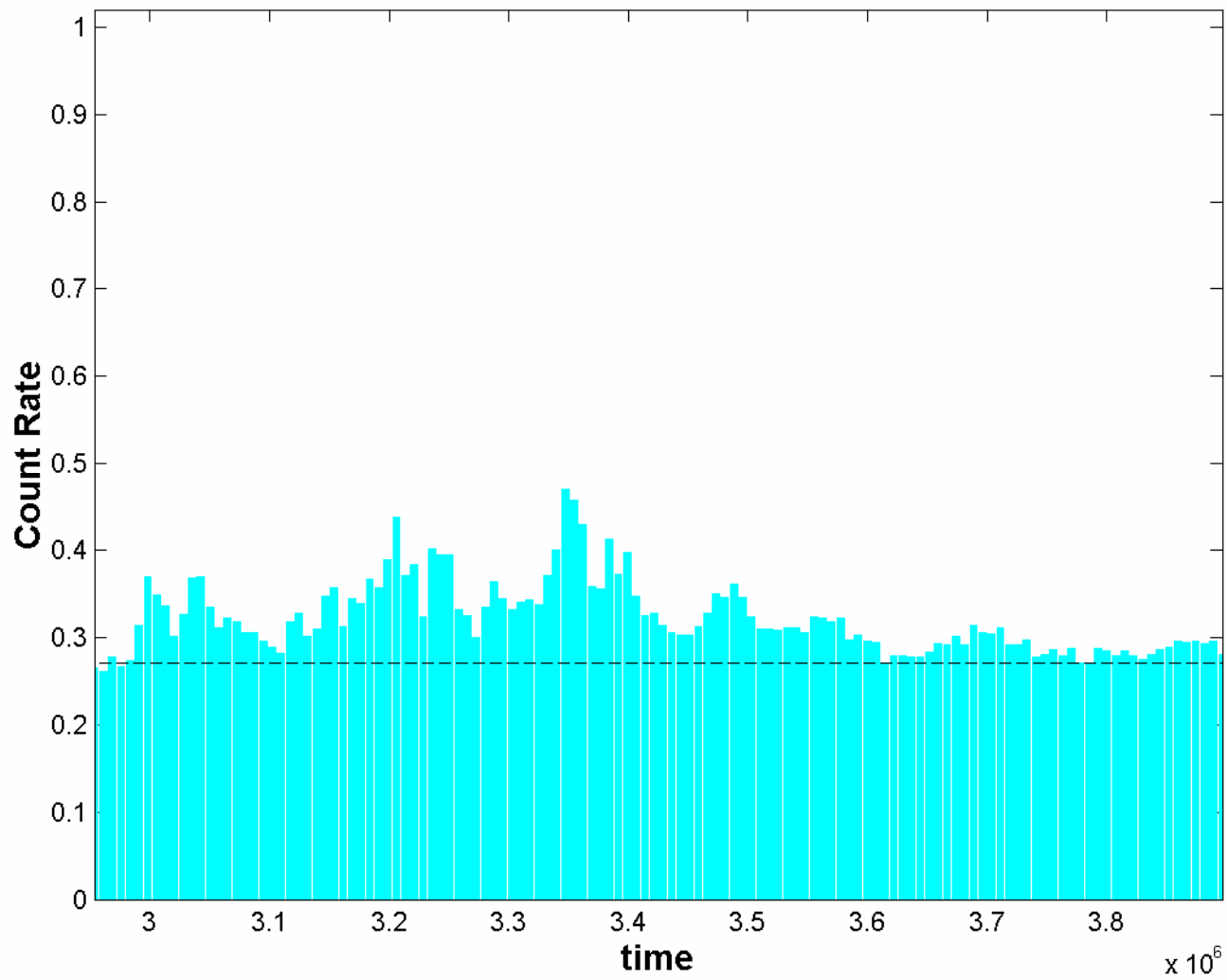
w(t) is the EM partitioning, or weighting, function

NB: the term  $\sum_i \log[ w(t_i) ]$  in the full log-likelihood is a constant, irrelevant for model fitting.

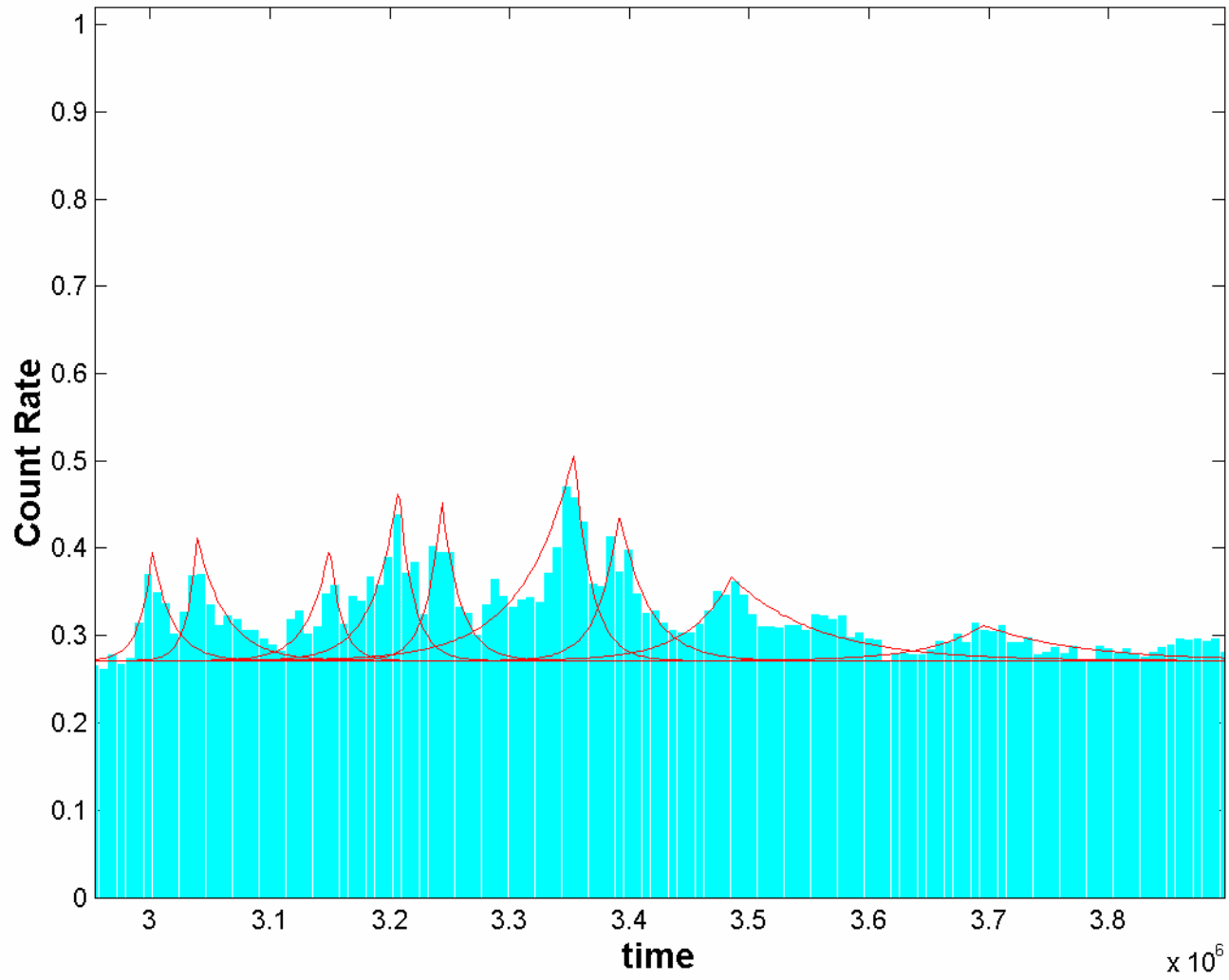
### Arbitrary Bins



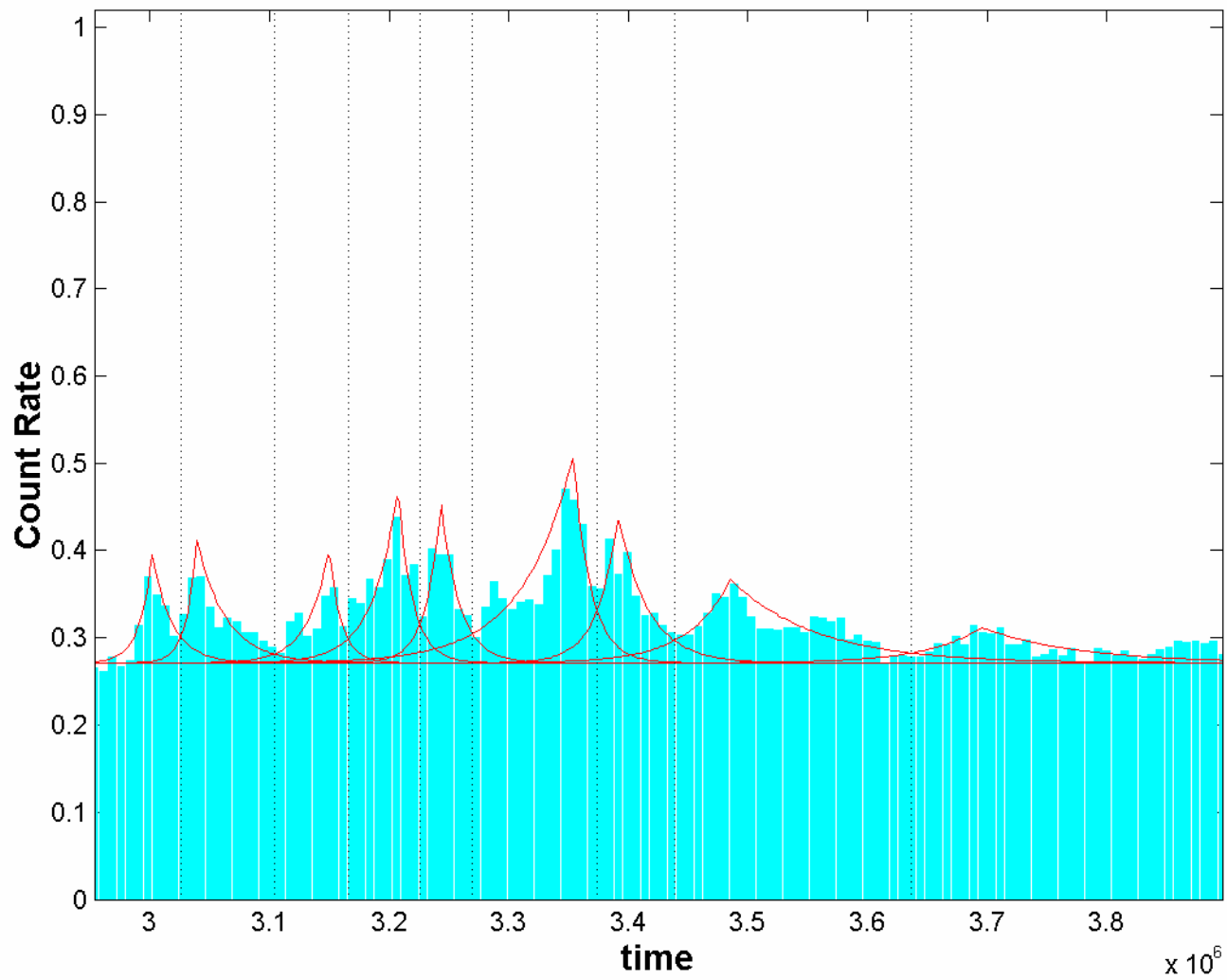
### Background Level



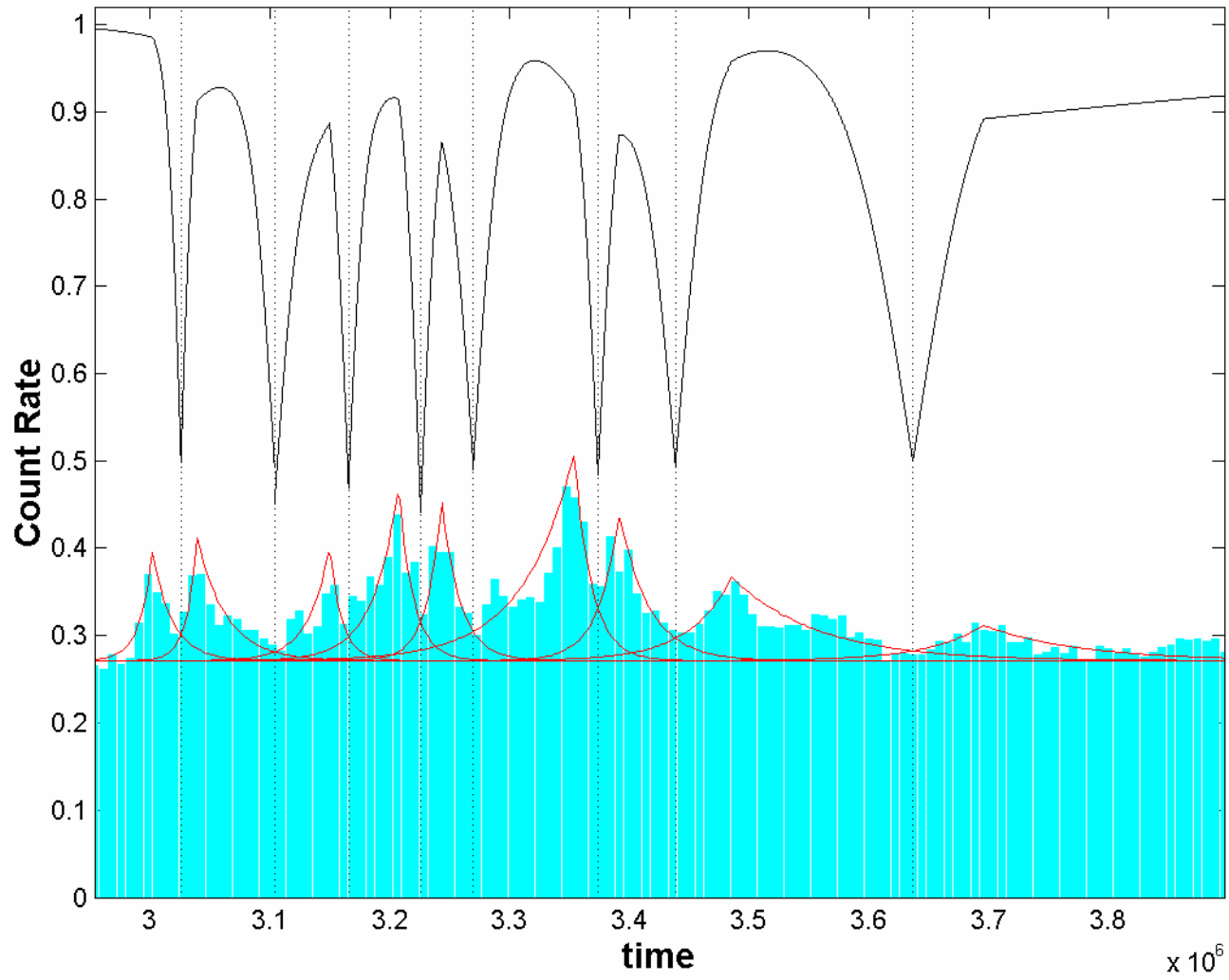
Initial Model from Optimal Segmentation

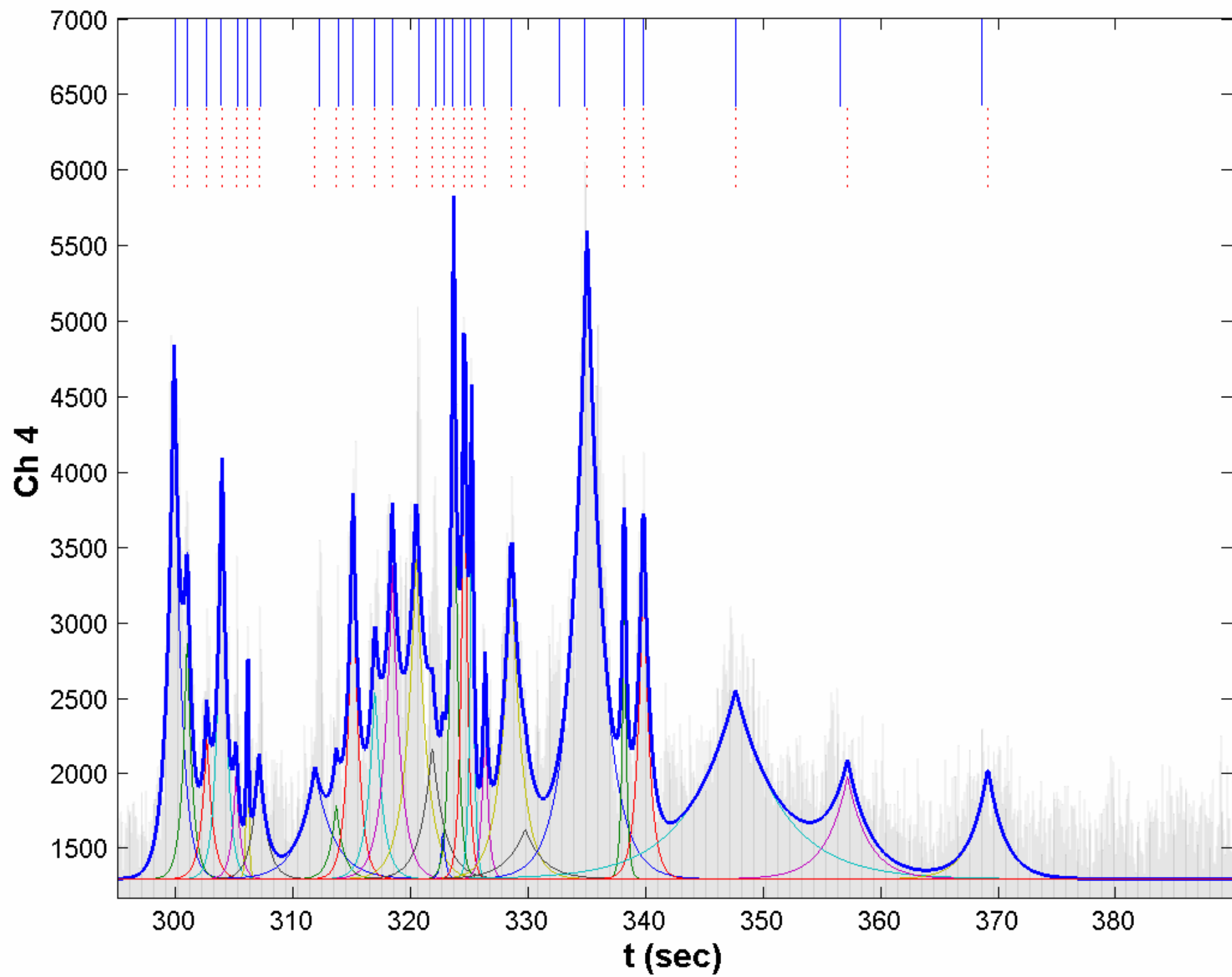


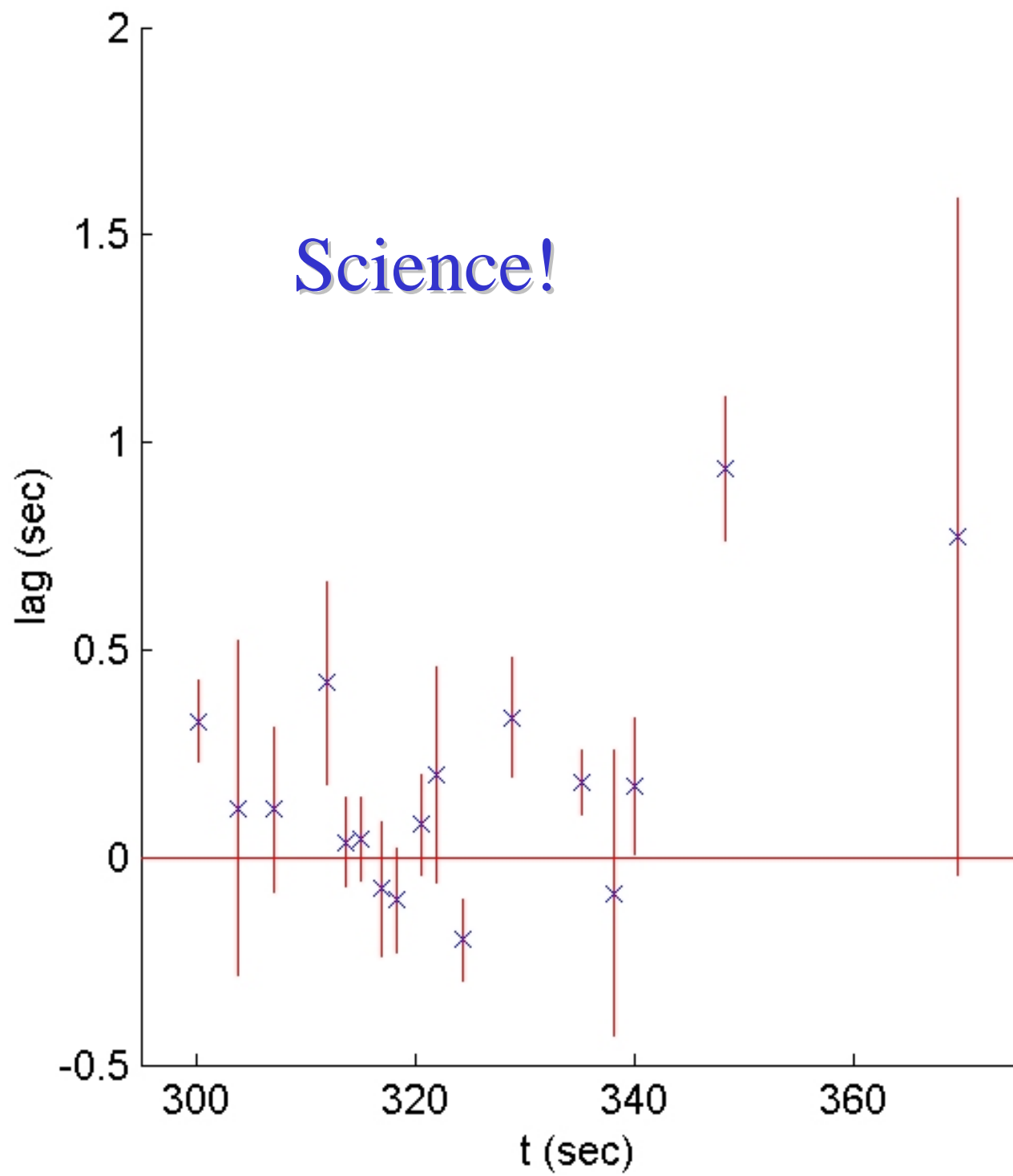
### Partition Interval



### Localized Weighting Functions









## Bootstrap Method: Time Series of N Discrete Events

For many iterations:

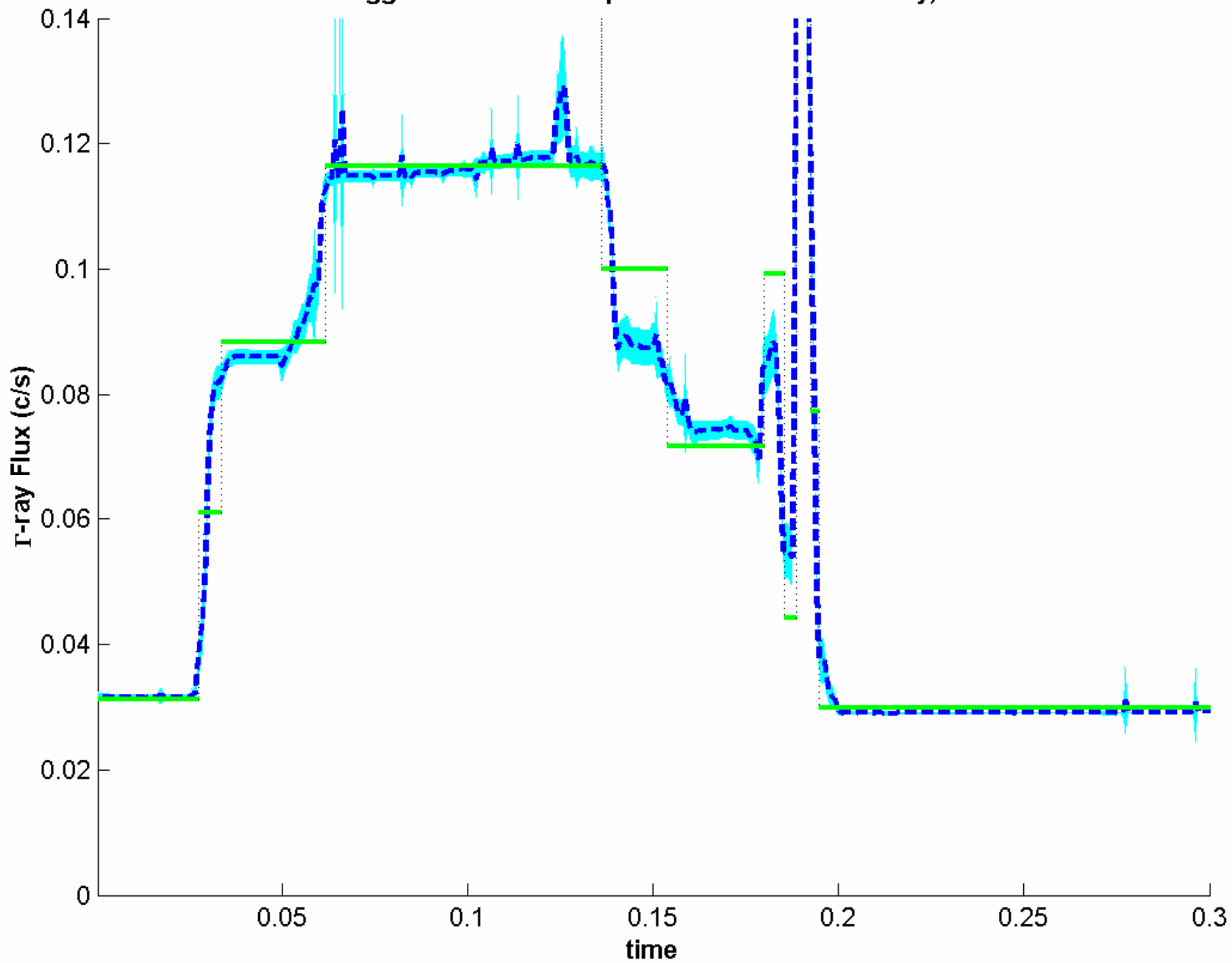
- Randomly select N of the observed events *with replacement*
- Analyze this sample just as if it were real data

Compute mean and variance of the bootstrap samples

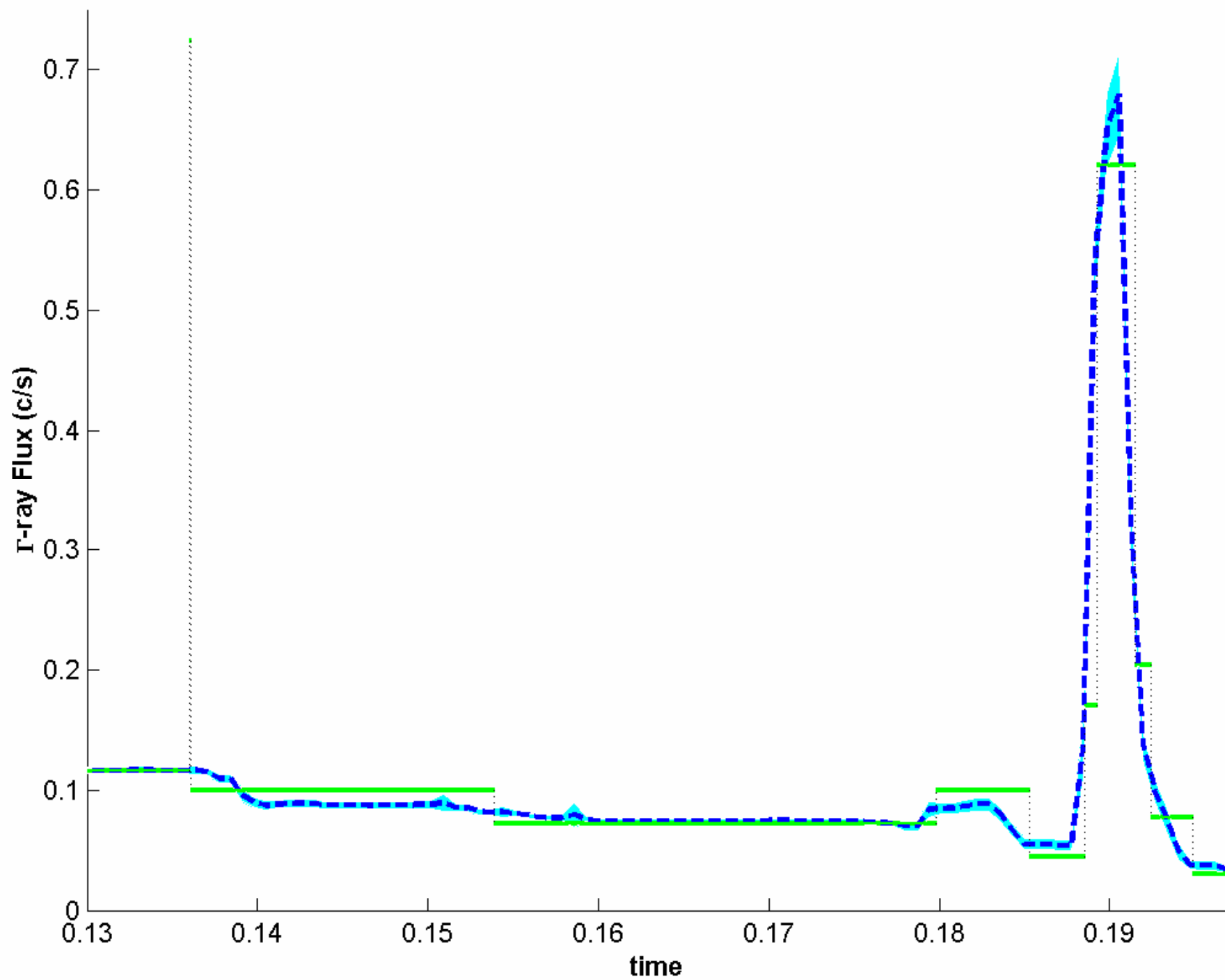
Bias = result for real data – bootstrap mean  
RMS error derived from bootstrap variance

Caveat: The real data does not have the repeated events in bootstrap samples. I am not sure what effect this has.

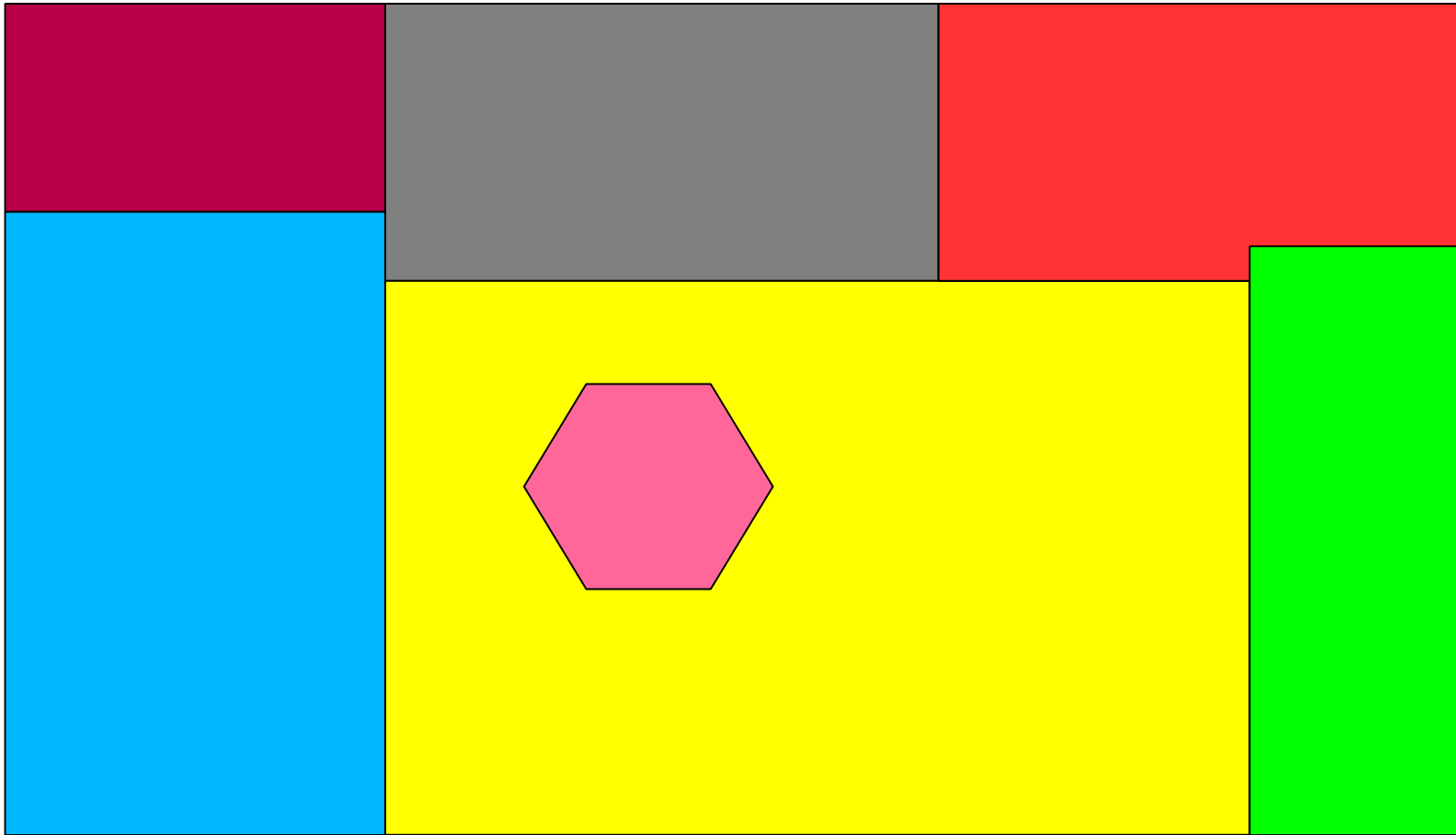
BATSE Trigger 1453: Bootstrap mean and  $5\sigma$  Uncertainty, ML Blocks



BATSE Trigger 1453: Bootstrap mean and  $5\sigma$  Uncertainty, ML Blocks



# Piecewise Constant Model (partitions the data space)



Signal modeled as constant over each partition element (block).

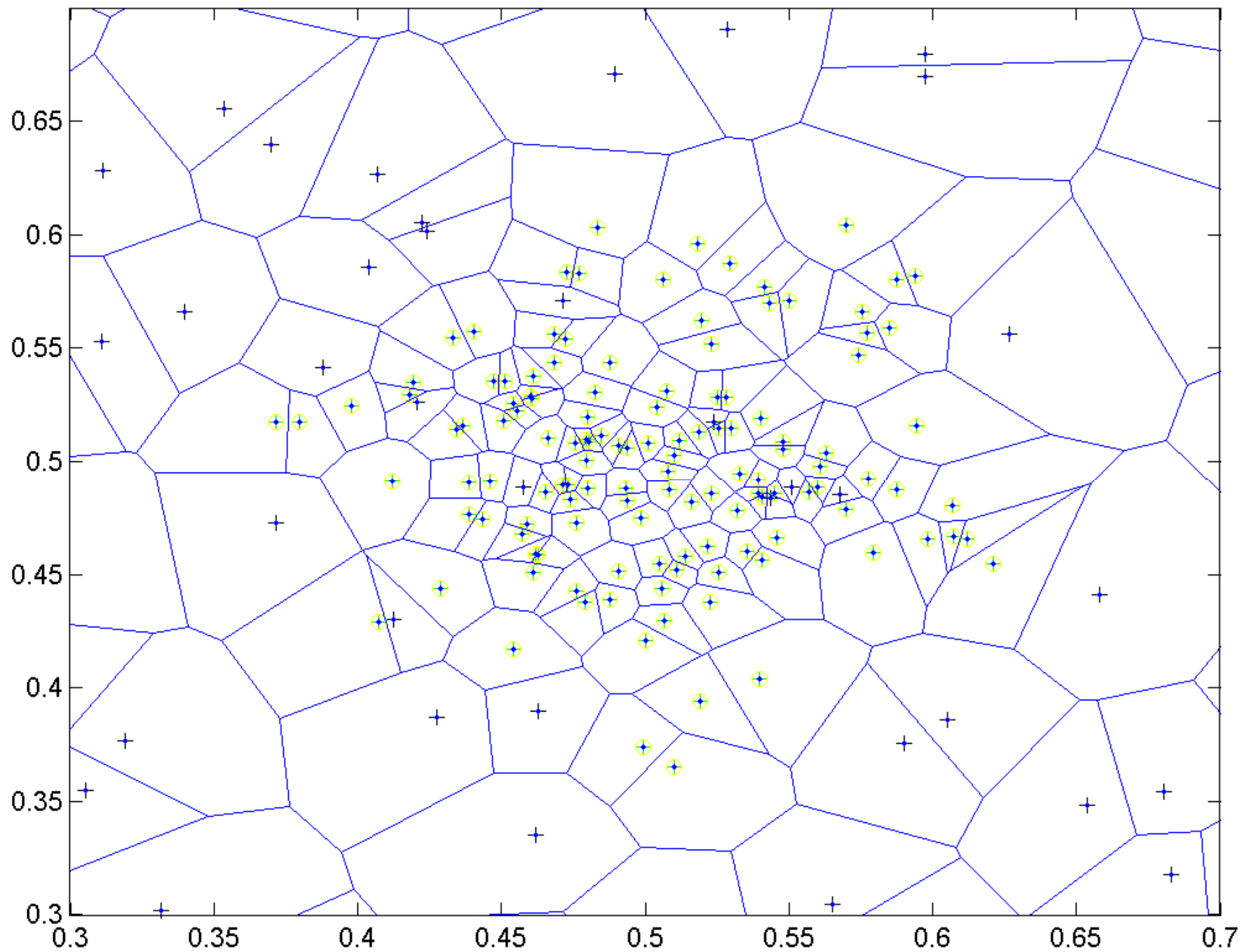
# Optimum Partitions in Higher Dimensions

- Blocks are collections of Voronoi cells (1D,2D,...)
- Relax condition that blocks be connected
- Cell location now irrelevant
- Order cells by volume

Theorem: Optimum partition consists of blocks  
that are connected in this ordering

- Now can use the 1D algorithm,  $O(N^2)$
- Postprocessing: identify connected block fragments

**Data: Voronoi Tessellation**



# Blocks

Block: a set of data cells

Two cases:

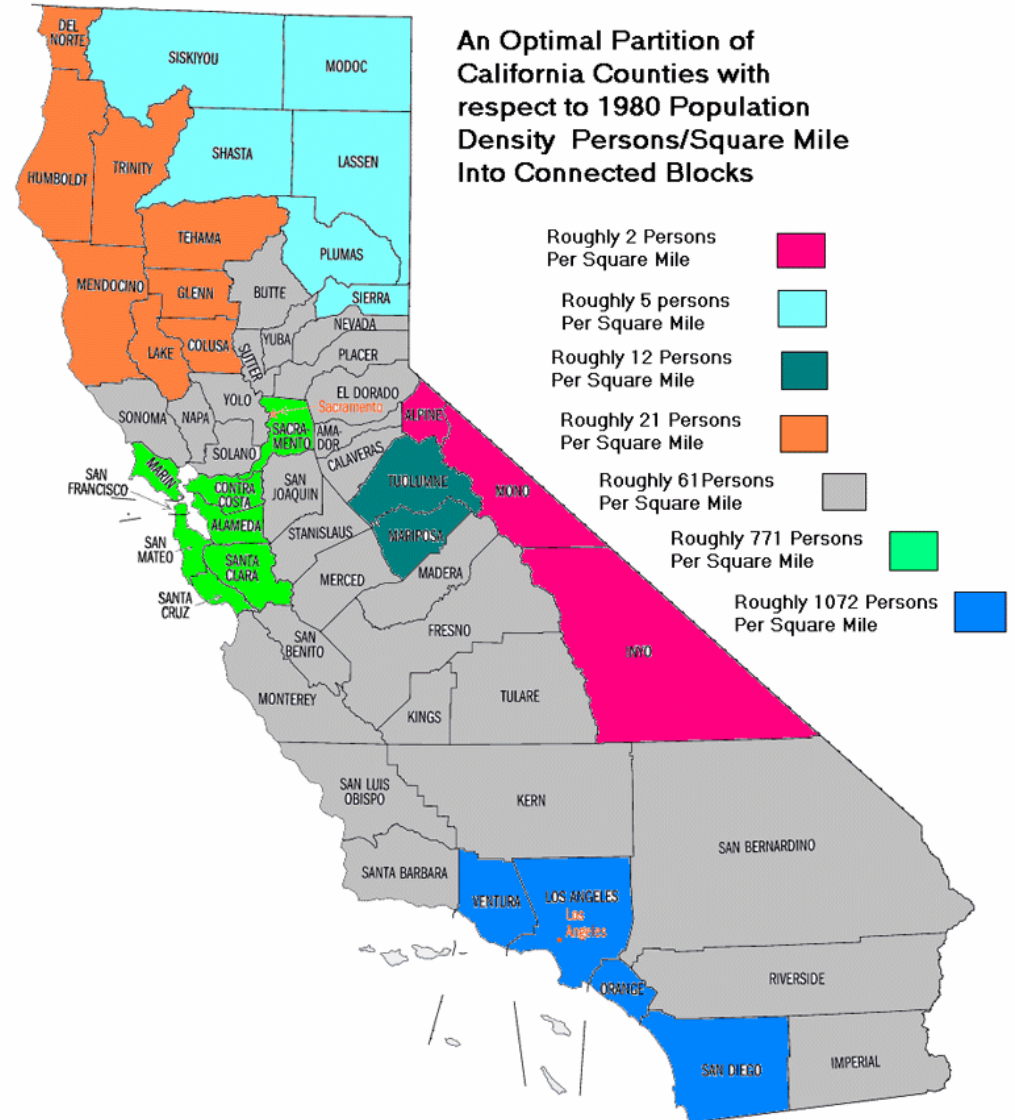
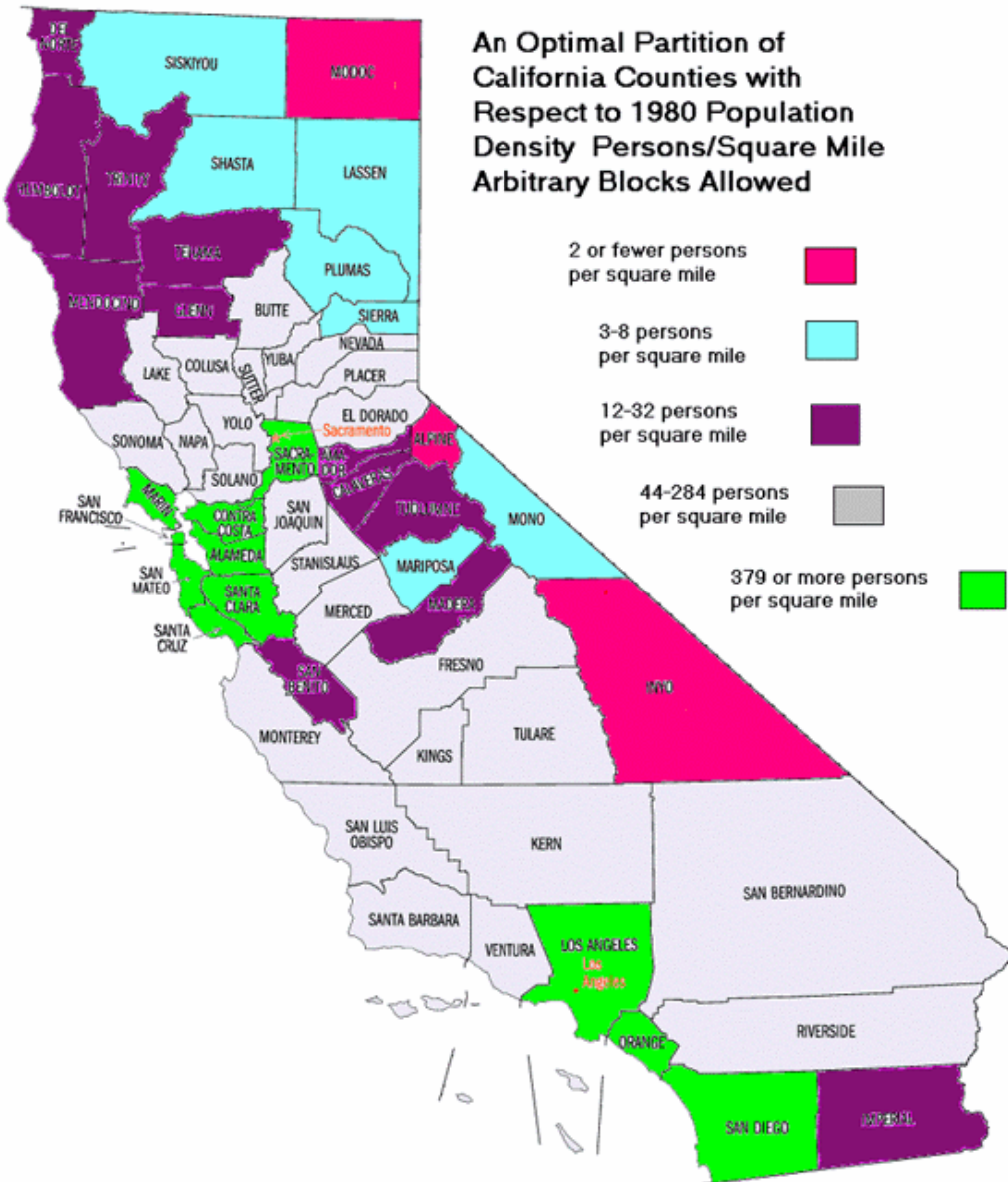
- Connected (can't break into distinct parts)
- Not constrained to be connected

Model = set of blocks

Fitness function:

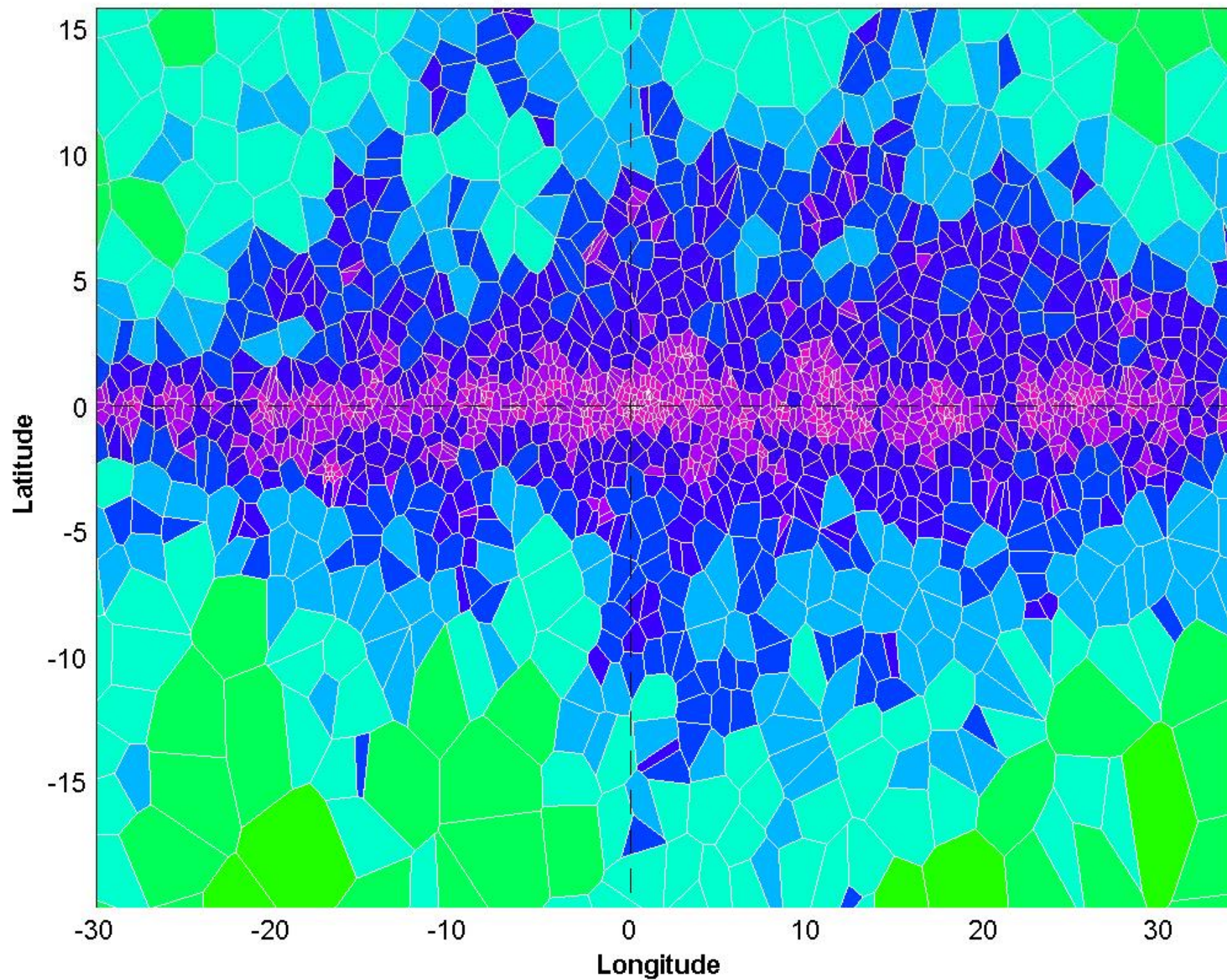
$F(\text{Model}) = \text{sum over blocks } F(\text{Block})$

# Connected vs. Arbitrary Blocks

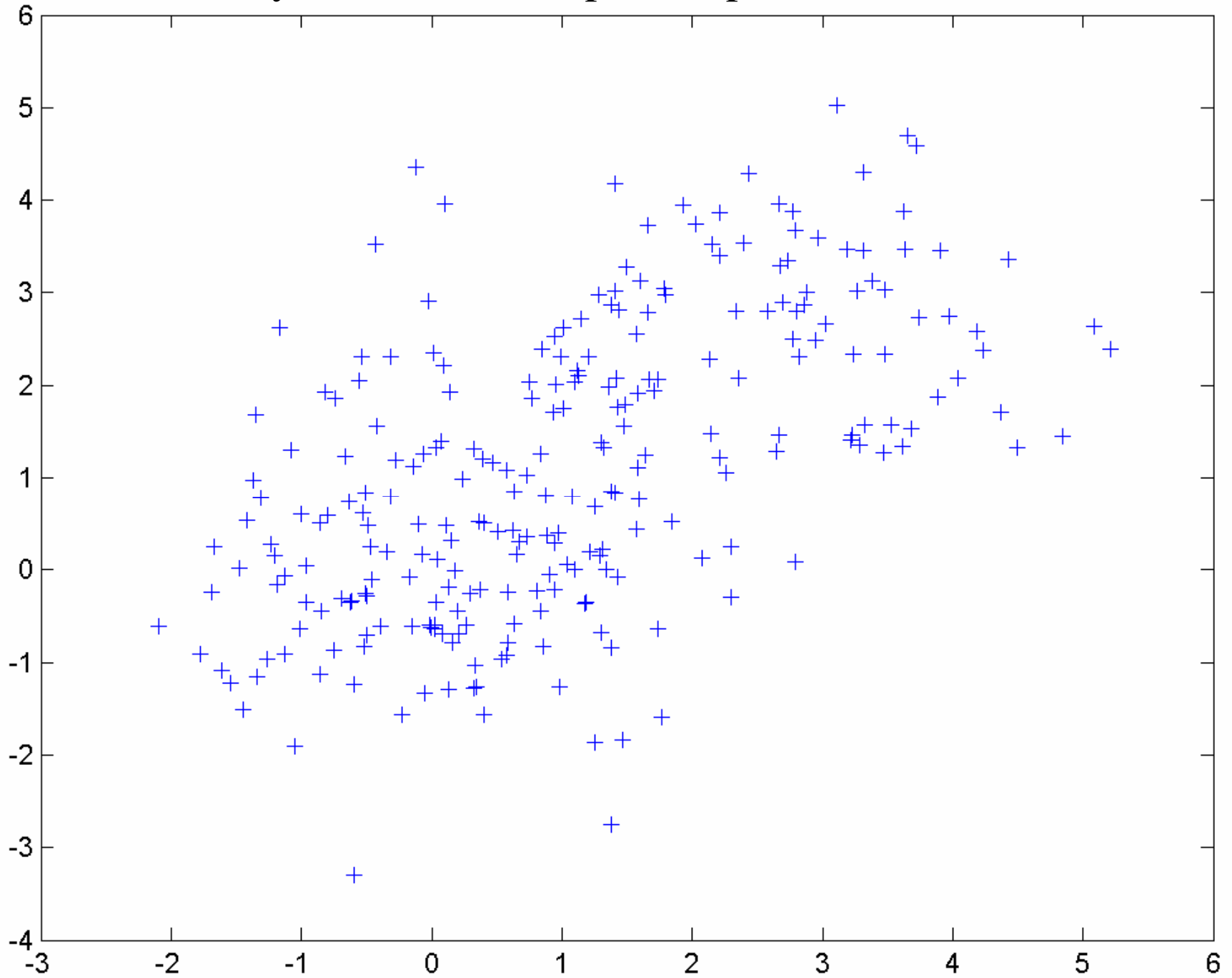


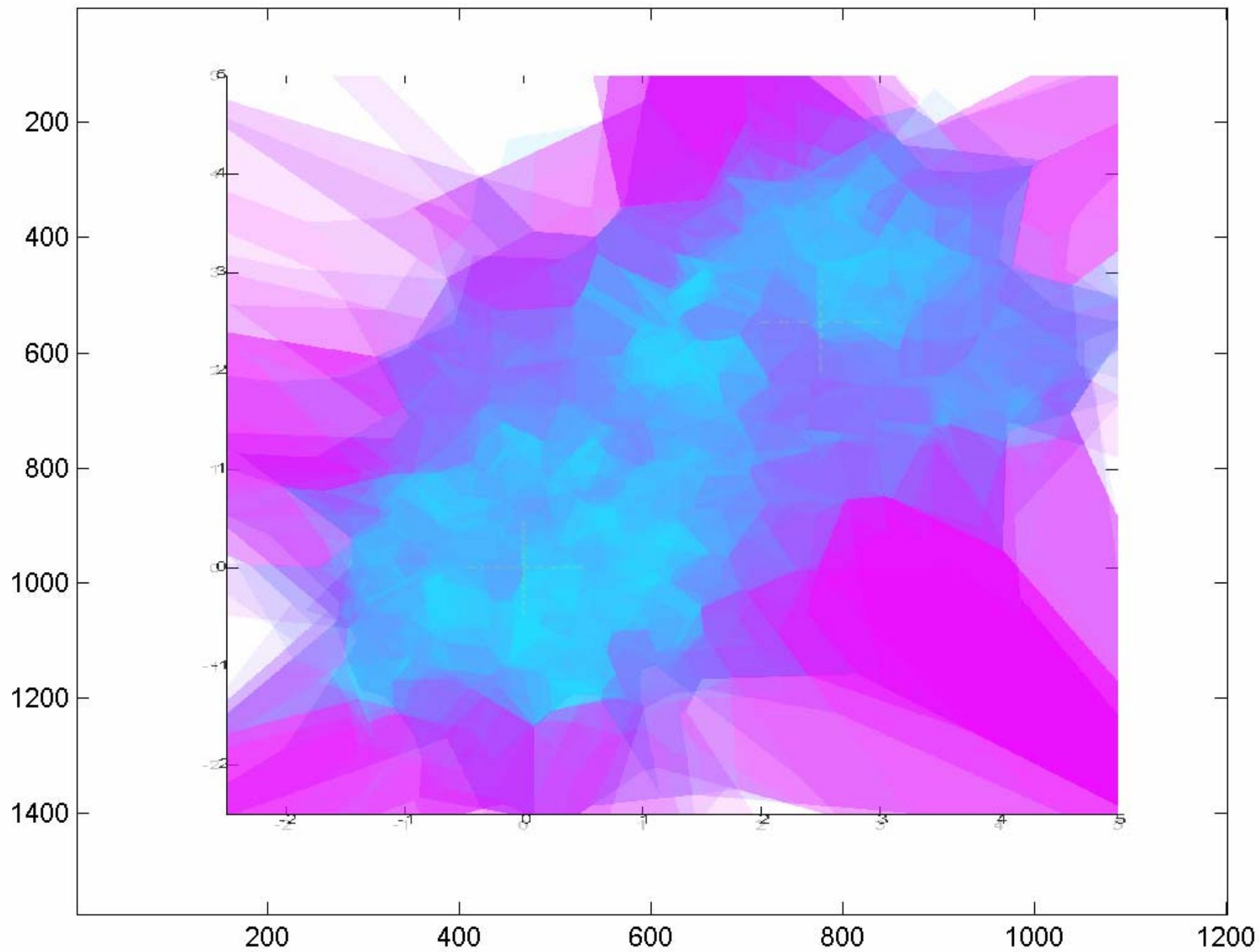


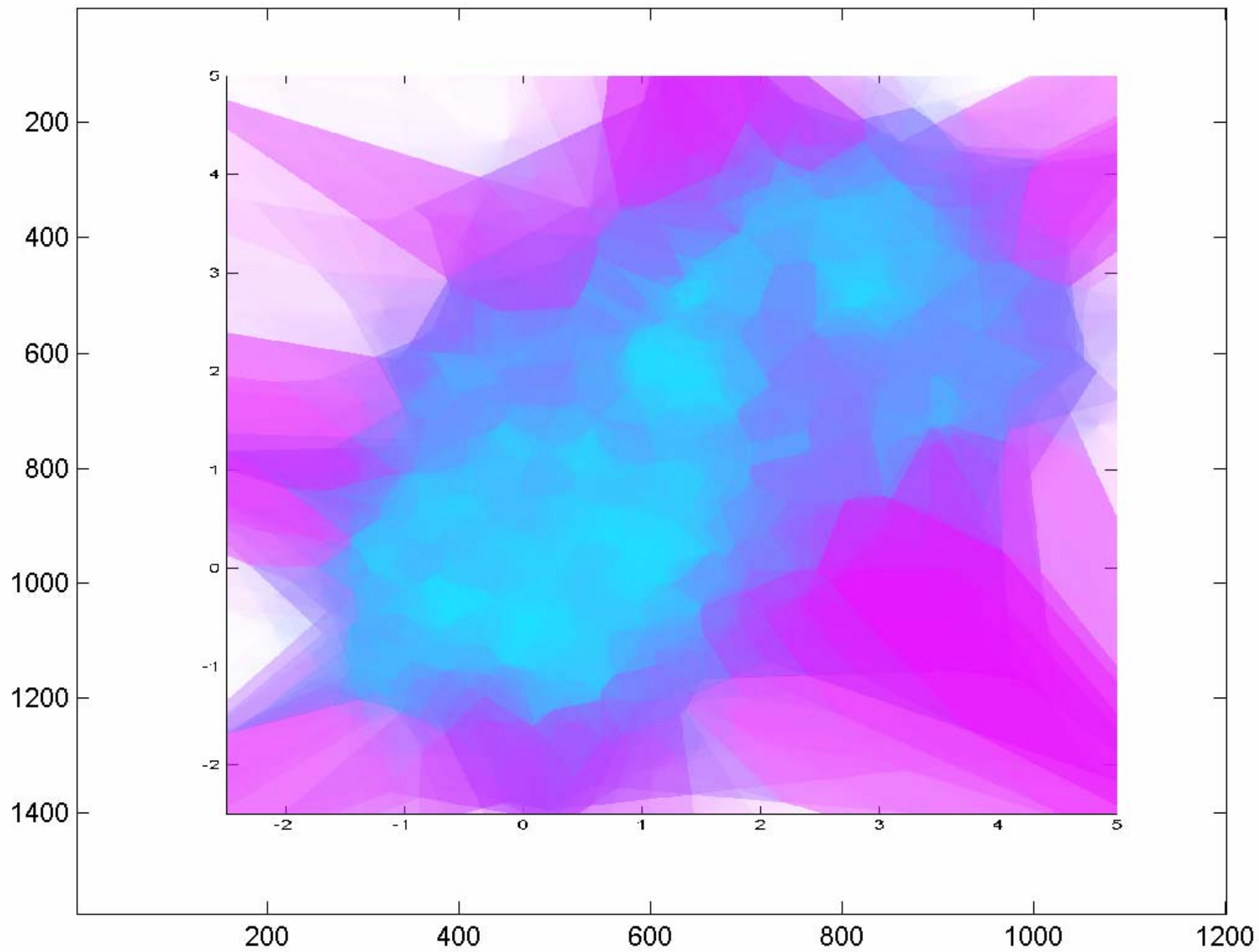
Gamma-rays from Galactic Center 12 blocks

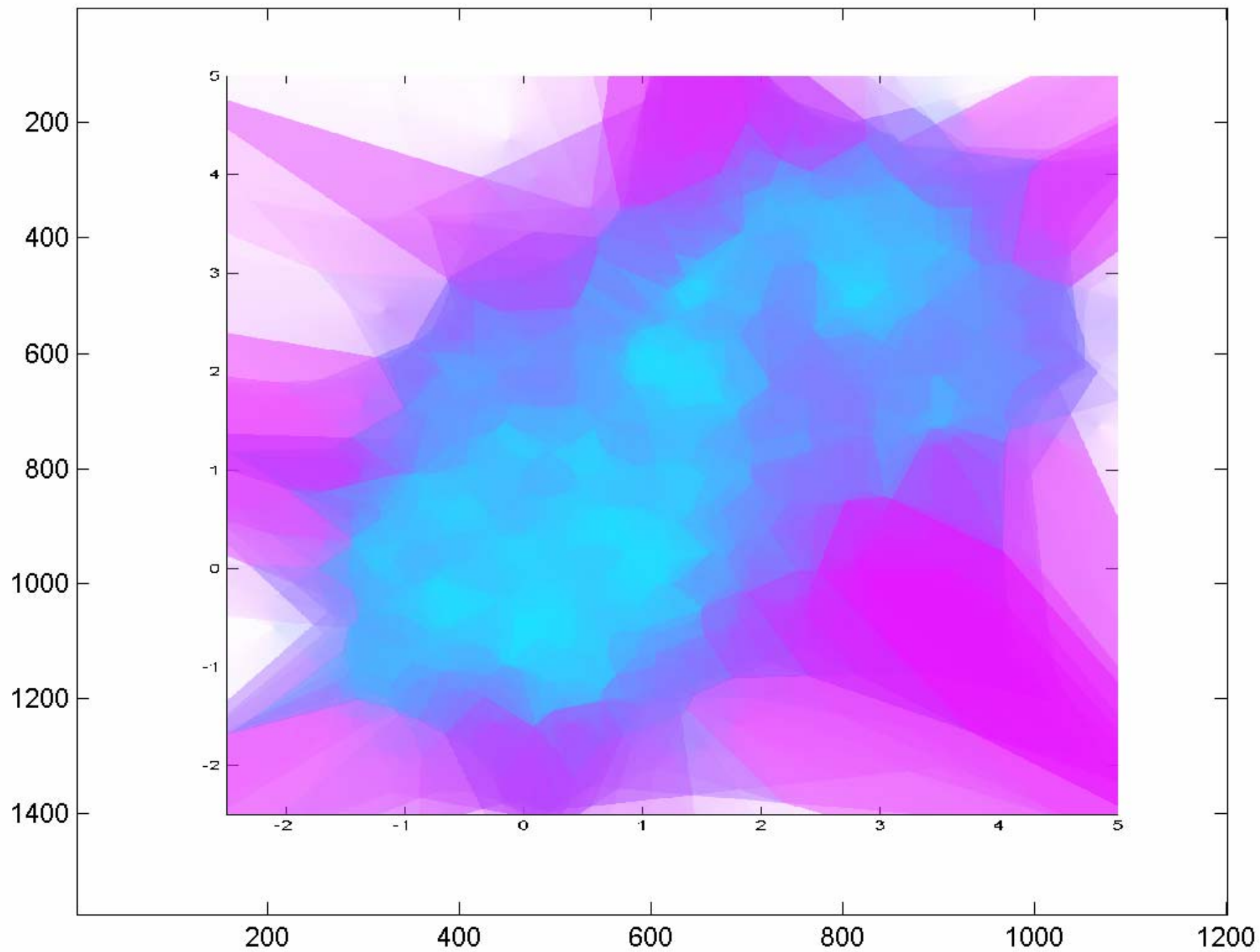


# 2D Synthetic Bootstrap Example: Raw Data









$E > 1 \text{ GEV} - \langle E \rangle / \langle \text{Area} \rangle$

# Local Mean & Variance of Area/Energy (idea due to Bill Atwood)

