

# MAPR<sup>®</sup>

## High Performance Time Series Databases



# Agenda

- What is anomaly detection?
- Some examples
- Some generalization
- Compression == Truth
- Deep dive into deep learning
- Why this matters for time series databases



# Who I am

- Ted Dunning, Chief Application Architect, MapR  
[tdunning@mapr.com](mailto:tdunning@mapr.com)  
[tdunning@apache.org](mailto:tdunning@apache.org)  
@ted\_dunning
- Committer, mentor, champion, PMC member on several Apache projects
- Mahout, Drill, Zookeeper others



# Who we are

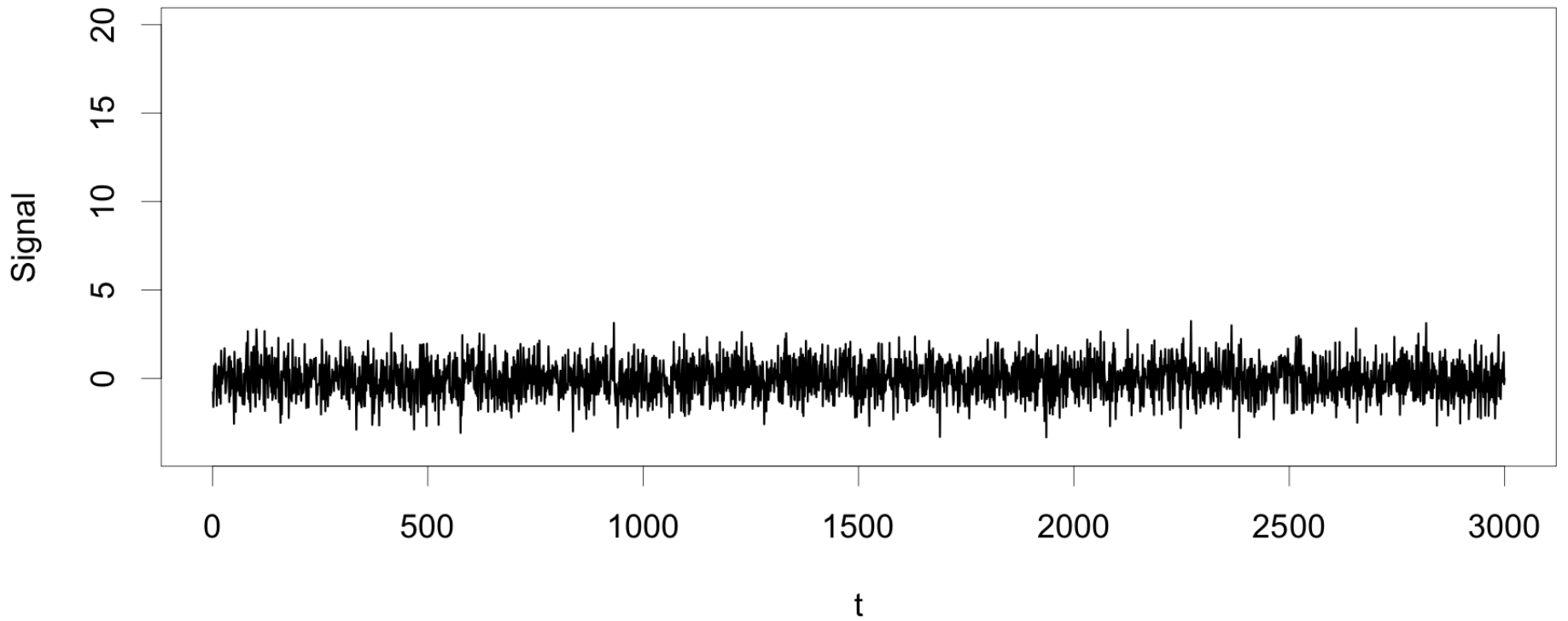
- MapR makes the technology leading distribution including Hadoop
- MapR integrates real-time data semantics directly into a system that also runs Hadoop programs seamlessly
- The biggest and best choose MapR
  - Google, Amazon
  - Largest credit card, retailer, health insurance, telco
  - Ping me for info

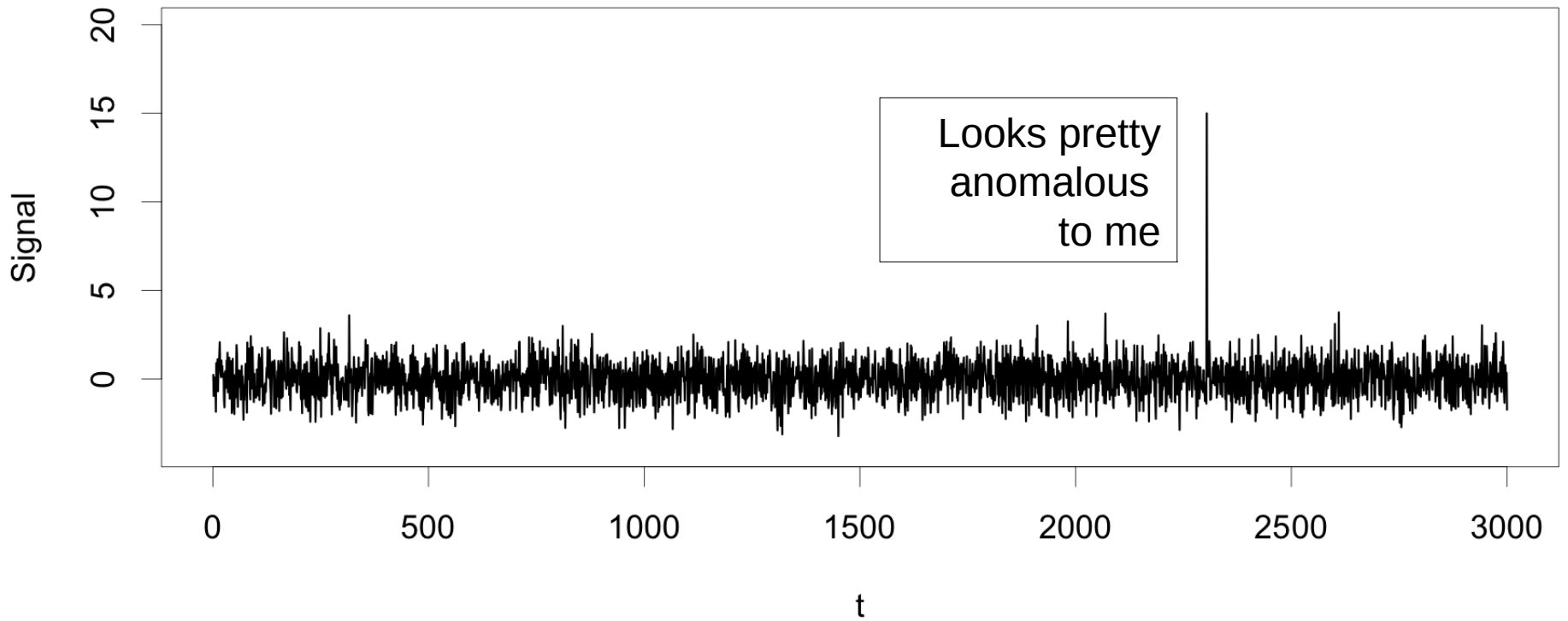


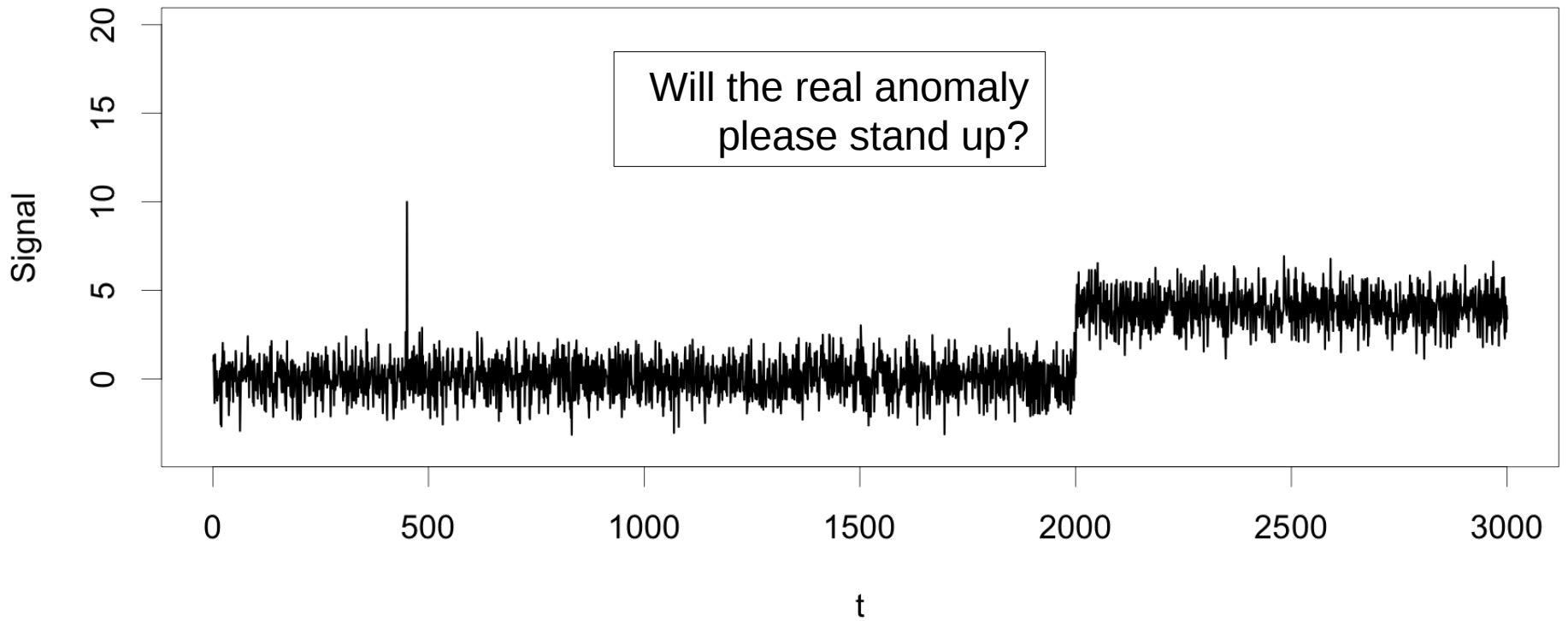
# What is Anomaly Detection?

- What just happened that shouldn't?
  - but I don't know what failure looks like (yet)
- Find the problem before other people see it
  - especially customers and CEO's
- But don't wake me up if it isn't really broken









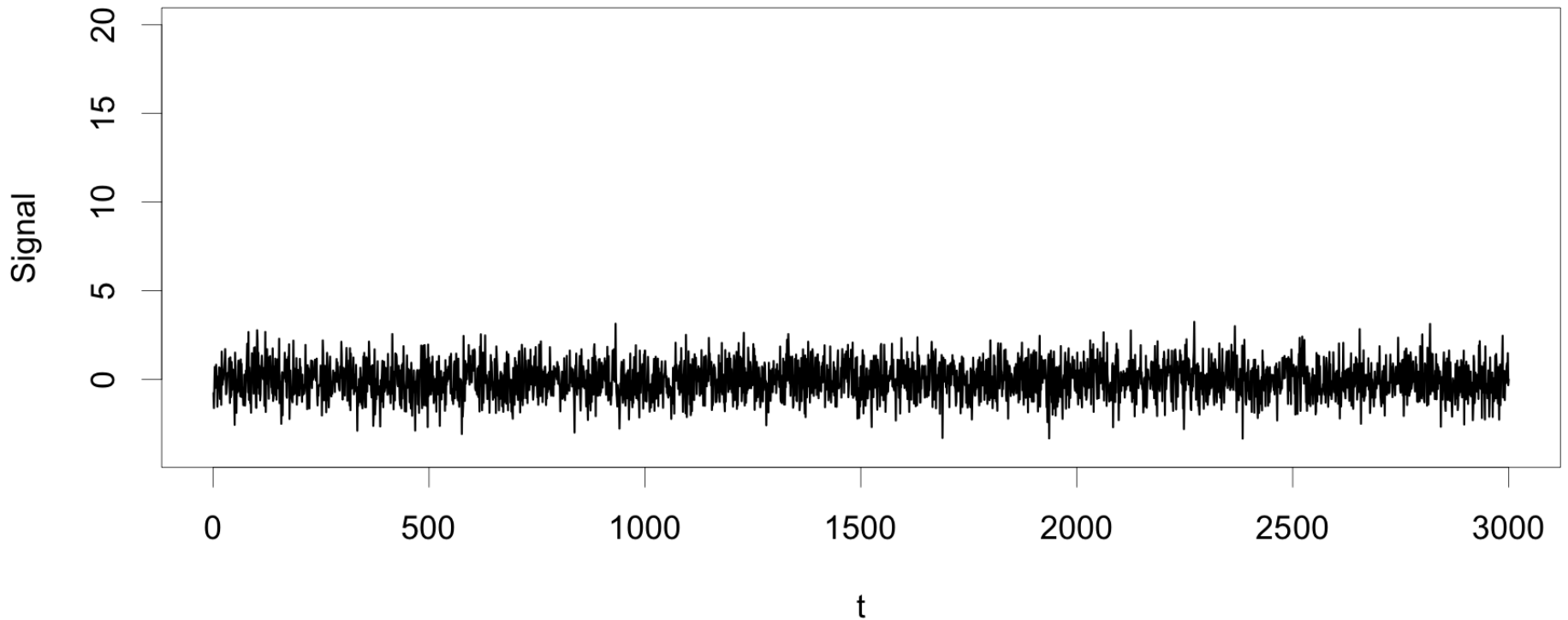


# What Are We Really Doing

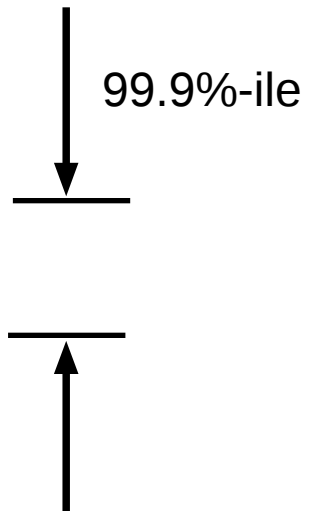
- We want action when something breaks  
(dies/falls over/otherwise gets in trouble)
- But action is expensive
- So we don't want false alarms
- And we don't want false negatives
  
- We need to trade off costs



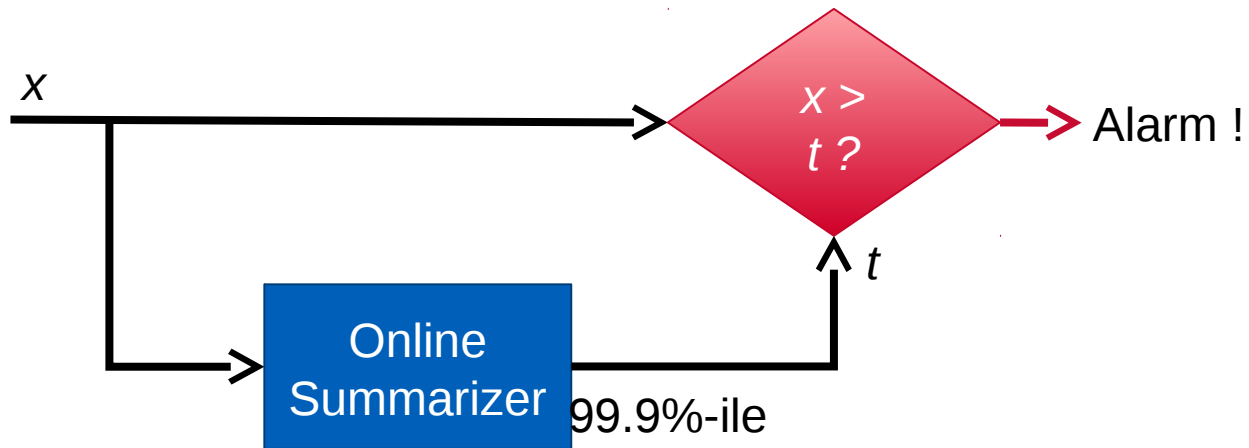
# A Second Look



# A Second Look



# How Hard Can it Be?

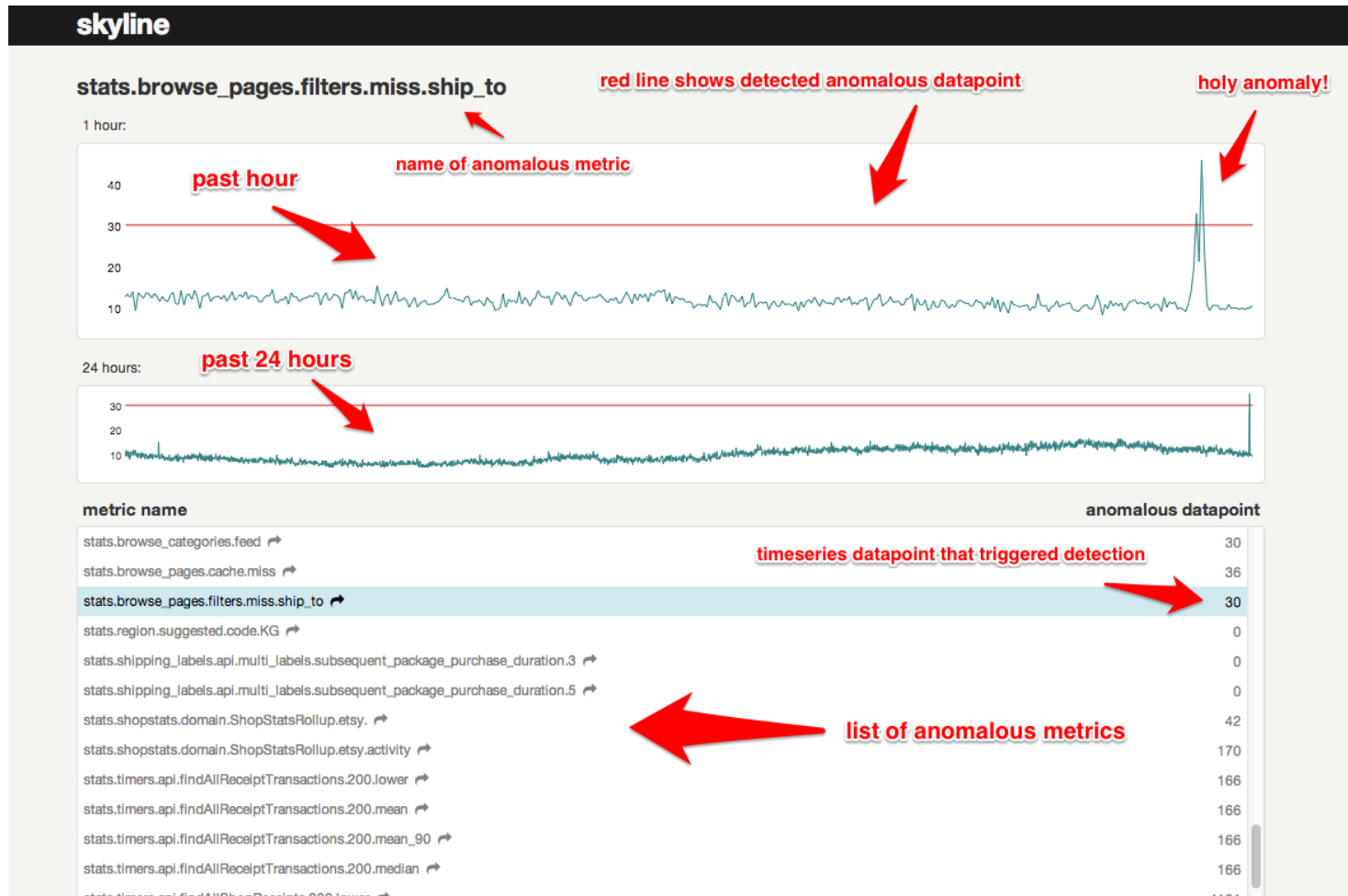


# On-line Percentile Estimates

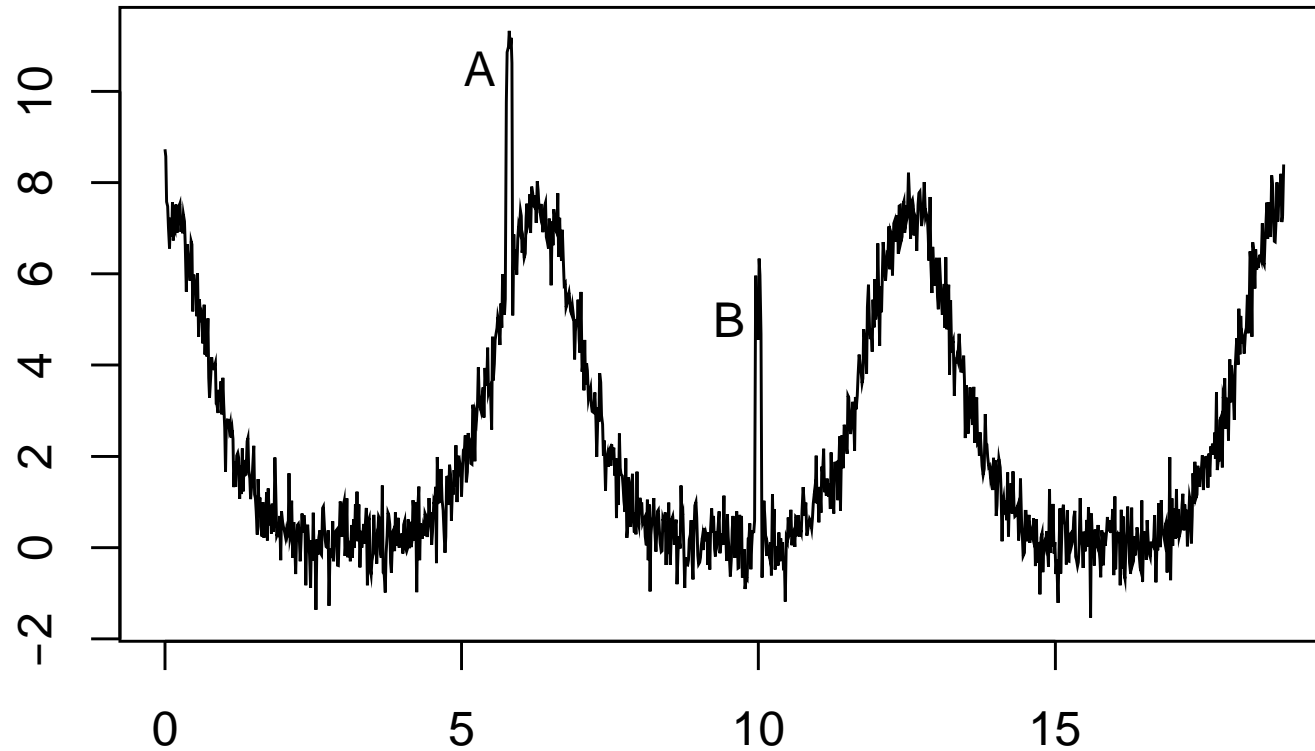
- Apache Mahout has on-line percentile estimator
  - very high accuracy for extreme tails
  - new in version 0.9 !!
- What's the big deal with anomaly detection?
- This looks like a solved problem



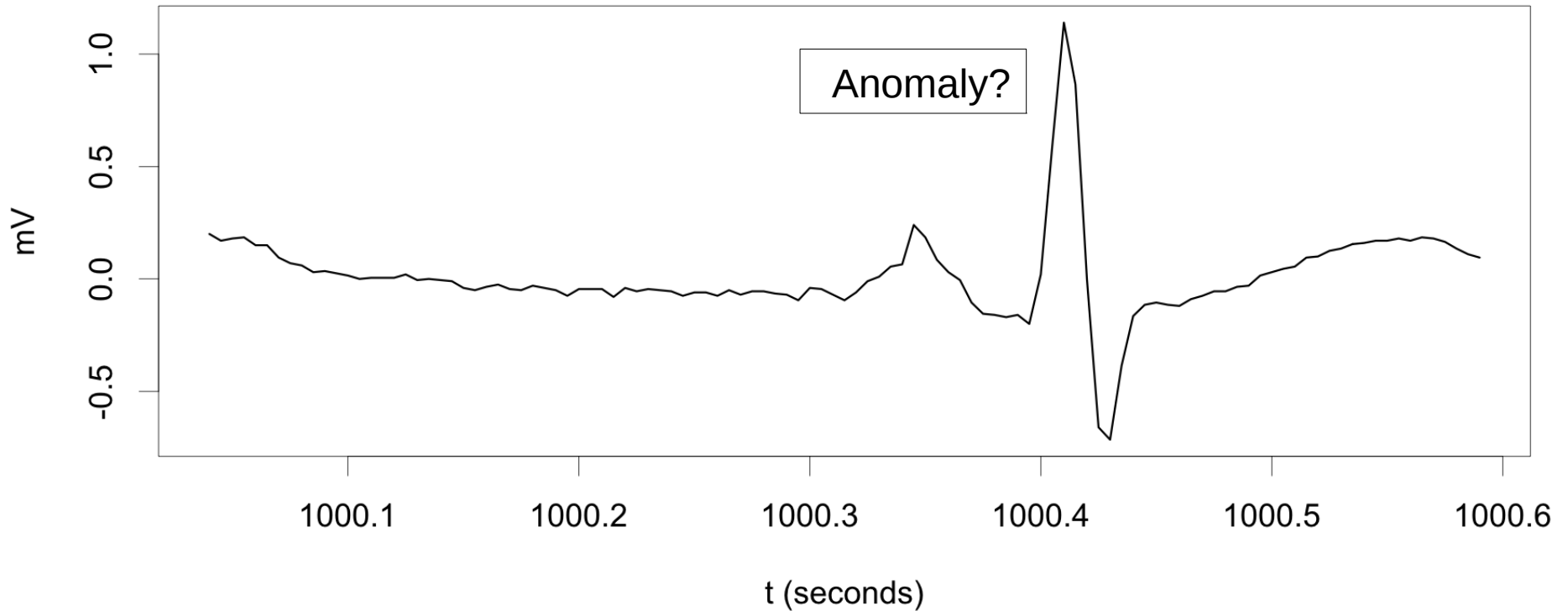
# Already Done? Etsy Skyline?



# What About This?

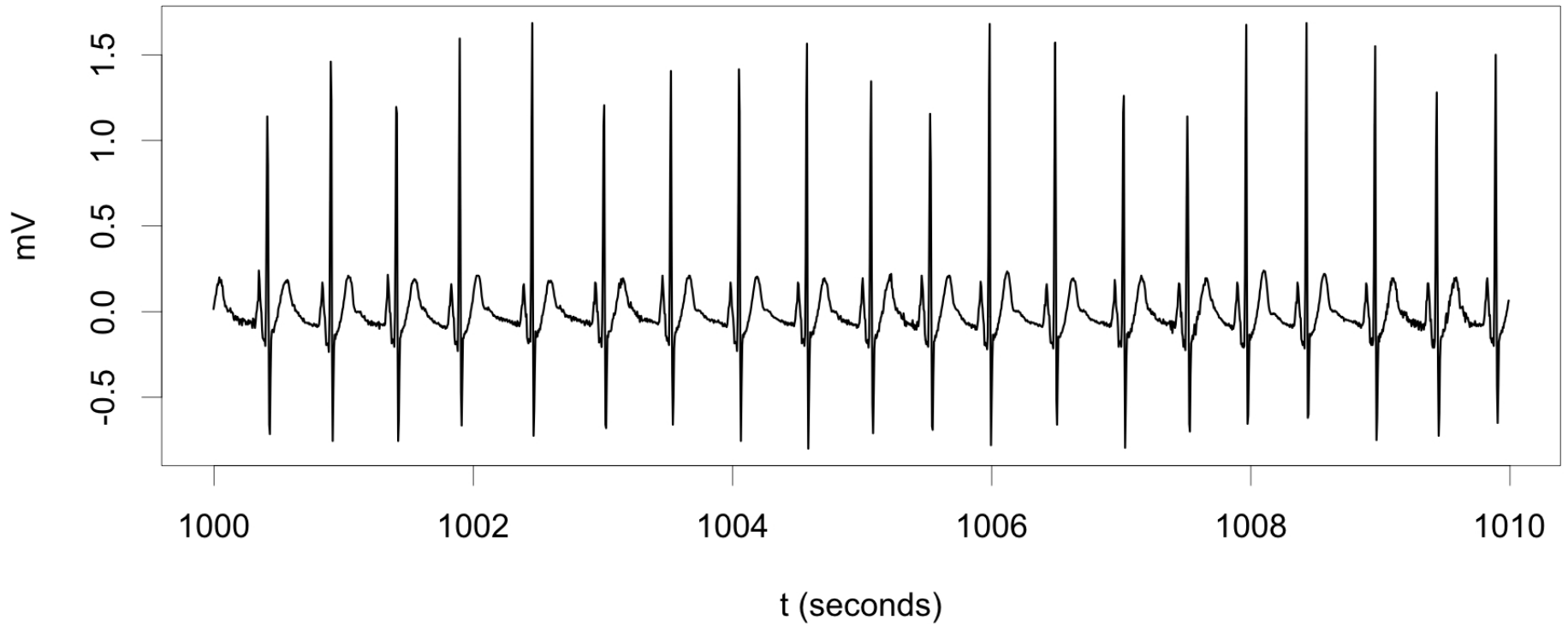


# Spot the Anomaly

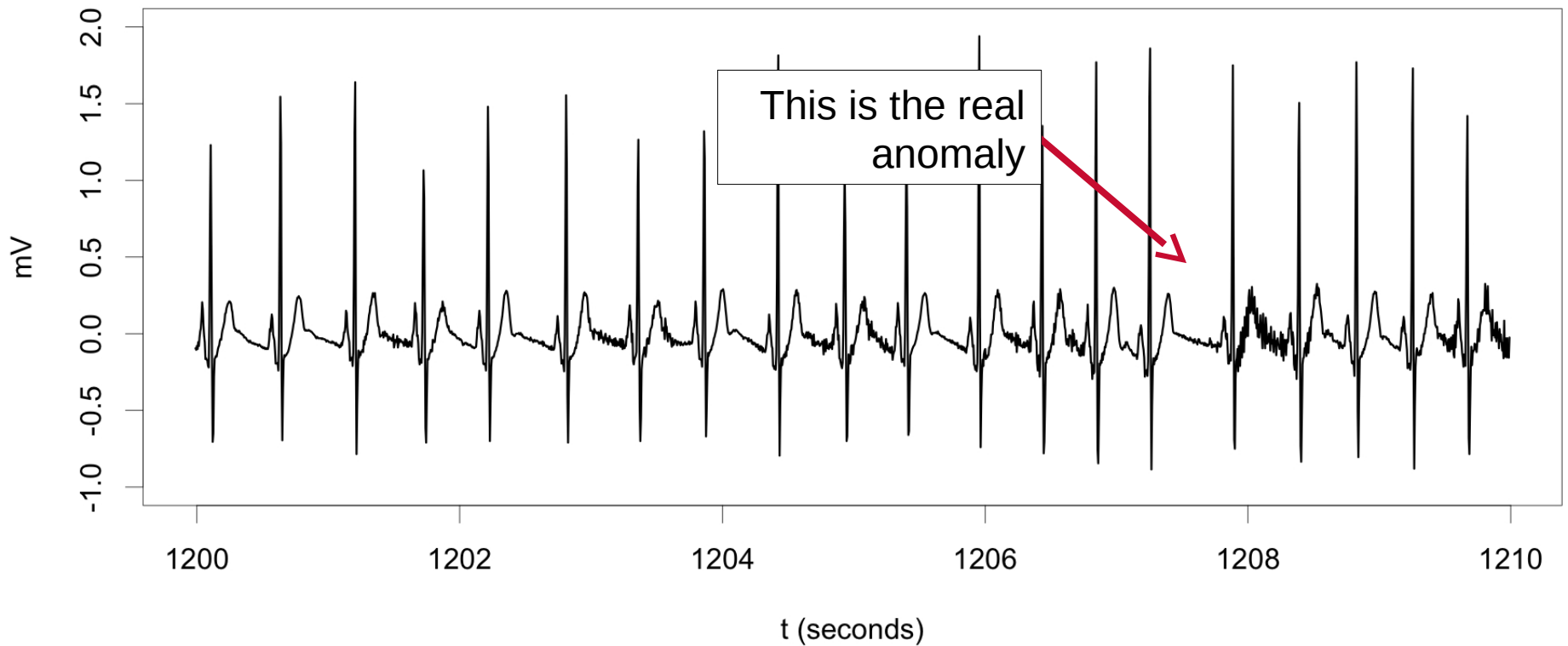




# Maybe not!



# Where's Waldo?



# Normal Isn't Just Normal

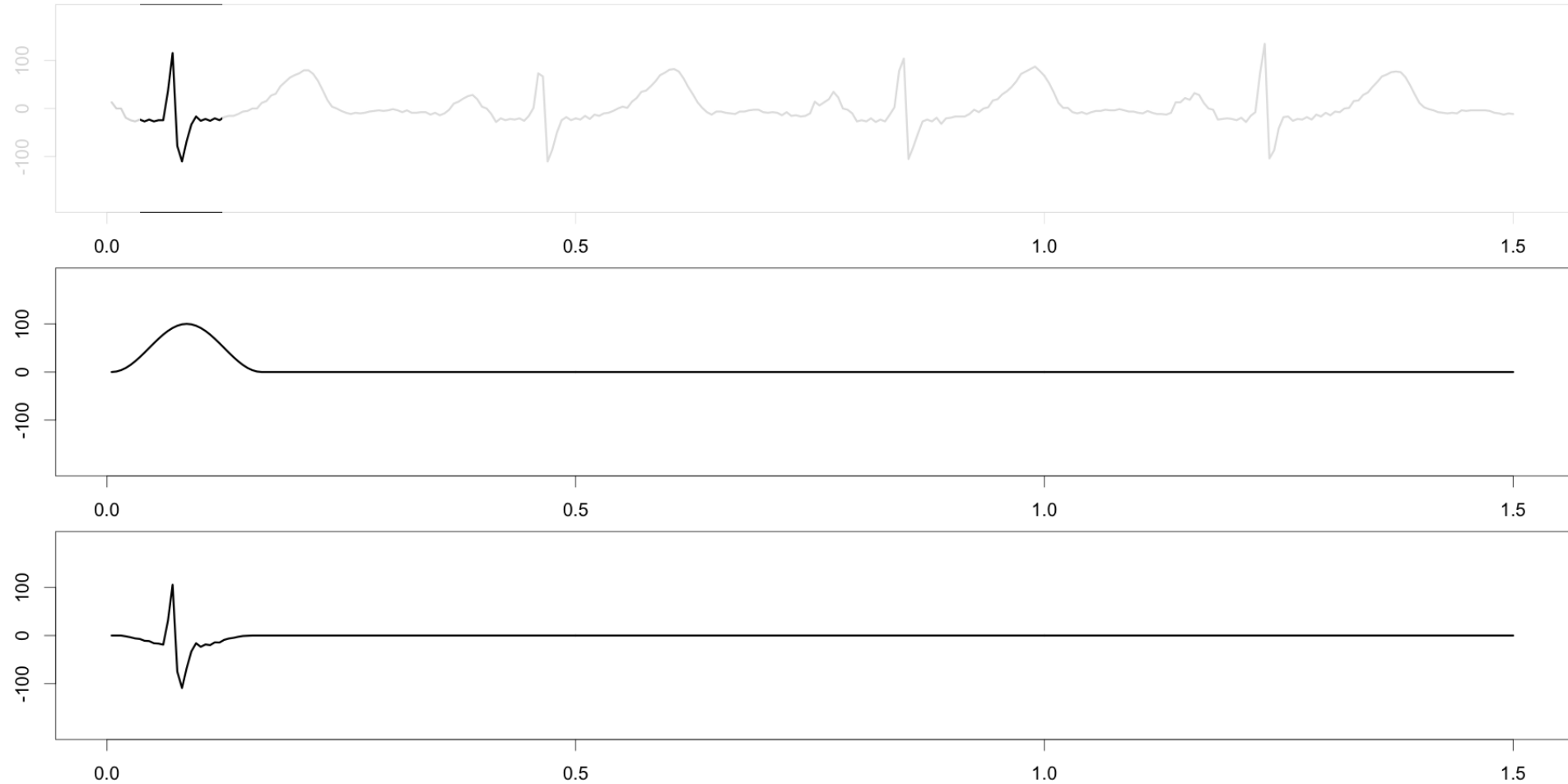
- What we want is a *model* of what is normal
- What doesn't fit the model is the *anomaly*
- For simple signals, the model can be simple ...

$$X \sim \mathcal{N}(0, \varepsilon)$$

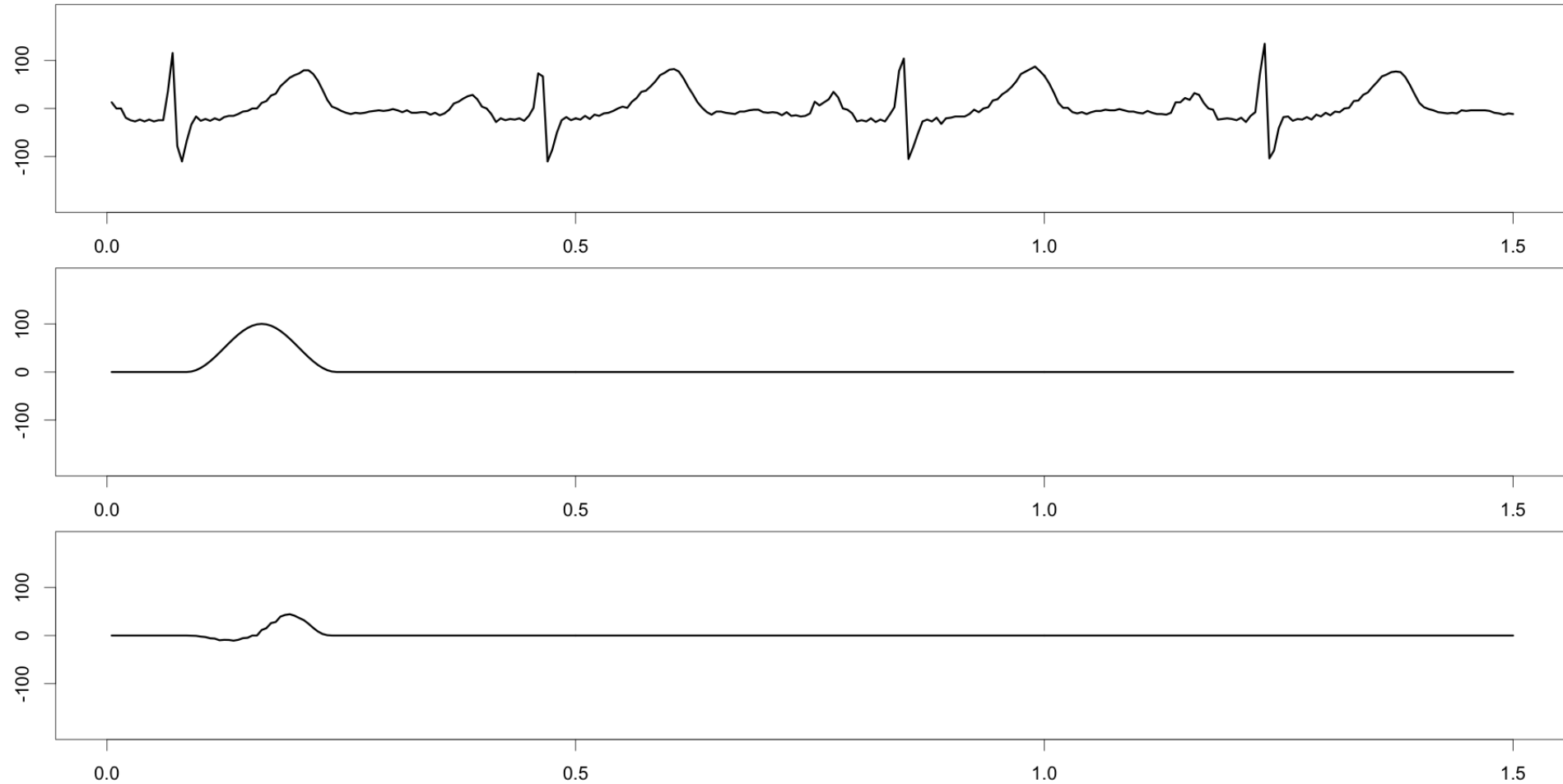
- The real world is rarely so accommodating



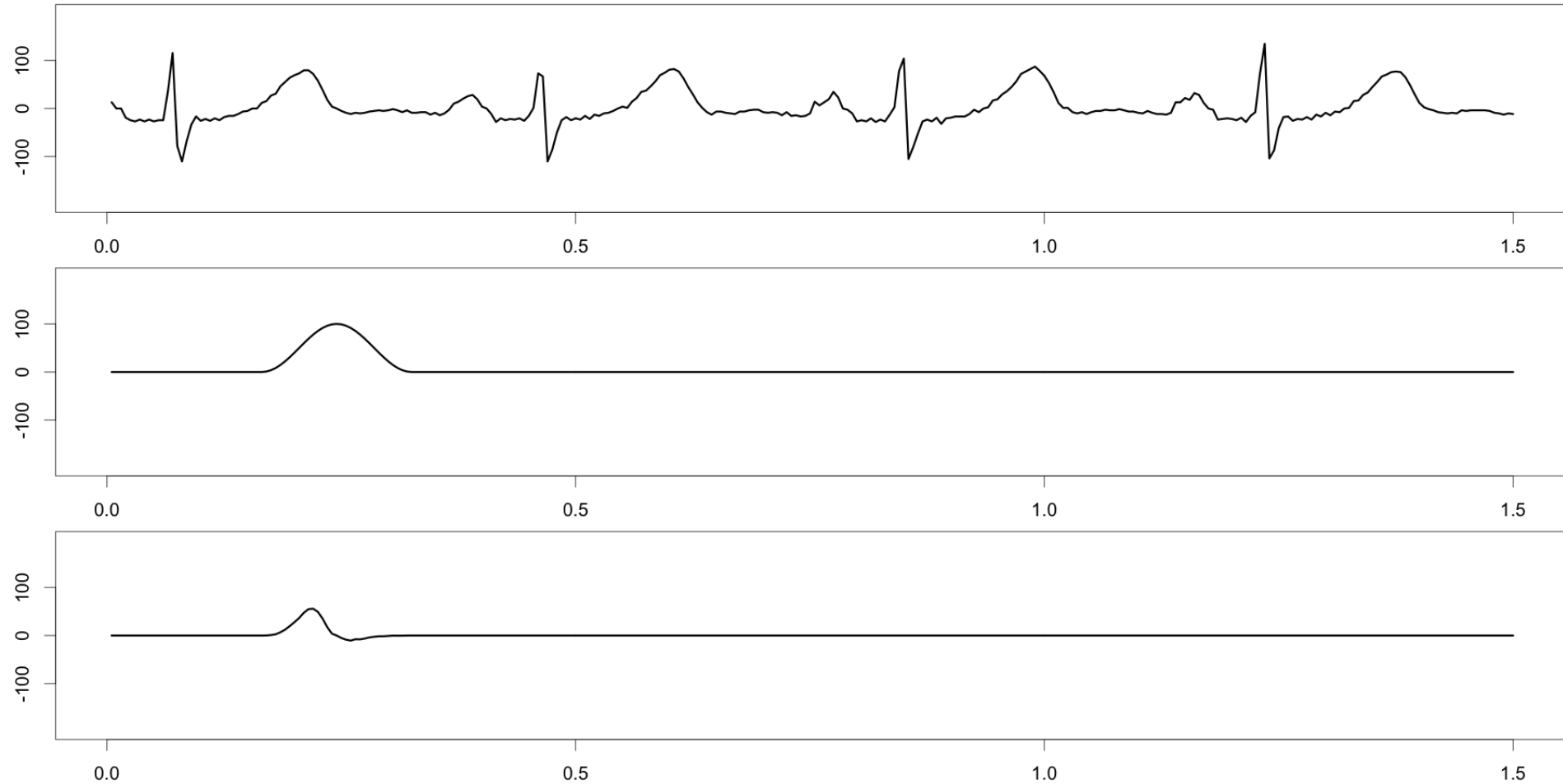
# We Do Windows



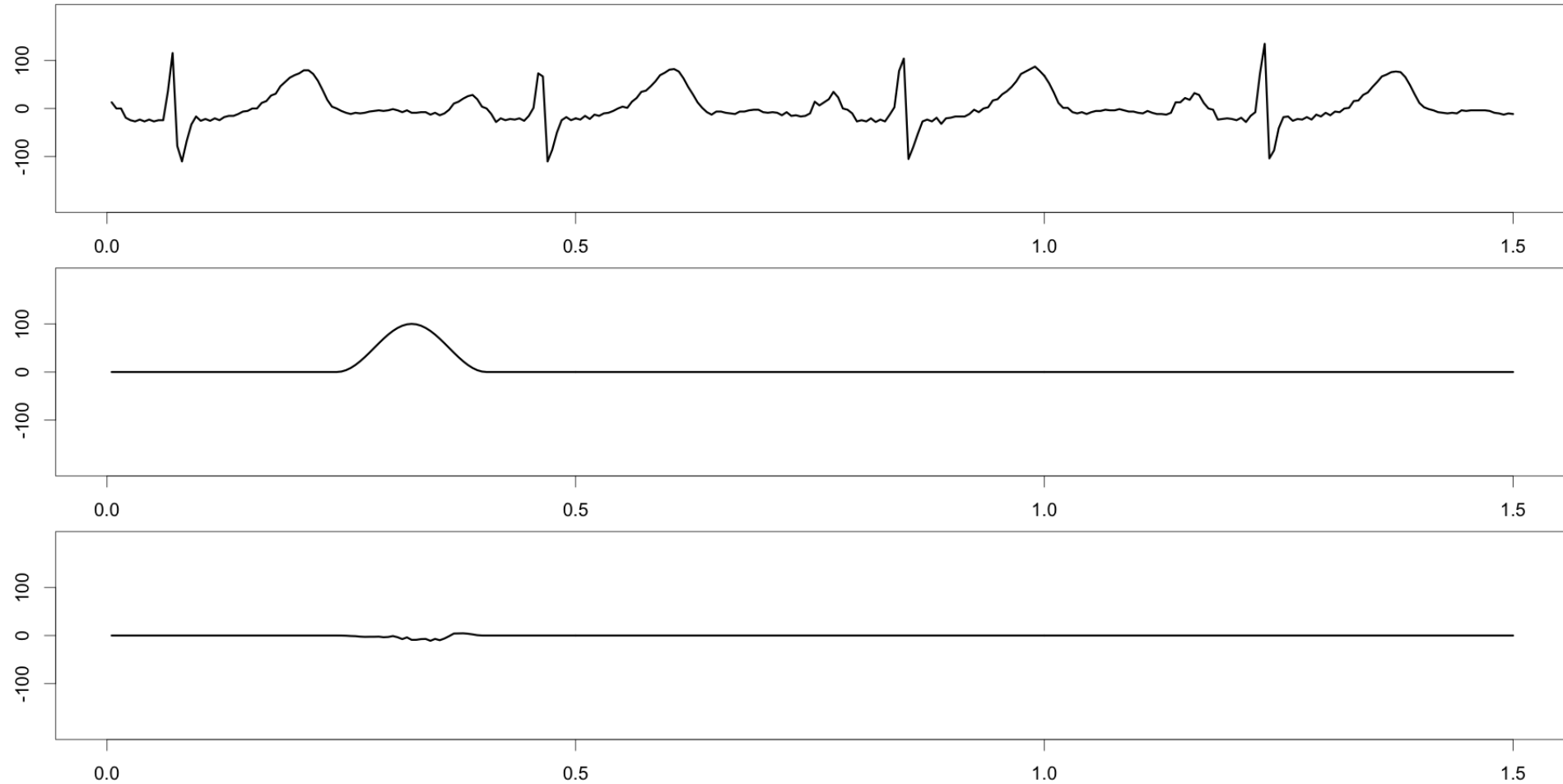
# We Do Windows



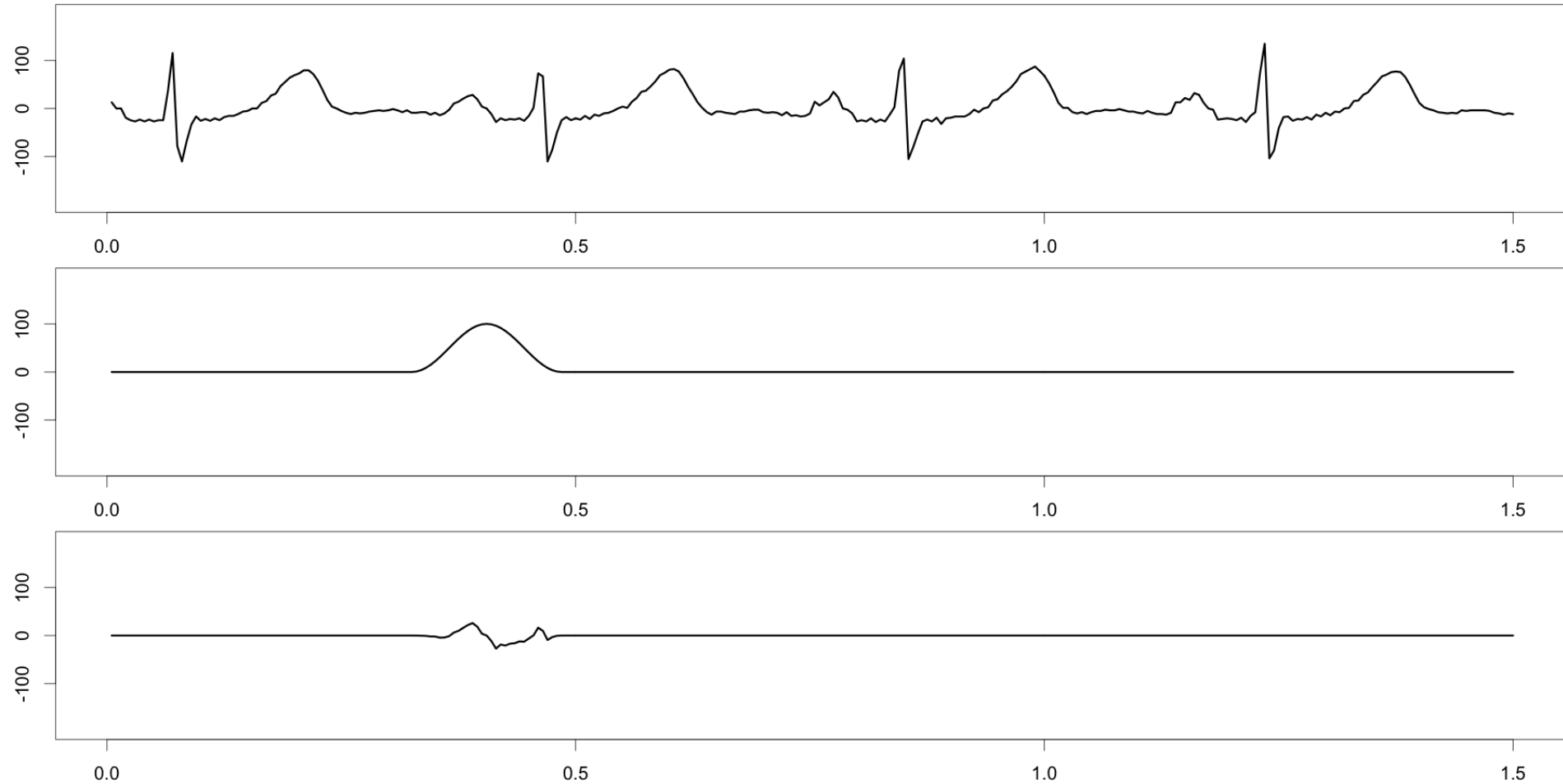
# We Do Windows



# We Do Windows

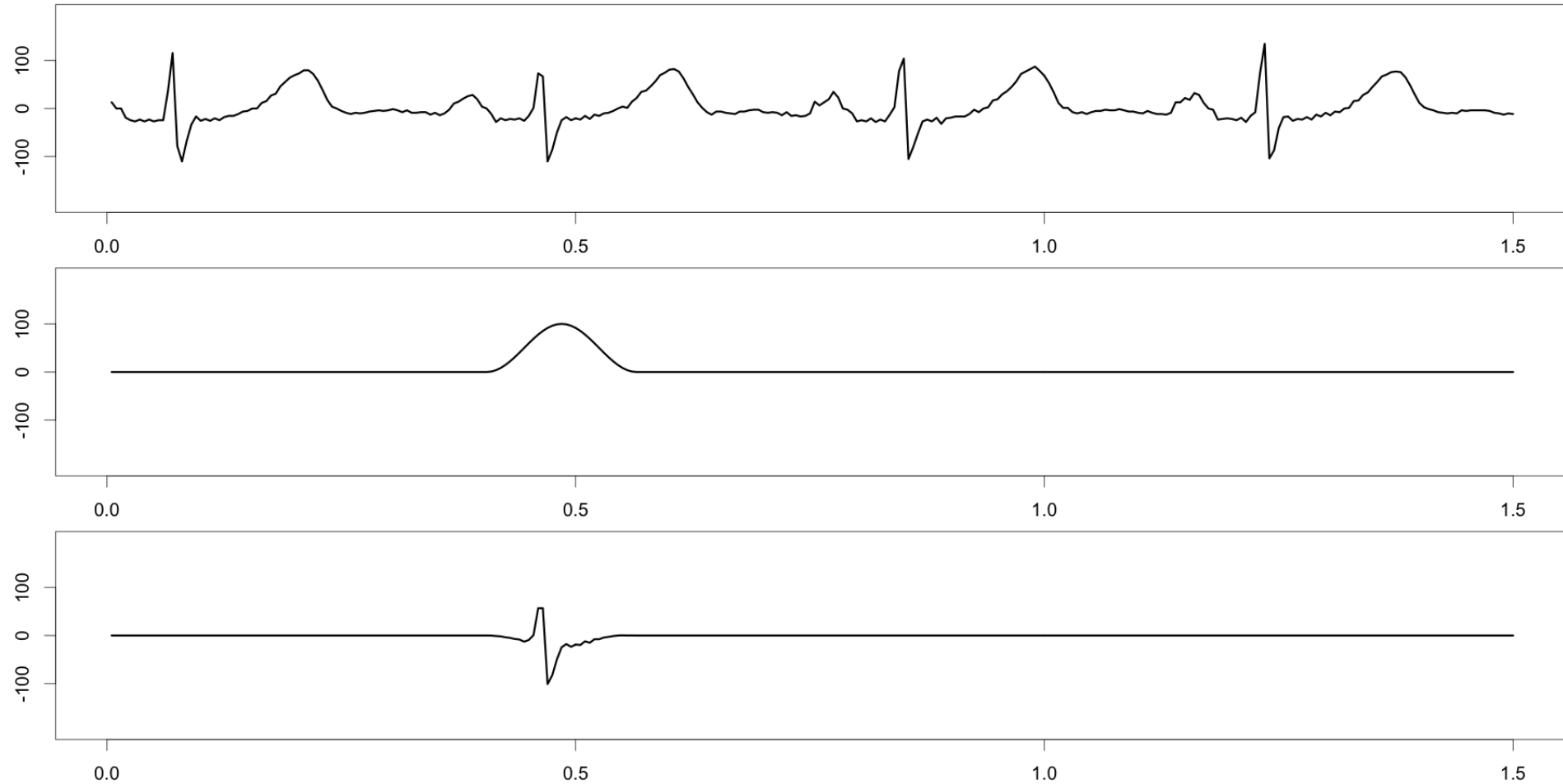


# We Do Windows

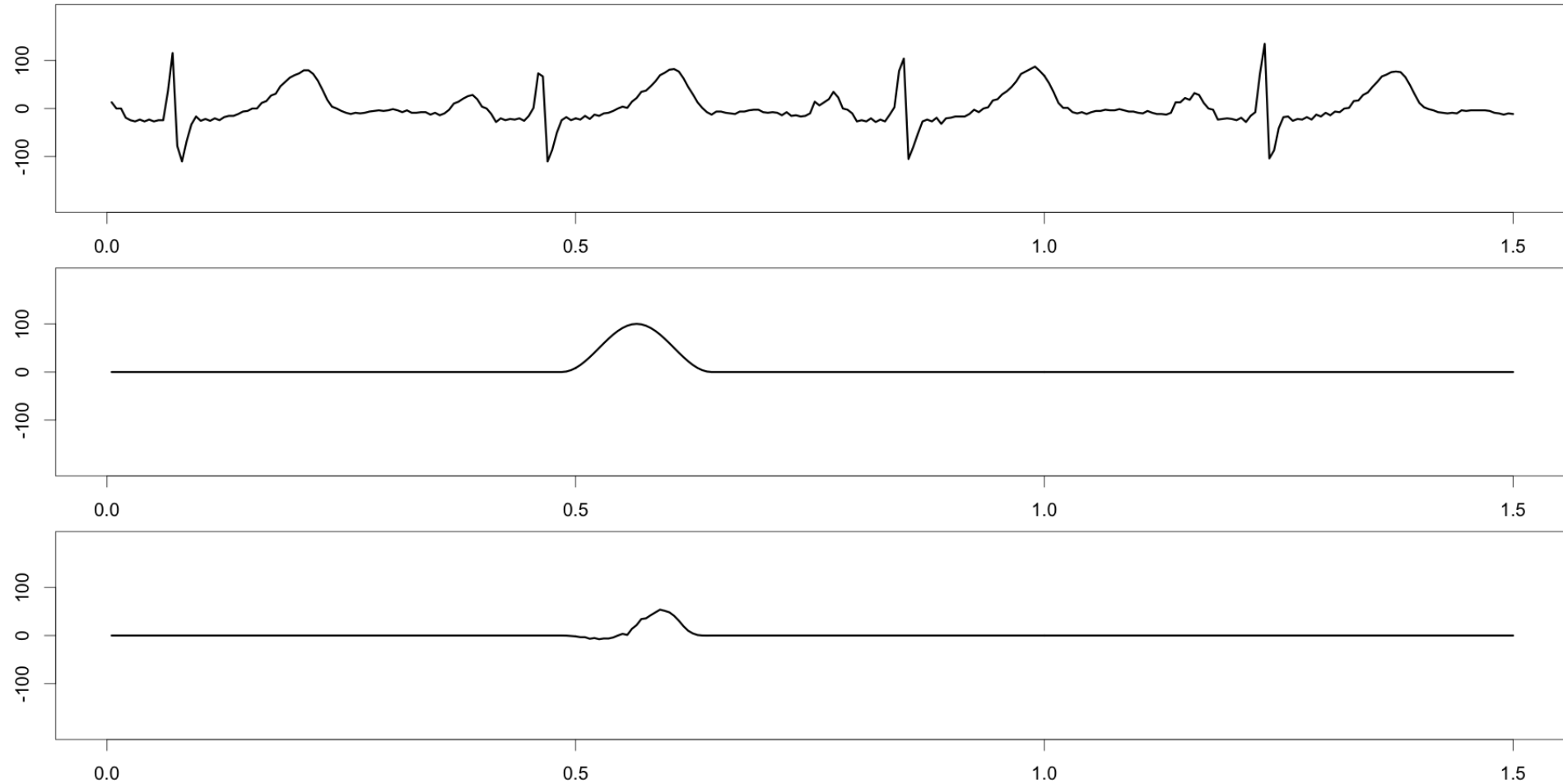




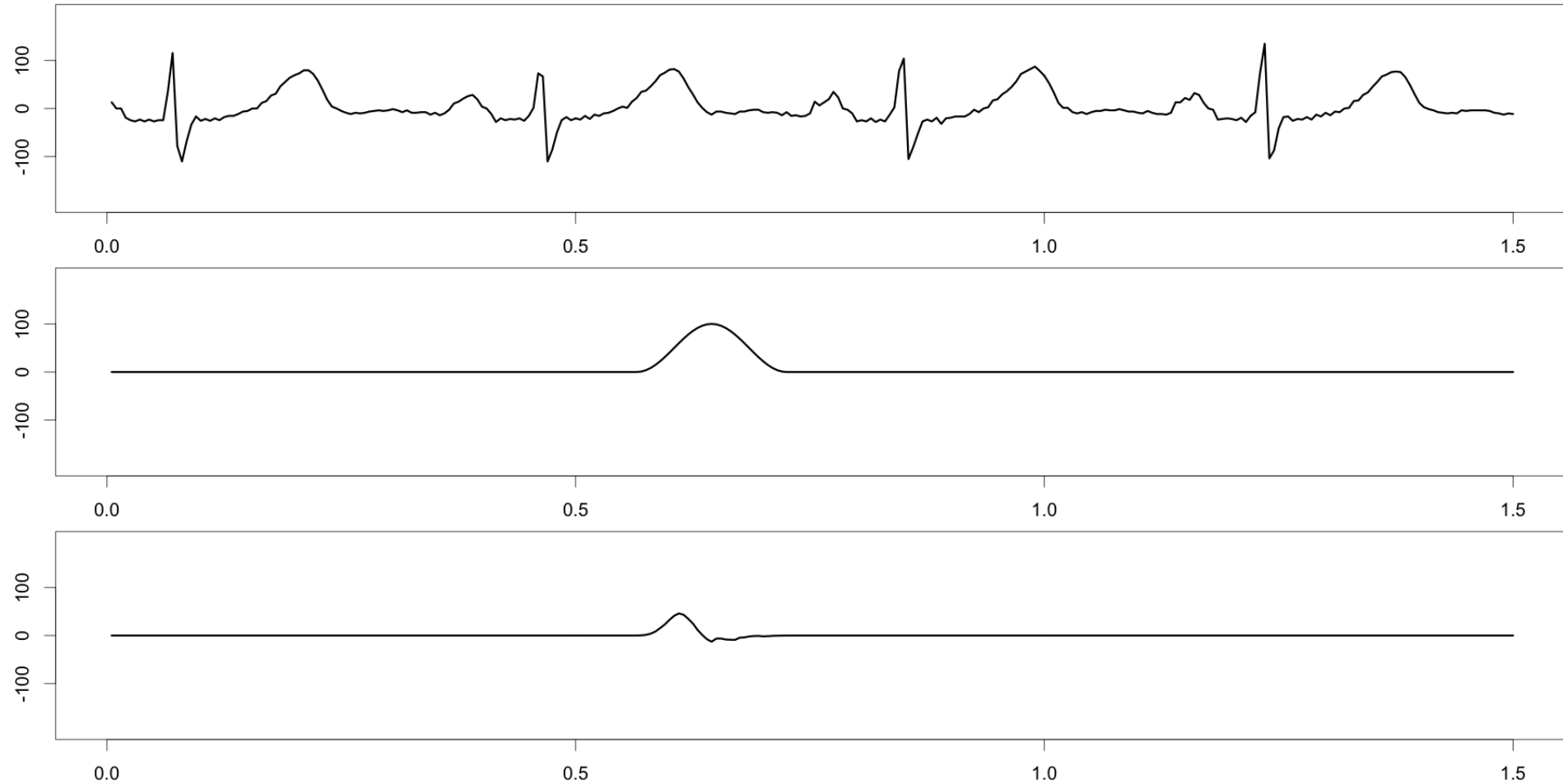
# We Do Windows



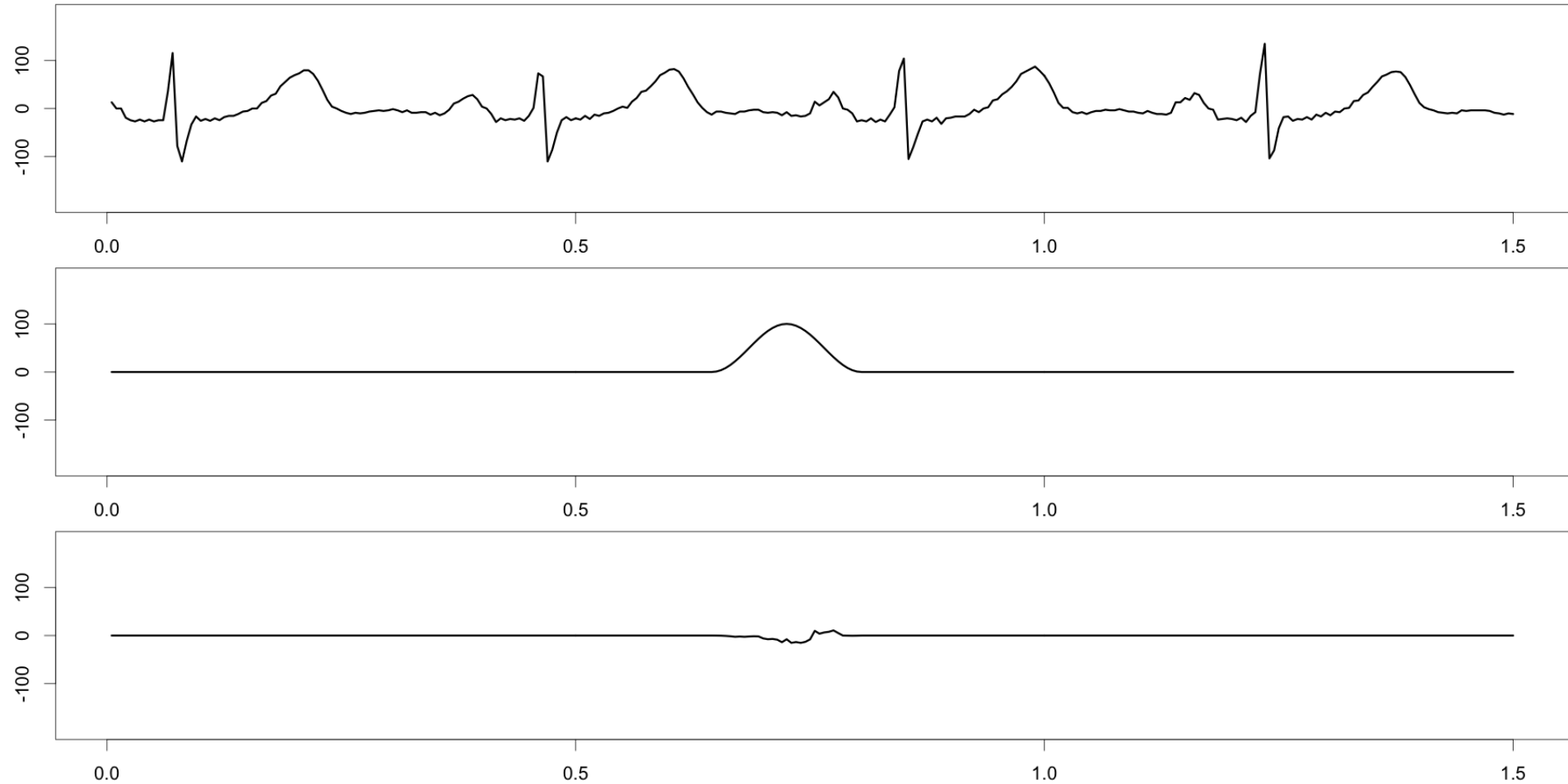
# We Do Windows



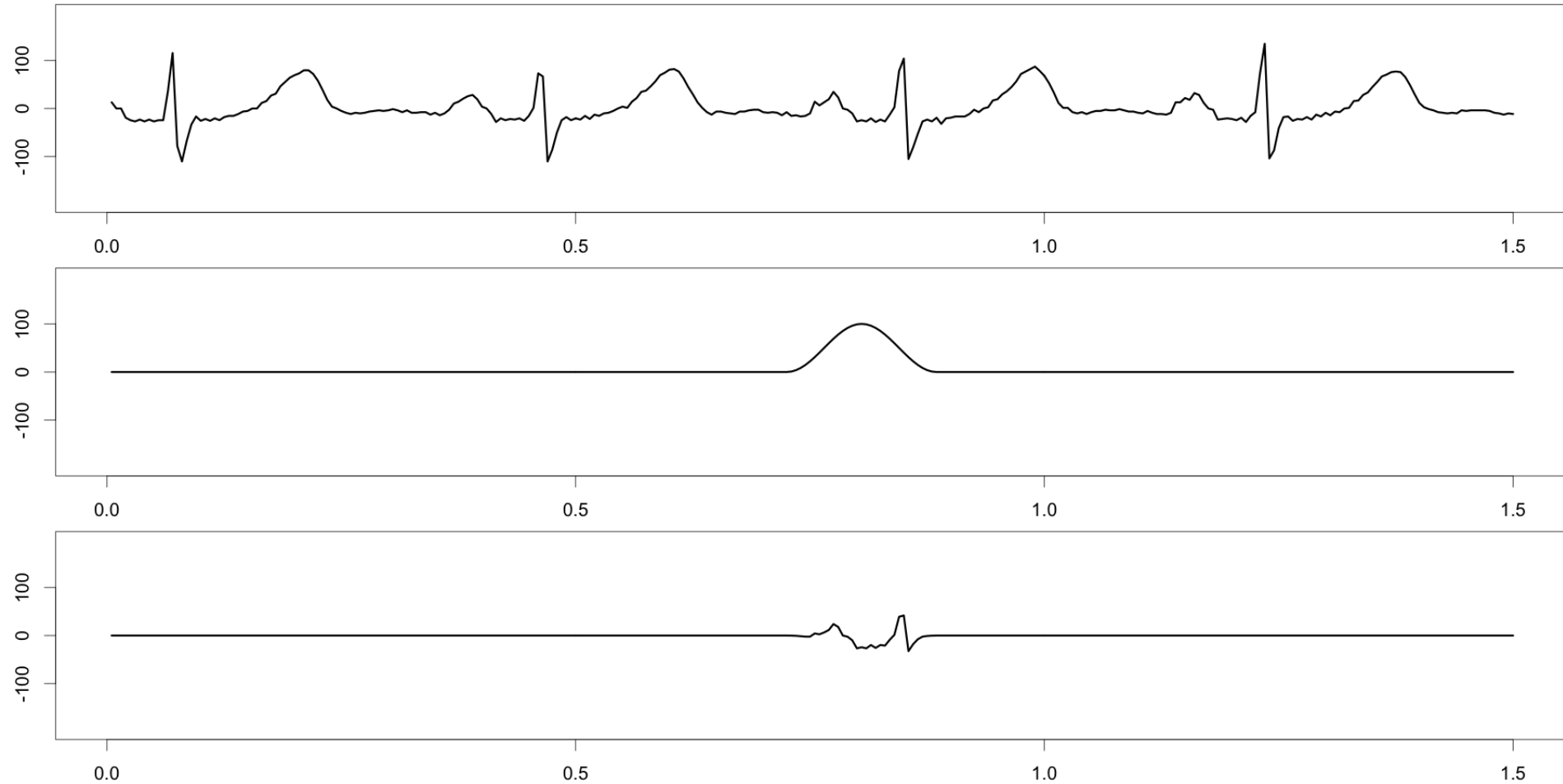
# We Do Windows



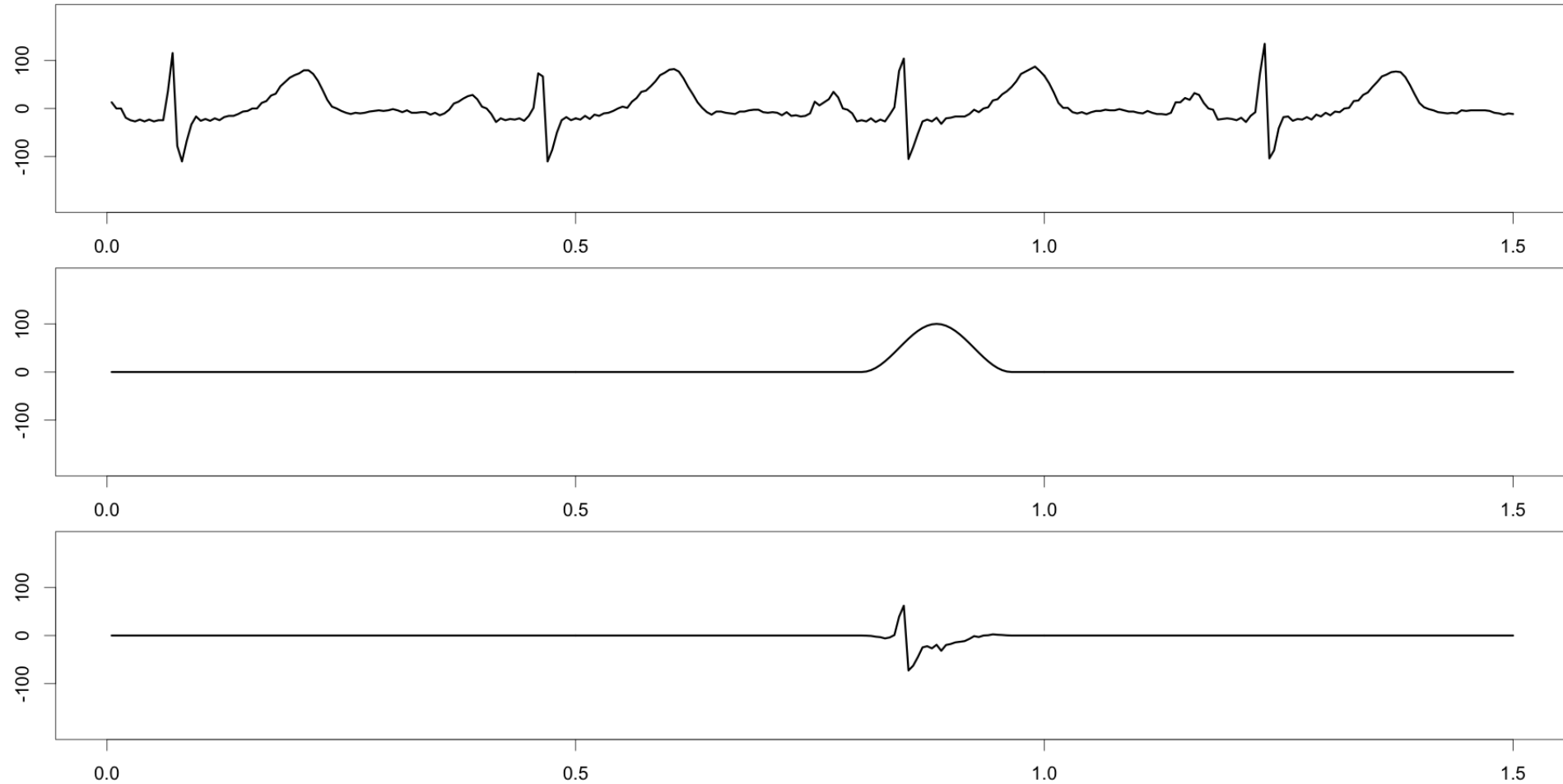
# We Do Windows



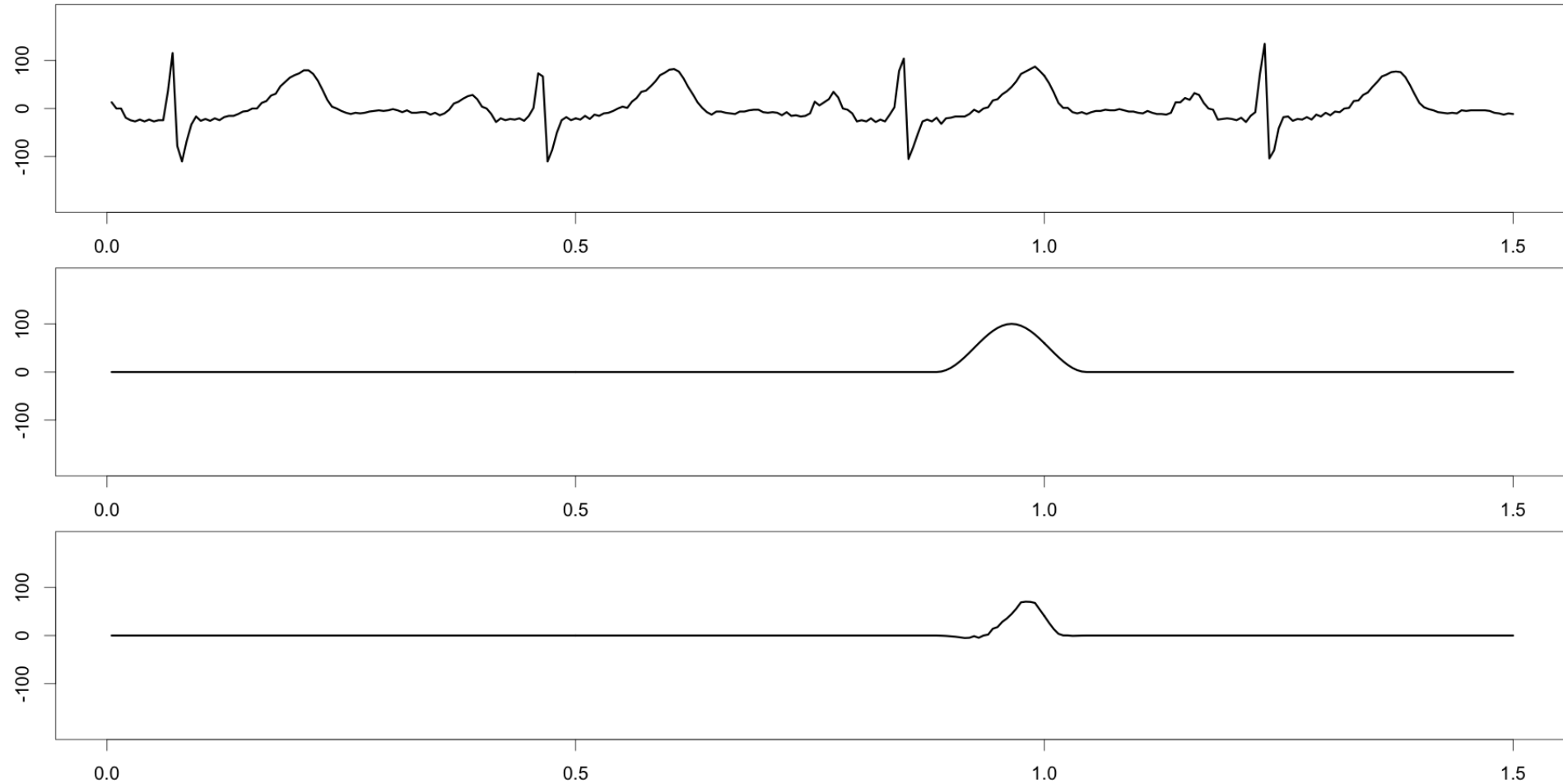
# We Do Windows



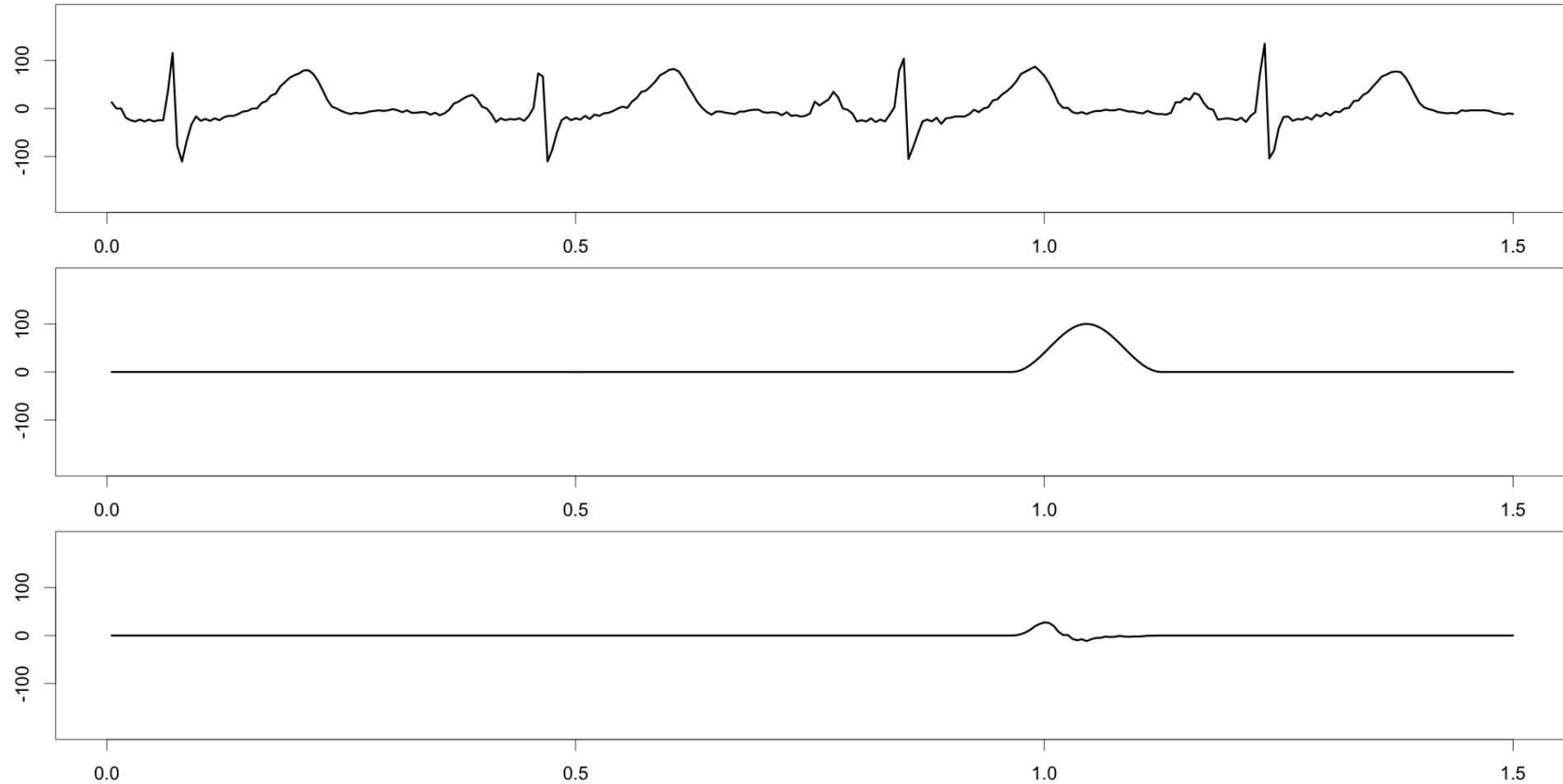
# We Do Windows



# We Do Windows

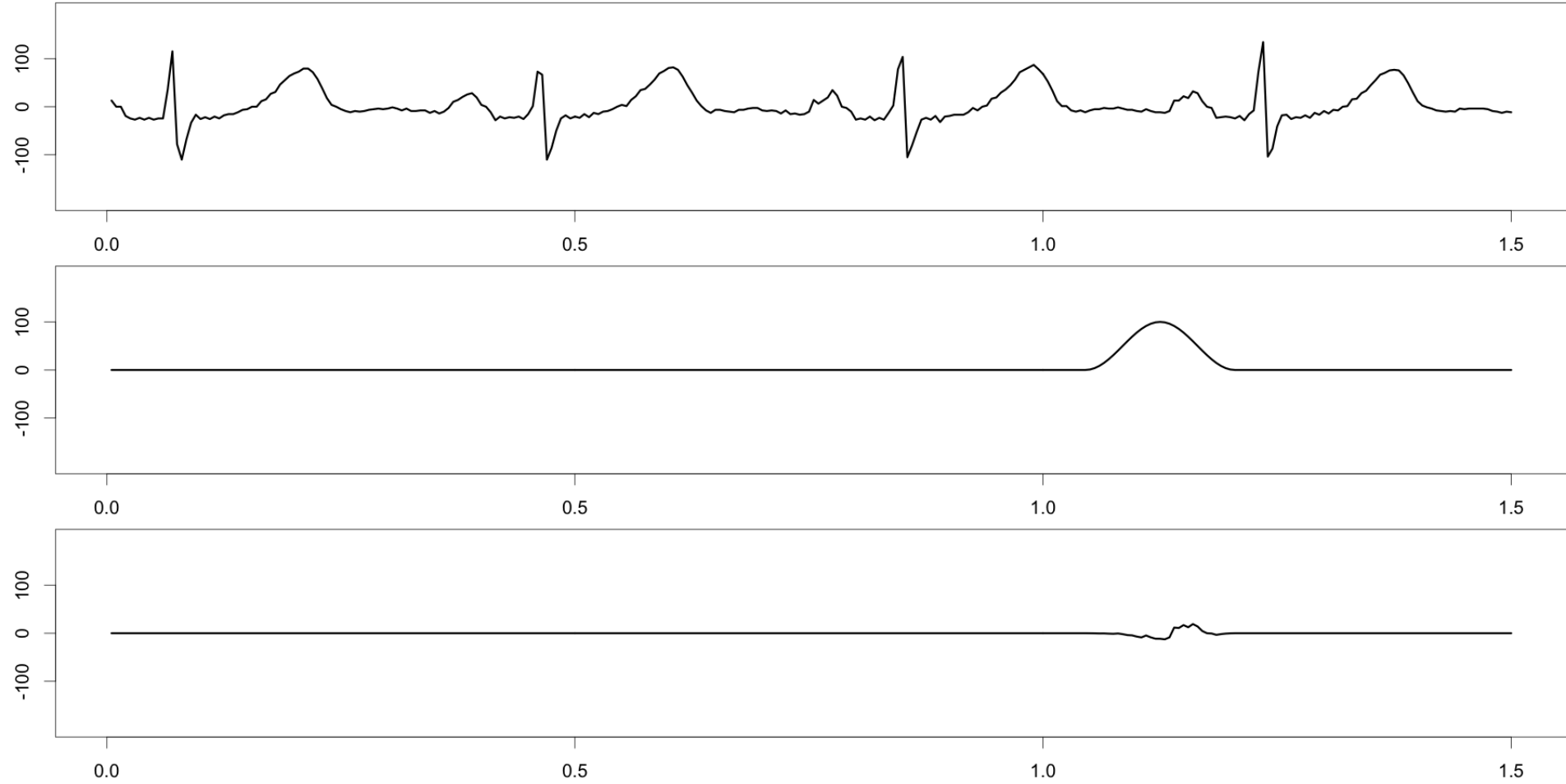


# We Do Windows

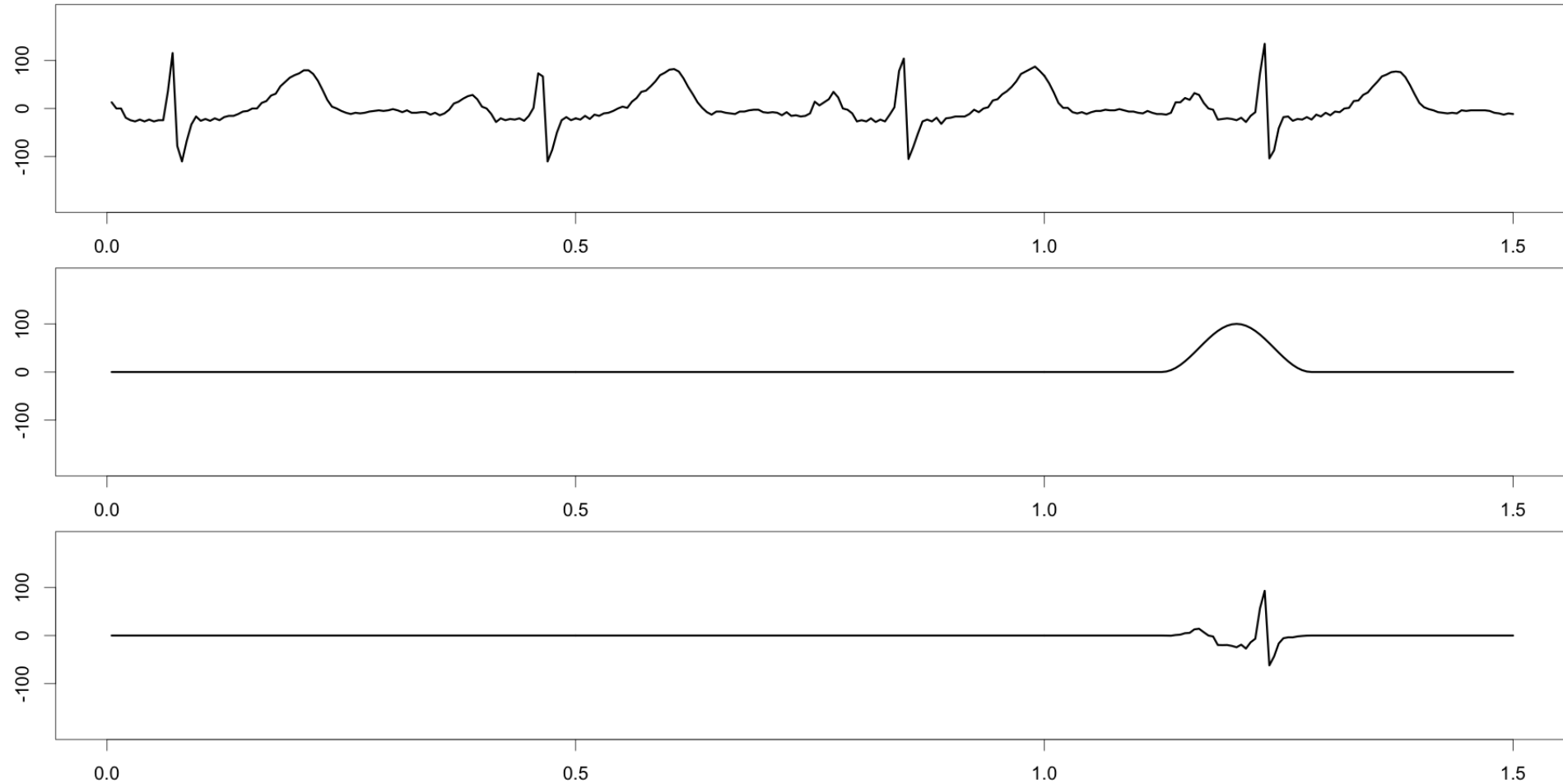




# We Do Windows



# We Do Windows

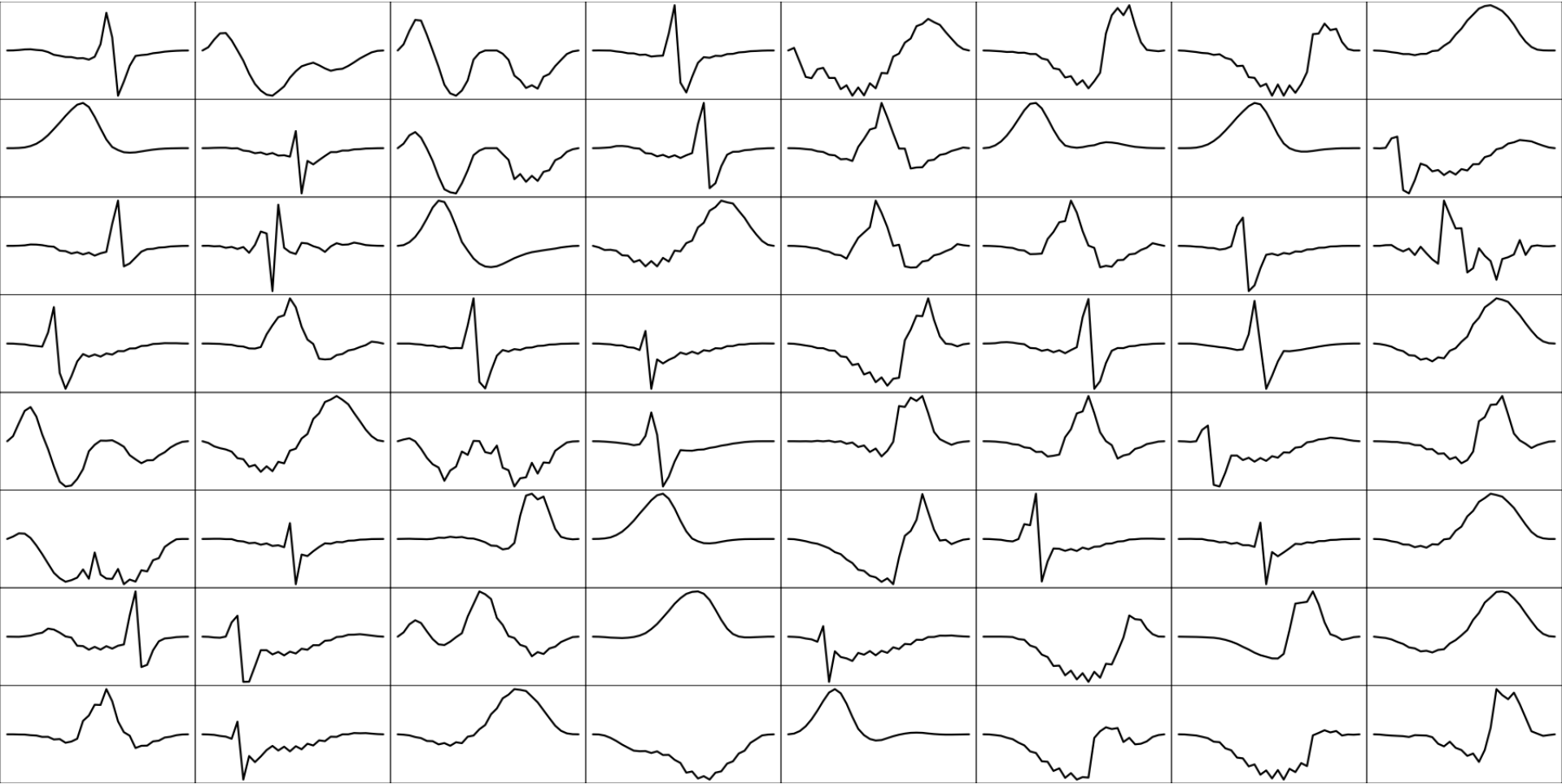


# Windows on the World

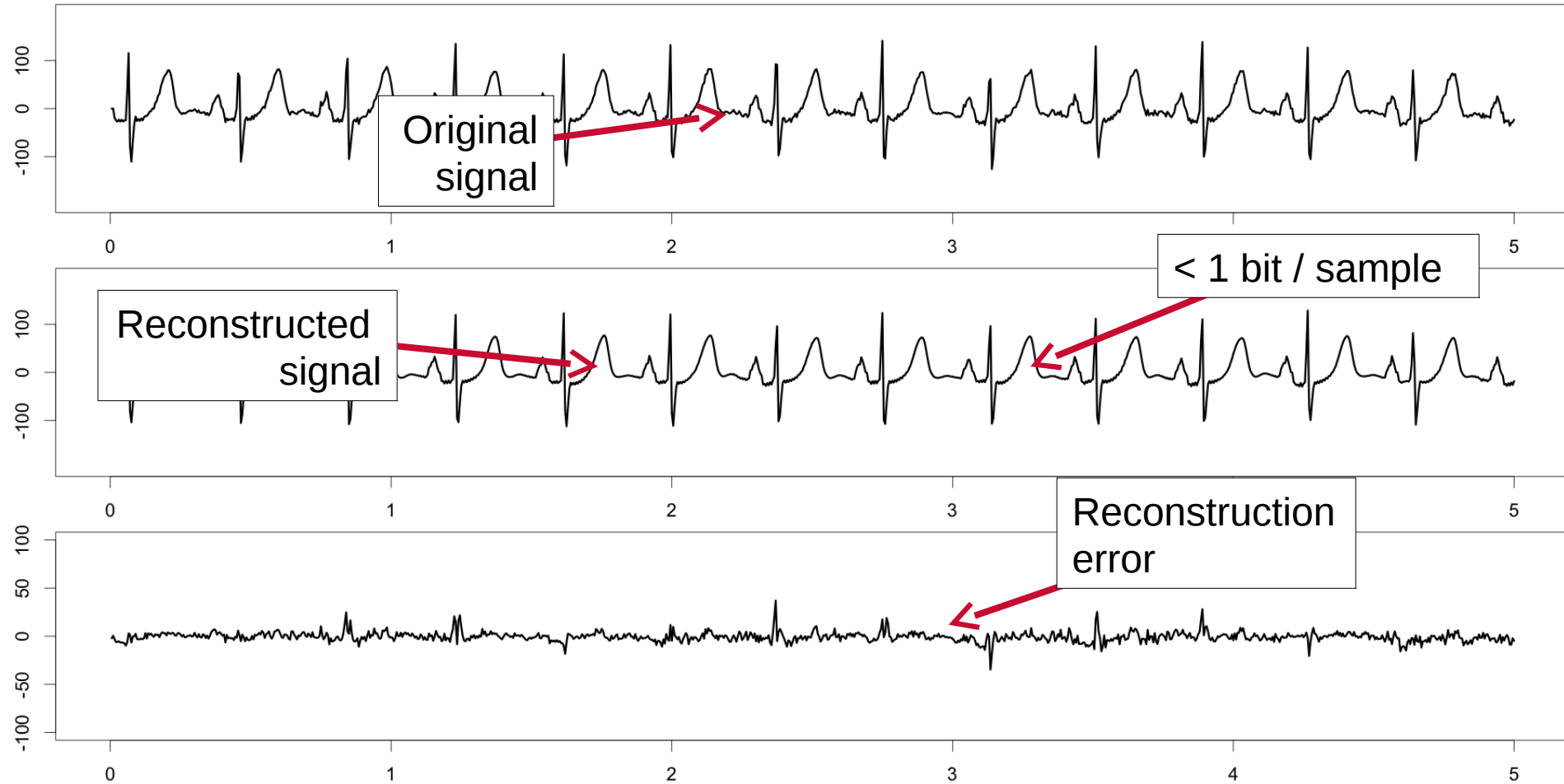
- The set of windowed signals is a nice model of our original signal
- Clustering can find the prototypes
  - Fancier techniques available using sparse coding
- The result is a dictionary of shapes
- New signals can be encoded by shifting, scaling and adding shapes from the dictionary



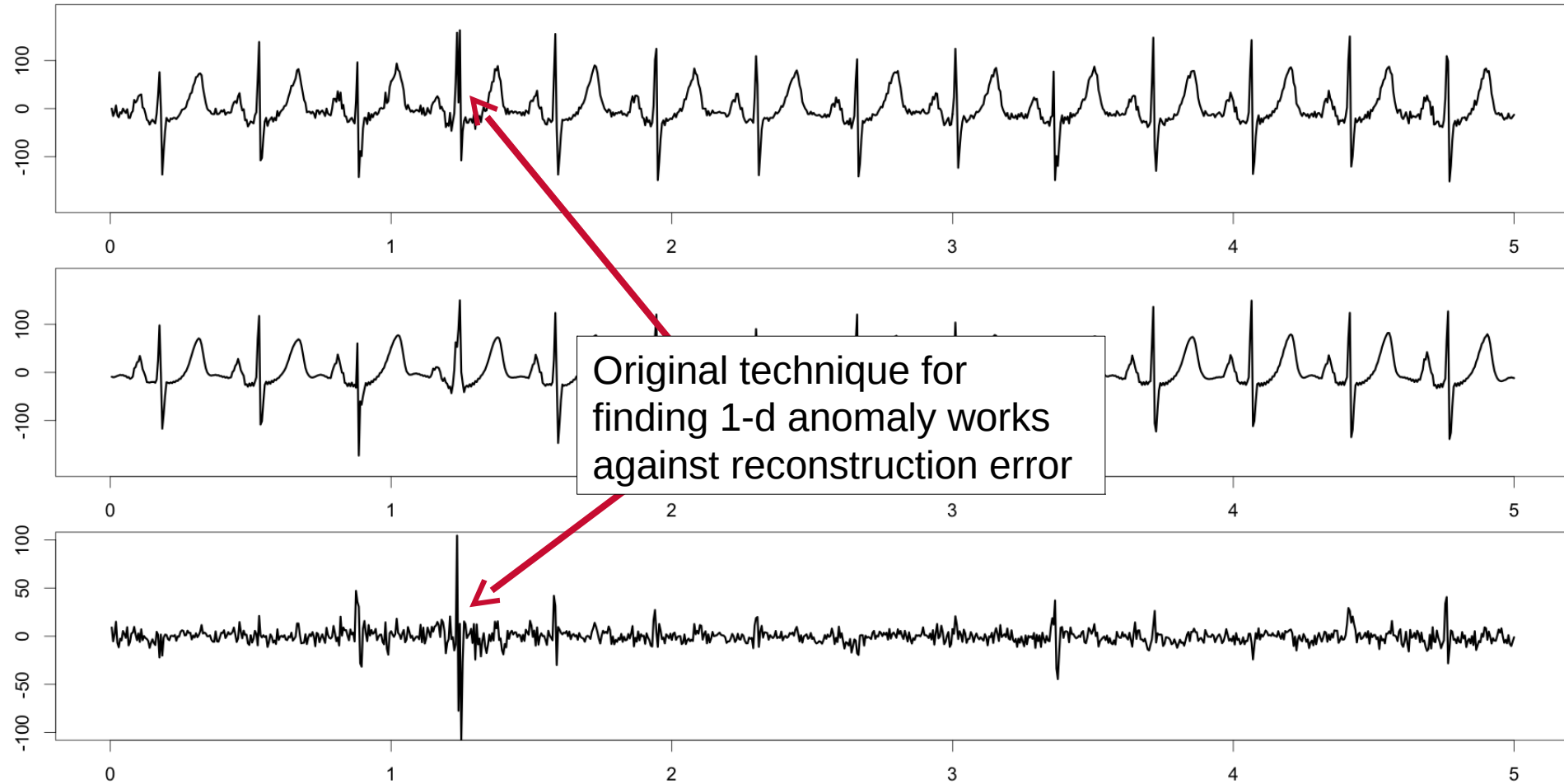
# Most Common Shapes (for EKG)



# Reconstructed signal



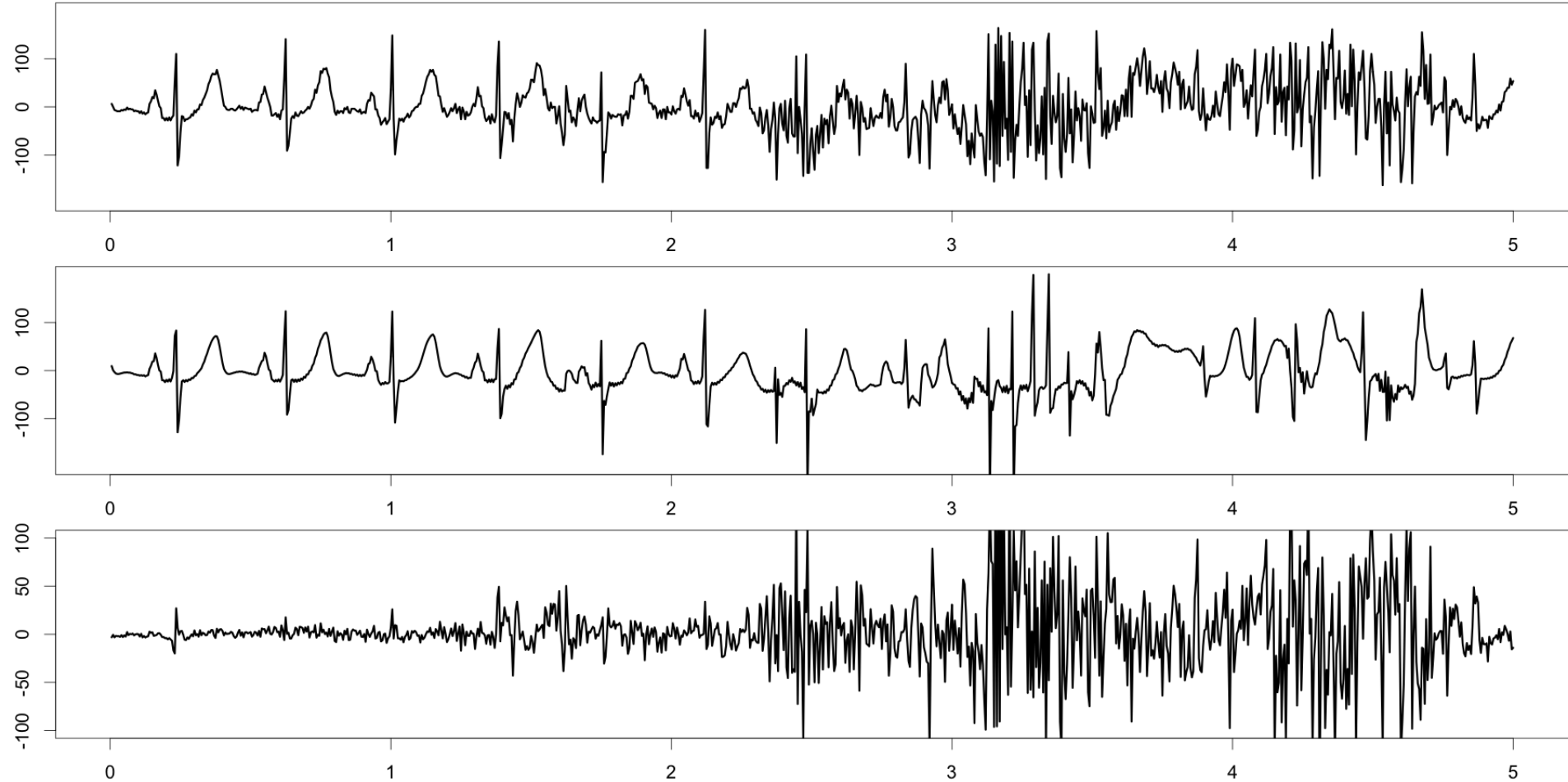
# An Anomaly



# Close-up of anomaly

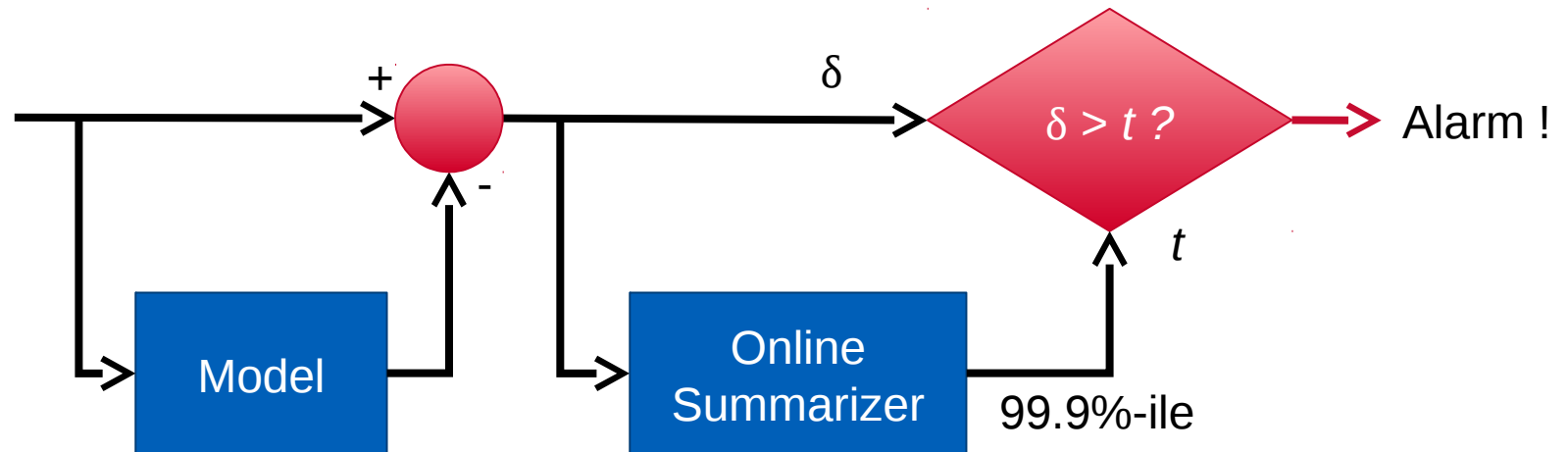


# A Different Kind of Anomaly





# Model Delta Anomaly Detection



# The Real Inside Scoop

- The model-delta anomaly detector is really just a sum of random variables
  - the model we know about already
  - and a normally distributed error
- The output (delta) is (roughly) the log probability of the sum distribution (really  $\delta^2$ )
- Thinking about probability distributions is good



# Example: Event Stream (timing)

- Events of various types arrive at irregular intervals
  - we can assume Poisson distribution
- The key question is whether frequency has changed relative to expected values
- Want alert as soon as possible



# Poisson Distribution

- Time between events is exponentially distributed

$$\Delta t \sim \lambda e^{-\lambda t}$$

- This means that long delays are exponentially rare

$$P(\Delta t > T) = e^{-\lambda T}$$
$$-\log P(\Delta t > T) = \lambda T$$

- If we know  $\lambda$  we can select a good threshold
  - or we can pick a threshold empirically



# Recap (out of order)

- Anomaly detection is best done with a probability model
- $-\log p$  is a good way to convert to anomaly measure
- Adaptive quantile estimation works for auto-setting thresholds



# Recap

- Different systems require different models
- Continuous time-series
  - sparse coding to build signal model
- Events in time
  - rate model base on variable rate Poisson
  - segregated rate model
- Events with labels
  - language modeling
  - hidden Markov models



# But Wait! Compression is Truth

- Maximizing  $\log \pi_k$  is minimizing compressed size
  - (each symbol takes  $-\log \pi_k$  bits on average)
- Maximizing  $\log \pi_k$  happens where  $\pi_k = p_k$ 
  - (maximum likelihood principle)



# But Auto-encoders Find Max Likelihood

- Minimal error  $\Rightarrow$  maximum likelihood
- Maximum likelihood  $\Rightarrow$  maximum compression
- So good anomaly detectors give good compression





# In Case You Want the Details

$$E[ p_k \log \pi_k ] = \sum_k p_k \log \pi_k$$

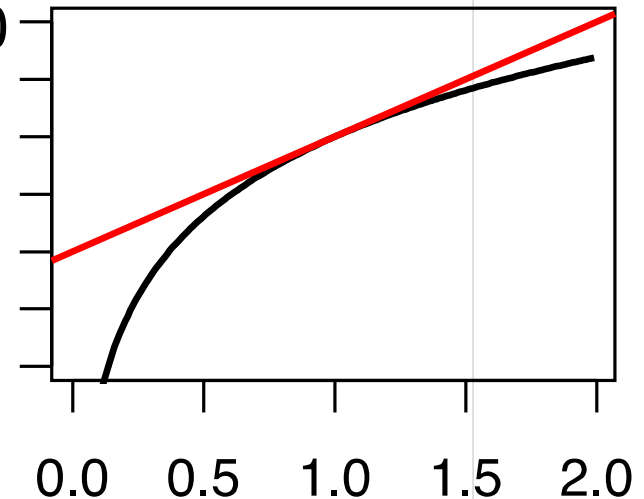
$$\log x \leq x - 1$$

$$\sum_k p_k \log \frac{\pi_k}{p_k} \leq \sum_k p_k \left( 1 - \frac{\pi_k}{p_k} \right) = \sum_k p_k - \sum_k \pi_k = 0$$

$$\sum_k p_k \log \pi_k - \sum_k p_k \log p_k \leq 0$$

$$\sum_k p_k \log \pi_k \leq \sum_k p_k \log p_k$$

$$E[ p_k \log \pi_k ] \approx \frac{1}{n} \sum_i \log \pi_{x_i}$$

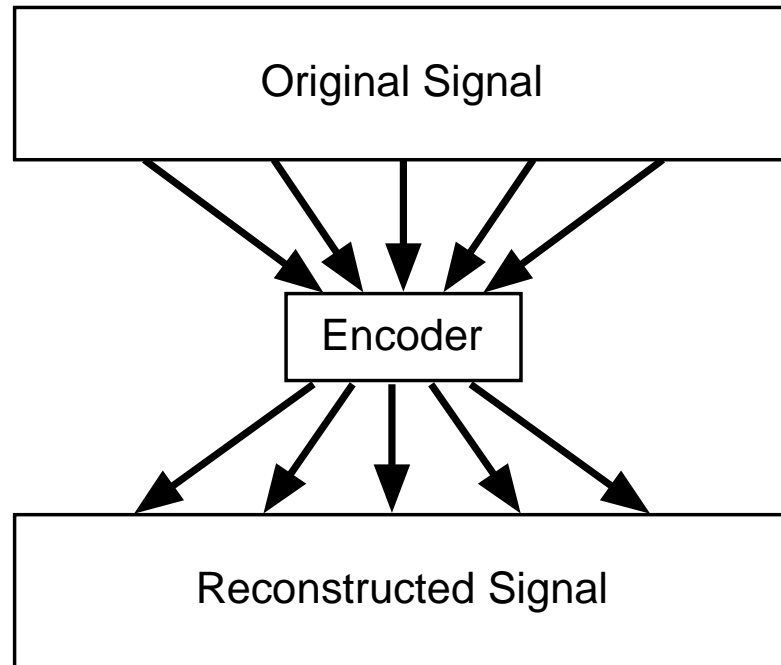


# Pause To Reflect on Clustering

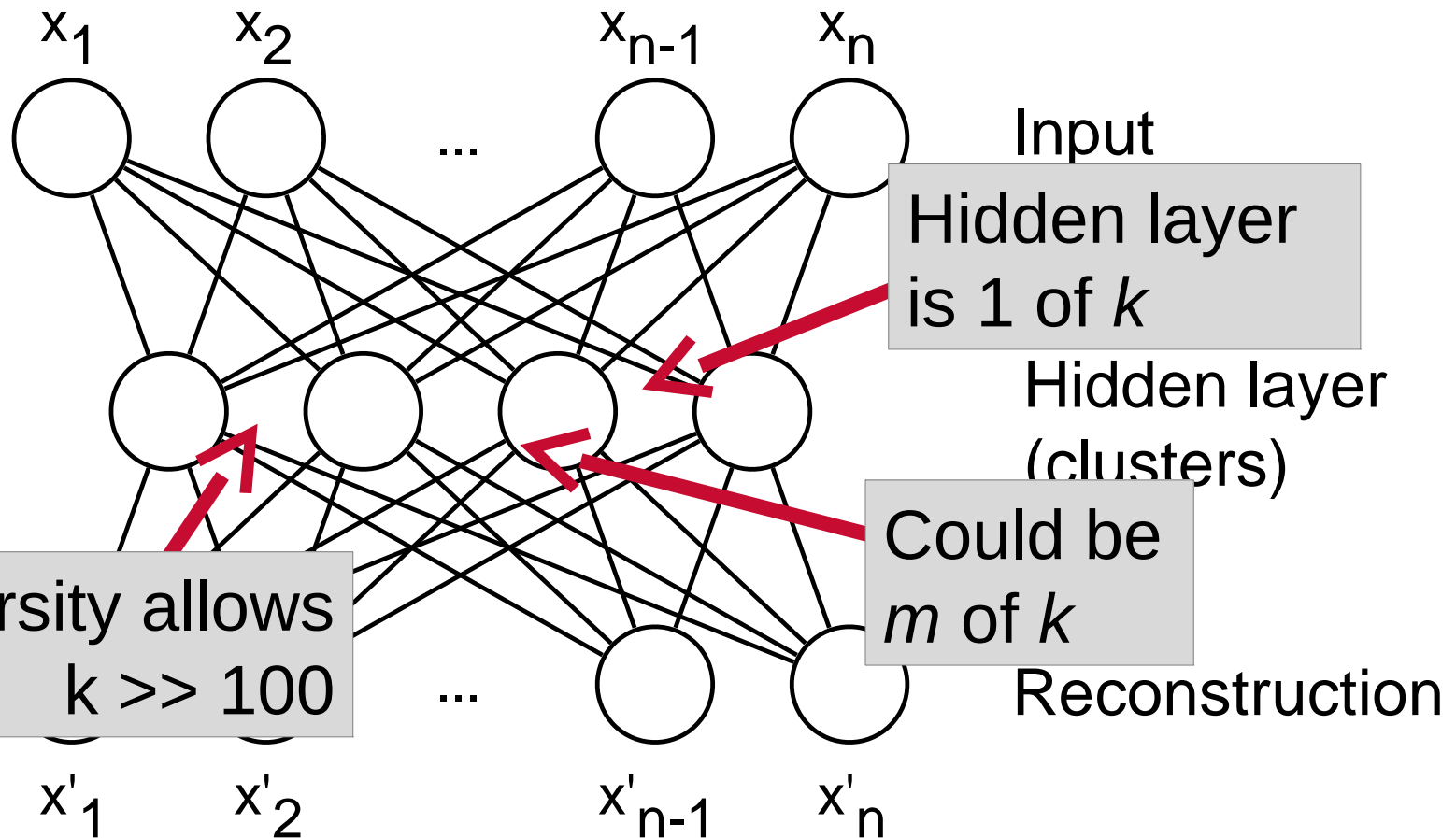
- Use windowing to apportion signal
  - Hamming windows add up to 1
- Find nearest cluster for each window
  - Can use dot product because all clusters normalized
- Scale cluster to right size
  - Dot product again
- Subtract from original signal



# Auto-encoding - Information Bottleneck

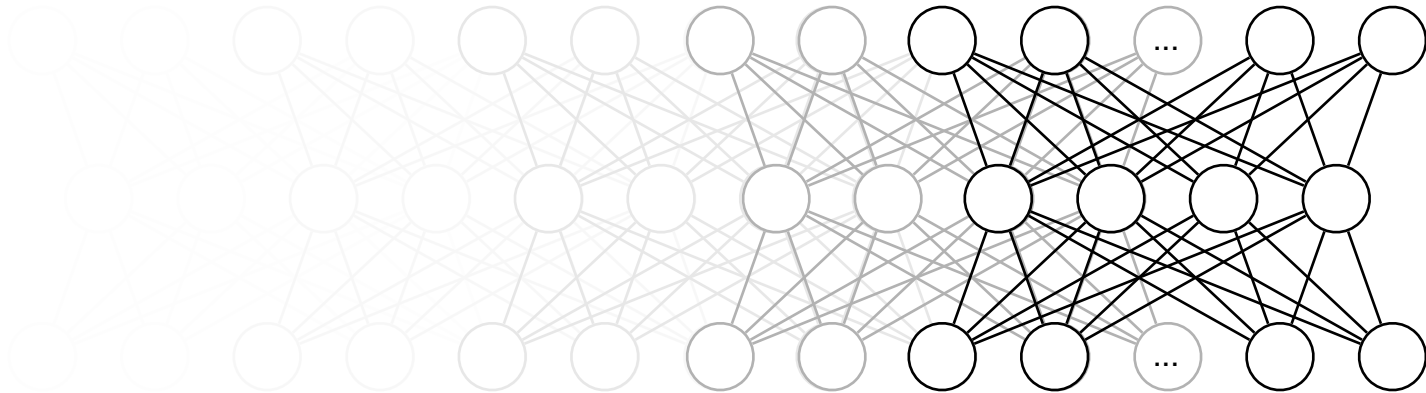


# Clustering as Neural Network



# Overlapping Networks

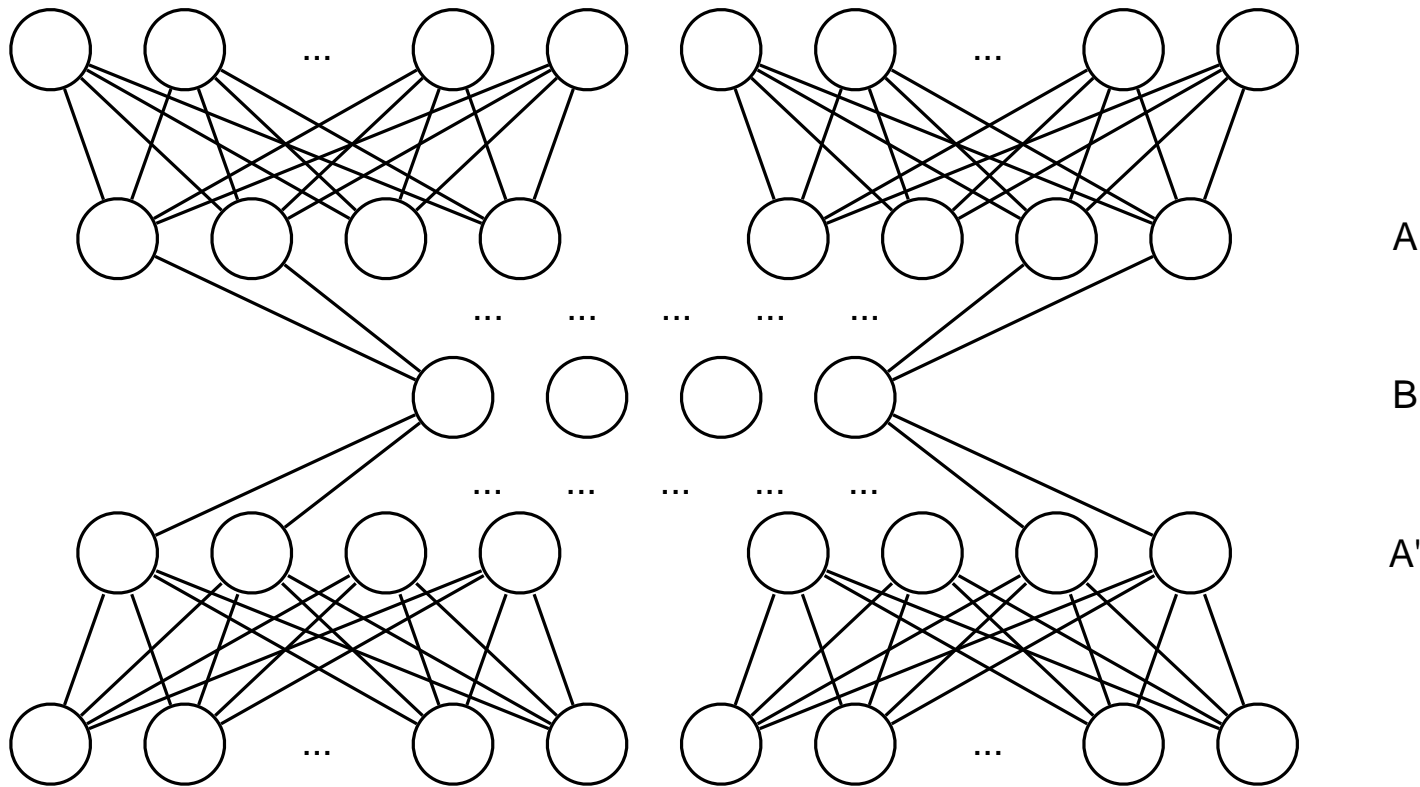
Time series input



Reconstructed time series



# Deep Learning



# What About the Database?

- We don't have to keep the reconstruction
- We can keep the first level nodes
  - And the reconstruction error
- To keep the first level nodes
  - We can keep the second level nodes
  - Plus the reconstruction error



# What Does it Matter?

- Even one level of auto-encoding compresses
  - 30-50x in EKG example with k-means
- Multiple levels compress more
  - Understanding => Truth => Compression
- Higher levels give semantic search



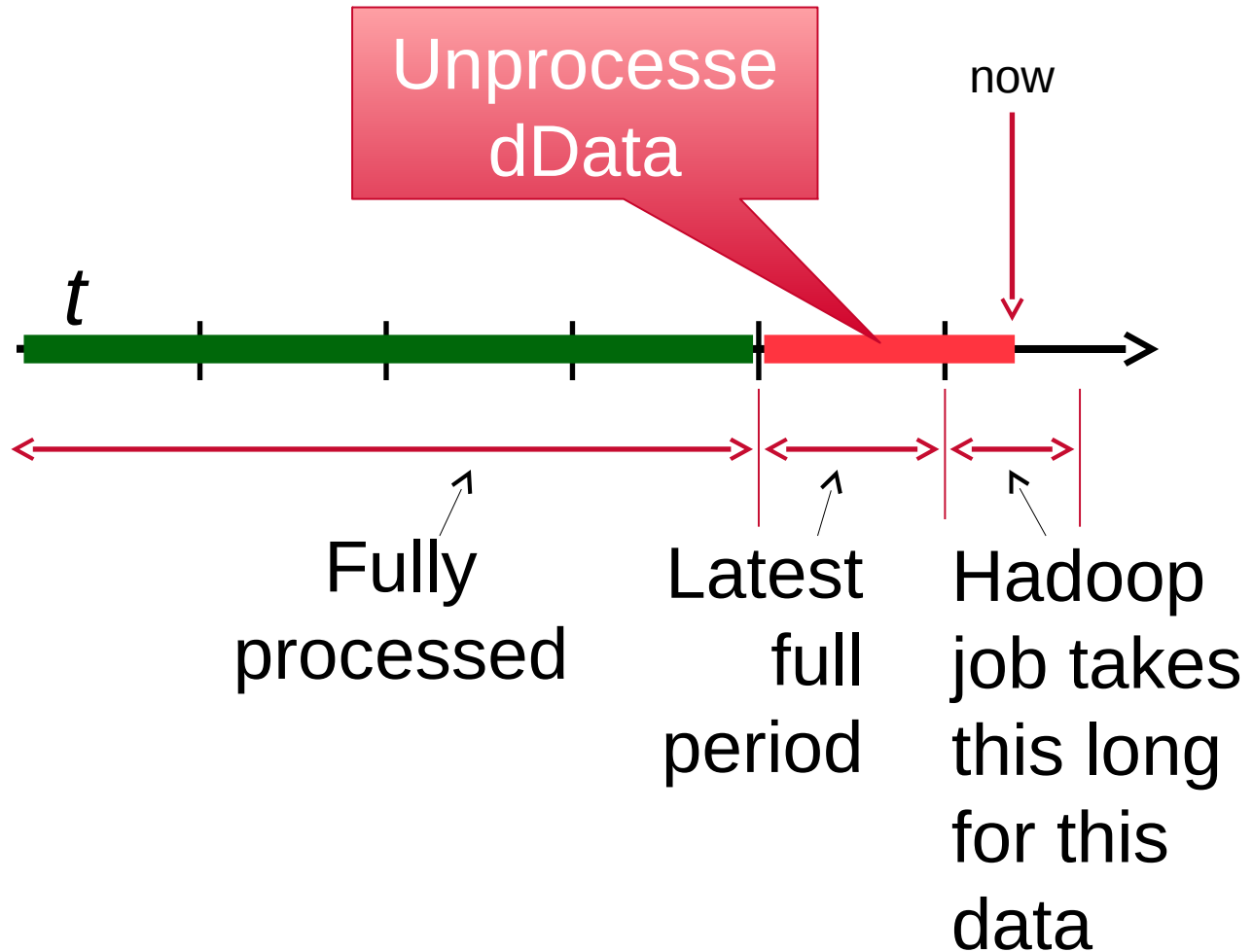


# How Do I Build Such a System

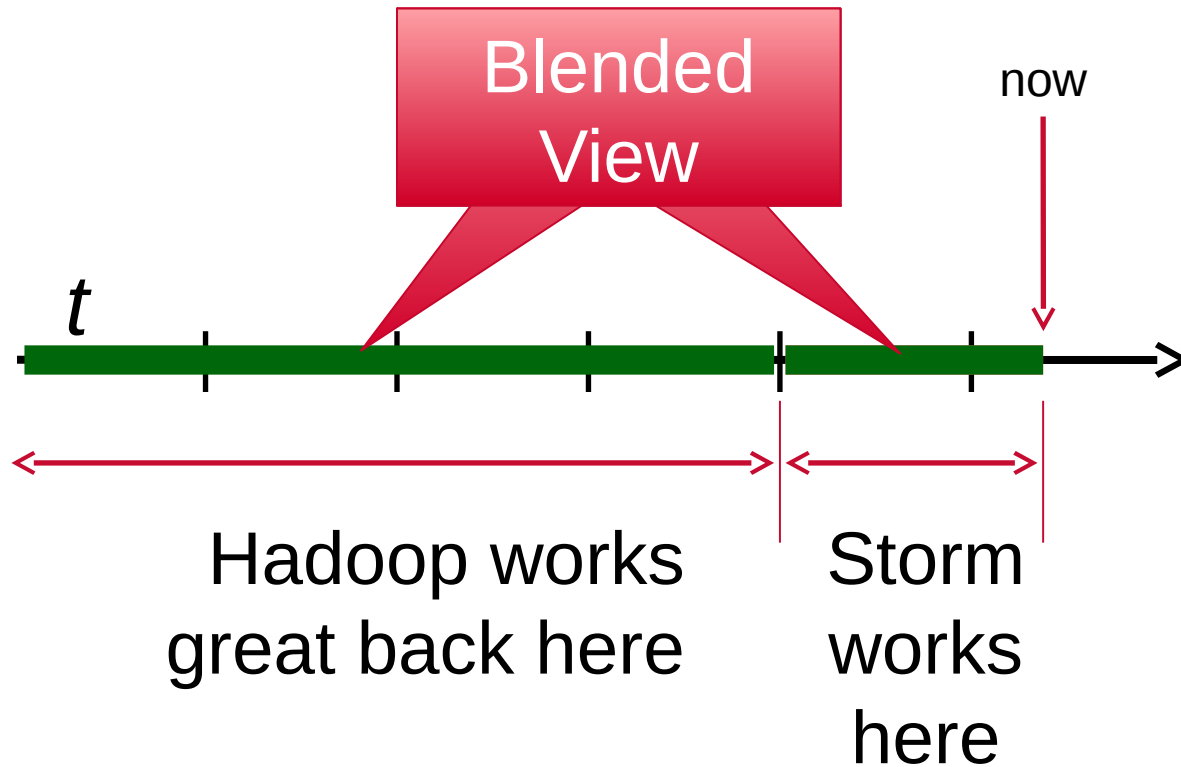
- The key is to combine real-time and long-time
  - real-time evaluates data stream against model
  - long-time is how we build the model
- Extended Lambda architecture is my favorite
- See my other talks on [slideshare.net](http://slideshare.net) for info
- Ping me directly



# Hadoop is Not Very Real-time



# Real-time and Long-time together



# Who I am

- Ted Dunning, Chief Application Architect, MapR  
[tdunning@mapr.com](mailto:tdunning@mapr.com)  
[tdunning@apache.org](mailto:tdunning@apache.org)  
@ted\_dunning
- Committer, mentor, champion, PMC member on several Apache projects
- Mahout, Drill, Zookeeper others



