

# Effective management of high volume numeric data with histograms

Fred Moyer @Circonus  
DataEngConf SF '18



[@phredmoyer](#)



- Engineer to Engineer @circonus
- Recovering C and Perl programmer
- Geeking out on histograms since 2015

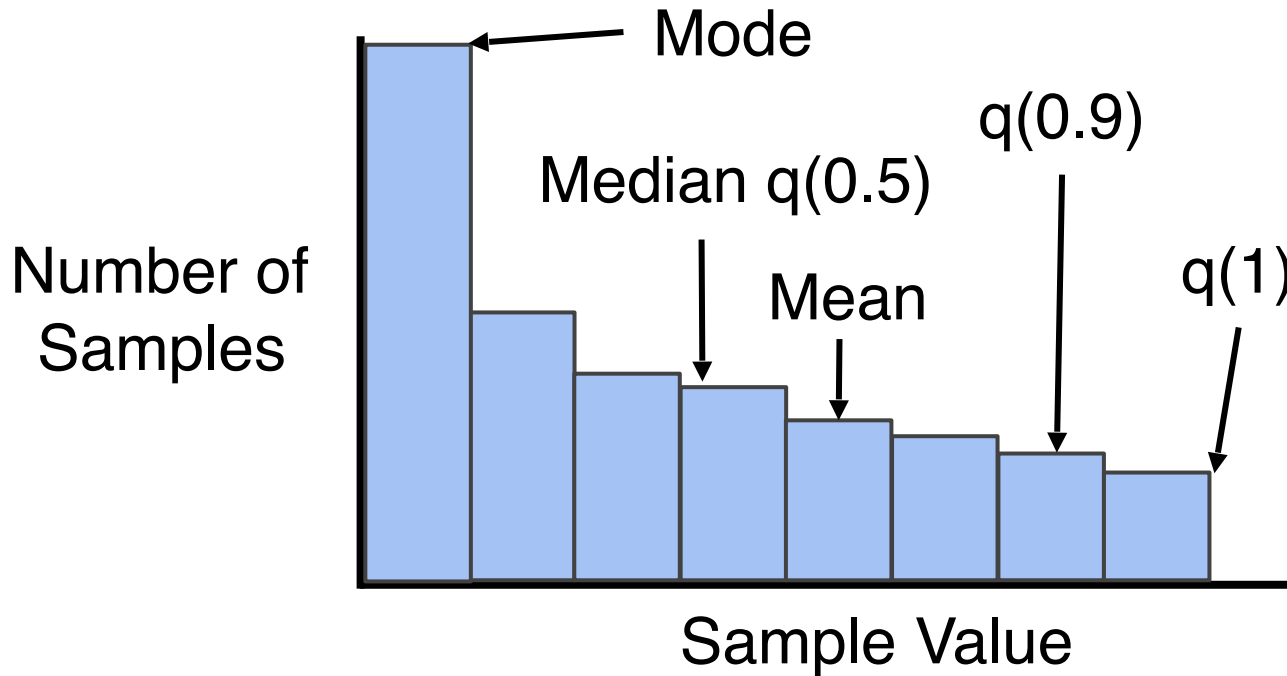
## Pain driven development

- Observability tools caused a telemetry firehose
- Existing monitoring systems got washed away
- Average based metrics gave limited insight

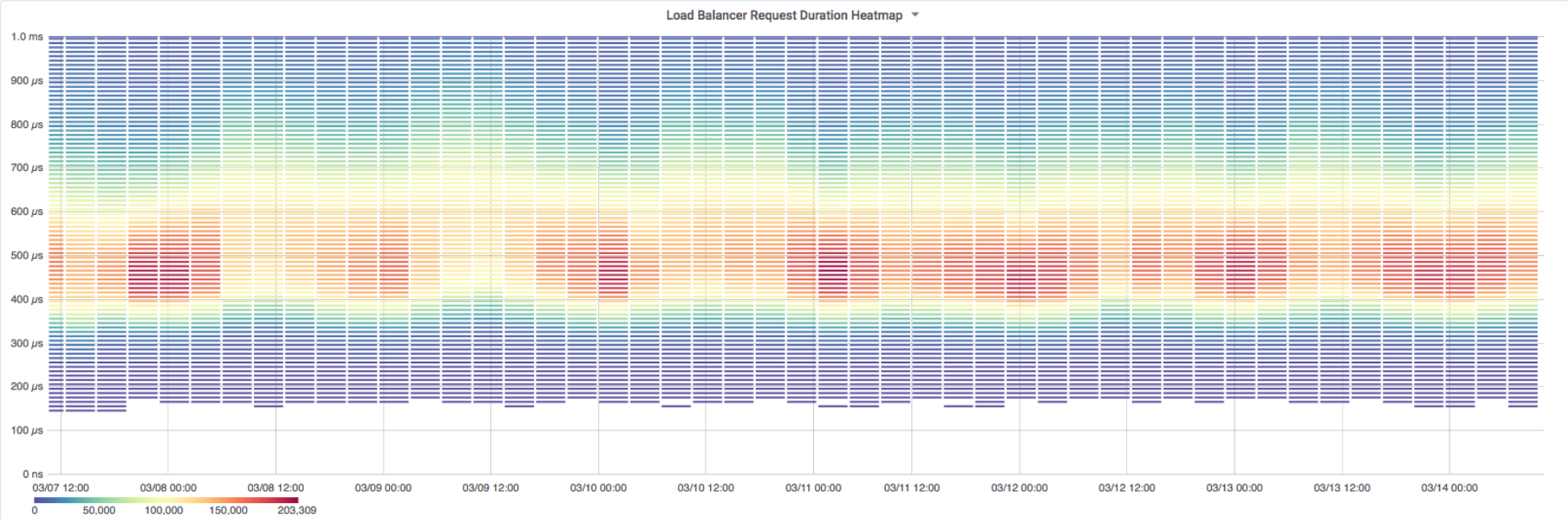
# “Effective Management”

- Performance AND scalability
- Avoid memory allocations, copies, locks, waits
- Persist data in size efficient structures

# Histogram Basics



# Heatmap Basics



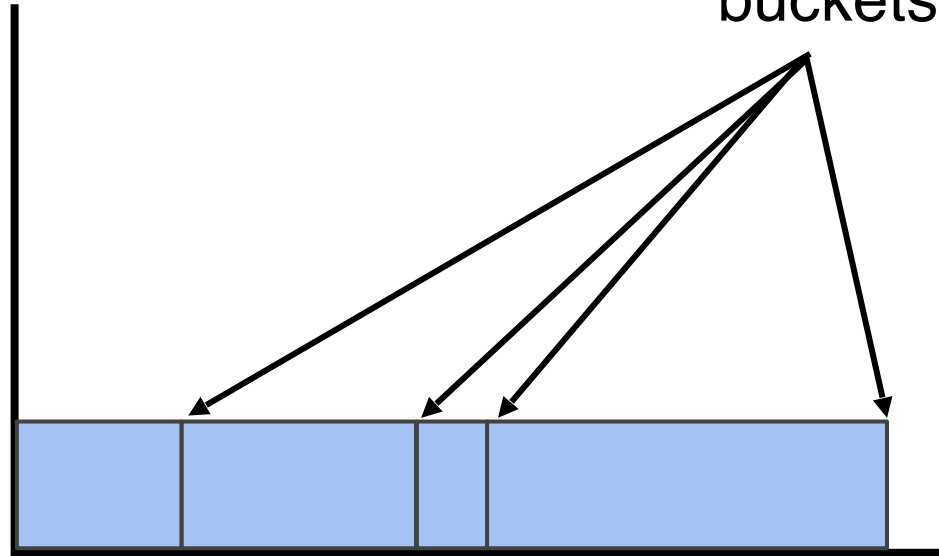
# Histogram Types

- Fixed Bucket
- Approximate
- Linear
- Log Linear
- Cumulative

# Fixed Bucket

User specified bins/  
buckets

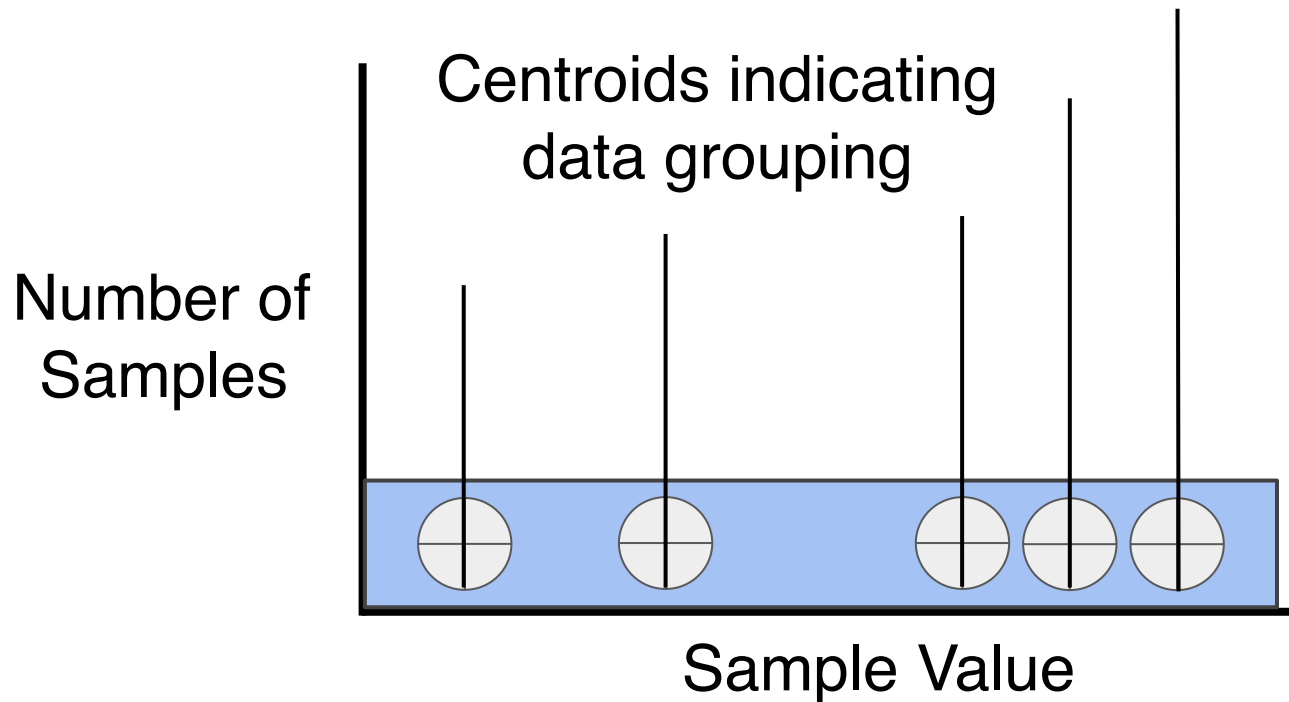
Number of  
Samples



Sample Value



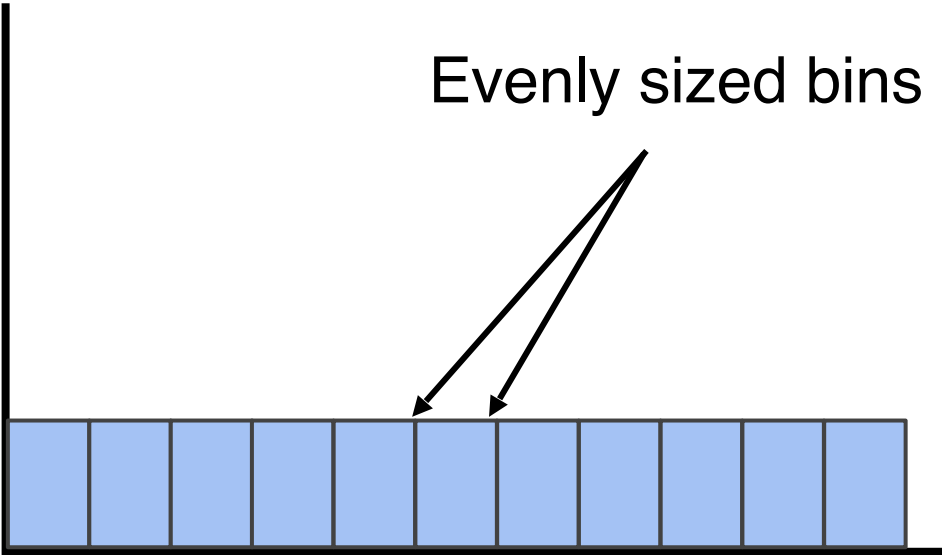
# Approximate



# Linear

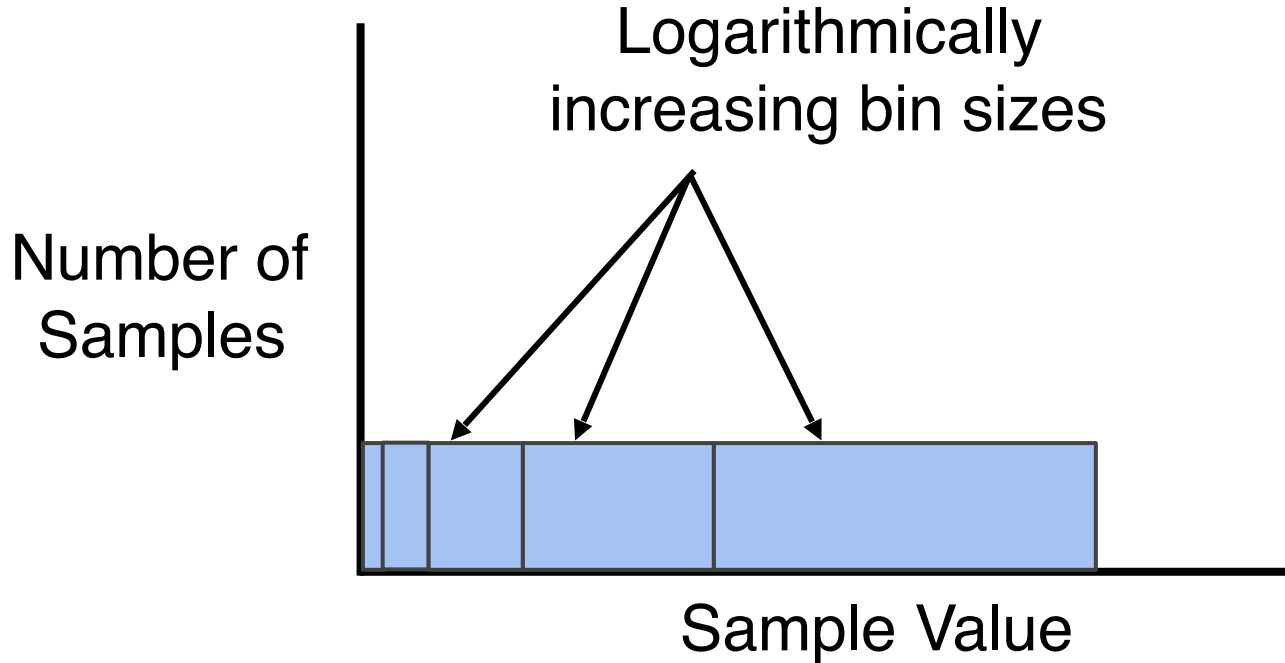


Number of Samples



Sample Value

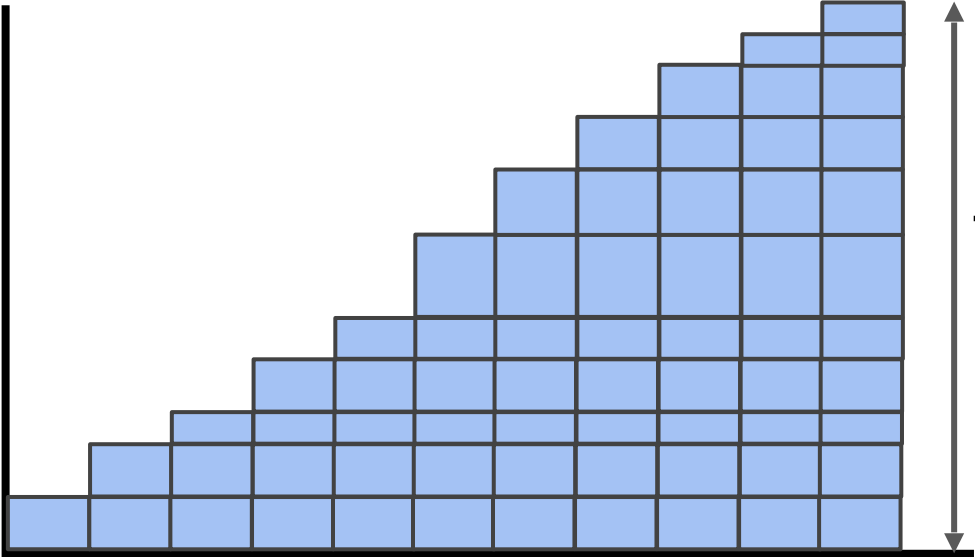
# Log Linear



# Cumulative



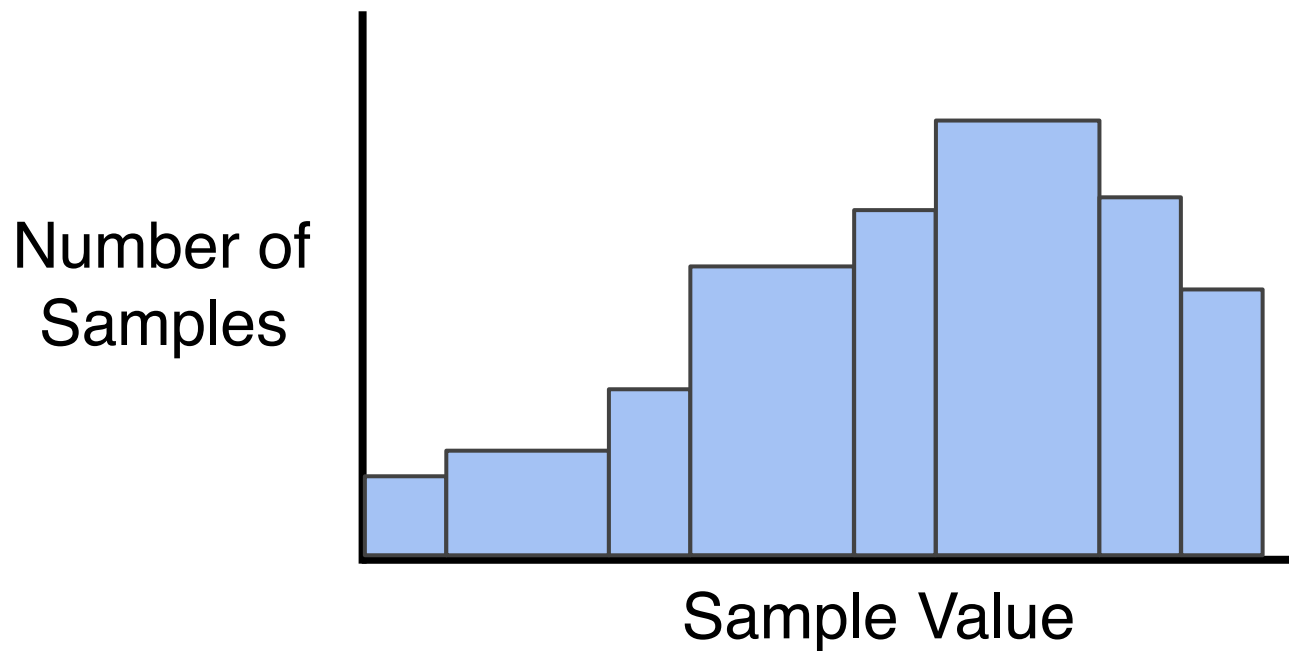
Number of Samples



Total Sample Count

Sample Value

Custom



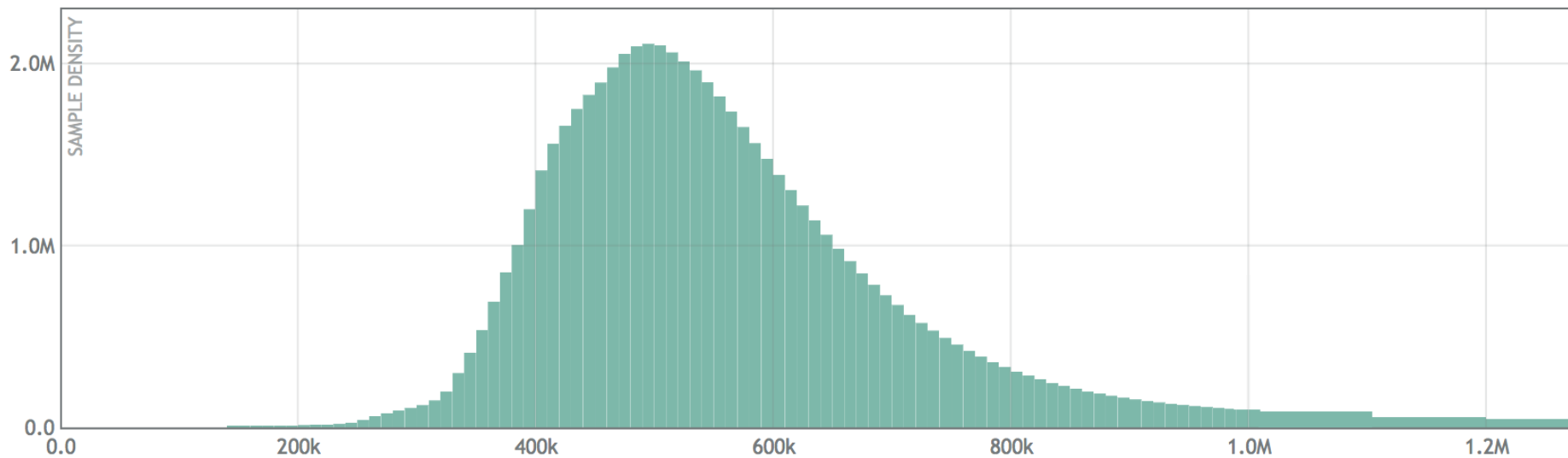
# Open Source Log Linear



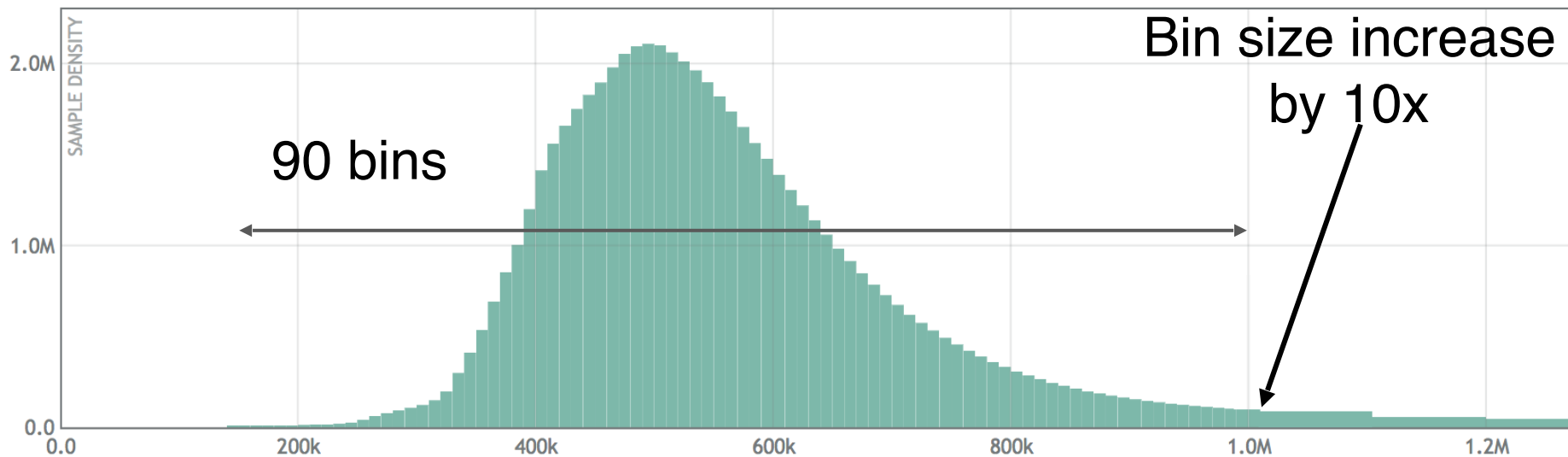
C - [github.com/circonus-labs/libcirclhist](https://github.com/circonus-labs/libcirclhist)

Go - [github.com/circonus-labs/circonuslhist](https://github.com/circonus-labs/circonuslhist)

# Open Source Log Linear

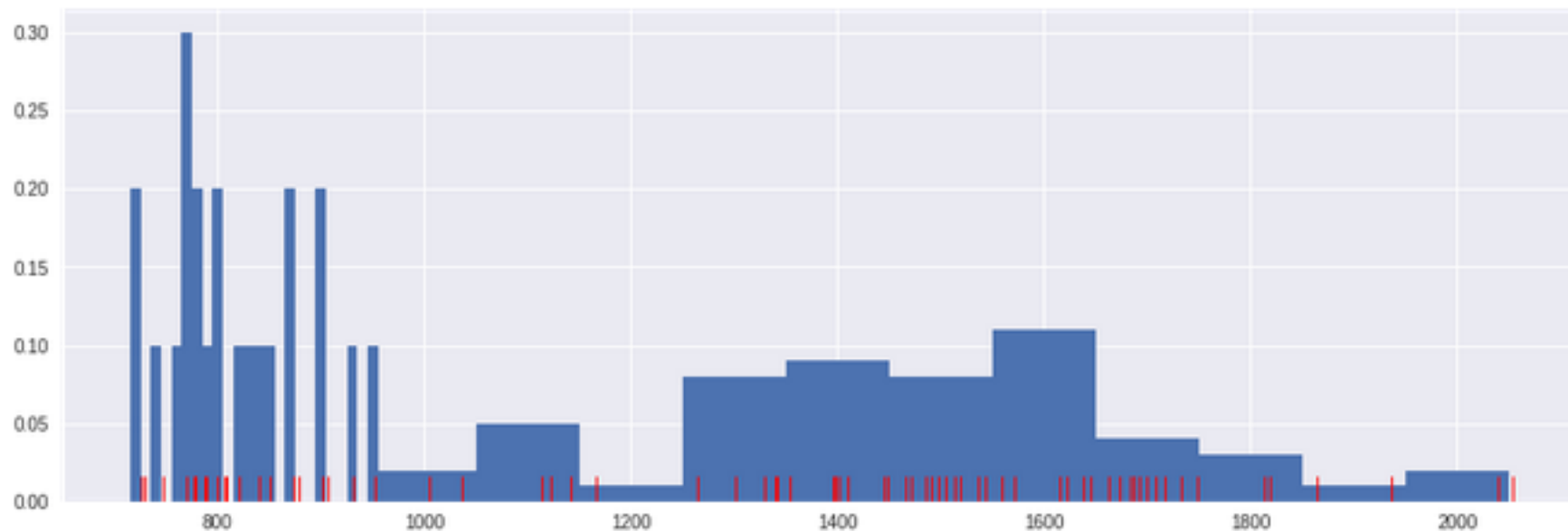


# Open Source Log Linear

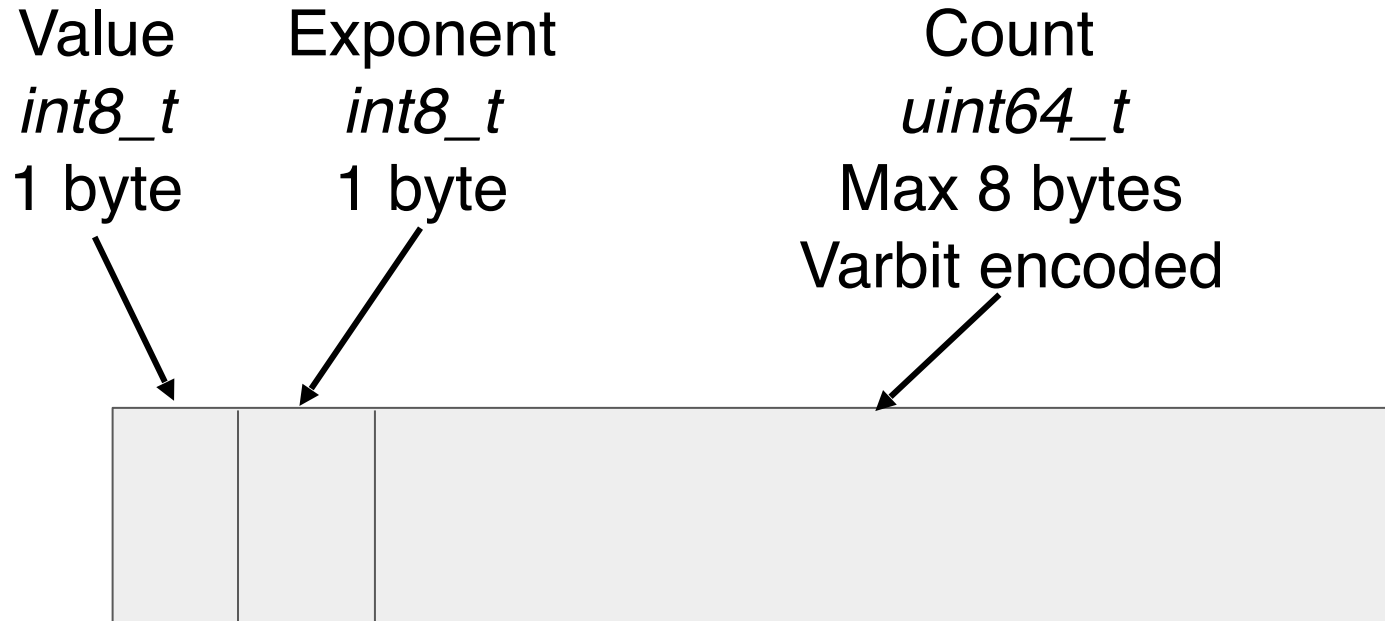




# Open Source Log Linear



# Bin data structure



Storage efficiency - 1 month

30 days of one minute histograms

$30 \text{ days} * 24 \text{ hours/day} * 60 \text{ bins/hour} * 300 \text{ bin span} * 10 \text{ bytes/bin} * 1 \text{ kB}/1,024 \text{ bytes} * 1 \text{ MB}/1024 \text{ kB} = 123.6 \text{ MB}$

Storage efficiency - 1 year

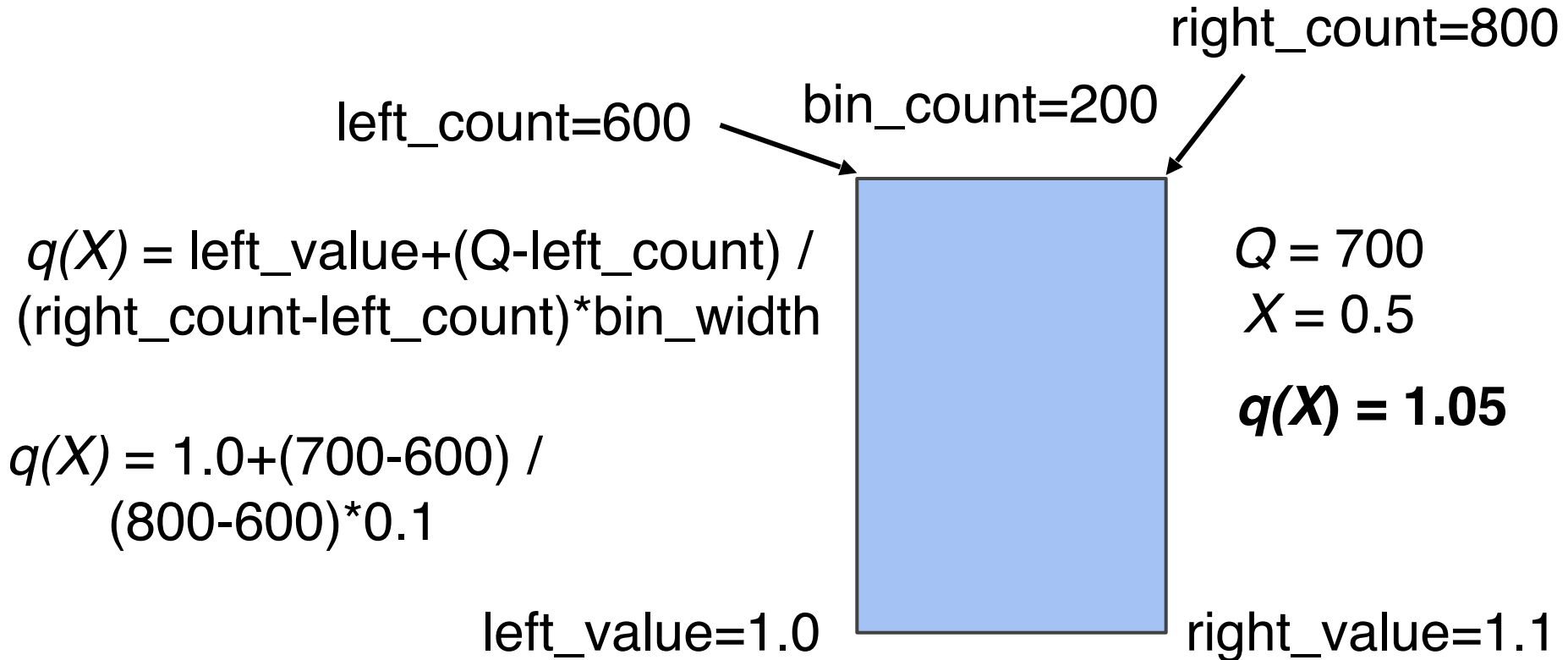
365 days of five minute histograms

$365 \text{ days} * 24 \text{ hours/day} * 12 \text{ bins/hour} * 300 \text{ bin span} * 10 \text{ bytes/bin} * 1\text{kB}/1,024\text{bytes} * 1\text{MB}/1024\text{kB} = 300.9 \text{ MB}$

# Quantile calculation

1. Given a quantile  $q(X)$  where  $0 < X < 1$
2. Sum up the counts of all the bins,  $C$
3. Multiply  $X * C$  to get count  $Q$
4. Walk bins, sum bin boundary counts until  $> Q$
5. Interpolate quantile value  $q(X)$  from bin

# Linear interpolation



$$q(X) = \text{left\_value} + (\text{Q} - \text{left\_count}) / (\text{right\_count} - \text{left\_count}) * \text{bin\_width}$$

$$q(X) = 1.0 + (700 - 600) / (800 - 600) * 0.1$$

$$Q = 700$$
$$X = 0.5$$

$$q(X) = 1.05$$

# Recap



- Several different types of histograms
- Highly space efficient
- $O(1)$  and  $O(n)$  complexity calculating quantiles
- What other fun things can we do?

# Inverse Quantiles

- What's the 95th percentile latency?
  - $q(0.95) = 10\text{ms}$
- What percent of requests exceeded 10ms?
  - 5% for this data set; what about others?



# Inverse Quantile calculation

1. Given a sample value  $X$ , locate its bin
2. Using the previous linear interpolation equation, solve for  $Q$  given  $X$

## Inverse Quantile calculation

$$X = \text{left\_value} + (Q - \text{left\_count}) / (\text{right\_count} - \text{left\_count}) * \text{bin\_width}$$

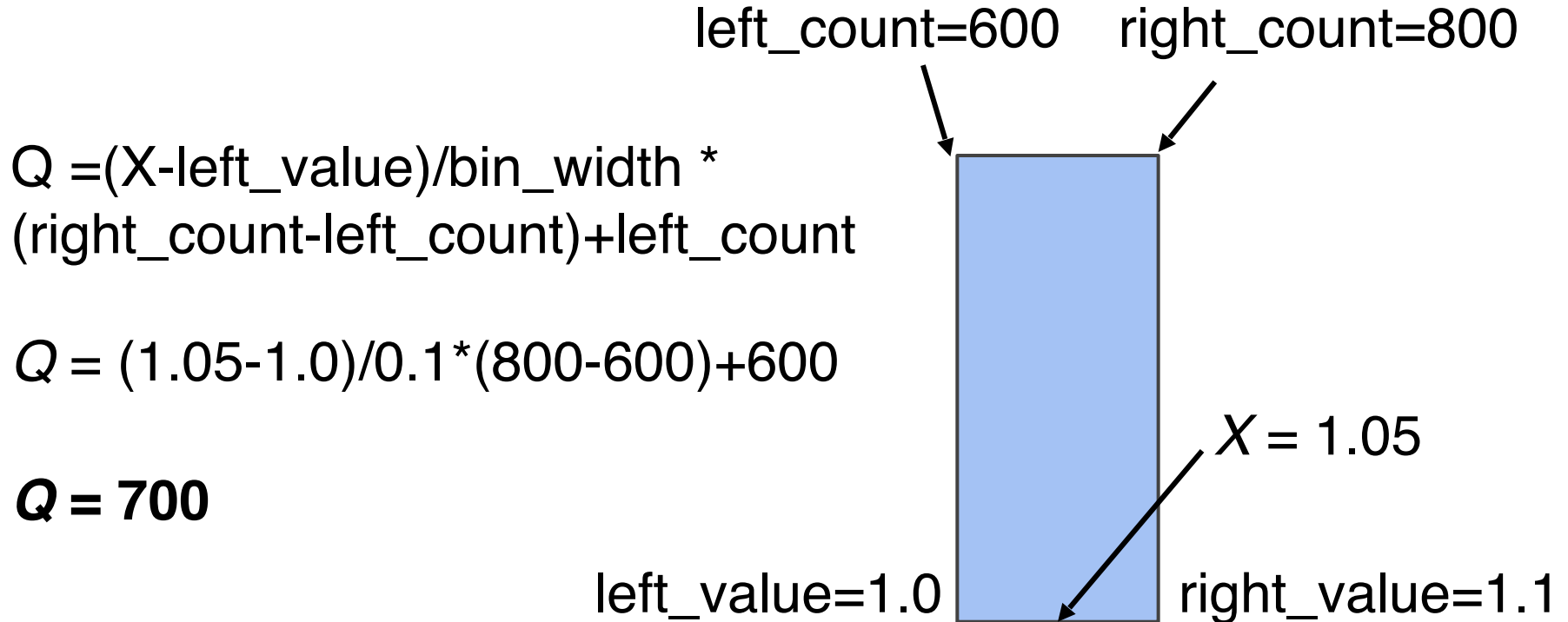
$$X - \text{left\_value} = (Q - \text{left\_count}) / (\text{right\_count} - \text{left\_count}) * \text{bin\_width}$$

$$(X - \text{left\_value}) / \text{bin\_width} = (Q - \text{left\_count}) / (\text{right\_count} - \text{left\_count})$$

$$(X - \text{left\_value}) / \text{bin\_width} * (\text{right\_count} - \text{left\_count}) = Q - \text{left\_count}$$

$$Q = (X - \text{left\_value}) / \text{bin\_width} * (\text{right\_count} - \text{left\_count}) + \text{left\_count}$$

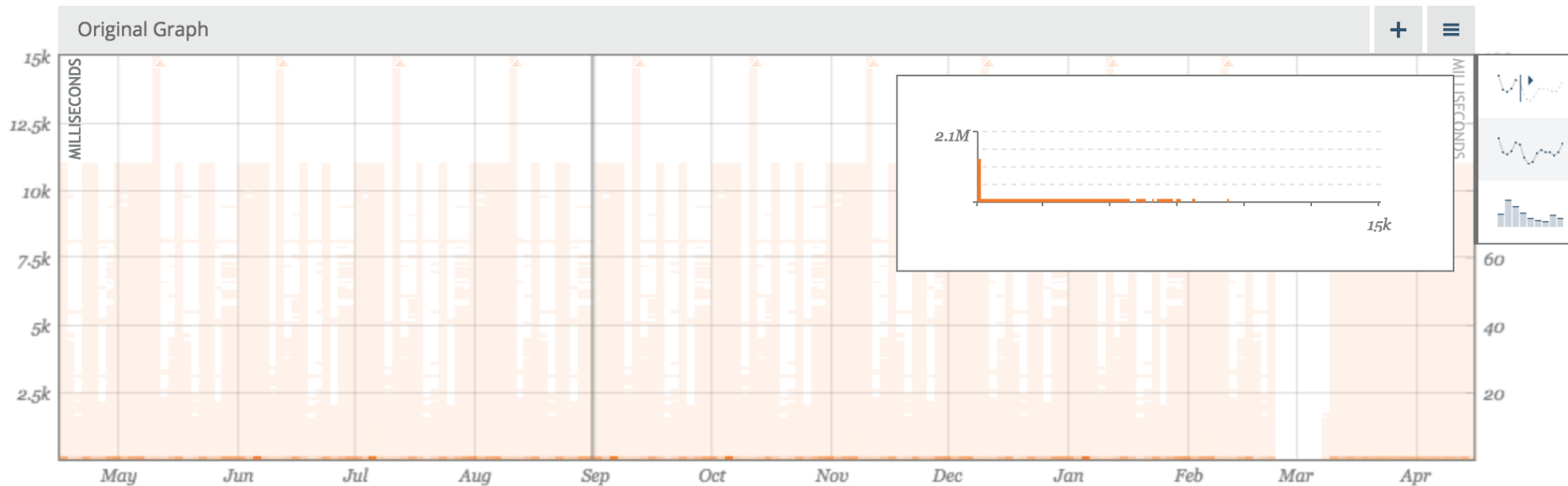
# Linear interpolation



# Inverse Quantile calculation

1. ~~Given a sample value  $X$ , locate its bin~~
2. ~~Using the previous linear interpolation equation, solve for  $Q$  given  $X$~~
3. Sum the bin counts up to  $Q$  as  $Q_{left}$
4. Inverse quantile  $q_{inv}(X) = (Q_{total} - Q_{left}) / Q_{total}$
5. For  $Q_{left} = 700$ ,  $Q_{total} = 1,000$ ,  $q_{inv}(X) = 0.3$
6. 30% of sample values exceeded  $X$

# Quantiles - Heatmap



Aug 29 2017, 20:00 (1d)



circonus-demo circonus.net json: api `GET` /getState (on demo-replay.circonus.net, from Chicago, IL, US) (ms)

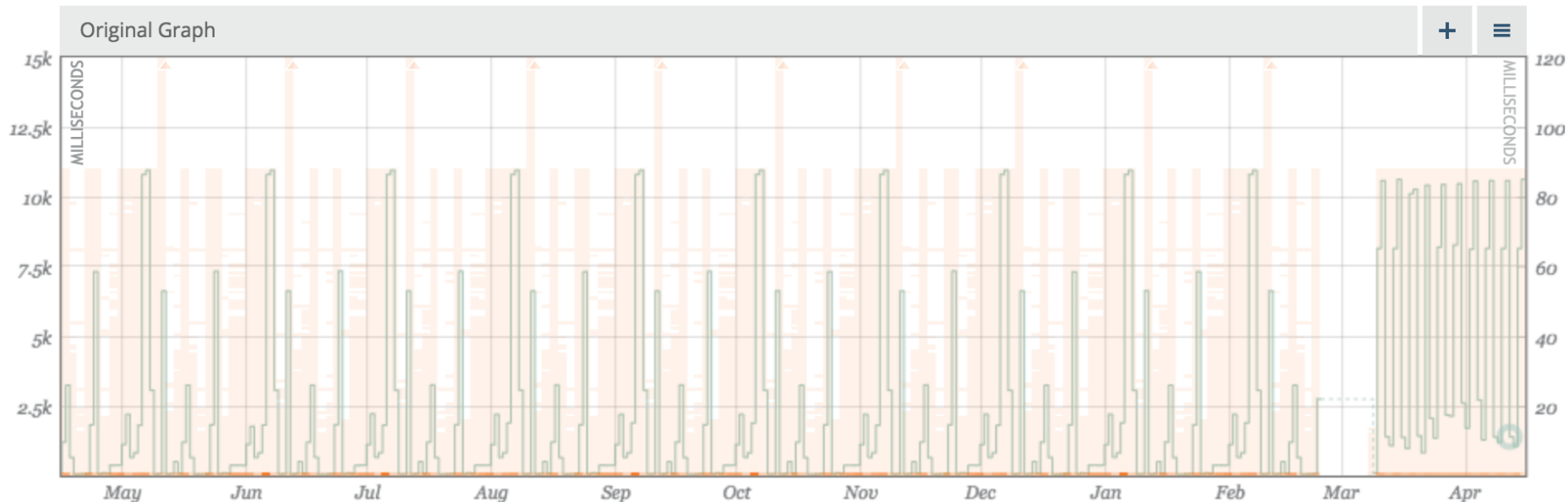
[4.3k - 4.4k] 6 of 1281020 samples

98%

1%

1%

# Quantiles - q(0.9)



Apr 11 2018, 20:00 (1d)

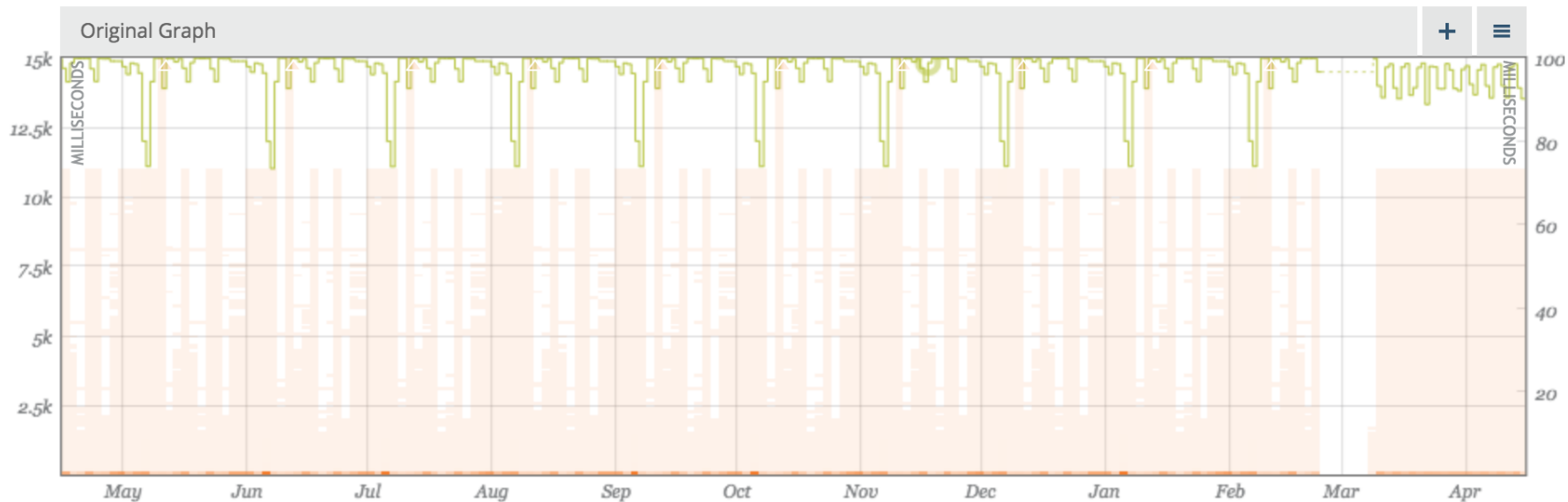
L R 99th percentile




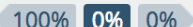
1.3913050314464999k

L R circonus-demo circonus.net json: api `GET` /getState (on demo-replay.circonus.net, from Chicago, IL, US)

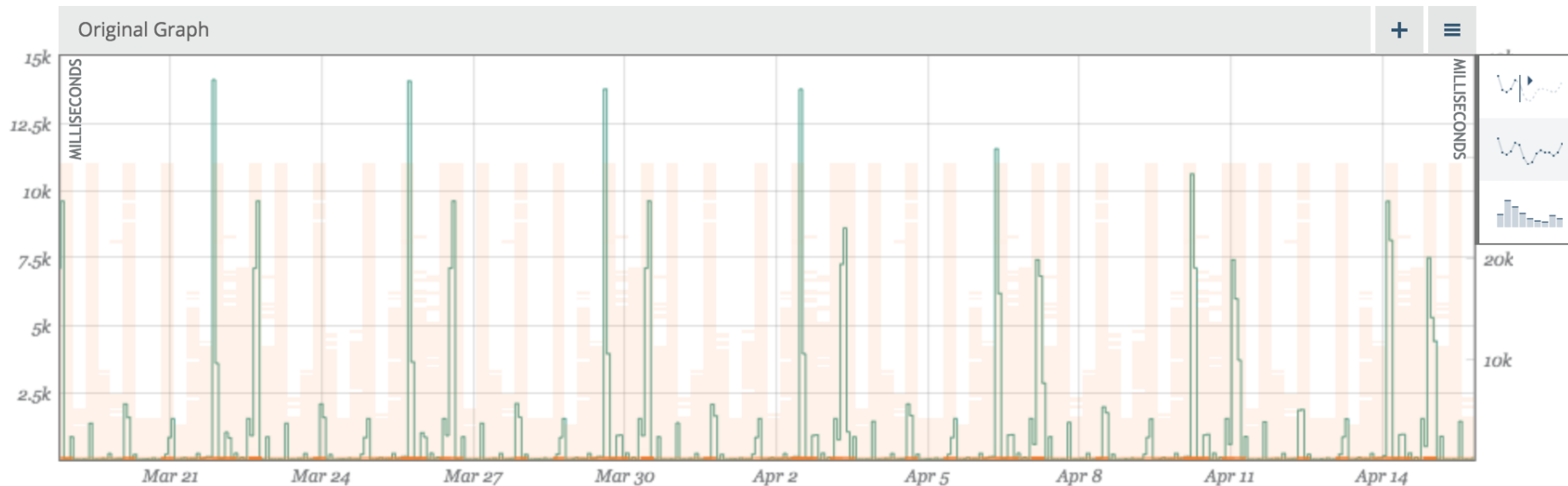
-

# Inverse Quantiles - SLO



 <i>Nov 17 2017, 19:00 (1d)</i>	
 Inverse quantile calculation for 500 milliseconds	98.783298034185
 circonus-demo circonus.net json: api `GET` /getState (on demo-replay.circonus.net, from Chicago, IL, US)	[12k - 13k] 0 of 361939 samples 

# Inverse Quantiles - SLO



Mar 18 2018, 18:30 (1h 30M)



Inverse quantile calculation count of requests that violated 500 millisecond SLO

18.959302325581k

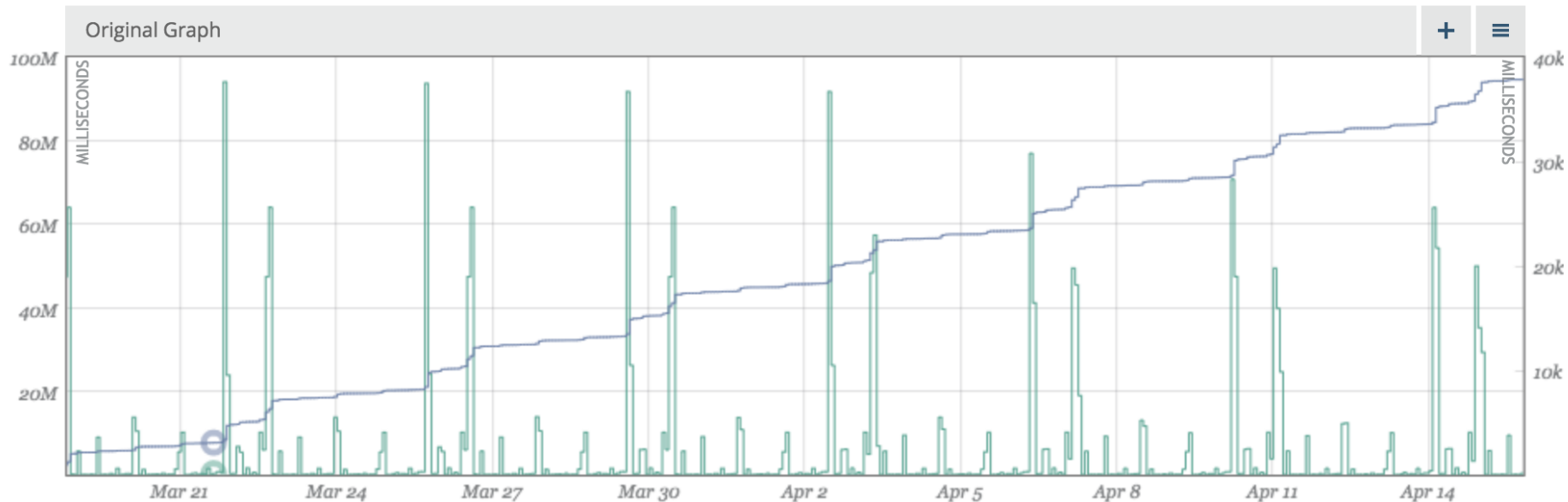


circonus-demo circonus.net json: api `GET` /getState (on demo-replay.circonus.net, from Chicago, IL, US)

-



# Inverse Quantiles - SLO



Mar 21 2018, 15:30 (1h 30M)

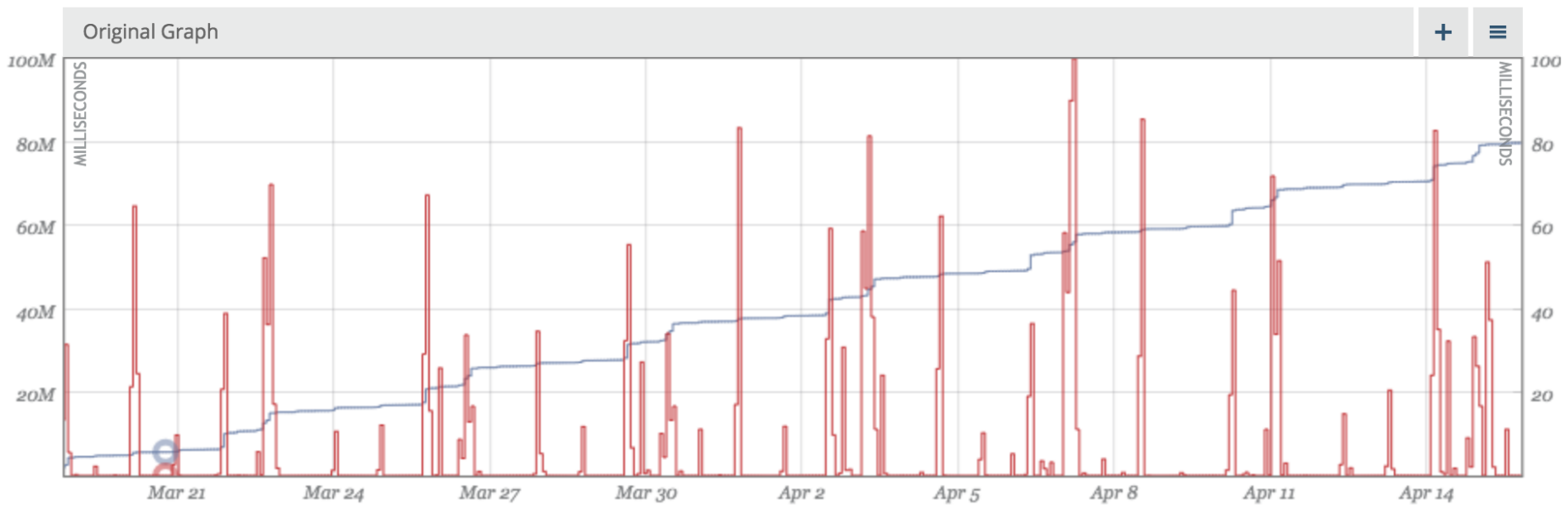
**L R** Inverse quantile calculation count of requests that violated 500 millisecond SLO

0.14487179487179k

**L R** Cumulative number of requests exceeded the 200 millisecond SLO

7.68593M

# Anomalies



		⚙️	Mar 20 2018, 18:30 (1h 30M)
L R	Cumulative number of requests exceeded the 500 millisecond SLO		5.663712222222222M
L R	SLO violation anomaly detection		0

Thank you!



Questions?

Bug me at the Circonus booth

Come to Office Hours

Tweet @phredmoyer or @circonus