

Anomaly Detection

Fault Tolerance

Anticipation

Patterns

John Allspaw
SVP, Tech Ops
Qcon London 2012

Etsy

Four Cornerstones

Erik Hollnagel

(Anticipation)

Knowing
What
To Expect

(Response)

Knowing
What
To Look For

Knowing
What
To Do

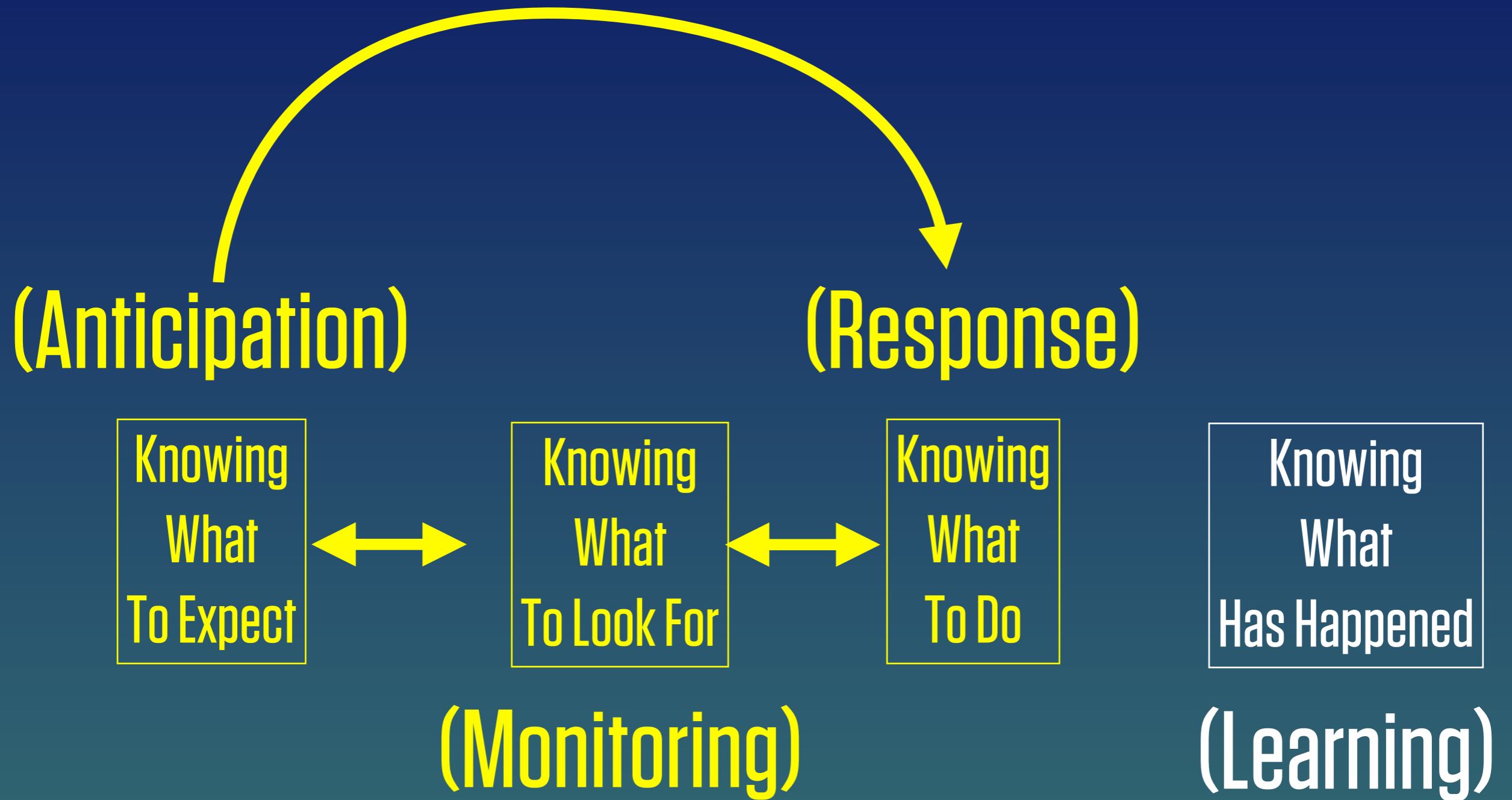
Knowing
What
Has Happened

(Monitoring)

(Learning)

Four Cornerstones

Erik Hollnagel



Four Cornerstones

Erik Hollnagel

(Anticipation)

Knowing
What
To Expect

(Response)

Knowing
What
To Do

(Monitoring)

Knowing
What
To Look For

(Learning)

Knowing
What
Has Happened

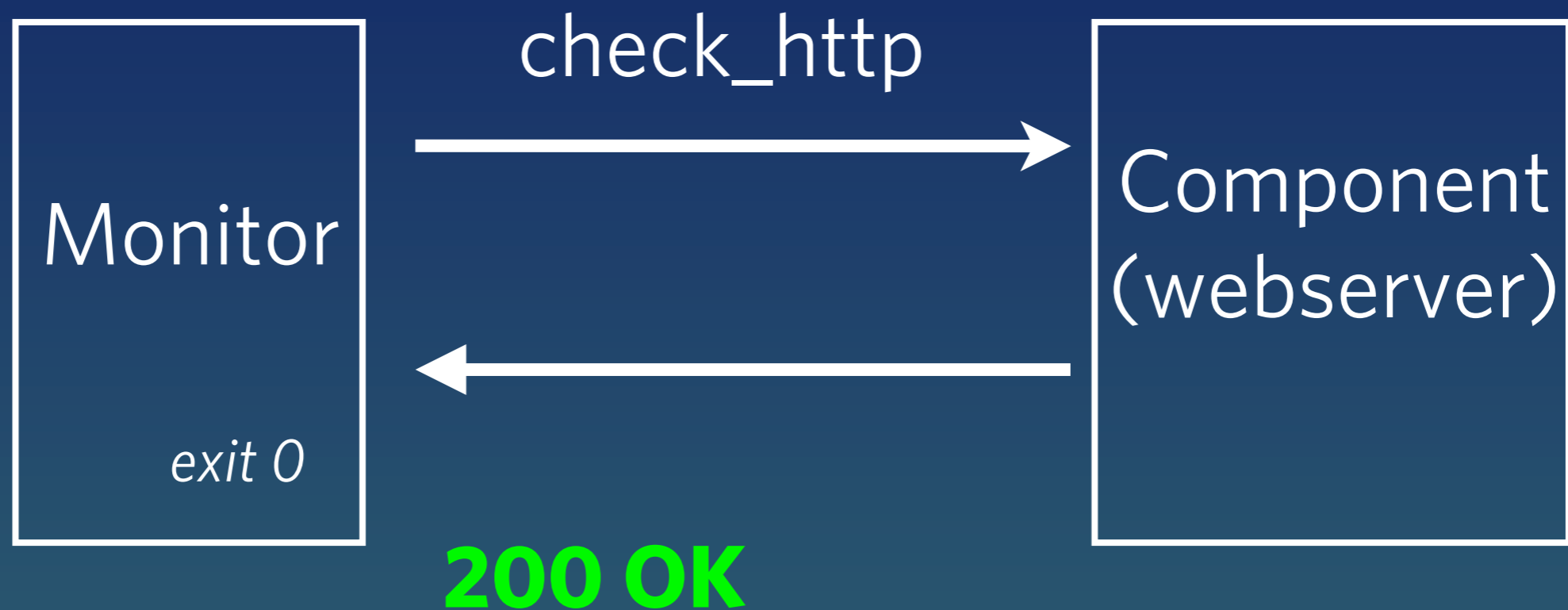
Anomaly Detection

Anomaly Detection

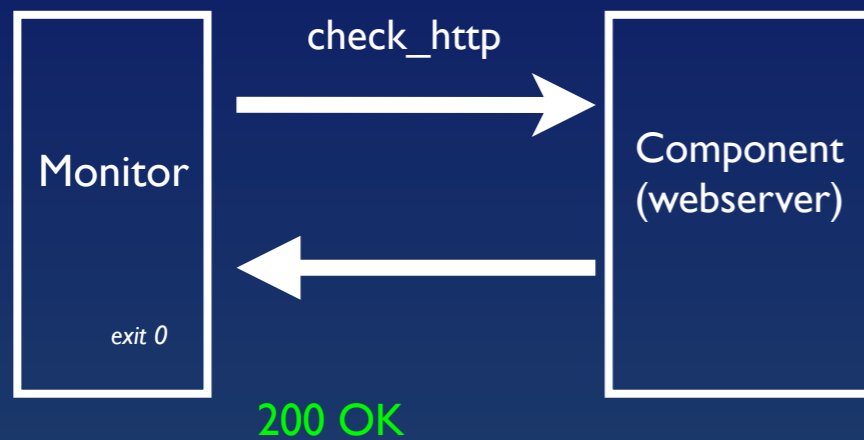
- **Getting at the state of health**
- **Evaluating the state of health**
- **Components AND systems**

Supervisory

Example: Active health check



Supervisory



Pros:

- Easy to implement**
- Easy to understand**
- Well-known pattern**

Cons:

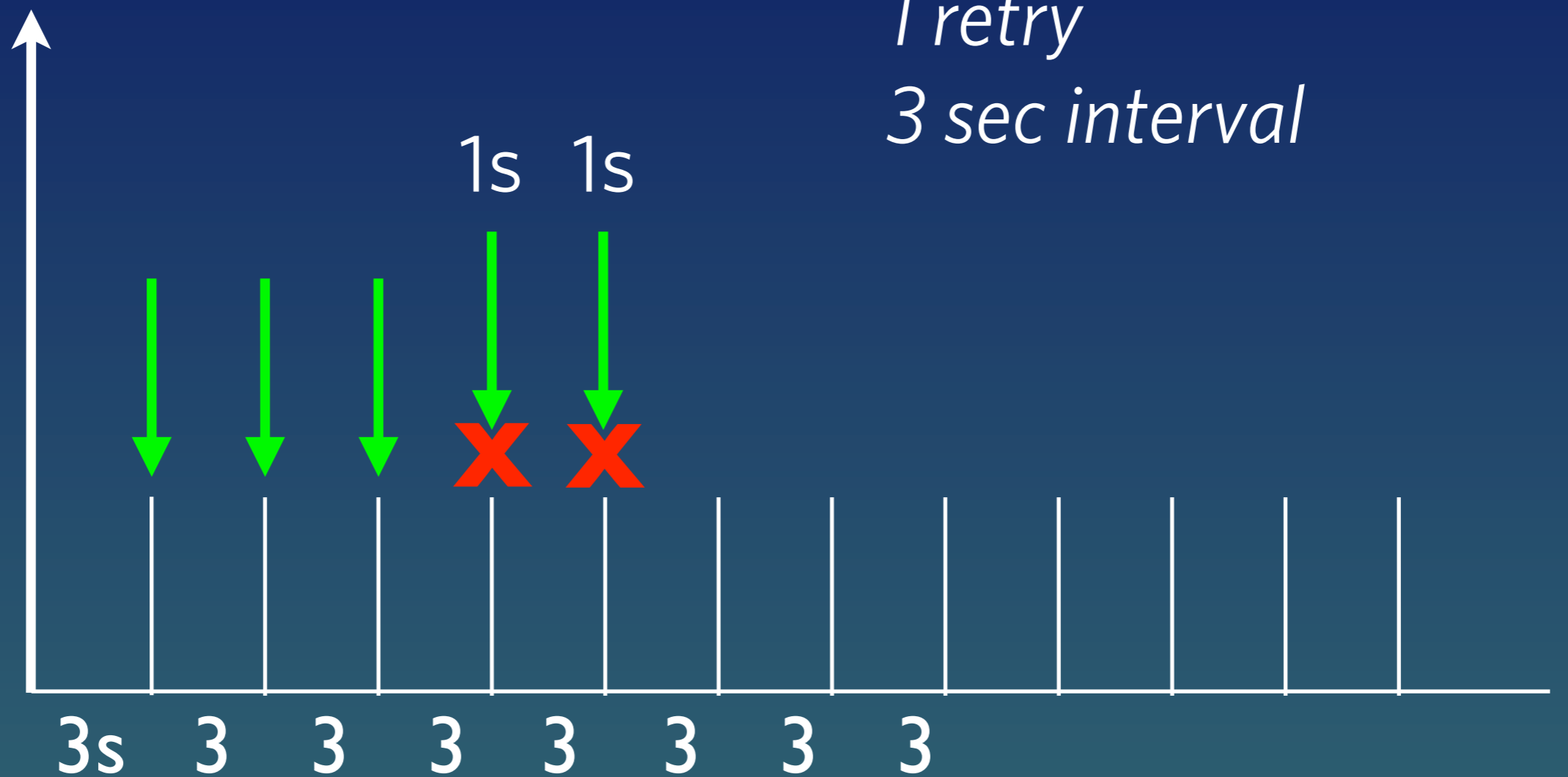
- Messaging can fail**
- Scalability is limited**

Supervisor Sensitivity

1 sec timeout

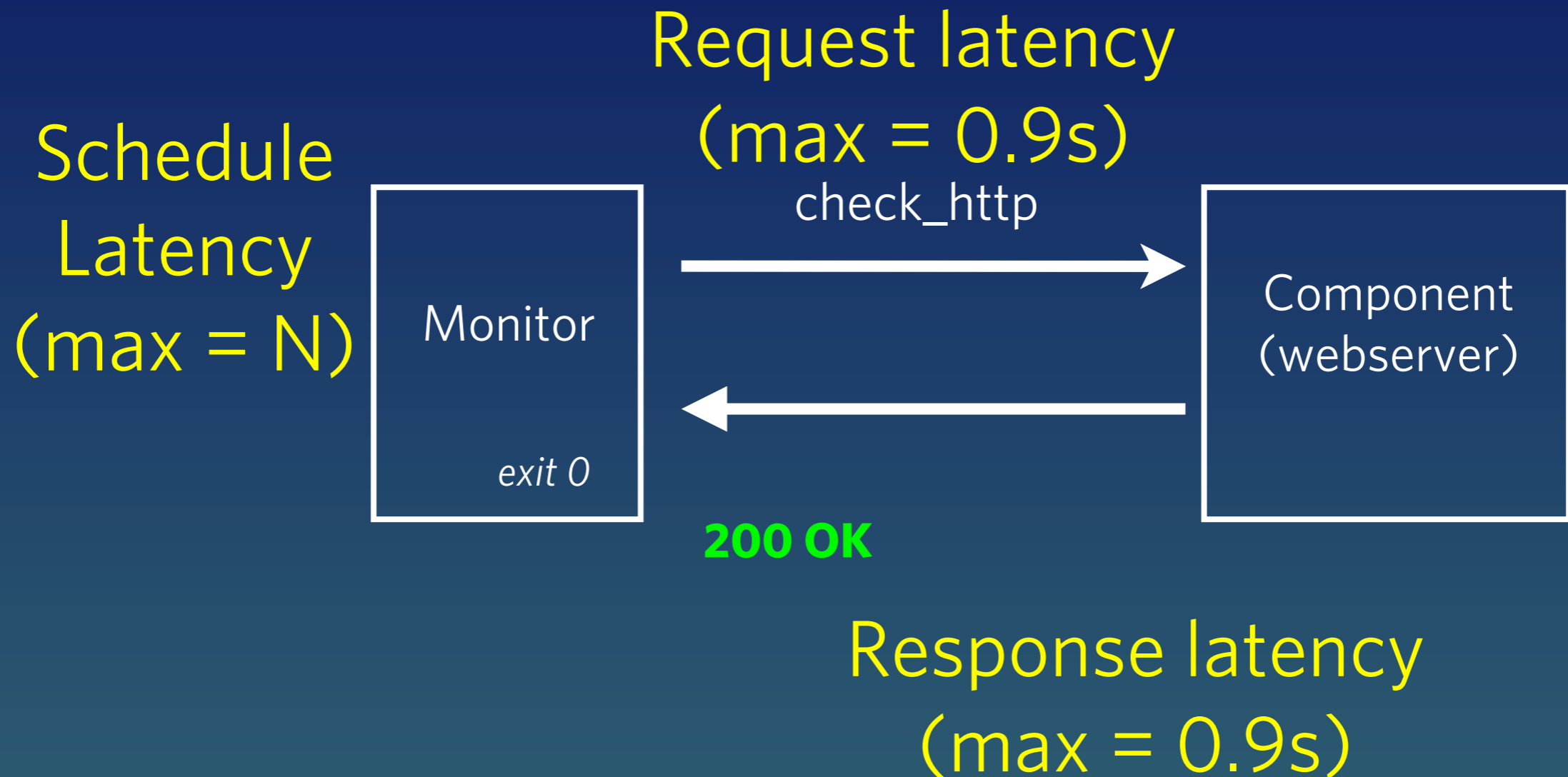
1 retry

3 sec interval



(7.9 sec exposure)
Up to ~2.9s for the previous interval

Supervisor Sensitivity



Supervisor Sensitivity

*How many seconds of errors
can you tolerate serving?*

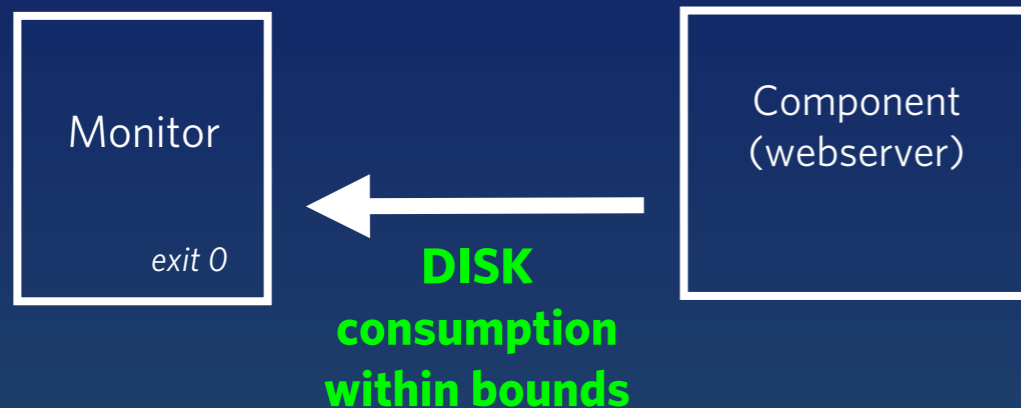
Supervisory

Example: Interval Passive health check



Supervisory

Example: Interval Passive health check



Pros:

Efficient

Scalability is different

Fewer moving parts

Less exposure

Can submit to multiple places

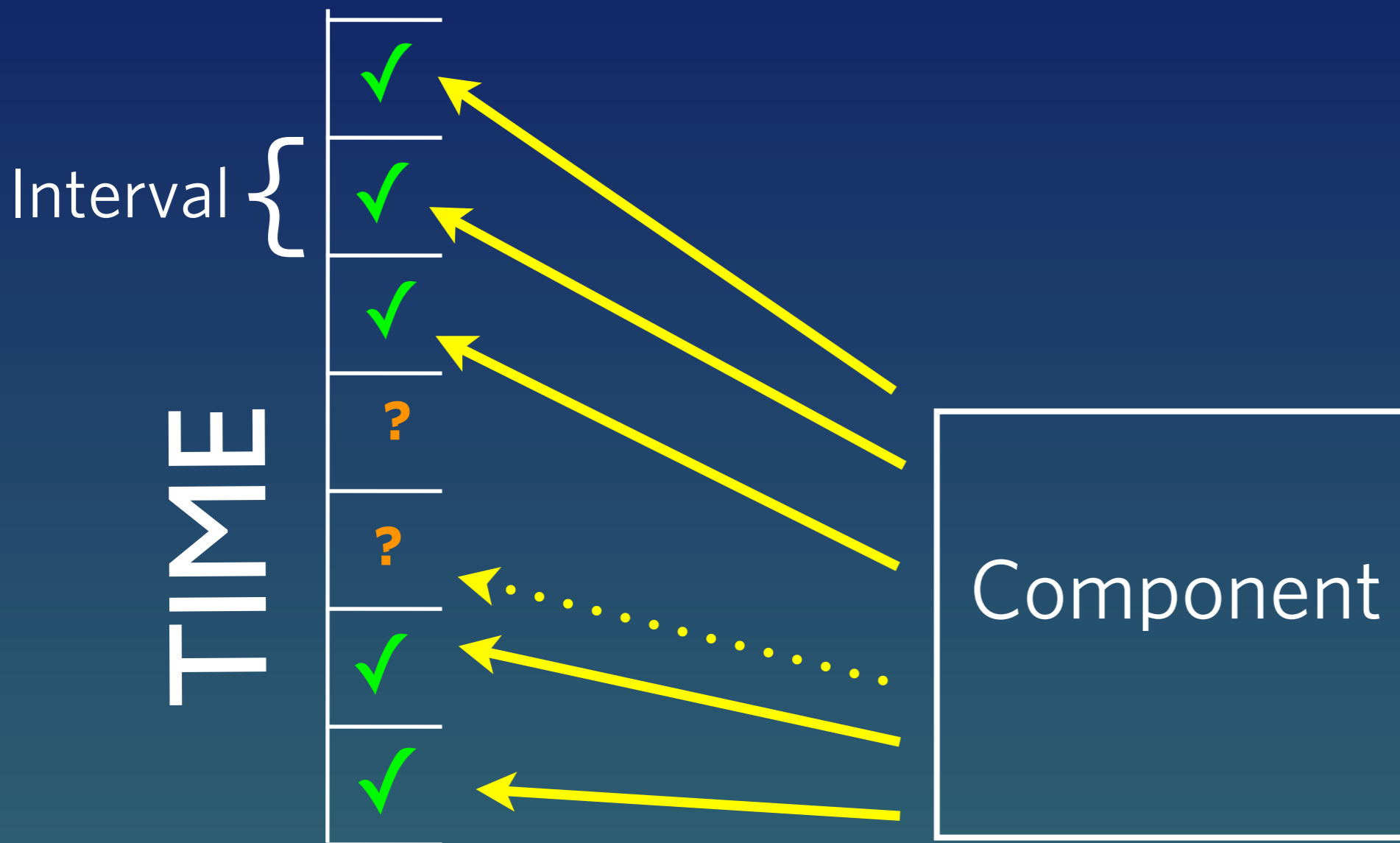
Cons:

Nonideal for network-based services

Different tuning (windowed expectation)

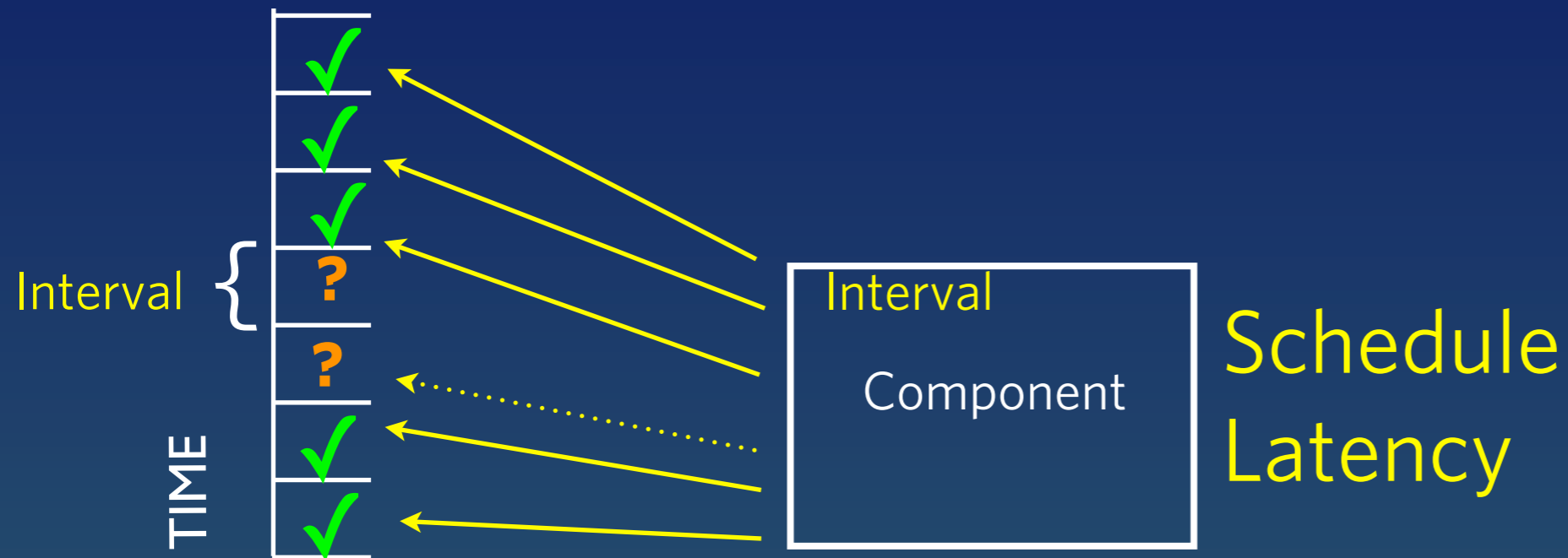
Supervisory

Example: Passive health check



Supervisory

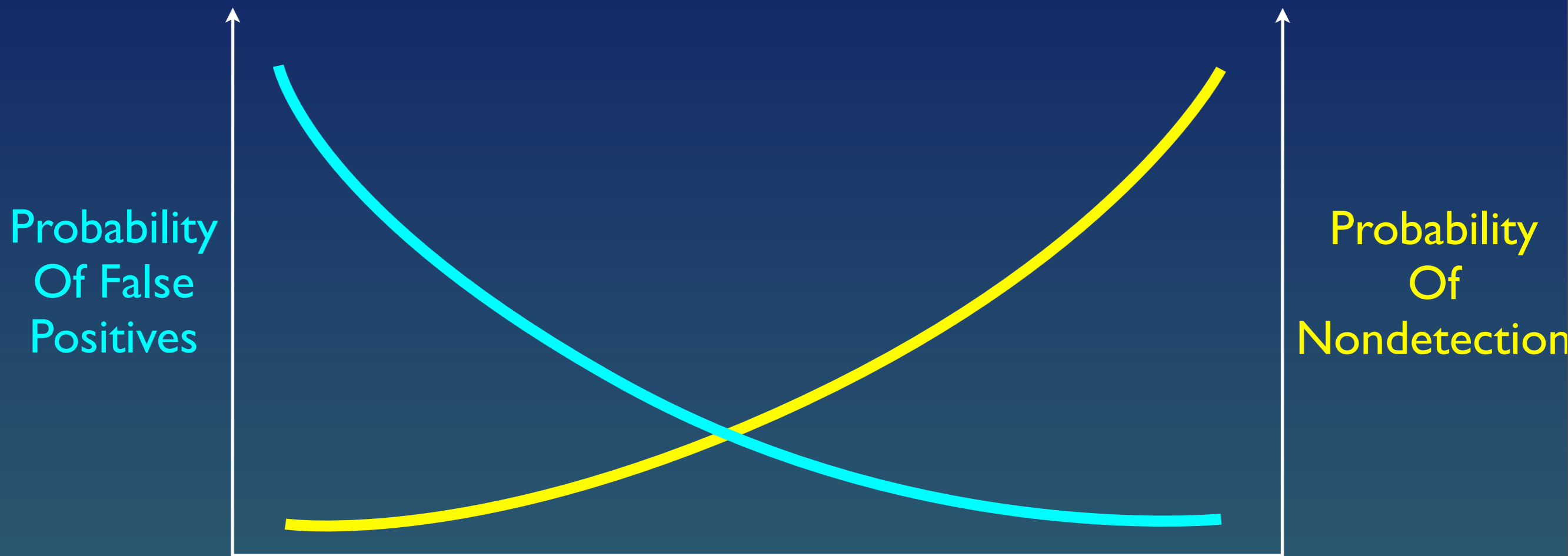
Example: Passive health check



Exposure =

$$(\text{Schedule} + \text{Interval}) * \text{UnknownConsecutiveIntervals} + 1$$

Frequency and Transience



Probability
Of False
Positives

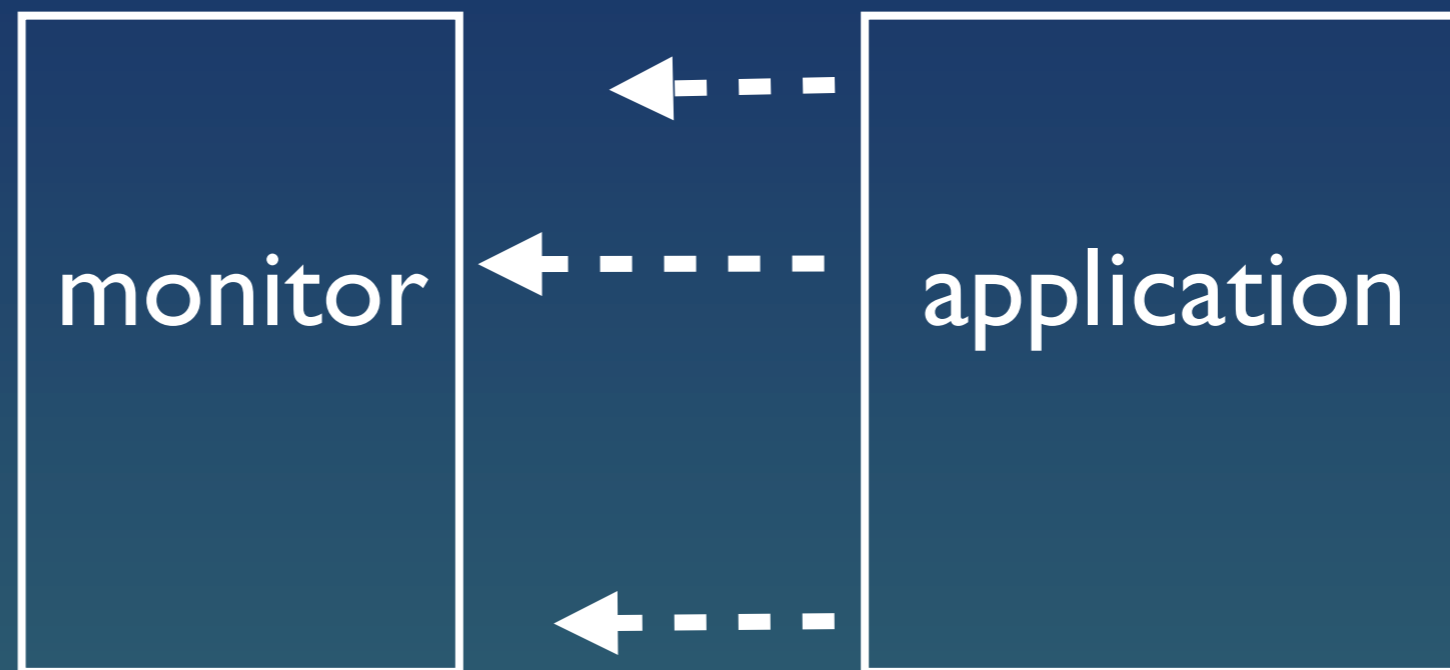
Probability
Of
Nondetection

Short intervals
Low # of retries
Short timeouts

Long intervals
High # of retries
Long timeouts

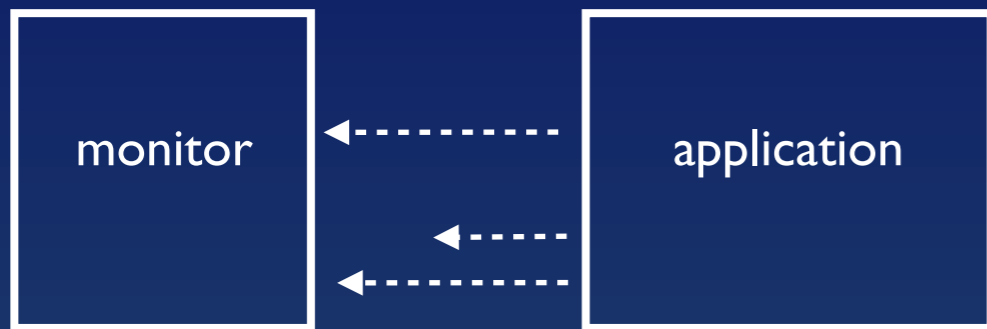
In-Line

Example: Passive application event logging



Supervisory

Example: Passive application event logging



Pros:

On-demand publish

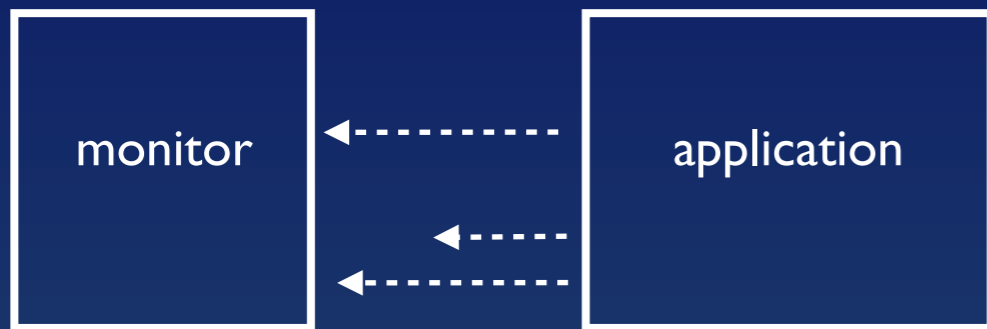
Cons:

Onus is on the app

Can't be 100% sure it's working

Supervisory

Example: Passive application event logging



Positive events (sales, registrations, etc.)

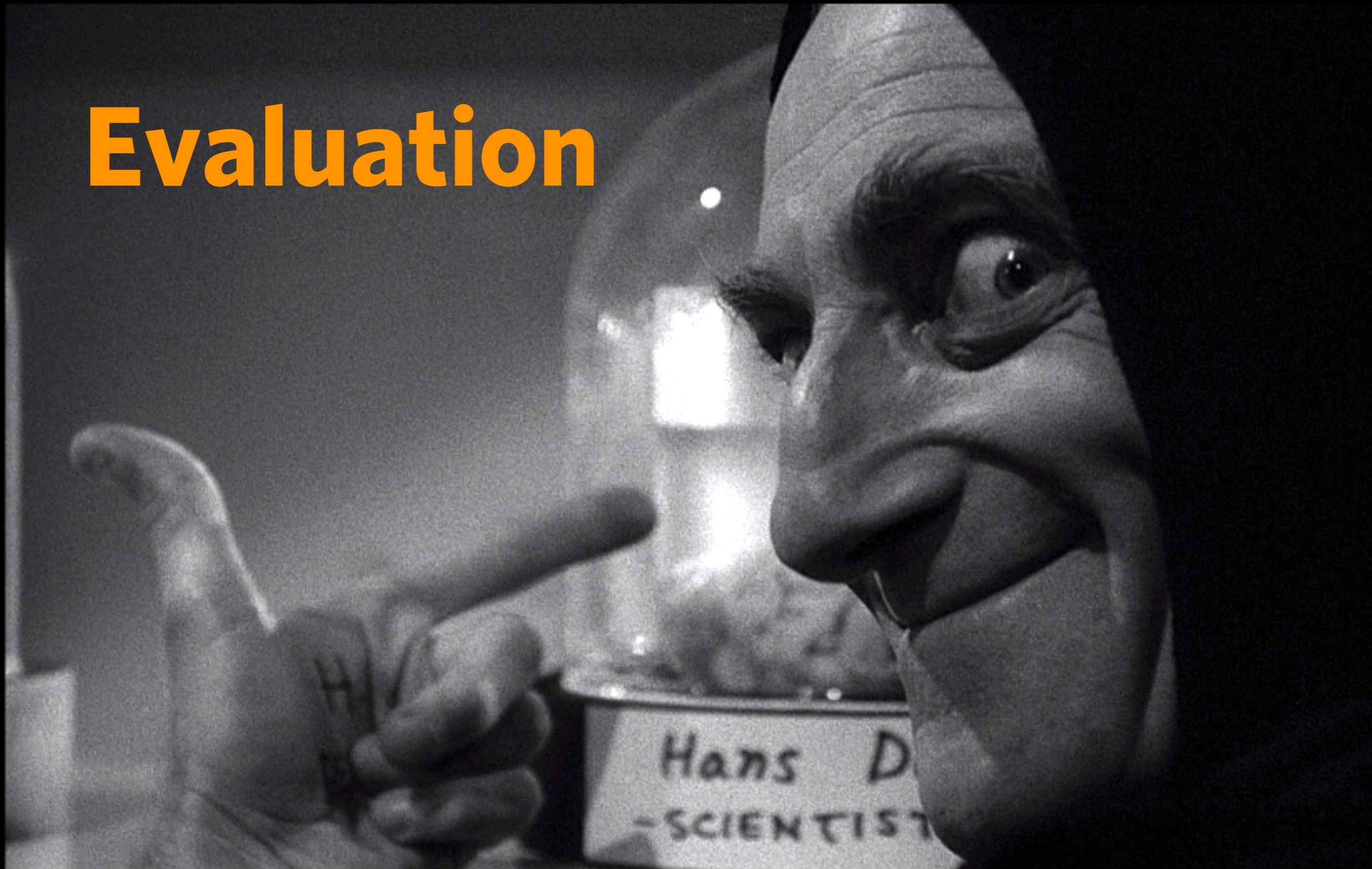
Negative events (errors, exceptions, etc.)

Lack or presence of data mean different things, so history is paramount.

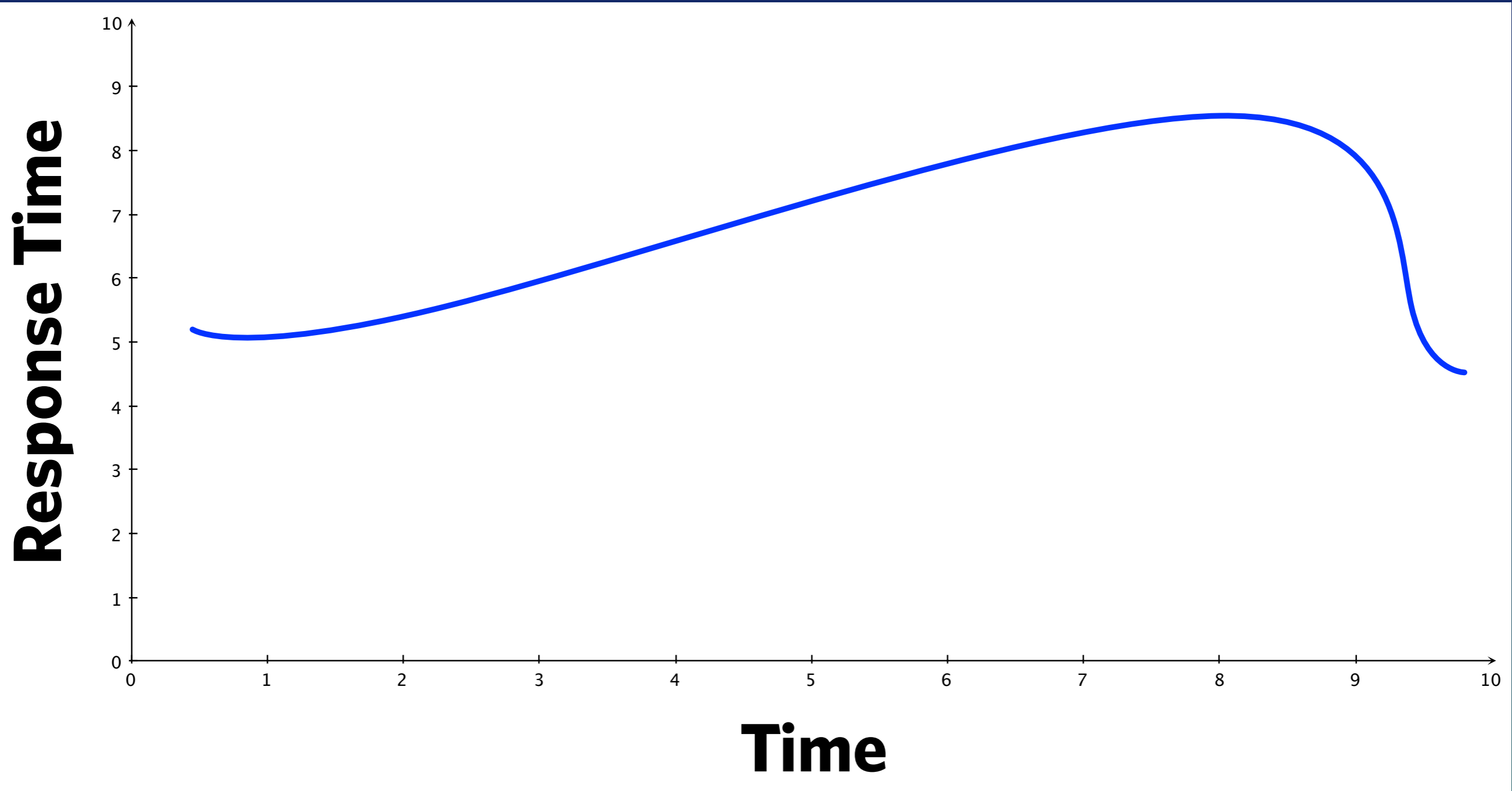


Context

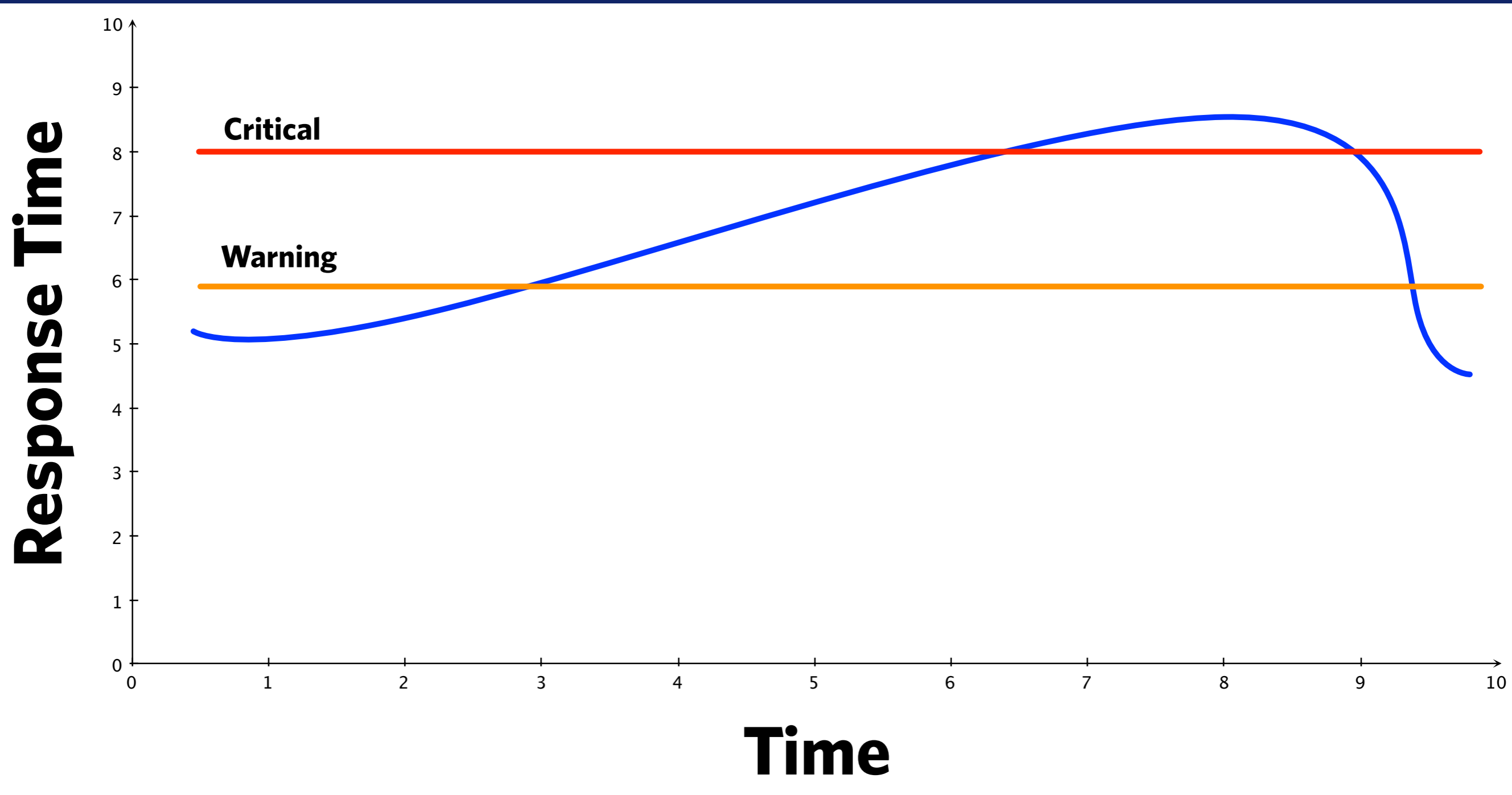
Evaluation



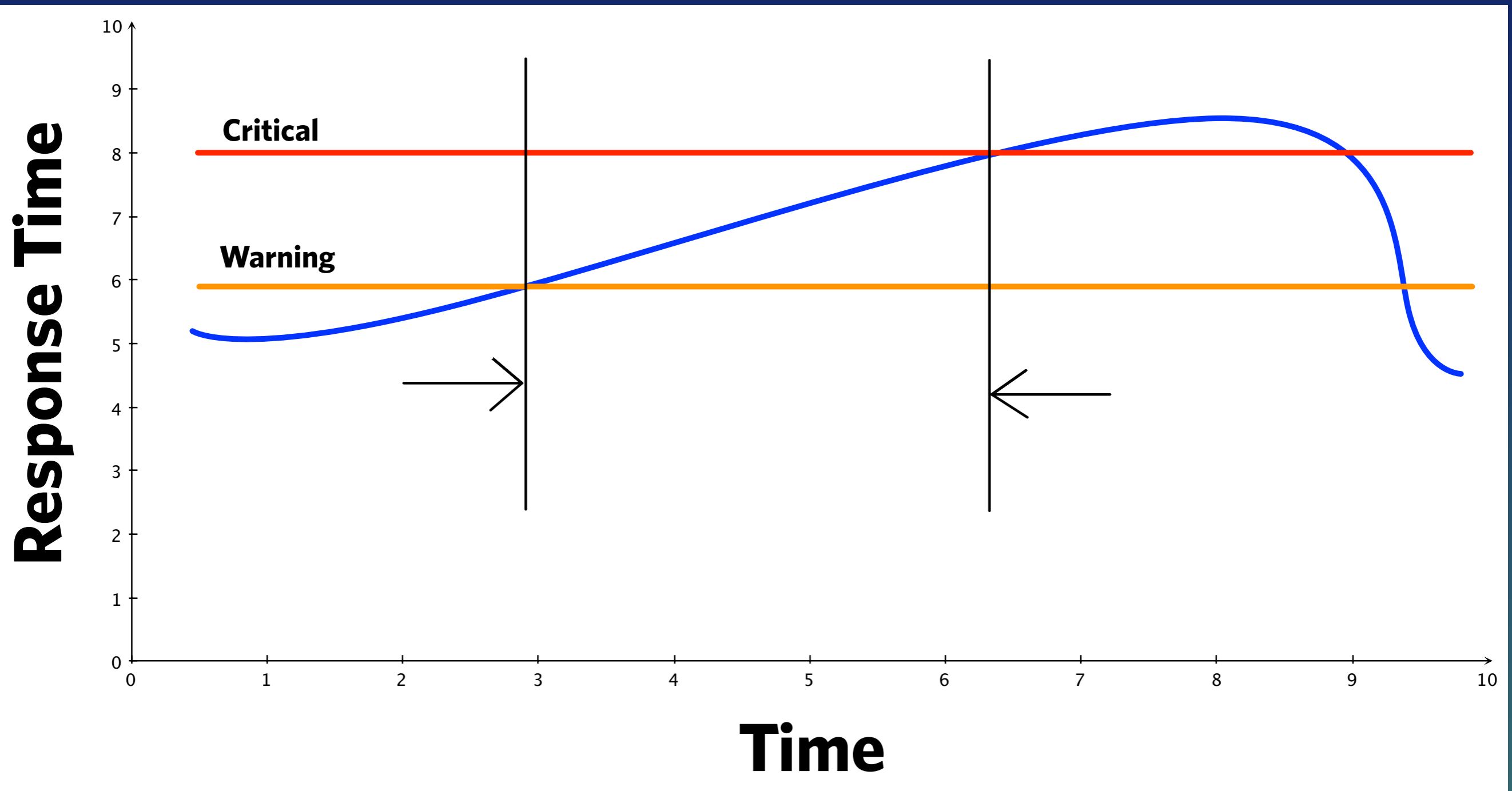
what is 'abnormal' ?



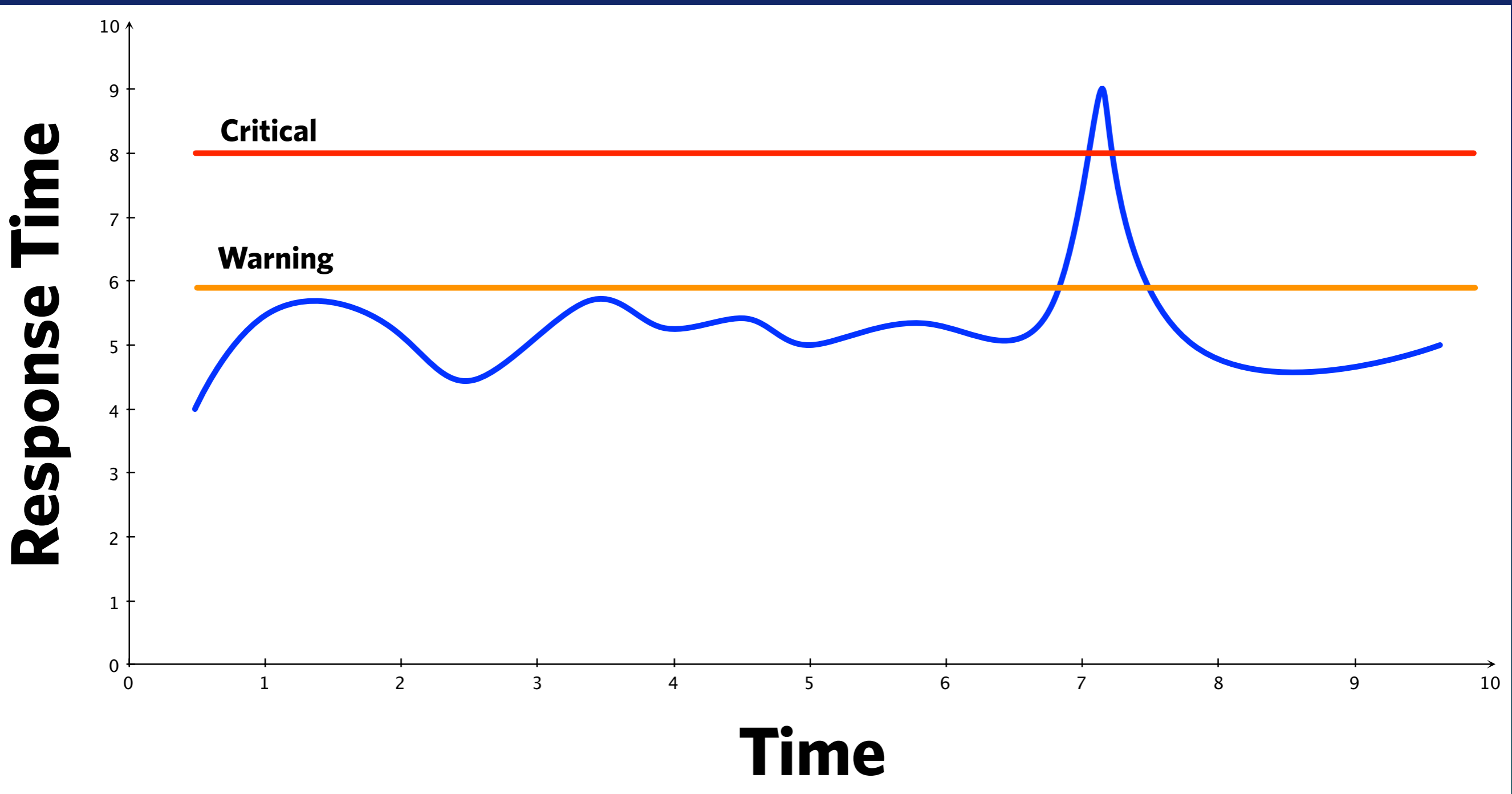
Static Thresholds



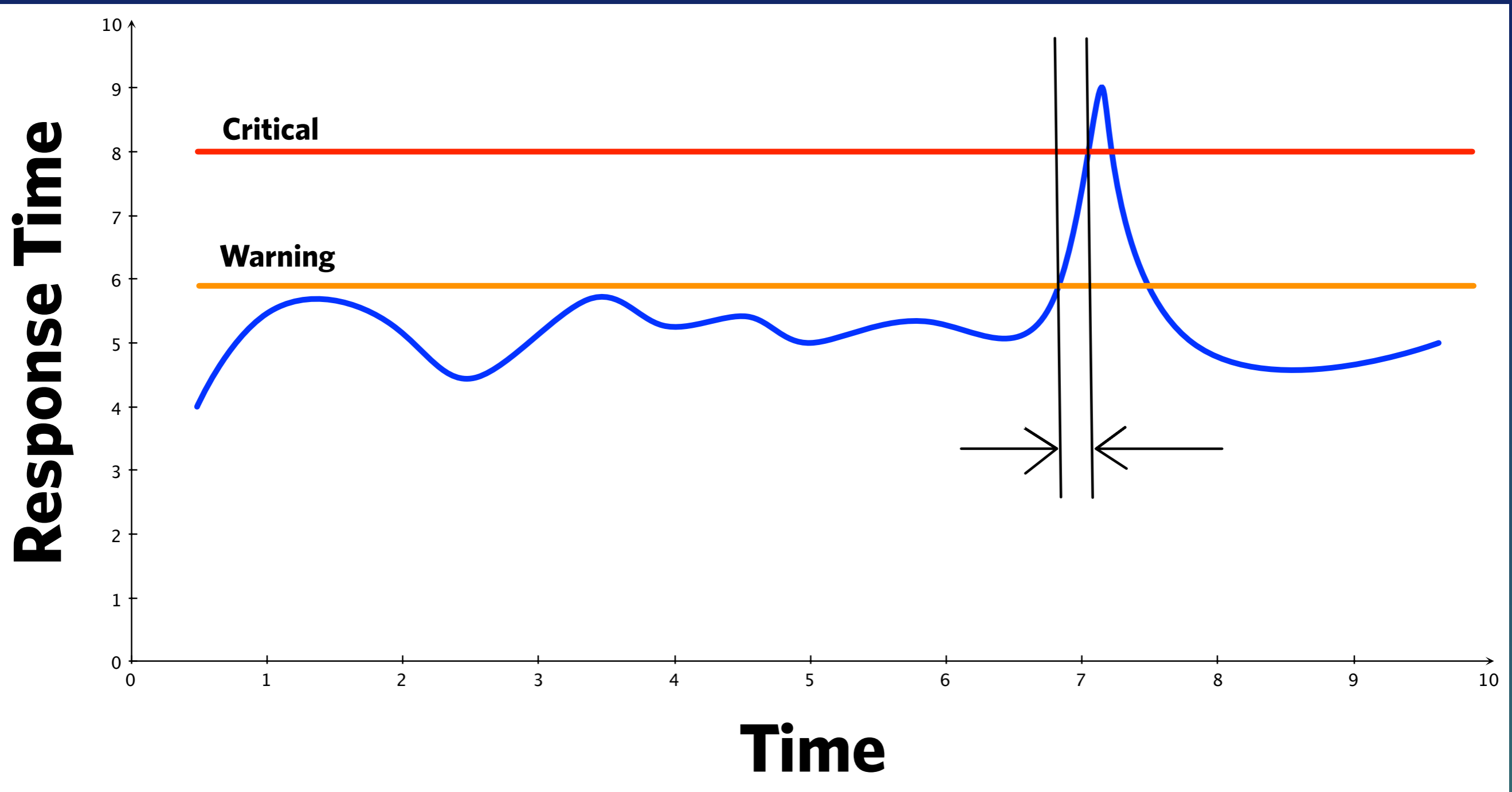
Static Thresholds



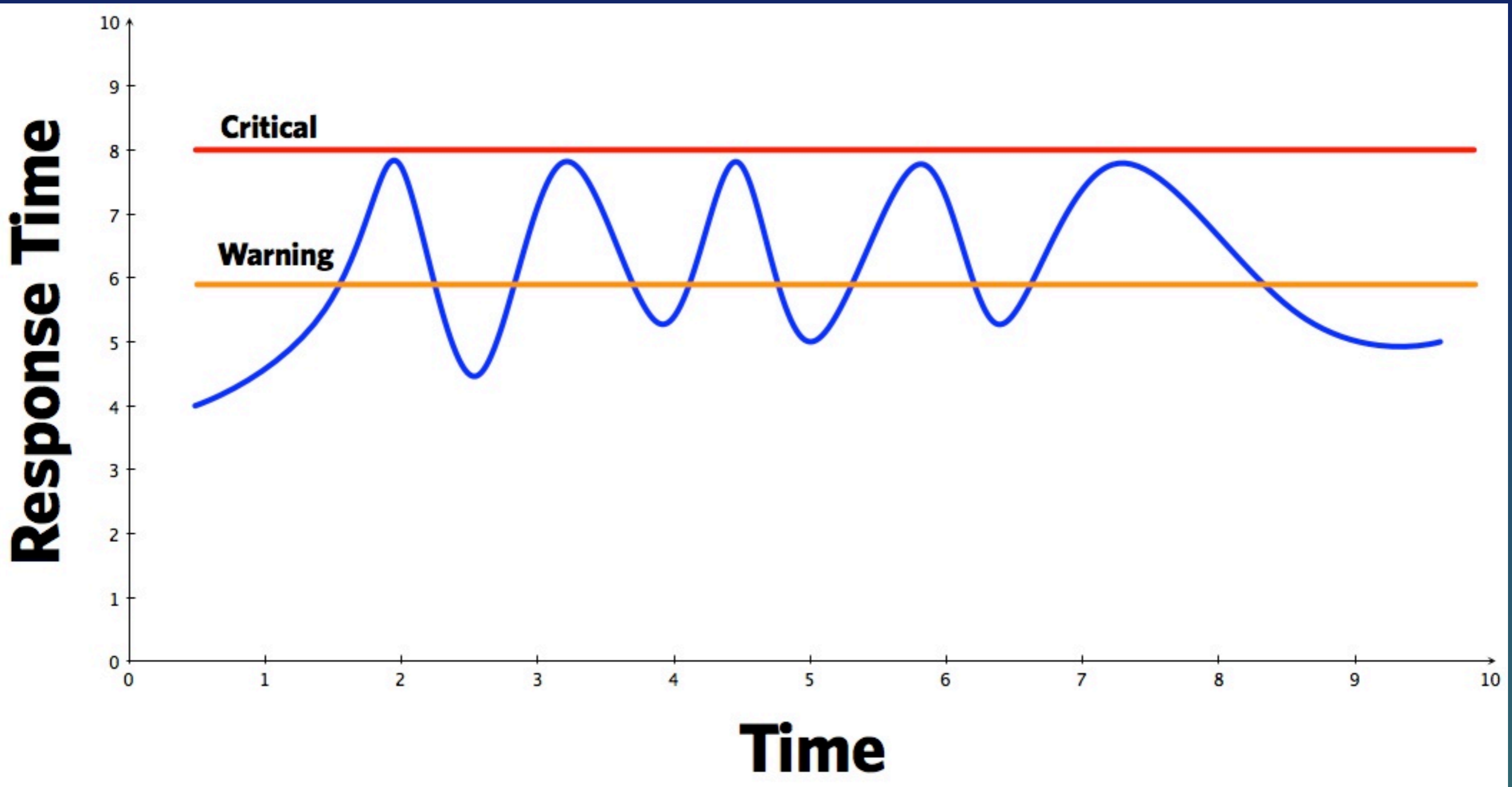
Static Thresholds



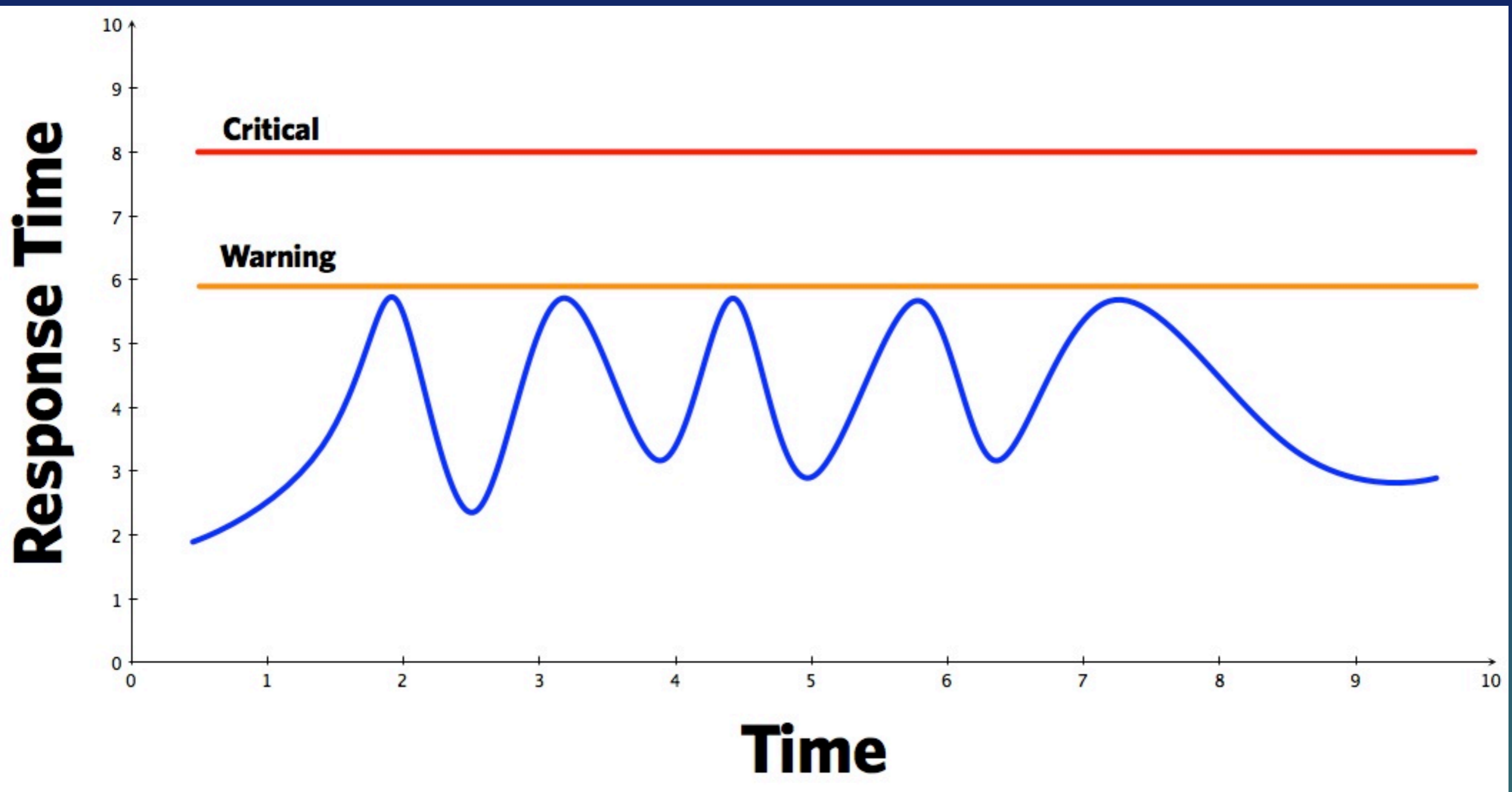
Static Thresholds



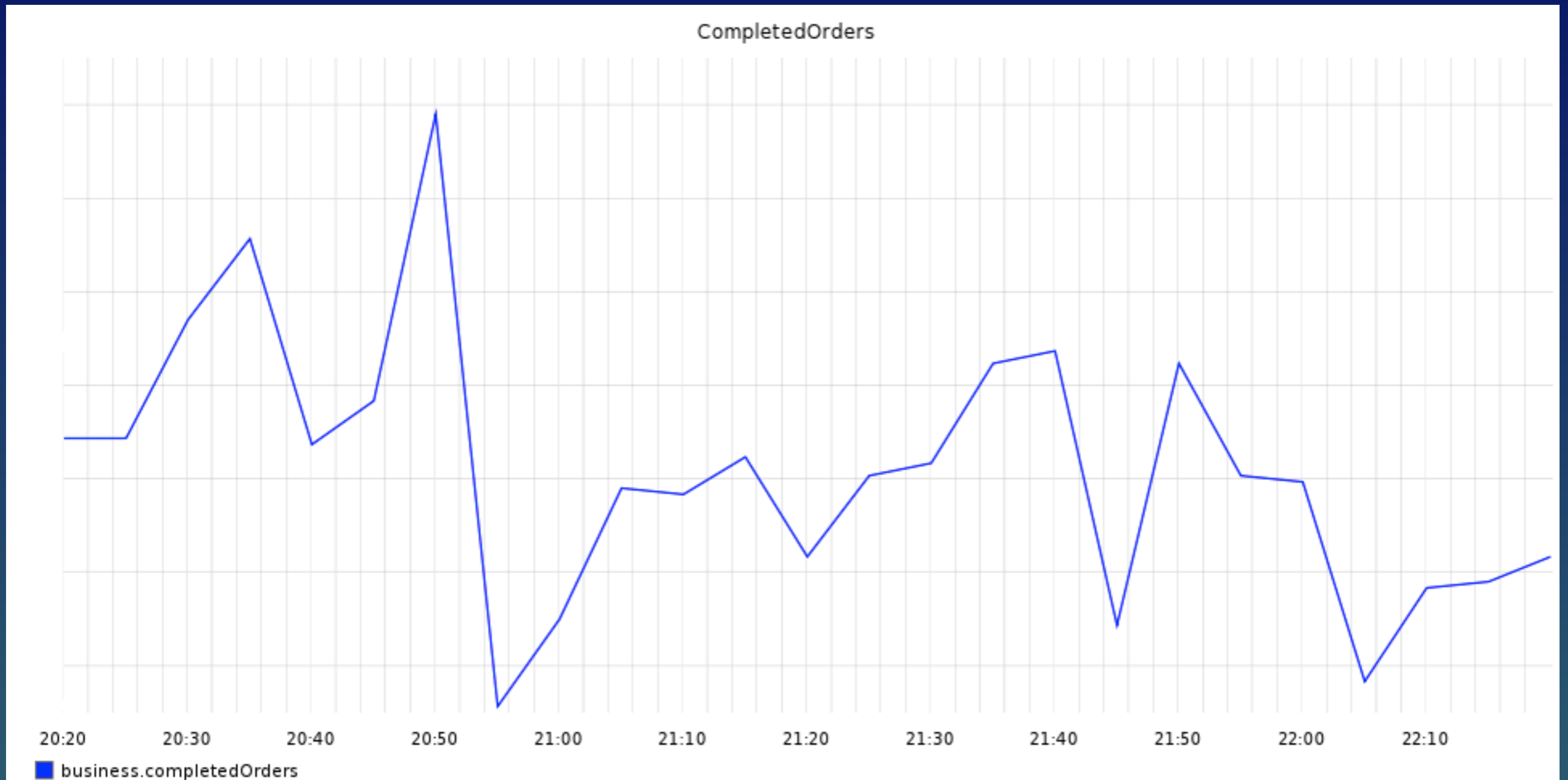
Static Thresholds



Static Thresholds

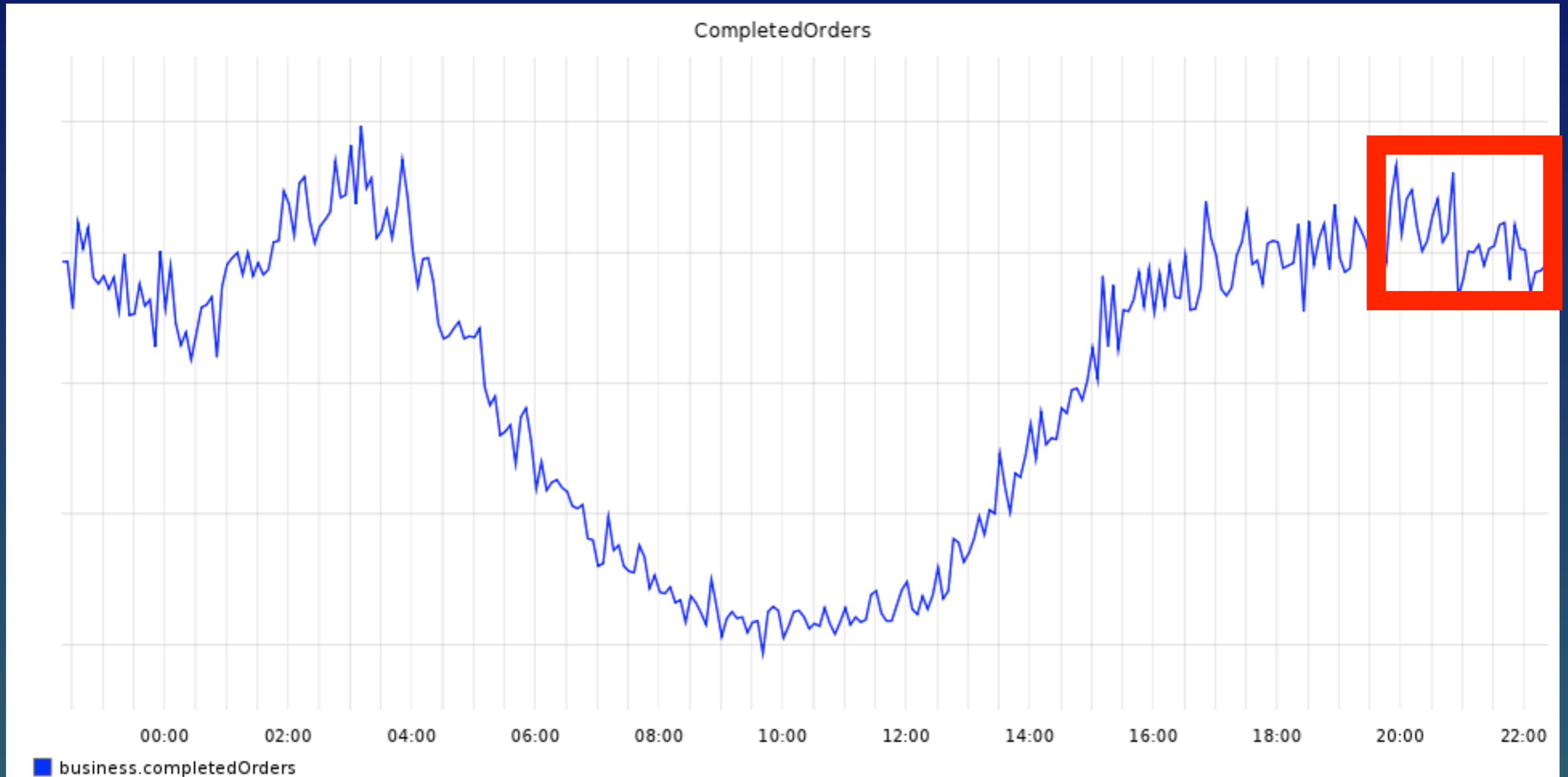


Context



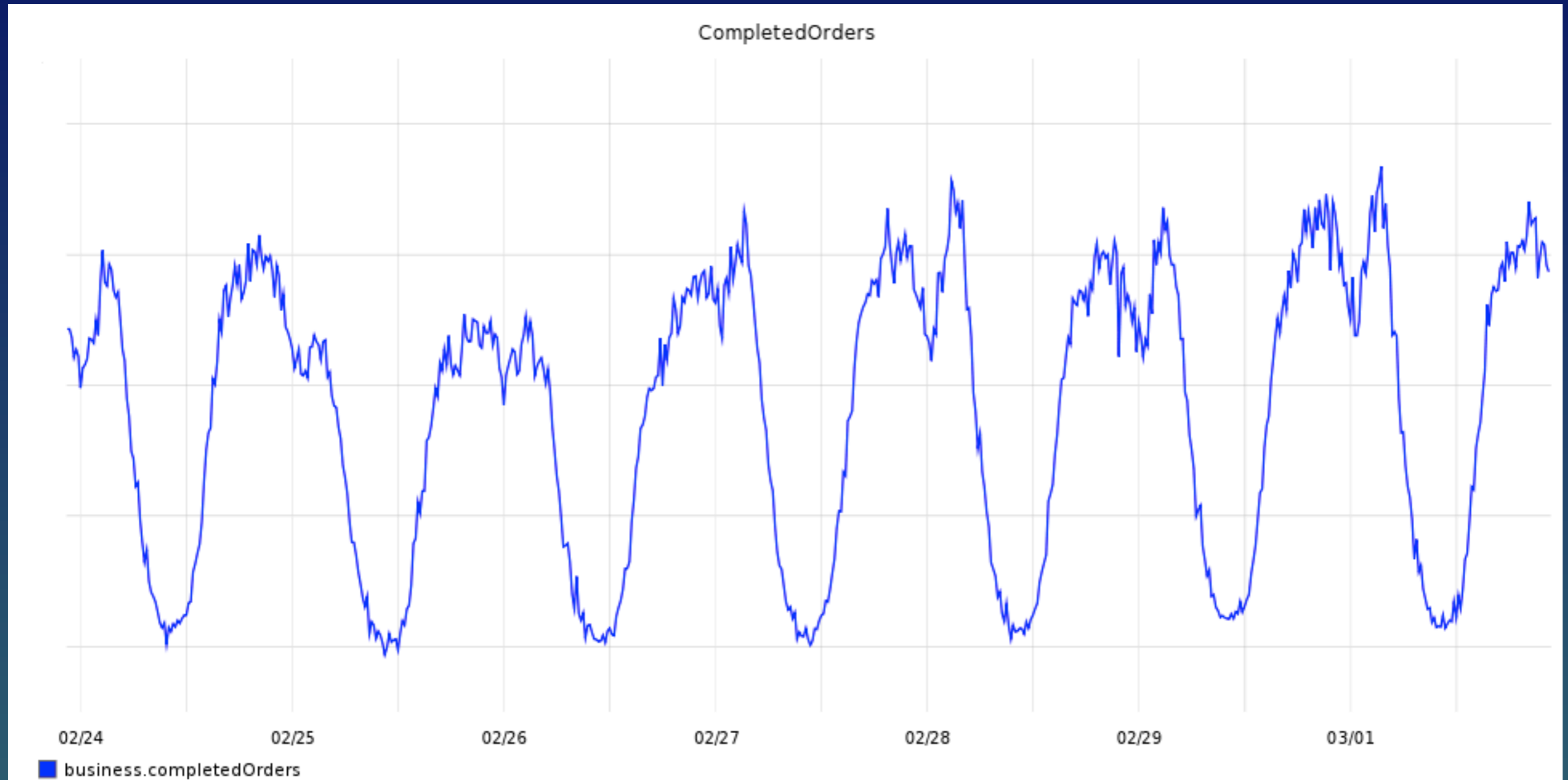
Normal?

Context



24 hours

Context

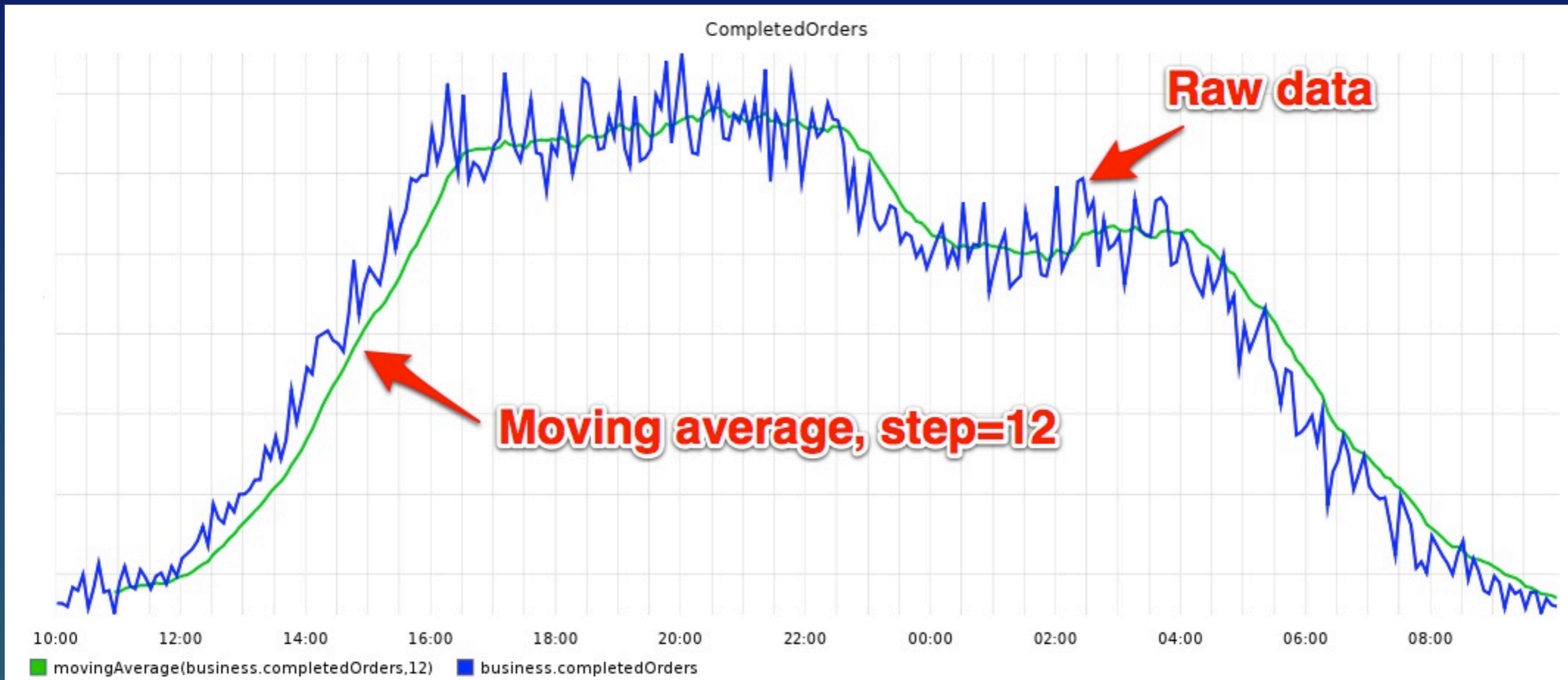


7 days

Context



Context



Smoothing?

Context

Holt-Winters Exponential Smoothing

Recent points influencing a forecast, exponentially decreasing influence backwards in time.

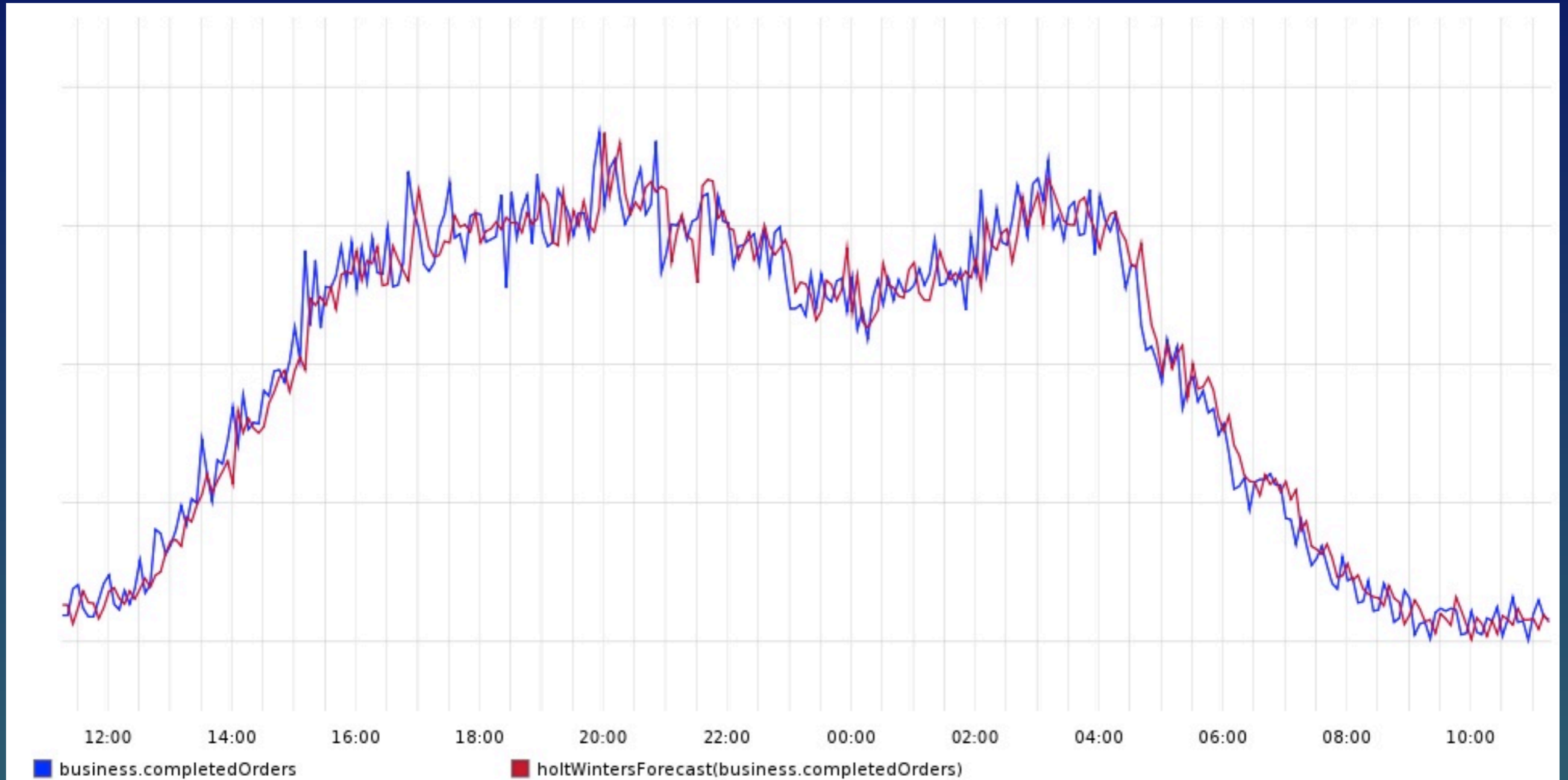
en.wikipedia.org/wiki/Exponential_smoothing

Context

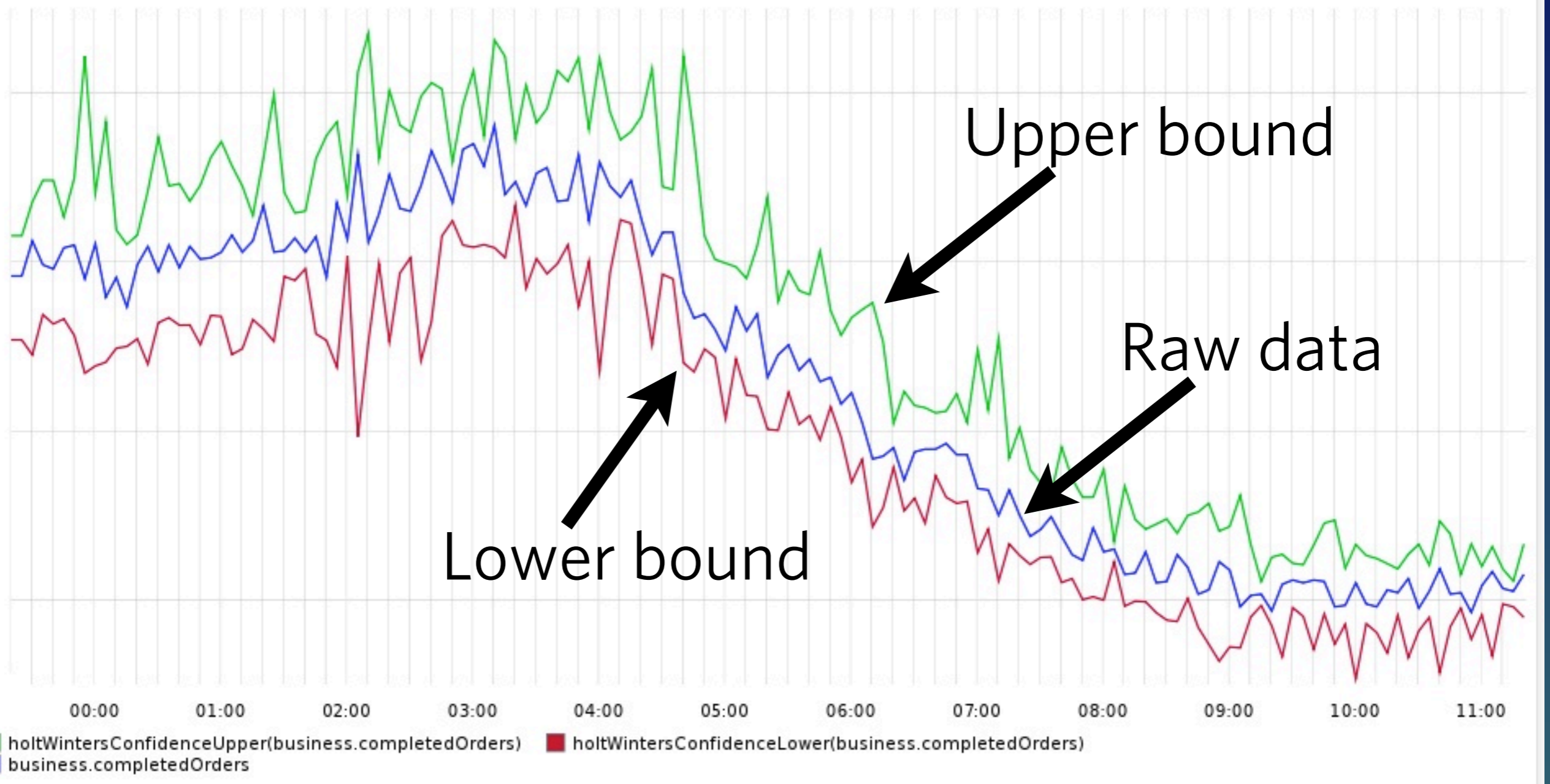
Aberrant Behavior Detection in Time Series for Network Monitoring

http://static.usenix.org/events/lisa00/full_papers/brutlag/brutlag_html/

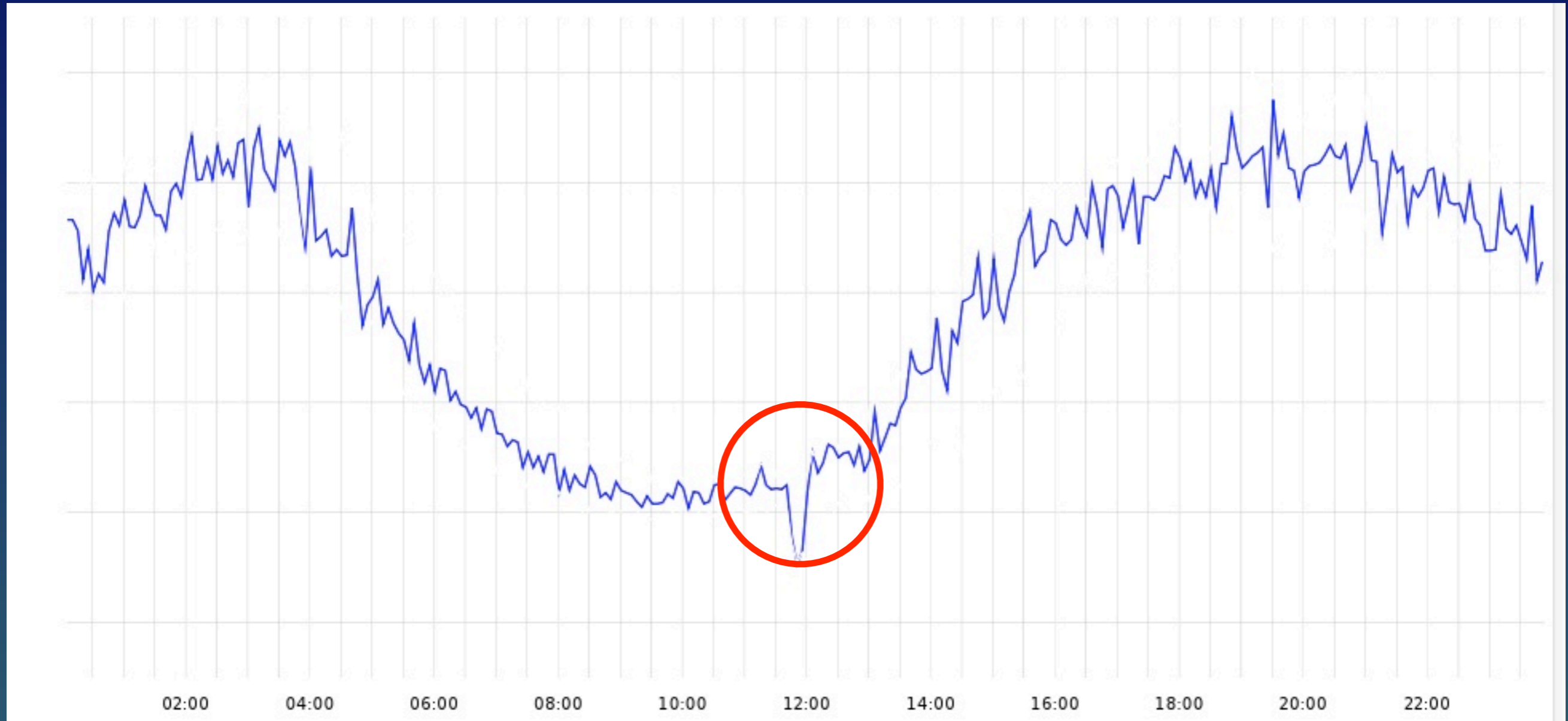
Dynamic Thresholds



Dynamic Thresholds

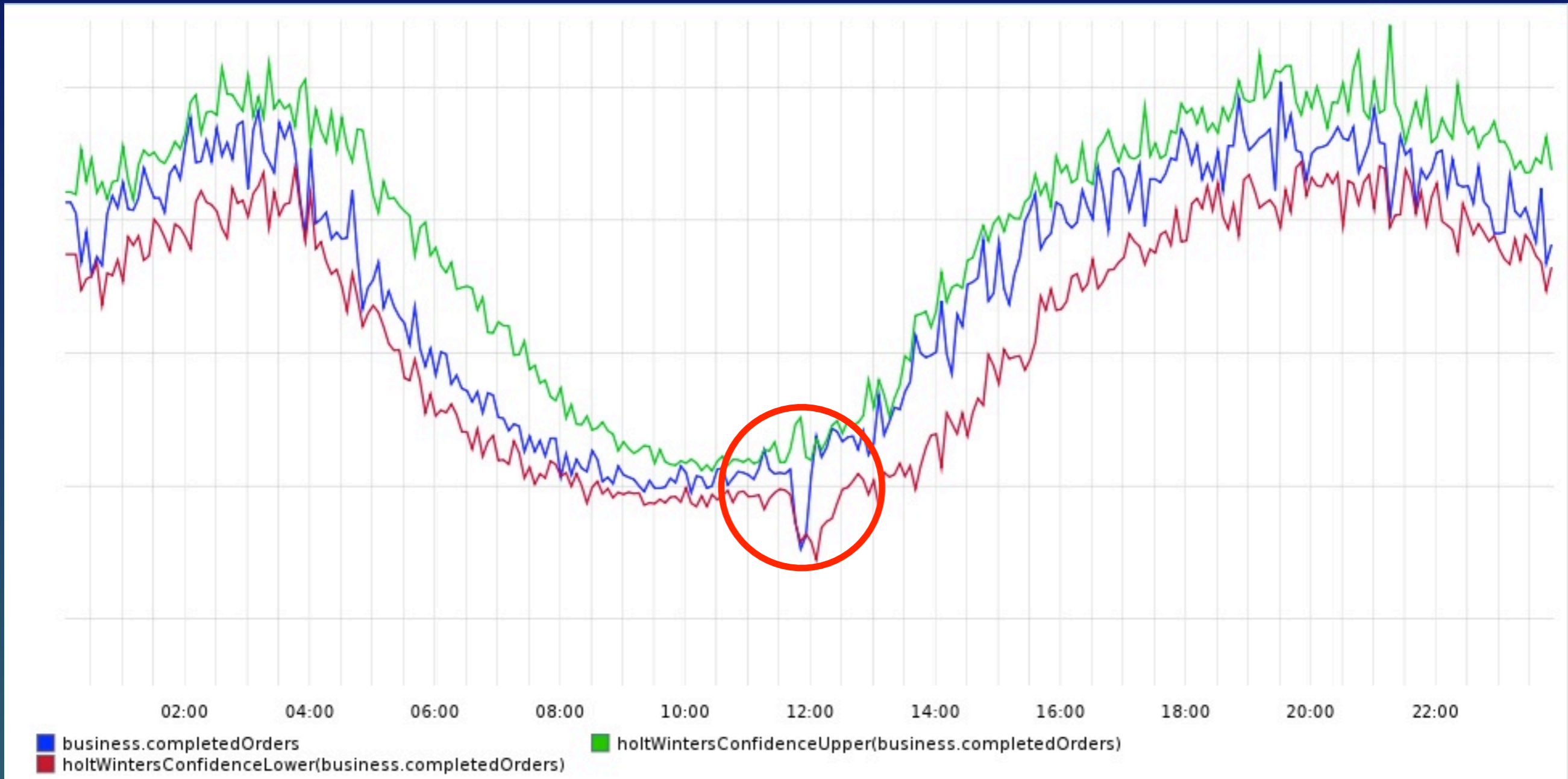


Dynamic Thresholds



Hrm.....

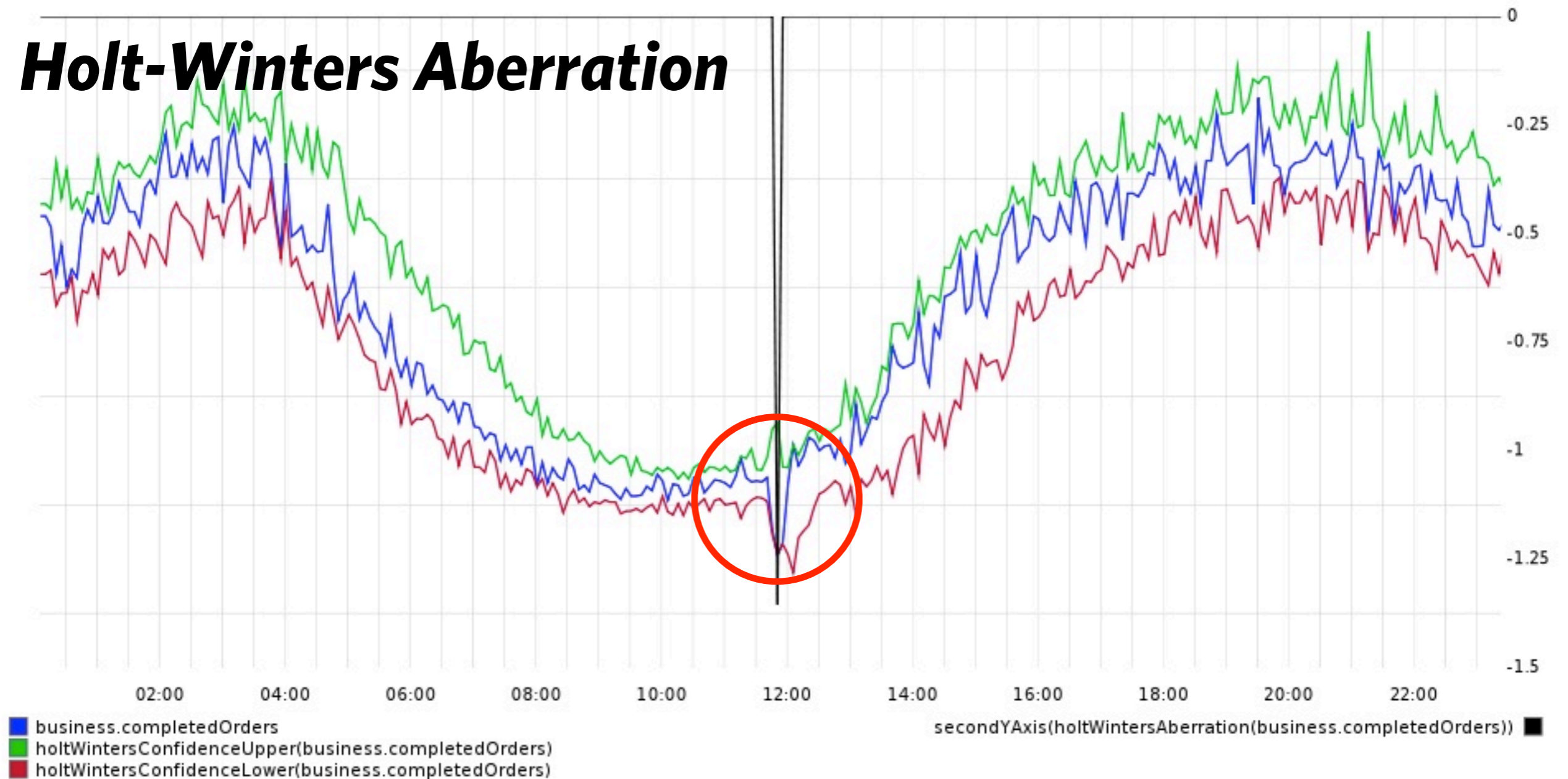
Dynamic Thresholds



Hrm.....

Dynamic Thresholds

Holt-Winters Aberration



Ah!

Dynamic Thresholds

**Graphite metrics collection w/Holt-Winters
abberations**

<http://graphite.wikidot.com/>

Nagios check for Graphite data

https://github.com/etsy/nagios_tools/blob/master/check_graphite_data

Four Cornerstones

Erik Hollnagel

(Anticipation)

Knowing
What
To Expect

(Response)

Knowing
What
To Do

(Monitoring)

Knowing
What
To Look For

(Learning)

Knowing
What
Has Happened

FAULT TOLERANCE

Detection of fault X



Triggers corrective action Y



Clean up, report back

(RECOVERY OR MASKING)

Variation Tolerance

Adaptive Systems

Expected
Variation



Adaptive Systems

Expected
Variation



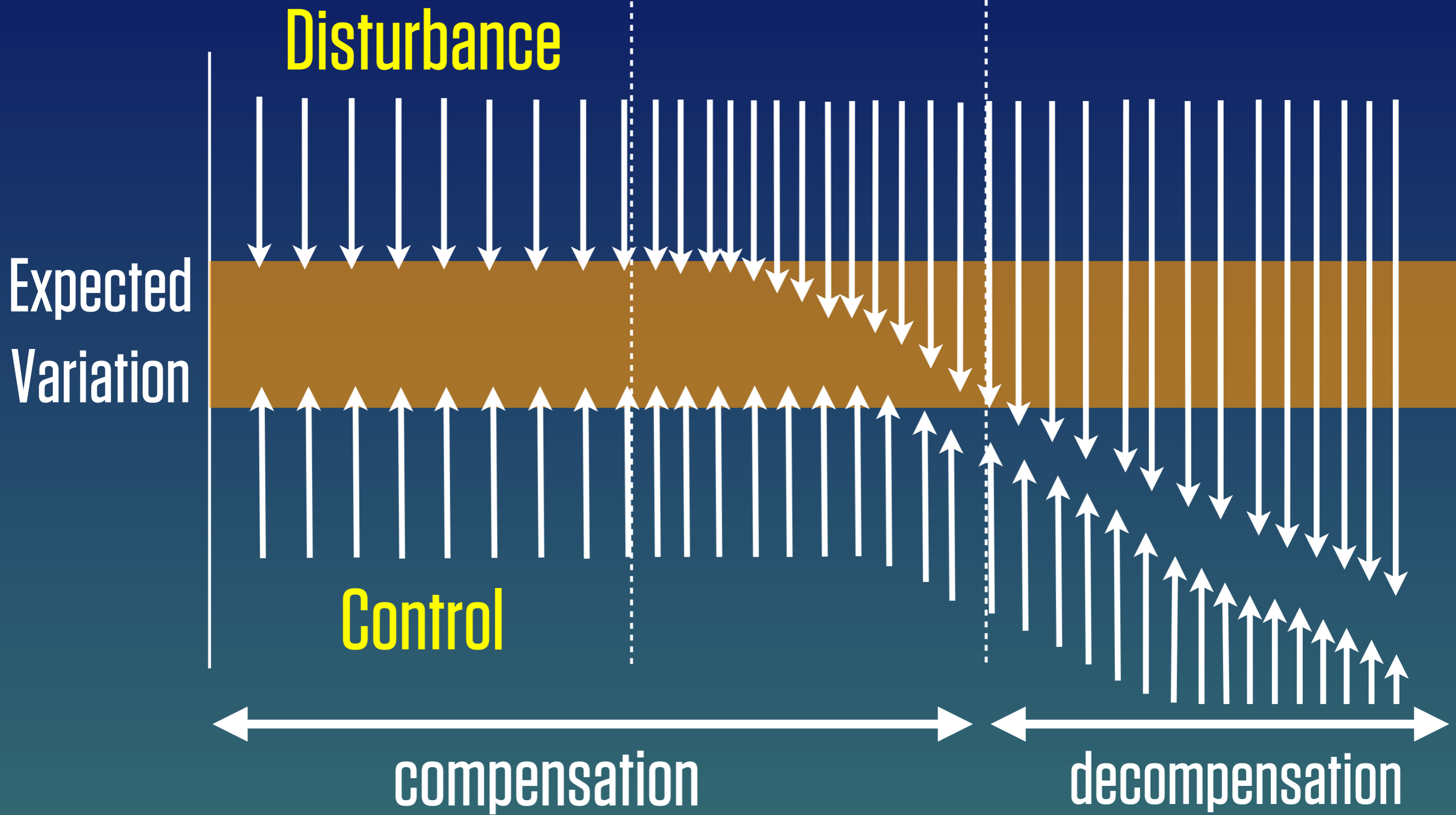
Adaptive Systems

Expected
Variation

A horizontal bar chart with a gold bar and a teal bar. A vertical white line is on the left side of the gold bar.

New Disturbances
Arise

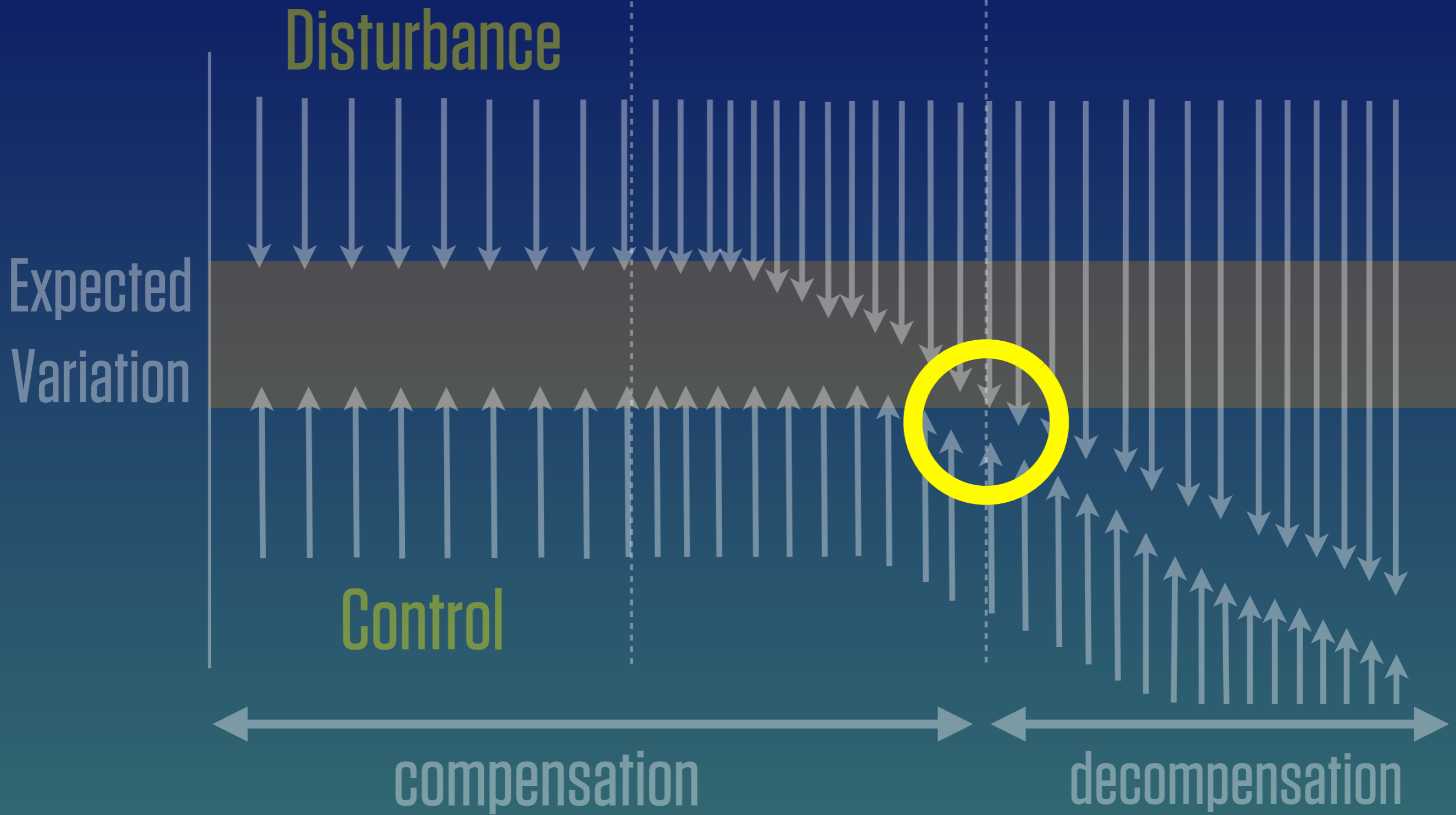
Compensation is
Exhausted



Woods, 2011

New Disturbances
Arise

Compensation is
Exhausted



New Disturbances
Arise

Compensation is
Exhausted

Variation

Expected
Variation

Fault

Control

compensation

decompensation



Variations \neq Faults

Dead
Corrupt
Late
Wrong

Fault Tolerance

Redundancy

Spatial (server, network, process)

Temporal (checkpoint, "rollback")

Informational (data in N locations)

Fault Tolerance

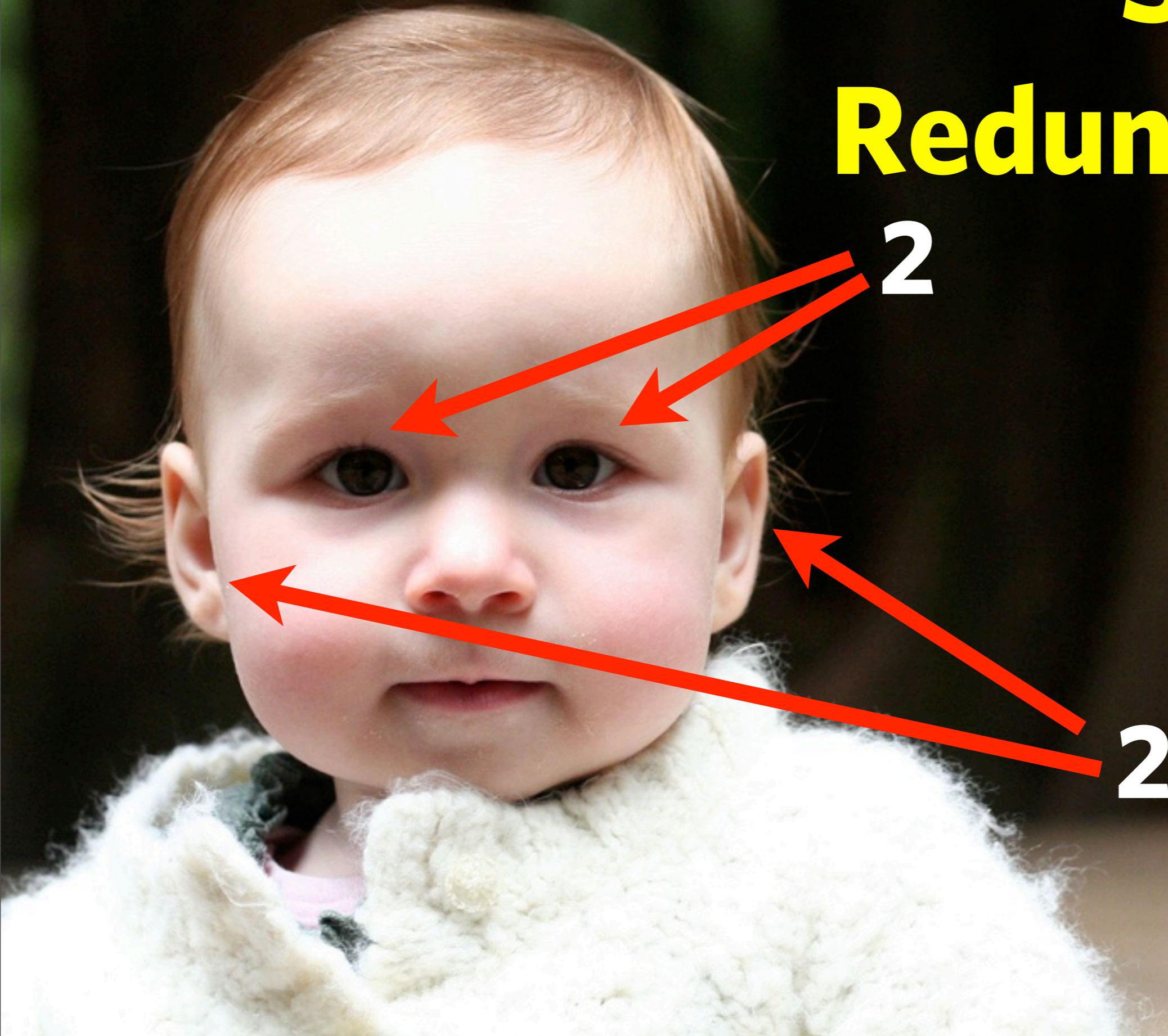
Redundancy

Spatial (server, network, process)

Temporal (checkpoint, "rollback")

Informational (data in N locations)

Spatial Redundancy



Spatial Redundancy

Active/Active



A close-up photograph of a light brown and white cat's face. The cat's eyes are green, and its mouth is slightly open, showing its teeth. The background is dark. Overlaid on the image are two text elements: 'Spatial Redundancy' in yellow and 'Active/Passive' in white.

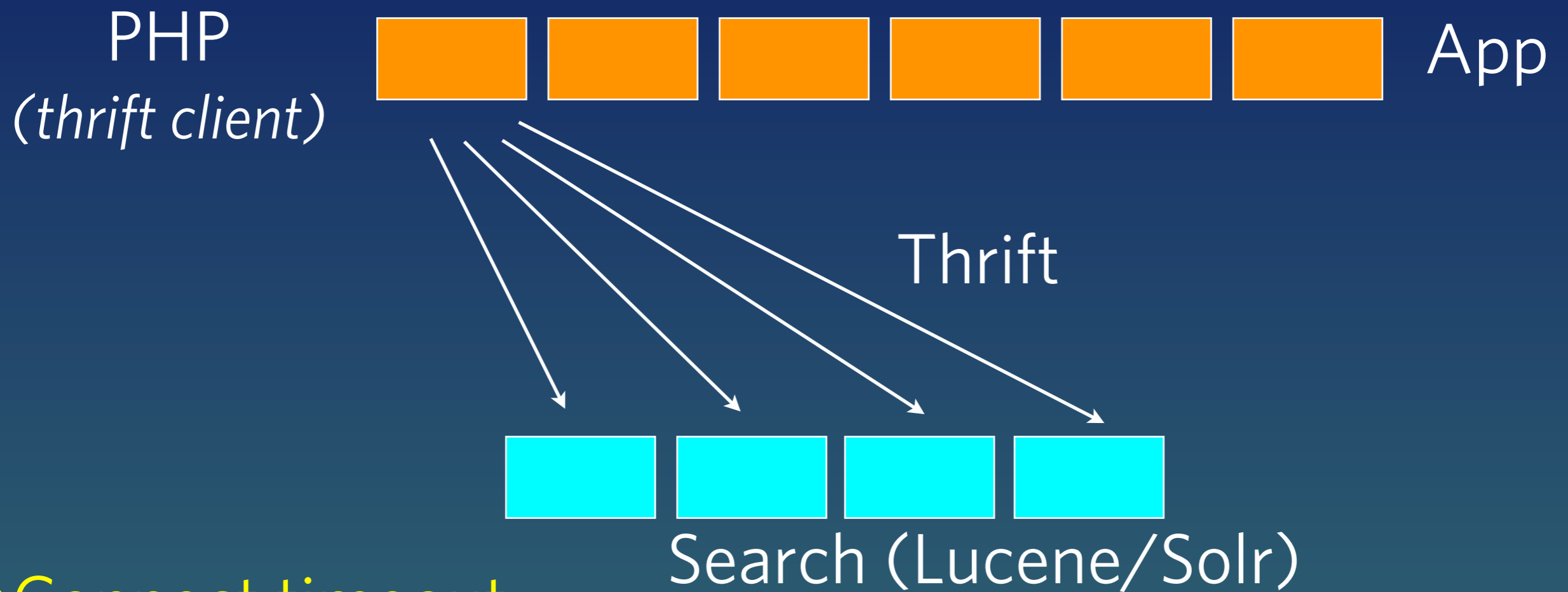
Spatial Redundancy

Active/Passive

Spatial Redundancy

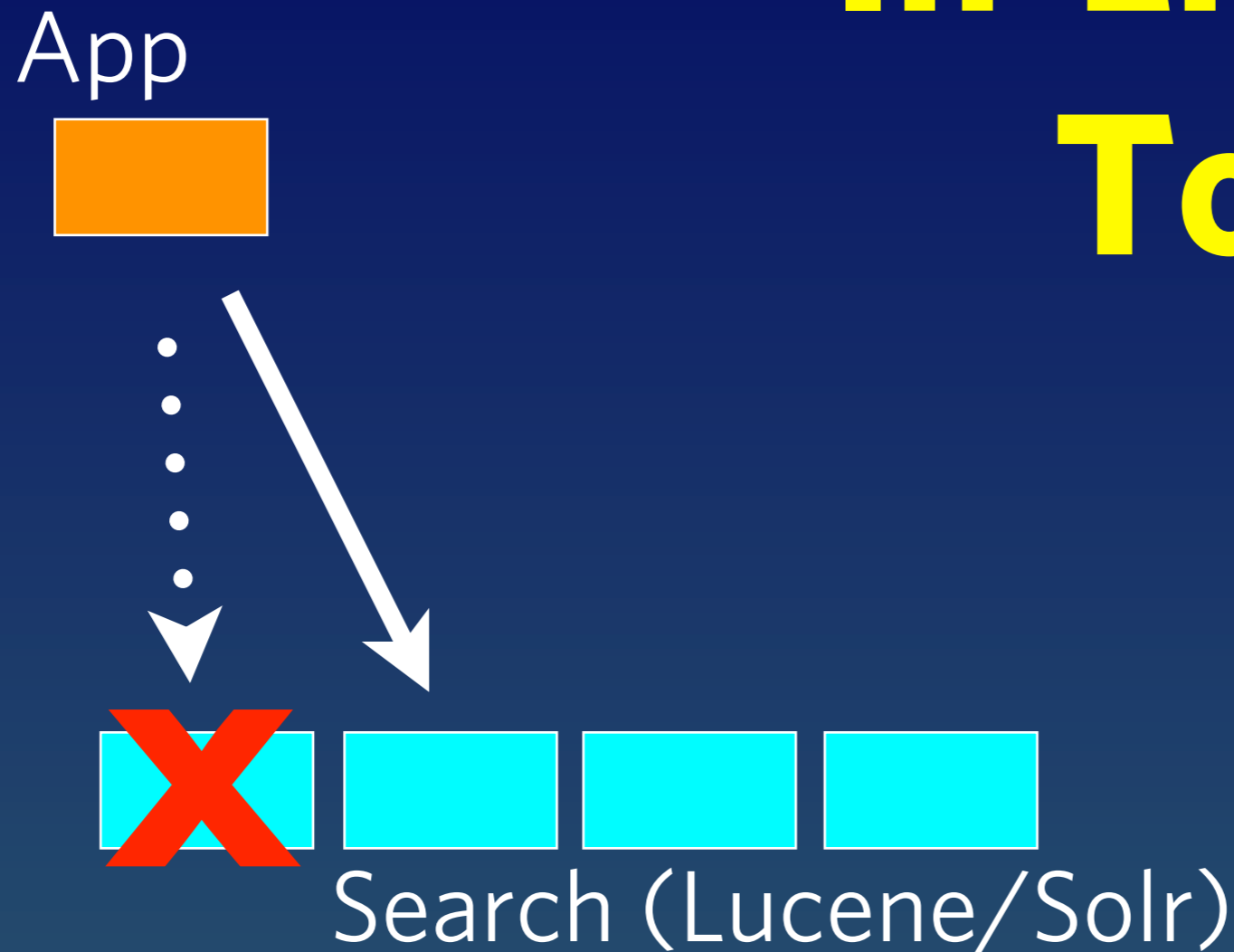
**Roaming Spare
Dedicated Spare**

In-Line Fault Tolerance



- Connect timeout
- Send timeout
- Receive timeout

In-Line Fault Tolerance



1. App attempts connection, can't
2. Caches APC user object with 60s "TTL" key=server:port
3. Moves to next server in rotation, skipping any found in APC

In-Line Fault Tolerance

APC
Opcode Cache

Refresh Data View Host Stats System Cache Entries Per-Directory Entries **User Cache Entries** Version Check Clear user Cache

Attribute	Value
Info	thrift_failtime:search06.ny4.etsy.com:9012~
Ttl	0 Hostname Port
Type	user
Num Hits	6 (0.00%)
Mtime	2010/10/27 19:05:36
Creation Time	2010/10/27 19:05:36
Deletion Time	None
Access Time	2010/10/27 19:05:40 Epoch+60 seconds
Ref Count	0
Mem Size	240
Stored Value	1288206336

<http://thrift.apache.org/>

</lib/php/src/TSocketPool.php>

In-Line Fault Tolerance

Pros:

- Distributed checking and perspective**
- Handles transient failures**
- Auto-recovery**

Cons:

- Onus is on the app for implementation**

Fault Tolerance

Nagios Event Handlers

Attempt to recover from specific conditions

Chain together recovery actions

[http://nagios.sourceforge.net/docs/3_0/
eventhandlers.html](http://nagios.sourceforge.net/docs/3_0/eventhandlers.html)

If (fault X) then

HUP process; re-check

If (OK) then notify+exit

ELSE

Hard restart process; re-check

If (OK) then notify+exit

ELSE

Remove from production; notify+exit

***How many seconds of errors
can you tolerate serving?***

Fail Closed

When fault is found, and can't be recovered or masked, operations cease to protect the rest of the system from damage.

Depth and Dependencies

Monitor



Load Balancers



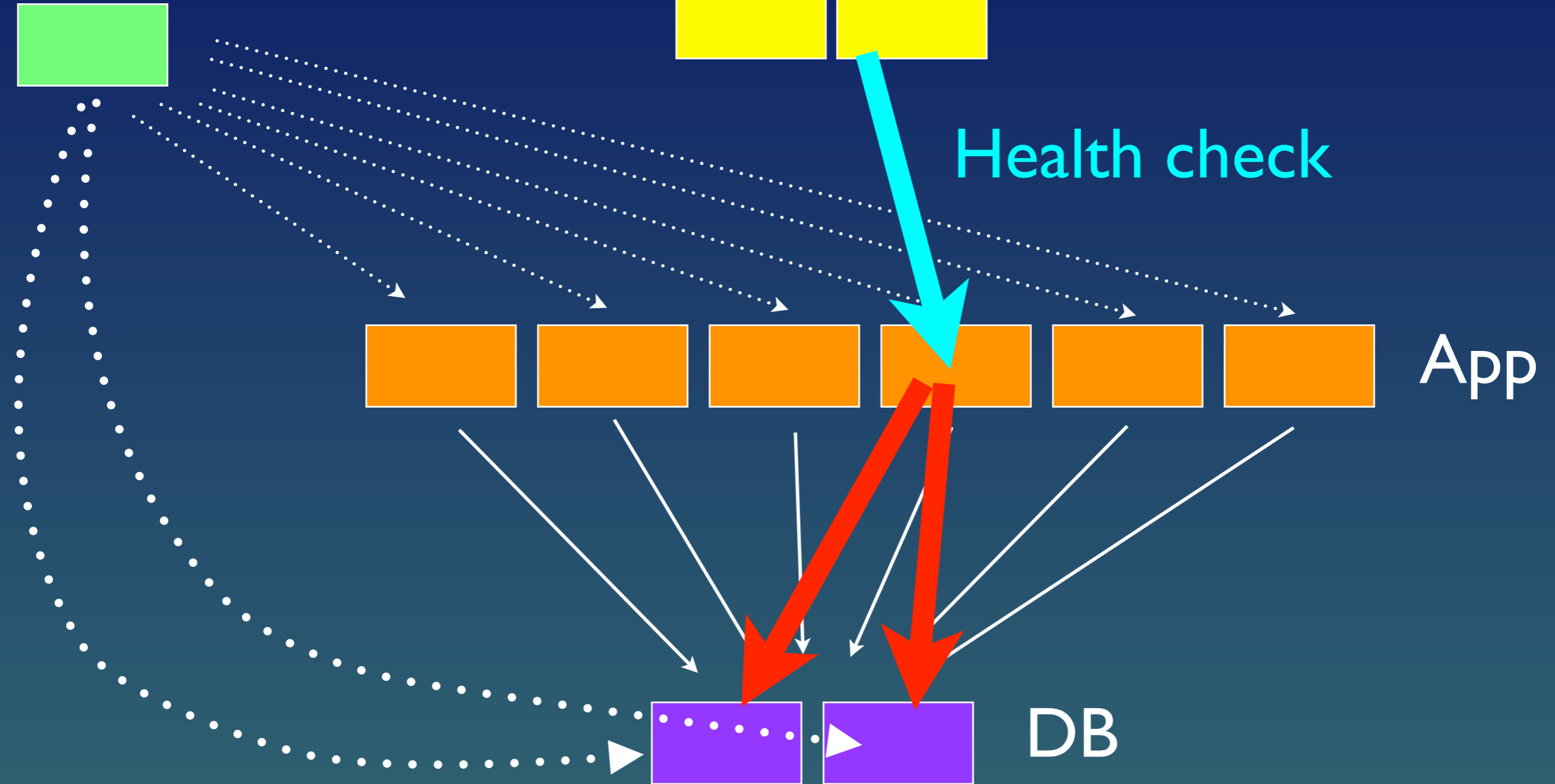
Health check



App



DB



Depth and Dependencies

WARNING:

Don't be too

crazy

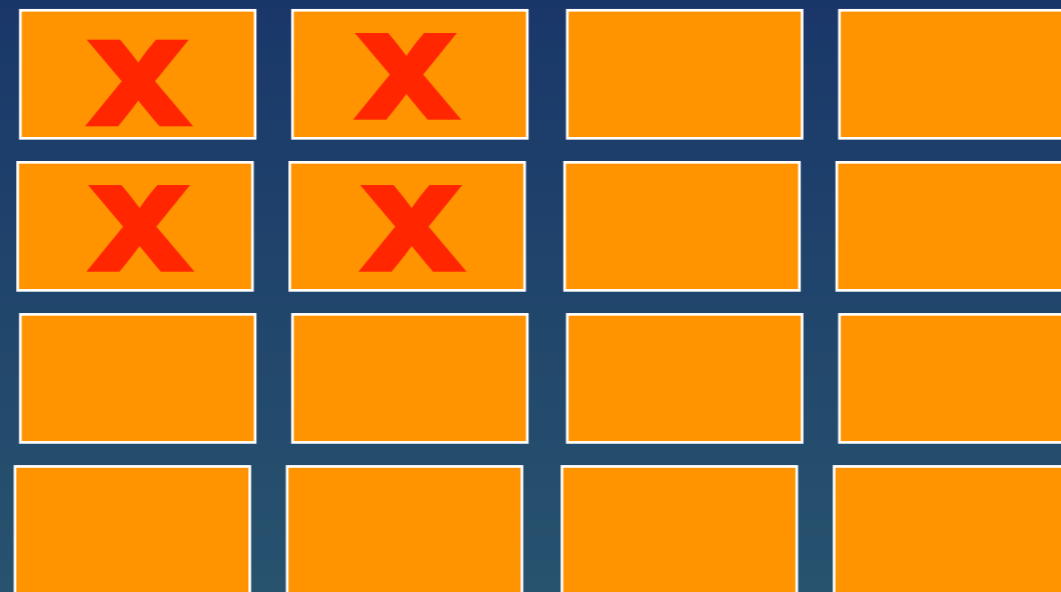
Health check

App

DB

Fail Closed

Aggregate Cluster Checking



```
If (clusterfail > 25%) then notify+exit  
ELSE OK
```

Fail Open

When a fault happens, and can't be masked or recovered, operations continue without the feature.

Fail Open

Example 1 at **Etsy**: Geo Targeting

Housewares 10,340
Jewelry 3,034
Glass 2,198
Art 2,104
Children 714
Furniture 550
Ceramics and Pottery 511
See more ▶

Shop Location
 All locations
 Brooklyn, NY 133
Change
 Ships to United States
Change

Link Lamp
Interactive 8-bit Question Blo...
\$74.99 USD

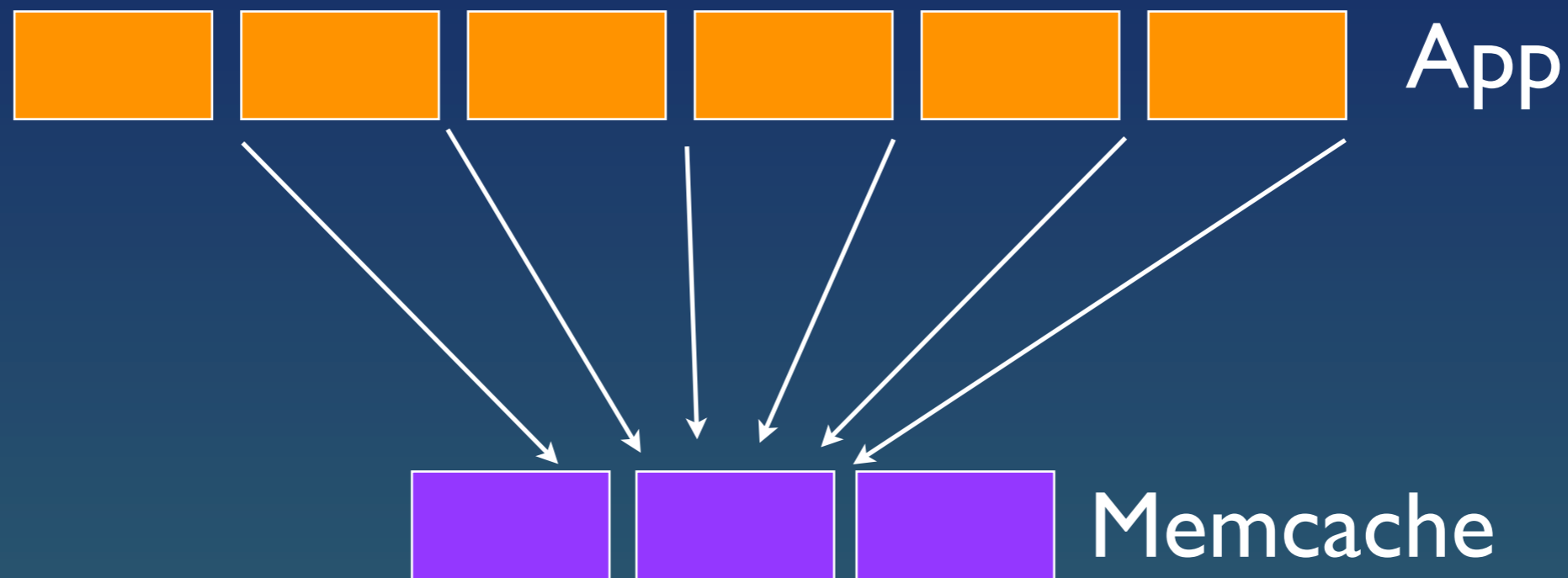
Lamp Shade Only
Unique handmade lamp mad...
Indu...

For sale near you

50ms Internal SLA on guessing location via client IP.
If >50ms, we just don't show local results.

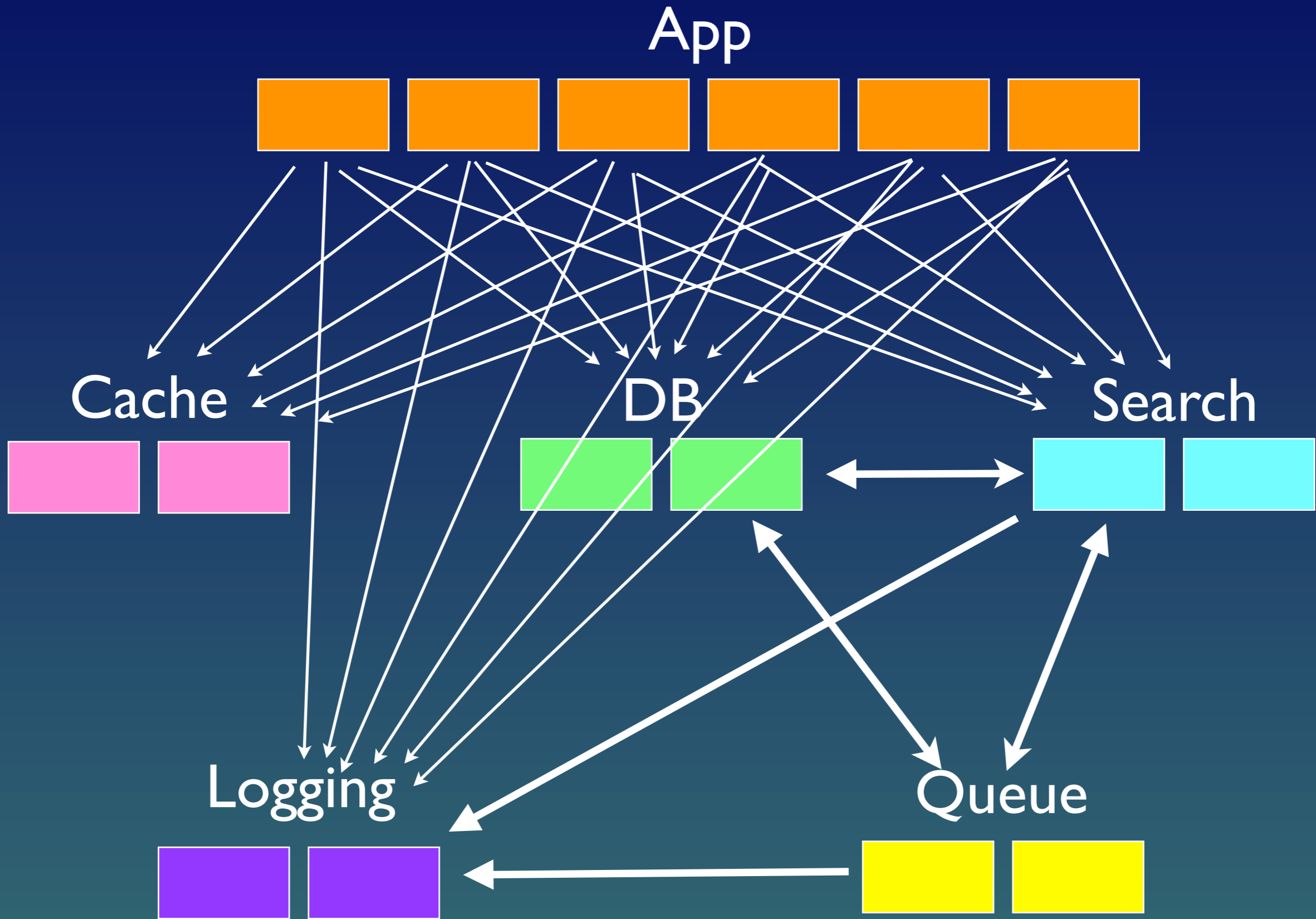
Fail Open

Example 2 at Etsy: Rate Limiting



Internal SLA on incrementing counters+checking totals.
If $>SLA$, we let the action continue, and throw fire-and-forget counter if we can.

SYSTEMIC



App



Cache



DB



Search



Logging



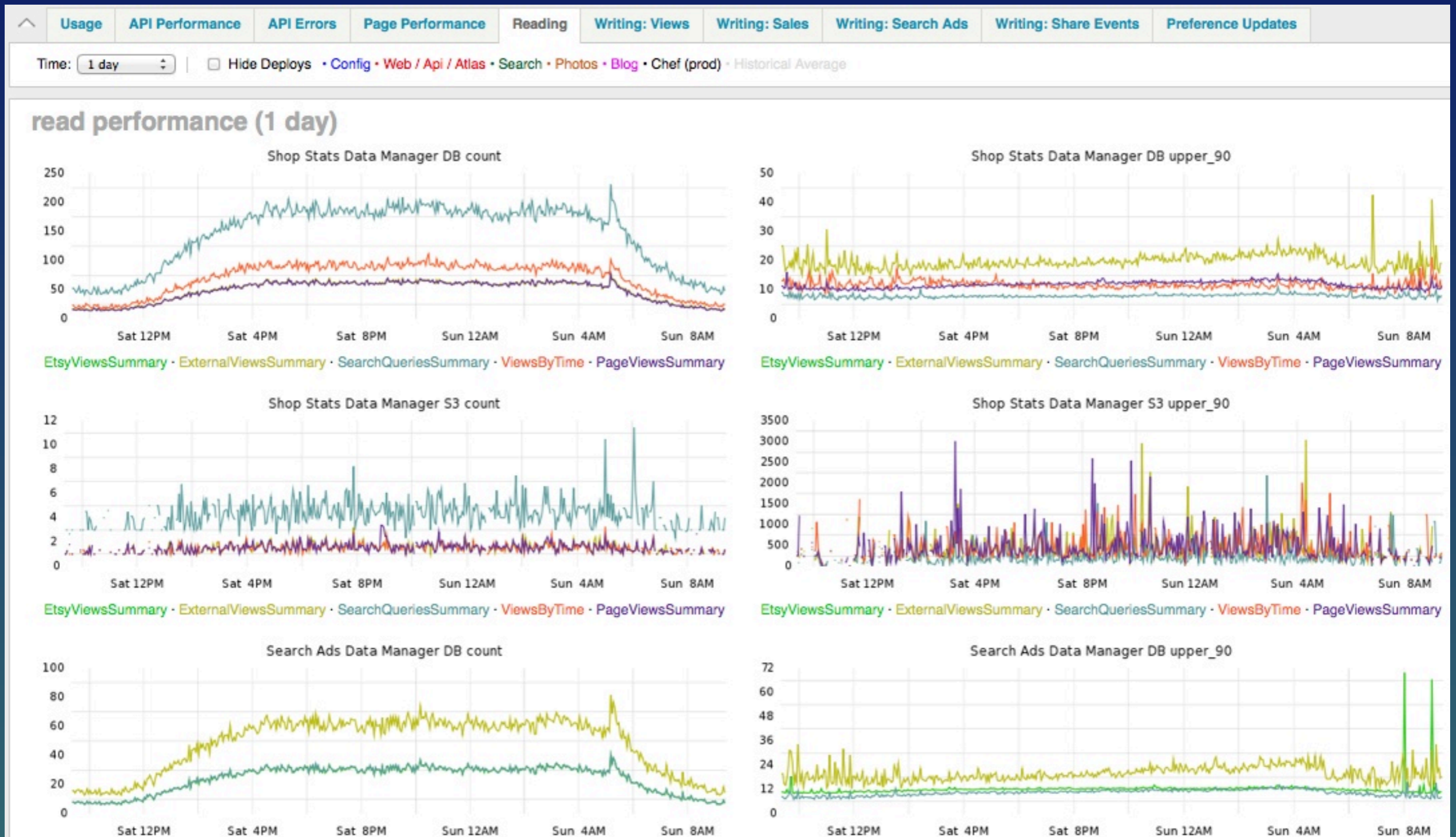
Queue



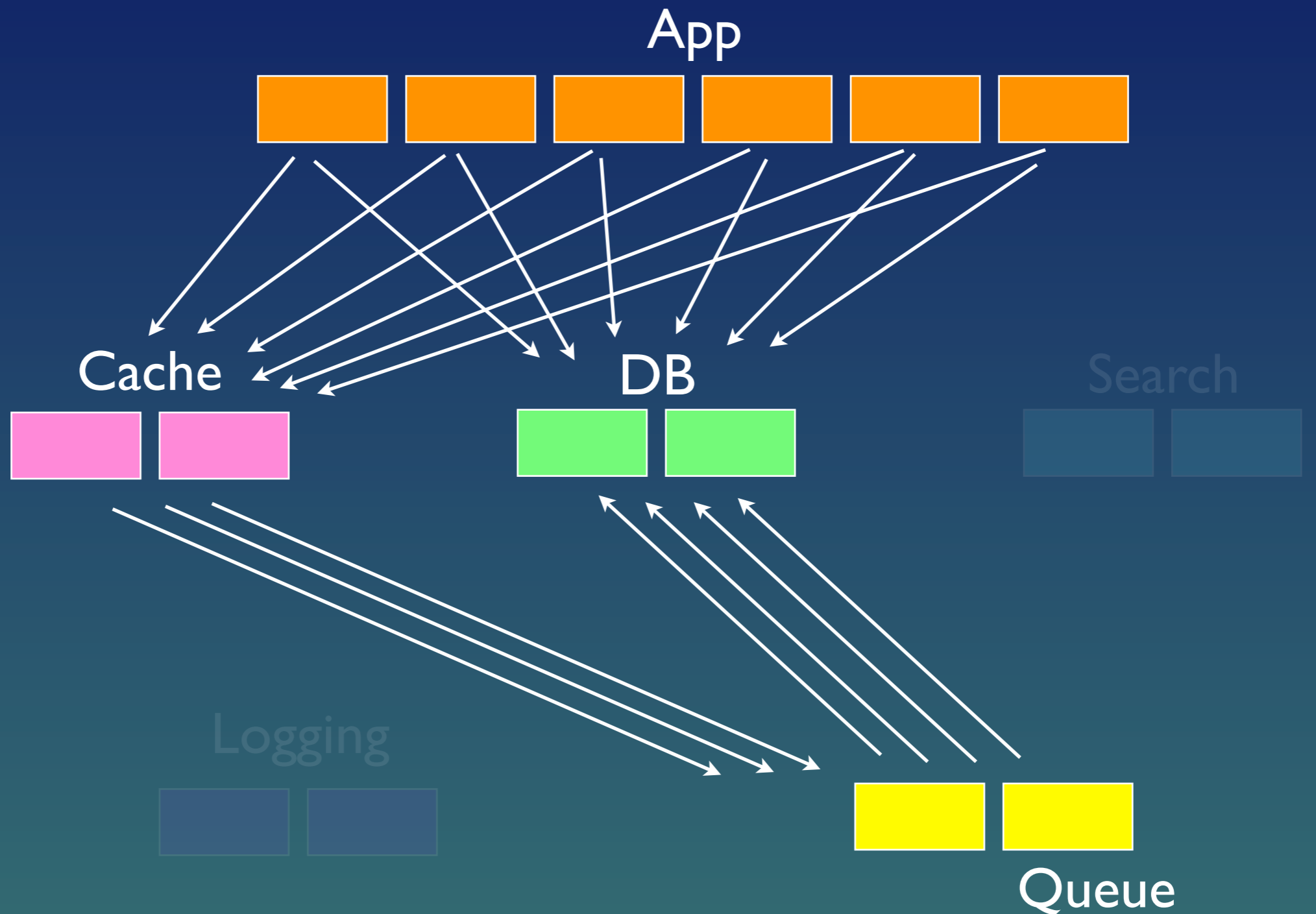
Functional Resonance



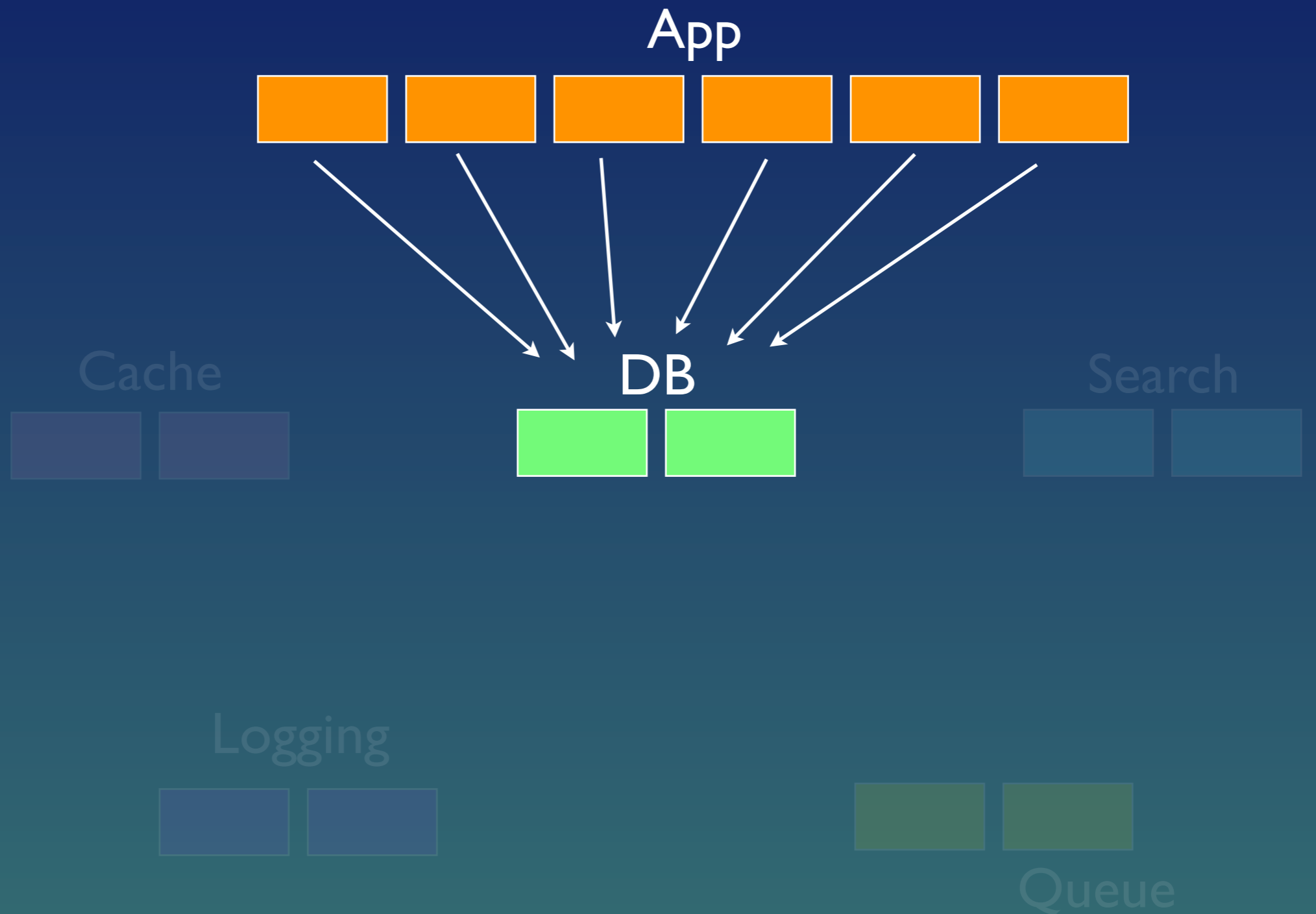
Shop Stats



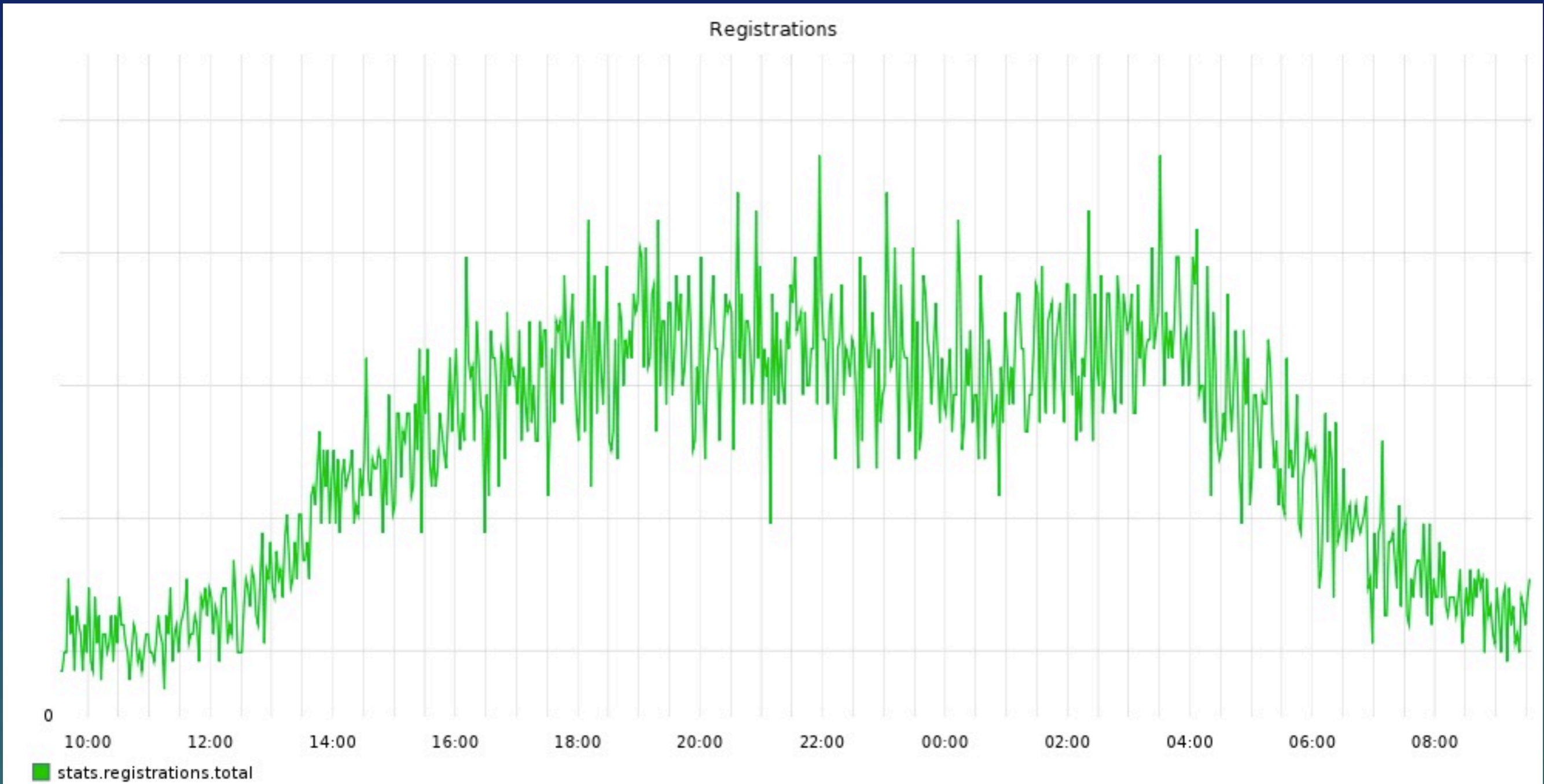
Shop Stats



Registration



Registration



Shop Stats
Logins
Registrations
Checkout
New Listings
Photos
Search
API
Rate limiting
Data Analysis
Search
A/B analysis

Page performance
Search Ads
Editorial content
Email systems
Feedback
Messaging/Convos
Activity Feeds
Circles
Shipping
Mobile
Internationalization
Testing
Fraud

Systemic

Application/Functionality Health



Componential/Resource Health

Four Cornerstones

Erik Hollnagel

(Anticipation)

Knowing
What
To Expect

(Response)

Knowing
What
To Look For

Knowing
What
To Do

Knowing
What
Has Happened

(Monitoring)

(Learning)

Anticipation

- **During design of architecture**
- **During choice of technologies**
- **During design of monitoring and metrics**

TRADE-OFFS

**“What could possibly
go wrong?”**

REQUISITE

IMAGINATION

Possible Foreseeable Situations

**Situations
Considered By
Novice Designer**

**Situations
Considered By
Average Designer**

**Situations
Considered By
Expert Designer**

Adamski and Westrum, 2003

Anticipation

Failure Mode Effects Analysis (FMEA)

http://en.wikipedia.org/wiki/Failure_mode_and_effects_analysis

Failure Mode Effects and Criticality Analysis (FMECA)

<http://en.wikipedia.org/wiki/>

[Failure_mode,_effects,_and_criticality_analysis](#)

Architectural reviews

Go or No-Go meetings

“Game Day” exercises

Anticipation

Servers

Networks

Software

Applications

Monitoring

Metrics

Traffic

PEOPLE

(Anticipation)

(Response)

Knowing
What
To Expect

Knowing
What
To Look For

Knowing
What
To Do

Knowing
What
Has Happened

(Monitoring)

(Learning)

THE END