# TIM GROUP

**UNLOCK YOUR ALPHA.**

## TIM IDEAS

The most comprehensive global alpha capture network for trade idea distribution, evaluation and management

### BUYSIDE
Identify exclusive, market-beating opportunities quickly and manage brokers wisely.

### SELLSIDE
Establish your value, broaden your network of connections and increase commissions.

## TIM FUNDS

Superior analytics and portfolio management for funds of hedge funds directly or working with administrators and custodians.

### OVERVIEW
TIM Funds transforms the data that administrators and custodians already hold into powerful, accurate, analytics.

### CONNECT
TIM Funds Connect is a series of APIs that enable users to interact seamlessly with in-house or third party technology or data systems.

## 2005

**TIM IDEAS**

**The most comprehensive global alpha capture network for trade idea distribution, evaluation and management**

### BUYSIDE
Identify exclusive, market-beating opportunities quickly and manage brokers wisely.

### SELLSIDE
Establish your value, broaden your network of connections and increase commissions.

## 2006

**TIM FUNDS**

**Superior analytics and portfolio management for funds of hedge funds directly or working with administrators and custodians.**

### OVERVIEW
TIM Funds transforms the data that administrators and custodians already hold into powerful, accurate, analytics.
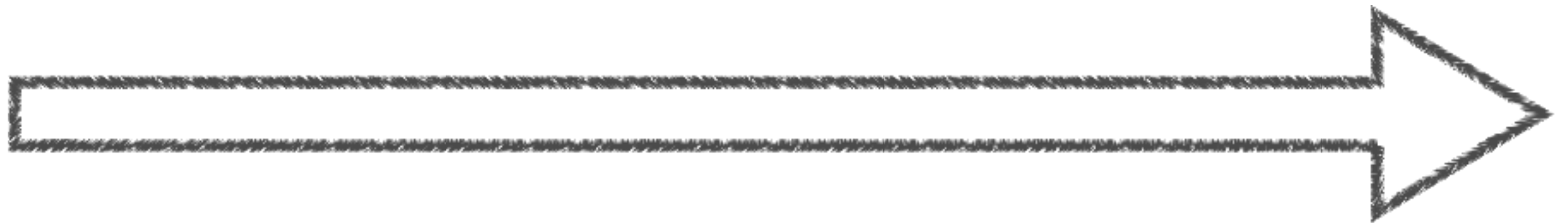
### CONNECT
TIM Funds Connect is a series of APIs that enable users to interact seamlessly with in-house or third party technology or data systems.

**2007** **2008** **2009**

- Adopt pair programming

- Add extensive testing including TDD
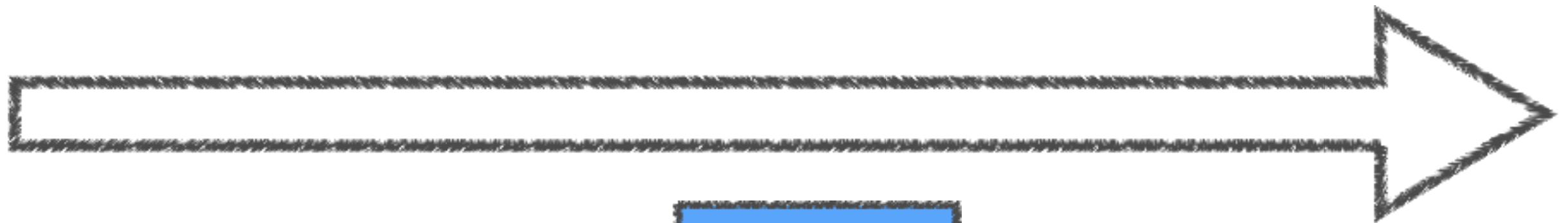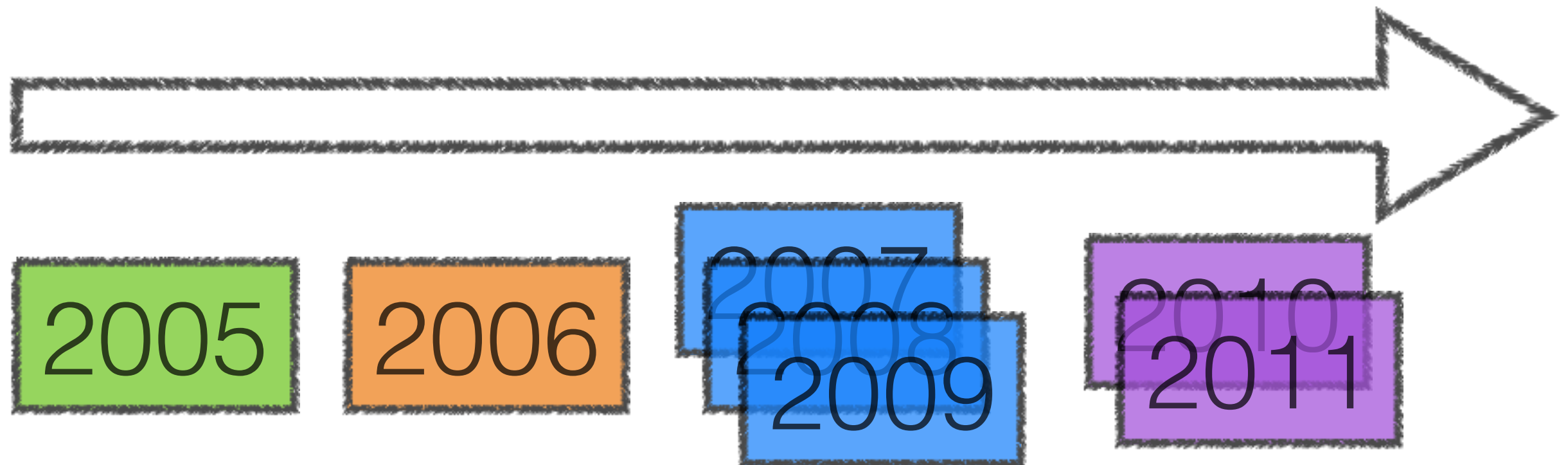
- Support with elaborate continuous integration

**2005** **2006**

2012

- At migration complete merged sysadmin project team with Ops and Tools team to create new Infrastructure team

- Went live from our new data centers in January

2005   2006   2007   2008   2009   2010   2011

- 20May12 3:18am: pg-timapp-0U2 TIM Tomcat Direct alert, resolved itself within 5 minutes. OoH effort: BoD: 5m.
- 18May12 8:15am: portal on macau alerting because it had hit the open files limit. Restarted portal to resolve the immediate issue. OoH effort: Bob: 45m → passed to Kei/Mehul.
- 18May12 5:20am: mail queue on pg-timapp-005 over threshold, 78/50 messages in queue. Resolved itself. OoH effort: Bob: 10m
- 17May12 1:18am: Nagios alerts for portal returning http 500. Investigated, restarted macau. OoH effort: Bob: 1h.
- 18May12 8:00am: installing MBA Test Environment TC certificate: OoH effort: Kei: ?
- 15May12 3:50am: HK office reported problems with delays in pricing for Japanese/Korean stocks—ended up being a TRKD problem. OoH effort: Bob: 1h, Jason: 30m
- 14May12 8:00pm: MBA TC certificate cutover: OoH effort: Kei: 1
- 12May12 8:00am: Standard TIM Funds and TIM Ideas release: OoH effort: Kei: 1
- 11May12: 0059: TRKD down. Resolved itself. Adam: ?
- 9May12: 0014: ApiRankingCyclicProcess failure. Resolved itself. Adam: 30m.
- 8May12: 0400: TradeRepricerOfOddSiblingEventsCyclicProcess commit failure: StaleStateException. Resolved itself. Adam: 30m.
- 7May12: 0850: Ideas Stop Loss cyclic process failure due to StaleObjectStateException Jtf: 15 min
- 6May12: 0640: TF Log Check failure due to Fx CP. Jtf: ? Tony: ? Isaac: ? Graham A.: ? Ian: ? JoeS: 1h
- 5May12 ?: *Scheduled Maintenence Jtf: ?*
- 5May12 01:00, 02:40, 06:30 : Merc and TF FX CP failures. Kei: ?
- 4May12 1700: prep for *Scheduled Maintenance. Jtf: 0.5*
- 28-29Apr12 ?: Ideas upgrade process broke replication. Bob: 10h Gary: 2h Adam: Kei: Jeff: Waseem: Joel: Jason:
- 28Apr12 1600: *Release of Funds & Ideas. Long release because trying to run both simultaneously botched the Ideas release. Jtf: 2.5 hrs, Waz: ? Kei: ?*
- 27Apr12 05:55pm: Gavin enabled CS restricted lists entitlement and wanted us to check if it was running. Bob: 30m
- 25Apr12 05:07am: TradeRepricerOfOddSiblingEventsCyclicProcess Failed: StaleStateException. Resolved itself. Bob: 30m
- 22Apr22 ?: TIM Connect failures, turned out to be remote end, not us. Adam: ? Waz: ? Jeff: ?
- 21Apr12 16:13: ExchangeRateUploader failure (Could not open a connection to SQL Server). Resolved itself. Adam: 30m
- 20Apr12 08:10am: High number of processes alert on proxy server. Adam: 2m
- 19Apr12 04:45am: PG TIM FW failure leading to lost packets and TIM outage. Adam: 1h
- 17Apr12 07:00pm: *TIM Funds maintenance release: OoH effort: Kei: 0.5*
- 15Apr12 09:50am: *TIM Ideas healthcheck: OoH effort: Kei: 1*
- 14Apr12 7:00am: *TIM Ideas database archiving: OoH effort: Kei: 4*
- 14Apr12 7:00am: *Standard TIM Funds and TIM Ideas release: OoH effort: Kei: 2*
- 14Apr12 21:00: Brixton upgrade OoH Effort: Richard: 2
- 10Apr12 10:30pm: *TIM Ideas maintenance release : OoH effort: Kei: 0.5*
- 9Apr12 Midnight: TIMConnect outages. OoH effort: Bob: 4h Adam: 1h Kei: 4 Jeff: 4 Jason: 2h Gavin: 2h
- 4Apr12 00:56: Cyclic Process failure for hourly idea rankings. OoH effort: Bob 4h (+ Brian Roberts 30m, Gavin 1h)
- 31Mar12 04:45: TradeRepricerOfEvenSiblingEventsCyclicProcess failure. Resolved itself. OoH effort: Adam: 10m
- 30Mar12 8:00pm: *TIM Funds maintenance release : OoH effort: Kei: 0.5*
- 29Mar12 04:48: TradeRepricerOfOddSiblingEventsCyclicProcess failure. Resolved itself. OoH effort: Adam: 10m
- 25Mar12 9:00am: *Standard TIM Funds and TIM Ideas healthcheck: OoH effort: Kei: 1*
- 24Mar12 07:00: *Standard TIM Funds and TIM Ideas releases + other maintenance: OoH effort: Kei: 3 GaryR: 2*
- 20Mar12 07:00: Abacus outage: OoH effort: Kei: 1
- 19Mar12 02:28: Today returns generation issue: OoH effort: Kei: 5h
- 19Mar12 00:54: Today returns generation issue: OoH effort: Adam: 4h
- 15Mar12 22:30: *TIM Ideas maintenance release: OoH effort: Kei: 0.5*
- 14Mar12 22:30: bounce TC/tomcat for Marshall Wace: OoH effort: Kei: 0.5
- 11Mar12 9:00am: *Standard TIM Funds and TIM Ideas healthcheck: OoH effort: Kei: 0.5*
- 10Mar12 9:00am: *Standard TIM Funds and TIM Ideas releases: OoH effort: Kei: 1*
- 06Mar12 1730: Email backups failing due to wrong VMware disk size. Implemented recovery solution after discussing with Paul Howard/Department 7. OoH effort: Craig + Richard: 2.5h
- 29Feb12 11pm: bos-f17-cloud-2 (fcsentiment + shortel) same issue as bos-f17-cloud-1. OoH effort: Bob: 0.5h
- 28Feb12 12am: bos-f17-cloud-1 (fcdomain + others) 208+ day uptime bug (cpu soft lockup/"divide by zero"). OoH effort: Bob: 1.5h
- 25Feb12 8:30am: *Standard TIM Funds and TIM Ideas releases: OoH effort: Kei: 0.5*
- 20Feb12 0530: broken replication on dominica... denmark running with old data. took site down. restore backup from tuesday and replayed binlog. Lost data since Saturday. OoH effort: Adam: 3h30m Gary: Kei:
- 19Feb12 2300 : broken replication on dominica. OoH effort: Adam: 1h
- 18Feb12 2001: broken replication on dominica. OoH effort: Adam: 2h
- 18Feb12 1311: Merc problem with interday pricing. Reuters IP addressed changed, making it unreachable due to our firewall rules. OoH effort: Adam: 1h30m
- 18Feb12 08:00am: *TIM Funds fall back to Powergate : OoH effort: Jarek: tbd , Kei: tbd*
- 16Feb12 10:30pm: *TC feature release : OoH effort: Kei: 0.5*
- 12Feb12 8:30am: TIM Funds production database server (Denmark) went down. After receiving alerts for an hour ForLinux called the Support number and reached Jarek. Denmark was not reachable over iLO (why?) and Kei followed the documented directions and switched Funds to DR site. Gary was able to connect to iLo and reboot Denmark but it did not come up cleanly. Jarek fixed some partition stuff (clean this part up?) and Denmark recovered. Though we did not have sync-binlog enabled (fix this?) the positions matched and Jarek was able to start replicating from OY back to Denmark. Other work: allowing mail out from OY; changing Nagios alerts. Root Cause: ? Outage duration: 2.5 hrs, 0800–1030 OoH effort: Jarek: 4, Kei: 1, Gary: , Craig: 1
- 11Feb12 9:00am: *Standard TIM Funds and TIM Ideas releases: OoH effort: Kei: 3*
- 10Feb12 10:00pm: *Scheduled SugarCRM and Drupal upgrades. OoH effort: Craig 1.5*

# Devops is About CAMS

- **Culture**
  People and process first. If you don't have culture, all automation attempts will be fruitless.

- **Automation**
  This is one of the places you start once you understand your culture. At this point, the tools can start to stitch together an automation fabric for Devops. Tools for release management, provisioning, configuration management, systems integration, monitoring and control, and orchestration become important pieces in building a Devops fabric.

- **Measurement**
  If you can't measure, you can't improve. A successful Devops implementation will measure everything it can as often as it can… performance metrics, process metrics, and even people metrics.

- **Sharing**
  Sharing is the loopback in the CAMS cycle. Creating a culture where people share ideas and problems is critical. Jody Mulkey, the CIO at Shopzilla, told me that they get in the war room the developers and operations teams describe the problem as the enemy, not each other. Another interesting motivation in the Devops movement is the way sharing Devops success stories helps others. First, it attracts talent, and second, there is a belief that by exposing ideas you can create a great open feedback that in the end helps them improve.

John Willis: What Devops Means to Me
http://is.gd/B6z7E7

- **Culture**
  People and process first.

- Ops was application support, the semi-technical buffer between production & users and the developers

- Company bottleneck seen as developer time

- Enabling architectural change was the motivating factor

John Willis: What Devops Means to Me
http://is.gd/B6z7E7

- Continuous deployment for new components put in place by development

## • Measurement
If you can't measure, you can't improve.

- "Metrics": a two week block on the Gantt chart

- Application metrics are (almost) all technical

- Business metrics come from the same nightly "ops project" scripts as before the transition

- "Shared problems are the basis of teamwork"

- Root Cause Analysis is our established mechanism for addressing problems

- Success on shared problems have been creating a virtuous cycle

# The wake up call

Fault tolerant infrastructure

**+**

Fault tolerant architecture

**=**

Outage within 15 minutes of going live

- New component introduced to have calculation code against a slave database

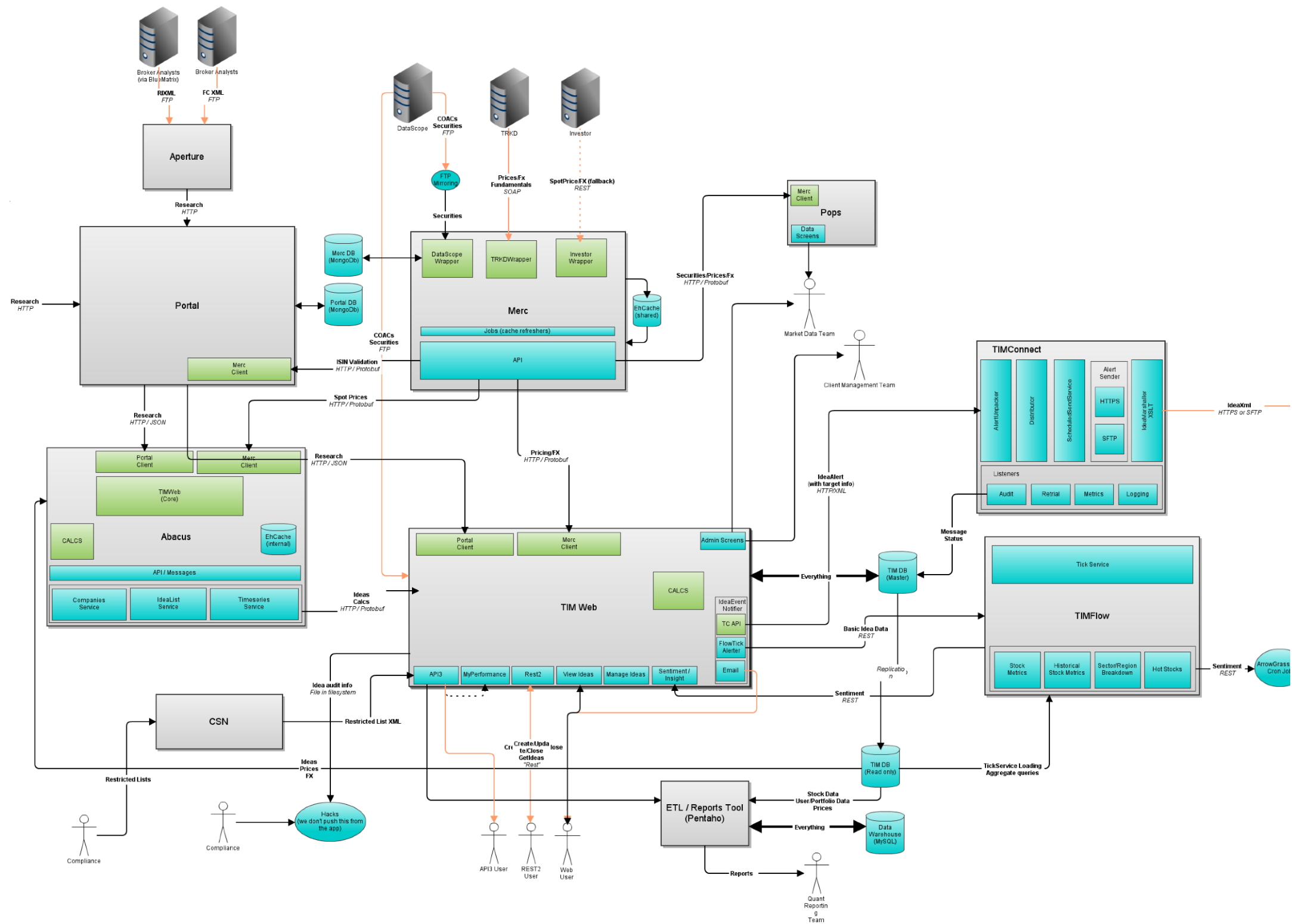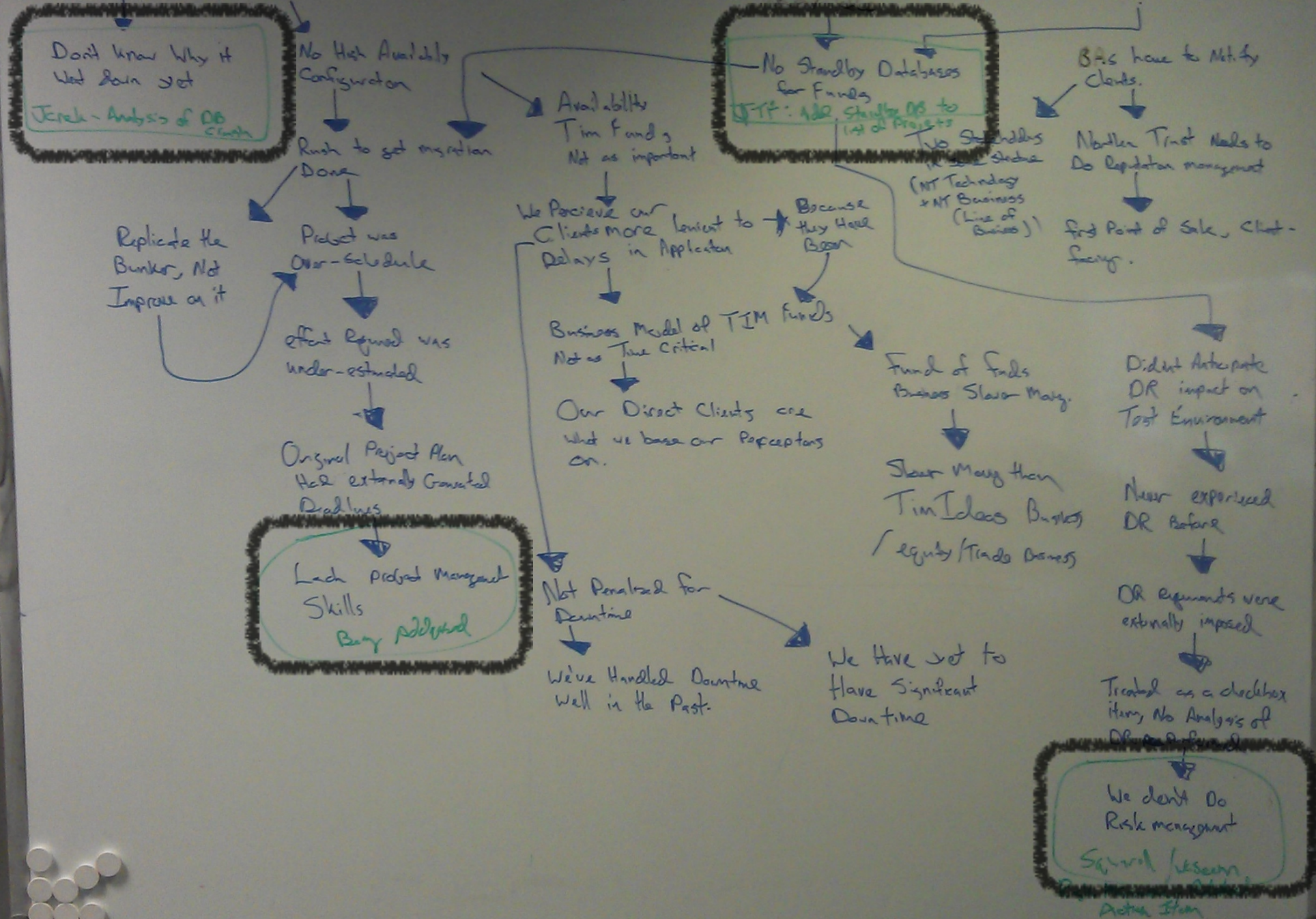- API code transparently calling new component; fallback to old code on failure

- Two component instances behind a load balancer

**Broker Analysts (via Blu Matrix)**
RIXML
*FTP*

**Broker Analysts**
FC XML
*FTP*

**DataScope**
COACs Securities
*FTP*

**TRKD**
Prices/Fx Fundamentals
*SOAP*

**Investor**
SpotPrice,FX (fallback)
*REST*

**Aperture**

Research
*HTTP*

**Portal**

Research
*HTTP*

Merc DB (MongoDb)

Portal DB (MongoDb)

Merc Client

**Merc**

DataScope Wrapper

TRKDWrapper

Investor Wrapper

Securities

EhCache (shared)

Jobs (cache refreshers)

API

COACs Securities
*FTP*

ISIN Validation
*HTTP / Protobuf*

Spot Prices
*HTTP / Protobuf*

Research
*HTTP / JSON*

Securities,Prices/Fx
*HTTP / Protobuf*

**Pops**
Merc Client
Data Screens

Market Data Team

Client Management Team

**Abacus**

Portal Client

Merc Client

TIMWeb (Core)

CALCS

EhCache (internal)

API / Messages

Companies Service

IdeaList Service

Timeseries Service

Research
*HTTP / JSON*

Ideas Calcs
*HTTP / Protobuf*

Pricing/FX
*HTTP / Protobuf*

**TIM Web**

Portal Client

Merc Client

CALCS

Admin Screens

IdeaEvent Notifier

TC API

FlowTick Alerter

Email

API3

MyPerformance

Rest2

View Ideas

Manage Ideas

Sentiment / Insight

**TIMConnect**

AlertDispatcher

Distributor

Scheduler/SendService

Alert Sender
HTTPS
SFTP

IdeaMarshaller XSLT

Listeners

Audit

Retrial

Metrics

Logging

IdeaXml
*HTTPS or SFTP*

IdeaAlert (with target info)
*HTTP/XML*

Message Status

TIM DB (Master)

Everything

**TIMFlow**

Tick Service

Stock Metrics

Historical Stock Metrics

Sector/Region Breakdown

Hot Stocks

Basic Idea Data
*REST*

Replication

Sentiment
*REST*

ArrowGrass Cron Job

Sentiment
*REST*

**CSN**

Restricted List XML

Idea audit info
*File in filesystem*

Ideas Prices FX

Restricted Lists

**Hacks** (we don't push this from the app)

Compliance

Compliance

Create/Update/Close GetIdeas "Rest"

TIM DB (Read only)

TickService Loading Aggregate queries

**ETL / Reports Tool (Pentaho)**

Stock Data
User/Portfolio Data
Prices

Everything

**Data Warehouse (MySQL)**

Reports

Quant Reporting Team

API3 User

REST2 User

Web User

grams at gliffy.com

wiki.timgroup.com/results_of_the_16_february_2012_rca_on_the_tim_fun   ↻    Google

# [[results_of_the_16_february_2012_rca_on_the_tim_funds_outage_16_february_2012]]

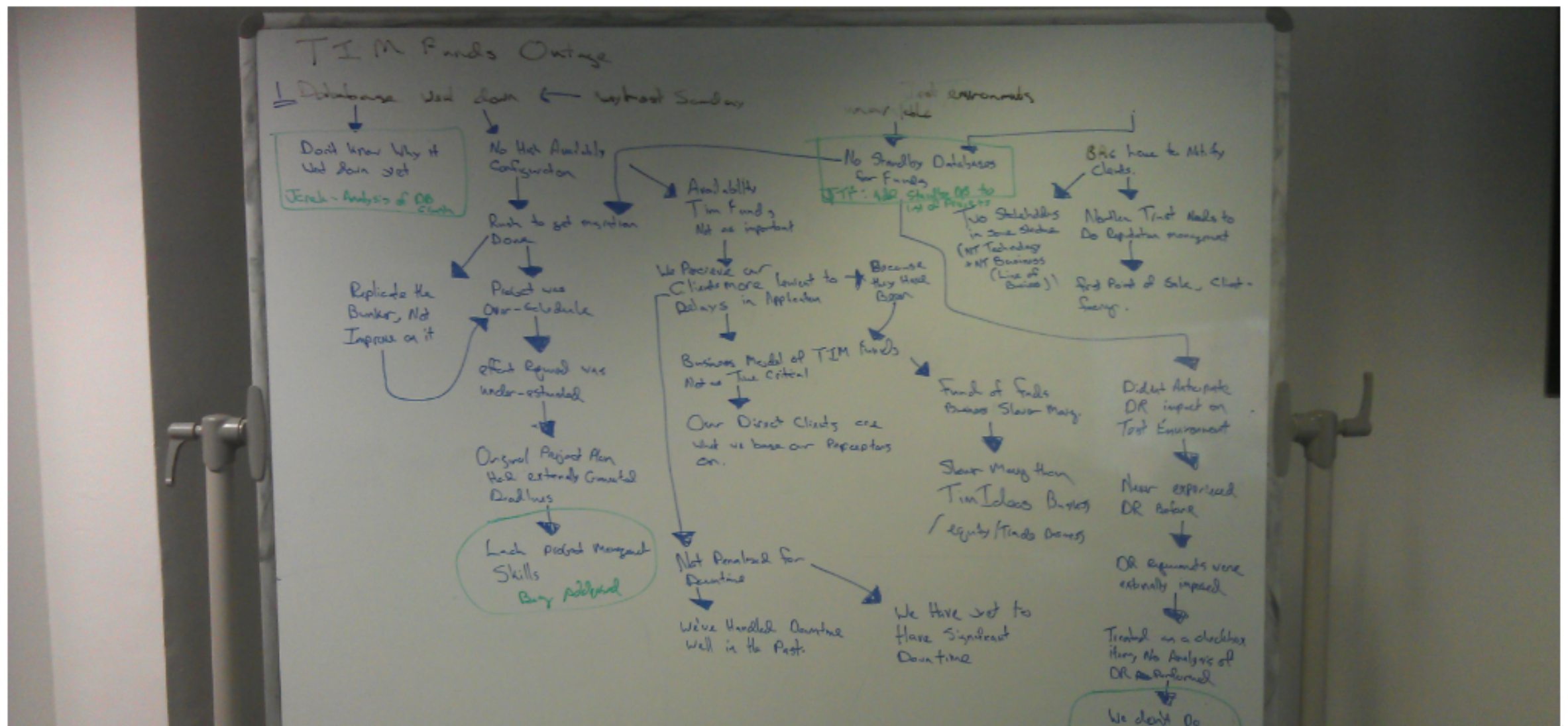## TIM GROUP WIKI

Edit this page   Old revisions         Recent changes   [    ]   Search

Trace: • start • addressing_some_limitations_of_the_5_whys • root-cause_analyses • results_of_the_16_february_2012_rca_on_the_tim_funds_outage_16_february_2012

## Actions

☐ Analysis of DB Crash - **Jarek**

☑ ~~Add standby databases to list of tasks~~ - **JTF**

☑ ~~generate Risk Management action item~~ - **Squirrel / Waseem**

Edit

## Picture

# TIM Ideas

CSN not working in production - 11 May 2012

Abacus was a version behind LATEST in production 08 May 2012

DB Slave Replication UUID problem in production - 01 May 201

Connect out of file handles 09 April 2012

Unpredictable race condition caused delivery failure of single message to multiple target systems (TIMConnec

NPE on My Performance caused by null value in portfolio fees 29March2012

Historical prices missing for non-trading days 24March 2012

Abacus caused site outage when turned on

Hourly rankings ran slowly following 2012Jan28 release 02February2012

Lock Currency bugs 01 February 2012

Incorrect time series data caused erroneous reports 31 January 2012

Dec2011Bug-Ranking cyclic processes ran all rankings, not limited to hourly or non-hourly per cyclic proces
18January 2012

RCA for table lock in production due to running SAC SQLs in production 17January2012

TechnicallyfocusedRCAdrivenbytheSACrelease17January2012

Tim funds Data Loss on Restore from Disaster Recovery Site

**Column 1:**
Customer Confidence/Satisfaction Impacted
↓
Customer was first to Discover problem
↓
Hard for us to know Data is New or Old
↓
Database monitoring Watches Database State But Not the Data itself
↓
[ Data Monitoring Pushed to Phase 2 RCA ]
Covered in Previous RCA
Tim funds 16 Feb 2012
↓
[ No High Availability Setup ]
16 Feb 2012 RCA
↓
No Resilience / Degraded Service Capacity
↓
Not Sure if it's Applicable to Tim funds Domain
↓
...the Possibility
↓
[ Risk Management Philosophy RCA ]

**Column 2:**
Cost to Customer
↓
Customer (key) Missed Deadline
↓
Could Not Work on Site while it was Down and/or Unsure for Potentially Missing Data
↓
Replication from DR to Production has Not Worked or has Never Worked
↓
Never tested DR to Production Replication
↓
DR testing Scenario Did Not include Replication Back to Prod
↓
Didn't Have a Realistic Scenario of Restoring Production
↓
[ Inherited the DR Scenario from the Bunker Without Reconsideration / Reevaluation Previous RCA ]

Re-iterate Squirrell Cultural Action from Previous RCA

– Schedule DR Scenario Review – JTF
– Schedule failure Analysis for TIM funds – JTF

**Column 3:**
Potential Loss To Tim Group
↓
SLA for 99% uptime
↓
Client Asked for SLA in contract
↓
Client's Business Negatively Impacted by Site Downtime
↓
Loss of Reputation

Wanted to Replicate Bunker Setup
↓
[ Replicated the Bunker Setup RCA 16/02/12 ]

Ask Squirrell to Explain the Reason function Accounts for Human Costs – JTF

**Column 4:**
Lost Weekend Time !
↓
Woken up by Duplicate Primary Key Alerts
↓
Replication Stopped
↓
Master DB was a Week out of Sync
↓
Database is Removed on Replication 😟 (extra Bad)
↓
Shared Slaves / Multiple-Purposed
↓
Not enough Hosts To Run Slaves as Single-Purposed
↓
Can't use Virtualization
↓
Current firewall configuration Makes implementing Virtualization Too expensive

**Column 5:**
Team Prioritized Work Disrupted !

We Relied on Bunker 24/7 Support to Handle Some of the out-of-hrs Load.
↓
Value Customers over employees
↓
We Still Have Startup Values
↓
founders No longer in the Line-of-fire

**Column 6:**
Monday Morning Stress
↓
failed to DR
↓
Single Point of failure at Database
↓
[ No Stand-by Database for TIM funds ]
↓
Did Not Consider Impact of outage on the Team
↓
Culturally We Have Dismissive of the Human Costs.
↓
We only feel the Pain When the People we Lean on are Gone.
↓
We would Rather Current Revenue Save Costs.

We don't consider the Benefit of Reducing friction

# Connect Failure Analysis

## questions from tim connect failure analysis

### email failures

- do we know if email is failing to send, to alert us of errors?
  - This is all handled by log4j and it should be logging email failures.

- can we fail back to detecting & notifiying about errors via log file
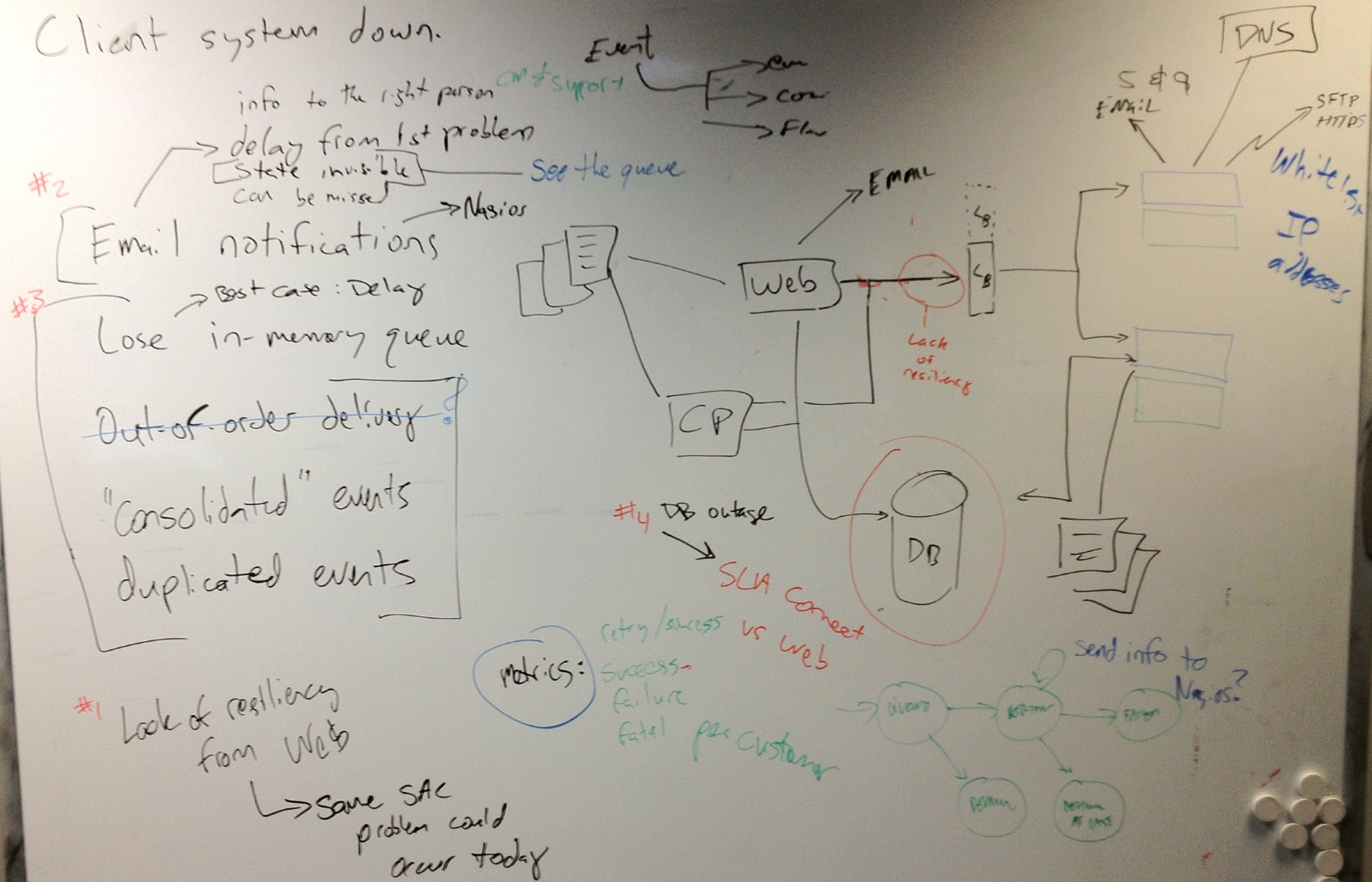  - Yes this can be configured in log4j.properties file.

?

### database failures

- if insert/update fails what happens to the connect message processing?
  - Message is processed and the persistence error is logged and emailed.
- can the insert fail without notification? (lost message)
  - If an insert fails, the error is logged and emailed.

### crash recovery

- if connect crashes and we lose the in-memory queue state do we know what queries we need to use to find (messages that were received by connect but final update is neither permentantly failed or successfully sent.)
  - A message arrives to TC and it's de-serialized. Once the message is de-serialized a new row added to TC table. However if TC crashes in this small time-frame, that message is lost.
- if we run that query now... do we have any unsent messages?
  - select count(*) from TC_MESSAGE_TARGET_STATUS where STATUS = 'SEND_FAILED' → 0 (answer is no.)

# Client system down.

info to the right person  cont support

→ delay from 1st problem

**#2** [ State invisible ] ───── See the queue
       can be missed → Nagios

[ Email notifications

**#3** ───→ Best case: Delay

Lose in-memory queue

~~Out-of-order delivery~~

"Consolidated" events

duplicated events ]

**#1** Look of resiliency
from Web

        └→ Same SAC
           problem could
           occur today

Event ──→ Car
       └─→ Con
         → Flu

metrics: retry/success vs Web
         success-
         failure
         fatal per Customer

**#y** DB outage
        ↘ SLA Connect
          vs Web

Web ──→ EMAIL

LB

Lack of resiliency

CP

DB

S & q
EMAIL

DNS

SFTP
HTTPS

White list
IP
addresses

send info to
Nagios?

# Connect Risk Analysis

Building on the Failure Analysis below we had a discussion of the risks we currently perceive with Connect.

Edit

## primary risks

### 1. Lack of resilliency in Web component

The Web component (including cyclic processes) don't retry if they fail to reach Connect. This means that transient failures re dropped messages that – at best – require manual intervention to resolve. Because any recovery attempt is manual it is not poss meet any sort of MTTR (Mean Time To Recover) SLA that is measured in minutes. It also means we are likely to end up with r messages (create is dropped, update/close goes through).

### 2. Email notifications as error reporting

### 3. Fragility of in-memory queue

### 4. Lack of high-availability database

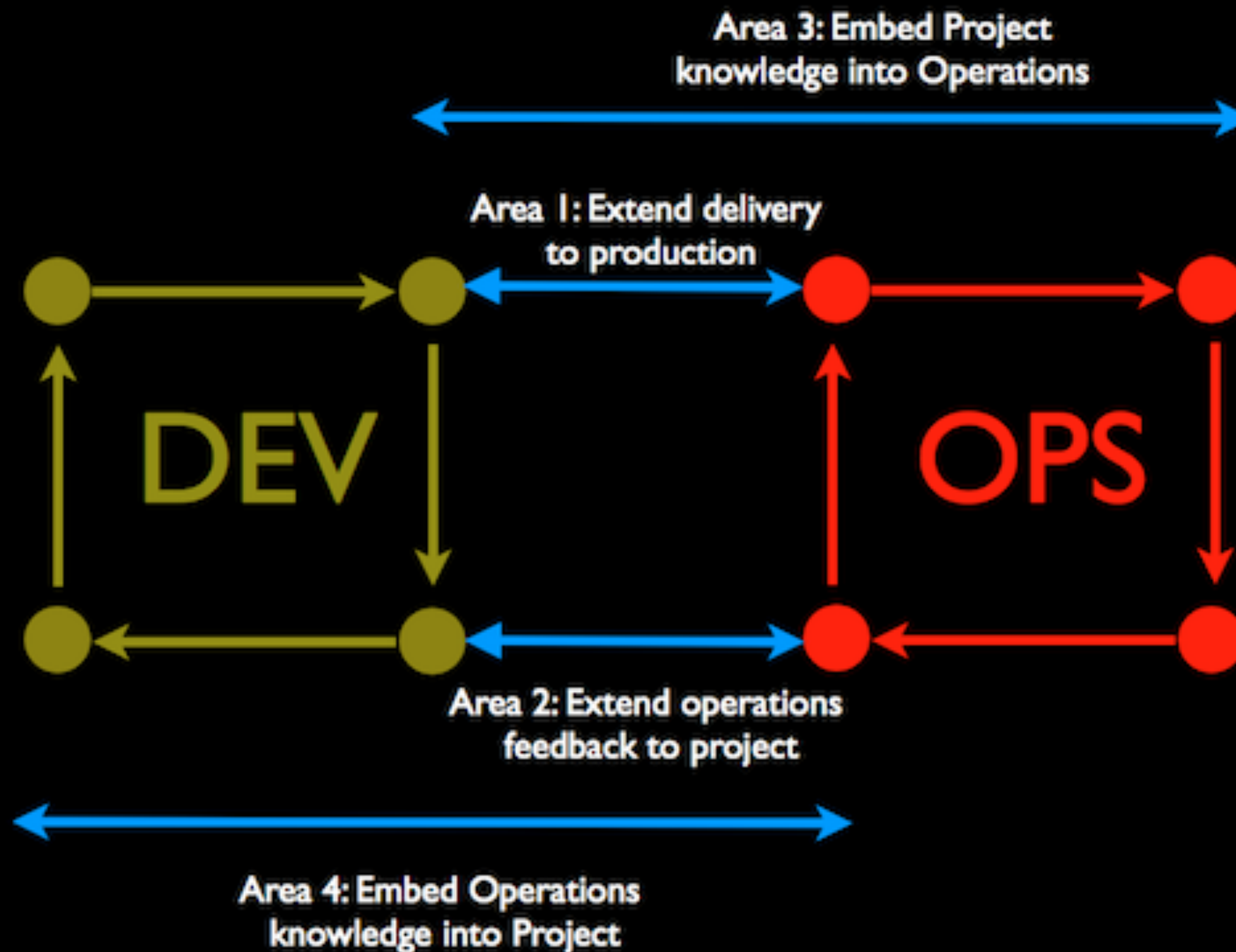## actions to add to current scope of work

**Add information on queue to the status page**

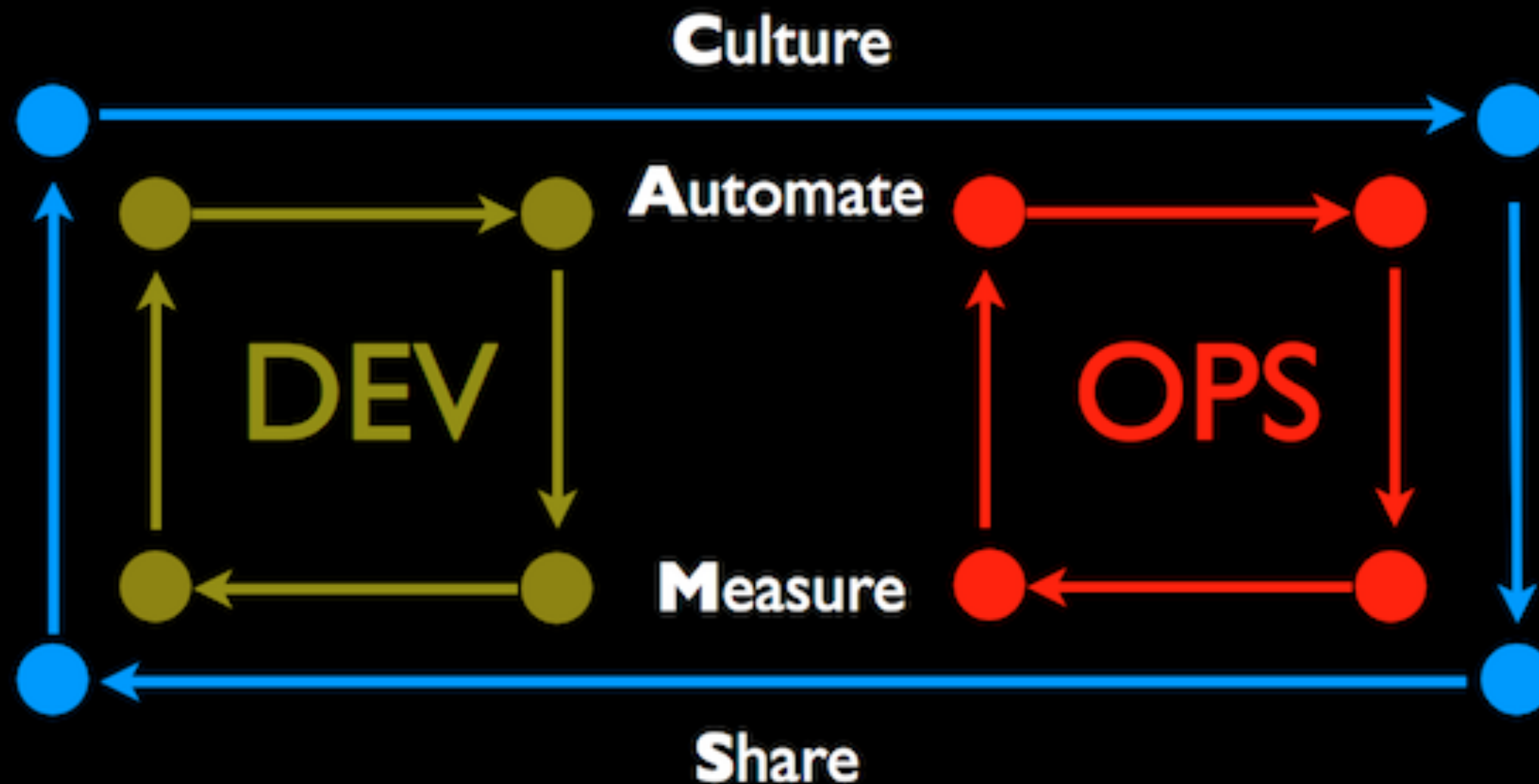**Have clients add the IP address of our new connect machine to their internal whitelist**

**Send warning/error information to Nagios**

**Add metrics around various Connect states**

Patrick Debois: Devops Areas - Codifying devops practices
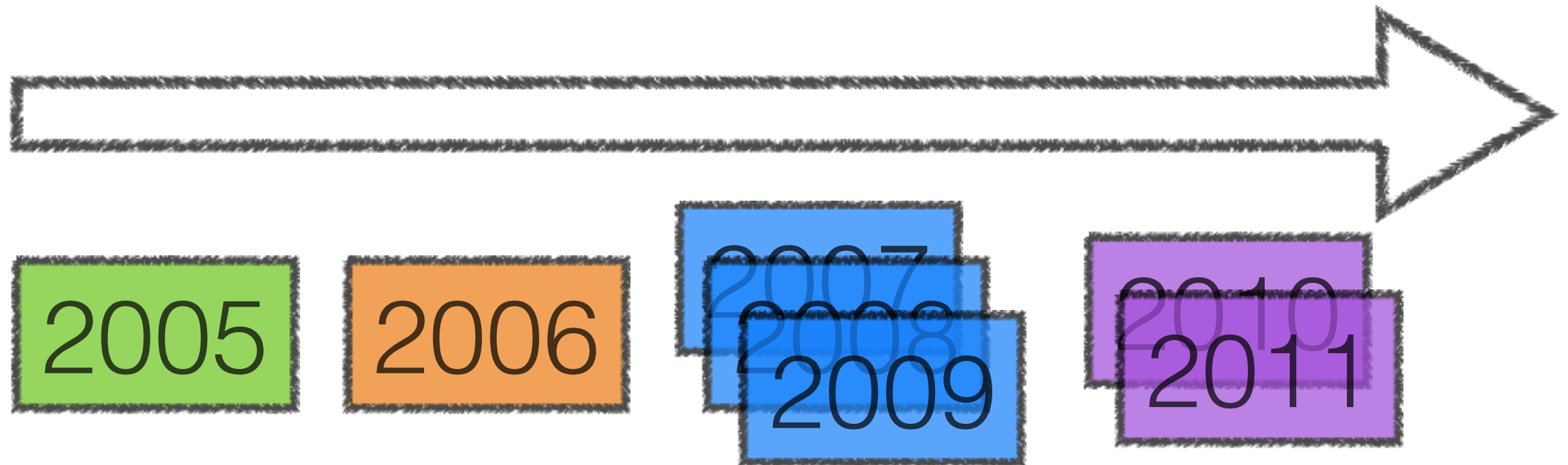http://is.gd/lyicB3

Patrick Debois: Devops Areas - Codifying devops practices
http://is.gd/IyicB3

# Questions?

jtf@jeffreyfredrick.com

2012

2005 2006 2007 2008 2009 2010 2011